

# What Does Flow-Matching Bring to TD-Learning?

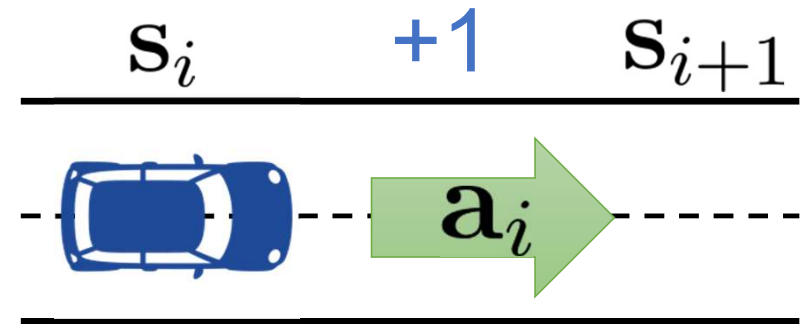
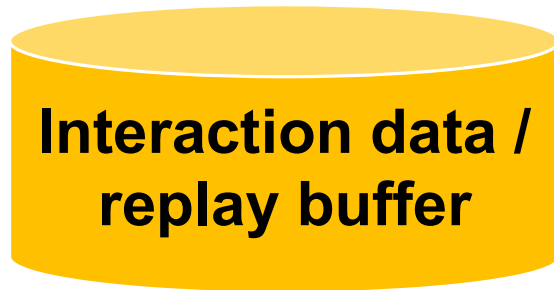
Bhavya Agrawalla, Michal Nauman, Aviral Kumar

**Carnegie Mellon University**  
School of Computer Science



UNIVERSITY  
OF WARSAW

# Primer: Notation and Terminology.



Action taken by data  
collection strategy

$$\{s_i, a_i, r(s_i, a_i), s_{i+1}\}_{i=1}^n$$

observation  
(state of the world)

reward value  
(how good this  
action was)

next  
observation

# Primer: TD-Learning.

**Bellman equation**

$$Q^\pi(\mathbf{s}_i, \mathbf{a}_i) := r(\mathbf{s}_i, \mathbf{a}_i) + \gamma Q^\pi(\mathbf{s}_{i+1}, \mathbf{a}_{i+1})$$

How can we learn to estimate  $Q$ ?

