



EarnHFT: Efficient Hierarchical Reinforcement Learning for High Frequency Trading

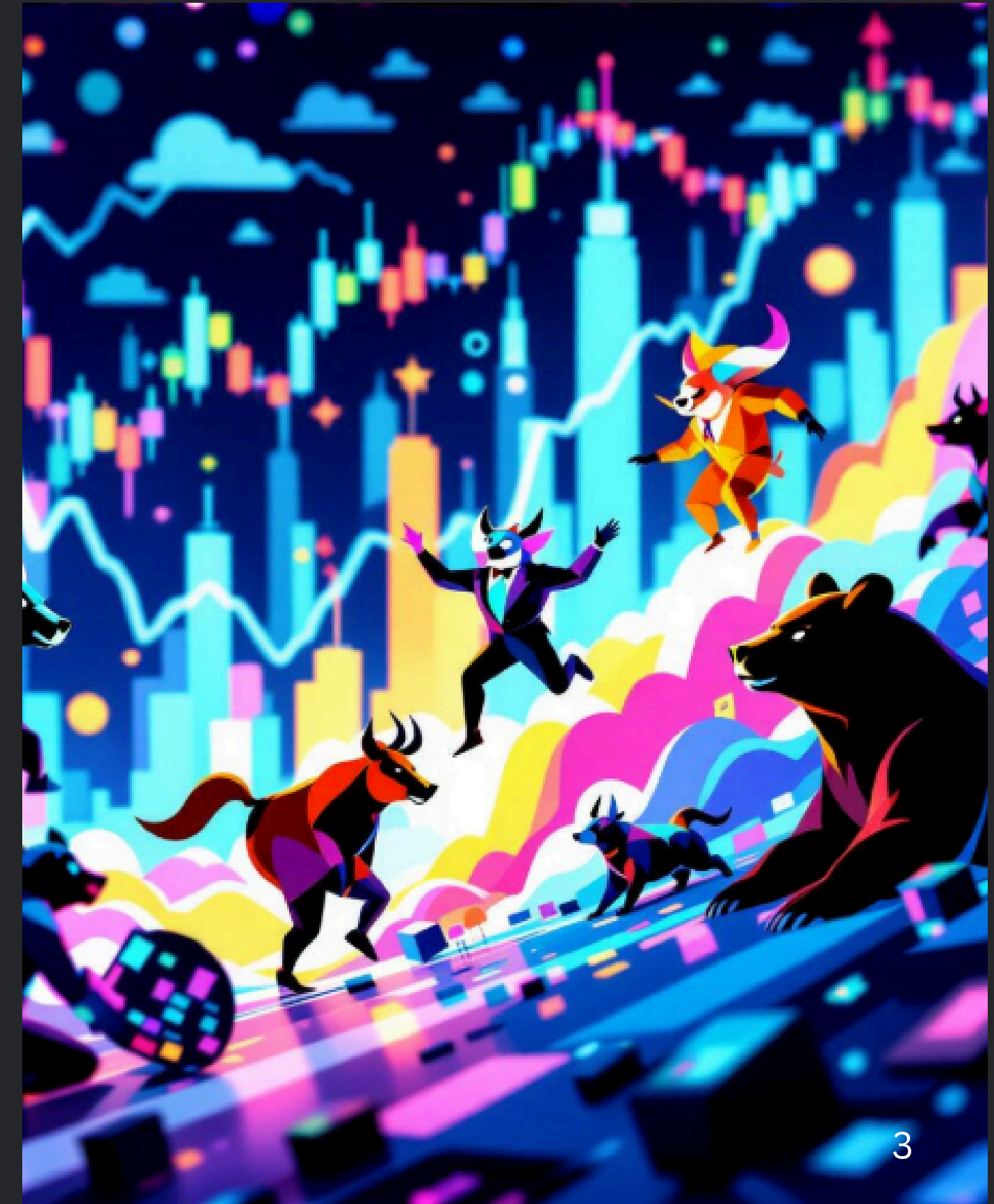
Shayan Kebriti

June 2025

Introduction

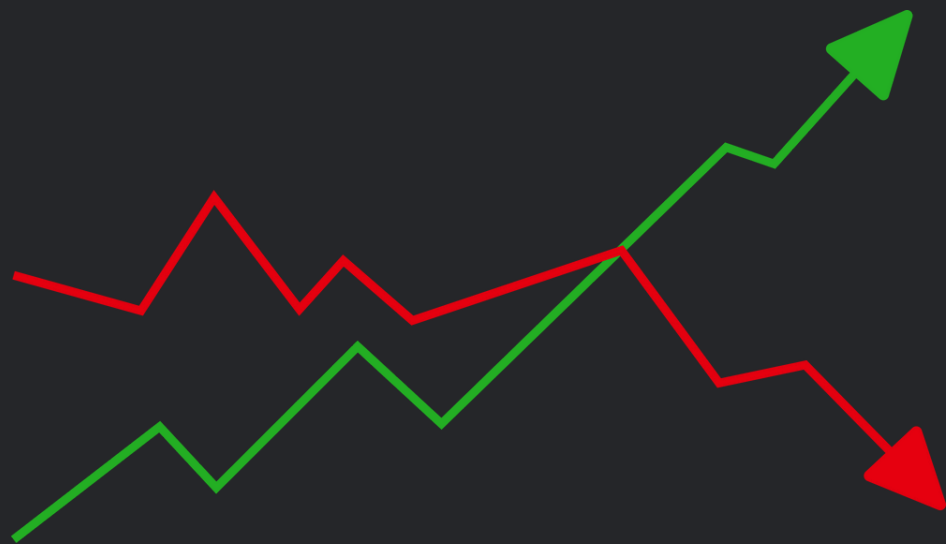
High Frequency Trading (HFT)

