

Select and Trade: Towards Unified Pair Trading with Hierarchical Reinforcement Learning

Weiguang Han

Boyi Zhang

Min Peng

han.wei.guang@whu.edu.cn

zhangby@whu.edu.cn

pengm@whu.edu.cn

Wuhan University

Wuhan, Hubei, China

Yanzhao Lai

laiyanzhao@swjtu.edu.cn

Southwest Jiaotong University

Chengdu, Sichuan, China

Qianqian Xie

qianqian.xie@manchester.ac.uk

University of Manchester

Manchester, United Kingdom

Jimin Huang*

jimin@chancefocus.com

Chancefocus AMC.

Shanghai, China

ABSTRACT

Pair trading is one of the most effective statistical arbitrage strategies which seeks a neutral profit by hedging a pair of selected assets. Existing methods generally decompose the task into two separate steps: pair selection and trading. However, the decoupling of two closely related sub-tasks can block information propagation and lead to limited overall performance. For pair selection, ignoring the trading performance results in the wrong assets being selected with irrelevant price movements, while the agent trained for trading can overfit to the selected assets without any historical information of other assets. To address it, in this paper, we propose a paradigm for automatic pair trading as a unified task rather than a two-step pipeline. We design a hierarchical reinforcement learning framework to jointly learn and optimize two sub-tasks. A high-level policy would select two assets from all possible combinations and a low-level policy would then perform a series of trading actions. Experimental results on real-world stock data demonstrate the effectiveness of our method on pair trading compared with both existing pair selection and trading methods.

CCS CONCEPTS

• **Applied computing** → **Forecasting**; • **Computing methodologies** → **Partially-observable Markov decision processes**.

KEYWORDS

pair trading, hierarchical reinforcement learning, pair selection, automatic trading

ACM Reference Format:

Weiguang Han, Boyi Zhang, Min Peng, Qianqian Xie, Yanzhao Lai, and Jimin Huang. 2018. Select and Trade: Towards Unified Pair Trading with Hierarchical Reinforcement Learning. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 14 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

1 INTRODUCTION

Since 1987, pair trading, a basic statistical arbitrage approach, has been intensively practised and studied [28]. It is an integral component of the financial market and plays a crucial role in enhancing market efficiency [46]. On worldwide markets and varied asset types, such as stocks, futures, and cryptocurrency, it has been argued that pair trading is effective in the long term [24]. In contrast to portfolio selection to find the optimal portfolio with the highest “risky profit” [20], it seeks “riskless profit” by performing arbitrage tradings on the abnormal price movements of two correlated assets [28]. It first picks two correlated assets and monitors the spread between their respective prices. If the spread widens abnormally, it will execute trading operations on two assets and gain a profit when the spread recovers to its usual value. For instance, if Google’s price is usually \$2 higher than Facebook’s and it suddenly rises to \$5 higher, the strategy will short Google (expect its price to fall) and long Facebook (expect its price to increase) and close two tradings when the spread returns to \$2. The total return of the strategy is the sum of the returns from the two transactions on Google and Facebook, which relies solely on the spread between the two assets. Therefore, it is irrelevant to the vast majority of common risks, such as market fluctuations, since the profit resulting from the risk on Google would compensate for the loss resulting from the risk on Facebook. Nonetheless, it depends on two crucial factors: (1) the chosen two assets should be beneficial for pair trading, with a spread that exhibits significant mean-reversion and high volatility; and (2) a flexible agent that can identify abnormal increases and falls of spread from normal fluctuations.

Generally, previous methods for pair trading divided the process into two discrete stages: **pair selection** and **trading**. For pair

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).
Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXXX.XXXXXXX>

selection, they generally employ predefined statistical tests or fundamental distance measurements to select two assets based on their historical price [5, 7, 8, 10, 11, 15, 16, 19, 30, 38, 39, 41, 42, 50]. For example, a number of previous researches apply the cointegration test [50] to determine if the historical price spread between two assets is stable. After selecting the pair, they would engage in trading using fixed-threshold-based strategies to generate the return in a subsequent period. Recently, inspired by the successful deployment of reinforcement learning in other areas [1, 14, 21, 32], there have been efforts to introduce reinforcement learning to train a flexible agent and report a significant improvement over traditional methods [13, 22, 51].

However, existing methods of automated pair trading are still confronted with drawbacks. Despite the fact that they can ensure the relevance between the selected two assets and perform tradings by decoupling pair selection and trading, it prevents the flow of information between them, which can be a significant concern since they are tightly coupled. For pair selection, existing methods would choose the wrong asset pairs since the employed model-free metrics are **target irrelevant**, which means, they consider no performance of candidate asset pairs during the following trading period. For example, the optimal asset pair with the lowest Euclidean distance would have zero spread and trading opportunities. It is fundamental to dynamically learn the measurement of the future profitability of asset pairs from the data. As for trading, existing methods can be **target overfitting** due to only observing the pre-selected asset pair during the training and ignoring other asset pairs and the market. Although reinforcement learning allows their methods to learn a flexible agent which can explore different trading actions during the training, the learned agent could show poor performance in the trading period with unseen market data since only partial historical information is leveraged.

Despite the critical necessity to jointly simulate the two phases of pair trading, there have been no prior attempts in this area. In this research, we propose a novel paradigm for automated pair trading, in which the two-step process is formulated as a unified task rather than a pipeline with two independently sorted sub-tasks. The approach must simultaneously choose the trading pair from candidate pairs in a formation period and trade it in a later trading period in order to optimize trading performance. Although the paradigm is straightforward for the task, it poses two challenges to the development of successful approaches. First, it is challenging to represent the sequential process of the paradigm in which trading occurs after pair selection, i.e., selecting two correlated assets and then trading on their anomalies. Second, there are complicated relationships between pair selection and trading that must be fully utilized to generate risk-free profits. Pair selection intuitively sets the input of trading, while trading offers the output of pair selection in the form of profit.

To address these issues, we design a new framework **TRIALS** that adopts feudal hierarchical reinforcement learning (FHRL) [37] to jointly learn and optimize two steps: a high-level reinforcement learning policy as the manager for pair selection, and a low-level reinforcement learning policy as the worker for trading. The agent in our proposed framework would first select an asset pair from all possible combinations of assets, and then perform a series of trading actions based on the selected pair. Given a set of assets,

for the high-level manager, **states** are the historical price features of these assets in the formation period; **options** are all possible combinations of these assets; **rewards** are the overall performance which is generated from the low-level worker on the trading period. As For the low-level worker, given the chosen option as two selected assets, i.e. Google and Facebook, **states** are the historical price features of two assets and the trading information of the agent such as historical actions, cash, and present net value; **actions** are three discrete trading actions including *long* (Buy Google to sell it later and sell Facebook to buy it back later), *short* (Sell Google and buy Facebook), and *clear* (Sale previous bought assets and buy previously sold assets to close tradings); and **rewards** are the overall performance of the agent in the formation period. Notice that the rewards for the high-level manager and the low-level worker are devised on the trading and formation period respectively, although they are both generated via the same low-level worker. This allows our method to optimize the agent at two levels jointly, which guides the high-level manager to select optimal asset pairs according to their trading performance on the unseen market data, and forces the low-level worker to consider different asset pairs and capture the common profitable patterns for pair trading. We further verify the effectiveness of our method on U.S. and Chinese stock datasets compared with both previous pair selection and trading methods. The experimental results prove that our proposed method can achieve the best performance, which attributes to the correct pair selection and corresponding precise trading actions.

In summary, our contributions can be listed as:

- (1) We are the first to introduce a new task for pair trading that combines the existing two tasks as pair selection and trading. In order to optimize the total trading performance, it is necessary for the approach to simultaneously consider these two steps, which were previously overlooked in both pair selection and trading.
- (2) We design a novel end-to-end hierarchical framework that introduces feudal hierarchical reinforcement learning to jointly optimize a high-level policy for pair selection and a low-level policy for trading.
- (3) Experimental results on both U.S. and Chinese stock markets demonstrate the effectiveness of our method compared with existing pair selection and trading methods.

2 RELATED WORK

2.1 Traditional Pair Selection

For pair selection, previous methods aim to find two assets whose prices have moved together historically in a formation period, and their future spread is assumed to be historical mean-reverted [24]. They generally adopted statistical or fundamental similarity measurements based on historical price information to perform asset pair selection before trading. The distance approach was first introduced [8, 10, 16, 19, 38, 39] for pair selection, which simply adopted distance metrics such as the sum of Euclidean squared distance (SSD) for the price time series to model the connection between two assets. However, an ideal asset pair in these model-free methods were expected to be two assets with exactly the same price movement in historical time, which have zero trading opportunities for no fluctuations of price spread. There were also

methods [5, 7, 12, 15, 30, 41, 42, 50] that directly model the tradability of a candidate pair based on the Engle-Granger cointegration test, which performs linear regression using the price series of two assets and expects the residual to be stationary.

