# JAX-LOB: A GPU-Accelerated limit order book simulator to unlock large scale reinforcement learning for trading

Sascha Frey[*][†]
Department of Computer Science
University of Oxford
UK

Kang Li[*]
Department of Statistics
University of Oxford
UK

Peer Nagy[*]
Oxford-Man Institute of Quantitative Finance
University of Oxford
UK

Silvia Sapora
Foerster Lab for AI Research
University of Oxford
UK

Chris Lu
Foerster Lab for AI Research
University of Oxford
UK

Stefan Zohren
Man-Group
Oxford-Man Institute of Quantitative Finance
University of Oxford
UK

Jakob Foerster
Foerster Lab for AI Research
University of Oxford
UK

Anisoara Calinescu
Department of Computer Science
University of Oxford
UK

## ABSTRACT

Financial exchanges across the world use limit order books (LOBs) to process orders and match trades. For research purposes it is important to have large scale efficient simulators of LOB dynamics. LOB simulators have previously been implemented in the context of agent-based models (ABMs), reinforcement learning (RL) environments, and generative models, processing order flows from historical data sets and hand-crafted agents alike. For many applications, there is a requirement for processing multiple books, either for the calibration of ABMs or for the training of RL agents. We showcase the first GPU-enabled LOB simulator designed to process thousands of books in parallel, with a notably reduced per-message processing time. The implementation of our simulator – JAX-LOB – is based on design choices that aim to best exploit the powers of JAX without compromising on the realism of LOB-related mechanisms. We integrate JAX-LOB with other JAX packages, to provide an example of how one may address an optimal execution problem with reinforcement learning, and to share some preliminary results from end-to-end RL training on GPUs.

## CCS CONCEPTS

• **Computing methodologies** → **Discrete-event simulation**; **Reinforcement learning**; *Massively parallel and high-performance simulations.*

## KEYWORDS

limit order books, reinforcement learning, high frequency trading, trade execution, market replay, order book simulator

## 1 INTRODUCTION

Markets, i.e. matching buyers and sellers while determining a price, are a crucial component of modern economies. In many instances these markets rely on *auction mechanisms* for their operation. At the most basic level, for a single good sold at a specific time, these mechanisms allow all buyers to state how much they are willing to pay for a given good, and the good then goes to the highest bidder.

In contrast, in financial markets the price finding process usually needs to happen *continuously* and for *any number* of shares during the operating hours of the market to ensure liquidity. To accomplish this, the limit order book (LOB) mechanism is used in most modern electronic exchanges as a price finding tool.

At a high level, LOBs implement an *any-time* auction mechanism that allows all market participants to submit *buy* and *sell* orders specifying a *quantity* of stocks and a *price* at which they are willing to trade. As soon as there are compatible buy and sell orders, i.e. where the buying price of one order at least matches the selling price of another order, a trade happens.

Due to their central role in the financial system, the ability to accurately and efficiently model LOB dynamics is extremely valuable. For example, it might allow a financial company to offer better services or may enable the government to predict the impact of financial regulation on the stability of the financial system.

Practically, there are a number of possible scientific approaches for taking advantage of such an LOB simulator, including agent-based models, reinforcement learning (RL), and generative models.

However, due to the low signal-to-noise ratio a common aspect of these approaches is the need for large scale simulations that take advantage of modern compute hardware and parallelism. Addressing this, **a GPU accelerated LOB simulator is the core**

---
[*]Equal contribution
[†]Correspondance address: sascha.frey@st-hughs.ox.ac.uk

**contribution** of our paper. One specific task of interest is the execution of orders by brokers or funds, where the goal is to sell/buy a specific quantity of stocks over a given time frame at the best possible price. This *optimal execution problem* is both a motivation and an example use-case for the work presented in this paper.

In Section 2 we provide a brief overview of the LOB and the trade execution problem, and in Section 3 we consider related works both in building LOB simulators and addressing the execution problem with reinforcement learning (RL). In Section 4 we present the first GPU-accelerated LOB simulator using JAX [9] and discuss the detailed workings of the LOB, both general and specific to our implementation. In Section 5 we lay the groundwork for the use of our simulator for RL, an application domain which benefits significantly from the parallelism enabled by our implementation.

In Section 5.2 we contribute to solving the optimal execution problem with RL by wrapping the JAX-LOB simulator in a JAX-native *gymnax* [25] execution environment. In our environment, both the experience rollout (i.e. agents interacting with the world to collect data) and the learning updates (i.e. agents training on the data) are conducted on the same GPU which avoids the GPU-CPU communication bottleneck. Finally, in Section 6 we integrate the environment with *PureJaxRL* [26] and use recurrent proximal policy optimization (PPO) to train an execution agent as an example application. Lastly, we compare the computational speed with equivalent CPU based implementations and show that the maximum speedup occurs when applying JAX-LOB to an RL task, taking full advantage of the parallelism.

Our main contributions are as follows:
**1)** The first GPU-accelerated LOB simulator, JAX-LOB
**2)** At least 7x speedup in training an RL agent, compared to equivalent CPU-based LOB simulators: 550 versus 74 steps per second during training, on the same hardware (Section 6.2).

We commit to open-sourcing JAX-LOB incl. the RL integration and strongly believe it will unlock a number of new research avenues for the community.

## 2 BACKGROUND

**The LOB**, the underlying data structure and associated mechanisms which power modern electronic exchanges, has been studied extensively [8, 17]. The LOB is the collection of all orders to buy or sell a security that have yet to be matched. Two types of orders may be submitted to the LOB at any time: limit orders which require the specification of a price and a quantity, and market orders, which require only a quantity. Traders are free to submit orders at any time, at which point they are either matching with a compatible opposite order in the book, or added to the book – in the case of a limit order. Further details of the LOB are discussed in Section 4, together with details specific to our implementation. LOB simulators have a number of application domains, including generative models [3]; in this paper we focus on their use to train trading agents by using RL for trading tasks such as market making [7, 15] or trade execution [27].

