

WHAT FLOW-MATCHING BRINGS TO TD-LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent work shows that flow-matching can be effective for value estimation in RL, but it remains unclear why they work well or whether flow-matching Q-functions differ fundamentally from standard critics. We show that their success is not explained by distributional RL: explicitly modeling return distributions often degrades performance. Instead, we argue that flow-matching Q-functions are effective because they couple a learned velocity field with an integration that is used both during training and to read out Q-values at inference. This coupling enables robust value prediction through *test-time recovery* from imperfect intermediate estimates where errors dampen out as more integration steps are performed. This mechanism is absent in monolithic critics. Beyond test-time recovery, training with the integration procedure induces more *plastic* representations, allowing critics to represent non-stationary future TD targets without overwriting previous features. We formalize these effects and validate them empirically, showing that flow-matching critics outperform monolithic critics by over $2\times$ in performance and achieve $5\text{--}10\times$ higher sample efficiency in high-UTD regimes.

1 INTRODUCTION

Recent works have demonstrated that flow-matching networks can be highly effective for value function estimation in off-policy reinforcement learning (RL) (Agrawalla et al., 2025; Espinosa-Dice et al., 2025a; Dong et al., 2025). These flow-matching critics depart from standard “monolithic” architectures, which map state-action pairs to scalar Q-values in a single forward pass, by instead estimating values via the iterative integration of a learned velocity field given a noise input. This approach yields substantial empirical gains, and is more robust than monolithic networks in offline (Levine et al., 2020) and offline-to-online RL settings (Nakamoto et al., 2024). However, while the gap is clear, the mechanism driving it is unexplained: do flow-matching critics succeed because they model a value distribution, or do they offer a different inductive bias that limits pathologies of temporal-difference (TD) learning?

A natural hypothesis is that flow-matching critics succeed because they implicitly model return distributions, similar to distributional RL (Bellemare et al., 2017). Indeed, much recent work applying flow-matching to value learning (Espinosa-Dice et al., 2025a; Dong et al., 2025) explicitly adopts distributional objectives. We test this hypothesis and find that it does not explain the observed gains: explicitly incorporating distributional updates often degrades performance relative to simple “expected-value” backups, which do not learn the return distribution. As such, flow-matching critics trained with standard TD consistently outperform strong distributional baselines. These results indicate that the advantage of flow matching does not come from modeling return distributions or distributional RL.

So why do flow-matching critics work? In this paper, we show that their effectiveness stems from training a velocity network “wrapped” by an integration procedure (Figure 7), which is also used at inference time to predict Q-values. This form of *iterative computation* plays a dual role. First, at inference time, the integration process enables the critic to recover from poor intermediate estimates and refine its prediction, yielding substantially more robust value estimates; we term this phenomenon *test-time recovery*. Second, beyond test-time recovery, training with the integration procedure induces more *plastic* features that allow better fitting subsequent TD targets. In particular, when trained with integration, the velocity network’s representations need not change substantially to absorb successive TD targets, as these shifts can instead be accommodated through integration. This mitigates pathologies associated with loss of plasticity (Lyle et al., 2023) or representational

capacity (Kumar et al., 2022). Put together, flow-matching critics outperform monolithic architectures with identical computation graphs that lack dense velocity supervision.

We support this argument empirically and theoretically. **Theoretically**, we formalize the notions of *test-time recovery* and the preservation of *plasticity* in simple settings, and show that flow-matching induces weight update dynamics that preserve and reweight previously learned features rather than overwriting them under non-stationary TD targets in linear settings. Such dynamics are absent in monolithic critics and their ensembles, and achieving similar behavior instead requires auxiliary objectives or explicit regularization that bias optimization and must be carefully tuned. **Empirically**, we show that flow-matching critics tolerate substantially higher noise, are more robust to interventions such as resetting or freezing network components, learn more isotropic features, and achieve much stronger performance, resulting in a $2\times$ performance gain and a $10\times$ improvement in sample efficiency in high update-to-data online RL with offline data. We further show that these gains arise from fitting velocities rather than directly regressing to TD targets.

2 FLOW-MATCHING ENABLES TEST-TIME RECOVERY

2.1 FLOW-MATCHING IMPLEMENTS AN ITERATIVE PROCESS

We begin by formalizing test-time recovery in Definition H.1. Intuitively, it describes the ability of a flow-matching Q-function to compensate for errors or inconsistencies introduced during inference, using a *single set of trained parameters*, by controlling the number of integration steps for computing Q-values. This mechanism is absent in monolithic critics that only query the network once.

From a control-theoretic perspective, this definition of test-time recovery corresponds to *incremental stability* of the inference-time dynamics induced by the learned velocity field, meaning that perturbations to intermediate states or velocity evaluations are progressively damped along the integration trajectory. We elaborate on this connection formally in Appendix N. Intuitively, TTR means that errors introduced at intermediate integration steps can be corrected by spending more integration steps. When β_K decreases with K (see Theorem N.1), increasing the number of integration steps improves robustness, explaining why flow-matching critics benefit from additional test-time compute.

2.2 ANALYSIS: TEST-TIME RECOVERY (TTR) IN PRACTICE

We now evaluate the ability of flow-matching critics to perform test-time recovery. To do so, we introduce controlled perturbations into the integration process and measure how sensitive the Q-value estimates are to these interventions. For comparison, we also construct analogous perturbations for monolithic critics and evaluate their behavior.

Experiment: Injecting staleness into earlier integration steps of the flow critic. In this experiment, we split the integration procedure into two phases. After training a flow-matching critic for $T = 250,000$ gradient steps, we evaluate a modified procedure in which the first $\kappa\%$ of the integration steps ($\kappa \in \{0, 25, 50, 75, 100\}$) are intentionally performed using a stale snapshot of the velocity field taken from the checkpoint at time T , while the remaining steps are completed using the current network parameters. We compare this procedure to a baseline which runs all integration steps with current parameters ($\kappa = 0$). If flow-matching exhibits test-time recovery, then later integration steps should be able to correct *at least some* errors arising from staleness of the early steps of the integration trajectory, even though the parameters of the velocity field cannot be updated.

Results. Table 1 provides empirical evidence of test-time recovery. On 2 of 3 environments, using stale velocity parameters for the first 25% or 50% of the integration steps in the critic yields policies with higher success rates than using no stale velocities at all, indicating that flow-matching critics can recover from errors introduced early during integration. This behavior is not universal, as for humanoidmaze-medium any amount of staleness substantially degrades performance.

In contrast, analogous interventions applied to monolithic feed-forward critics consistently result in pronounced performance drops, as shown in Figure 3. We observe similar degradation for ResNet-

Table 1: *Effect of injecting staleness into early integration steps* of the flow-matching critic. Entries are shaded by performance degradation relative to the best value within each environment as $\kappa = \{0, 25, 50, 75, 100\}\%$ increases left to right. A flow-matching critic can still succeed when the first 25%–50% of integration uses a stale velocity field (best in each row shown in **bold**).

Env.	Success as κ increases
antsoccer-arena	45 → 50 → 47 → 46 → 35
hmmaze-medium	98 → 43 → 63 → 62 → 39
cube-double	72 → 84 → 82 → 58 → 64

style critics with residual connections (Figure 2), suggesting that a computation graph resembling an integration process is insufficient to induce TTR. Instead, recovery requires training with a flow-matching loss itself. Finally, as we discuss in the next section, flow matching also induces more robust feature representations, explaining why flow-matching critics remain resilient not only to stale velocities but also to stale internal features.

Experiment: Robustness to noisy TD supervision.

In the previous experiment, we perturbed the integration procedure itself. Here, we examine a complementary probe in which we corrupt the *training-time supervision*, i.e., TD targets that supervise the velocity field. Of course, this corruption would degrade the performance of any approach, but if flow-matching does indeed exhibit TTR, it should exhibit a more graceful degradation as the magnitude of corruption increases. We add noise to the velocity network targets as per the procedure in Appendix M.1. This noise directly affects the local supervision signal at each integration step. We also train a monolithic Q-network baseline with target noise and present its results below.

Results. Figure 1 shows that performance degradation is more graceful for the flow-matching critic, which consistently maintains higher performance (and even no degradation) than the monolithic approach as noise magnitude increases. These results suggest that robustness to noisy supervision can arise in flow-matching critics. Overall, the learned integration dynamics with flow-matching remains more stable, allowing later integration steps to partially attenuate the effect of noisy supervision.

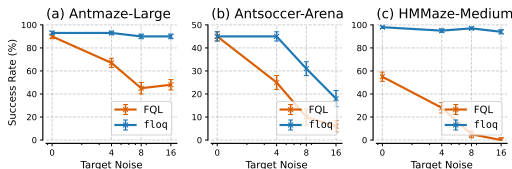


Figure 1: **Performance of fLoq and FQL critics under target noise.** Flow-matching critics are substantially more robust to noise in TD targets, while FQL performance degrades much faster, even when starting from similar initial performance (antmaze/antsoccer).

3 FLOW-MATCHING LEARNS PLASTIC FEATURES

So far we found that flow-matching critics exhibit test-time recovery, and that this behavior arises from supervising velocities along the integration trajectory rather than from the computation graph alone. We now ask whether this same structure also shapes training dynamics. In this section, we show that training a velocity field jointly with the integration procedure leads to more *plastic* representations, which are crucial for mitigating the pathologies induced by non-stationary TD targets.

3.1 EMPIRICAL EVIDENCE OF PLASTIC FEATURE LEARNING

To empirically test this intuition, we examine whether flow-matching critics learn more plastic and qualitatively different feature representations than monolithic networks. In particular, we ask whether the interaction between flow matching and TD-learning leads to differences in learned representations relative to monolithic critics, and whether such differences are specific to TD-learning. We run several experiments that we describe in the paragraphs below and Appendix A.

Experiment: Measuring feature plasticity by freezing features.

To assess whether these representational differences are consequential, we probe feature plasticity by freezing the early layers of the critic at an intermediate point during offline training and continuing TD-learning on the offline dataset. If the learned features are sufficiently expressive to support future TD targets, the impact of freezing should diminish with further training. As shown in Figure 3, freezing features causes monolithic critics to collapse to near-zero performance across almost all environments, with little to no recovery. In contrast, flow-matching critics recover to performance comparable to the unfrozen baseline on all but one environment, where performance remains competitive. This indicates that features learned by flow-matching critics remain useful for representing *future value functions* that will be encountered, even when these features are no longer updated. This supports the notion that flow-matching learns features that help model the value improvement path (Dabney et al., 2020). We further confirm this behavior by applying the same intervention to a monolithic critic with a ResNet-style architecture

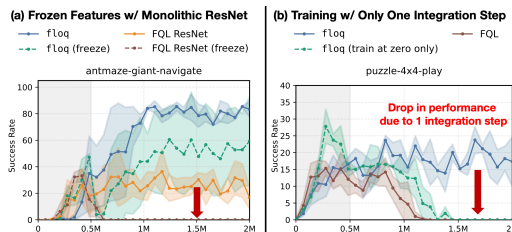


Figure 2: **Frozen features and integration.** (a) A monolithic ResNet critic collapses when features are frozen despite a similar computation graph. (b) A single integration step improves stability over monolithic critics but remains less stable than full flow matching (drop indicated by red arrow), highlighting the role of multi-step integration in preserving plasticity.