↓  
Convert into a frozen target

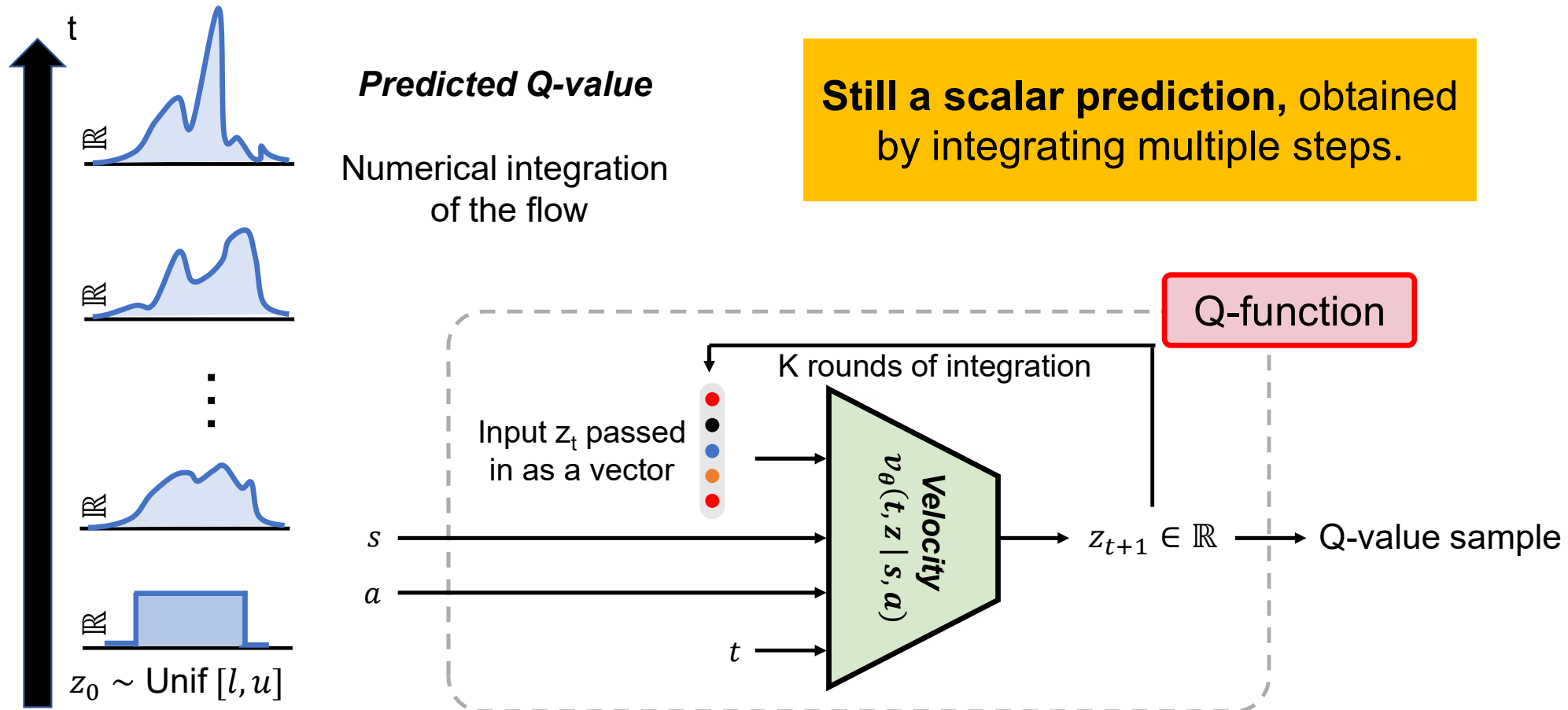
$$y(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma Q_{\theta^{\text{old}}}(\mathbf{s}_{i+1}, \mathbf{a}_{i+1})$$

**Training loss**

$$\min_{\theta} \mathbb{E}_{\text{data}} \left[ (Q_{\theta}(\mathbf{s}_i, \mathbf{a}_i) - y(\mathbf{s}_i, \mathbf{a}_i))^2 \right]$$

*(and periodically update the target network)*

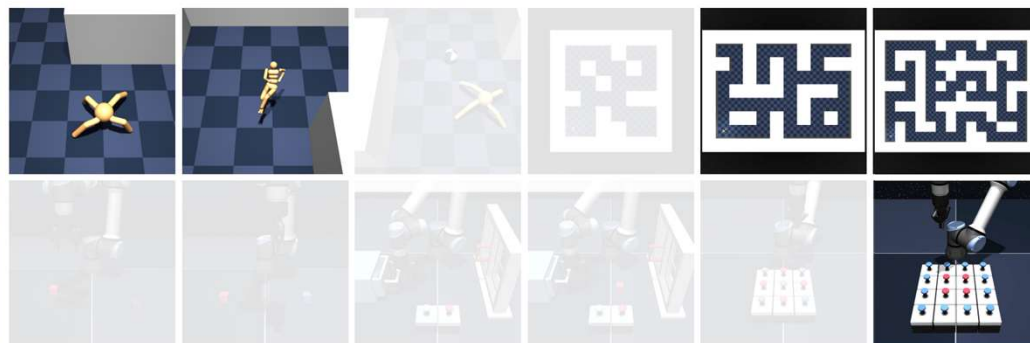
# Primer: Flow-Matching Q-Functions (f1oq).



# f1oq Performance.

## Simulated offline RL benchmarks

(*Hard tasks: precision is key*)