**JAX** [9] is an accelerator agnostic framework that enables just-in-time (JIT) compilation using accelerated linear algebra (XLA) [16], automatic differentiation, and automatic vectorization which can easily be executed on the GPU. The framework is designed for high-performance machine learning research and forms the basis

of both the *gymnax* [25] and *PureJaxRL* [26] frameworks which we leverage, using optimal trade execution as an example problem. JAX has also been used for other applications where massive parallelism is advantageous, such as fluid dynamics [23] or molecular dynamics [32]. Whilst the focus of this paper is on the JIT compilation and GPU vectorization, the differentiation of LOB models for calibration [31] is a topic that is thus far underexplored in the literature.

**Optimal trade execution** is a well-studied problem in finance, with traditional approaches using methods from stochastic optimal control and analytic market-impact functions [1, 6, 24, 29] to derive optimal policies. The aim of optimal trade execution is to buy or sell a specified number of stocks in a specified time frame, obtaining the best possible average price at low risk. This usually requires minimizing the cost due to market impact, especially when the trades are large compared to the liquidity of the market [1].

## 3 RELATED WORK

There are a number of implementations that reconstruct the mechanics of the LOB [5, 10, 11, 19, 30]. The implementations may be further subdivided by the application they are used for: agent-based models (ABMs) and market replay. Frameworks built for ABMs, such as the agent-based interactive discrete event simulation (ABIDES) framework [10] or the multi-agent exchange environment (MAXE) [5] implement a heterogeneous set of agents to submit fictional orders to the order book. ABIDES is implemented in standard Python and has an extension, ABIDES-gym [2], which creates a gym-like interface for reinforcement learning. MAXE is implemented in C++ and benefits from increased speed due to compilation [5]. In market replay models, such as the reinforcement learning for market making (RL4MM) framework [19], historical orders based on message data, such as the "limit order book system - the efficient reconstructor" (LOBSTER) [18] data set are submitted. The primary limitation of ABMs is the difficulty in creating sufficiently realistic agent behaviors. Despite significant recent efforts [36], this remains a limitation, especially when compared to replayed market data, which directly replicates the data distribution. However, the strength of ABMs is also a key limitation of market replay [20]. The implementation of strategic agents results in a dynamically reactive order flow based on outside intervention (e.g. a RL agent), whereas replayed historical order flow will remain unchanged. JAX-LOB, while originally designed to process data from LOBSTER [18], may process messages from any source, including strategic agents. In contrast to the discussed implementations, JAX-LOB runs on the GPU, allowing for large-scale parallelization.

Ning et al. [28] use double deep Q Networks (DDQN) [35] to solve the optimal execution problem using only market orders. These are, however, generally sub-optimal as trades may "walk the book", i.e. consume the entire available volume on one or more levels of the LOB, and eliminate potential gains from the spread by using passive limit orders. Furthermore, market impact is stylized by using a penalty function for larger orders. In this setting, they outperform the time-weighted average price (TWAP) benchmark in 7 out of 9 evaluated stocks. Dabérius et al. [13] compares the use of DDQN [35] and PPO [34] for optimal execution using only market orders in a stochastic model of price dynamics without price impact. In building a LOB simulator, we provide the tools to

move away from having to stylize the market impact and create an environment that allows for the submission of limit orders at various prices, thus avoiding the risk of inadvertently consuming volume on multiple levels of the book. Nevertheless, the exclusive use of historical LOB data means that only direct market impact is accounted for. Such historical messages have no strategic behavior that may react to an RL agent's actions.

Addressing this drawback, Karpe et al. [21] uses DDQN in an agent-based simulation environment based on the ABIDES simulator [10], which replays real LOB message data while adding additional heuristic agents, which react to changes in the LOB state caused by the RL agent. Their simulation setup allows addressing optimal execution and optimal level placement jointly. They find that their agent converges to a TWAP strategy. While additional agents allow for more realistic market impact in principle, the ABIDES heuristics use a somewhat simplistic momentum strategy based on aggregate summary statistics, instead of a learned rational strategy. Similarly, Fang et al. [14] use PPO to address the combined optimal execution and placement problem using the approach of replaying historical data in the ABIDES [10] simulator. Their agent learns placement of up to three levels of the LOB concurrently. They use a large number of engineered features, including technical indicators, an attention-based network architecture, dubbed Dual-Window Denoise PPO, and a reward formulation based on execution cost and an imitation term, steering the policy towards a TWAP policy. Results indicate a potential improvement over TWAP but suffer from high standard errors. The environment we set up aims to move away from such engineered feature spaces and reward functions, and focuses on using recurrent neural networks exclusively to automatically extract features and account for the longer-term memory of the agent.

## 4 THE JAX-LOB SIMULATOR

Some of the key challenges in applying deep reinforcement learning to trade execution and other high-frequency tasks are the low signal-to-noise ratio, the risk of overfitting to specific training days, and building a simulator with realistic market impact. A straightforward way to compensate for the first two issues is to increase the number of state-action transitions available for training. To speed up the generation of observations of faithful LOB representations using high-frequency data we use JAX [9].

We make a number of design decisions to address some of the constraints of JIT-compiled JAX code:

- Pure functions without side effects (e.g. global variables)
- Fixed array size and type
- Control-flow that can be compiled and parallelized efficiently (See Section 4.3).

Most CPU implementations of LOBs are based on hash-maps, queues, doubly-linked lists and sorted dictionaries [5, 10, 19] which enable quick access to data and maintain sorting of the book throughout. Given the JAX requirement of compile-time fixed-size arrays, implementing a similar structure means memory must be pre-allocated for all price-levels and orders. The use of arrays means that re-ordering upon removal of entries is far more costly than with linked lists. We, therefore, opt for an architecture that does not use a tree-like structure, nor do we keep orders sorted at all times.