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

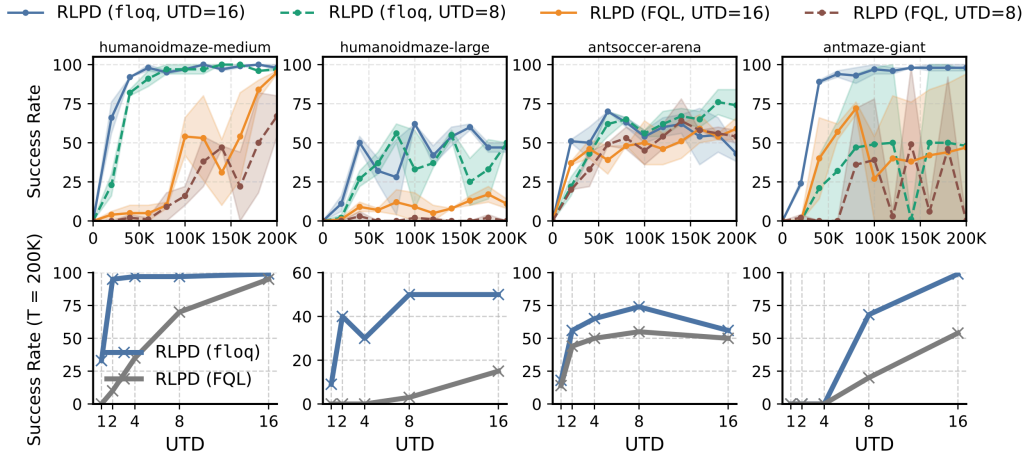


Figure 4: **High-UTD online RL with offline data.** Under RLPD with a 50:50 mix of offline and online data, flow-matching critics (f1oq) consistently outperform monolithic critics (FQL) across all UTD values and substantially improve sample efficiency, achieving gains of up to 5–10× in some environments.

(Figure 2), which also collapses once intermediate layers are frozen. To isolate the role of integration, we additionally evaluate a flow-matching critic trained with only a single integration step (“train at zero only”), where no integration is performed at inference time. While this variant outperforms a monolithic critic (FQL) after the point of intervention (Figure 2; Figure 10 for more tasks), it is substantially less stable than full flow matching, highlighting the crucial role of integration in enabling plastic representations. Finally, we note that unlike in online RL, where learning can recover from frozen/reset features by collecting new data online (Nikishin et al., 2022), such mechanisms are unavailable in offline RL.

4 HIGH-UTD ONLINE RL WITH PRIOR DATA

Our analysis shows that the gains of flow matching arise from training the velocity field jointly with the integration procedure. Given non-stationary TD targets, flow-matching enables test-time recovery and learns more plastic features. Here, we examine whether these properties improve performance in high update-to-data (UTD) regimes. While aggressive data reuse is known to improve efficiency, it often degrades feature plasticity and destabilizes learning (Chen et al., 2021; D’Oro et al., 2022). If a flow-matching critic indeed maintains adaptability under non-stationary TD targets, its performance should scale faster when increasing UTD. To test this, we incorporate flow-matching critics into the RLPD framework (Ball et al., 2023), enabling more aggressive data reuse. As shown in Figure 4, high-UTD training with f1oq, a flow-matching critic, achieves substantially stronger performance at large UTD ratios, with roughly 2× higher final return and a 5-10× improvement in efficiency compared to an RLPD baseline on top of FQL, using monolithic critics. Flow-matching critics are more stable and do not destabilize, even at the highest UTD value we tested. These results show that the plastic features learned via flow-matching enable efficient learning in high UTD settings.

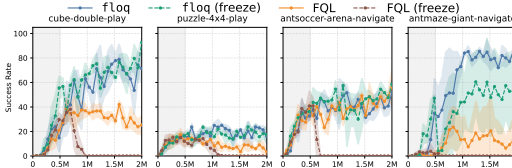


Figure 3: **Measuring feature plasticity.** Feature plasticity on four OG-bench tasks, measured by freezing all but the final two layers at $T=0.5M$ steps (gray region). Solid curves denote fully trained runs; dashed curves show performance after freezing penultimate features. Monolithic FQL exhibits a sharp collapse after freezing, whereas flow-matching critics remain stable and continue improving, indicating greater plasticity.

216 REFERENCES

- 217 Bhavya Agrawalla, Michal Nauman, Khush Agrawal, and Aviral Kumar. floq: Training critics via
218 flow-matching for scaling compute in value-based rl, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2509.06863)
219 [2509.06863](https://arxiv.org/abs/2509.06863).
- 220
- 221 Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants.
222 In *International Conference on Learning Representations (ICLR)*, 2023.
- 223
- 224 Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit
225 acceleration by overparameterization. *arXiv preprint arXiv:1802.06509*, 2018.
- 226
- 227 Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *ArXiv*, abs/1607.06450,
228 2016.
- 229 Philip J Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient online reinforcement learning
230 with offline data. In *International Conference on Machine Learning*, pp. 1577–1594. PMLR, 2023.
- 231
- 232 Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement
233 learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*,
234 pp. 449–458. JMLR. org, 2017.
- 235
- 236 Onur Celik, Zechu Li, Denis Blessing, Ge Li, Daniel Palenicek, Jan Peters, Georgia Chalvatzaki,
237 and Gerhard Neumann. Dime: Diffusion-based maximum entropy reinforcement learning. *arXiv*
preprint arXiv:2502.02316, 2025.
- 238
- 239 Xinyue Chen, Che Wang, Zijian Zhou, and Keith W. Ross. Randomized ensembled double q-learning:
240 Learning fast without a model. In *International Conference on Learning Representations*, 2021.
241 URL <https://openreview.net/forum?id=AY8zfZm0tDd>.
- 242
- 243 Will Dabney, Georg Ostrovski, David Silver, and Rémi Munos. Implicit quantile networks for
244 distributional reinforcement learning. *arXiv preprint arXiv:1806.06923*, 2018.
- 245
- 246 Will Dabney, André Barreto, Mark Rowland, Robert Dadashi, John Quan, Marc G Bellemare, and
247 David Silver. The value-improvement path: Towards better representations for reinforcement
248 learning. *arXiv preprint arXiv:2006.02243*, 2020.
- 249
- 250 Perry Dong, Chongyi Zheng, Chelsea Finn, Dorsa Sadigh, and Benjamin Eysenbach. Value flows,
251 2025. URL <https://arxiv.org/abs/2510.07650>.
- 252
- 253 Pierluca D’Oro, Max Schwarzer, Evgenii Nikishin, Pierre-Luc Bacon, Marc G Bellemare, and Aaron
254 Courville. Sample-efficient reinforcement learning by breaking the replay ratio barrier. In *The*
255 *Eleventh International Conference on Learning Representations*, 2022.
- 256
- 257 Nicolas Espinosa-Dice, Kianté Brantley, and Wen Sun. Expressive value learning for scalable offline
258 reinforcement learning. *arXiv Preprint*, 2025a.
- 259
- 260 Nicolas Espinosa-Dice, Yiyi Zhang, Yiding Chen, Bradley Guo, Owen Oertel, Gokul Swamy, Kianté
261 Brantley, and Wen Sun. Scaling offline rl via efficient and expressive shortcut models. *arXiv*
262 *preprint arXiv:2505.22866*, 2025b.
- 263
- 264 Jesse Farebrother, Jordi Orbay, Quan Vuong, Adrien Ali Taïga, Yevgen Chebotar, Ted Xiao, Alex
265 Irpan, Sergey Levine, Pablo Samuel Castro, Aleksandra Faust, et al. Stop regressing: Training
266 value functions via classification for scalable deep rl. *arXiv preprint arXiv:2403.03950*, 2024.
- 267
- 268 Justin Fu, Aviral Kumar, Ofir Nachum, G. Tucker, and Sergey Levine. D4rl: Datasets for deep
269 data-driven reinforcement learning. *ArXiv*, abs/2004.07219, 2020.
- 265
- 266 Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in
267 actor-critic methods. In *International Conference on Machine Learning (ICML)*, pp. 1587–1596,
268 2018.
- 269
- Hado van Hasselt. Double q-learning. In *Proceedings of the 23rd International Conference on Neural*
Information Processing Systems - Volume 2, 2010.