| Environment (5 tasks each)        | Gaussian Policy |        | Flow Policy  |          |         | Flow Q-function (Ours) |               |
|-----------------------------------|-----------------|--------|--------------|----------|---------|------------------------|---------------|
|                                   | BC              | ReBRAC | SORL         | FQL (1M) | FQL(2M) | f1oq (Def.)            | f1oq (Best)   |
| antmaze-large                     | 11 ±1           | 81 ±5  | 89 ±2        | 79 ±3    | 83 ±5   | <b>91 ±5</b>           | <b>91 ±5</b>  |
| antmaze-giant ( <i>Hard</i> )     | 0 ±0            | 26 ±8  | 9 ±6         | 22 ±19   | 27 ±23  | 36 ±21                 | <b>51 ±12</b> |
| hmmaze-medium                     | 2 ±1            | 22 ±8  | 64 ±4        | 57 ±5    | 69 ±20  | <b>82 ±10</b>          | <b>82 ±10</b> |
| hmmaze-large ( <i>Hard</i> )      | 1 ±0            | 2 ±1   | 5 ±2         | 9 ±6     | 16 ±9   | <b>28 ±9</b>           | <b>28 ±9</b>  |
| antsoccer-arena                   | 1 ±0            | 0 ±0   | <b>69 ±2</b> | 60 ±2    | 61 ±10  | 65 ±12                 | 65 ±12        |
| cube-single                       | 5 ±1            | 91 ±2  | 97 ±1        | 96 ±1    | 94 ±5   | <b>98 ±3</b>           | <b>98 ±3</b>  |
| cube-double ( <i>Hard</i> )       | 2 ±1            | 12 ±1  | 25 ±3        | 29 ±2    | 25 ±6   | 47 ±15                 | <b>47 ±15</b> |
| scene                             | 5 ±1            | 41 ±3  | 57 ±2        | 56 ±2    | 57 ±4   | <b>58 ±6</b>           | <b>58 ±6</b>  |
| puzzle-3x3 ( <i>Hard</i> )        | 2 ±0            | 21 ±1  | — ±—         | 30 ±1    | 29 ±5   | <b>37 ±7</b>           | <b>37 ±7</b>  |
| puzzle-4x4 ( <i>Hard</i> )        | 0 ±0            | 14 ±1  | — ±—         | 17 ±2    | 9 ±3    | 21 ±5                  | <b>28 ±6</b>  |
| Average Score (All Environments)  | 3               | 31     | —            | 46       | 47      | 56                     | <b>59</b>     |
| Average Score (Hard Environments) | 1               | 15     | —            | 21       | 21      | 34                     | <b>38</b>     |

~2x

~1.75x

~2x

~3x

# Why is Flow-Matching Effective for Q-Learning?

**Recent works [1,2,3,..] show significant gains with flow-matching Q-functions.**

**Why are flow-matching Q-functions effective? Our findings:**

- 1. Distribution modelling** does **NOT** explain their success.
- 2. Integration** enables **test-time recovery (TTR)** from imperfect value estimates.
- 3. Flow-matching loss** enables **plasticity** under non-stationary TD targets.

**[1]** Agrawalla et. al, floq: Training Critics via Flow-Matching for scaling compute in Value-Based RL.

**[2]** Dong et. al, Value flows.

**[3]** Espinosa-Dice et. al, Expressive Value-Learning for scalable offline RL.

# Distributional RL is an Insufficient Explanation.

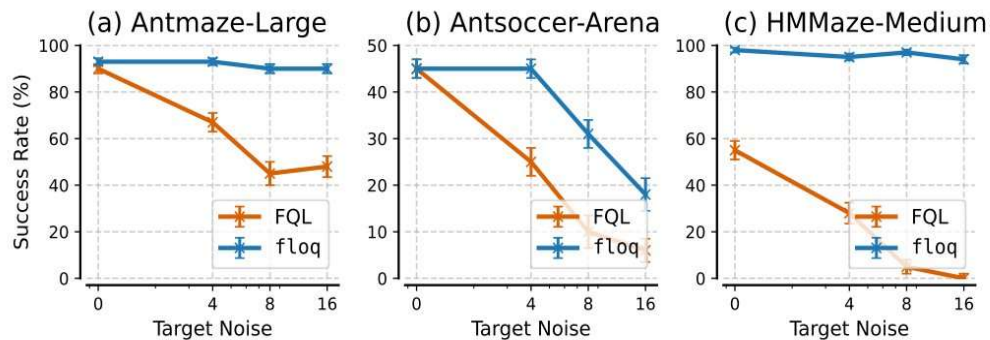
1. Expected-value TD backups (E) often outperform distributional RL backups (D).
2. Examples: Success rate is 52 v/s 30 (hmmaze-large), 86 v/s 74 (antmaze-giant).
3. Variance of final Q-values is much lower for E (last column), further supporting the claim.

