

# 000 001 002 003 004 005 DIFFUSION MODULATION VIA ENVIRONMENT MECH- 006 ANISM MODELING FOR PLANNING 007 008 009

010 **Anonymous authors**  
011 Paper under double-blind review  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023

## ABSTRACT

024 Diffusion models have shown promising capabilities in trajectory generation for  
025 planning in offline reinforcement learning (RL). However, conventional diffusion-  
026 based planning methods often fail to account for the fact that generating trajectories  
027 in RL requires unique consistency between transitions to ensure coherence  
028 in real environments. This oversight can result in considerable discrepancies  
029 between the generated trajectories and the underlying mechanisms of a real envi-  
030 ronment. To address this problem, we propose a novel diffusion-based planning  
031 method, termed as Diffusion Modulation via Environment Mechanism Modeling  
032 (DMEMM). DMEMM modulates diffusion model training by incorporating key  
033 RL environment mechanisms, particularly transition dynamics and reward func-  
034 tions. Experimental results demonstrate that DMEMM achieves state-of-the-art  
035 performance for planning with offline reinforcement learning.  
036  
037

## 1 INTRODUCTION

038 Offline reinforcement learning (RL) has garnered significant attention for its potential to leverage  
039 pre-collected datasets to learn effective policies without requiring further interaction with the envi-  
040 ronment (Levine et al., 2020). One emerging approach within this domain is the use of diffusion  
041 models for trajectory generation (Janner et al., 2022). Diffusion models (Sohl-Dickstein et al., 2015;  
042 Ho et al., 2020), initially popularized for tasks such as image synthesis, have demonstrated promis-  
043 ing capabilities in generating coherent and diverse trajectories for planning in offline RL settings  
044 (Janner et al., 2022; Ni et al., 2023; Li, 2024; Goyal and Grand-Clement, 2023). Nevertheless, the  
045 essential differences between mechanisms in image synthesis and RL necessitate specific consider-  
046 ations for the effective application of diffusion models in RL.  
047

048 In image synthesis (Ho et al., 2020), diffusion models primarily aim to produce visually coherent  
049 outputs consistent in style and structure, while RL tasks demand environment and task oriented  
050 consistency between transitions in the generated trajectories (Janner et al., 2022) to ensure that  
051 the generated sequences are not only plausible but also effective for policy learning (Kumar et al.,  
052 2020). This consistency is essential for ensuring that the sequence of actions within the generated  
053 trajectories can successfully guide the RL agent from the current state to the target state. However,  
054 conventional diffusion-based planning methods often overlook this need for transition coherence  
055 (Janner et al., 2022). By simply adopting traditional diffusion models like DDPM, which utilize  
056 a fixed isotropic variance for Gaussian distributions, such diffusion-based planning models may  
057 fail to adequately capture the transition dynamics necessary for effective RL, leading to inaccurate  
058 trajectories and suboptimal learned policies (Wu et al., 2019; Niu et al., 2024).  
059

060 To address this problem, we introduce a novel diffusion-based planning method called Diffusion  
061 Modulation via Environment Mechanism Modeling (DMEMM). This method modulates the diffu-  
062 sion process by integrating RL-specific environment mechanisms, particularly transition dynamics  
063 and reward functions, directly into the diffusion model training process on offline data, thereby en-  
064 hancing the diffusion model to better capture the underlying transition and reward structures of the  
065 offline data. Specifically, we modify the diffusion loss by weighting it with the cumulative reward,  
066 which biases the diffusion model towards high-reward trajectories, and introduce two auxiliary mod-  
067 ulation losses based on empirical transition and reward models to regularize the trajectory diffu-  
068 sion process, ensuring that the generated trajectories are not only plausible but also reward-optimized.  
069 Additionally, we also utilize the transition and reward models to guide the sampling process dur-  
070

054  
 055  
 056  
 057  
 058  
 059  
 060  
 ing planning trajectory generation from the learned diffusion model, further aligning the outputs  
 with the desired transition dynamics and reward structures. We conducted experiments on multiple  
 RL environments. Experimental results indicate that our proposed method achieves state-of-the-art  
 performance compared to previous diffusion-based planning approaches.

061  
 062  
 063  
 064  
 065  
 066  
 067  
 068  
 This work presents a significant step forward in the application of diffusion models for trajectory  
 generation in offline RL. The main contributions can be summarized as follows:

069  
 070  
 071  
 072  
 073  
 074  
 075  
 076  
 077  
 078  
 079  
 080  
 081  
 082  
 083  
 084  
 085  
 086  
 087  
 088  
 089  
 • We identify a critical problem in conventional diffusion model training for offline RL plan-  
 ning, where fixed isotropic variance and disregard for rewards may lead to a mismatch be-  
 tween generated trajectories and those desirable for RL. To address this, we propose a novel  
 method called Diffusion Modulation via Environment Mechanism Modeling (DMMEM).

• We incorporate RL-specific environment mechanisms, including transition dynamics and  
 reward functions, into diffusion model training through loss modulation, enhancing the  
 quality and consistency of the generated trajectories in a principled manner and providing  
 a fundamental framework for adapting diffusion models to offline RL tasks.

• Our results on multiple RL environments show that the proposed method achieves state-of-  
 the-art performance in offline RL planning, validating the effectiveness of our approach.

## 090 2 RELATED WORKS

### 091 2.1 OFFLINE REINFORCEMENT LEARNING

092  
 093  
 094  
 095  
 096  
 097  
 098  
 099  
 100  
 101  
 102  
 103  
 104  
 105  
 106  
 107  
 Offline reinforcement learning (RL) has gained significant traction in recent years, with various  
 approaches proposed to address the challenges of learning from static datasets without online envi-  
 ronment interactions. Fujimoto et al. (2019) introduced Batch Constrained Q-Learning (BCQ) that  
 learns a perturbation model to constrain the policy to stay close to the data distribution, mitigating the  
 distributional shift issue. Wu et al. (2019) conducted Behavior Regularized Offline Reinforcement  
 Learning (BRAC) that incorporates behavior regularization into actor-critic methods to prevent the  
 policy from deviating too far from the data distribution. Conservative Q-Learning (CQL) by Kumar  
 et al. (2020) uses a conservative Q-function to underestimate out-of-distribution actions, preventing  
 the policy from exploring unseen state-action regions. Kostrikov et al. (2022) conducted Implicit  
 Q-Learning (IQL) to directly optimize the policy to match the expected Q-values under the data  
 distribution. Goyal and Grand-Clement (2023) introduce Robust MDPs to formulate offline RL as a  
 robust optimization problem over the uncertainty in the dynamics model. Planning has emerged as a  
 powerful tool for solving offline RL tasks. MOPO (Yu et al., 2020) incorporates uncertainty-aware  
 planning into offline RL by penalizing simulated trajectories that deviate from the offline dataset.