- 270 Marcel Hussing, Claas Voelcker, Igor Gilitschenski, Amir-massoud Farahmand, and Eric Eaton.
271 Dissecting deep rl with high update ratios: Combatting value divergence. *arXiv preprint*
272 *arXiv:2403.05996*, 2024.
- 273 Aviral Kumar, Rishabh Agarwal, Dibya Ghosh, and Sergey Levine. Implicit under-parameterization
274 inhibits data-efficient deep reinforcement learning. In *International Conference on Learning*
275 *Representations*, 2021.
- 276 Aviral Kumar, Rishabh Agarwal, Tengyu Ma, Aaron Courville, George Tucker, and Sergey Levine.
277 DR3: Value-based deep reinforcement learning requires explicit regularization. *International*
278 *Conference on Learning Representations*, 2022.
- 279 Aviral Kumar, Rishabh Agarwal, Xinyang Geng, George Tucker, and Sergey Levine. Offline Q-
280 learning on diverse multi-task data both scales and generalizes. In *International Conference on*
281 *Learning Representations*, 2023.
- 282 Hojoon Lee, Youngdo Lee, Takuma Seno, Donghu Kim, Peter Stone, and Jaegul Choo. Hyperspherical
283 normalization for scalable deep reinforcement learning. *arXiv preprint arXiv:2502.15280*, 2025.
- 284 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,
285 review, and perspectives on open problems. *arXiv preprint*, 2020.
- 286 Tianhong Li and Kaiming He. Back to basics: Let denoising generative models denoise. *arXiv*
287 *preprint arXiv:2511.13720*, 2025.
- 288 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching
289 for generative modeling. In *International Conference on Learning Representations (ICLR)*, 2023.
- 290 Lei Lv, Yunfei Li, Yu Luo, Fuchun Sun, Tao Kong, Jiafeng Xu, and Xiao Ma. Flow-based policy for
291 online reinforcement learning. *arXiv preprint arXiv:2506.12811*, 2025.
- 292 Clare Lyle, Mark Rowland, Will Dabney, Marta Kwiatkowska, and Yarin Gal. Learning dynamics and
293 generalization in deep reinforcement learning. In *International Conference on Machine Learning*,
294 pp. 14560–14581. PMLR, 2022.
- 295 Clare Lyle, Zeyu Zheng, Evgenii Nikishin, Bernardo Avila Pires, Razvan Pascanu, and Will Dabney.
296 Understanding plasticity in neural networks. In *International Conference on Machine Learning*,
297 2023.
- 298 Clare Lyle, Zeyu Zheng, Khimya Khetarpal, James Martens, Hado P van Hasselt, Razvan Pascanu,
299 and Will Dabney. Normalization and effective learning rates in reinforcement learning. *Advances*
300 *in Neural Information Processing Systems*, 37:106440–106473, 2024.
- 301 Haitong Ma, Tianyi Chen, Kai Wang, Na Li, and Bo Dai. Efficient online reinforcement learning for
302 diffusion policy. *arXiv preprint arXiv:2502.00361*, 2025.
- 303 Mitsuhiro Nakamoto, Simon Zhai, Anikait Singh, Max Sobol Mark, Yi Ma, Chelsea Finn, Aviral
304 Kumar, and Sergey Levine. Cal-ql: Calibrated offline rl pre-training for efficient online fine-tuning.
305 *Advances in Neural Information Processing Systems*, 36, 2024.
- 306 Michal Nauman, Michał Borkiewicz, Piotr Miłoś, Tomasz Trzcinski, Mateusz Ostaszewski, and
307 Marek Cygan. Overestimation, overfitting, and plasticity in actor-critic: The bitter lesson of
308 reinforcement learning. In *International Conference on Machine Learning*, 2024a.
- 309 Michal Nauman, Mateusz Ostaszewski, Krzysztof Jankowski, Piotr Miłoś, and Marek Cygan. Bigger,
310 regularized, optimistic: Scaling for compute and sample-efficient continuous control. *Advances in*
311 *Neural Information Processing Systems*, 2024b.
- 312 Michal Nauman, Marek Cygan, Carmelo Sferrazza, Aviral Kumar, and Pieter Abbeel. Bigger,
313 regularized, categorical: High-capacity value functions are efficient multi-task learners. *arXiv*
314 *preprint arXiv:2505.23150*, 2025.
- 315 Evgenii Nikishin, Max Schwarzer, Pierluca D’Oro, Pierre-Luc Bacon, and Aaron Courville. The
316 primacy bias in deep reinforcement learning. In *International conference on machine learning*, pp.
317 16828–16847. PMLR, 2022.

324 Evgenii Nikishin, Junhyuk Oh, Georg Ostrovski, Clare Lyle, Razvan Pascanu, Will Dabney, and André
325 Barreto. Deep reinforcement learning with plasticity injection. *Advances in Neural Information*
326 *Processing Systems*, 36:37142–37159, 2023.

327 Daniel Palenicek, Florian Vogt, Joe Watson, Ingmar Posner, and Jan Peters. Xqc: Well-conditioned
328 optimization accelerates deep reinforcement learning. *arXiv preprint arXiv:2509.25174*, 2025.

329
330 Seohong Park, Kevin Frans, Sergey Levine, and Aviral Kumar. Is value learning really the main
331 bottleneck in offline rl? *arXiv preprint arXiv:2406.09329*, 2024.

332
333 Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. Ogbench: Benchmarking
334 offline goal-conditioned rl. In *International Conference on Learning Representations (ICLR)*,
335 2025a.

336
337 Seohong Park, Qiyang Li, and Sergey Levine. Flow q-learning. *arXiv preprint arXiv:2502.02538*,
338 2025b.

339
340 Rafael Rafailov, Kyle Beltran Hatch, Anikait Singh, Aviral Kumar, Laura Smith, Ilya Kostrikov,
341 Philippe Hansen-Estruch, Victor Kolev, Philip J Ball, Jiajun Wu, et al. D5rl: Diverse datasets for
data-driven deep reinforcement learning. In *Reinforcement Learning Conference (RLC)*, 2024.

342
343 Allen Z Ren, Justin Lidard, Lars L Ankile, Anthony Simeonov, Pulkit Agrawal, Anirudha Majumdar,
344 Benjamin Burchfiel, Hongkai Dai, and Max Simchowitz. Diffusion policy policy optimization.
arXiv preprint arXiv:2409.00588, 2024.

345
346 Denis Tarasov, Vladislav Kurenkov, Alexander Nikulin, and Sergey Kolesnikov. Revisiting the
347 minimalist approach to offline reinforcement learning. In *Neural Information Processing Systems*
348 *(NeurIPS)*, 2023.

349
350 Andrew Wagenmaker, Mitsuhiko Nakamoto, Yunchu Zhang, Seohong Park, Waleed Yagoub, Anusha
351 Nagabandi, Abhishek Gupta, and Sergey Levine. Steering your diffusion policy with latent space
reinforcement learning. *arXiv preprint arXiv:2506.15799*, 2025.

352
353 Zhendong Wang, Jonathan J Hunt, and Mingyuan Zhou. Diffusion policies as an expressive policy
354 class for offline reinforcement learning.

355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

Appendices

A ADDITIONAL EMPIRICAL EVIDENCE FOR PLASTIC FEATURE LEARNING

Experiment: Measuring properties of learned features.

We measure the ℓ_2 -norm of *post-layernorm* features learned by the velocity network (for flow critics) and the critic network (for monolithic critics) for three learning algorithms: (a) TD-learning, (b) SARSA (using the dataset action for the TD backup), and (c) regression to pre-computed Monte Carlo (MC) returns. As shown in Figure 5 (ref. Fig 8 for more tasks), across several tasks, flow-matching critics trained with TD-learning learn much lower-norm penultimate hidden-layer features than monolithic critic networks, despite the absence of any explicit regularization, which has in fact appeared as a desirable property (Kumar et al., 2023; Hussing et al., 2024). This also indicates that a significant burden of modeling the Q-value scale is deferred to the final layer and the integration, rather than encoded through the network. As such, the model is less likely to learn spurious features to explain the changes in target magnitudes.

Notably, this effect is absent when training with SARSA or Monte-Carlo (MC) regression (Figure 5) (ref. Fig 9 for more tasks), where both flow-matching and monolithic critics exhibit similar feature statistics. This suggests that the representational differences arise specifically from the interaction between flow matching and TD-learning, and not from architectural differences alone. In particular, TD-learning relies on bootstrapping from policy actions that may be out-of-distribution for the offline dataset, introducing a higher degree of non-stationarity in the targets than SARSA.

Experiment: Intermediate velocity supervision is crucial.

Finally, we test whether the specific choice of supervising the *velocity field*, rather than absolute TD targets, is crucial for obtaining benefits of flow matching discussed so far. Motivated by recent work (Li & He, 2025), we implement a variant of flow matching in which each integration step is supervised to directly predict the TD target value itself, rather than the velocity. As shown in Figure 6, although this variant still uses integration and receives dense supervision at every integration step, it fails to retain the benefits of flow matching degrading to a performance comparable to that of a monolithic critic. Empirically, we find that the network learns to ignore the interpolant and instead fit the target values independently at each step, effectively collapsing to an ensemble of monolithic critics. As a result, performance degrades and the gains from iterative computation disappear. This behavior highlights that training a velocity field together with the integration “wrapper” is critical for attaining better performance.

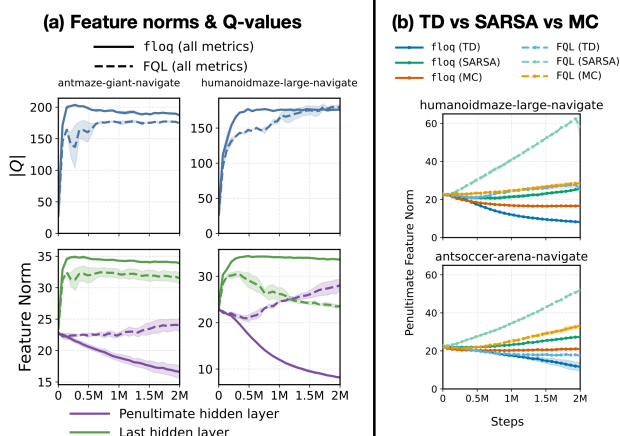


Figure 5: **Feature norms.** (a) Learned feature norms and average Q-values for monolithic critics (FQL) and flow-matching critics (fLoq) in the penultimate and last hidden layers. While the last hidden layer adapts to the scale of Q-values for both methods, the penultimate layer in fLoq exhibits a much more rapid decrease in feature norms, indicating more flexible features largely decoupled from Q-value magnitude. (b) Penultimate-layer feature norms for fLoq trained with TD, SARSA, and MC targets. TD yields the fastest decrease, whereas SARSA and MC resemble monolithic FQL, suggesting that flow-matching critics, especially under TD, learn more robust representations under non-stationary targets.

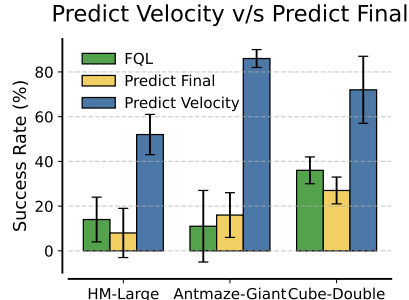


Figure 6: **Impact of predicting velocities versus final TD-targets** at each integration step. Predicting the final TD-target substantially degrades the performance of a flow-matching critic, yielding performance comparable to monolithic critics. This underscores the importance of dense velocity supervision in flow-matching critics.

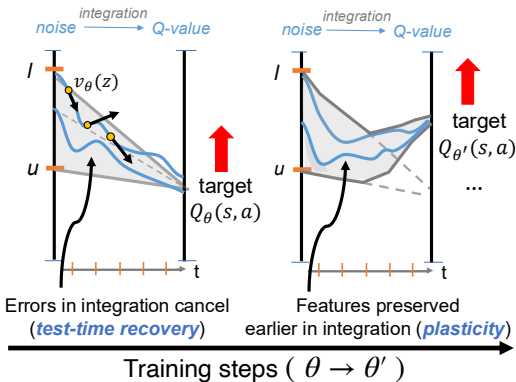
432 B INTUITION FOR PLASTIC FEATURE
433 LEARNING
434

435
436 **Intuition.** Non-stationary TD targets computed on unseen actions can lead to pathologies such as
437 rank collapse, exploding feature norms and dead neurons, as critics must repeatedly overwrite their
438 features to track moving targets, eventually exhausting representational capacity (Kumar et al., 2021;
439 Lyle et al., 2022; Nikishin et al., 2022). We hypothesize that flow-matching critics alleviate these
440 issues due to their structure. Rather than requiring the velocity network to directly update its features
441 to chase each new TD target, flow matching allows the integration process to absorb large changes
442 in value estimates, reducing the need for major changes to the network’s internal representations.
443 Thus, the learned features retain greater plasticity over the course of training. *Viewed through this*
444 *lens, integration plays a dual role:* at test time, it enables recovery and robust Q-value estimates; at
445 training time, it acts as a buffer between changing TD targets and learned representations, promoting
446 more stable feature learning.