Instead, we define two arrays $\mathcal{A}$ and $\mathcal{B}$ to represent the sides of the order book containing all active ask and bid orders, respectively.

$$\mathcal{A} = [\boldsymbol{a_1}, ..., \boldsymbol{a_N}] \quad \mathcal{B} = [\boldsymbol{b_1}, ..., \boldsymbol{b_N}] \tag{1}$$

$$\boldsymbol{o_i} = [P_i, Q_i, OID_i, TID_i, Ts_i, Tns_i] \in \mathcal{A} \cup \mathcal{B}, \quad i \in [1, N] \tag{2}$$

Each side of the book has a fixed capacity for $N$ orders, where each order $\boldsymbol{o}$ has six features (2): Price $P$, Quantity $Q$, Order ID $OID$, Trader ID $TID$, Time (seconds) $Ts$, and Time (nanoseconds) $Tns$. Empty positions in $\mathcal{A}$ or $\mathbf{B}$ are indicated by setting all features to -1. The size $N$ must be chosen such that the book is not saturated during a given experiment.

### 4.1 Basic operations

There are three operations that can be applied to either side of the order book; the *addition* of a new order, the *cancellation* of an existing order, and the *matching* of an existing order with an incoming order on the other side of the book followed by its removal from the book.

(1) The addition of an order requires the identification of an empty position ($\boldsymbol{o_i} = -1$) in the array, and the insertion of the order-specific data into the correct fields.

(2) Cancellation requires locating the order ID ($OID$) of the order to be canceled and removing the corresponding quantity from the book.

(3) During the matching operation, an incoming aggressive order, denoted $\boldsymbol{o_a}$, is matched against an existing order on the other side of the book, the standing order $\boldsymbol{o_s}$. The matching logic aims to find the remaining quantities for both the aggressive ($Q_a \in \boldsymbol{o_a}$) and standing ($Q_s \in \boldsymbol{o_s}$) orders using the following operations

$$Q'_s = max(0, Q_s - Q_a), \quad Q'_a = Q_a - Q_s.$$

When two orders are matched, a trade $\boldsymbol{t_j}$ is logged (3).

$$\boldsymbol{t_j} = [P_j, Q_j, OID_{a,j}, OID_{s,j}, Ts_j, Tns_j] \tag{3}$$

Whereby:

$$P_j = P_s$$
$$Q_j = Q_s - Q'_s$$
$$OID_{a,j} = OID_a \in \boldsymbol{o_a}$$
$$OID_{s,j} = OID_s \in \boldsymbol{o_s}$$
$$Ts_j, Tns_j = Ts_a, Tns_a$$

Up to N trades are logged in a fixed-size array $\mathcal{T}$ due to constraints of the XLA compiler (4).

$$\mathcal{T} = [\boldsymbol{t_1}, ..., \boldsymbol{t_N}]^T \tag{4}$$

Again, $\boldsymbol{t_j} = -1$ indicates an empty slot.

For all operations, upon completion, both *Asks* and *Bids* are checked for orders $\boldsymbol{o_i}$ where $Q_i \leq 0$, in which case we set $\boldsymbol{o_i} = -1$.

A single incoming order may match with more than just one standing order. The matching logic contains a while loop that repeatedly attempts to match the incoming order against the next-best standing order on the opposite side of the book. The best standing order $Best(\boldsymbol{o_s})$ is defined by the price-time order priority algorithm, the most commonly used LOB matching algorithm [17]. The best

order is the one with the ask price or bid price (5). If multiple orders share this price, the one with the earliest arrival time is considered.

$$P_{ask} = \min_{o \in \mathcal{A}}(P_i) \quad P_{bid} = \min_{o \in \mathcal{B}}(P_i) \qquad (5)$$

The while-loop continues as long as the book is non-empty, $Q_a > 0$, and the prices overlap as follows

$$P_a \leq P_s \qquad \text{if marketable sell order,}$$
$$P_a \geq P_s \qquad \text{if marketable buy order.}$$

Based on the above descriptions it is intuitive to expect varying computational complexity across these basic operations. To validate this intuition, we time the basic operations in order books of different maximum capacities $N$ (Table 1). The operations take longer with increasing capacity in the book $N$, and as expected, the slowest operation (*Match*) takes more than twice as long to execute as the fastest (*Cancel*).

**Table 1: Average run-time for the three basic operations that may be applied to the LOB for different sizes $N$. The tests are conducted on an Nvidia 2080 Ti and averaged across 1000 serial trials. The book is initially filled to one-third capacity with random orders.**

| Capacity $N$ | Add (ms) | Cancel (ms) | Match (ms) |
|---|---|---|---|
| 10 | 0.115 | 0.081 | 0.184 |
| 100 | 0.117 | 0.095 | 0.206 |
| 1000 | 0.162 | 0.093 | 0.243 |

The relative size of the incoming order $Q_a \in o_a$ will have a significant impact on the computation time for the matching operation. All else being equal, a larger quantity requires the consideration of more standing orders to match the arriving order, thus computing more iterations of the matching logic. A simple way to illustrate this is to consider an order book with capacity $N = 100$, filled to one-third capacity, and to submit market orders of varying sizes (Table 2). This increase in time is not drastic until particularly large orders are submitted, but shows the importance of constraining the action space to reasonable quantities to ensure quick execution.

**Table 2: Analysis of the effect of incoming order size on the time required to complete the matching operation. As the size increases, more standing orders need to be considered to fully match the incoming order. The capacity for the book is set to $N = 100$. Testing is done on an Nvidia 2080 Ti GPU.**

| Market $Q_a$ | Time to match (ms) |
|---|---|
| 0 | 0.132 |
| 10 | 0.206 |
| 500 | 0.271 |
| 1,000 | 0.336 |
| 10,000 | 2.326 |

## 4.2 Message types

Which of the three basic operations (Section 4.1) is called depends on the type $T$ of message $m$ (6) passed to the order book.

$$m = [T, S, Q, P, OID, TID, Ts, Tns] \qquad (6)$$

The structure of data contained in messages $m$ is:

(1) $T$ - The type of message: limit orders ($T = 1$), cancel orders ($T = 2$), delete orders ($T = 3$), and market orders ($T = 4$).

(2) $S$ - The side of the order: bid ($S = 1$) or ask ($S = -1$).
(3) $P$ - The price at which the order should be submitted or the price of the order to be cancelled/deleted. This is ignored for market orders.
(4) $Q$ - The size of the order, i.e. the quantity to buy/sell or the to remove from an existing order for cancel/delete orders.
(5) $OID$ - The unique order ID.
(6) $TID$ - The trader ID, an identification marking the source of the order. Used to identify orders from individual agents.
(7) $Ts, Tns$ - The time of receipt of the order, split into two fields representing the time as full seconds and fractional digits as nanoseconds represented by 32-bit integers.