However, the mean-reversion properties of the spread of an asset pair in the future can be irrelevant to their mean-reversion strength in history, which limits the trading performance of the selected pair from these parameter-free methods. Although there were also methods [25] that integrated neural networks to learn the metrics, they proposed to measure the profit of assets rather than the asset pairs and selected the top and bottom assets to form the trading pair, which is difficult to find two matched assets. [54] was the most similar study which considers pair trading as a unified portfolio management task. Their methods using historical price spread as the metric nevertheless suffer from the same issue, even though they can dynamically learn the trading and allocation ratios of each pair.

2.2 Reinforcement Learning for Pair Trading

After pair selection, previous methods generate trading signals which trigger contradictory actions on two assets during the trading period. Based on the assumption that the spread of the selected pair would still revert to its historical mean value, previous methods generally employ simple threshold-based rules that they would long the undervalued and short the overvalued asset when the spread is higher or lower than the historical mean by pre-defined thresholds [24]. However, it requires expert knowledge to identify the optimal trading thresholds in the time-varying market.

Inspired by the success of applying reinforcement learning (RL) in financial trading problems [14], previous attempts generally focused on introducing RL methods to develop flexible trading agents after pair selection via traditional methods. [13] used the cointegration method to select trading pairs, and adopted Q-Learning [52] to select optimal trading parameters. Kim and Kim introduced a deep Q-network [35] to select the best trading threshold for cointegration approaches [23]. [31] proposed to detect structural changes and improve reinforcement learning trading methods. Brim, Wang et al. directly utilized the RL methods to train an agent for trading [6, 51]. [22] further introduced stop-loss boundaries to control the risk. Although these methods have shown the benefits of the integration of RLs as a smart trading agent, they still adopt traditional methods for pair selection which only consider the historical performance of the trading pair. Moreover, their trading agent can easily overfit to the only observable asset pair and show limited performance on the unseen future market. However, there were no previous efforts to address the problem, which requires the method to jointly learn how to select and trade asset pairs.

2.3 Hierarchical Reinforcement Learning

Many approaches have been proposed for building agents within the context of hierarchical reinforcement learning (HRL) [40, 43, 53]. The feudal framework is one popular approach for HRL, in which the action space of a higher-level policy consists of sub-goals corresponding to various sub-tasks and the objective of this lower-level policy is to achieve the input sub-goal [9]. In HRL, different levels of temporal abstraction enable efficient credit assignment

over longer timescales [49]. At the same time, a subtask may itself be easier to learn and the learned sub-tasks lead to more structured exploration over the course of training of the HRL agent [36]. In previous works, the low-level policy generally learned handcrafted sub-goals [26], discovered options [3] or intrinsic rewards [49], while the high-level policy is learned using extrinsic rewards from the environment. The decomposition of feudal HRL can also help to model complex tasks that are difficult for normal RL methods.

As one of the most challenging applications, pair trading consisting of two separate steps requires the method to optimize two related but different sub-tasks. Existing methods generally deem the process as a two-step pipeline and apply different methods for each step respectively. It inevitably blocks the information propagation between these two steps and introduces extra noise due to the error accumulation step-by-step. To the best of our knowledge, our work is the first one that applies HRL in pair trading to end-to-end learning and inference.

3 HIERARCHICAL PAIR TRADING FRAMEWORK

In this section, we illustrate the detail of our proposed hierarchical pair trading framework, as shown in Fig. 1.

3.1 Formalization

Generally, pair trading consists of two steps: pair selection and trading. In pair selection, it would select two correlated assets from all possible combinations of assets to form a trading pair. Given the trading pair, it would perform a series of trading actions to earn market-neutral profit in a subsequent period. The task aims to maximize the trading profit of the selected asset pair, which requires selecting the optimal trading pair and choosing correct trading actions during the trading period. Different from previous approaches that generally take two steps separately, in this paper, we propose to jointly learn to select and trade the pair in a unified hierarchical framework. Therefore, given a formation period with $T_{\mathcal{F}}$ time points consisting of $\{0, 1, \dots, T_{\mathcal{F}} - 1\}$, a subsequent trading period with $T_{\mathcal{T}}$ time points consisting of $\{0, 1, \dots, T_{\mathcal{T}} - 1\}$, and selected N assets $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$, there are formation price series $\{p_0^x, p_1^x, \dots, p_{T_{\mathcal{F}}-1}^x\}$ and trading price series $\{p_0^x, p_1^x, \dots, p_{T_{\mathcal{T}}-1}^x\}$ for each asset $x \in \mathcal{X}$ that is associated with each time point in formation period and trading period respectively.

Formally, we formulate the pair trading process as the feudal hierarchical reinforcement learning framework [37]. As shown in Fig.2, a feudal hierarchical reinforcement learning framework consists of two controllers: a high-level controller called *manager* and a low-level controller as *worker*. The manager is designed to set the option which aims to maximize the *extrinsic* reward or the goal of the task. By selecting an option, the manager would trigger the worker which is guided by the *intrinsic* reward. Different from the extrinsic reward as the overall target of the task, the intrinsic reward is a sub-goal of the manager given the selected option. Therefore, the decomposition allows the method to satisfy requirements at multiple levels to solve complex tasks that are infeasible for centralized reinforcement learning.

To this end, we design a high-level controller as the manager for pair trading, which aims to select two assets as a pair and

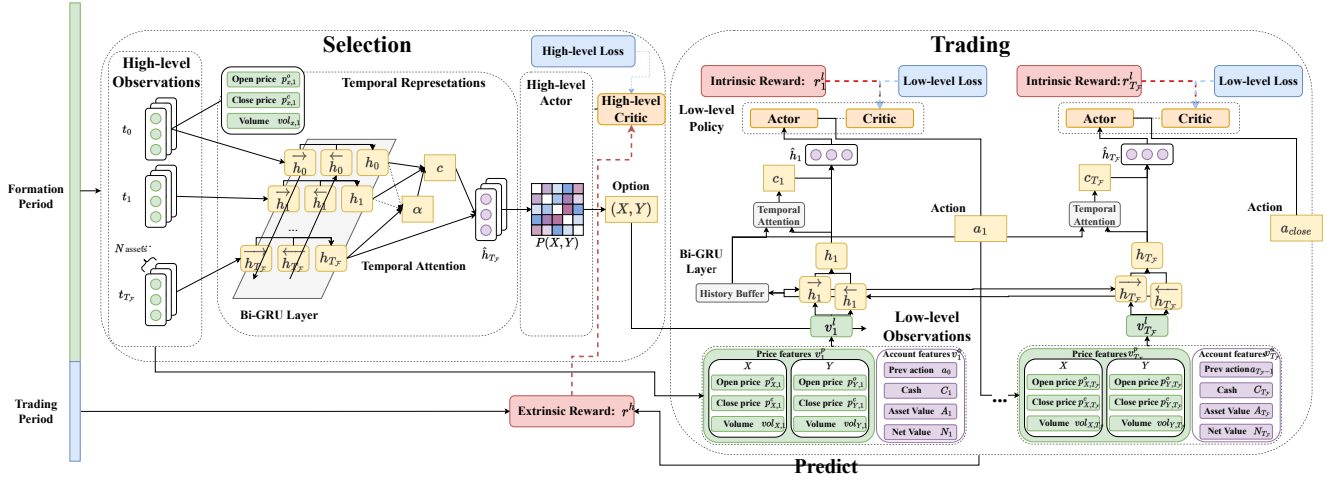


Figure 1: The hierarchical framework for pair trading.

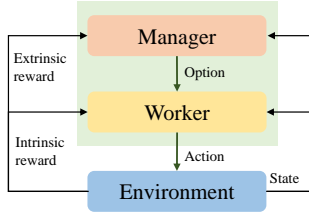


Figure 2: The feudal hierarchical reinforcement learning framework.

maximize their trading performance via pair trading. It is expected that the pair should possess the highest profit in the subsequent trading period among all possible combinations of assets. Thus the extrinsic reward for the manager is the profit of the selected pair in the trading period.

However, to achieve the optimal trading profit of the selected pair, it is required to consider a different sub-task where the agent is supposed to perform a series of sequential trading decisions on the selected pair. Since the target is different from the selection, we derive a low-level controller as the worker which only focuses on learning a flexible and profitable trading policy. We adopt the profit of the selected pair in the formation period to guide the learning of the worker as the intrinsic reward. After the worker is fully trained with the historical formation data, it is utilized to yield the trading performance of the selected pair with the unseen market data in the trading period, which is further taken as the extrinsic reward.

3.2 Pair Selection with High-Level Controller

For pair selection, we aim to select the optimal asset pair from all possible pairs of assets. It can be deemed as a contextual bandit [27] $M = (S^h, O, T^h, R^h, \Psi, Q^h)$ over options, where S^h refers to the state space, O is the option space, T^h is the transitions among states, R^h is the designed reward, Ψ is the observation state which is generated from the current state $s^h \in S^h$ and the option referring to the

high-level action $o \in O$ according to the probability distribution $Q^h(s^h, o)$. Different from trading, the pair selection process is a one-step decision process that the agent would perform an option $o_0 \in O$ under the current state s_0^h , resulting in the transition from s_0^h to s_1^h with the probability $T(s_1^h | s_0^h, o_1)$. After the option is selected, a low-level POMDP as the worker would be triggered to perform trading according to the selected option.