447
448
449 C WHY DOES TTR OCCUR WITH FLOW-MATCHING CRITICS BUT IS ABSENT IN
450 MONOLITHIC NETWORKS?
451

452
453 **Why does TTR occur with flow-matching critics but is absent in monolithic networks?**

454 Although Definition H.1 characterizes an inference-time condition, flow-matching critics
455 train a velocity network using TD-style supervision applied densely along the integration trajectory,
456 across many interpolant inputs z and state-action pairs. As a result, the learned dynamics
457 are explicitly shaped to correct local deviations at each integration step. Consequently, errors
458 incurred early in the integration can be progressively attenuated, leading to more accurate Q-
459 value estimates. This intuition is supported by Definition H.1, which shows that even when the
460 per-step TD errors of a flow-matching critic are comparable to those of a monolithic critic, the
461 iterative integration mechanism can yield a more robust final value estimate. In contrast, mono-
462 lithic critics, including ResNet-style architectures with similar inference-time computation graphs,
463 are have only terminal supervision and lack any mechanism to control intermediate representations or
464 suppress error propagation. As a result, errors introduced at intermediate layers tend to accumulate.
465 Finally, we show that this same dense supervision also improves the learned features of the velocity
466 network when chasing non-stationary TD targets, allowing the model to adapt by preserving previous
467 features, which strengthens the TTR phenomenon.



468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
Figure 7: *Flow-matching critics under non-stationary TD-targets.* Over the course of training ($\theta \rightarrow \theta'$), flow-matching critics adapt through a combination of integration (t) and updating the weights of the velocity network v_θ , resulting in **plastic** features and **test-time recovery** from imperfect value estimates.

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
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

Flow-matching critics learn plastic features.

- Flow-matching critics learn more plastic representations than monolithic critics: freezing features severely degrades monolithic critics but has little effect on flow-matching critics.
- Supervising *velocities* is essential for these benefits; directly supervising TD targets collapses flow matching to monolithic behavior, eliminating test-time recovery and plasticity.

Test-time recovery in flow-matching critics.

- Flow-matching critics exhibit *test-time recovery*: perturbations to integration or TD targets dampen as more integration steps are performed.
- Supervising velocities is crucial for effective TTR; monolithic architectures with similar computation graphs (e.g., ResNets) do not exhibit TTR.

E RELATED WORK

Flow-matching (Albergo & Vanden-Eijnden, 2023; Lipman et al., 2023) and diffusion are used to represent policies in RL (Ren et al., 2024; Celik et al., 2025; Ma et al., 2025; Lv et al., 2025; Wang et al.; Park et al., 2025b). More recently, flow-matching has been applied to value learning by parameterizing critics as velocity fields and training them via TD learning using both standard and distributional objectives (Agrawalla et al., 2025; Espinosa-Dice et al., 2025a; Dong et al., 2025). While these works demonstrate strong empirical results, it remains unclear whether these gains stem from capacity, distributional modeling, or the learning dynamics induced by flow-matching. We provide evidence for the latter, showing that flow-matching yields more plastic features and enables robust value estimation via test-time recovery.

Prior work has identified several pathologies in TD learning, including value overestimation (Hasselt, 2010; Fujimoto et al., 2018), growth in parameter norms (Nikishin et al., 2022; Nauman et al., 2024a), and rapid loss of plasticity (Nikishin et al., 2023; Lyle et al., 2023) as training progresses. These issues are commonly attributed to bootstrapping and target non-stationarity, motivating a range of stabilization techniques, including architectural modifications such as layer normalization (Ba et al., 2016; Ball et al., 2023; Nauman et al., 2024b), weight or feature normalization (Kumar et al., 2022; Hussing et al., 2024; Lee et al., 2025), and alternative objectives such as cross-entropy style losses (Bellemare et al., 2017; Farebrother et al., 2024). Collectively, these approaches improve representational capacity and help maintain plasticity (Kumar et al., 2021; 2022; Lyle et al., 2024), and are widely used in modern RL algorithms (Kumar et al., 2023; Lee et al., 2025; Nauman et al., 2025; Palenicek et al., 2025). We show that flow-matching critics stabilize TD learning by addressing plasticity and robustness through the implicit bias of flow matching, without other regularization.

F PRELIMINARIES, NOTATION, AND SETUP

Under usual notation of states and actions, we train a policy $\pi(a|s)$ that induces a distribution over the return random variable $Z^\pi(s, a) \triangleq \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$. The expectation of $Z^\pi(s, a)$ is the Q-function of the policy: $Q^\pi(s, a) = \mathbb{E}[Z^\pi(s, a)]$. We focus on RL training with an offline dataset $\mathcal{D} = \{(s, a, r, s')\}$, in the fully offline RL (Levine et al., 2020) setting for most analysis. Value-based methods learn a network $Q_\theta(\mathbf{s}, \mathbf{a})$ by minimizing TD error.

Flow-matching critics. Flow-matching value functions depart from monolithic models that directly map (\mathbf{s}, \mathbf{a}) to a scalar. Instead, they represent values via a learned transformation of a random noise $\mathbf{z} \in \mathbb{R}$. Concretely, these methods parameterize a time-dependent velocity field $v_\theta(\mathbf{z}, t | \mathbf{s}, \mathbf{a})$ that defines an ODE over \mathbf{z} . Starting from an initial noise sample $\mathbf{z}_0 \sim p_0(\mathbf{z})$ at $t = 0$, numerical integration of this ODE ($\psi(t, \mathbf{z} | \mathbf{s}, \mathbf{a})$) produces a value sample at $t = 1$. Several recent works use this parameterization to learn *distributional* value functions. In particular, Espinosa-Dice et al. (2025a); Dong et al. (2025) train the velocity field so that integrating over the initial noise recovers the full return distribution $Z^\pi(\mathbf{s}, \mathbf{a})$. Given a transition $(\mathbf{s}, \mathbf{a}, r, s') \sim \mathcal{D}$, a distributional TD target is constructed by pushing forward noise through the target flow at the next state: $\mathbf{z}' \sim p_0(\mathbf{z})$, $\tilde{Z}(\mathbf{s}, \mathbf{a}; \mathbf{z}') = r(\mathbf{s}, \mathbf{a}) + \gamma \psi_{\tilde{\theta}}(1, \mathbf{z}' | s', \mathbf{a}')$, where $\psi_{\tilde{\theta}}$ denotes the integrated target flow, $\mathbf{a}' \sim \pi(\cdot | s')$. Flow matching is then applied to align the velocity field with the transport from \mathbf{z} (distinct from \mathbf{z}') to samples from $\tilde{Z}(\mathbf{s}, \mathbf{a}; \mathbf{z}')$. A typical objective is given by

$$\mathcal{L}_{\text{dist}}(\theta) := \mathbb{E}_{(\mathbf{s}, \mathbf{a}, r, s') \sim \mathcal{D}, \mathbf{z}, \mathbf{z}', t \sim \text{Unif}(0, 1)} \left[\left\| v_\theta(\mathbf{z}(t), t | \mathbf{s}, \mathbf{a}) - \bar{s}_{\mathbf{z}, \mathbf{z}'}(\mathbf{s}, \mathbf{a}) \right\|_2^2 \right], \quad (1)$$

where the target velocity is $\bar{s}_{\mathbf{z}, \mathbf{z}'}(\mathbf{s}, \mathbf{a}) := \tilde{Z}(\mathbf{s}, \mathbf{a}; \mathbf{z}') - \mathbf{z}$, and the interpolant is $\mathbf{z}(t) = t \cdot \mathbf{z} + (1 - t) \cdot \tilde{Z}(\mathbf{s}, \mathbf{a}; \mathbf{z}')$. This objective explicitly matches the *entire return distribution*.

f1oq. In contrast, Agrawalla et al. (2025) use flow matching to represent a Q-function while targeting only the *expected* return. Although f1oq integrates noise into a Q-value “sample” and therefore produces stochastic outputs during inference, its TD target collapses the next-state flow to a scalar

expectation. Consequently, this approach does not learn or enforce a return distribution. Specifically, given m i.i.d. noise samples $\{\mathbf{z}'_j\}_{j=1}^m$, define $y(\mathbf{s}, \mathbf{a}) := r(\mathbf{s}, \mathbf{a}) + \gamma \frac{1}{m} \sum_{j=1}^m \psi_{\bar{\theta}}(1, \mathbf{z}'_j | \mathbf{s}', \mathbf{a}')$, which estimates the expected-value TD target $r(\mathbf{s}, \mathbf{a}) + \gamma Q_{\bar{\theta}}(\mathbf{s}', \mathbf{a}')$. Flow-matching is then used to regress from initial noise $\mathbf{z} \sim \text{Unif}[l, u]$ to target $y(\mathbf{s}, \mathbf{a})$ via:

$$\mathcal{L}_{\text{floq}}(\theta) := \mathbb{E}_{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}') \sim \mathcal{D}, \mathbf{z} \sim \text{Unif}[l, u], t \sim \text{Unif}(0, 1)} \left[\left\| v_{\theta}(\mathbf{z}(t), t | \mathbf{s}, \mathbf{a}) - (y(\mathbf{s}, \mathbf{a}) - \mathbf{z}) \right\|_2^2 \right]. \quad (2)$$

where the interpolant is $\mathbf{z}(t) = t \cdot \mathbf{z} + (1 - t) \cdot y(\mathbf{s}, \mathbf{a})$. Thus, unlike the distributional objective above, `floq` does not learn or enforce a distributional Bellman equation; it uses flow matching purely as a *parameterization* of the expected Q-function. Equivalently, the target used to train the velocity field in Equation 2 can be viewed as the expectation of $s_{\mathbf{z}, \mathbf{z}'}(\mathbf{s}, \mathbf{a})$ over the random variable \mathbf{z}' alone.