While the reader familiar with the LOBSTER [18] data set will recognize significant similarities to the message structure in this data set, they are not identical. We take the following steps to process orders of different types:

(1) Limit orders are first matched against the opposite side of the book, the remaining quantity is then added as a single new order to the corresponding side.
(2) Cancel and delete orders are treated identically, by calling a cancellation (Section 4.1).
(3) Market orders are matched against the opposite side with a price $P_m = 0$ if an ask order, or $P_m = max\_int$ if a bid order. Any unmatched quantity is disregarded.

Considering each type and side as eight separate cases allows defining explicit functions for each case, enabling the use of a single conditional statement upon receipt of the message, rather than multiple branches within the matching logic, which we find improves performance under `vmap` (Section 4.3). The compute time (Table 3) to process each of the three message types varies and is related to the basic operations required. When compared to the results in Table 1 the time is only marginally increased by the branching statement, and there is still a notable difference across different order types.

**Table 3: Time to process different message types. The last column differs from the second column in that the price crosses the spread, requiring a matching operation. Tests are done on an Nvidia 2080 Ti GPU, the book capacity is $N = 100$.**

| Capacity $N$ | Limit Order (ms) | Cancel Order (ms) | Limit Order across book (ms) |
|---|---|---|---|
| 10 | 0.108 | 0.077 | 0.163 |
| 100 | 0.157 | 0.097 | 0.203 |
| 1000 | 0.195 | 0.095 | 0.247 |

## 4.3 Vectorising map - `vmap`

The benefit of implementing the LOB in JAX is derived from parallelizing computation on the GPU. However, the processing of messages sent to the order book in a continuous double-auction is intrinsically a serial process: both the treatment of messages and the internal matching logic must be strictly ordered. This is in opposition to a classical auction with, say, a central clearing process [33], or even stochastic models of the limit order book [12]. Therefore, to achieve parallelism we process multiple books in parallel using the `vmap` operator.

One particularity of the `vmap` operator is that it transforms a number of control-flow statements into `select` statements which

execute all of the conditional branches at run time. In the case of limit order books, this means that each of the possible eight cases is computed and the runtime is therefore limited by the slowest branch. This can be observed when looking at the timed results (Table 4) when processing the same message across $N\_books = 1000$ identical order books in parallel. The time no longer varies across different types of messages but only depends on the capacity. The reason we avoid continuous sorting is that when removing orders in an array-based view, the entire array needs to be re-sorted. This order sorting is comparatively costly, which would not be problematic if called rarely. However, the aforementioned behavior of the vmap operator means that it would be called for every message, even if no order needs to be removed. Even though we have done extensive tests, it is possible that there is still a bottleneck of this nature which would explain the worst-case message processing time of over 2 seconds in Table 4.

**Table 4: Time to process different message types with vmap. The difference to table 3 is that messages are processed by $N\_Book = 1000$ books in parallel. Times measure the processing of a message in all books. The effective processing time for a single message can therefore be interpreted as being in $\mu s$ and compared directly to the times in Table 2.**

| Capacity | Limit Order (ms) | Cancel Order (ms) | Limit Order across book (ms) |
|---|---|---|---|
| 10 | 1.402 | 1.407 | 1.393 |
| 100 | 2.649 | 2.614 | 2.612 |
| 1000 | 11.75 | 11.43 | 11.47 |

The resulting times per message type are compared to two CPU implementations (Table 5): one implementation of the LOB using RB trees and linked lists as well as a similar implementation named RL4MM [19] using numpy. The CPU environments are described in more detail in sections 5 and 6. When parallelizing to more than 1000 books, our JAX-LOB has a faster per-message processing time.

**Table 5: Processing times for messages on the CPU order book implementation based on *linked list* and *RB trees*, and the implementation used in RL4MM [19] based on *numpy* arrays, and for the JAX-LOB implementation based on *jax arrays* with capacity $N = 100$ and $N\_Books = 1000$ books in parallel. A key difference between the CPU implementation and that in RL4MM is the fact that data is pre-loaded, as opposed to being accessed from a database. Both CPU implementations are tested on a single core of an Apple M1 chip and the JAX-LOB on an Nvidia 2080 Ti.**

| Order type | RL4MM Time ($\mu s$) | CPU Time ($\mu s$) | JAX-LOB Time ($\mu s$) |
|---|---|---|---|
| Limit | 4.91 | 5.3 | 2.6 |
| Cancel | 16.12 | 3.6 | 2.6 |
| Limit Cross | 37.94 | 7 | 2.6 |

## 5 *GYMNAX* ENVIRONMENTS FOR THE LOB

Our JAX-LOB is well suited for the *gymnax* reinforcement learning environment framework designed around JAX to produce parallel environment roll-outs on the GPU. In the following sections we briefly describe the construction of a general base environment that can be used to derive environments for various tasks, such as market-making, trade execution, or other intraday tasks. We further describe our example order execution environment, which we use in Section 6 to begin training a recurrent PPO-based RL agent.

### 5.1 Base environment

The primary roles of the base environment are to load LOBSTER [18] data for replay, provide an interface to the JAX order book, and set up a skeleton trading environment based on JAX-LOB.

*5.1.1 Data loading.* We opt to pre-load message data in fixed-size, non-overlapping time windows (figure 1) to save time during execution. The number of messages per step is constant, with a variable time duration per step. Conversely, the duration of a single time window is fixed, implying a variable horizon [2, 19], i.e. maximum number of steps. This allows the execution of trades with a strict time constraint. To satisfy the fixed-size array constraints, additional steps (indicated by the gray borders in Figure 1) with all values set to zero are added. These steps are never processed as the episode termination condition is satisfied beforehand.