3.2.1 Observation. For the agent in the high-level contextual bandit, only limited information of the market state can be observed, i.e., the price features of the assets in history, which means the agent can only receive the observation $v_0 \in \Psi$ with probability as $Q(s_1^h, o_0)$. The observation $v_0 \in \Psi$ is the price features for all assets $x \in \mathcal{X}$ associated with each time step $t \in T_f$ in the formation period, including the open price $p_{x,t}^o$, the close price $p_{x,t}^c$, and the volume $vol_{x,t}$.

3.2.2 Option. The option o is a pair (x_i, x_j) selected from all possible combinations of assets in \mathcal{X} . When the low-level POMDP ended, the agent would select the next option according to the high-level contextual bandit.

3.2.3 State. Given the observation v_0^h consisting of the open price $p_{x,t}^o$, the close price $p_{x,t}^c$, and the volume at each time point $vol_{x,t}$ for each asset $x \in \mathcal{X}$ in the formation period $t \in T_f$, we adopt the Bi-directional GRU (Bi-GRU) [18] to capture the temporal correlations between historical price features. Our method takes the previous hidden state h_{t-1} as the hidden state of the forward GRU and the next state h_t as the hidden state of the backward GRU. Since the asset prices possess strong auto-correlation effects [33], it is fundamental to model the relationships from both history and future, which helps the method to capture salient information embedded in the asset price fluctuations. Therefore, we represent our latent state h_t as:

$$\vec{h}_t = \text{GRU}(v_{0,t}^h, \vec{h}_{t-1}), \overleftarrow{h}_t = \text{GRU}(v_{0,t}^h, \overleftarrow{h}_{t+1}), h_t = [\vec{h}_t, \overleftarrow{h}_t] \quad (1)$$

where $h_t \in \mathbb{R}^{d_h}$ is the concatenation of the forward hidden state and backward hidden state, d_h is the hidden dimension, and $v_{0,t}^h \in \mathbb{R}^{N \times 3}$

are the price features of all assets at the time step $t \in T_{\mathcal{F}}$ of the formation period.

As a matter of fact, Bi-GRU has the long-distance forgetting problem [4], especially when there are thousands of time steps in the formation period. We further introduce a temporal attention mechanism to dynamically select salient information from all historical time steps by:

$$\alpha_k = \frac{\exp(\text{score}(h_{T_{\mathcal{F}}}, h_k))}{\sum_{k'=0}^{T_{\mathcal{F}}-1} \exp(\text{score}(h_{T_{\mathcal{F}}}, h_{k'}))}, c_{T_{\mathcal{F}}} = \sum_k \alpha_k h_k \quad (2)$$

$$\hat{h}_{T_{\mathcal{F}}} = \text{LayerNorm}(\text{LeakyRelu}(W_c[h_{T_{\mathcal{F}}}, c_{T_{\mathcal{F}}}])) \quad (3)$$

where $\text{score}(h_{T_{\mathcal{F}}}, h_k) = \frac{h_{T_{\mathcal{F}}} h_k}{\sqrt{d_h}}$ is the scaled dot-product attention score [48]. We also adopt LeakyRelu and LayerNorm [2] to stabilize the hidden state dynamics. We adopt the final output $\hat{h}_{T_{\mathcal{F}}} \in \mathbb{R}^{N \times d_h}$ as the state $s_0^h \in \mathbb{R}^{N \times d_h}$ of the high-level contextual bandit.

3.2.4 Policy. The stochastic policy for pair selection $\mu : \mathcal{S} \rightarrow \mathcal{O}$ refers to a probability distribution over options:

$$o_0 \sim \mu(o_0 | s_0^h) = \text{softmax}(\text{triu}(s_0^h s_0^{hT})) \quad (4)$$

where *triu* is to extract and return the flattened upper triangular part of the given matrix.

3.2.5 Reward. The reward of the high-level contextual bandit is the same as the target of the task, which is to maximize the profit of the trading period given the option o_0 . However, realizing the optimal trading profit requires the method to learn a different sub-task. Therefore, we propose to utilize a low-level POMDP triggered by the selected option. It is first trained with the intrinsic reward in the formation period and then utilized to perform tradings to yield the trading profit in the trading period. Following previous RL-based trading methods, we also maximize the cumulative profit over the trading period with $T_{\mathcal{T}}$ time points:

$$R^h = \prod_{t \in T_{\mathcal{T}}} (1 + R_t^h) \quad (5)$$

where R_t^h is the return of the low-level policy. We would provide further details in the following subsections.

3.3 Trading with Low-Level Controller

When the high-level controller has selected a trading pair as the option, the low-level controller will perform trading based on the given trading pair as a series of trading actions in a subsequent trading period to achieve the trading profit. Formally, we formulate the decision process of the trading as a Partially Observable Markov Decision Process (POMDP) [17] $M = (S^l, A, T^l, R^l, \Omega, Q^l)$, where S^l refers to the state space, A is the action space, T^l is the transitions among states, R^l is the designed reward, Ω is the partial observation state which is generated from the current state $s^l \in S^l$ and action $a \in A$ according to the probability distribution $Q^l(s^l, a)$. At each time point, the agent would perform an action $a_t \in A$ under the current state s_t^l , resulting in the transition from s_t^l to s_{t+1}^l with the probability $T^l(s_{t+1}^l | s_t^l, a_t)$. Similar to pair selection, the actual market states are partially observed and only the historical prices and volumes of assets, along with the historical account information

of the agent such as actions, amounts of cash, and returns can be leveraged, while other information is ignored. In detail, the agent can only receive the observation $v_{t+1}^l \in \Omega$ with probability as $Q(s_{t+1}^l, a_t)$, which requires the agent to fully exploit the historical observations up to present time point.

3.3.1 Observation. The observation $v_t^l \in \Omega$ consists of two different feature sets, including: (1) the account features $v_t^a \in \Omega^a$ as previous action a_{t-1} , present cash C_t , present asset value V_t , and cumulative profit as the net value N_t ; (2) the price features $v_t^p \in \Omega^p$ as the open price $p_{i,t}^o$, the close price $p_{i,t}^c$, and the volume $\text{vol}_{i,t}$ of for each asset $i \in \{X, Y\}$. Following previous work [29], we simplify the impact of tradings performed by our agent on the market state as a constant loss to each trading. Therefore the action of our agent would not affect the state and price features of assets in our observation.

3.3.2 Action. The action in each time step is to perform a pair of contradictory trading actions on two assets respectively. The action space $A = \{L, C, S\} = \{1, 0, -1\}$ consists of three discrete actions each of which involves two trading actions for two assets $\{X, Y\}$ respectively. In detail, the *L* action represents the long trading action which means to long asset X and short asset Y at the same time, the *C* action for clear referring to clear two assets if longed or shorted any before, and the *short* action for short which is to short asset X and long asset Y . Notice that for different asset pairs, the trading action at each time step could be assigned to other actions.

3.3.3 State. Different from the policy for pair selection, the agent is required to estimate the latent market state s_t^l according to the history $H_t = \{v_1^l, a_1, v_2^l, \dots, a_{t-1}, v_t^l\}$. Although the market state cannot be directly observed, the historical information embedded in H_t , especially the sequential dependencies can help the agent to generate better estimation.

Therefore, we also introduce Bi-GRU to encode the history. The previous hidden state h_{t-1} is deemed as the hidden state of the forward GRU and the next state h_t as the hidden state of the backward GRU:

$$\vec{h}_t = \text{GRU}(v_{0,t}^l, \vec{h}_{t-1}), \overleftarrow{h}_t = \text{GRU}(v_{0,t}^l, \overleftarrow{h}_{t+1}), h_t = [\vec{h}_t, \overleftarrow{h}_t] \quad (6)$$

where $h_t \in \mathbb{R}^{d_h}$ is also the concatenation of the forward hidden state and backward hidden state, d_h is the hidden dimension, and $v_{0,t}^h \in \mathbb{R}^{2 \times M}$ are the input features of the selected two assets and M is the feature dimension. We transform discrete variables such as previous action a_{t-1} in v_t^a into continuous embeddings via an embedding layer $E_a \in \mathbb{R}^{3 \times d_a}$, where d_a is the corresponding embedding size.

We also introduce a temporal attention mechanism:

$$\alpha_k = \frac{\exp(\text{score}(h_t, h_k))}{\sum_{k'=0}^{t-1} \exp(\text{score}(h_t, h_{k'}))}, c_t = \sum_k \alpha_k h_k \quad (7)$$

$$\hat{h}_t = \text{LayerNorm}(\text{LeakyRelu}(W_c[h_t, c_t])) \quad (8)$$

where $\text{score}(h_t, h_k) = \frac{h_t h_k}{\sqrt{d_h}}$ is the scaled dot-product attention score. The output $\hat{h}_t \in \mathbb{R}^{d_h}$ as the $s_t^l \in \mathbb{R}^{d_h}$ of the low-level POMDP.