Agrawalla et al. (2025) argue that `floq` is effective due to *iterative computation*, rather than distributional RL, a perspective also reflected in their training objective. Rather than fitting a Q-function in a single pass, they suggest that integration steps enable a gradual refinement of value estimates. While this explanation is appealing at a high level, it leaves open what iterative computation provides *formally*. In particular, if the benefits of flow matching arise solely from iterative computation at inference time, it is unclear why monolithic architectures fail to exhibit similar behavior.

Our goal. Given these different design choices underlying flow-matching critics, our goal is to identify the mechanisms by which these critics improve TD-learning. How does iterative integration interact with TD bootstrapping? Why does `floq` improve performance despite targeting only expected values? We argue that the answer lies not in distributional modeling, but in *test-time recovery* and improved *representational plasticity* enabled by flow-based critics.

Table 2: *Expected-value (E) vs. distributional (D)* `floq` on representative OGBench tasks. Entries report **E / D**. While both variants learn similar expected Q-values, **D** produces higher-variance estimates without improving performance.

Env.	Succ. (%)	Q_{θ}	$\text{Var}_{\mathbf{z}}(Q)$
hmmaze-large	52 / 30	-180 / -170	0.2/4.5
antmaze-giant	86 / 74	-190 / -200	0.1/0.7
cube-double	72/72	-130 / -130	1.1/6.3
hmmaze-medium	94/94	-170 / -170	0.3/2.3

G DO FLOW-MATCHING CRITICS WORK BECAUSE OF DISTRIBUTIONAL RL?

We now explicitly test whether distributional RL is necessary for strong performance of flow-matching critics. To do so, we use the TD-update from `floq` (Agrawalla et al., 2025) and modify it to use a distributional backup (Equation 1), also following Espinosa-Dice et al. (2025a). Importantly, we keep the velocity field architecture and hyperparameters identical across the two variants. We compare expected-value `floq` and its distributional counterpart on four representative OG-Bench tasks (Park et al., 2025a), along different axes: **(a)** the expected Q-value recovered on the offline dataset; **(b)** variance of the learned Q-value distribution; and **(c)** the performance of the learned policy.

Results. Observe in Table 2, that both expected and distributional variants recover nearly expected Q-values. However, the statistics of the learned Q-value distributions differ substantially. For instance, the standard deviation of the expected variant is significantly *lower* than that of the distributional variant. This behavior is consistent with the training objective. Since `floq` does not attempt to match the return distribution, and OG-Bench tasks are expected to exhibit high-variance, multi-modal returns (Espinosa-Dice et al., 2025a; Dong et al., 2025), regressing to the expected-value target yields substantially lower-variance estimates. These results confirm that while both algorithms learn stochastic Q-function predictions, `floq` does not model the return distribution fully.

Despite this, in Table 2 we see that the distributional variant offers no benefits in performance and is often worse than expected-value `floq`. In addition, Agrawalla et al. (2025) also shows that `floq` outperforms strong distributional RL baselines such as C51 (Bellemare et al., 2017) and IQN (Dabney et al., 2018). Taken together, these findings demonstrate that flow-matching critics can perform extremely well even in the absence of distributional RL training, ruling it out as an explanation. These results also motivate our use of `floq` for the rest of the analysis.

Distributional RL is an insufficient explanation.

- Standard `floq` outperforms its distributional variant although it does not fit the return distribution, as it learns lower variance Q-value distributions.

Definition H.1 (Test-Time Recovery). Fix a state-action pair (\mathbf{s}, \mathbf{a}) and let $\psi_\theta(t, \mathbf{z} | \mathbf{s}, \mathbf{a})$ denote the flow interpolant obtained by integrating the velocity field v_θ . Consider a K -step numerical integrator with step size $\eta = 1/K$ and discrete times $t_k = k/K$. Let $\{\psi^k\}_{k=0}^K$ be the unperturbed trajectory defined by

$$\psi^{k+1} = \psi^k + \eta v_\theta(\psi^k, t_k | \mathbf{s}, \mathbf{a}),$$

and let $\{\tilde{\psi}^k\}_{k=0}^K$ be a perturbed trajectory satisfying

$$\tilde{\psi}^{k+1} = \tilde{\psi}^k + \eta(v_\theta(\tilde{\psi}^k, t_k | \mathbf{s}, \mathbf{a}) + \xi_k), \quad \tilde{\psi}^0 = \psi^0 = \mathbf{z},$$

where $\{\xi_k\}_{k=0}^{K-1}$ are *arbitrary* perturbations or errors incurred in velocity evaluations. We say that a trained flow-matching critic exhibits *test-time recovery* if the terminal error of the flow $\Delta_K(\mathbf{z}) := \tilde{\psi}^K - \psi^K$ satisfies: $\|\Delta_K(\mathbf{z})\| \leq \beta_K \sum_{k=0}^{K-1} \eta \|\xi_k\|$, for a stability factor $\beta_K < 1$ that decreases with K .

H DEFINITIONS

I BROADER IMPACTS

This work contributes to a better understanding of flow-matching methods in reinforcement learning (RL). We do not foresee any direct societal impacts specific to this work beyond those generally associated with advances in reinforcement learning research.

J THEORETICAL ANALYSIS OF PLASTICITY IN THE LINEAR SETTING

Theoretical analysis in a linear setting. To isolate the source of plasticity under non-stationary TD targets, we analyze a linear setting comparing a monolithic critic to a linear flow-matching critic. Although both represent linear predictors with comparable computation graphs, their adaptation dynamics differ fundamentally. Our analysis (Theorem J.1) shows that even when feature directions are frozen, flow-matching critics can continue adapting purely by reweighting existing features through gain dynamics induced by the integration process. This enables motion within the span of learned features without overwriting them, providing a mechanistic explanation for plasticity under non-stationary learning. In contrast, monolithic critics necessarily modify their feature representation to adapt. Full derivations, extensions to ensembles, and discussion are provided in Appendix J.

To understand why flow-matching critics preserve plasticity under non-stationary TD targets, we analyze a simple yet informative linear setting. We consider TD-learning with a *linear* flow-matching critic and compare it to a corresponding “deep” linear ResNet (Arora et al., 2018). Although both represent linear predictors with comparable computation graphs, their learning dynamics differ.

Linear model setup. We consider learning predictors on inputs $\mathbf{x} \in \mathbb{R}^d$. A *monolithic* critic is given by $f_{\text{mono}}(\mathbf{x}; t) \triangleq w(t)^\top \mathbf{x}$, where $w(t) \in \mathbb{R}^d$ is trained directly against a non-stationary TD target, and may be composed of several linear residual blocks. Since the input interpolant in each step of integration in a *flow-matching* critic is the output from the previous step, it naturally makes the resulting function recursive in the input \mathbf{x} . A compact representation of the (expected) output of the linear flow network can be represented via two sets of parameters: **a)** linear velocity slices, $\{u_m(t)\}_{m=1}^{T-1}$, and **b)** gains $\{v_k(t)\}_{k=1}^{T-1}$ as $f_{\text{FM}}(\mathbf{x}; t) \triangleq \sum_{m=1}^{T-1} \beta_m(t) u_m(t)^\top \mathbf{x}$, $\beta_m(t) \triangleq \sum_{k=m}^{T-1} \alpha_k v_k(t)$. Here, $\{\alpha_k\}$ denotes the step sizes of the integration process. The gain $v_k(t)$ intuitively captures the contribution of the k -th integration step to the final Q-value prediction as we unroll the integration process. Although both predictors represent linear functions of \mathbf{x} , we show that the monolithic network can adapt only by updating $w(t)$, whereas the flow-matching model can adapt by reweighting existing features $\{u_m\}$ through the gain dynamics $\{v_k(t)\}$. Our theoretical result for this linear setting is shown in Theorem J.1. A discussion of Theorem J.1 is provided in Appendix O.

This theorem explains the mechanism underlying plasticity in the linear setting. Specifically, when the feature directions $\{u_m\}$ are held fixed over an interval, a flow-matching critic can still track changes in the regression target by adapting the predictor exclusively through the coefficients $\beta_m(t)$. These coefficients are driven by the gain dynamics $\{v_k(t)\}$ (Appendix O) induced by the integration process, rather than by changes to the features themselves. Thus, adaptation is mediated by a redistribution

Theorem J.1 (Feature-reweighting can happen in flow critics but not in monolithic critics). Fix a time interval $[t_0, t_1]$ and suppose $\dot{u}_m(t) = 0$ for all m and $t \in [t_0, t_1]$. Then, the following conditions hold for monolithic and flow-matching critics:

1. (**Monolithic**). If $\partial_t f_{\text{mono}}(\cdot; t) \neq 0$ on $[t_0, t_1]$, then necessarily $\dot{w}(t) \neq 0$.
2. (**Flow-matching**). The Euler flow-matching predictor satisfies

$$\partial_t f_{\text{FM}}(\mathbf{x}; t) = \left(\sum_{m=1}^{T-2} \dot{\beta}_m(t) u_m \right)^\top \mathbf{x},$$

with $\dot{\beta}_m(t)$ driven by the gain dynamics $\{\dot{v}_k(t)\}$.

3. (**Time-selective responsiveness**). The sensitivity of $\dot{v}_k(t)$ to changes (drift) in target satisfies

$$\frac{\partial \dot{v}_k(t)}{\partial b(t)} = -2(1 - \alpha_k)u_k(t),$$

so later slices (small α_k) respond strongly to drift, while early slices respond weakly.

of contributions across integration slices, allowing the predictor to reallocate mass among existing features without overwriting them $u_m(t)$.

In contrast, a predictor $f_{\text{mono}}(\mathbf{x}; t) = w(t)^\top \mathbf{x}$ admits no such decomposition: any change in the predictor necessarily requires $\dot{w}(t) \neq 0$, and thus directly rewriting the feature vector. Equivalently, **while flow matching allows motion in the predictor space within the span of previously learned features, monolithic training must modify the feature span itself to track non-stationary targets.** This quantifies the intuition about plasticity discussed at the beginning of this section. Although derived in a linear setting, this analysis offers insight into the training dynamics observed in our nonlinear experiments. For completeness, we extend Theorem J.1 to ensembles of monolithic networks in Appendix O.8, and identify a similar limitation.

K ADDITIONAL RESULTS FOR FLOW-MATCHING CRITICS

Post-layerNorm feature Norms for flow-matching critics (f1oq) vs monolithic critics (FQL).

We show in Figure 8 that while the last hidden layer adapts to the Q -scale for both f1oq and monolithic critics, the penultimate hidden layer in f1oq exhibits a rapid decrease in feature norms compared to FQL. This indicates that f1oq learns more adaptive representations in the penultimate hidden layer that are largely decoupled from the Q -value scale.