Further, the order book state (Level-2) of the first ten levels of the book at the beginning of each time window is used to initialize the book. Unlike the message data, this does not provide insight into whether price levels contain one or multiple orders, nor does it provide any order IDs. We work around this by assuming that there is exactly one order per price level containing the sum of the listed volumes. We use $OID_i \in [-inf, -9000]$ starting at $-9000$ and descending for these initial orders. We allow for the cancellation of these orders if a cancellation message satisfies $P_i = P_{cancel}$ and $OID_i < -9000$.
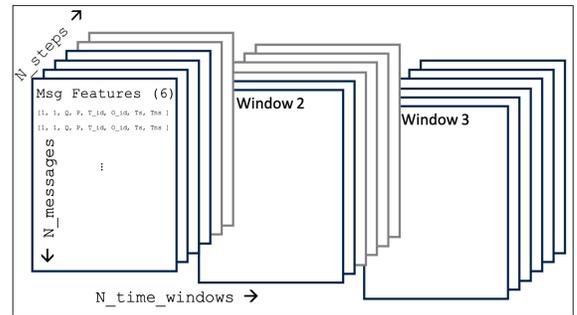


**Figure 1: Structure of data loaded into the *gymnax* environments. In our experiments, each step processes N_messages = 100 messages. The number of N_steps is variable with a fixed time of 30 minutes but the data is padded with zeros (grey boxes). There are N_time_windows per trading day of data.**

*5.1.2 Core state-space requirements.* The base environment defines the core features required in the state space, additional elements may be added for derived environments:

- Order book:
  - Ask side: Array($N_{orders} \times N_{features}$)

– Bid side: Array($N_{orders} \times N_{features}$)
- Trades: Array($N_{trades} \times N_{features}$)
- Initial time: $[T_s, T_{ns}]$
- Current time: $[T_s, T_{ns}]$
- ID counter
- Window Index
- Step counter

The ID counter is used to generate orders submitted by the agent(s), and is incremented accordingly during an episode.

*5.1.3 Functionality.* We keep the description of the base environment brief, leaving consideration of the action and observation spaces, and the reward function to the execution environment in Section 5.2. The core functionality of the environment during a step is as follows:

(1) The received actions (environment dependent) are transformed into messages to be processed by the JAX-LOB.
(2) The data messages are obtained from the loaded data in the environment parameters using the window index and step counter variables (Section 5.1.1) for indexing.
(3) The current time is updated based on the last data message.
(4) All messages are processed by the JAX-LOB.

The episode terminates when the elapsed time is longer than the episode time.

## 5.2 Execution environment

The execution environment extends the state space and defines the observation and action space, reward function, and a new termination condition.

*5.2.1 Augmented state space.* The state space in 5.1.2 is augmented with intra-step data and agent-specific information:

- Level 1 data for every processed message from data since the last step:
  – Best ask price: Array($N_{messages}$)
  – Best ask volume: Array($N_{messages}$)
  – Best bid price: Array($N_{messages}$)
  – Best bid volume: Array($N_{messages}$)
- the initial mid-price $\frac{P_{ask}+P_{bid}}{2}$ at the beginning of the episode
- The size of the execution task
- The quantity executed thus far
- The total revenue due to execution thus-far

*5.2.2 Observation space.* The state space is internal to the environment. The observation space is accessible by the RL agent and is designed to contain information that is accessible in real-life and avoids redundancy. The structure is a single one-dimensional array of size $N_{messages} \times 6 + 2 \times 2 + 4$ containing the following metrics:

- Best bid prices between steps: Array($N_{messages}$)
- Best ask prices between steps: Array($N_{messages}$)
- Mid-prices between steps: Array($N_{messages}$)
- The prices $n$ ticks deep on the side of the book of the task between steps

$$P_{passive} = \left\{ \begin{array}{ll} P_{ask} + n \times ticksize, & \text{if } task \text{ is } sell \\ P_{bid} - n \times ticksize, & \text{if } task \text{ is } buy \end{array} \right\} \quad (7)$$

- The spreads between steps: $spread = P_{ask} - P_{bid}$: Array($N_{messages}$)
- The current time: $[T_s, T_{ns}]$
- The elapsed time in the episode: $[T_s, T_{ns}]$
- The initial mid price at the episode start: $P_{init}$
- The price drift: $P_{mid} - P_{init}$
- The size of the execution task: $Q_{task}$
- The quantity executed thus far: $Q_{exec}$
- The Level 1 imbalances between steps: $Q_{ask} - Q_{bid}$: Array($N_{messages}$)

*5.2.3 Action space.* In order to have a rich set of actions whilst maintaining a reasonably simple learning task, we opt for a continuous action space of four dimensions. Similar to [14], at every step, the agent must choose the size of the order to place at four pre-defined prices:

(1) The "far touch" price: the best price on the opposite side of the book

$$P_{far\_touch} = \left\{ \begin{array}{ll} P_{bid} & \text{if } task \text{ is } sell \\ P_{ask} & \text{if } task \text{ is } buy \end{array} \right\}. \quad (8)$$

(2) The Mid-price (Section 5.2.2).
(3) The "near touch" price: best price on the side of the book of the task.
(4) The "passive" price, as defined in Section 5.2.2.

*5.2.4 Reward function.* We propose a reward function that combines the agent's ability to outperform the baseline volume-weighted average price (VWAP) strategy and the ability to use predicted price trends. It is composed of two parts, the advantage and the drift. The former represents the advantage gained over the average price of execution, weighted by volume, during the step. The latter represents the effect of price movements of the average traded price with respect to the mid-price at the beginning of the episode. The parameter $\lambda$ can be adjusted to weigh the effect of this drift. In equations (9) and (10), the indices $i$ and $j$ respectively refer to the set of executed trades during a step, and the set of executed trades of the RL agent.

$$R = \sum_j Q_j(P_j - P_{VWAP}) + \lambda \sum_j Q_j(P_{VWAP} - P_{init}) \quad (9)$$

$$P_{VWAP} = \frac{\sum_i Q_i P_i}{\sum_i Q_i}, \quad \sum_j Q_j < \text{ task size} \quad (10)$$

*5.2.5 Termination condition and market order submission.* In comparison to the base environment, the execution environment allows for an additional termination condition: the completion of the task, i.e. the execution of the desired quantity. In practice, the step function is modified such that this becomes the only termination condition, by submitting a market order for the remaining quantity one minute before the episode time ends.