**Table 1: Comparing expected-value (E) vs. distributional (D) flow-matching critics on representative OGBench tasks. Each entry reports E / D. While both variants learn similar expected Q-values, D produces higher-variance estimates (that more closely reflect variation in the return distribution) but does not outperform E.**

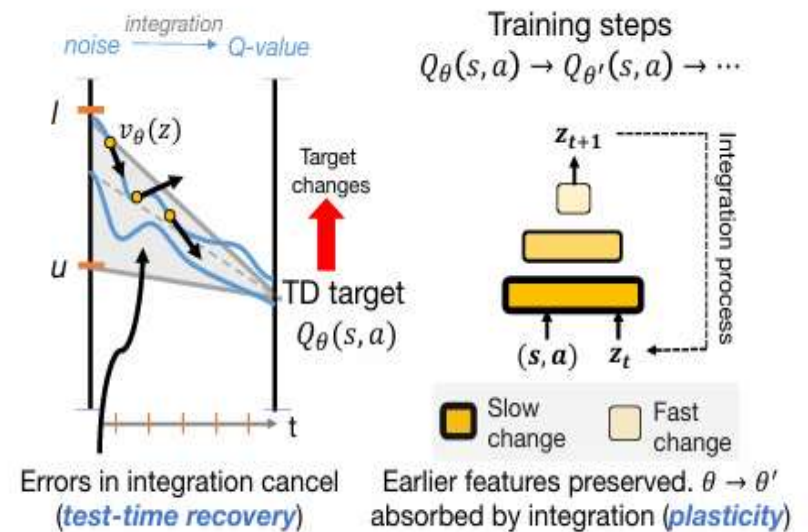
| Env.          | Success (%)    | $Q_\theta(\mathbf{s}, \mathbf{a})$ | $\text{Var}_z(Q)$ |
|---------------|----------------|------------------------------------|-------------------|
| hmmaze-large  | <b>52 / 30</b> | -180/ - 170                        | 0.2/ <b>4.5</b>   |
| antmaze-giant | <b>86 / 74</b> | -190/ - 200                        | 0.1/ <b>0.7</b>   |
| cube-double   | 72 / 72        | -130/ - 130                        | 1.1/ <b>6.3</b>   |
| hmmaze-medium | 94 / 94        | -170/ - 170                        | 0.3/ <b>2.3</b>   |

# f1oq Learns to “Recover” from Noisy Q-Values

1. When training against moving TD-targets, flow-critics learn to use the integration to correct errors. (“test-time recovery”).
2. As a result, they are more robust to noisy targets and learn plastic features.



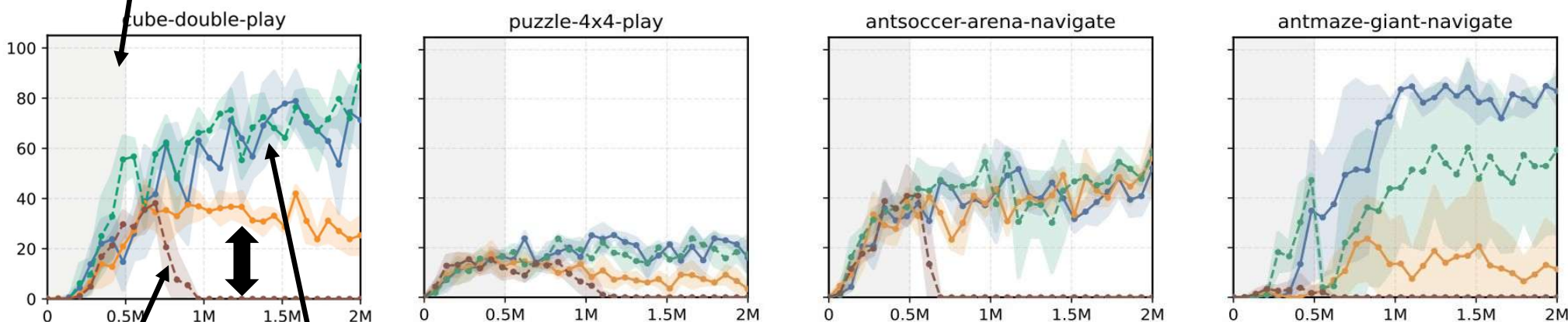
**Figure 2: Performance of flow-matching (f1oq) and monolithic (FQL) critics when trained with target noise.** Observe that flow-matching critics are much more robust to noise in TD targets, while performance of FQL (monolithic critic) degrades substantially faster, even when they start at a similar point (antmaze/antsoccer).