Penultimate hidden layer feature norms comparison for f1oq (TD), f1oq (SARSA) and f1oq (MC).

We show in Figure 9 that f1oq trained with TD shows the fastest decrease in penultimate hidden layer feature norms, whereas f1oq (SARSA) and f1oq (MC) trends resemble those of the monolithic FQL critics. This suggests that flow-matching critics develop better representations under non-stationary TD targets.

Frozen features with a single integration step. We show in Figure 10 that flow-matching critics with one integration step are more stable than monolithic critics, but fall short in learning plastic features than full flow matching critics. This highlights that integration plays a crucial role in preserving feature plasticity.

L BENCHMARKS

Following evaluation protocols from recent work in offline RL (Park et al., 2025b; Wagenmaker et al., 2025; Espinosa-Dice et al., 2025b;a; Dong et al., 2025), we use the **OGBench** task suite (Park et al., 2025a) as our main evaluation benchmark (see Figure 11). OGBench provides a number of diverse, challenging tasks across robotic locomotion and manipulation, where these tasks are generally more challenging than standard D4RL tasks (Fu et al., 2020), which have been saturated as of 2024 (Tarasov et al., 2023; Rafailov et al., 2024; Park et al., 2024). While OGBench was originally designed for benchmarking offline goal-conditioned RL, we use its reward-based single-task variants (“-singletask” from Park et al. (2025b)).

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

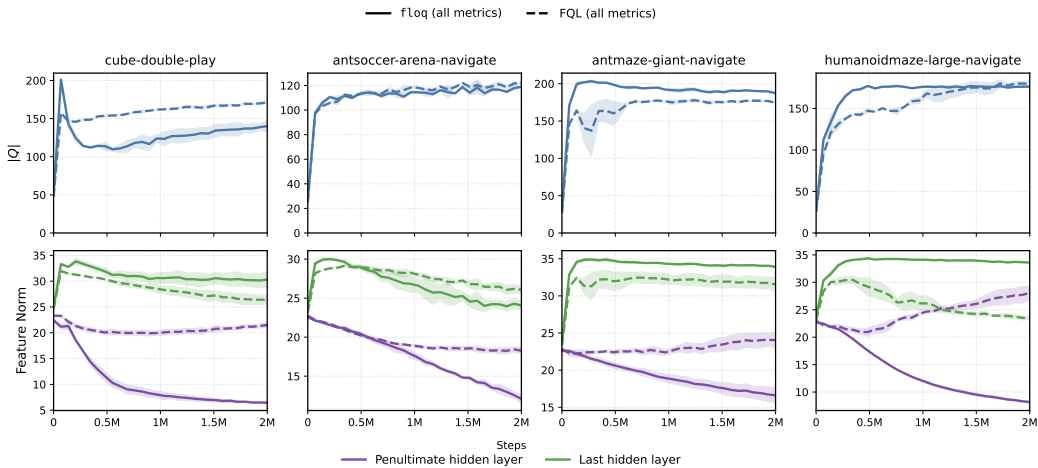


Figure 8: While the last hidden layer adapts to the Q-value scale for both `fLoq` and monolithic critics, the penultimate hidden layer layer in `fLoq` exhibits a qualitatively distinct, rapid decrease in feature norms. This indicates that `fLoq` learns more adaptive representations in the penultimate hidden layer that are largely decoupled from the Q-value scale.

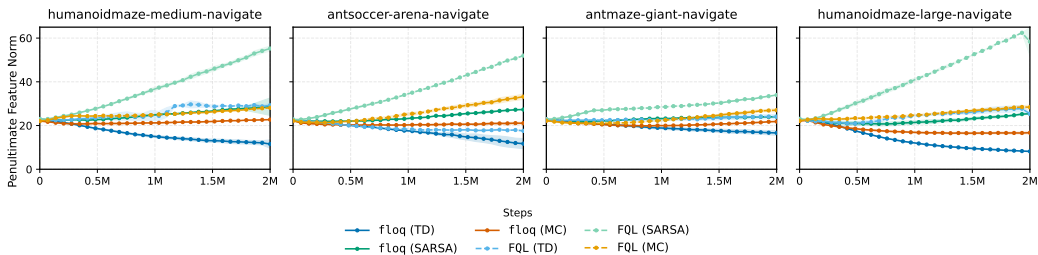


Figure 9: Flow-matching critics trained with TD-learning (`fLoq`) shows the fastest decrease in penultimate hidden layer feature norms, whereas `fLoq` (SARSA) and `fLoq` (MC) trends resemble those of the monolithic FQL critics. Thus flow-matching critics develop particularly better representations under non-stationary TD targets.

We used the `default` task in each OGBench environment for our analysis and RLPD experiments, following the protocol in recent works (Park et al., 2025b; Wagenmaker et al., 2025; Espinosa-Dice et al., 2025b;a; Dong et al., 2025).

M EXPERIMENTAL DETAILS AND HYPERPARAMETERS

M.1 OFFLINE RL ANALYSIS EXPERIMENTS (SECTION G, SECTION 2, SECTION 3).

Experimental details for robustness to noisy TD supervision (Figure 1, Section 2.2). We added $\text{Unif}[-\kappa, \kappa]$ (for $\kappa \in \{0, 4, 8, 16\}$) noise to the Q-network targets for FQL (monolithic critics) and to the velocity-network targets for `fLoq` (flow-matching critics).

Hyperparameters. We used the hyper-parameters from Agrawalla et al. (2025) for both `fLoq` (flow-matching critics) and FQL (monolithic critics) for all offline RL analysis experiments.

M.2 RLPD EXPERIMENTS (SECTION 4).

We tuned the BC-regularization coefficient (α) in the range $[10, 100]$ (step size 10) for both `fLoq` and FQL. All other hyper-parameters were kept to the same values from Agrawalla et al. (2025).

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

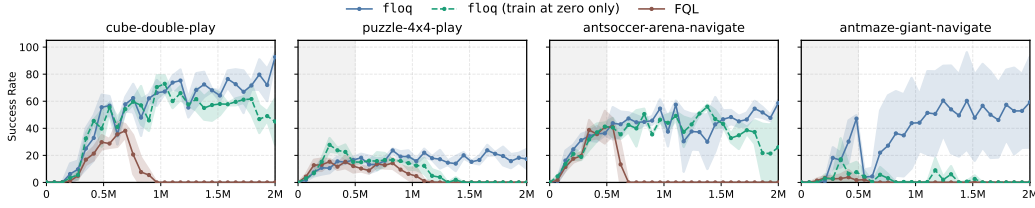


Figure 10: Flow-matching critics with one integration step are more stable than monolithic critics but less stable than full flow matching. This emphasizes that test-time recovery in the form of integration plays a crucial role in preserving feature plasticity.

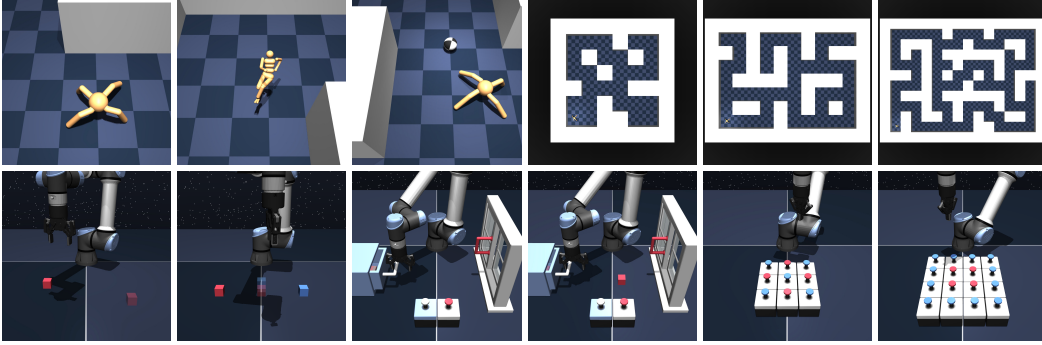


Figure 11: **OGBench** (Park et al., 2025a) domains. These tasks include high-dimensional state and action spaces, sparse rewards, stochasticity, as well as hierarchical structure.

N THEORETICAL RESULT FOR TEST-TIME RECOVERY (SECTION 2.1).

Theorem N.1. Fix a state-action pair (\mathbf{s}, \mathbf{a}) and let $\psi_\theta(t, \mathbf{z} | \mathbf{s}, \mathbf{a})$ denote the flow interpolant obtained by integrating the velocity field v_θ . Assume that the velocity field v_θ satisfies

$$\frac{\partial v_\theta(x, t; \mathbf{s}, \mathbf{a})}{\partial x} \leq -\frac{c_1}{1-t}.$$

for some constant $0 < c_1 < 1$ and all $t \in (0, 1 - \frac{1}{K})$.

Consider a K -step numerical integrator with step size $\eta = 1/K$ and discrete times $t_k = k/K$. Let $\{\psi^k\}_{k=0}^K$ be the unperturbed trajectory defined by

$$\psi^{k+1} = \psi^k + \eta v_\theta(\psi^k, t_k | \mathbf{s}, \mathbf{a}),$$

and let $\{\tilde{\psi}^k\}_{k=0}^K$ be a perturbed trajectory satisfying

$$\tilde{\psi}^{k+1} = \tilde{\psi}^k + \eta(v_\theta(\tilde{\psi}^k, t_k | \mathbf{s}, \mathbf{a}) + \xi_k), \quad \tilde{\psi}^0 = \psi^0 = \mathbf{z},$$

where $\{\xi_k\}_{k=0}^{K-1}$ are arbitrary perturbations (or errors) in velocity evaluations. Assume further that each $\|\xi_k\| \leq \varepsilon$ for some $\varepsilon > 0$.

Then the terminal error $\Delta_K(\mathbf{z}) := \tilde{\psi}^K - \psi^K$ satisfies

$$\|\Delta_K(\mathbf{z})\| \leq \beta_K \varepsilon$$

where

$$\beta_K = c_2 K^{-c_1}$$

for some constant $c_2 > 0$.

Note. We argue here that assumption on v_θ is reasonable for flow-matching networks (by choosing any appropriate $c_1 < 1$).

This is because the learned trajectories are straight to the first order (the curvature, while being important to distinguish actions, appears at the scale of advantages; which are typically much smaller

than the Q -value - see Figure 16 in Agrawalla et al. (2025) for an example). Moreover, the variance in flow outputs is much smaller than the range of initial noise for $\text{fl}\circ\text{q}$ style expected-value training (see Table 2).

This implies that to the first order, we have

$$v_\theta(x, t; s, a) \approx \frac{Q(s, a) - x}{1 - t} \implies \frac{\partial v_\theta(x, t; s, a)}{\partial x} \approx -\frac{1}{1 - t} \leq -\frac{c_1}{1 - t}$$

for an appropriate constant $c_1 < 1$.