### 108 2.2 DIFFUSION MODEL IN REINFORCEMENT LEARNING

109  
 110  
 111  
 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  
 138  
 139  
 140  
 141  
 142  
 143  
 144  
 145  
 146  
 147  
 148  
 149  
 150  
 151  
 152  
 153  
 154  
 155  
 156  
 157  
 158  
 159  
 160  
 161  
 162  
 163  
 164  
 165  
 166  
 167  
 168  
 169  
 170  
 171  
 172  
 173  
 174  
 175  
 176  
 177  
 178  
 179  
 180  
 181  
 182  
 183  
 184  
 185  
 186  
 187  
 188  
 189  
 190  
 191  
 192  
 193  
 194  
 195  
 196  
 197  
 198  
 199  
 200  
 201  
 202  
 203  
 204  
 205  
 206  
 207  
 208  
 209  
 210  
 211  
 212  
 213  
 214  
 215  
 216  
 217  
 218  
 219  
 220  
 221  
 222  
 223  
 224  
 225  
 226  
 227  
 228  
 229  
 230  
 231  
 232  
 233  
 234  
 235  
 236  
 237  
 238  
 239  
 240  
 241  
 242  
 243  
 244  
 245  
 246  
 247  
 248  
 249  
 250  
 251  
 252  
 253  
 254  
 255  
 256  
 257  
 258  
 259  
 260  
 261  
 262  
 263  
 264  
 265  
 266  
 267  
 268  
 269  
 270  
 271  
 272  
 273  
 274  
 275  
 276  
 277  
 278  
 279  
 280  
 281  
 282  
 283  
 284  
 285  
 286  
 287  
 288  
 289  
 290  
 291  
 292  
 293  
 294  
 295  
 296  
 297  
 298  
 299  
 300  
 301  
 302  
 303  
 304  
 305  
 306  
 307  
 308  
 309  
 310  
 311  
 312  
 313  
 314  
 315  
 316  
 317  
 318  
 319  
 320  
 321  
 322  
 323  
 324  
 325  
 326  
 327  
 328  
 329  
 330  
 331  
 332  
 333  
 334  
 335  
 336  
 337  
 338  
 339  
 340  
 341  
 342  
 343  
 344  
 345  
 346  
 347  
 348  
 349  
 350  
 351  
 352  
 353  
 354  
 355  
 356  
 357  
 358  
 359  
 360  
 361  
 362  
 363  
 364  
 365  
 366  
 367  
 368  
 369  
 370  
 371  
 372  
 373  
 374  
 375  
 376  
 377  
 378  
 379  
 380  
 381  
 382  
 383  
 384  
 385  
 386  
 387  
 388  
 389  
 390  
 391  
 392  
 393  
 394  
 395  
 396  
 397  
 398  
 399  
 400  
 401  
 402  
 403  
 404  
 405  
 406  
 407  
 408  
 409  
 410  
 411  
 412  
 413  
 414  
 415  
 416  
 417  
 418  
 419  
 420  
 421  
 422  
 423  
 424  
 425  
 426  
 427  
 428  
 429  
 430  
 431  
 432  
 433  
 434  
 435  
 436  
 437  
 438  
 439  
 440  
 441  
 442  
 443  
 444  
 445  
 446  
 447  
 448  
 449  
 450  
 451  
 452  
 453  
 454  
 455  
 456  
 457  
 458  
 459  
 460  
 461  
 462  
 463  
 464  
 465  
 466  
 467  
 468  
 469  
 470  
 471  
 472  
 473  
 474  
 475  
 476  
 477  
 478  
 479  
 480  
 481  
 482  
 483  
 484  
 485  
 486  
 487  
 488  
 489  
 490  
 491  
 492  
 493  
 494  
 495  
 496  
 497  
 498  
 499  
 500  
 501  
 502  
 503  
 504  
 505  
 506  
 507  
 508  
 509  
 510  
 511  
 512  
 513  
 514  
 515  
 516  
 517  
 518  
 519  
 520  
 521  
 522  
 523  
 524  
 525  
 526  
 527  
 528  
 529  
 530  
 531  
 532  
 533  
 534  
 535  
 536  
 537  
 538  
 539  
 540  
 541  
 542  
 543  
 544  
 545  
 546  
 547  
 548  
 549  
 550  
 551  
 552  
 553  
 554  
 555  
 556  
 557  
 558  
 559  
 560  
 561  
 562  
 563  
 564  
 565  
 566  
 567  
 568  
 569  
 570  
 571  
 572  
 573  
 574  
 575  
 576  
 577  
 578  
 579  
 580  
 581  
 582  
 583  
 584  
 585  
 586  
 587  
 588  
 589  
 590  
 591  
 592  
 593  
 594  
 595  
 596  
 597  
 598  
 599  
 600  
 601  
 602  
 603  
 604  
 605  
 606  
 607  
 608  
 609  
 610  
 611  
 612  
 613  
 614  
 615  
 616  
 617  
 618  
 619  
 620  
 621  
 622  
 623  
 624  
 625  
 626  
 627  
 628  
 629  
 630  
 631  
 632  
 633  
 634  
 635  
 636  
 637  
 638  
 639  
 640  
 641  
 642  
 643  
 644  
 645  
 646  
 647  
 648  
 649  
 650  
 651  
 652  
 653  
 654  
 655  
 656  
 657  
 658  
 659  
 660  
 661  
 662  
 663  
 664  
 665  
 666  
 667  
 668  
 669  
 670  
 671  
 672  
 673  
 674  
 675  
 676  
 677  
 678  
 679  
 680  
 681  
 682  
 683  
 684  
 685  
 686  
 687  
 688  
 689  
 690  
 691  
 692  
 693  
 694  
 695  
 696  
 697  
 698  
 699  
 700  
 701  
 702  
 703  
 704  
 705  
 706  
 707  
 708  
 709  
 710  
 711  
 712  
 713  
 714  
 715  
 716  
 717  
 718  
 719  
 720  
 721  
 722  
 723  
 724  
 725  
 726  
 727  
 728  
 729  
 730  
 731  
 732  
 733  
 734  
 735  
 736  
 737  
 738  
 739  
 740  
 741  
 742  
 743  
 744  
 745  
 746  
 747  
 748  
 749  
 750  
 751  
 752  
 753  
 754  
 755  
 756  
 757  
 758  
 759  
 760  
 761  
 762  
 763  
 764  
 765  
 766  
 767  
 768  
 769  
 770  
 771  
 772  
 773  
 774  
 775  
 776  
 777  
 778  
 779  
 780  
 781  
 782  
 783  
 784  
 785  
 786  
 787  
 788  
 789  
 790  
 791  
 792  
 793  
 794  
 795  
 796  
 797  
 798  
 799  
 800  
 801  
 802  
 803  
 804  
 805  
 806  
 807  
 808  
 809  
 8010  
 8011  
 8012  
 8013  
 8014  
 8015  
 8016  
 8017  
 8018  
 8019  
 8020  
 8021  
 8022  
 8023  
 8024  
 8025  
 8026  
 8027  
 8028  
 8029  
 8030  
 8031  
 8032  
 8033  
 8034  
 8035  
 8036  
 8037  
 8038  
 8039  
 8040  
 8041  
 8042  
 8043  
 8044  
 8045  
 8046  
 8047  
 8048  
 8049  
 8050  
 8051  
 8052  
 8053  
 8054  
 8055  
 8056  
 8057  
 8058  
 8059  
 8060  
 8061  
 8062  
 8063  
 8064  
 8065  
 8066  
 8067  
 8068  
 8069  
 8070  
 8071  
 8072  
 8073  
 8074  
 8075  
 8076  
 8077  
 8078  
 8079  
 8080  
 8081  
 8082  
 8083  
 8084  
 8085  
 8086  
 8087  
 8088  
 8089  
 8090  
 8091  
 8092  
 8093  
 8094  
 8095  
 8096  
 8097  
 8098  
 8099  
 80100  
 80101  
 80102  
 80103  
 80104  
 80105  
 80106  
 80107  
 80108  
 80109  
 80110  
 80111  
 80112  
 80113  
 80114  
 80115  
 80116  
 80117  
 80118  
 80119  
 80120  
 80121  
 80122  
 80123  
 80124  
 80125  
 80126  
 80127  
 80128  
 80129  
 80130  
 80131  
 80132  
 80133  
 80134  
 80135  
 80136  
 80137  
 80138  
 80139  
 80140  
 80141  
 80142  
 80143  
 80144  
 80145  
 80146  
 80147  
 80148  
 80149  
 80150  
 80151  
 80152  
 80153  
 80154  
 80155  
 80156  
 80157  
 80158  
 80159  
 80160  
 80161  
 80162  
 80163  
 80164  
 80165  
 80166  
 80167  
 80168  
 80169  
 80170  
 80171  
 80172  
 80173  
 80174  
 80175  
 80176  
 80177  
 80178  
 80179  
 80180  
 80181  
 80182  
 80183  
 80184  
 80185  
 80186  
 80187  
 80188  
 80189  
 80190  
 80191  
 80192  
 80193  
 80194  
 80195  
 80196  
 80197  
 80198  
 80199  
 80200  
 80201  
 80202  
 80203  
 80204  
 80205  
 80206  
 80207  
 80208  
 80209  
 80210  
 80211  
 80212  
 80213  
 80214  
 80215  
 80216  
 80217  
 80218  
 80219  
 80220  
 80221  
 80222  
 80223  
 80224  
 80225  
 80226  
 80227  
 80228  
 80229  
 80230  
 80231  
 80232  
 80233  
 80234  
 80235  
 80236  
 80237  
 80238  
 80239  
 80240  
 80241  
 80242  
 80243  
 80244  
 80245  
 80246  
 80247  
 80248  
 80249  
 80250  
 80251  
 80252  
 80253  
 80254  
 80255  
 80256  
 80257  
 80258  
 80259  
 80260  
 80261  
 80262  
 80263  
 80264  
 80265  
 80266  
 80267  
 80268  
 80269  
 80270  
 80271  
 80272  
 80273  
 80274  
 80275  
 80276  
 80277  
 80278  
 80279  
 80280  
 80281  
 80282  
 80283  
 80284  
 80285  
 80286  
 80287  
 80288  
 80289  
 80290  
 80291  
 80292  
 80293  
 80294  
 80295  
 80296  
 80297  
 80298  
 80299  
 80300  
 80301  
 80302  
 80303  
 80304  
 80305  
 80306  
 80307  
 80308  
 80309  
 80310  
 80311  
 80312  
 80313  
 80314  
 80315  
 80316  
 80317  
 80318  
 80319  
 80320  
 80321  
 80322  
 80323  
 80324  
 80325  
 80326  
 80327  
 80328  
 80329  
 80330  
 80331  
 80332  
 80333  
 80334  
 80335  
 80336  
 80337  
 80338  
 80339  
 80340  
 80341  
 80342  
 80343  
 80344  
 80345  
 80346  
 80347  
 80348  
 80349  
 80350  
 80351  
 80352  
 80353  
 80354  
 80355  
 80356  
 80357  
 80358  
 80359  
 80360  
 80361  
 80362  
 80363  
 80364  
 80365  
 80366  
 80367  
 80368  
 80369  
 80370  
 80371  
 80372  
 80373  
 80374  
 80375  
 80376  
 80377  
 80378  
 80379  
 80380  
 80381  
 80382  
 80383  
 80384  
 80385  
 80386  
 80387  
 80388  
 80389  
 80390  
 80391  
 80392  
 80393  
 80394  
 80395  
 80396  
 80397  
 80398  
 80399  
 80400  
 80401  
 80402  
 80403  
 80404  
 80405  
 80406  
 80407  
 80408  
 80409  
 80410  
 80411  
 80412  
 80413  
 80414  
 80415  
 80416  
 80417  
 80418  
 80419  
 80420  
 80421  
 80422  
 80423  
 80424  
 80425  
 80426  
 80427  
 80428  
 80429  
 80430  
 80431  
 80432  
 80433  
 80434  
 80435  
 80436  
 80437  
 80438  
 80439  
 80440  
 80441  
 80442  
 80443  
 80444  
 80445  
 80446  
 80447  
 80448  
 80449  
 80450  
 80451  
 80452  
 80453  
 80454  
 80455  
 80456  
 80457  
 80458  
 80459  
 80460  
 80461  
 80462  
 80463  
 80464  
 80465  
 80466  
 80467  
 80468  
 80469  
 80470  
 80471  
 80472  
 80473  
 80474  
 80475  
 80476  
 80477  
 80478  
 80479  
 80480  
 80481  
 80482  
 80483  
 80484  
 80485  
 80486  
 80487  
 80488  
 80489  
 80490  
 80491  
 80492  
 80493  
 80494  
 80495  
 80496  
 80497  
 80498  
 80499  
 80500  
 80501  
 80502  
 80503  
 80504  
 80505  
 80506  
 80507  
 80508  
 80509  
 80510  
 80511  
 80512  
 80513  
 80514  
 80515  
 80516  
 80517  
 80518  
 80519  
 80520  
 80521  
 80522  
 80523  
 80524  
 80525  
 80526  
 80527  
 80528  
 80529  
 80530  
 80531  
 80532  
 80533  
 80534  
 80535  
 80536  
 80537  
 80538  
 80539  
 80540  
 80541  
 80542  
 80543  
 80544  
 80545  
 80546  
 80547  
 80548  
 80549  
 80550  
 80551  
 80552  
 80553  
 80554  
 80555  
 80556  
 80557  
 80558  
 80559  
 80560  
 80561  
 80562  
 80563  
 80564  
 80565  
 80566  
 80567  
 80568  
 80569  
 80570  
 80571  
 80572  
 80573  
 80574  
 80575  
 80576  
 80577  
 80578  
 80579  
 80580  
 80581  
 80582  
 80583  
 80584  
 80585  
 80586  
 80587  
 80588  
 80589  
 80590  
 80591  
 80592  
 805

108 **3 PRELIMINARIES**

110 Reinforcement learning (RL) (Sutton and Barto, 2018) can be modeled as a Markov Decision Process (MDP)  $M = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$  in a given environment, where  $\mathcal{S}$  denotes the state space,  $\mathcal{A}$  corresponds to the action space,  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  defines the transition dynamics,  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  represents the reward function, and  $\gamma$  is a discount factor. Offline RL aims to train an RL agent from an offline dataset  $\mathcal{D}$ , consisting of a collection of trajectories  $\{\tau_1, \tau_2, \dots, \tau_i, \dots\}$ , with each trajectory  $\tau_i = (s_0^i, a_0^i, r_0^i, s_1^i, a_1^i, r_1^i, \dots, s_T^i, a_T^i, r_T^i)$  sampled from the underlying MDP in the given environment. In particular, the task of planning in offline RL aims to generate planning trajectories from an initial state  $s_0$  by simulating action sequences  $a_{0:T}$  and predicting future states  $s_{0:T}$  based on those actions. The objective is to learn an optimal plan function such that the cumulative reward can be maximized when executing the plan under the underlying MDP of the given environment.

120 **3.1 PLANNING WITH DIFFUSION MODEL**

122 Diffusion probabilistic models, commonly known as “diffusion models” (Sohl-Dickstein et al., 123 2015; Ho et al., 2020), are a class of generative models that utilize a unique Markov chain 124 framework. When applied to planning in offline RL, the objective is to generate best planning 125 trajectories  $\{\tau\}$  by learning a diffusion model on the offline RL dataset  $\mathcal{D}$ .