3.3.4 *Policy.* The stochastic policy for trading $\pi : \mathcal{S} \rightarrow \mathcal{A}$ yields a probability distribution over actions given the low-level state s_t^l and the high-level option o_0 :

$$a_t \sim \pi(a_t | s_t^l; o_0) = \text{softmax}(W_\pi s_t^l) \quad (9)$$

3.3.5 *Reward.* The intrinsic reward for the low-level controller is also the cumulative profit over a period with T time points:

$$R_t = \prod_{t \in T} (1 + R_t) \quad (10)$$

where R_t is the return of the agent given the action a_t :

$$\begin{aligned} R_t &= a_{t-1} r_{X,t} - a_{t-1} r_{Y,t} - c |a_t - a_{t-1}| \\ &= a_{t-1} (r_{X,t} - r_{Y,t}) - c |a_t - a_{t-1}| \end{aligned} \quad (11)$$

Notice that the return of the agent is irrelevant to the market for hedging the return of two assets as $r_{X,t} - r_{Y,t}$. To yield a positive return, it is required to select the optimal trading pair and precisely trading actions according to the historical performance of the trading pair. For training, we use the formation period to guide the learning of the worker, and the trading period to yield a high-level extrinsic reward with the fully trained worker.

3.4 Hierarchical Policy Learning

For high-level policy updating, following the Advantage Actor-Critic method (A2C) [34], We update the policy and the value function every step as:

$$\begin{aligned} \nabla_{\theta_h^P} \log \mu(o_0 | s_0^h; \theta_h^P) A(s_0^h; \theta_h^A) \\ \nabla_{\theta_h^A} \frac{1}{2} A^2(s_0^h; \theta_h^A) \end{aligned} \quad (12)$$

where $A(s_0^h; \theta_h^A) = r_1^h + \gamma V(s_0^h; \theta_h^A) - V(s_0^h; \theta_h^A)$ is the estimation of the advantage function for the high-level controller and the option o_0 is sampled from the option distribution $\mu(o_0 | s_0^h; \theta_h^P)$.

As for low-level policy updating, we apply A2C update similarly,

$$\begin{aligned} \nabla_{\theta_l^P} \log \pi(a_t | s_t^l; o_0, \theta_l^P) A(s_t^l; \theta_l^A) \\ \nabla_{\theta_l^A} \frac{1}{2} A^2(s_t^l; \theta_l^A) \end{aligned} \quad (13)$$

where $A(s_t^l; \theta_l^A) = r_{t+1}^l + \gamma V(s_{t+1}^l; \theta_l^A) - V(s_t^l; \theta_l^A)$ is the estimation of the advantage function for the low-level controller and action a_t is sampled from the action distribution $\pi(a_t | s_t^l; o_0, \theta_l^P)$.

For training, we adopt the same formation period data as the input to train both high-level and low-level policy, where the performance of the low-level policy during the trading period would be considered as the reward of the high-level policy. As for evaluation and testing, we directly infer the option and corresponding actions without exploration.

Algorithm 1 Training

Require: N assets \mathcal{X} , loop conditions M, N

Ensure: Model parameters $\theta_h = \{\theta_h^P, \theta_h^A\}, \theta_l = \{\theta_l^P, \theta_l^A\}$

- 1: Initialize parameters θ_h, θ_l for the high-level controller and low-level controller respectively
 - 2: **for** iteration=1, 2, 3, ..., M **do**
 - 3: Sample option o_0 from $\mu(o_0 | s_0^h)$
 - 4: Select pair from \mathcal{X} and initialize the trading environment
 - 5: **for** iteration=1, 2, 3, ..., N **do**
 - 6: **while** not reach termination condition **do**
 - 7: Sample action a_t from $\pi(a_t | s_t^l; o_0)$
 - 8: Execute action, then obtain the next state and intrinsic reward from the trading environment
 - 9: Update θ_l by Eq (13)
 - 10: **end while**
 - 11: **end for**
 - 12: Obtain extrinsic reward from pair selection environment
 - 13: Update θ_h by Eq (12)
 - 14: **end for**
-

4 EXPERIMENTS

4.1 Dataset

Following previous methods [24], we build a dataset based on a pool of real stocks from S&P 500¹ for recent 21 years from 01/02/2000 to 12/31/2020. We filter stocks that have missing data throughout the whole period, resulting in 150 stocks with 5,284 trading days. To support the evaluation and development of pair trading, we introduce a new daily emerging stock market dataset (Chinese CSI 300 dataset) including 300 stocks and 5,088 time steps from the CSMAR database². Similar to previous work [24], we construct our stock dataset using a pool of stocks from the CSI 300 index for the last 21 years, from 01/02/2000 to 12/31/2020. Instead of all stocks in the market, we select the stocks that used to belong to the major market index CSI 300 and filter out stocks that have missing price data over the period. We compare our dataset and the U.S. stock market dataset S&P 500 in Table 1.

Dataset	Market	Period	Assets	Time Steps
S&P 500	U.S	2000 - 2020	150	5284
CSI 300	China	2000 - 2020	300	5088

Table 1: The statistics of datasets.

For each trading day, we use the fundamental price features as the features of stocks, including open price, close price, and volume. Additionally, we normalize price features such as open price and close price with logarithm.

Different from previous methods, we randomly split stocks into five non-overlapping sub-datasets, as shown in Appendix A. For each subset with, we perform experiments of our method and baselines to evaluate their generality. We use the first 90% trading days

¹Tiingo. Tiingo stock market tools. <https://api.tiingo.com/documentation/iex>

²www.gtarsc.com

as train data, the following 5% as validation data, and the rest 5% as test data. For training, we further use the first 85% trading days to train our methods to simultaneously select the optimal trading pair from possible combinations and perform optimal trading actions based on the optimal trading pair in the rest of 5% trading days. The trained model is evaluated on the validation data to select the best hyperparameters based on which the performance of the model among the test data is reported. We independently evaluate and report the performance of all methods on each subset, along with the mean and standard deviation over all subsets.

For our method and ablations, we use the RMSProp optimizer [47] and perform a bayesian parameter search [44] for each subset to set the optimal hyper-parameters respectively. We implement our method based on Pytorch and stable-baselines, and conduct all our experiments on a server with 2 NVIDIA Tesla V100 GPUs.

4.2 Baselines

We compare our methods with the following baselines: (1) **Pair selection methods**: they mainly focus on selecting the optimal asset pair which is expected to yield the best performance with threshold-based trading rules, such as **GGR** [16] which uses average Euclidean distance to select pairs, **Cointegration** [50] which adopts the augmented Engle-Granger two-step cointegration test to select the trading pair, and **Correlation** [11] which selects two assets that have the highest correlation. (2) **Trading methods**: they generally aim to train an agent to perform optimal trading actions with the asset pair which is generally selected using the augmented Engle-Granger two-step cointegration test, i.e. **Wang et al.** [51] that adopting the reinforcement learning to maximize the overall profit.

4.3 Metrics

As previous trading methods [51], we first evaluate our method along with baselines with their trading performance on the test data using (1) **Sharpe ratio (SR)** is the ratio of the profit to the risk [45], which is calculated as $(E(R_t) - R_f)/V(R_t)$, where R_t is the daily return and R_f is a risk-free daily return that is set to 0.00085 as previous methods. (2) **Annualized return (AR)** is the expected profit of the agent when trading for a year. (3) **Maximum drawdown (MDD)** measures the risk as the maximum potential loss from a peak to a trough during the trading period. (4) **Annualized Volatility (AV)** measures the risk as the volatility of return over the course of a year.

We also employ fundamental measurements to measure the selected pair of all methods: average **Euclidean distance (ED)** [16] which is the average euclidean distance of the historical price series of two assets.

4.4 Main Results

As shown in Table 2, our method TRIALS achieves the best performance among all methods in all metrics and most stock subsets of both S&P 500 and CSI 300. The detailed performance of each subset in two datasets is presented in Appendix C. It demonstrates the effectiveness of our method for simultaneously learning the pair selection and trading with a unified hierarchical reinforcement learning framework. In detail, TRIALS has the highest average SR

and AR, which indicates that our trained trading agent can yield a remarkable profit with controlled risks based on the selected pair of our method. Our method presents a consistently high performance in two datasets, which clearly shows the two tasks in our unified framework are complementary since the pair selection task requires the trading performance of assets while the trading task depends on the selected pair. This is further proved when TRIALS also yields the lowest average MDD in S&P 500 and a relatively low MDD in CSI 300. Since AV indicates both the fluctuations during rises and falls, our method presents a relatively high average AV in both datasets.

In contrast, previous pair selection methods such as GGR, Cointegration, and Correlation, underperform our method. The average SR of GGR is -1.37 and -1.19, and its AR in both two datasets are negative. Cointegration has an even worse average SR of -1.83 and -1.50, and also negative AR in both datasets. Similarly, correlation shows an average SR as -1.41 and -1.37 along with a negative SR. It clearly shows that pair selection methods based on pre-defined statistical tests or fundamental measurements would fail to select the optimal trading pair without considering the trading performance of the asset pairs. Moreover, the statistical tests and fundamental measurements adopted in their methods cannot measure the profitability of the asset pair even with the test data. For example, our method has a higher ED compared with existing methods, despite the fact that our method yields a significant profit in both datasets.