- Uses algorithms to make trading decisions in short time scales (e.g., second-level)
- Widely used in Cryptocurrency (e.g. Bitcoin)
 - High potential for profitability due to high volatility.
 - Lower Risk
 - Run 24/7, No overnight risk
- Helps the discovery of the true price of an asset



Challenges

- An extremely large time horizon
 - Atari games → ~ 6,000 steps
 - HFT → ~ 1 million steps (agents need to be evaluated in dozens of days)
 - Large time horizons need more data to converge
- Dramatic market changes
 - Agents trained on historical data fail in maintaining performance over long periods



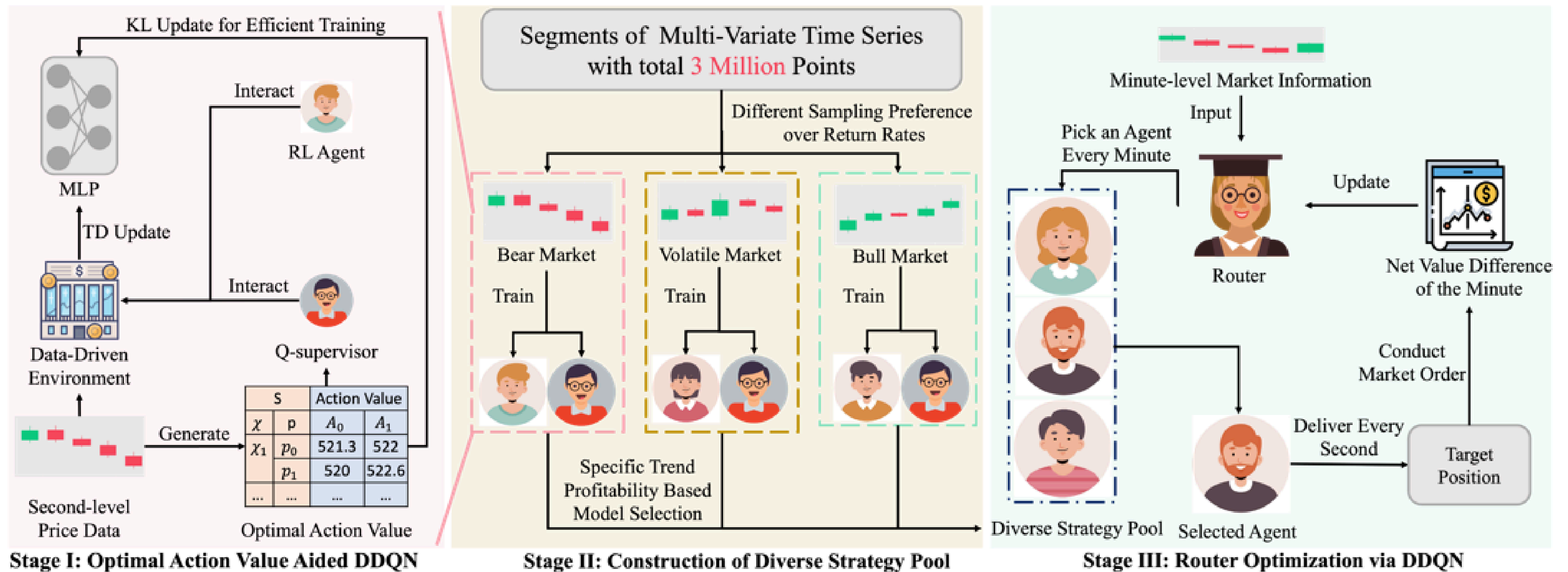
Proposed Method

Hierarchical MDP Formulation

- Low-level MDP (second-level operation)
 - Action: Choose a target position
 - State: Second-level market features
 - Reward: Net value change from the last second
- High-level MDP (minute-level operation):
 - Action: Choose a low-level policy (agent)
 - State: Minute-level market characteristics
 - Reward: Total profit from the chosen agent over the past minute
- Goal
 - Learn a set of diverse low-level policies for different market trends
 - Train a high-level policy to dynamically select the right agent based on current market conditions