## 5.3 Rollout cost comparison

Before considering the RL problem in Section 6, we consider the run-time improvements of the execution environment based on the JAX-LOB over comparable gym environments running on the CPU.

We compare our execution environment with the RL4MM [19] environment designed for market making, but with comparable

core logic. All orderbooks are run with data from January 2nd 2015 for ten levels for the stock TSLA. In the CPU and JAX-LOB environments, 100 messages are processed per step. RL4MM measures step duration in time, but we find that 100 messages are processed roughly every 2 minutes.

In the case of the JAX-LOB environment, we run different numbers of environments in parallel on an Nvidia 2080 Ti GPU, whilst the RL4MM and the CPU environments are run in series on an Apple M1 CPU. The improvement (Table 6) is roughly a factor 10 over the RL4MM environment and a factor 5 over our CPU implementation for 1000 parallel environments. Increasing the number of environments further is a feasible possibility, especially on GPUs with larger memory, but such economies of scale are eventually limited due to memory.

**Table 6: Average time per step comparison between the RL4MM environment, our CPU-based *gym* environment, and the JAX-LOB based *gymnax* environment for a varying number of environments in parallel. The CPU environments are run on an Apple M1 processor and the *gymnax* environment is run on an Nvidia 2080 Ti.**

| - | RL4MM | CPU | *gymnax* | | |
|---|---|---|---|---|---|
| N_Envs | - | - | 10 | 1000 | 10000 |
| Total time (s) | 29.6 | 33.8 | 69 | 329 | 2306 |
| N_steps | 5k | 9k | 5k | 500k | 5000k |
| Time/step (ms) | 5.9 | 3.5 | 13.8 | 0.66 | 0.46 |

## 6 USING THE *GYMNAX* ENVIRONMENT FOR REINFORCEMENT LEARNING

### 6.1 Training loop

We re-implement Recurrent-PPO [14] to allow for a continuous action space. The full architecture of the actor and critic networks is given in Figure 2. For every network update step, 10 steps are collected across 1,000 environments to create a batch size of 10,000 transitions. For each of the four epochs per update step, this batch is subdivided into four mini-batches which are used to calculate the loss functions, take the gradients and update the network through gradient descent using the Adam [22] optimizer implemented in optax [4].

### 6.2 JAX enabled training speedup

Whilst we do not have benchmarks for the exact same problem definition in other approaches, we nevertheless compare the training speeds for another environment running on the CPU. This is in addition to the reported data in Section 5.3 as the transfer of data from the CPU to GPU for the network update adds significant overhead to the run-time. With *gymnax* running on the GPU, this is not necessary and we see further improvements in Figure 3. The RL4MM package, which we use for comparison in Sections 4.3 and 5.3, is omitted in this comparison as we could not replicate fair experimental conditions due to issues arising with the use of the GPU. To qualify the comparison, we give a very brief overview of the CPU-based order book we use for the comparison. It is based on a LOB simulator using the traditional tree-like structure and uses the Stable-Baselines3 implementation of Recurrent PPO to solve
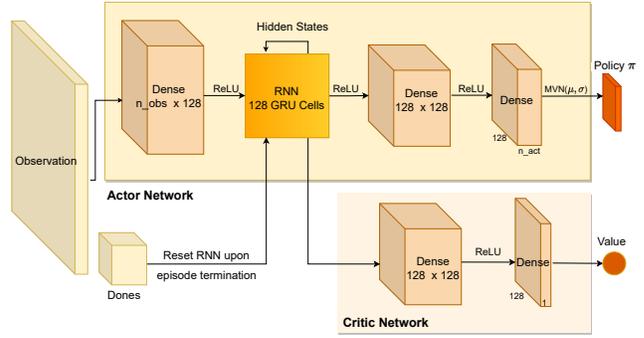


**Figure 2: The architecture of the RNN-PPO actor and critic networks. The hidden states of the RNN layer add the necessary recurrence for the network to have memory of previous states. The Boolean termination values ("Dones") are required to reset the memory of the RNN at episode termination. The policy $\pi$ is a multivariate normal distribution $mvn(\mu, \sigma)$ whereby $\sigma$ is a trained network parameter and $\mu$ is the output vector from the final dense network layer.**

an execution task. There are a number of key variables which are common across both experiments which we aim to control as far as possible to make the comparison equitable:

- *Messages processed per step* - 100 messages
- *Levels of LOBSTER data considered* - 10 levels
- *Episode length* - 30 minutes
- *Hardware* - GPU: Nvidia A40, CPU: 32 core AMD EPYC 7513, both experiments run network updates on the GPU.
- *Parallelization* - 1000 environments in parallel for JAX-LOB, 1000 vectorized environments for CPU but collected serially.
- *Batch size* - 10,000, 10 steps across 1000 environments
- *Mini-batch size* - 2,500, 4 per batch
- *Number of Epochs* - 4 epochs
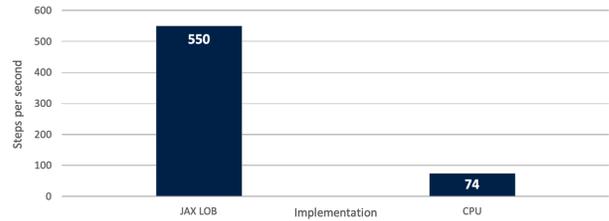- *Data* - April 2021 LOBSTER high frequency data for AMZN



**Figure 3: Training steps per second for gymnax environment based on JAX-LOB and our CPU gym environment implementation. The former is more than 7 times faster.**

We observe a 7x speedup of the training loop based on the JAX-LOB simulator compared to that using the CPU implementation of the LOB. This is larger than the reported 5x speedup in Table 2. We attribute this to the lack of data transfer across processing units. This does not yet consider hyper-parameter tuning which would further benefit from parallelization, and is trivial to achieve with our fully JAX-based implementation.