As for trading methods such as Wang et al. which adopt reinforcement learning to train a flexible agent, it shows a better performance than pair trading methods such as Cointegration with the same selected pair. However, it is limited by the selected pair based on the cointegration test which is shown to be ineffective in capturing the profitability of the asset pair, resulting in a lower SR compared with our method.

4.5 Ablation Study

To evaluate the contributions of two tasks in our unified framework, we further propose an ablation of our proposed method to compare with our method, as shown in Table 2, which is **TRIALS wo TR** that adopts a fixed trading agent with predefined thresholds after our RL based pair selection.

Compared with TRIALS wo TR, our method which jointly optimizes the two tasks presents the best performance with the highest average SR and lowest average ED. In contrast, TRIALS wo TR is misguided by the trading performance of fixed trading rules which strongly relies on the wrong estimation of the mean and standard deviation of the price spread as the historical mean and standard deviation, resulting in a worse result in comparison to our method.

However, TRIALS wo TR still outperforms existing parameter-free methods, which proves the importance of dynamically learning the measurement of the future profitability according to their trading performance in pair selection.

We also display the visualizations of the learned pair selection probability of our method in Fig.3. It clearly shows that our method can precisely capture the complex connections between asset pairs. For example, EQR engages in the real estate investment and ABT engages in chemicals. Although there are no direct connections, our method finds that they are strongly and consistently correlated,

Model		GGR	Cointegration	Correlation	Wang	TRIALS	TRIALS wo TR
S&P 500	SR↑	-1.37 (0.79)	-1.83 (0.27)	-1.41 (0.21)	1.18 (0.43)	1.84 (0.24)	0.07 (0.16)
	AR↑	-0.15 (0.09)	-0.36 (0.20)	-0.14 (0.05)	0.21 (0.11)	0.50 (0.14)	0.01 (0.20)
	MDD↑	-0.20 (0.08)	-0.37 (0.20)	-0.20 (0.04)	-0.09 (0.05)	-0.09 (0.01)	-0.25 (0.07)
	AV↓	0.13 (0.03)	0.27 (0.20)	0.12 (0.02)	0.16 (0.06)	0.22 (0.04)	0.22 (0.04)
	ED↓	0.014 (5e-3)	0.021 (0.02)	0.007 (0.002)	0.021 (0.02)	0.037 (0.01)	0.01 (4e-3)
CSI 300	SR↑	-1.19 (0.74)	-1.50 (0.97)	-1.37 (0.25)	0.75 (0.68)	1.91 (0.88)	0.95 (0.88)
	AR↑	-0.17 (0.11)	-0.25 (0.17)	-0.21 (0.07)	0.24 (0.23)	0.68 (0.51)	0.13 (0.12)
	MDD↑	-0.29 (0.06)	-0.29 (0.13)	-0.25 (0.06)	-0.18 (0.09)	-0.14 (0.07)	-0.12 (0.09)
	AV↓	0.18 (0.03)	0.19 (0.03)	0.17 (0.05)	0.25 (0.07)	0.26 (0.09)	0.17 (0.07)
	ED↓	0.013 (6e-3)	0.017 (8e-3)	0.015 (8e-3)	0.017 (8e-3)	0.046 (0.02)	0.02 (8e-3)

Table 2: Mean(Standard Deviation) of all metrics on S&P 500 and CSI 300.

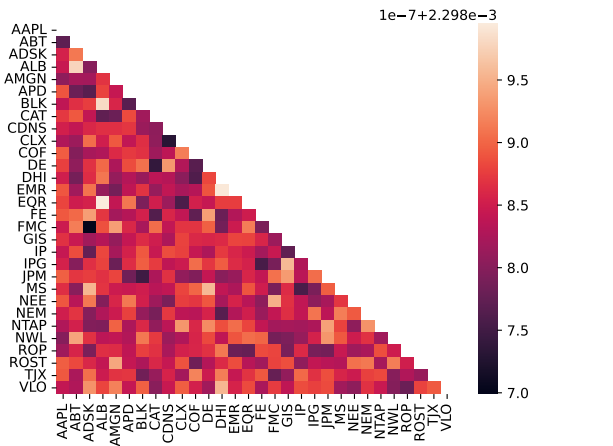


Figure 3: The pair selection probabilities of TRIALS.

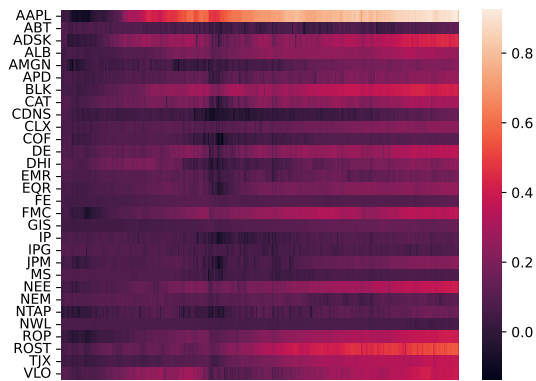
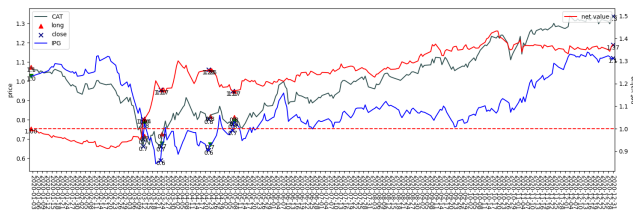


Figure 4: The visualization of temporal attentions

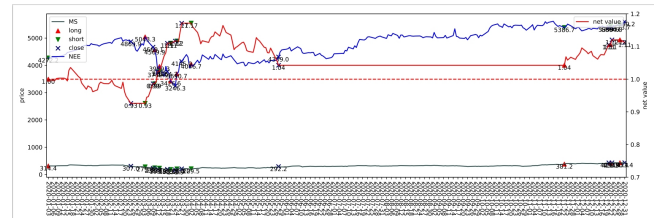
which indicates a more complex multi-hop relationship between these two stocks such as industry spillover. Besides, we display the temporal attention of our method in Fig.4. As shown in Fig.4, our method can fully exploit the temporal information by temporal attention, which means, for AAPL, we would focus more on the latest features.

4.6 Case Study

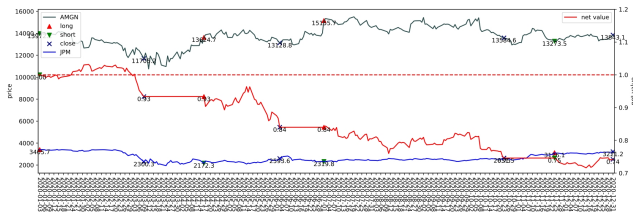
To further verify the profitability of the selected pair and learned trading agent of our method, we show the detailed trading actions, positions, and profit during the trading period of the selected pair by TRIALS, TRIALS wo TR, GGR, and Wang et al. in Set 2, as shown in Fig.5. The larger version of these figures is presented in Appendix B. Since GGR and Wang et al. both ignore the trading performance of assets, there are irrelevant movements of the prices of the selected pair such as NEE and MS, resulting in wrong tradings with great



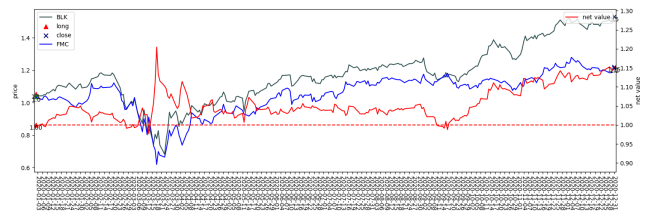
(a) The trading detail of TRIALS.



(b) The trading detail of TRIALS wo TR.



(c) The trading detail of GGR.



(d) The trading detail of Wang et al.

Figure 5: The trading details

loss. As for GGR and TRIALS w/o TR, they show irrational trading decisions due to the wrongly-estimated thresholds for trading, also leading to poor trading performance.

Compared with them, our method that jointly learns to select optimal pair and trade can simultaneously consider the information of all assets and dynamically learns the measurement according to the optimal trading performance based on a flexible agent. It allows our method to select the profitable pair CAT and IPG, which shows multiple trading opportunities in the trading period which are precisely captured by our trained trading method. Besides, our method can observe multiple asset pairs, which forces the worker to capture the consistent pattern for pair trading instead of overfitting to only one selected asset pair. Thus the worker in our method can precisely capture the trading opportunities and yield significant profit.

In contrast, although TRIALS w/o TR also learns to dynamically select asset pairs according to their trading performance, the fixed-threshold-based trading method can only provide biased information, resulting in less profitable pair selection.