Overview of EarnHFT

- Three Stages for Training
- Two-level Inference



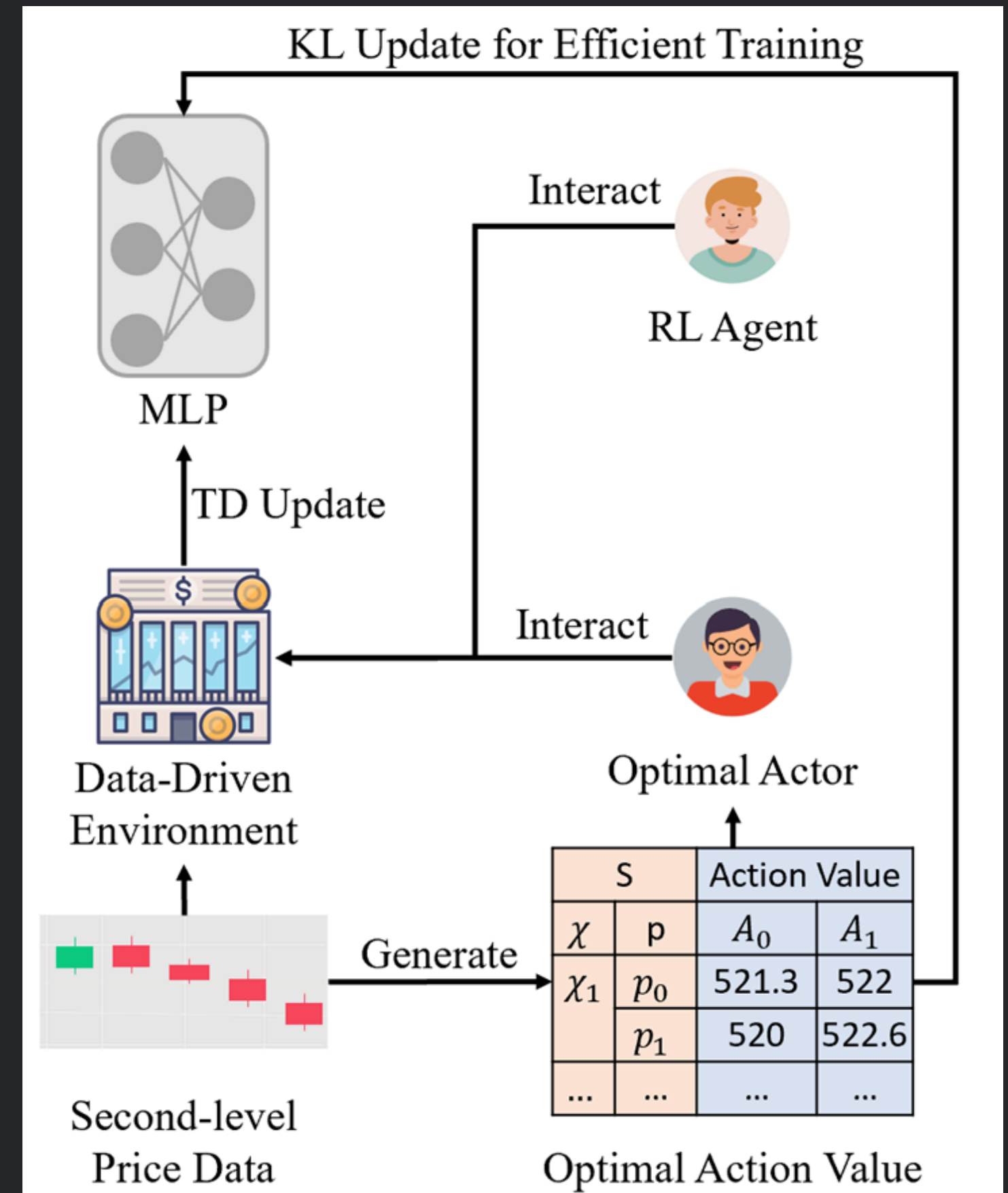
Stage 1: Efficient RL with Q-Teacher

- Computing Q^* using Dynamic Programming
- Using DDQN for Q
- Overall loss:

$$L(\theta_i) = L_{td} + \alpha KL(Q_t(\chi, p, \cdot; \theta_i) || Q^*(\chi, p, \cdot))$$

$$L_{td} = (r + \gamma \max Q_t(\chi', a, \cdot; \theta'_i) - Q_t(\chi, p, a; \theta_i))^2$$

- The *Optimal Actor* helps overcome the drift due to exploration of the *RL Agent*