127 **Trajectory Representation** In the diffusion model applied to RL planning, it is necessary to 128 predict both states and actions. Therefore, the trajectory representation in the model is in an image-like 129 matrix format. In particular, trajectories are represented as two-dimensional arrays (Janner et al., 130 2022), where each column corresponds to a state-action pair  $(s_t, a_t)$  of the trajectory:

$$\tau = \begin{bmatrix} s_0 & s_1 & \cdots & s_T \\ a_0 & a_1 & \cdots & a_T \end{bmatrix}$$

135 **Trajectory Diffusion** The diffusion model (Ho et al., 2020) comprises two primary processes: the 136 forward process and the reverse process. The forward process (diffusion process) is a Markov chain 137 characterized by  $q(\tau^k | \tau^{k-1})$  that gradually adds Gaussian noise at each time step  $k \in \{1, \dots, K\}$ , 138 starting from an initial clean trajectory sample  $\tau^0 \sim \mathcal{D}$ . The conditional probability is particularly 139 defined as a Gaussian probability density function, such as:

$$q(\tau^k | \tau^{k-1}) := \mathcal{N}(\tau^k; (1 - \beta_k)\tau^{k-1}, \beta_k \mathbf{I}), \quad (1)$$

140 with  $\{\beta_1, \dots, \beta_K\}$  representing a predefined variance schedule. By introducing  $\alpha_k := 1 - \beta_k$  and 141  $\bar{\alpha}_k := \prod_{i=1}^k \alpha_i$ , one can succinctly express the diffused sample at any time step  $k$  as:

$$\tau^k = \sqrt{\bar{\alpha}_k} \tau^0 + \sqrt{1 - \bar{\alpha}_k} \epsilon, \quad (2)$$

146 where  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . The reverse diffusion process is an iterative denoising procedure, and can be 147 modeled as a parametric Markov chain characterized by  $p_\theta(\tau^{k-1} | \tau^k)$ , starting from a Gaussian 148 noise prior  $\tau^K \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , such that:

$$p_\theta(\tau^{k-1} | \tau^k) = \mathcal{N}(\tau^{k-1}; \mu_\theta(\tau^k, k), \sigma_k^2 \mathbf{I}), \quad (3)$$

$$\text{with } \mu_\theta(\tau^k, k) = \frac{1}{\sqrt{\alpha_k}} \left( \tau^k - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k}} \epsilon_\theta(\tau^k, k) \right). \quad (4)$$

154 **Training** In the literature, the diffusion model is trained by predicting the additive noise  $\epsilon$  (Ho 155 et al., 2020) using the noise network  $\epsilon_\theta(\tau^k, k) = \epsilon_\theta(\sqrt{\bar{\alpha}_k} \tau^0 + \sqrt{1 - \bar{\alpha}_k} \epsilon, k)$ . The training loss is 156 expressed as the mean squared error between the additive noise  $\epsilon$  and the predicted noise  $\epsilon_\theta(\tau^k, k)$ :

$$L_{\text{diff}} = \mathbb{E}_{k \sim \mathcal{U}(1, K), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \tau^0 \sim \mathcal{D}} \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_k} \tau^0 + \sqrt{1 - \bar{\alpha}_k} \epsilon, k)\|^2 \quad (5)$$

160 where  $\mathcal{U}(1, K)$  denotes a uniform distribution over numbers in  $[1, 2, \dots, K]$ . With the trained noise 161 network, the diffusion model can be used to generate RL trajectories for planning through the reverse 162 diffusion process characterized by Eq.(3).

162 4 METHOD  
163

164 In this section, we present our proposed diffusion approach, Diffusion Modulation via Environment  
165 Mechanism Modeling (DMEMM), for planning in offline RL. This method integrates the essential  
166 transition and reward mechanisms of reinforcement learning into an innovative modulation-based  
167 diffusion learning framework, while maintaining isotropic covariance matrices for the diffusion  
168 Gaussian distributions to preserve the benefits of this conventional setup—simplifying model com-  
169 plexity, stabilizing training and enhancing performance. Additionally, the transition and reward  
170 mechanisms are further leveraged to guide the planning phase under the trained diffusion model,  
171 aiming to generate optimal planning trajectories that align well with both the underlying MDP of  
172 the environment and the objectives of RL.

173 4.1 MODULATION OF DIFFUSION TRAINING  
174

175 In an RL environment, the transition dynamics and reward function are two fundamental components  
176 of the underlying MDP. Directly applying conventional diffusion models to offline RL can lead to  
177 a mismatch between the generated trajectories and those optimal for the underlying MDP in RL.  
178 This is due to the use of isotropic covariance and the disregard for rewards in traditional diffusion  
179 models. To tackle this problem, we propose to modulate the diffusion model training by deploying  
180 a reward-aware diffusion loss and enforcing auxiliary regularizations on the generated trajectories  
181 based on environment transition and reward mechanisms.

182 Given the offline data  $\mathcal{D}$  collected from the RL environment, we first learn a probabilistic transition  
183 model  $\hat{T}(s_t, a_t)$  and a reward function  $\hat{\mathcal{R}}(s_t, a_t)$  from  $\mathcal{D}$  as regression functions to predict the next  
184 state  $s_{t+1}$  and the corresponding reward  $r_t$  respectively. These models can serve as estimations of  
185 the underlying MDP mechanisms. In order to regularize diffusion model training for generating  
186 desirable trajectories, using the learned transition model and reward function, we need to express  
187 the output trajectories of the reverse diffusion process in terms of the diffusion model parameters,  $\theta$ .  
188 To this end, we present the following proposition.

189 **Proposition 1.** *Given the reverse process encoded by Eq.(3) and Eq.(4) in the diffusion model,  
190 the output trajectory  $\hat{\tau}^0$  denoised from an intermediate trajectory  $\tau^k$  at step  $k$  has the following  
191 Gaussian distribution:*

$$192 \quad \hat{\tau}^0 \sim \mathcal{N}(\hat{\mu}_\theta(\tau^k, k), \hat{\sigma}^2 \mathbf{I}), \quad (6)$$

$$194 \quad \text{where } \hat{\mu}_\theta(\tau^k, k) = \frac{1}{\sqrt{\bar{\alpha}_k}} \tau^k - \sum_{i=1}^k \frac{1 - \alpha_i}{\sqrt{(1 - \bar{\alpha}_i) \prod_{j=1}^i \alpha_j}} \epsilon_\theta(\tau^i, i). \quad (7)$$

198 Conveniently, we can use the mean of the Gaussian distribution above directly as the most likely  
199 output trajectory, denoted as  $\hat{\tau}^0 = \hat{\mu}_\theta(\tau^k, k)$ . This allows us to express the denoised output trajec-  
200 tory explicitly in terms of the parametric noise network  $\epsilon_\theta$ , and thus the parameters  $\theta$  of the diffusion  
201 model. Moreover, by deploying Eq.(2), we can get rid of the latent  $\{\tau^1, \dots, \tau^k\}$  and re-express  $\hat{\tau}^0$   
202 as the following function of a sampled clean trajectory  $\tau^0$  and some random noise  $\epsilon$ :

$$203 \quad \hat{\tau}_\theta^0(\tau^0, k, \epsilon) = \tau^0 + \frac{\sqrt{1 - \bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} \epsilon - \sum_{i=1}^k \frac{1 - \alpha_i}{\sqrt{(1 - \bar{\alpha}_i) \prod_{j=1}^i \alpha_j}} \epsilon_\theta(\sqrt{\bar{\alpha}_i} \tau^0 + \sqrt{1 - \bar{\alpha}_i} \epsilon, i). \quad (8)$$

206 Next, we leverage this output trajectory function to modulate diffusion model training by developing  
207 novel auxiliary modulation losses.

209 4.1.1 TRANSITION-BASED DIFFUSION MODULATION  
210

211 As previously discussed, the deployment of a fixed isotropic variance in conventional diffusion  
212 models has the potential drawback of overlooking the underlying transition mechanisms of the RL  
213 environment. As a result, there can be potential mismatches between the transitions of generated  
214 trajectories and the underlying transition dynamics. Consequently, the RL agent may diverge from  
215 the expected states when executing the planning actions generated by the diffusion model, leading to  
poor planning performance. To address this problem, the first auxiliary modulation loss is designed

216 to minimize the discrepancy between the transitions in the generated trajectories from the diffusion  
 217 model and those predicted by the learned transition model  $\widehat{\mathcal{T}}$ , which encodes the underlying transi-  
 218 tion mechanism. Specifically, for each transition  $(s_t, a_t, s_{t+1})$  in a generated trajectory  $\widehat{\tau}_\theta^0(\tau^0, k, \epsilon)$ ,  
 219 we minimize the mean squared error between  $s_{t+1}$  and the predicted next state using the transition  
 220 model  $\widehat{\mathcal{T}}$ . This leads to the following transition-based diffusion modulation loss:  
 221

$$222 \quad L_{\text{tr}} = \mathbb{E}_{k \sim \mathcal{U}(1, K), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \tau^0 \sim \mathcal{D}} \left[ \sum_{(s_t, a_t, s_{t+1}) \in \widehat{\tau}_\theta^0(\tau^0, k, \epsilon)} \|s_{t+1} - \widehat{\mathcal{T}}(s_t, a_t)\|^2 \right] \quad (9)$$

226 Here, the expectation is taken over the uniform sampling of time step  $k$  from  $[1 : K]$ , the random  
 227 sampling of noise  $\epsilon$  from a standard Gaussian distribution, and the random sampling of input tra-  
 228 jectories from the offline training data  $\mathcal{D}$ . Through function  $\widehat{\tau}_\theta^0$ , this loss  $L_{\text{tr}}$  is a function of the  
 229 diffusion model parameters  $\theta$ . By minimizing this transition-based modulation loss, we enforce that  
 230 the generated trajectories from the diffusion model are consistent with the transition dynamics ex-  
 231 pressed in the offline dataset. This approach enhances the fidelity of the generated trajectories and  
 232 improves the overall performance of the diffusion model in offline reinforcement learning tasks.  
 233

#### 4.1.2 REWARD-BASED DIFFUSION MODULATION

235 The goal of planning is to generate trajectories that maximize cumulative rewards when executed  
 236 under the underlying MDP of the given environment. Thus, focusing solely on the fit of the planning  
 237 trajectories to the transition dynamics is insufficient. It is crucial to guide the diffusion model train-  
 238 ing to directly align with the planning objective. Therefore, the second auxiliary modulation loss is  
 239 designed to maximize the reward induced in the generated trajectories. As the trajectories generated  
 240 from diffusion models do not have reward signals, we predict the reward scores of the state-action  
 241 pairs  $\{(s_t, a_t)\}$  in each trajectory generated through function  $\widehat{\tau}_\theta^0(\cdot, \cdot, \cdot)$  using the learned reward  
 242 function  $\widehat{\mathcal{R}}(\cdot, \cdot)$ . Specifically, we formulate the reward-based diffusion modulation loss function as  
 243 the following negative expected trajectory-wise cumulative reward from the generated trajectories:

$$244 \quad L_{\text{rd}} = -\mathbb{E}_{k \sim \mathcal{U}(1, K), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \tau^0 \sim \mathcal{D}} \left[ \sum_{(s_t, a_t) \in \widehat{\tau}_\theta^0(\tau^0, k, \epsilon)} \widehat{\mathcal{R}}(s_t, a_t) \right] \quad (10)$$

247 Through function  $\widehat{\tau}_\theta^0$ , this loss  $L_{\text{rd}}$  again is a function of the diffusion model parameters  $\theta$ . By  
 248 computing the expected loss over different time steps  $k \in [1 : K]$ , different random noise  $\epsilon$ , and all  
 249 input trajectories from the offline dataset  $\mathcal{D}$ , we ensure that the modulation is consistently enforced  
 250 across all instances of diffusion model training.  
 251