5 CONCLUSION

In this paper, we proposed a novel paradigm for automatic pair trading that unifies the two sub-tasks: pair selection and trading. Based on it, we designed a feudal hierarchical reinforcement learning method consisting of a high-level manager for selection and a low-level worker for trading. The manager focused on selecting a pair as the option from all possible combinations of assets to maximize its trading performance, while the worker was to achieve the option set by the manager and yield the trading performance of the selected pair after training on historical data. Experimental results on the real-world stock data prove that the two steps in pair trading are closely related and complementary, which our method can fully exploit and jointly optimize to generate a significant improvement compared with existing pair selection and trading methods. In the future, we would further integrate more

representation methods for learning the representations of assets and consider other information such as natural language texts and macroeconomic variables.

REFERENCES

- [1] Saud Almahdi and Steve Y Yang. 2019. A constrained portfolio trading system using particle swarm algorithm and recurrent reinforcement learning. *Expert Systems with Applications* 130 (2019), 145–156.
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. Layer Normalization. <https://doi.org/10.48550/ARXIV.1607.06450>
- [3] Pierre-Luc Bacon, Jean Harb, and Doina Precup. 2017. The Option-Critic Architecture. *ArXiv abs/1609.05140* (2017).
- [4] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).
- [5] William K. Bertram. 2009. Analytic Solutions for Optimal Statistical Arbitrage Trading. *ERN: Optimization Techniques; Programming Models; Dynamic Analysis (Topic)* (2009).
- [6] Andrew Brim. 2020. Deep Reinforcement Learning Pairs Trading with a Double Deep Q-Network. *2020 10th Annual Computing and Communication Workshop and Conference (CCWC)* (2020), 0222–0227.
- [7] Cathy W. S. Chen, Max Chen, and Shu-Yu Chen. 2014. Pairs Trading via Three-Regime Threshold Autoregressive GARCH Models. In *TES*.
- [8] Huafeng (Jason) Chen, Shaojun Chen, Zhuo Chen, and Feng Li. 2019. Empirical Investigation of an Equity Pairs Trading Strategy. *American Finance Association Meetings (AFA)* (2019).
- [9] Peter Dayan and Geoffrey E Hinton. 1992. Feudal Reinforcement Learning. In *Advances in Neural Information Processing Systems*, S. Hanson, J. Cowan, and C. Giles (Eds.), Vol. 5. Morgan-Kaufmann. <https://proceedings.neurips.cc/paper/1992/file/d14220e66aee73c49038385428ec4c-Paper.pdf>
- [10] Binh Huu Do and Robert W. Faff. 2012. Are Pairs Trading Profits Robust to Trading Costs. *Journal of Financial Research* 35 (2012), 261–287.
- [11] Robert J Elliott, John Van Der Hoek*, and William P Malcolm. 2005. Pairs trading. *Quantitative Finance* 5, 3 (2005), 271–276.
- [12] Robert J R Elliott, John Van Der Hoek*, and W. Paul Malcolm. 2005. Pairs trading. *Quantitative Finance* 5 (2005), 271 – 276.
- [13] Saeid Fallahpour, Hasan Hakimian, Khalil Taheri, and Ehsan Ramezani. 2016. Pairs trading strategy optimization using the reinforcement learning method: a cointegration approach. *Soft Computing* 20, 12 (2016), 5051–5066.
- [14] Thomas G Fischer. 2018. *Reinforcement learning in financial markets—a survey*. Technical Report. FAU Discussion Papers in Economics.
- [15] Alexander Galenko, Elmira Popova, and Ivilina Popova. 2012. Trading in the Presence of Cointegration. *The Journal of Alternative Investments* 15 (2012), 85 – 97.
- [16] Evan Gatev, William N Goetzmann, and K Geert Rouwenhorst. 2006. Pairs trading: Performance of a relative-value arbitrage rule. *The Review of Financial Studies* 19, 3 (2006), 797–827.

- [17] Matthew Hausknecht and Peter Stone. 2015. Deep recurrent q-learning for partially observable mdp. In *2015 aaai fall symposium series*.
- [18] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [19] Nicolas Huck and Komivi Afawubo. 2015. Pairs trading and selection methods: is cointegration superior? *Applied Economics* 47 (2015), 599 – 613.
- [20] Gevorg Hunanyan. 2019. Portfolio Selection. *Finanzwirtschaft, Banken und Bankmanagement I Finance, Banks and Bank Management* (2019).
- [21] Musonda Katongo and Ritabrata Bhattacharyya. 2021. The use of deep reinforcement learning in tactical asset allocation. Available at SSRN 3812609 (2021).
- [22] Sang-Ho Kim, Deog-Yeong Park, and Ki-Hoon Lee. 2022. Hybrid Deep Reinforcement Learning for Pairs Trading. *Applied Sciences* (2022).
- [23] Taewook Kim and Ha Young Kim. 2019. Optimizing the pairs-trading strategy using deep reinforcement learning with trading and stop-loss boundaries. *Complexity* 2019 (2019).
- [24] Christopher Krauss. 2017. Statistical Arbitrage Pairs Trading Strategies: Review and Outlook. *Journal of Economic Surveys* 31, 2 (2017), 513–545. <https://doi.org/10.1111/joes.12153>
- [25] Christopher Krauss, Xuan Anh Do, and Nicolas Huck. 2017. Deep neural networks, gradient-boosted trees, random forests: Statistical arbitrage on the S&P 500. *Eur. J. Oper. Res.* 259 (2017), 689–702.
- [26] Tejas D. Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Joshua B. Tenenbaum. 2016. Hierarchical Deep Reinforcement Learning: Integrating Temporal Abstraction and Intrinsic Motivation. In *NIPS*.
- [27] John Langford and Tong Zhang. 2007. The Epoch-Greedy algorithm for contextual multi-armed bandits. In *NIPS 2007*.
- [28] John P. Lehoczky and Mark J. Schervish. 2018. Overview and History of Statistics for Equity Markets.
- [29] David A. Lesmond, Michael J. Schill, and Chunsheng Zhou. 2003. The Illusory Nature of Momentum Profits. *AFIA 2002 Atlanta Meetings (Archive)* (2003).
- [30] Yan-Xia Lin, Michael McCrae, and Chandra Gulati. 2006. Loss protection in pairs trading through minimum profit bounds: A cointegration approach. *Adv. Decis. Sci.* 2006 (2006), 73803:1–73803:14.
- [31] Jing-You Lu, Hsu-Chao Lai, Wen-Yueh Shih, Yi-Feng Chen, Shengkai Huang, Hao-Han Chang, Jun-Zhe Wang, Jun-Long Huang, and Tian-Shyr Dai. 2022. Structural break-aware pairs trading strategy using deep reinforcement learning. *The Journal of Supercomputing* 78 (2022), 3843 – 3882.
- [32] Giorgio Lucarelli and Matteo Borrotti. 2019. A deep reinforcement learning approach for automated cryptocurrency trading. In *IFIP International Conference on Artificial Intelligence Applications and Innovations*. Springer, 247–258.
- [33] Ian Martin. 2021. On the Autocorrelation of the Stock Market*. *Journal of Financial Econometrics* 19, 1 (01 2021), 39–52. <https://doi.org/10.1093/jfinc/ebaa033> arXiv:<https://academic.oup.com/jfec/article-pdf/19/1/39/36732400/nbaa033.pdf>
- [34] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*. PMLR, 1928–1937.
- [35] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. *nature* 518, 7540 (2015), 529–533.
- [36] Ofir Nachum, Haoran Tang, Xingyu Lu, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. 2019. Why Does Hierarchy (Sometimes) Work So Well in Reinforcement Learning? *ArXiv abs/1909.10618* (2019).
- [37] Shubham Pateria, Budhitama Subagdja, Ah-Hwee Tan, and Hiok Chai Quek. 2021. Hierarchical Reinforcement Learning. *ACM Computing Surveys (CSUR)* 54 (2021), 1 – 35.
- [38] Marcelo Scherer Perlin. 2007. M of a Kind: A Multivariate Approach at Pairs Trading. *Emerging Markets: Finance* (2007).
- [39] Andy Pole. 2007. Statistical Arbitrage: Algorithmic Trading Insights and Techniques.
- [40] Adrian Pope, Jaime Shinsuke Ide, Daria Mićović, Henry Diaz, David Rosenbluth, Lee Ritholtz, Jason C. Twedt, Thayne T. Walker, Kevin Alcedo, and Daniel Javorsek. 2021. Hierarchical Reinforcement Learning for Air-to-Air Combat. *2021 International Conference on Unmanned Aircraft Systems (ICUAS)* (2021), 275–284.
- [41] Heni Puspasingrum, Yan-Xia Lin, and Chandra Gulati. 2010. Finding the Optimal Pre-set Boundaries for Pairs Trading Strategy Based on Cointegration Technique. *Journal of Statistical Theory and Practice* 4 (2010), 391–419.
- [42] Hossein Rad, Rand Kwong Yew Low, and Robert W. Faff. 2015. The Profitability of Pairs Trading Strategies: Distance, Cointegration, and Copula Methods. *Wharton Research Data Services (WRDS) Research Paper Series* (2015).
- [43] Abdelrhman Saleh, Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, and Rosalind W. Picard. 2020. Hierarchical Reinforcement Learning for Open-Domain Dialog. In *AAAI*.
- [44] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. 2015. Taking the human out of the loop: A review of Bayesian optimization. *Proc. IEEE* 104, 1 (2015), 148–175.
- [45] William F Sharpe. 1994. The sharpe ratio. *Journal of portfolio management* 21, 1 (1994), 49–58.
- [46] Megan Shearer, David Byrd, Tucker Hybinette Balch, and Michael P. Wellman. 2021. Stability effects of arbitrage in exchange traded funds: an agent-based model. *Proceedings of the Second ACM International Conference on AI in Finance* (2021).
- [47] Tijmen Tieleman and Geoffrey Hinton. 2012. Rmsprop: Divide the gradient by a running average of its recent magnitude. coursera: Neural networks for machine learning. *COURSERA Neural Networks Mach. Learn* (2012).
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*. I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>
- [49] Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Manfred Otto Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. 2017. FeUdal Networks for Hierarchical Reinforcement Learning. *ArXiv abs/1703.01161* (2017).
- [50] Ganapathy Vidyamurthy. 2004. *Pairs Trading: quantitative methods and analysis*. Vol. 217. John Wiley & Sons.
- [51] Cheng Wang, Patrik Sandås, and Peter Beling. 2021. Improving Pairs Trading Strategies via Reinforcement Learning. In *2021 International Conference on Applied Artificial Intelligence (ICAAI)*. IEEE, 1–7.
- [52] Christopher JCH Watkins and Peter Dayan. 1992. Q-learning. *Machine learning* 8, 3-4 (1992), 279–292.
- [53] Ruobing Xie, Shaoliang Zhang, Rui Wang, Feng Xia, and Leyu Lin. 2021. Hierarchical Reinforcement Learning for Integrated Recommendation. In *AAAI*.
- [54] Fucui Xu and Shannon Siew Ngee. Tan. 2020. Dynamic Portfolio Management Based on Pair Trading and Deep Reinforcement Learning. *2020 The 3rd International Conference on Computational Intelligence and Intelligent Systems* (2020).