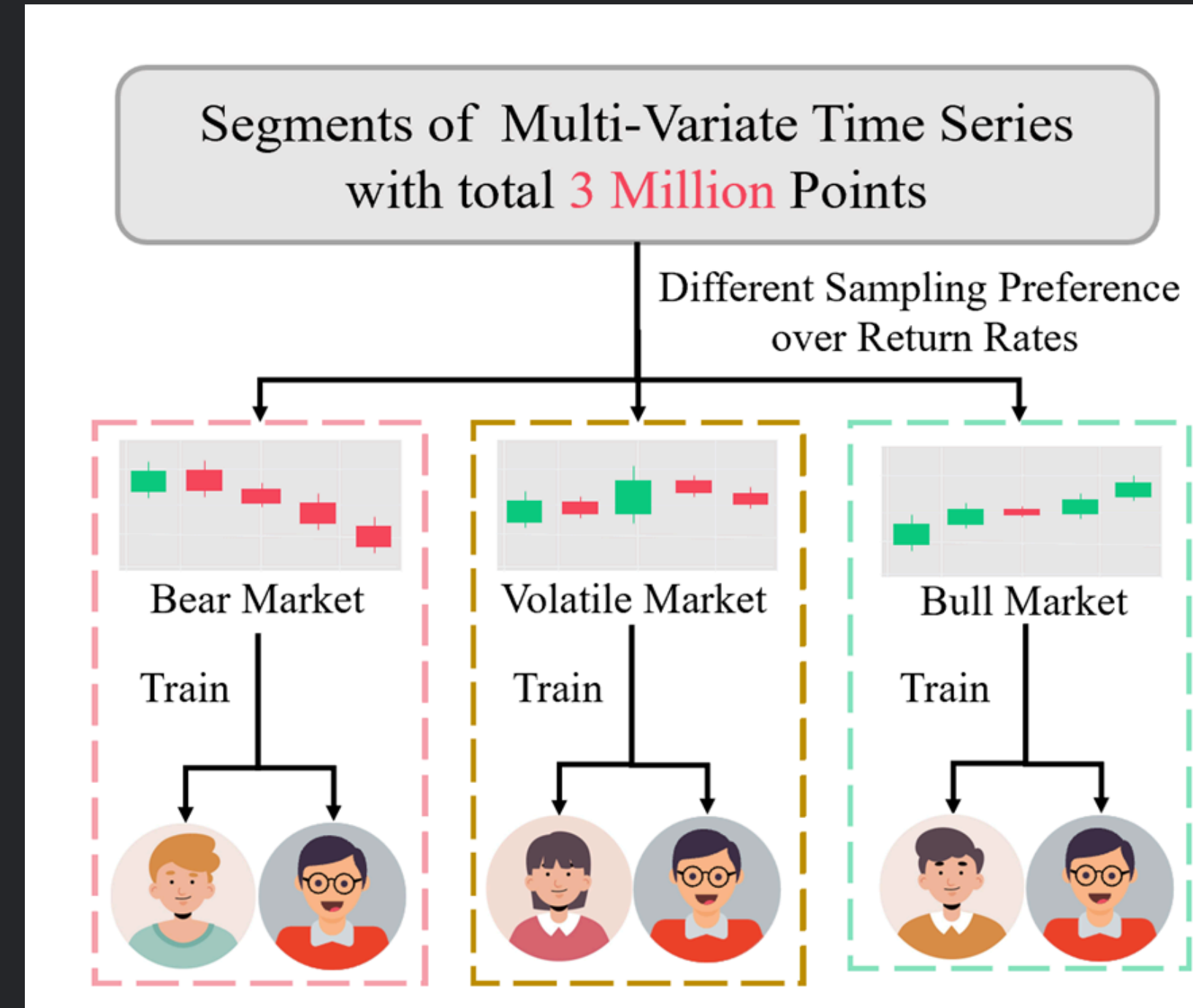
Stage 2: Building Diverse Agents

1. Seprate dataset into chunks

2. Train agents with different market trends by giving different priorities to each chunk:

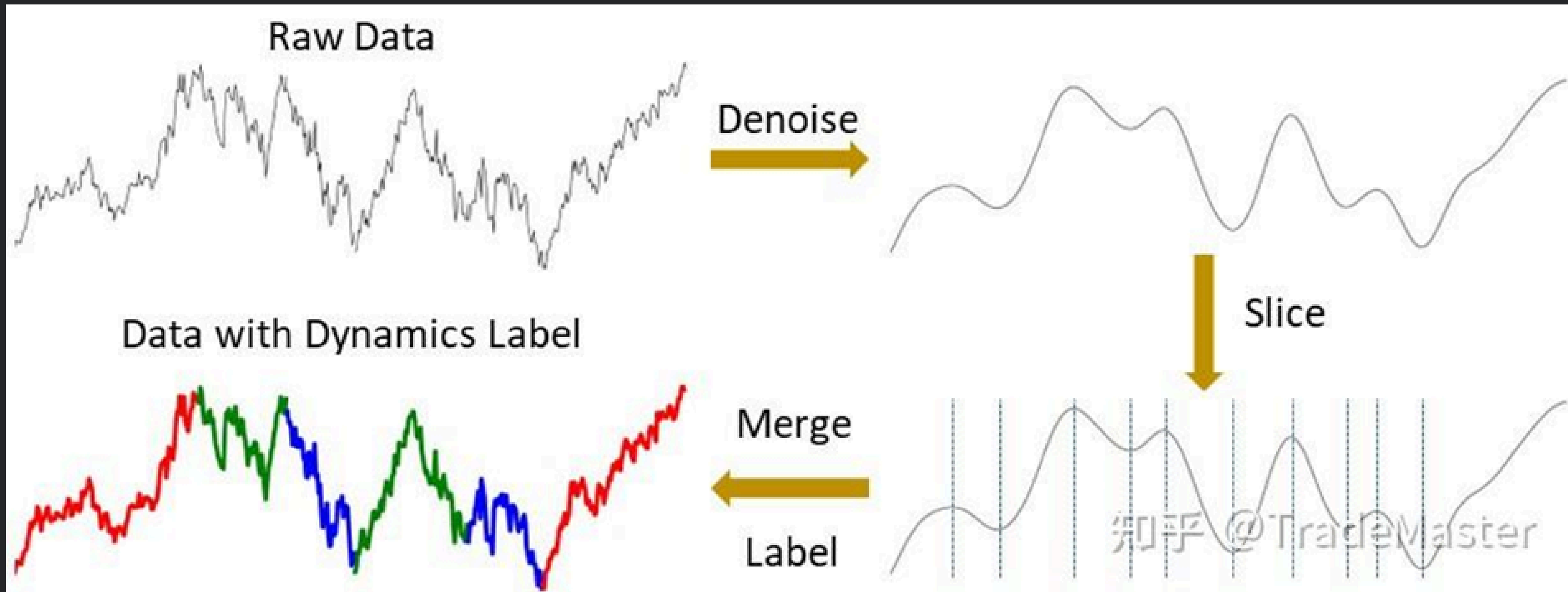
$$f(x) = \begin{cases} \frac{e^{\beta r}}{pdf(r)} & \text{if } Q_{\frac{\theta}{2}}(R) \leq r \leq Q_{1-\frac{\theta}{2}}(R) \\ e^{\beta r} & \text{if } r \geq Q_{1-\frac{\theta}{2}}(R) \vee r \leq Q_{\frac{\theta}{2}}(R) \end{cases}$$

- $r \rightarrow$ return
- $\beta \rightarrow$ Preference parameter (high/low return)
- $pdf(r) \rightarrow$ Probability density of return r
- $Q_{\theta}(R) \rightarrow$ Quantile function
- $\theta \rightarrow$ Risk threshold