252 By minimizing this reward-based loss, we ensure that the generated trajectories are not only plau-  
 253 sible but also reward-optimized to align with the reward structure inherent in the offline data. This  
 254 approach improves the quality of the trajectories generated from the diffusion model and enhances  
 255 the overall policy learning process in offline reinforcement learning tasks.  
 256

#### 4.1.3 REWARD-AWARE DIFFUSION LOSS

258 In addition to the auxiliary modulation losses, we propose to further align diffusion model training  
 259 with the goal of RL planning by devising a novel reward-aware diffusion loss to replace the origi-  
 260 nal one. The original diffusion loss (shown in Eq.(5)) minimizes the expected per-trajectory mean  
 261 squared error between the true additive noise and the predicted noise, which gives equal weights  
 262 to different training trajectories without differentiation. In contrast, we propose to weight each tra-  
 263 jectory instance  $\tau^0$  from the offline dataset  $\mathcal{D}$  using its normalized cumulative reward, so that the  
 264 diffusion training can focus more on the more informative trajectory instances with larger cumula-  
 265 tive rewards. Specifically, we weight each training trajectory  $\tau^0$  using its normalized cumulative  
 266 reward and formulate the following reward-aware diffusion loss:  
 267

$$268 \quad L_{\text{wdiff}} = \mathbb{E}_{k \sim \mathcal{U}(1, K), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \tau^0 \sim \mathcal{D}} \left[ \left( \sum_{(s_t, a_t) \in \tau^0} \frac{\mathcal{R}(s_t, a_t)}{T_{\max} \cdot r_{\max}} \right) \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_k} \tau^0 + \sqrt{1 - \bar{\alpha}_k} \epsilon, k)\|^2 \right] \quad (11)$$

---

270           **Algorithm 1** Diffusion Training  
 271           **Require:** Offline data  $\mathcal{D} = \{(s_0^i, a_0^i, r_0^i, \dots, s_T^i, a_T^i, r_T^i)\}$ .  
 272           Learn transition model  $\hat{T}(s_t, a_t)$  and reward function  $\hat{\mathcal{R}}(s_t, a_t)$  from offline data  $\mathcal{D}$ .  
 273           Initialize noise network  $\epsilon_\theta(\tau^k, k)$ .  
 274           **while** not converged **do**  
 275            Sample a trajectory from offline data  $\tau^0 \sim \mathcal{D}$ .  
 276            Sample a random diffusion step  $k \sim \mathcal{U}(1, K)$ .  
 277            Sample a random noise  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ .  
 278            Calculate the gradient  $\nabla_\theta L_{\text{total}}$  of Eq. (12) and take gradient descent step.  
 279           **end while**  
 280  
 281

---

282           Here,  $\sum_{(s_t, a_t) \in \tau^0} \mathcal{R}(s_t, a_t)$  is the trajectory-wise cumulative reward on the original offline data  
 283           instance  $\tau^0 \in \mathcal{D}$ ;  $T_{\text{max}}$  denotes the largest trajectory length and  $r_{\text{max}}$  denotes the maximum possible  
 284           per-step reward. By using  $T_{\text{max}} \cdot r_{\text{max}}$  as the normalizer, we scale the cumulative reward to a ratio  
 285           within  $(0, 1]$  to weight the corresponding per-trajectory diffusion loss. This weighting mechanism  
 286           biases the diffusion model toward high-reward trajectories, ensuring that those trajectories yielding  
 287           higher cumulative rewards are more accurately represented, thus aligning diffusion training with  
 288           the planning objectives in offline RL. This approach improves the model’s performance on rare but  
 289           valuable trajectories, which are crucial for effective RL.

290           4.1.4 FULL MODULATION FRAMEWORK

292           The proposed full modulated diffusion model comprises all of the three loss components presented  
 293           above: the reward-aware diffusion loss  $L_{\text{wdiff}}$ , the transition-based auxiliary modulation loss  $L_{\text{tr}}$ , and  
 294           the reward-based auxiliary modulation loss  $L_{\text{rd}}$ . By integrating these loss terms together, we have  
 295           the following total loss for modulated diffusion training:

296           
$$L_{\text{total}} = L_{\text{wdiff}} + \lambda_{\text{tr}} L_{\text{tr}} + \lambda_{\text{rd}} L_{\text{rd}}, \quad (12)$$

298           where  $\lambda_{\text{tr}}$  and  $\lambda_{\text{rd}}$  are trade-off parameters that balance the contributions of the transition-based and  
 299           reward-based auxiliary losses, respectively. Standard diffusion training algorithm can be utilized to  
 300           train the model  $\theta$  by minimizing this total loss function. By employing this integrated loss function,  
 301           we establish a comprehensive modulation framework that incorporates essential domain and task  
 302           knowledge into diffusion model training, offering a general capacity of enhancing the adaptation  
 303           and broadening the applicability of diffusion models.

304           4.1.5 DIFFUSION TRAINING ALGORITHM

306           The complete training process of the diffusion model is presented in Algorithm 1. Prior to training  
 307           the diffusion model, a probabilistic transition model  $\hat{T}(s_t, a_t)$  and a reward model  $\hat{\mathcal{R}}(s_t, a_t)$  are  
 308           learned from the offline dataset  $\mathcal{D}$ . Afterward, the noise network is initialized and iteratively trained.  
 309           During each iteration, an original trajectory  $\tau^0$  is sampled from the offline dataset  $\mathcal{D}$ , along with a  
 310           randomly selected diffusion step  $k$  and noise sample  $\epsilon$ . Gradient descent is then applied to minimize  
 311           the total loss  $L_{\text{total}}$ .

312           4.2 PLANNING WITH DUAL GUIDANCE

315           Once trained, the diffusion model can be used to generate trajectories for planning during an RL  
 316           agent’s online interactions with the environment. The generation procedure starts from an initial  
 317           noise trajectory  $\tau^K \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and gradually denoises it by following the reverse diffusion process  
 318            $\tau^{k-1} \sim \mathcal{N}(\mu^{k-1}, \sigma_k^2 \mathbf{I})$  for each time step  $k \in \{K, K-1, \dots, 1\}$ , where  $\mu^{k-1}$  is estimated through  
 319           Eq. (4). In each diffusion time step  $k$ , the first state  $s_0$  of the trajectory  $\tau^k$  is fixed to the current  
 320           state  $s$  of the RL agent in the online environment to ensure the plan starts from it. The denoised  
 321           trajectory  $\tau^0$  after  $K$  diffusion time steps is treated as the plan for the RL agent, which is intended  
 322           to maximize the RL agent’s long-term performance without extra interaction with the environment.

323           To further enhance the objective of planning, some previous work (Janner et al., 2022) has utilized  
 324           the learned reward function to guide the sampling process of planning. In this work, we propose to

---

324 **Algorithm 2** Planning with Dual Guidance325 **Require:** Noise network  $\epsilon_\theta$ , tradeoff parameter  $\alpha$ , environment ENV, covariances  $\{\sigma_k^2\}$ .326 Initialize environment step  $t = 0$ .327 **while** not finished **do**328 Initialize noise trajectory  $\tau_t^K$ :  $\tau_t^K \sim \mathcal{N}(0, \mathbf{I})$ .329 **for** diffusion step  $k = K, \dots, 1$  **do**330 Compute the mean  $\mu^{k-1}$  using Eq. (4).331 Compute the guidance  $\mathbf{g}$  using Eq. (14).332 Sample next trajectory  $\tau_t^{k-1}$  with Eq.(13)333 Set current state  $s_t$  to the trajectory:  $\tau_t^{k-1}(s_0) = s_t$ .334 **end for**335 Execute the first action of plan  $\tau_t^0(a_0)$ :  $s_{t+1} = \text{ENV}(s_t, \tau_t^0(a_0))$ 336 Increment environment step by 1:  $t = t + 1$ 337 **end while**338  
339 deploy dual guidance for each reverse diffusion step  $k$  by exploiting both the reward function  $\widehat{\mathcal{R}}$  and  
340 the transition model  $\widehat{\mathcal{T}}$  learned from the offline dataset  $\mathcal{D}$ . Following previous works on conditional  
341 reverse diffusion (Dhariwal and Nichol, 2021), we incorporate the dual guidance by perturbing the  
342 mean of the Gaussian distribution  $\mathcal{N}(\mu^{k-1}, \sigma_k^2 \mathbf{I})$  used for reverse diffusion sampling. Specifically,  
343 we integrate the gradient  $\mathbf{g}$  of the linear combination of the reward function and transition function  
344 w.r.t the trajectory into  $\mu^{k-1}$ , such that:

345  
346 
$$\tau^{k-1} \sim \mathcal{N}(\mu^{k-1} + \alpha \sigma_k^2 \mathbf{I} \mathbf{g}, \sigma_k^2 \mathbf{I}) \quad (13)$$

347 where  $\mathbf{g}$  is computed as:

348  
349 
$$\mathbf{g} = \sum_{t=0}^T \nabla_{(s_t, a_t)} \widehat{\mathcal{R}}(s_t, a_t) + \lambda \sum_{t=0}^{T-1} \nabla_{(s_t, a_t)} \log \widehat{\mathcal{T}}(s_{t+1} | s_t, a_t) \quad (14)$$

350 where  $\alpha$  is a tradeoff parameter that controls the degree of guidance. By incorporating both the  
351 reward and transition guidance, we aim to enhance the planning process to generate high-quality  
352 trajectories that are both reward-optimized and transition-consistent, improving the overall planning  
353 performance. The details of the proposed planning procedure is summarized in Algorithm 2.354 

## 5 EXPERIMENT

355 In this section, we present the experimental setup and results for evaluating our proposed method,  
356 DMEMM, across various offline RL tasks. We conduct experiments on the D4RL locomotion suite  
357 and Maze2D environments to assess the performance of DMEMM compared to several state-of-the-  
358 art methods. The experiments are designed to demonstrate the effectiveness of our approach across  
359 different tasks, expert levels, and complex navigation scenarios.360 **Environments** We conduct our experiments on D4RL (Fu et al., 2020) tasks to evaluate the performance  
361 of planning in offline RL settings. Initially, we focus on the D4RL locomotion suite to assess  
362 the general performance of our planning methods across different tasks and expert levels of demon-  
363 strations. The RL agents are tested on three different tasks: HalfCheetah, Hopper, and Walker2d,  
364 and three different levels of expert demonstrations: Med-Expert, Medium, and Med-Replay. We use  
365 the normalized scores provided in the D4RL (Fu et al., 2020) benchmarks to evaluate performance.  
366 Subsequently, we conduct experiments on Maze2D (Fu et al., 2020) environments to evaluate per-  
367 formance on maze navigation tasks.368  
369 **Comparison Methods** We benchmark our methods against several leading approaches in each  
370 task domain, including Model Predictive Path Integral (MPPI) (Williams et al., 2016), Batch-  
371 Constrained Deep Q-Learning (BCQ) (Fujimoto et al., 2019), Conservative Q-Learning (CQL) (Ku-  
372 mar et al., 2020), Implicit Q-Learning (IQL) (Kostrikov et al., 2022), and Decision Transformer  
373 (DT) (Chen et al., 2021). Additionally, we compare our methods with the state-of-the-art offline  
374 RL approach, Selecting from Behavior Candidates (SfBC) (Chen et al., 2023), as well as several

378

379  
Table 1: This table presents the scores on D4RL locomotion suites for various comparison methods.  
380  
Results are averaged over 5 seeds.