# Flow-Matching Enables Plastic Feature Learning.

Freeze features of the network here

floq floq (freeze) FQL FQL (freeze)



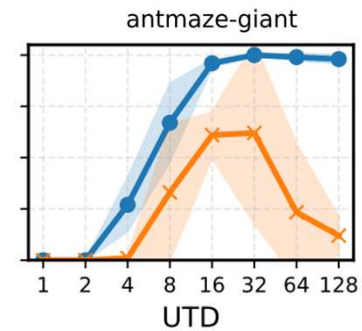
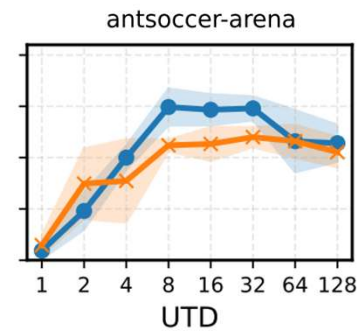
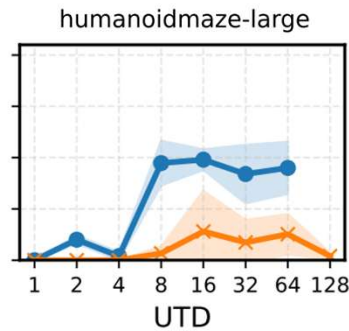
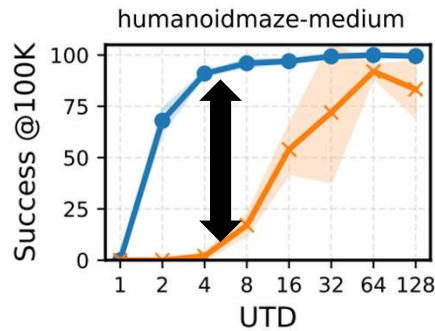
monolithic net

floq net

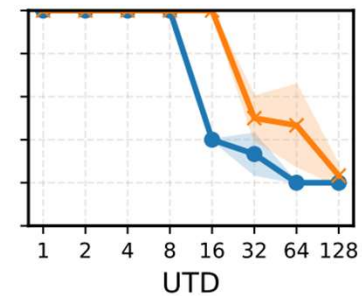
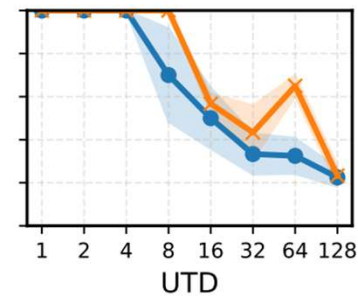
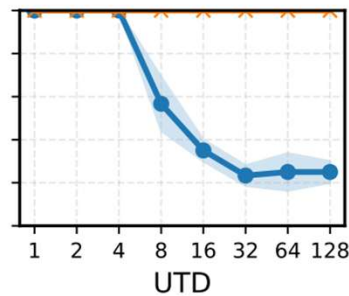
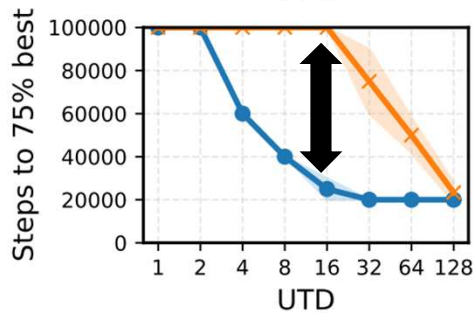
**Key finding:** floq features can support future targets (i.e., “*plastic*” features)