Stage 2 (Cont.): Construction of Agent Pool

1. Market Segmentation & Labelling



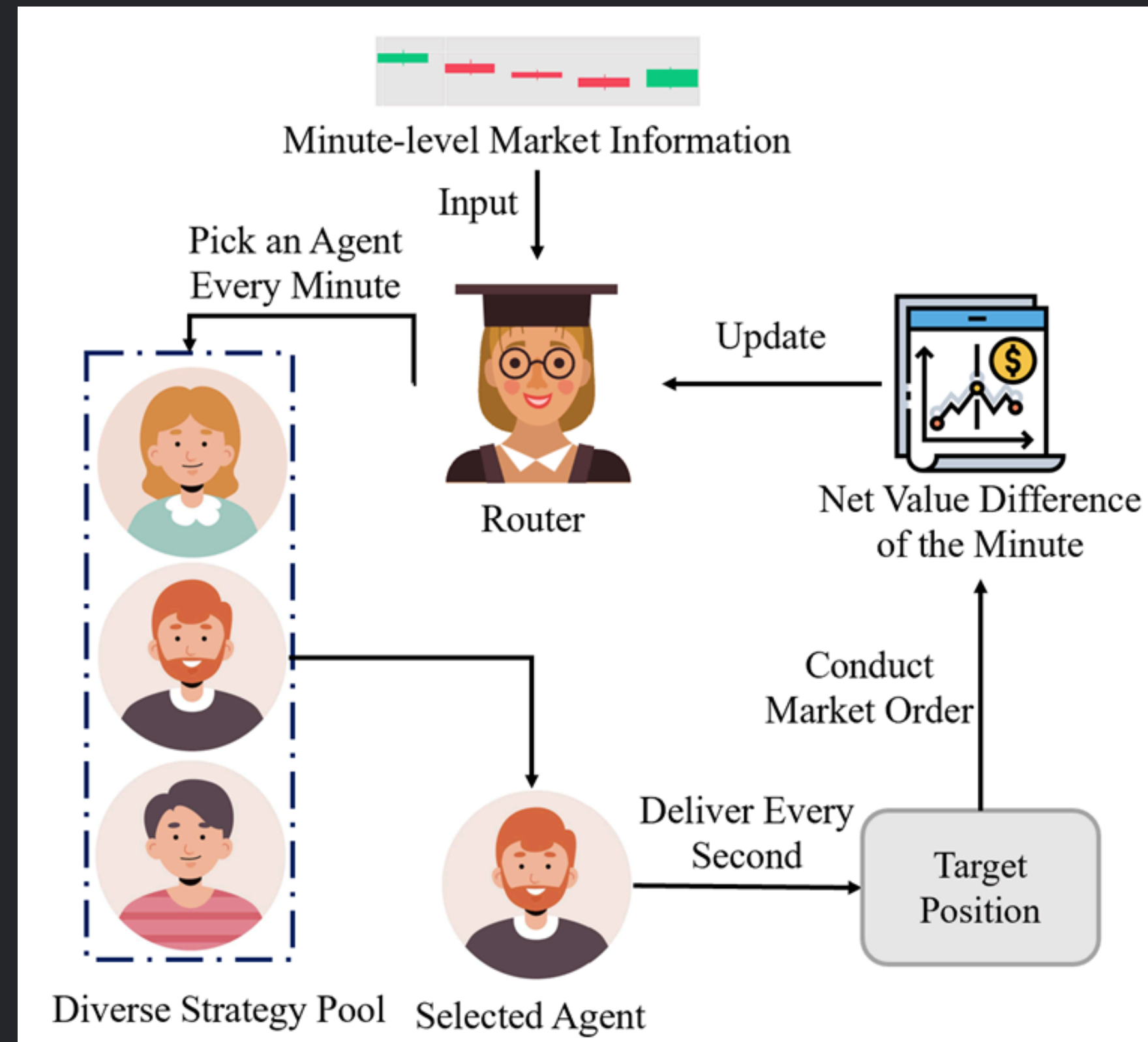
2. Evaluate the trained agents and choose the top ones to construct an agent pool of (m, n) .