Gym Tasks	BC	DT	IQL	CQL	SfBC	LDCQ	Diffuser	DD	HDMI	HD-DA	DMEMM (Ours)
HalfCheetah (Med-Expert)	55.2	86.8	86.7	91.6	92.6 $\pm$ 0.5	90.2 $\pm$ 0.9	88.9 $\pm$ 0.3	90.6 $\pm$ 1.3	92.1 $\pm$ 1.4	92.5 $\pm$ 0.3	<b>94.6<math>\pm</math>1.2</b>
Hopper (Med-Expert)	52.5	107.6	91.5	105.4	108.6 $\pm$ 2.1	109.3 $\pm$ 0.4	103.3 $\pm$ 1.3	111.8 $\pm$ 1.8	113.5 $\pm$ 0.9	115.3 $\pm$ 1.1	<b>115.9<math>\pm</math>1.6</b>
Walker2d (Med-Expert)	107.5	108.1	109.6	108.8	109.8 $\pm$ 0.2	111.3 $\pm$ 0.2	106.9 $\pm$ 0.2	108.8 $\pm$ 1.7	107.9 $\pm$ 1.2	107.1 $\pm$ 0.1	<b>111.6<math>\pm</math>1.1</b>
HalfCheetah (Medium)	42.6	42.6	47.4	44.0	45.9 $\pm$ 2.2	42.8 $\pm$ 0.7	42.8 $\pm$ 0.3	49.1 $\pm$ 1.0	48.0 $\pm$ 0.9	46.7 $\pm$ 0.2	<b>49.2<math>\pm</math>0.8</b>
Hopper (Medium)	52.9	67.6	66.3	58.5	57.1 $\pm$ 4.1	69.4 $\pm$ 3.5	74.3 $\pm$ 1.4	79.3 $\pm$ 3.6	76.4 $\pm$ 2.6	99.3 $\pm$ 0.3	<b>101.2<math>\pm</math>1.4</b>
Walker2d (Medium)	75.3	74.0	78.3	72.5	77.9 $\pm$ 2.5	66.2 $\pm$ 1.7	79.6 $\pm$ 0.6	82.5 $\pm$ 1.4	79.9 $\pm$ 1.8	84.0 $\pm$ 0.6	<b>86.5<math>\pm</math>1.5</b>
HalfCheetah (Med-Replay)	36.6	36.6	44.2	45.5	37.1 $\pm$ 1.7	41.8 $\pm$ 0.4	37.7 $\pm$ 0.5	39.3 $\pm$ 4.1	44.9 $\pm$ 2.0	38.1 $\pm$ 0.7	<b>46.1<math>\pm</math>1.3</b>
Hopper (Med-Replay)	18.1	82.7	94.7	95.0	86.2 $\pm$ 9.1	68.5 $\pm$ 4.3	93.6 $\pm$ 0.4	100.0 $\pm$ 0.7	99.6 $\pm$ 1.5	94.7 $\pm$ 0.7	<b>100.6<math>\pm</math>0.9</b>
Walker2d (Med-Replay)	26.0	66.6	73.9	77.2	65.1 $\pm$ 5.6	86.2 $\pm$ 2.5	70.6 $\pm$ 1.6	75.0 $\pm$ 4.3	80.7 $\pm$ 2.1	84.1 $\pm$ 2.2	<b>85.8<math>\pm</math>2.6</b>
<b>Average</b>	51.9	74.7	77.0	77.6	75.6	76.2	77.5	81.8	82.6	84.6	<b>87.9</b>

389  
Table 2: This table presents the scores on Maze2D navigation tasks for various comparison methods.  
390  
Results are averaged over 5 seeds.

Environment	MPPI	IQL	Diffuser	HDMI	HD-DA	DMEMM (Ours)
Maze2D U-Maze	33.2	47.4	113.9 $\pm$ 3.1	120.1 $\pm$ 2.5	128.4 $\pm$ 3.6	<b>132.4<math>\pm</math>3.0</b>
Maze2D Medium	10.2	34.9	121.5 $\pm$ 2.7	121.8 $\pm$ 1.6	135.6 $\pm$ 3.0	<b>138.2<math>\pm</math>2.2</b>
Maze2D Large	5.1	58.6	123.0 $\pm$ 6.4	128.6 $\pm$ 2.9	<b>155.8<math>\pm</math>2.5</b>	153.2 $\pm$ 3.3
Multi2D U-Maze	41.2	24.8	128.9 $\pm$ 1.8	131.3 $\pm$ 1.8	144.1 $\pm$ 1.2	<b>145.6<math>\pm</math>2.6</b>
Multi2D Medium	15.4	12.1	127.2 $\pm$ 3.4	131.6 $\pm$ 1.9	140.2 $\pm$ 1.6	<b>140.8<math>\pm</math>2.2</b>
Multi2D Large	8.0	13.9	132.1 $\pm$ 5.8	135.4 $\pm$ 2.5	<b>165.5<math>\pm</math>0.6</b>	159.6 $\pm$ 3.8
AntMaze U-Maze	–	62.2	76.0 $\pm$ 7.6	86.1 $\pm$ 2.4	94.0 $\pm$ 4.9	<b>96.2<math>\pm</math>5.5</b>
AntMaze Medium	–	70.0	31.9 $\pm$ 5.1	–	88.7 $\pm$ 8.1	<b>90.1<math>\pm</math>6.4</b>
AntMaze Large	–	47.5	0.0 $\pm$ 0.0	71.5 $\pm$ 3.5	<b>83.6<math>\pm</math>5.8</b>	79.6 $\pm$ 7.7

400 diffusion-based offline RL methods, including Diffuser (Janner et al., 2022), Decision Diffuser  
401 (DD) (Ajay et al., 2023), Latent Diffusion-Constrained Q-learning (LDCQ / LDGC) (Venkatraman et al., 2024), Hierarchical Diffusion for Offline Decision Making (HDMI) (Li et al., 2023), and  
402 Hierarchical Diffuser with Dense Actions (HD-DA) (Chen et al., 2024).

404  
405 

## 5.1 EXPERIMENTAL RESULTS ON D4RL

407 The experimental results summarized in Table 1 highlight the performance of various comparison  
408 methods across different Gym tasks, with scores averaged over 5 seeds. Our proposed method,  
409 DMEMM, consistently outperforms other methods across all tasks. Notably, in the HalfCheetah  
410 environments, DMEMM achieves a 2.0-point improvement on the Med-Expert dataset, and an 6.8-  
411 point improvement on the Med-Replay dataset compared to the previous best results. Additionally,  
412 DMEMM shows a 2.8-point increase on the Med-Expert Walker2D task, demonstrating that  
413 DMEMM effectively extracts valuable information, particularly from data that is not purely expert-  
414 level.

415 In most tasks, DMEMM outperforms the previous state-of-the-art method HD-DA, another variant  
416 of a Diffuser based planning method, by more than 2.0 points on average. Compared to Diffuser,  
417 DMEMM shows superior performance on all tasks, indicating that our method improves the  
418 consistency and optimality of diffusion model training in offline RL planning.

419 Overall, DMEMM achieves outstanding performance. With an average score of 87.9, DMEMM  
420 yields a substantial improvement over the second-highest average score of 84.6 achieved by HD-  
421 DA. These results clearly demonstrate the robustness and superiority of DMEMM in enhancing  
422 performance across various Gym tasks.

423  
424 

## 5.2 EXPERIMENTAL RESULTS ON MAZE2D

425 We present our experimental results on the Maze2D navigation tasks in Table 2, where the results  
426 are averaged over 5 seeds. The table shows that in all three environments, particularly at the U-Maze  
427 and Medium difficulty levels, our proposed DMEMM method significantly outperforms other com-  
428 parison methods. Specifically, on Maze2D tasks, DMEMM achieves a 4.0 point improvement over  
429 the state-of-the-art HD-DA method on the U-Maze task, and a 2.6 point increase on the Medium-  
430 sized maze. Compared to Diffuser, DMEMM shows an almost 20-point improvement. These results  
431 indicate that our method performs exceptionally well in generating planning solutions for navigation  
432 tasks.

432

433 Table 3: the scores on the Walker2D environment at three different levels for all four ablation vari-  
434 ants. Results are averaged over 5 seeds.

435

Gym Tasks	DMEMM	DMEMM-w/o-weighting	DMEMM-w/o- $\lambda_{tr}$	DMEMM-w/o- $\lambda_{rd}$	DMEMM-w/o-tr-guide
Med-Expert	<b>111.6±1.1</b>	110.4±0.8	108.4±1.2	110.4±0.6	109.9±1.0
WMedium	<b>86.5±1.5</b>	85.6±1.2	82.8±1.4	84.4±0.9	83.0±1.8
Med-Replay	<b>85.8±2.6</b>	84.6±2.2	82.2±1.7	83.7±2.5	82.6±3.2

439

440 However, HD-DA shows better performance on the large maze tasks. This is likely due to the hi-  
441 erarchical structure of HD-DA, which offers an advantage in larger, more complex environments  
442 by breaking long-horizon planning into smaller sub-tasks, an area where our method is not speci-  
443 fically designed to excel. Nevertheless, DMEMM remains competitive in larger environments, while  
444 demonstrating superior performance in smaller and medium-sized tasks.

445

## 5.3 ABLATION STUDY

447

448 We conduct an ablation study to evaluate the effectiveness of the key components in our DMEMM  
449 framework. We compare the full DMEMM model with four ablated variants: (1) DMEMM-w/o-  
450 weighting, which removes the weighting function in the reward-aware diffusion loss; (2) DMEMM-  
451 w/o- $\lambda_{tr}$ , which omits the transition-based diffusion modulation loss; (3) DMEMM-w/o- $\lambda_{rd}$ , which  
452 omits the reward-based diffusion modulation loss; and (4) DMEMM-w/o-tr-guide, which removes  
453 transition guidance in the dual-guided sampling procedure. The ablation study is conducted on all  
454 locomotion tasks across three levels of expert demonstrations, while only the Walker2D results are  
455 reported in the main paper. The complete results are provided in Table 5 of the Appendix. Table 3  
456 summarizes the performance of all four ablation variants on the D4RL locomotion benchmarks,  
457 averaged over five random seeds.

458

459 The results highlight the contribution of each component in DMEMM. Across all three difficulty  
460 levels, the full DMEMM model consistently achieves the best performance. In particular, remov-  
461 ing transition-related components, either the transition-based modulation loss (DMEMM-w/o- $\lambda_{tr}$ ) or the  
462 transition guidance (DMEMM-w/o-tr-guide), leads to substantial performance drops, under-  
463 scoring the importance of explicitly modeling transition dynamics in our approach. Incorporating  
464 transition information significantly improves the consistency and fidelity of generated trajectory  
465 plans. Moreover, DMEMM-w/o- $\lambda_{rd}$  and DMEMM-w/o-weighting yield comparable results, with  
466 DMEMM-w/o- $\lambda_{rd}$  showing a slightly larger degradation. This indicates that the designed reward  
467 model and its weighting mechanism play a key role in improving the optimality of planned trajec-  
468 tories.