Proof. We analyze the evolution of the error between the unperturbed trajectory ψ^k and the perturbed trajectory $\tilde{\psi}^k$.

1. Error Dynamics Let $\Delta_k = \tilde{\psi}^k - \psi^k$. By subtracting the update equations for the two trajectories, we obtain:

$$\begin{aligned} \Delta_{k+1} &= \tilde{\psi}^{k+1} - \psi^{k+1} \\ &= (\tilde{\psi}^k - \psi^k) + \eta \left[v_\theta(\tilde{\psi}^k, t_k | \mathbf{s}, \mathbf{a}) - v_\theta(\psi^k, t_k | \mathbf{s}, \mathbf{a}) \right] + \eta \xi_k \\ &= \Delta_k + \eta \left[v_\theta(\tilde{\psi}^k, t_k | \mathbf{s}, \mathbf{a}) - v_\theta(\psi^k, t_k | \mathbf{s}, \mathbf{a}) \right] + \eta \xi_k \end{aligned}$$

2. Applying the Mean Value Theorem By the Mean Value Theorem (MVT) for vector-valued functions, the difference in the velocity fields can be expressed as:

$$v_\theta(\tilde{\psi}^k, t_k) - v_\theta(\psi^k, t_k) = \left(\int_0^1 \nabla_x v_\theta(\psi^k + \tau \Delta_k, t_k) d\tau \right) \Delta_k$$

Applying the assumption $\frac{\partial v_\theta}{\partial x} \leq -\frac{c_1}{1-t}$, and substituting $t_k = k/K$ and $\eta = 1/K$:

$$\begin{aligned} \Delta_{k+1} &\leq \left(1 - \frac{\eta c_1}{1 - t_k} \right) \Delta_k + \eta \xi_k \\ &= \left(1 - \frac{(1/K)c_1}{1 - (k/K)} \right) \Delta_k + \eta \xi_k \\ &= \left(1 - \frac{c_1}{K - k} \right) \Delta_k + \eta \xi_k \end{aligned}$$

3. Unrolling the Recurrence Since $\tilde{\psi}^0 = \psi^0$, we have $\Delta_0 = 0$. Unrolling the recurrence from $k = 0$ to $K - 1$:

$$\begin{aligned} \|\Delta_K\| &\leq \sum_{k=0}^{K-1} \left(\prod_{j=k+1}^{K-1} \left(1 - \frac{c_1}{K-j} \right) \right) \eta \|\xi_k\| \\ &\leq \frac{\varepsilon}{K} \left[[(1 - c_1)] + [(1 - c_1)(1 - \frac{c_1}{2})] + [(1 - c_1)(1 - \frac{c_1}{2})(1 - \frac{c_1}{3})] + \dots + [(1 - c_1)(1 - \frac{c_1}{2})(1 - \frac{c_1}{3}) \dots (1 - \frac{c_1}{K})] \right] \\ &\leq \frac{\varepsilon(c_2 \cdot K^{1-c_1})}{K} \\ &\leq (c_2 K^{-c_1}) \varepsilon \end{aligned}$$

for a constant $c_2 > 0$, as desired. \square

O COMPARING FLOW-MATCHING AND MONOLITHIC NETWORKS IN THE LINEAR SETTING (SECTION J)

Fix $T \geq 3$ and step size $h := \frac{1}{T-1}$. Consider the linear Euler recursion

$$s_{i+1}(x, z; t) = (1 + h v_i(t)) s_i(x, z; t) + h u_i(t)^\top x, \quad i = 1, \dots, T-1, \quad (3)$$

$$s_1(x, z; t) = z, \quad \hat{y}(x, z; t) := s_T(x, z; t), \quad (4)$$

with parameters $u_i(t) \in \mathbb{R}^d$, $v_i(t) \in \mathbb{R}$.

864 **Gain parameters.** In the Euler recursion

$$865 s_{i+1} = (1 + hv_i)s_i + hu_i^\top x,$$

866 the scalar parameter v_i is a *gain*: it multiplicatively scales the incoming state before propagation.
867 Positive gain amplifies the accumulated state, while negative gain attenuates it. Because the end-to-
868 end contribution of an injected feature is

$$869 \beta_m = h \prod_{j>m} (1 + hv_j),$$

870 downstream gains determine how much signal injected at earlier slices survives to the output.
871 Learning the gains therefore corresponds to learning how strongly the model should trust previously
872 accumulated information.

873 O.1 UNROLLED PREDICTOR

874 **Lemma O.1** (Unrolled form). *Define*

$$875 P_{>m}(t) := \prod_{j=m+1}^{T-1} (1 + hv_j(t)), \quad \beta_m(t) := h P_{>m}(t).$$

876 Then

$$877 \hat{y}(x, z; t) = \left(\prod_{j=1}^{T-1} (1 + hv_j(t)) \right) z + \sum_{m=1}^{T-1} \beta_m(t) u_m(t)^\top x. \quad (5)$$

878 Consequently, if $\mathbb{E}[z] = 0$,

$$879 f(x; t) := \mathbb{E}_z[\hat{y}(x, z; t)] = \sum_{m=1}^{T-1} \beta_m(t) u_m(t)^\top x. \quad (6)$$

880 O.2 FLOW-MATCHING GRADIENT FLOW

881 Let $Y(t)$ be a time-dependent target and Z independent noise with $\mathbb{E}[Z] = 0$, $\text{Var}(Z) = \sigma_z^2$. Define

$$882 \alpha_i := \frac{T-i}{T-1}, \quad S_i(t) := \alpha_i Z + (1 - \alpha_i)Y(t).$$

883 Each slice is trained with the local quadratic loss

$$884 \mathcal{L}_i(t, w_i) = \mathbb{E}[(w_i^\top [X, S_i(t)] - (Y(t) - Z))^2], \quad w_i := \begin{bmatrix} u_i \\ v_i \end{bmatrix}.$$

885 Define the moment matrices

$$886 A_i(t) := \mathbb{E}[\tilde{X}_i(t)\tilde{X}_i(t)^\top], \quad b_i(t) := \mathbb{E}[\tilde{X}_i(t)(Y(t) - Z)], \quad \tilde{X}_i(t) := [X, S_i(t)].$$

887 Then gradient flow is

$$888 \dot{w}_i(t) = -2(A_i(t)w_i(t) - b_i(t)). \quad (7)$$

889 O.3 GRADIENT FLOW OF THE PREDICTION

890 **Theorem O.2** (Exact decomposition). *The time derivative of the mean predictor $f(x; t) = \sum_m \beta_m(t) u_m(t)^\top x$ satisfies*

$$891 \frac{d}{dt} f(x; t) = \underbrace{\sum_{m=1}^{T-1} \beta_m(t) \dot{u}_m(t)^\top x}_{\text{Term 2a}} + \underbrace{\sum_{k=2}^{T-1} s_k(x, 0; t) \frac{h \dot{v}_k(t)}{1 + hv_k(t)}}_{\text{Term 2b}}, \quad (8)$$

892 where $s_k(x, 0; t) = \sum_{m<k} \beta_m(t) u_m(t)^\top x$.

893 Equivalently,

$$894 \frac{d}{dt} f(x; t) = \sum_{m=1}^{T-1} \beta_m(t) \dot{u}_m(t)^\top x + \sum_{m=1}^{T-2} \dot{\beta}_m(t) u_m(t)^\top x, \quad \dot{\beta}_m(t) = \beta_m(t) \sum_{k=m+1}^{T-1} \frac{h \dot{v}_k(t)}{1 + hv_k(t)}.$$

918 O.4 MOVING TARGETS: EXPLICIT GAIN DYNAMICS

919 Let

$$920 \quad \Sigma := \mathbb{E}[XX^\top], \quad b(t) := \mathbb{E}[XY(t)], \quad q(t) := \mathbb{E}[Y(t)^2].$$

921 No distributional assumptions are made beyond existence of these moments.

922 **Lemma O.3** (Gain dynamics). *For each slice k ,*

$$923 \quad \dot{v}_k(t) = -2 \left((1 - \alpha_k) b(t)^\top u_k(t) + ((1 - \alpha_k)^2 q(t) + \alpha_k^2 \sigma_z^2) v_k(t) - (1 - \alpha_k) q(t) + \alpha_k \sigma_z^2 \right). \quad (9)$$

924 O.5 COMPARISON WITH MONOLITHIC TRAINING

925 A monolithic linear predictor $f_{\text{mono}}(x; t) = w(t)^\top x$ trained by squared loss satisfies

$$926 \quad \dot{w}(t) = -2(\Sigma w(t) - b(t)). \quad (10)$$

927 Thus its instantaneous movement is

$$928 \quad \partial_t f_{\text{mono}}(x; t) = \dot{w}(t)^\top x,$$

929 and any adaptation requires direct motion of $w(t)$.

930 O.6 MAIN RESULT

931 **Theorem O.4** (Feature-reweighting advantage under moving targets). *Assume $\Sigma \succ 0$ and finite second moments. Fix a time interval $[t_0, t_1]$ and suppose*

$$932 \quad \dot{u}_m(t) = 0 \quad \text{for all } m \text{ and } t \in [t_0, t_1].$$

- 933 1. (Monolithic) *If $\partial_t f_{\text{mono}}(\cdot; t) \neq 0$ on $[t_0, t_1]$, then necessarily $\dot{w}(t) \neq 0$.*
- 934 2. (Flow matching) *The Euler flow-matching predictor satisfies*

$$935 \quad \partial_t f_{\text{FM}}(x; t) = \left(\sum_{m=1}^{T-2} \dot{\beta}_m(t) u_m \right)^\top x,$$

936 *with $\dot{\beta}_m(t)$ driven entirely by the gain dynamics $\{\dot{v}_k(t)\}$.*

- 937 3. (Time-selective responsiveness) *The sensitivity of $\dot{v}_k(t)$ to target drift satisfies*

$$938 \quad \frac{\partial \dot{v}_k(t)}{\partial b(t)} = -2(1 - \alpha_k) u_k(t),$$

939 *so later slices (α_k small) respond strongly to drift, while early slices respond weakly.*

940 Hence, under moving targets, flow matching admits predictor motion via reweighting existing features (Term 2b), whereas monolithic training can adapt only by rewriting its feature vector.

941 O.6.1 INTERPRETATION

942 Although both models represent linear predictors, flow matching decomposes learning into two distinct channels: (i) feature learning through \dot{u}_m (Term 2a), and (ii) feature reweighting through downstream gain updates \dot{v}_k (Term 2b). Under moving targets, the gain dynamics are slice-dependent and explicitly modulated by the target moments, enabling rapid reallocation of importance among previously learned features without altering them. Monolithic training lacks this mechanism and must instead directly overwrite its parameter vector to track drift.