$m \rightarrow$ number of market trends

$n \rightarrow$ initial positions

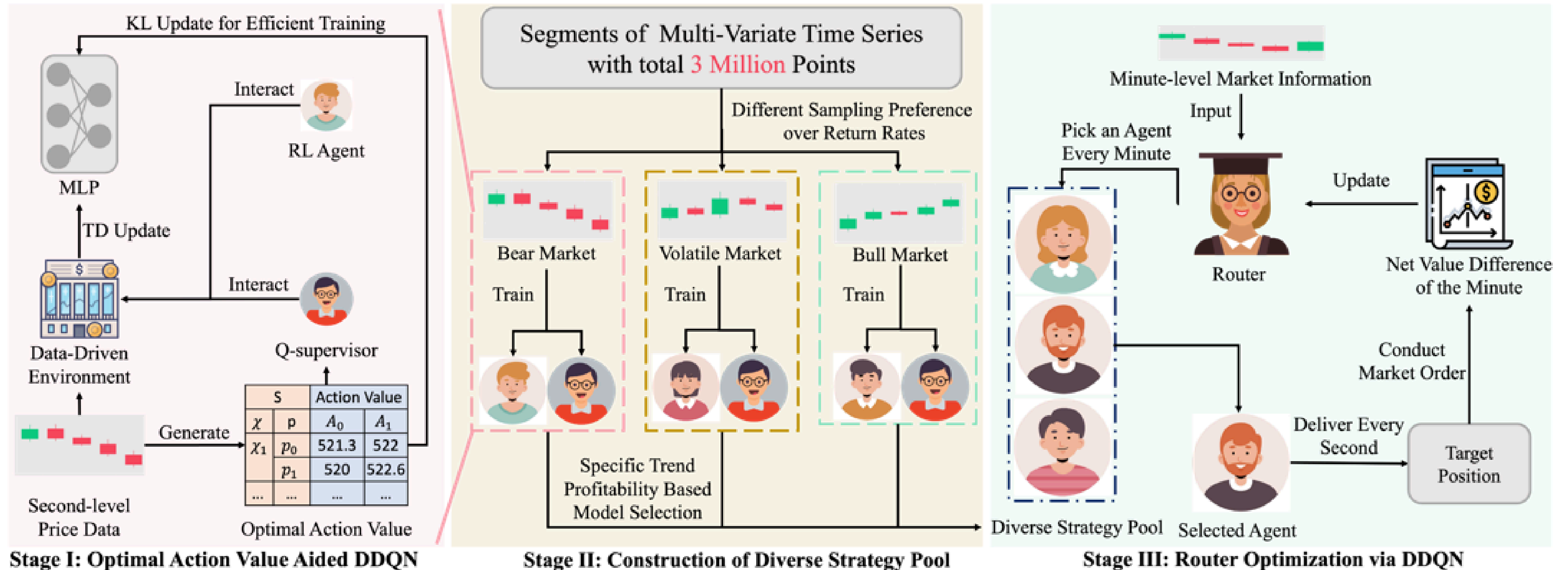
Stage 3: Dynamic Routing Optimization

- Train a DDQN for *Router*
- We reduce the number of possible low-level agents to m
 - using the agents with the same initial position



Summary

Summary



Results

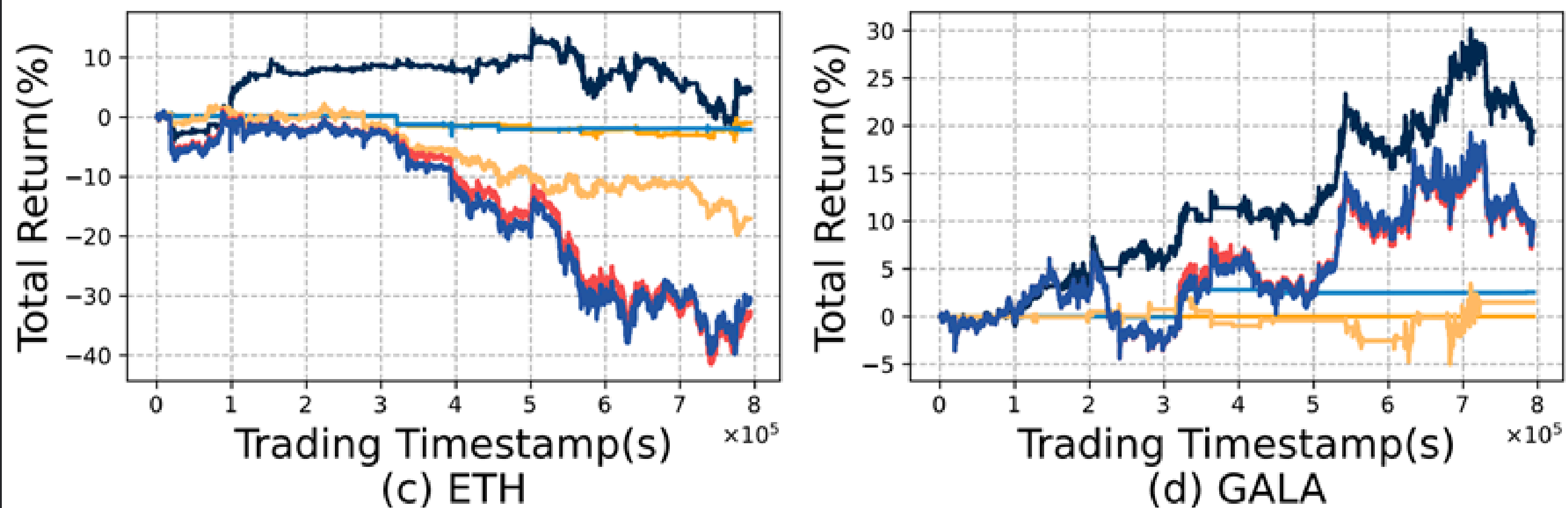
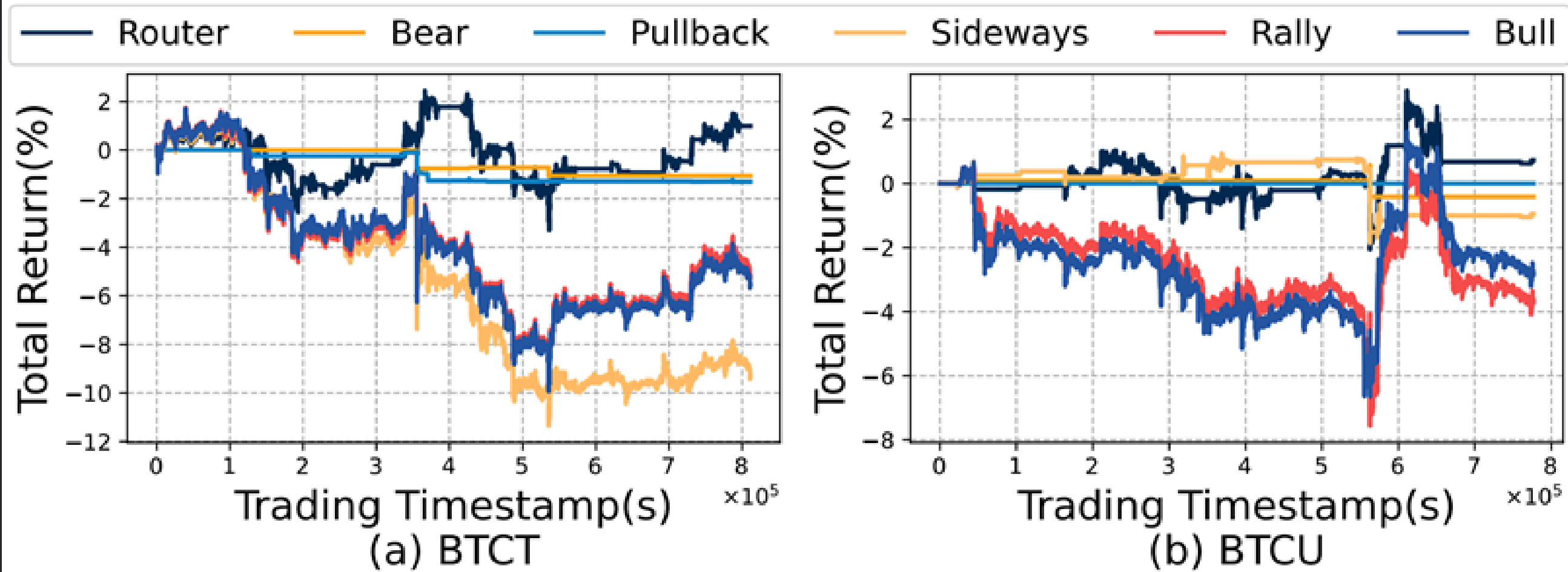
Experiments Results