469

470 Overall, the ablation study demonstrates that each component of our DMEMM method contributes  
471 significantly to its performance. Removing any of these components results in a noticeable de-  
472 crease in performance, highlighting the importance of the weighting function, transition-based and  
473 reward-based diffusion modulation loss, and transition guidance in achieving optimal results in of-  
474 fline reinforcement learning tasks.

475

## 6 CONCLUSION

476

477 In this work, we addressed a critical limitation of conventional diffusion-based planning methods in  
478 offline RL, which often overlook the consistency of transition dynamics in planned trajectories. To  
479 overcome this challenge, we proposed Diffusion Modulation via Environment Mechanism Model-  
480 ing (DMEMM), a novel approach that integrates RL-specific environment mechanisms, particularly  
481 transition dynamics and reward functions, into the diffusion model training process. By modulating  
482 the diffusion loss with cumulative rewards and introducing auxiliary losses based on transition  
483 dynamics and reward functions, DMEMM enhances both the coherence and quality of the generated  
484 trajectories, ensuring they are plausible and optimized for policy learning. Our experimental  
485 results across multiple offline RL environments demonstrate the effectiveness of DMEMM, achieving  
486 state-of-the-art performance compared to previous diffusion-based planning methods. The proposed  
487 approach significantly improves the alignment of generated trajectories, addressing the discrepan-  
488 cies between offline data and real-world environments. This provides a promising framework for  
489 further exploration of diffusion models in RL and their potential practical applications.

486 REFERENCES  
487

488 S. Levine, A. Kumar, G. Tucker, and J. Fu, “Offline reinforcement learning: Tutorial, review, and  
489 perspectives on open problems,” *arXiv preprint arXiv:2005.01643*, 2020.

490 M. Janner, Y. Du, J. B. Tenenbaum, and S. Levine, “Planning with diffusion for flexible behavior  
491 synthesis,” in *International Conference on Machine Learning (ICML)*, 2022.

492

493 J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, “Deep unsupervised learning  
494 using nonequilibrium thermodynamics,” in *Proceedings of the 32nd International Conference on  
495 Machine Learning*. PMLR, 2015.

496

497 J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” in *Advances in Neural  
498 Information Processing Systems (NeurIPS)*, 2020.

499

500 F. Ni, J. Hao, Y. Mu, Y. Yuan, Y. Zheng, B. Wang, and Z. Liang, “Metadiffuser: Diffusion model  
501 as conditional planner for offline meta-rl,” in *International Conference on Machine Learning  
502 (ICML)*. PMLR, 2023.

503

504 W. Li, “Efficient planning with latent diffusion,” in *The Twelfth International Conference on Learn-  
505 ing Representations (ICLR)*, 2024.

506

507 V. Goyal and J. Grand-Clement, “Robust markov decision process: Beyond rectangularity,” *Mathe-  
508 matics of Operations Research*, 2023.

509

510 A. Kumar, A. Zhou, G. Tucker, and S. Levine, “Conservative q-learning for offline reinforcement  
511 learning,” *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.

512

513 Y. Wu, G. Tucker, and O. Nachum, “Behavior regularized offline reinforcement learning,” *arXiv  
514 preprint arXiv:1911.11361*, 2019.

515

516 H. Niu, Q. Chen, T. Liu, J. Li, G. Zhou, Y. Zhang, J. Hu, and X. Zhan, “xted: Cross-domain  
517 adaptation via diffusion-based trajectory editing,” *arXiv preprint arXiv:2409.08687*, 2024.

518

519 S. Fujimoto, D. Meger, and D. Precup, “Off-policy deep reinforcement learning without explo-  
520 ration,” in *International Conference on Machine Learning (ICML)*. PMLR, 2019.

521

522 I. Kostrikov, A. Nair, and S. Levine, “Offline reinforcement learning with implicit q-learning,” in  
523 *International Conference on Learning Representations (ICLR)*, 2022.

524

525 T. Yu, G. Thomas, L. Yu, S. Ermon, J. Y. Zou, S. Levine, C. Finn, and T. Ma, “Mopo: Model-based  
526 offline policy optimization,” *Advances in Neural Information Processing Systems (NeurIPS)*,  
527 2020.

528

529 Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, “Score-based genera-  
530 tive modeling through stochastic differential equations,” in *International Conference on Learning  
531 Representations (ICLR)*, 2021.

532

533 S. Venkatraman, S. Khaitan, R. T. Akella, J. Dolan, J. Schneider, and G. Berseth, “Reasoning with  
534 latent diffusion in offline reinforcement learning,” in *The Twelfth International Conference on  
535 Learning Representations*, 2024.

536

537 A. Ajay, Y. Du, A. Gupta, J. Tenenbaum, T. Jaakkola, and P. Agrawal, “Is conditional genera-  
538 tive modeling all you need for decision making?” in *The Eleventh International Conference on  
539 Learning Representations (ICLR)*, 2023.

540

541 C. Chen, F. Deng, K. Kawaguchi, C. Gulcehre, and S. Ahn, “Simple hierarchical planning with  
542 diffusion,” in *The Twelfth International Conference on Learning Representations (ICLR)*, 2024.

543

544 R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.

545

546 P. Dhariwal and A. Q. Nichol, “Diffusion models beat GANs on image synthesis,” in *Advances in  
547 Neural Information Processing Systems (NeurIPS)*, 2021.

540 J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine, “D4rl: Datasets for deep data-driven rein-  
541forcement learning,” *arXiv preprint arXiv:2004.07219*, 2020.

542

543 G. Williams, P. Drews, B. Goldfain, J. M. Rehg, and E. A. Theodorou, “Aggressive driving with  
544 model predictive path integral control,” in *2016 IEEE International Conference on Robotics and  
545 Automation (ICRA)*. IEEE, 2016.

546

547 L. Chen, K. Lu, A. Rajeswaran, and J. Lee, Pieter Abbeel, “Decision transformer: Reinforcement  
548 learning via sequence modeling,” *Advances in Neural Information Processing Systems (NeurIPS)*,  
549 2021.

550

551 H. Chen, C. Lu, C. Ying, H. Su, and J. Zhu, “Offline reinforcement learning via high-fidelity gener-  
552 ative behavior modeling,” in *The Eleventh International Conference on Learning Representations  
(ICLR)*, 2023.

553

554 W. Li, X. Wang, B. Jin, and H. Zha, “Hierarchical diffusion for offline decision making,” in *Inter-  
555 national Conference on Machine Learning (ICML)*. PMLR, 2023.

556

557 D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” in *International Conference on  
558 Learning Representations (ICLR)*, 2014.

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594 **A PROOF OF PROPOSITION 1**  
 595

596 In this section, we present the proof of Proposition 1.  
 597

598 *Proof.* To incorporate key RL mechanisms into the training of the diffusion model, we explore  
 599 the denoising process and trace the denoised data through the reverse diffusion process. Let  $\hat{\tau}^0$   
 600 represent the denoised output trajectory. It can be gradually denoised using the reverse process,  
 601 following the chain rule:  $\hat{\tau}^0 \sim p_\theta(\tau^K) \prod_{k=1}^K p_\theta(\tau^{k-1} | \tau^k)$ , where the detailed reverse process is  
 602 defined in Eq. (3) and Eq. (4). Starting from an intermediate trajectory  $\tau^k$  at step  $k$ , by combining  
 603 these two equations, the trajectory at the next diffusion step,  $k-1$ , can be directly sampled from the  
 604 distribution:

$$605 \hat{\tau}^{k-1} \sim \mathcal{N} \left( \frac{1}{\sqrt{\alpha_k}} \left( \tau^k - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k}} \epsilon_\theta(\tau^k, k) \right), \sigma_k^2 \mathbf{I} \right). \quad (15)$$

607 By applying the reparameterization trick (Kingma and Welling, 2014), we can derive a closed-form  
 608 solution for the above distribution. Let  $\epsilon_k$  represent the noise introduced in the reverse process  
 609  $p_\theta(\tau^{k-1} | \tau_k)$ , and the denoised trajectory can then be formulated as:

$$610 \hat{\tau}^{k-1} = \frac{1}{\sqrt{\alpha_k}} \left( \tau^k - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k}} \epsilon_\theta(\tau^k, k) \right) + \sigma_k \epsilon_k \\ 611 = \frac{1}{\sqrt{\alpha_k}} \tau^k - \frac{1 - \alpha_k}{\sqrt{(1 - \bar{\alpha}_k) \alpha_k}} \epsilon_\theta(\tau^k, k) + \sigma_k \epsilon_k. \quad (16)$$

615 In the following diffusion step  $k-2$ , the denoised data  $\hat{\tau}^{k-2}$  is sampled from a similar Gaussian dis-  
 616 tribution. By the Central Limit Theorem,  $\hat{\tau}^{k-1}$  serves as an unbiased estimate of  $\tau^{k-1}$ . Therefore,  
 617 the denoised data  $\hat{\tau}^{k-2}$  can be expressed as follows:

$$618 \hat{\tau}^{k-2} \sim \mathcal{N} \left( \frac{1}{\sqrt{\alpha_{k-1}}} \left( \tau^{k-1} - \frac{1 - \alpha_{k-1}}{\sqrt{1 - \bar{\alpha}_{k-1}}} \epsilon_\theta(\tau^{k-1}, k-1) \right), \sigma_{k-1}^2 \mathbf{I} \right) \\ 619 = \frac{1}{\sqrt{\alpha_{k-1}}} \left( \hat{\tau}^{k-1} - \frac{1 - \alpha_{k-1}}{\sqrt{1 - \bar{\alpha}_{k-1}}} \epsilon_\theta(\tau^{k-1}, k-1) \right) + \sigma_{k-1} \epsilon_{k-1} \\ 620 = \frac{1}{\sqrt{\alpha_{k-1}}} \left( \frac{1}{\sqrt{\alpha_k}} \tau^k - \frac{1 - \alpha_k}{\sqrt{(1 - \bar{\alpha}_k) \alpha_k}} \epsilon_\theta(\tau^k, k) - \frac{1 - \alpha_{k-1}}{\sqrt{1 - \bar{\alpha}_{k-1}}} \epsilon_\theta(\tau^{k-1}, k-1) + \sigma_k \epsilon_k \right) \\ 621 + \sigma_{k-1} \epsilon_{k-1} \\ 622 = \frac{1}{\sqrt{\alpha_k \alpha_{k-1}}} \tau^k - \frac{1 - \alpha_k}{\sqrt{(1 - \bar{\alpha}_k) \alpha_k \alpha_{k-1}}} \epsilon_\theta(\tau^k, k) - \frac{1 - \alpha_{k-1}}{\sqrt{(1 - \bar{\alpha}_{k-1}) \alpha_{k-1}}} \epsilon_\theta(\tau^{k-1}, k-1) \\ 623 + \frac{1}{\sqrt{\alpha_{k-1}}} \sigma_k \epsilon_k + \sigma_{k-1} \epsilon_{k-1}. \quad (17)$$