A STOCK SUBSET

In this section, we present the stocks of all 5 randomly split subgroups in the U.S. and Chinese stock markets in Table 3 and Table 4 respectively.

Set	Stocks
1	AMAT, AXP, BA, BAX, EA, EBAY, ED, EOG, GLW, IBM, IRM, LMT, MAS, MCO, MMM, MOS, NUE, PFE, PG, PPL, QCOM, RTX, SLB, SPG, SWKS, TGT, TXT, UNH, USB, WY
2	AAPL, ABT, ADSK, ALB, AMGN, APD, BLK, CAT, CDNS, CLX, COF, DE, DHI, EMR, EQR, FE, FMC, GIS, IP, IPG, JPM, MS, NEE, NEM, NTAP, NWL, ROP, ROST, TJX, VLO
3	ADBE, AES, AVY, BSX, C, CAH, CCL, CL, CMI, CTSH, DOV, DUK, EXC, F, GE, HSY, KO, KR, LUV, MRO, MSFT, NKE, PEAK, PLD, PNC, SCHW, SYY, UPS, VFC, YUM
4	A, ADM, ALL, ATVI, AZO, BMY, COST, CSCO, CVX, FCX, FDX, GS, HAL, HD, INTC, K, KIM, LEN, LOW, MCD, MMC, MRK, MSI, NVDA, PHM, STT, T, WMB, XOM, XRAY
5	AMZN, AON, APA, BAC, BBY, BEN, BK, CMCSA, CPRT, CVS, DHR, EIX, ETN, FAST, HON, HUM, MCK, MO, MTB, NLOK, PCAR, PGR, SBUX, TER, TRV, UNP, VZ, WFC, WHR, WMT

Table 3: The stocks of subsets for S&P 500.

Set	Stocks
1	000012,000016,000022,000024,000027,000029, 000046,000059, 000068,000088,000422,000518, 000520,000539,000541,000559, 000568,000598, 000625,000627,000671,000698,000708,000712, 000717,000729,000776,000825,000831,000878, 000916,000920, 000933,600061,600108,600109, 600118,600125,600138,600151, 600183,600228, 600239,600602,600611,600641,600642,600688, 600710,600717,600744,600770,600795,600811, 600820,600832, 600875,600884,600886,600893
2	000009,000089,000401,000402,000420,000425,000503,000507, 000536,000538,000553,000623,000709,000758,000806,000900, 000912,000927,000937,000938,600005,600009,600057,600058, 600066,600078,600111,600115,600123,600150,600153,600157, 600170,600177,600190,600196,600266,600600,600601,600606, 600639,600643,600649,600662,600724,600726,600780,600783, 600797,600804,600809,600812,600835,600838,600863,600879, 600880,600881,600887,600894
3	000031,000069,000423,000532,000562,000571,000599,000601, 000630,000631,000651,000661,000666,000690,000738,000750, 000767,000768,000778,000780,000793,000823,000828,000898, 000921,000930,000932,600072,600079,600085,600088,600091, 600096,600100,600104,600117,600135,600169,600171,600198, 600200,600210,600216,600608,600630,600635,600657,600663, 600704,600707,600718,600737,600739,600747,600823,600839, 600866,600868,600873,600874
4	000021,000036,000061,000063,000066,000408,000413,000498, 000528,000533,000550,000596,000612,000667,000682,000686, 000703,000707,000718,000728,000735,000737,000786,000883, 000886,000897,000939,000951,600006,600062,600074,600098, 600103,600110,600121,600126,600176,600215,600219,600220, 600221,600621,600633,600648,600654,600655,600660,600675, 600705,600748,600757,600761,600779,600808,600816,600827, 600837,600854,600867,600895
5	000001,000002,000039,000060,000400,000415,000429,000543, 000573,000581,000607,000629,000636,000652,000656,000659, 000680,000685,000727,000733,000783,000792,000800,000807, 000822,000826,000839,000858,000876,000895,000917,000949, 000959,600000,600007,600060,600068,600073,600089,600132, 600161,600188,600207,600208,600637,600638,600652,600653, 600664,600674,600690,600694,600703,600733,600741,600760, 600790,600805,600851,600871

Table 4: The stocks of subsets for CSI 300.

B RESULT PRESENTATION

We present a larger version of the detailed trading actions, positions, and profit during the trading period of the selected pair by BanditPair and GGR in Set 2 of S&P 500, as shown in Fig.6.

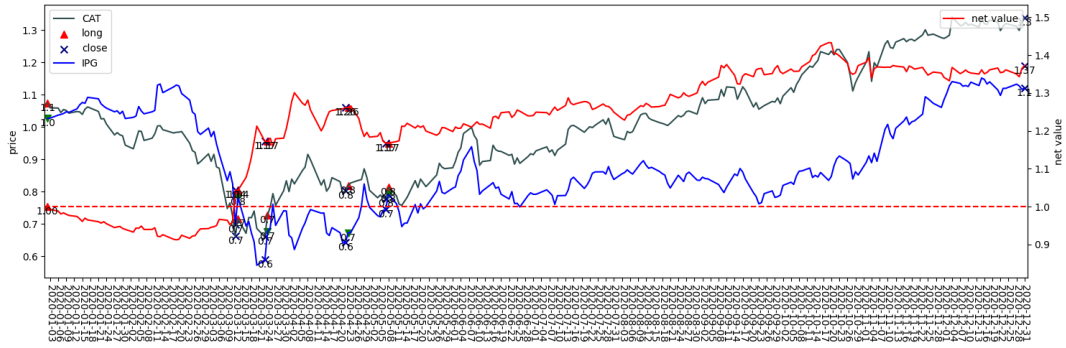
C DETAIL PERFORMANCE

As shown in Table 5 and Table 6, we report in this section the performance of each approach over all 5 subgroups of the U.S. and Chinese stock markets.