		Prof↑	RAP↑	Risk↓			Prof↑	RAP↑	Risk↓
Market	Model	TR(%)	SR	MDD(%)	Market	Model	TR(%)	SR	MDD(%)
BTCU	DRA	-4.56	-4.28	9.24	BTCT	DRA	-2.65	<u>-4.82</u>	5.84
	PPO	-3.61	-5.25	<u>6.41</u>		PPO	-0.60	-14.74	<i>0.65</i>
	CDQNRP	-2.83	<i>-2.91</i>	7.38		CDQNRP	<u>-0.60</u>	-19.52	0.61
	DQN	-3.48	-12.37	<i>4.09</i>		DQN	<i>0.47</i>	4.21	<u>0.66</u>
	MACD	-6.07	-10.11	9.98		MACD	-4.02	-5.80	6.44
	IV	<u>-2.99</u>	<u>-3.78</u>	8.32		IV	-12.01	-17.83	12.66
	EarnHFT	0.72	1.22	3.07		EarnHFT	0.99	<i>1.34</i>	5.61
ETH	DRA	-33.37	<u>-9.06</u>	45.88	GALA	DRA	10.56	4.77	10.60
	PPO	-22.61	-10.11	31.17		PPO	<u>10.56</u>	<u>4.77</u>	10.60
	CDQNRP	<u>-6.82</u>	-24.41	6.96		CDQNRP	5.22	4.51	<i>5.41</i>
	DQN	-11.02	-9.47	<i>13.79</i>		DQN	2.94	3.55	3.78
	MACD	-4.29	-1.78	16.35		MACD	2.37	1.79	9.84
	IV	-27.42	-12.27	33.96		IV	<i>13.95</i>	<i>6.74</i>	9.91
	EarnHFT	4.52	2.92	<u>13.89</u>		EarnHFT	19.41	9.77	<u>9.26</u>