632 The introduced noise  $\epsilon_{k-1}$  in diffusion step  $k-1$  can be combined with the noise  $\epsilon_k$  at diffu-  
 633 sion step  $k$  into a joint noise term,  $\bar{\epsilon}_{k-1}$ , by merging two Gaussian distributions,  $\mathcal{N}(0, \frac{\sigma_k^2}{\alpha_{k-1}} \mathbf{I})$  and  
 634  $\mathcal{N}(0, \sigma_{k-1}^2 \mathbf{I})$ , into  $\mathcal{N}(0, (\frac{\sigma_k^2}{\alpha_{k-1}} + \sigma_{k-1}^2) \mathbf{I})$ . Consequently, we obtain the distribution for the denoised  
 635 data  $\hat{\tau}^{k-2}$  with only directly computable terms, where  
 636

$$637 \hat{\tau}^{k-2} = \frac{1}{\sqrt{\alpha_k \alpha_{k-1}}} \tau^k - \frac{1 - \alpha_k}{\sqrt{(1 - \bar{\alpha}_k) \alpha_k \alpha_{k-1}}} \epsilon_\theta(\tau^k, k) - \frac{1 - \alpha_{k-1}}{\sqrt{(1 - \bar{\alpha}_{k-1}) \alpha_{k-1}}} \epsilon_\theta(\tau^{k-1}, k-1) \\ 638 + \sqrt{\frac{\sigma_k^2}{\alpha_{k-1}} + \sigma_{k-1}^2} \bar{\epsilon}_{k-1} \\ 639 \sim \mathcal{N} \left( \frac{1}{\sqrt{\alpha_k \alpha_{k-1}}} \tau^k - \frac{1 - \alpha_k}{\sqrt{(1 - \bar{\alpha}_k) \alpha_k \alpha_{k-1}}} \epsilon_\theta(\tau^k, k) - \frac{1 - \alpha_{k-1}}{\sqrt{(1 - \bar{\alpha}_{k-1}) \alpha_{k-1}}} \epsilon_\theta(\tau^{k-1}, k-1), \right. \\ 640 \left. \left( \frac{\sigma_k^2}{\alpha_{k-1}} + \sigma_{k-1}^2 \right) \mathbf{I} \right). \quad (18)$$

648 By repeating the denoising process for  $k$  iterations, we can ultimately obtain a closed-form representation of the denoised data  $\hat{\tau}^0$ .  
 649

$$\begin{aligned}
 651 \quad \hat{\tau}^0 &= \frac{1}{\sqrt{\prod_{i=1}^k \alpha_i}} \tau^k - \sum_{i=1}^k \frac{1 - \alpha_i}{\sqrt{(1 - \bar{\alpha}_i) \prod_{j=1}^i \alpha_j}} \epsilon_\theta(\tau^i, i) + \sqrt{\sigma_1^2 + \sum_{i=2}^k \frac{\sigma_i^2}{\prod_{j=1}^{i-1} \alpha_j} \bar{\epsilon}_1} \\
 652 \\
 653 \quad &= \frac{1}{\sqrt{\bar{\alpha}_k}} \tau^k - \sum_{i=1}^k \frac{1 - \alpha_i}{\sqrt{(1 - \bar{\alpha}_i) \bar{\alpha}_i}} \epsilon_\theta(\tau^i, i) + \sqrt{\sigma_1^2 + \sum_{i=2}^k \frac{\sigma_i^2}{\bar{\alpha}_{i-1}} \bar{\epsilon}_1}.
 \end{aligned} \tag{19}$$

654 Using the closed-form representation of the reparameterization trick, the final denoised data  $\hat{\tau}^0$  follows a Gaussian distribution, expressed as  $\hat{\tau}^0 \sim \mathcal{N}(\hat{\mu}_\theta(\tau^k, k), \hat{\sigma}^2 \mathbf{I})$ . The mean  $\hat{\mu}_\theta(\tau^k, k)$  captures  
 655 the denoising trajectory and is formulated as:  
 656

$$\hat{\mu}_\theta(\tau^k, k) = \frac{1}{\sqrt{\bar{\alpha}_k}} \tau^k - \sum_{i=1}^k \frac{1 - \alpha_i}{\sqrt{(1 - \bar{\alpha}_i) \bar{\alpha}_i}} \epsilon_\theta(\tau^i, i). \tag{20}$$

663 Similarly, the covariance  $\hat{\sigma}^2$  accounts for the accumulation of noise over all diffusion steps and is  
 664 written as:  
 665

$$\hat{\sigma}^2 = \sigma_1^2 + \sum_{i=2}^k \frac{\sigma_i^2}{\bar{\alpha}_{i-1}}. \tag{21}$$

□

## 669 B IMPLEMENTATION DETAILS

670 **Reward model and transition model pretraining** Before training the diffusion model, we first  
 671 pretrain both the reward model and the transition model on the concatenated inputs  $(s_t, a_t)$  from  
 672 the same dataset used for diffusion training (e.g., D4RL (Fu et al., 2020)), with the same Gaussian  
 673 normalization. The transition model is implemented as a MLP with two hidden layers of 512 units,  
 674 ReLU activations, and a linear output head predicting the next-state mean  $\mu(s_t, a_t) \in \mathbb{R}^{s_{\text{dim}}}$ , trained  
 675 with mean squared error against  $s_{t+1}$ . The reward model is also an MLP with two hidden layers  
 676 of 256 units, ReLU activations, and a linear output head, trained by regression to the per-timestep  
 677 rewards in the dataset. Both models are pretrained using the Adam optimizer (learning rate  $3 \times 10^{-4}$ ),  
 678 batch size 64, for  $5 \times 10^5$  training steps. After pretraining, the reward and transition models are  
 679 frozen during diffusion model training.  
 680

681 **Diffusion training** We adopt the core diffusion model and reward guidance implementation from  
 682 Diffuser (Janner et al., 2022). Both the diffusion backbone and the reward-guidance network use a  
 683 temporal U-Net trained on length- $T$  trajectories of concatenated  $(s_t, a_t)$ , with hard conditioning on  
 684 the initial observation  $s_0$ . We set the planning horizon to  $T = 32$  for locomotion tasks,  $T = 128$   
 685 for the three U-Maze tasks,  $T = 256$  for the three Medium-Maze tasks, and  $T = 384$  for the  
 686 three Large-Maze tasks. Observations and actions are Gaussian-normalized using statistics from the  
 687 offline dataset.  
 688

689 During training, the diffusion model is optimized with our designed total loss  $L_{\text{total}} = L_{\text{wdiff}} +$   
 690  $\lambda_{\text{tr}} L_{\text{tr}} + \lambda_{\text{rd}} L_{\text{rd}}$ , as defined in Eq. (12). The weights in the reward-aware diffusion loss  $L_{\text{wdiff}}$  are  
 691 clipped by  $r_{\text{max}}$ , which we set to 1 in practice. We use  $\lambda_{\text{tr}} = 0.1$  for the transition-based auxiliary  
 692 modulation loss  $L_{\text{tr}}$  and  $\lambda_{\text{rd}} = 0.05$  for the reward-based auxiliary modulation loss  $L_{\text{rd}}$ . When a  
 693 domain lacks stepwise rewards, the reward bias term is omitted. The diffusion backbone is trained  
 694 with the Adam optimizer (learning rate  $2 \times 10^{-4}$ ), batch size 32, gradient accumulation factor 2,  
 695 and EMA decay 0.995.  
 696

697 **Diffusion sampling** Following Janner et al. (2022), we use  $N = 20$  reverse diffusion steps and  
 698 apply reward-only gradient guidance with scale  $\alpha = 10^{-3}$  at each step, re-imposing conditioning  
 699 after every step. We report the top-scoring trajectories under reward guidance.  
 700

701 **Computational resources** All experiments were conducted on a cluster of 10 nodes, each  
 702 equipped with four Intel Xeon CPUs, 32 GB of RAM, and an NVIDIA GeForce RTX 2080 GPU  
 703 with 11 GB of VRAM.  
 704

702  
 703 Table 4: Quarter-wise transition mismatch  $E_i$  (mean squared state prediction error) for Diffuser and  
 704 DMEMM on D4RL locomotion tasks (Medium-Expert). Results are averaged over 5 runs.

Method	Environment	Quarter 1	Quarter 2	Quarter 3	Quarter 4
Diffuser	HalfCheetah	$29.10 \pm 0.32$	$34.67 \pm 0.40$	$34.96 \pm 0.37$	$38.71 \pm 0.42$
DMEMM	HalfCheetah	<b><math>9.74 \pm 0.05</math></b>	<b><math>9.88 \pm 0.04</math></b>	<b><math>9.91 \pm 0.09</math></b>	<b><math>9.99 \pm 0.06</math></b>
Diffuser	Hopper	$1.01 \pm 0.06$	$0.96 \pm 0.05$	$1.35 \pm 0.08$	$2.96 \pm 0.12$
DMEMM	Hopper	<b><math>0.82 \pm 0.07</math></b>	<b><math>0.75 \pm 0.06</math></b>	<b><math>1.02 \pm 0.09</math></b>	<b><math>1.85 \pm 0.15</math></b>
Diffuser	Walker2D	$1.42 \pm 0.07$	$3.21 \pm 0.06$	$3.52 \pm 0.11$	$4.03 \pm 0.13$
DMEMM	Walker2D	<b><math>1.12 \pm 0.08</math></b>	<b><math>2.19 \pm 0.10</math></b>	<b><math>2.43 \pm 0.11</math></b>	<b><math>2.61 \pm 0.14</math></b>

## 714 C EVALUATION OF TRANSITION MISMATCH

716 To investigate the potential transition mismatch problem in conventional Diffuser models (Janner  
 717 et al., 2022), we conduct an experiment to quantify the discrepancy between trajectories predicted  
 718 by the diffusion model and the actual environment rollouts. This analysis highlights how traditional  
 719 Diffuser suffers from model–environment dynamics gaps, and how our proposed DMEMM effec-  
 720 tively tackles this issue. We compare Diffuser and DMEMM on D4RL locomotion tasks (Fu et al.,  
 721 2020) at the Medium-Expert level.

722 We use the trained diffusion models with reward guidance from both Diffuser and DMEMM for  
 723 evaluation. During the online sampling stage, the planner at each timestep  $t$  generates an imagined  
 724 future sequence  $\{\hat{s}_{t+1}^t, \hat{s}_{t+2}^t, \dots, \hat{s}_{t+H}^t\}$ , where  $H$  is the planning horizon. Each predicted state  $\hat{s}_{t+h}^t$   
 725 is compared to the corresponding ground-truth state  $s_{t+h}$  collected from the environment for  $h \in$   
 726  $\{1, \dots, H\}$ . This measures the discrepancy caused by the mismatch between the diffusion model  
 727 and the true environment dynamics. At each prediction step we compute the L2-norm error  $e_{t,h} =$   
 728  $\|\hat{s}_{t+h}^t - s_{t+h}\|_2^2$ . To reduce computation, we reuse previously generated plans by backtracking from  
 729 stored plans at earlier states rather than recomputing forward rollouts from scratch.

730 To better analyze prediction quality over different time scales, we divide the planning horizon  
 731 into four equal-length quarters and report the average error in each quarter. Early quarters reflect  
 732 short-term prediction accuracy, while later quarters capture long-horizon stability. Formally, let  
 733 the trajectory length be  $T$ , the horizon  $H$ , and let quarter  $i \in \{1, 2, 3, 4\}$  cover prediction steps  
 734  $h \in \left[\frac{(i-1)H}{4} + 1, \frac{iH}{4}\right]$ . The average squared error for  $i$ -th quarter is:

$$737 \quad E_i = \frac{1}{(T-H) \cdot \frac{H}{4}} \sum_{t=0}^{T-H-1} \sum_{h=(i-1)H/4+1}^{iH/4} \|\hat{s}_{t+h}^t - s_{t+h}\|_2^2. \quad (22)$$

741 These quarter-wise errors quantify how transition mismatch accumulates along the trajectory:  
 742 smaller  $E_i$  in later quarters indicates better long-horizon predictive ability.