## 6.3 Training execution task

We begin to train an execution agent using the proposed setup on one day (April 1st, 2021) of data, chosen to be this narrow as the scope of this paper is primarily to showcase the functionality of the presented framework. The training curve, with a moving average window of 30 steps is shown in Figure 4. We set $\lambda = 0$ to make training easier, as the drift part of the reward function is much harder to learn, and compare the training curve to the TWAP benchmark strategy, which is to execute the desired volume linearly over the episode duration. Though the TWAP strategy is shown to be outperformed in training, we cannot make any claims on the success of this agent without conducting complete out-of-sample tests. Nevertheless, it indicates a promising direction of research based on the presented framework.
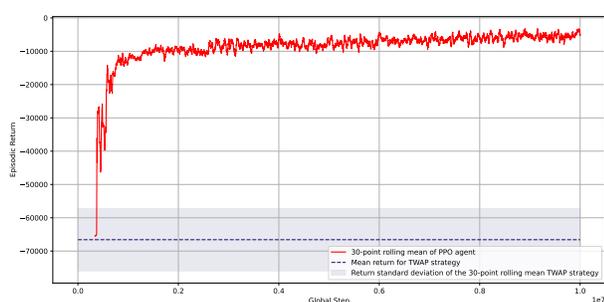


**Figure 4: Episodic return with $\lambda = 0$ (9) for a PPO agent during training and the benchmark TWAP strategy mean return on a single day of data with 30-point rolling mean.**

## 7 CONCLUSIONS AND FUTURE WORK

We present the first implementation of a limit order book simulator for GPUs, called JAX-LOB. Using this as part of a *gymnax* environment for reinforcement learning we obtain at least a 5x speedup compared to a comparable CPU implementation. When using the environment to train a RL agent with *PureJaxRL* we even see a 7x speed-up with respect to our comparable CPU implementation. This speedup due to parallelization is expected to contribute to research in applying RL to high-frequency trading and execution problems that require a reactive LOB simulator.

As part of this paper, we provide an example use case of training an RL agent for the execution task on a single day of message data. We plan to extend this work by following a rigorous RL training pipeline for the execution problem amongst others, including out-of-sample testing and hyper-parameter tuning. The latter is expected to benefit further from the parallel nature of the JAX-LOB implementation. Detailed experiments with respect to the action-space, feature-space, reward shaping, and network architecture will be required to train a robust and cost-effective RL execution agent.

## REFERENCES

[1] Robert Almgren and Neil Chriss. 2001. Optimal execution of portfolio transactions. *The Journal of Risk* 3, 2 (Jan. 2001), 5–39. https://doi.org/10.21314/jor.2001.041 Publisher: Infopro Digital Services Limited.
[2] Selim Amrouni, Aymeric Moulin, Jared Vann, Svitlana Vyetrenko, Tucker Balch, and Manuela Veloso. 2022. ABIDES-gym: gym environments for multi-agent discrete event simulation and application to financial markets. In *Proceedings of the Second ACM International Conference on AI in Finance (ICAIF '21)*. Association

for Computing Machinery, New York, NY, USA, 1–9. https://doi.org/10.1145/3490354.3494433
[3] Anonymous Author(s). forthcoming 2023. Generative AI for End-to-End Limit Order Book Modelling. (forthcoming 2023).
[4] Igor Babuschkin, Kate Baumli, Alison Bell, Surya Bhupatiraju, Jake Bruce, Peter Buchlovsky, David Budden, Trevor Cai, Aidan Clark, Ivo Danihelka, Antoine Dedieu, Claudio Fantacci, Jonathan Godwin, Chris Jones, Ross Hemsley, Tom Hennigan, Matteo Hessel, Shaobo Hou, Steven Kapturowski, Thomas Keck, Iurii Kemaev, Michael King, Markus Kunesch, Lena Martens, Hamza Merzic, Vladimir Mikulik, Tamara Norman, George Papamakarios, John Quan, Roman Ring, Francisco Ruiz, Alvaro Sanchez, Rosalia Schneider, Eren Sezener, Stephen Spencer, Srivatsan Srinivasan, Wojciech Stokowiec, Luyu Wang, Guangyao Zhou, and Fabio Viola. 2020. *The DeepMind JAX Ecosystem.* http://github.com/deepmind
[5] Peter Belcak, Jan-Peter Calliess, and Stefan Zohren. 2022. Fast Agent-Based Simulation Framework with Applications to Reinforcement Learning and the Study of Trading Latency Effects. In *Multi-Agent-Based Simulation XXII (Lecture Notes in Computer Science)*, Koen H. Van Dam and Nicolas Verstaevel (Eds.). Springer International Publishing, Cham, 42–56. https://doi.org/10.1007/978-3-030-94548-0_4
[6] Dimitris Bertsimas and Andrew W. Lo. 1998. Optimal control of execution costs. *Journal of Financial Markets* 1, 1 (April 1998), 1–50. https://doi.org/10.1016/S1386-4181(97)00012-8
[7] Taweh Beysolow II and Taweh Beysolow II. 2019. Market making via reinforcement learning. *Applied Reinforcement Learning with Python: With OpenAI Gym, Tensorflow, and Keras* (2019), 77–94.
[8] Jean-Philippe Bouchaud, Julius Bonart, Jonathan Donier, and Martin Gould. 2018. *Trades, Quotes and Prices: Financial Markets Under the Microscope.* Cambridge University Press. Google-Books-ID: dPRQDwAAQBAJ.
[9] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. 2018. *JAX: composable transformations of Python+NumPy programs.* http://github.com/google/jax
[10] David Byrd, Maria Hybinette, and Tucker Hybinette Balch. 2020. ABIDES: Towards High-Fidelity Multi-Agent Market Simulation. In *Proceedings of the 2020 ACM SIGSIM Conference on Principles of Advanced Discrete Simulation (SIGSIM-PADS '20)*. Association for Computing Machinery, New York, NY, USA, 11–22. https://doi.org/10.1145/3384441.3395986
[11] Philippe Casgrain. 2020. *LimitOrderBook.jl.* https://github.com/p-casgrain/LimitOrderBook.jl
[12] Rama Cont, Sasha Stoikov, and Rishi Talreja. 2010. A Stochastic Model for Order Book Dynamics. *Operations Research* 58, 3 (June 2010), 549–563. https://doi.org/10.1287/opre.1090.0780
[13] Kevin Dabérius, Elvin Granat, and Patrik Karlsson. 2019. Deep execution-value and policy based reinforcement learning for trading and beating market benchmarks. *Available at SSRN 3374766* (2019).
[14] Jin Fang, Jiacheng Weng, Yi Xiang, and Xinwen Zhang. 2022. Imitate then Transcend: Multi-Agent Optimal Execution with Dual-Window Denoise PPO. *arXiv preprint arXiv:2206.10736* (2022).
[15] Sumitra Ganesh, Nelson Vadori, Mengda Xu, Hua Zheng, Prashant Reddy, and Manuela Veloso. 2019. Reinforcement learning for market making in a multi-agent dealer market. *arXiv preprint arXiv:1911.05892* (2019).
[16] Google. 23. XLA: Optimizing Compiler for machine learning  :  tensorflow. https://www.tensorflow.org/xla
[17] Martin D. Gould, Mason A. Porter, Stacy Williams, Mark McDonald, Daniel J. Fenn, and Sam D. Howison. 2013. Limit order books. *Quantitative Finance* 13, 11 (Nov. 2013), 1709–1742. https://doi.org/10.1080/14697688.2013.803148 Publisher: Routledge _eprint: https://doi.org/10.1080/14697688.2013.803148.
[18] Ruihong Huang and Tomas Polak. 2011. LOBSTER: Limit Order Book Reconstruction System. https://doi.org/10.2139/ssrn.1977207
[19] Joseph Jerome, Gregory Palmer, and Rahul Savani. 2022. Market Making with Scaled Beta Policies. In *Proceedings of the Third ACM International Conference on AI in Finance*. ACM, New York NY USA, 214–222. https://doi.org/10.1145/3533271.3561745
[20] Joseph Jerome, Leandro Sanchez-Betancourt, Rahul Savani, and Martin Herdegen. 2022. Model-based gym environments for limit order book trading. http://arxiv.org/abs/2209.07823 arXiv:2209.07823 [cs, q-fin].
[21] Michäel Karpe, Jin Fang, Zhongyao Ma, and Chen Wang. 2020. Multi-agent reinforcement learning in a realistic limit order book market simulation. In *Proceedings of the First ACM International Conference on AI in Finance*. ACM.
[22] Diederik Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *International Conference on Learning Representations (ICLR)*. San Diega, CA, USA.
[23] Dmitrii Kochkov, Jamie A. Smith, Ayya Alieva, Qing Wang, Michael P. Brenner, and Stephan Hoyer. 2021. Machine learning–accelerated computational fluid dynamics. *Proceedings of the National Academy of Sciences* 118, 21 (May 2021), e2101784118. https://doi.org/10.1073/pnas.2101784118 Publisher: Proceedings of the National Academy of Sciences.