# Consequences of Plastic Feature Learning.

Success w/  
limited data



Sample  
efficiency



— floq — FQL

**Key finding:** Can improve with online experience in a data constrained regime *very fast*, with aggressive updates!

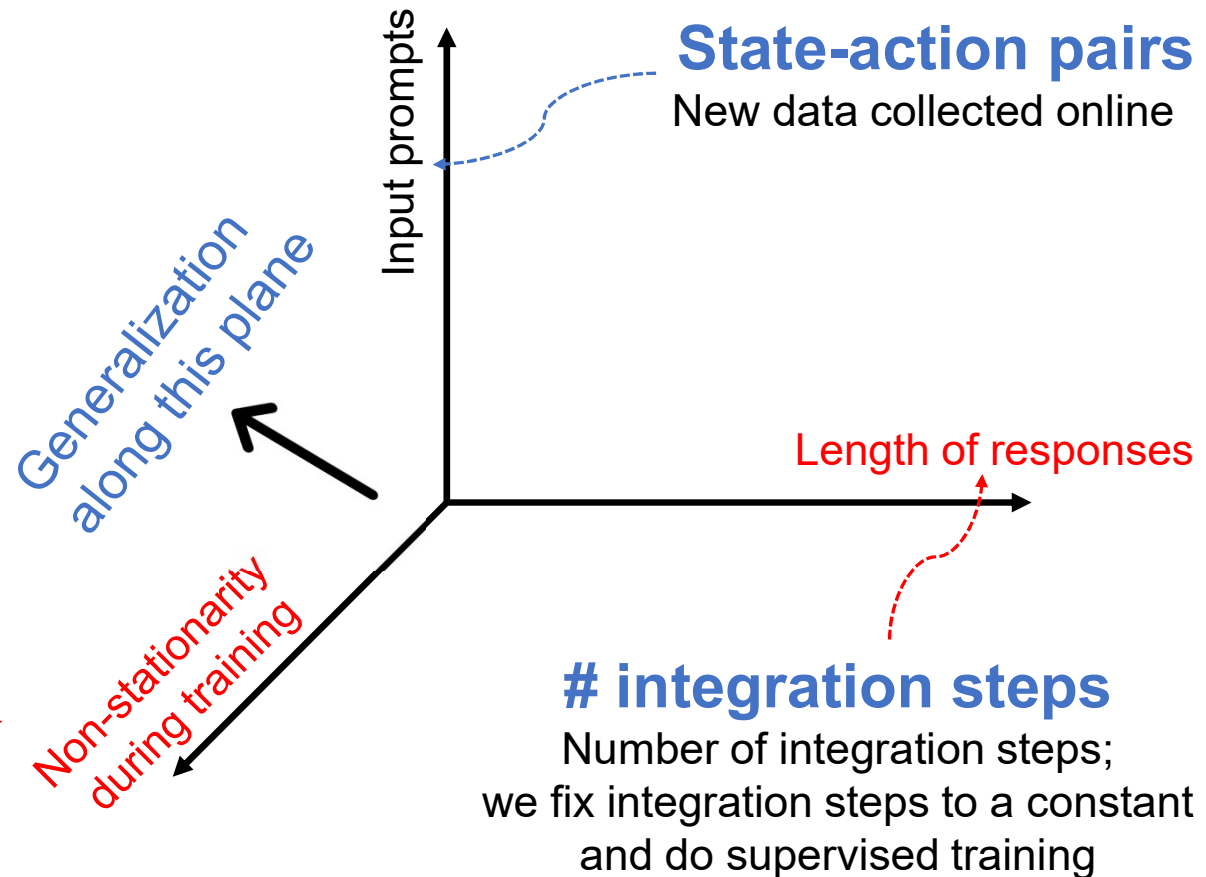
# Axes of Generalization in **f1oq** Networks

For **f1oq** nets

- New (s, a) pairs (yes!)
- ~~Longer lengths (no!)~~
- Training time (yes!)

## Non-stationarity

Moving targets and new data collected online in the environment



# Summary & Takeaways

- Training with a flow-matching loss presents a useful way to utilize **iterative test-time compute** for modeling non-stationary objects.
- Connections to test-time compute and ***plastic features*** from LLM land.
- **But still some open questions:**
  - What does exploration in a CoT correspond to for flows?
  - RL vs SFT gap in LLMs – what does it mean here?

Thank You!