744 The results are summarized in Table 4. We observe that DMEMM greatly outperforms Diffuser  
 745 in terms of transition mismatch on the HalfCheetah environment. On the other two environments,  
 746 DMEMM still surpasses Diffuser, though with a smaller relative improvement. We hypothesize that  
 747 transition mismatch is more severe in complex dynamical systems such as HalfCheetah, suggesting  
 748 that our method is particularly beneficial for environments with more challenging dynamics.

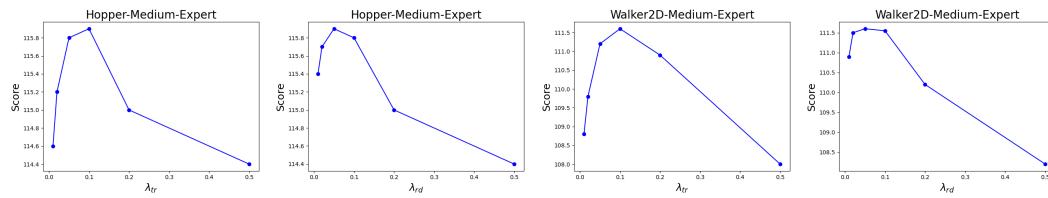
749 Furthermore, as the planning horizon increases, Diffuser’s prediction errors grow substantially, in-  
 750 dicating poor generalization to real online interactions despite strong offline fitting. In contrast,  
 751 DMEMM consistently maintains lower transition mismatch, especially in the later quarters, with-  
 752 out exhibiting the pronounced error escalation seen in Diffuser. Interestingly, although HalfCheetah  
 753 shows the highest absolute prediction errors, neither Diffuser nor DMEMM displays a sharp error  
 754 increase across quarters in this task.

755 Overall, these findings demonstrate the superior long-horizon prediction quality and robustness of  
 the proposed DMEMM method.

756

757 Table 5: This table presents the scores on D4RL locomotion suites for all four ablation variants.  
758 Results are averaged over 5 seeds.

Gym Tasks	DMEMM	DMEMM-w/o-weighting	DMEMM-w/o- $\lambda_{tr}$	DMEMM-w/o- $\lambda_{rd}$	DMEMM-w/o-tr-guide
HalfCheetah (Med-Expert)	<b>94.6±1.2</b>	93.8±0.9	92.2±0.6	92.8±1.2	92.5±1.3
Hopper (Med-Expert)	<b>115.9±1.6</b>	115.2±0.4	114.4±0.8	115.0±0.4	114.8±0.2
Walker2d (Med-Expert)	<b>111.6±1.1</b>	110.4±0.8	108.4±1.2	110.4±0.6	109.9±1.0
HalfCheetah (Medium)	<b>49.2±0.8</b>	48.0±1.1	46.3±0.4	47.1±0.6	46.9±0.9
Hopper (Medium)	<b>101.2±1.4</b>	100.4±1.2	98.6±1.8	100.1±1.1	99.8±1.6
Walker2d (Medium)	<b>86.5±1.5</b>	85.6±1.2	82.8±1.4	84.4±0.9	83.0±1.8
HalfCheetah (Med-Replay)	<b>46.1±1.3</b>	44.7±1.7	42.5±2.9	44.2±1.4	43.6±2.5
Hopper (Med-Replay)	<b>100.6±0.9</b>	98.8±1.2	97.0±0.9	98.2±0.6	96.2±1.2
Walker2d (Med-Replay)	<b>85.8±2.6</b>	84.6±2.2	82.2±1.7	83.7±2.5	82.6±3.2

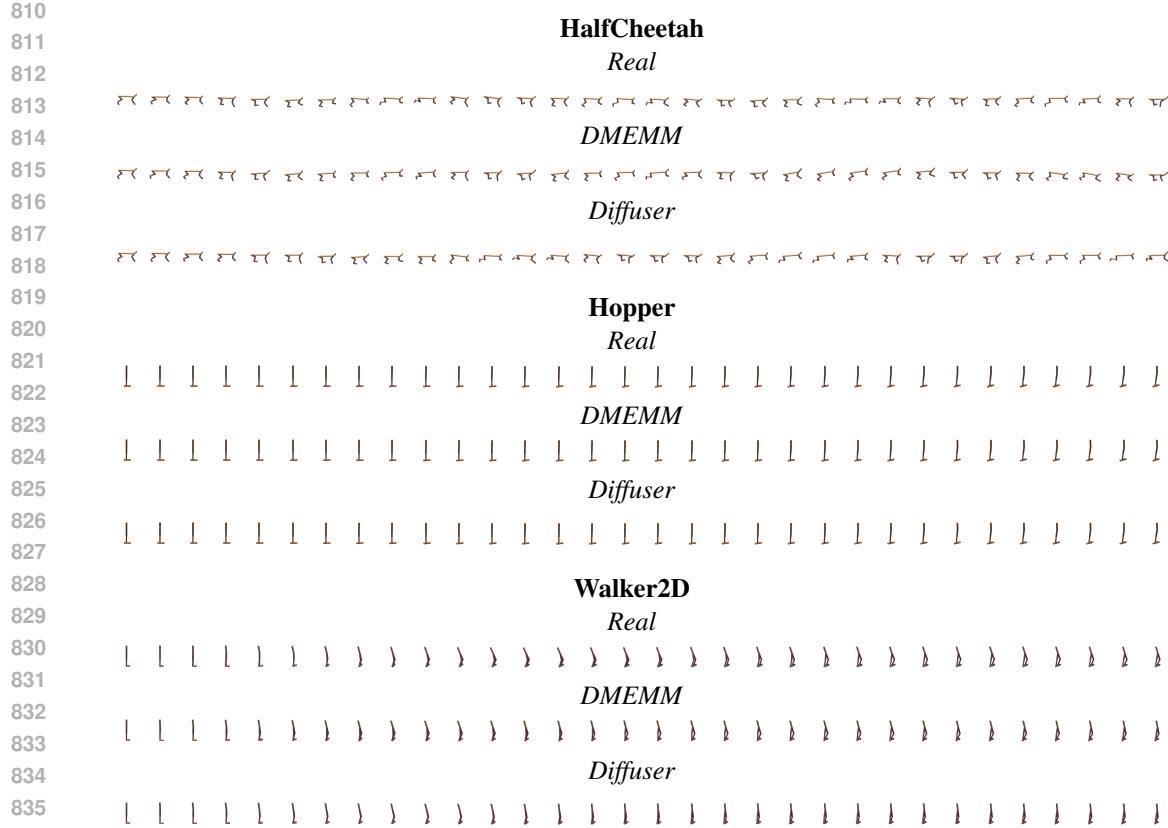
755 Figure 1: Hyperparameter sensitivity analysis of the tradeoff parameters for transition-based diffu-  
756 sion modulation loss ( $\lambda_{tr}$ ) and reward-based diffusion modulation loss ( $\lambda_{rd}$ ) on Hopper-Medium-  
757 Expert and Walker2D-Medium-Expert environments.778 

## D FULL ABLATION RESULTS

781 The complete ablation results of our DMEMM method compared with the four ablated variants (1)  
782 DMEMM-w/o-weighting, which removes the weighting function in the reward-aware diffusion loss;  
783 (2) DMEMM-w/o- $\lambda_{tr}$ , which omits the transition-based diffusion modulation loss; (3) DMEMM-  
784 w/o- $\lambda_{rd}$ , which omits the reward-based diffusion modulation loss; and (4) DMEMM-w/o-tr-guide,  
785 which removes transition guidance in the dual-guided sampling procedure are presented in Table 5.786 The conclusions drawn from the full ablation results are consistent with those reported in the  
787 main paper. Across all three environments and all expert demonstration levels, the performance  
788 of DMEMM is substantially degraded when either the transition-based diffusion modulation loss or  
789 the transition guidance is removed, highlighting the critical role of explicitly modeling transition dy-  
790 namics. Dropping the weighting function or the reward-based diffusion modulation loss also harms  
791 performance, with the reward auxiliary loss  $L_{rd}$  appearing relatively more important between the  
792 two. Overall, removing any single component leads to a noticeable performance drop compared to  
793 the full DMEMM model, demonstrating the effectiveness and necessity of each key component in  
794 our approach.795 

## E HYPERPARAMETER SENSITIVITY ANALYSIS

798 In this section, we analyze the sensitivity of the tradeoff parameters  $\lambda_{tr}$  (transition-based diffusion  
799 modulation loss) and  $\lambda_{rd}$  (reward-based diffusion modulation loss) to understand their impact on  
800 performance in offline RL tasks. The analysis is conducted on two environments: Hopper-Medium-  
801 Expert and Walker2D-Medium-Expert.802 Figures 1 illustrate the performance sensitivity to the tradeoff parameters. For  $\lambda_{tr}$ , the per-  
803 formance peaks at approximately  $\lambda_{tr} = 0.1$  in both the Walker2D-Medium-Expert and Hopper-  
804 Medium-Expert environments. Beyond this optimal point, performance declines notably, regard-  
805 less of whether  $\lambda_{tr}$  is increased or decreased. Similarly, for  $\lambda_{rd}$ , the performance also peaks around  
806  $\lambda_{rd} = 0.05$  in both environments. However, unlike  $\lambda_{tr}$ , performance shows little change when  $\lambda_{rd}$   
807 is adjusted within a small range, indicating that  $\lambda_{rd}$  is less sensitive than  $\lambda_{tr}$ . Overall, the hyper-  
808 parameter sensitivity analysis shows that both  $\lambda_{rd}$  and  $\lambda_{tr}$  have similar effects on performance and  
809 are robust across different tasks. Additionally, it confirms that the selected hyperparameters for our  
810 experiments are optimal.



838 Figure 2: Visualization of trajectories generated by Real environment rollouts, DMEMM, and Diffuser on HalfCheetah, Hopper, and Walker2D (32-long horizon per method). Figures are slightly  
839 overlapped horizontally to highlight differences in posture and stability.  
840

## 842 F TRAJECTORY CONSISTENCY VISUALIZATION

844 We provide a qualitative visualization of the trajectories generated during online planning to demon-  
845 strate that our proposed DMEMM method improves transition consistency. Starting from a fixed  
846 initial state, we sample trajectories using DMEMM and the standard Diffuser, and compare them  
847 against ground-truth trajectories collected from the offline dataset. To ensure a fair comparison, we  
848 visualize 32 consecutive frames for each method. The results are shown in Figure E.

849 From the figure, we observe that DMEMM produces trajectories that closely resemble the real tra-  
850 jectories across all environments, especially on HalfCheetah. For Hopper, the trajectories generated  
851 by all methods are visually similar, which is consistent with the quantitative results in Table 4. On  
852 Walker2D, trajectories generated by DMEMM remain closer to the real trajectories than those pro-  
853 duced by Diffuser, indicating that DMEMM can generate more consistent and stable motions than  
854 the base Diffuser method.

855  
856  
857  
858  
859  
860  
861  
862  
863