[24] Albert S. Kyle. 1985. Continuous Auctions and Insider Trading. *Econometrica* 53, 6 (1985), 1315–1335. https://doi.org/10.2307/1913210 Publisher: [Wiley, Econometric Society].

[25] Robert Tjarko Lange. 2022. *gymnax: A JAX-based Reinforcement Learning Environment Library*. http://github.com/RobertTLange/gymnax

[26] Chris Lu, Jakub Kuba, Alistair Letcher, Luke Metz, Christian Schroeder de Witt, and Jakob Foerster. 2022. Discovered policy optimisation. *Advances in Neural Information Processing Systems* 35 (2022), 16455–16468.

[27] Yuriy Nevmyvaka, Yi Feng, and Michael Kearns. 2006. Reinforcement learning for optimized trade execution. In *Proceedings of the 23rd international conference on Machine learning*. 673–680.

[28] Brian Ning, Franco Ho Ting Lin, and Sebastian Jaimungal. 2021. Double deep Q-learning for optimal execution. *Applied Mathematical Finance* 28, 4 (2021), 361–380.

[29] Anna A. Obizhaeva and Jiang Wang. 2013. Optimal trading strategy and supply/demand dynamics. *Journal of Financial Markets* 16, 1 (Feb. 2013), 1–32. https://doi.org/10.1016/j.finmar.2012.09.001

[30] James Paulin. 2019. Understanding flash crash contagion and systemic risk: a calibrated agent-based approach. https://ora.ox.ac.uk/objects/uuid:929fa3fe-4e5f-4cef-ad9f-03eb40110818

[31] Arnau Quera-Bofarull, Joel Dyer, Anisoara Calinescu, and Michael Wooldridge. 2023. Some challenges of calibrating differentiable agent-based models. http://arxiv.org/abs/2307.01085 arXiv:2307.01085 [cs, q-fin, stat].

[32] Samuel S. Schoenholz and Ekin D. Cubuk. 2020. JAX M.D. A Framework for Differentiable Physics. In *Advances in Neural Information Processing Systems*, Vol. 33. Curran Associates, Inc. https://papers.nips.cc/paper/2020/file/83d3d4b6c9579515e1679aca8cbc8033-Paper.pdf

[33] Maarten P. Scholl, Anisoara Calinescu, and J. Doyne Farmer. 2021. How market ecology explains market malfunction. *Proceedings of the National Academy of Sciences* 118, 26 (June 2021), e2015574118. https://doi.org/10.1073/pnas.2015574118 Publisher: Proceedings of the National Academy of Sciences.

[34] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal Policy Optimization Algorithms. *arXiv preprint arXiv:1707.06347* (2017).

[35] Hado Van Hasselt, Arthur Guez, and David Silver. 2016. Deep reinforcement learning with double Q-learning. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 30.

[36] Svitlana Vyetrenko, David Byrd, Nick Petosa, Mahmoud Mahfouz, Danial Dervovic, Manuela Veloso, and Tucker Balch. 2020. Get real: realism metrics for robust limit order book market simulations. In *Proceedings of the First ACM International Conference on AI in Finance (ICAIF '20)*. Association for Computing Machinery, New York, NY, USA, 1–8. https://doi.org/10.1145/3383455.3422561