943 O.7 FEATURE REWEIGHTING VIA DOWNSTREAM GAINS

944 We formalize the notion of *feature reweighting* induced by Term 2b.

972 O.7.1 DEFINITION (FEATURE REPRESENTATION)

973 At any training time t , define the collection of slice features

$$974 \mathcal{U}(t) := \{u_1(t), \dots, u_{T-1}(t)\} \subset \mathbb{R}^d,$$

975 and the corresponding effective predictor

$$976 f(x; t) = w_{\text{eff}}(t)^\top x, \quad w_{\text{eff}}(t) := \sum_{m=1}^{T-1} \beta_m(t) u_m(t),$$

977 where

$$978 \beta_m(t) := h \prod_{j=m+1}^{T-1} (1 + hv_j(t)).$$

979 Thus the model represents a linear function as a weighted combination of fixed feature vectors $u_m(t)$.

980 O.7.2 DEFINITION (FEATURE REWEIGHTING)

981 We say that the predictor undergoes *feature reweighting* on a time interval $I \subset \mathbb{R}_+$ if

$$982 \dot{u}_m(t) = 0 \quad \forall m, \forall t \in I, \quad \text{but} \quad \dot{\beta}_m(t) \neq 0 \text{ for some } m \text{ and } t \in I.$$

983 Equivalently, feature reweighting means that the predictor $f(\cdot; t)$ changes over time while the feature set $\mathcal{U}(t)$ remains fixed.

984 O.7.3 PROPOSITION (EXACT REWEIGHTING IDENTITY)

985 Under the Euler flow-matching dynamics,

$$986 \dot{w}_{\text{eff}}(t) = \sum_{m=1}^{T-1} \beta_m(t) \dot{u}_m(t) + \sum_{m=1}^{T-2} \dot{\beta}_m(t) u_m(t),$$

987 with

$$988 \dot{\beta}_m(t) = \beta_m(t) \sum_{k=m+1}^{T-1} \frac{h \dot{v}_k(t)}{1 + hv_k(t)}.$$

989 In particular, if $\dot{u}_m(t) = 0$ for all m on an interval I , then

$$990 \dot{w}_{\text{eff}}(t) = \sum_{m=1}^{T-2} \dot{\beta}_m(t) u_m \quad \text{for all } t \in I,$$

991 so the predictor changes *purely* by reweighting existing features.

992 O.7.4 INTERPRETATION

993 The decomposition above shows that the Euler flow-matching model separates learning into two distinct mechanisms:

- 994 • **Feature learning (Term 2a):** updates of $u_m(t)$ change the feature directions available to the model.
- 995 • **Feature reweighting (Term 2b):** updates of downstream gains $\{v_k(t)\}$ modify the coefficients $\{\beta_m(t)\}$, thereby changing the relative importance of previously learned features without altering them.

996 Crucially, feature reweighting is mediated by downstream gains: changing a single $v_k(t)$ rescales *all* features injected at earlier slices $m < k$. This creates a hierarchical structure in which earlier features may be introduced tentatively and later amplified or suppressed depending on downstream dynamics.

997 This mechanism is absent in monolithic linear training. For a monolithic predictor $f_{\text{mono}}(x; t) = w(t)^\top x$, any change in the predictor necessarily requires $\dot{w}(t) \neq 0$, meaning that features themselves must be rewritten. In contrast, flow matching admits predictor adaptation through coefficient updates alone, enabling the model to track changes in the target by reallocating importance among existing features.

Summary. Feature reweighting is the ability of the flow-matching model to change its prediction by adjusting downstream gains, while keeping the feature directions fixed. Term 2b is exactly the mathematical expression of this mechanism.

O.8 WHY AN ENSEMBLE OF MONOLITHIC NETWORKS CANNOT REPLICATE TERM 2B

In this section we formalize the statement that an ensemble of independently trained monolithic predictors cannot realize the *feature reweighting* mechanism provided by Term 2b in flow matching.

O.8.1 SETUP: MONOLITHIC ENSEMBLE

Let $(X, Y(t))$ be a time-indexed (moving-target) data process with finite second moments, and let π_1, \dots, π_K be fixed ensemble weights with $\pi_k \geq 0$ and $\sum_{k=1}^K \pi_k = 1$.

Each ensemble member is a monolithic linear predictor

$$f^{(k)}(x; t) := w_k(t)^\top x,$$

trained by gradient flow on the time-varying squared loss

$$\mathcal{L}^{(k)}(t, w) := \mathbb{E}[(w^\top X - Y(t))^2].$$

Then each member evolves according to

$$\dot{w}_k(t) = -2(\Sigma w_k(t) - b(t)), \quad \Sigma := \mathbb{E}[X X^\top], \quad b(t) := \mathbb{E}[X Y(t)].$$

The ensemble predictor is defined as the weighted average

$$f_{\text{ens}}(x; t) := \sum_{k=1}^K \pi_k f^{(k)}(x; t) = \sum_{k=1}^K \pi_k w_k(t)^\top x.$$

O.8.2 THEOREM: NO FEATURE REWEIGHTING CHANNEL IN A FIXED-WEIGHT ENSEMBLE

Theorem O.5 (Ensembles do not create a Term-2b-like reweighting mechanism). *Assume Σ exists and is finite and that $b(t)$ exists for all t . Define the ensemble-averaged parameter*

$$\bar{w}(t) := \sum_{k=1}^K \pi_k w_k(t).$$

Then:

1. (**Function-level collapse**) *The ensemble predictor is exactly a single monolithic predictor:*

$$f_{\text{ens}}(x; t) = \bar{w}(t)^\top x.$$

2. (**Dynamics collapse**) *The averaged parameter follows the same monolithic gradient flow:*

$$\dot{\bar{w}}(t) = -2(\Sigma \bar{w}(t) - b(t)).$$

3. (**No reweighting without parameter motion**) *For every t ,*

$$\partial_t f_{\text{ens}}(x; t) = \dot{\bar{w}}(t)^\top x = \sum_{k=1}^K \pi_k \dot{w}_k(t)^\top x.$$

In particular, if $\dot{w}_k(t) = 0$ for all k on a time interval I , then $f_{\text{ens}}(\cdot; t)$ is constant on I .

Consequently, a fixed-weight ensemble cannot change its predictor by reweighting previously learned features while keeping member parameters fixed; any change in f_{ens} requires direct motion of at least one member parameter $w_k(t)$.

1080 *Proof.* (1) Linearity gives

$$1081$$

$$1082 f_{\text{ens}}(x; t) = \sum_{k=1}^K \pi_k w_k(t)^\top x = \left(\sum_{k=1}^K \pi_k w_k(t) \right)^\top x = \bar{w}(t)^\top x.$$

$$1083$$

$$1084$$

1085 (2) Differentiate $\bar{w}(t) = \sum_k \pi_k w_k(t)$ and substitute the gradient flows:

$$1086$$

$$1087 \dot{\bar{w}}(t) = \sum_{k=1}^K \pi_k \dot{w}_k(t) = \sum_{k=1}^K \pi_k (-2(\Sigma w_k(t) - b(t))) = -2(\Sigma \sum_{k=1}^K \pi_k w_k(t) - b(t) \sum_{k=1}^K \pi_k),$$

$$1088$$

$$1089$$

1090 which equals $-2(\Sigma \bar{w}(t) - b(t))$ since $\sum_k \pi_k = 1$.

1091 (3) Differentiate $f_{\text{ens}}(x; t) = \bar{w}(t)^\top x$ to obtain $\partial_t f_{\text{ens}}(x; t) = \dot{\bar{w}}(t)^\top x$ and substitute $\dot{\bar{w}}(t) = \sum_k \pi_k \dot{w}_k(t)$. If $\dot{w}_k(t) = 0$ for all k on I , then $\dot{\bar{w}}(t) = 0$ on I , hence $f_{\text{ens}}(\cdot; t)$ is constant. \square

1092 O.8.3 INTERPRETATION AND COMPARISON TO TERM 2B

1093 Theorem O.5 shows that a fixed-weight ensemble does not introduce an additional adaptation channel: its mean prediction is governed by the same single-parameter monolithic dynamics.

1094 By contrast, in the Euler flow-matching model the effective parameter admits the decomposition

$$1095$$

$$1096 w_{\text{eff}}(t) = \sum_{m=1}^{T-1} \beta_m(t) u_m(t), \quad \dot{w}_{\text{eff}}(t) = \sum_m \beta_m(t) \dot{u}_m(t) + \sum_m \dot{\beta}_m(t) u_m(t),$$

$$1097$$

$$1098$$

1099 so the predictor can move even when $\dot{u}_m(t) = 0$ via the coefficient dynamics $\dot{\beta}_m(t)$ induced by downstream gains (Term 2b). The ensemble lacks any analog of coefficients $\{\beta_m(t)\}$ that can evolve independently of the base feature parameters, hence it cannot realize feature reweighting in this sense.

1100 **Remark.** If one allows the ensemble weights π_k themselves to vary with t or to depend on x (a learned gating network), then additional mechanisms become possible. However, such a model is no longer an ensemble of *independent monolithic networks with fixed averaging*; it introduces an explicit mixer/gate, which is precisely the kind of additional structure that Term 2b provides in the flow-matching construction.

1101 P DISCUSSION AND PERSPECTIVES ON FUTURE WORK

1102 This work explains the effectiveness of flow-matching critics in off-policy reinforcement learning by identifying dense, trajectory-level supervision as the key mechanism. Rather than gains from distributional modeling or expressivity, flow-matching critics learn a velocity field jointly with an integration procedure, enabling test-time recovery and preserving representation plasticity under non-stationary TD targets. In contrast to monolithic critics, this coupling between training and iterative computation allows robust adaptation to noise and target drift, highlighting dense intermediate supervision as a powerful inductive bias for stabilizing TD learning.

1103 Our findings open several directions for future work. On the practical side, it would be natural to explore alternative design choices for training flow-matching critics, such as representing velocities with higher-dimensional vectors rather than scalars, as well as to evaluate their effectiveness in settings that place stronger demands on plasticity, including continual learning with shifting task distributions. From a theoretical perspective, extending our analysis beyond the linear setting to nonlinear function approximation remains an important open problem. More broadly, our results point to interesting connections between plasticity under non-stationarity, flow matching, and test-time computation. While we study these interactions in the context of TD learning, the underlying principles governing flow-matching dynamics may extend to other domains, such as time-series modeling, suggesting a broader set of applications and theoretical questions for future investigation.