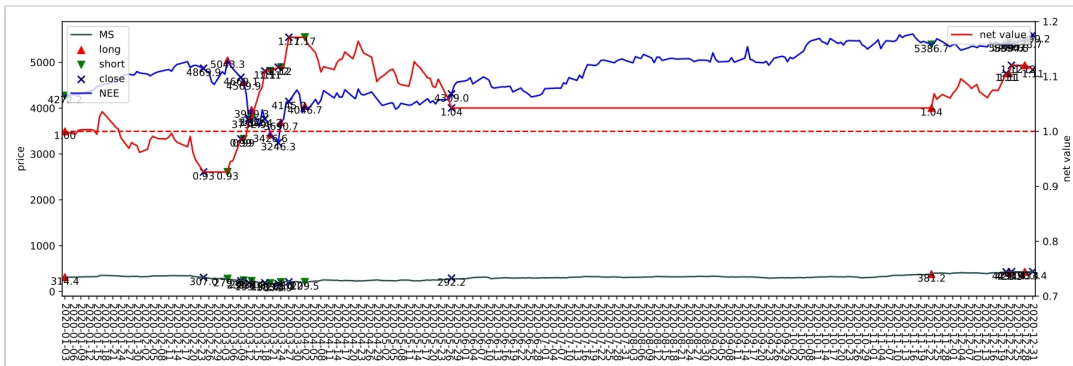
Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

Model		GGR	Cointegration	Correlation	Wang	TRIALS	TRIALS wo TR
Set 1	SR↑	-0.39	-1.56	-1.26	1.32	1.57	-0.08
	AR↑	-0.03	-0.39	-0.14	0.32	0.51	-0.04
	MDD↑	-0.15	-0.41	-0.23	-0.11	-0.09	-0.32
	AV↓	0.12	0.30	0.13	0.21	0.27	0.28
	ED↓	0.008	0.035	0.010	0.036	0.027	0.011
Set 2	SR↑	-2.01	-1.84	-1.58	0.73	2.10	-0.08
	AR↑	-0.29	-0.20	-0.19	0.14	0.64	-0.01
	MDD↑	-0.33	-0.25	-0.22	-0.08	0.11	-0.21
	AV↓	0.17	0.13	0.14	0.17	0.24	0.19
	ED↓	0.015	0.008	0.006	0.008	0.037	0.010
Set 3	SR↑	-1.49	-1.74	-1.30	1.70	1.86	0.02
	AR↑	-0.20	-0.74	-0.10	0.36	0.34	0.10
	MDD↑	-0.24	-0.74	-0.13	-0.10	-0.08	-0.20
	AV↓	0.16	0.65	0.09	0.18	0.15	0.20
	ED↓	0.022	0.044	0.005	0.005	0.029	0.008
Set 4	SR↑	-0.55	-1.66	-1.16	1.51	2.10	0.15
	AR↑	-0.06	-0.18	-0.08	0.10	0.67	0.03
	MDD↑	-0.10	-0.20	-0.16	0	-0.10	-0.33
	AV↓	0.14	0.13	0.09	0.05	0.25	0.24
	ED↓	0.016	0.006	0.006	0.006	0.060	0.302
Set 5	SR↑	-2.42	-2.35	-1.73	0.61	1.56	0.35
	AR↑	-0.16	-0.27	-0.20	0.13	0.34	0.07
	MDD↑	-0.19	-0.27	-0.25	-0.14	-0.09	-0.18
	AV↓	0.08	0.13	0.13	0.19	0.18	0.18
	ED↓	0.010	0.009	0.007	0.007	0.033	0.009
Mean (Std)	SR↑	-1.37 (0.79)	-1.83 (0.27)	-1.41 (0.21)	1.18 (0.43)	1.84 (0.24)	0.07 (0.16)
	AR↑	-0.15 (0.09)	-0.36 (0.20)	-0.14 (0.05)	0.21 (0.11)	0.50 (0.14)	0.01 (0.20)
	MDD↑	-0.20 (0.08)	-0.37 (0.20)	-0.20 (0.04)	-0.09 (0.05)	-0.09 (0.01)	-0.25 (0.07)
	AV↓	0.13 (0.03)	0.27 (0.20)	0.12 (0.02)	0.16 (0.06)	0.22 (0.04)	0.22 (0.04)
	ED↓	0.014 (5e-3)	0.021 (0.02)	0.007 (0.002)	0.021 (0.02)	0.037 (0.01)	0.01 (4e-3)

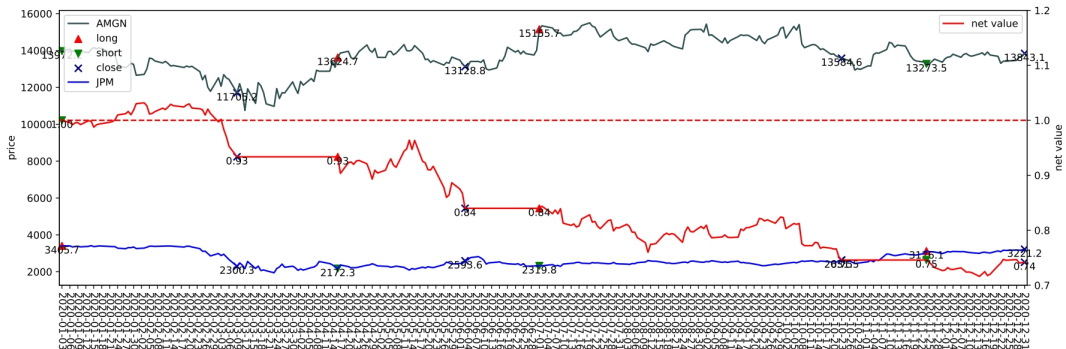
Table 5: Trading performance on S&P 500.



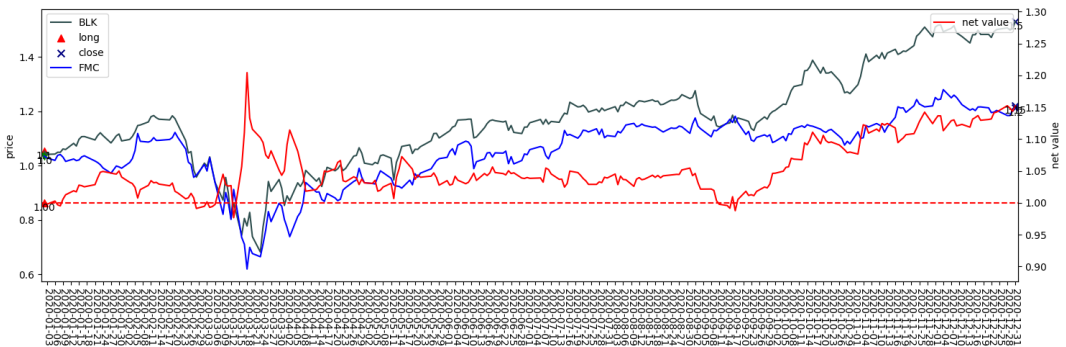
(a) The trading detail of TRIALS.



(b) The trading detail of TRIALS w/o TR.



(c) The trading detail of GGR.



(d) The trading detail of Wang et al.

Figure 6: The trading details

Model		GGR	Cointegration	Correlation	Wang	TRIALS	TRIALS wo TR
Set 1	SR↑	0.22	-2.19	-0.20	0.19	0.75	0.73
	AR↑	0.05	-0.31	-7e-3	0.05	0.13	0.14
	MDD↑	-0.27	-0.30	-0.11	-0.22	-0.08	-0.16
	AV↓	0.20	0.17	0.11	0.16	0.15	0.17
	ED↓	0.007	0.008*	0.011	0.008*	0.006	0.015
Set 2	SR↑	-1.18	-1.55	-1.27	0.17	2.41	1.11
	AR↑	-0.18	-0.21	-0.16	0.04	0.31	0.11
	MDD↑	-0.31	-0.27	-0.21	-0.18	-0.14	-0.10
	AV↓	0.17	0.16	0.14	0.12	0.10	0.19
	ED↓	0.016	0.007	0.031	0.007	0.023	0.008
Set 3	SR↑	-1.54	-0.60	-0.62	1.61	0.31	0.56
	AR↑	-0.15	-0.11	-0.08	0.24	0.06	0.08
	MDD↑	-0.15	-0.18	-0.18	-0.29	-0.12	-0.12
	AV↓	0.12	0.20	0.16	0.10	0.17	0.12
	ED↓	0.007*	0.018	0.009	0.018	0.011	0.007
Set 4	SR↑	-1.90	-2.68	-1.06	0.79	2.57	2.36
	AR↑	-0.29	-0.47	-0.22	0.29	0.53	0.45
	MDD↑	-0.08	-0.10	-0.07	-0.70	0	0
	AV↓	0.18	0.23	0.23	0.24	0.16	0.15
	ED↓	0.013	0.028*	0.014	0.028*	0.056	0.018
Set 5	SR↑	-0.68	-0.03	-1.76	1.91	1.03	1.03
	AR↑	-0.14	-6e-4	-0.24	0.69	0.17	0.17
	MDD↑	-0.32	-0.14	-0.23	-0.79	-0.11	-0.11
	AV↓	0.22	0.18	0.16	0.20	0.15	0.15
	ED↓	0.024	0.024	0.011	0.024	0.024	0.014
Mean (Std)	SR↑	-1.19 (0.74)	-1.50 (0.97)	-1.37 (0.25)	0.75 (0.68)	1.91 (0.88)	0.95 (0.88)
	AR↑	-0.17 (0.11)	-0.25 (0.17)	-0.21 (0.07)	0.24 (0.23)	0.68 (0.51)	0.13 (0.12)
	MDD↑	-0.29 (0.06)	-0.29 (0.13)	-0.25 (0.06)	-0.18 (0.09)	-0.14 (0.07)	-0.12 (0.09)
	AV↓	0.18 (0.03)	0.19 (0.03)	0.17 (0.05)	0.25 (0.07)	0.26 (0.09)	0.17 (0.07)
	ED↓	0.013 (6e-3)	0.017 (8e-3)	0.015 (8e-3)	0.017 (8e-3)	0.046 (0.02)	0.02 (8e-3)

Table 6: Trading performance on CSI 300.