Experiments Results

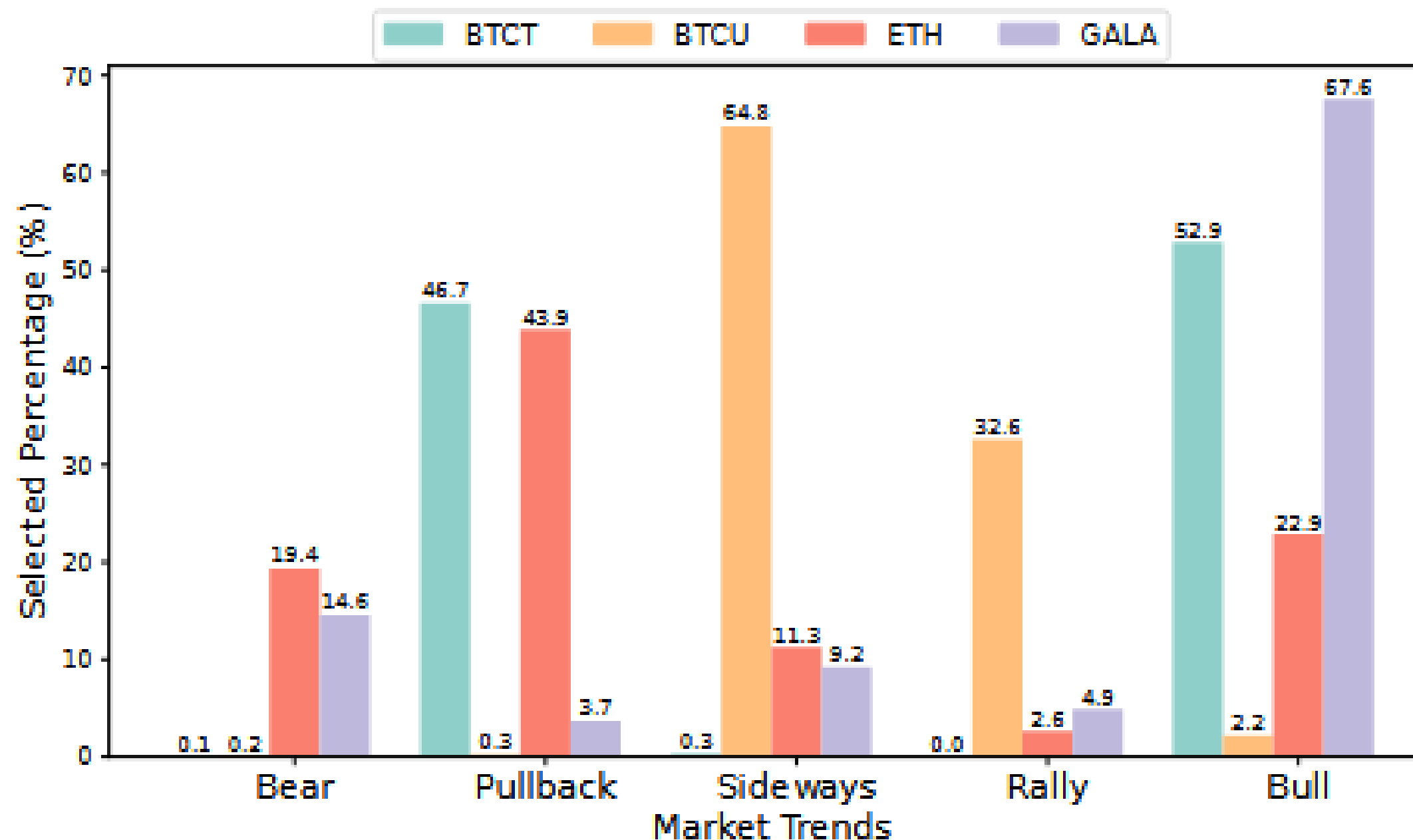


Experiments Results



Experiments Results

Dataset	Dynamics	Seconds	From	To
BTC/TUSD	Sideways	4057140	23/03/30	23/05/15
BTC/USDT	Sideways	3884400	22/09/01	22/10/15
ETH/USDT	Bear	3970800	22/05/01	22/06/15
GALA/USDT	Bull	3970740	22/07/01	22/08/15



Thanks!

Appendix

Algorithm 1: Construction of Optimal Action Value

Input: Multivariate Time Series \mathcal{D} with Length N , Commission Fee Rate δ , Action Space A

Output: A Table Q^* Indicating Optimal Action Value at Time t , Position p and Action a .

- 1: Initialize Q^* with shape $(N, |A|, |A|)$ and all elements 0.
 - 2: **for** $t \leftarrow N - 1$ to 1 **do**
 - 3: **for** $p \leftarrow 1$ to $|A|$ **do**
 - 4: **for** $a \leftarrow 1$ to $|A|$ **do**
 - 5: $Q^*[t, p, a] \leftarrow \max_{a'} Q^*[t+1, a, a'] + a \times p_{t+1}^{b1} -$
 $(p \times p_t^{b1} + E_t(p - a)).$
 - 6: **end for**
 - 7: **end for**
 - 8: **end for**
 - 9: **return** Q^*
-

Algorithm 2: Efficient RL with Q-Teacher

Input: Multivariate Time Series \mathcal{D} with Length N , Commission Fee Rate δ , Action Space A

Output: Network Parameter θ

- 1: Initialize experience replay R , network Q_θ , target network $Q_{\theta'}$ and construct the optimal action value using Algorithm 1 and trading environment Env .
 - 2: Initialize trading environment Env
 - 3: **for** $t = 1$ to $N - 1$ **do**
 - 4: Choose action a_ϵ using ϵ -greedy policy.
 - 5: Store transition $(s, a_\epsilon, r, s', Q^*)$ in D
 - 6: **end for**
 - 7: Reinitialize trading environment Env
 - 8: **for** $t = 1$ to $N - 1$ **do**
 - 9: Choose action a_o that $\text{argmax}_a Q^*[t, p, a]$.
 - 10: Store transition (s, a_o, r, s', Q^*) in R
 - 11: **end for**
 - 12: Sample transitions $(s_j, a_j, r_j, s'_j, Q_j^*)$
 - 13: Calculate L following Equation 2, do its gradient descent on θ and update $\theta' = \tau\theta + (1 - \tau)\theta'$.
 - 14: **return** Q_θ
-

Algorithm 3: Market Segmentation & Labelling

Input: A Time Series \mathcal{D} with Length N

Parameter: Risk threshold θ , Label number M

Output: Labels indicating the trend they belong to for every point in time series D

- 1: $D' \leftarrow$ denoising high frequency noise D .
 - 2: Divide D' according to its extrema into segments S .
 - 3: Merge adjacent segments in S if DTW (Muda, Begam, and Elamvazuthi 2010) and slop difference are small enough until S is stable.
 - 4: Calculate threshold $H = Q_{1-\frac{\theta}{2}}(R)$, $L = Q_{\frac{\theta}{2}}(R)$
 - 5: Calculate the upper bond and lower bonds of slopes for each label based on the quantile and the threshold.
 - 6: Label each segment based on the bonds.
 - 7: Return the label corresponding to each segment.
-

The Effectiveness of optimal value supervisor (OS) & optimal actor (OA)

OA	OS	GALA			ETH		
		CS	RS	AHL	CS	RS	AHL
✓	✓	78848	4.43	448	102400	12.32	81.3
✓		102400	0.24	38.7	102400	-1.40	4.15
	✓	4608	2.89	147	30720	4.87	35.8
		30720	-0.01	284	30720	-29.6	39.1

CS → Convergence Steps

RS → Converged Reward Sum

AHL → Average Holding Length