**⟨§⟩ ChatGPT**

# Micro-Structure Sniper: Capturing the Spread Using LOB Dynamics

## 1. Theoretical Framework & Financial Logic

### Order Book Imbalance (OBI) and Short-Term Price Pressure

**Order Book Imbalance (OBI)** is a metric that quantifies the difference in buy vs. sell volume at the top of the order book. A common definition (for best bid volume $Q_b$ and best ask volume $Q_a$) is:

$$\text{OBI} = \frac{Q_b - Q_a}{Q_b + Q_a},$$

which ranges from -1 to 1 [1] . An OBI near +1 indicates the bid side is much larger than the ask side (heavy buying interest), while -1 means the ask side dominates. This imbalance has a **predictive correlation with short-term price movements**. In fact, it's considered the "worst-kept secret" of high-frequency trading that volume imbalance strongly predicts the next price move [1] [2] . Empirical studies show that when OBI is close to +1 (excess bid liquidity), the mid-price is likely to **jump up**, and when OBI is near -1, the mid-price often **drops down** shortly after [2] . In essence, immediate price pressure is skewed toward the side with less liquidity – the market tends to move where there are fewer orders (since that side can be **"consumed"** more easily by aggressive trades) [3] [4] . Traders monitor OBI as an early signal: a lopsided book (e.g. bids far outweigh asks) suggests short-term upward pressure, whereas a surplus of asks signals downward pressure.

### Micro-Price vs. Mid-Price – Volume-Weighted Fair Value

The **mid-price** (midpoint of best bid and ask) is the simplest "fair value" estimator: $\text{Mid} = (P_b + P_a)/2$. However, mid-price is *naive* – it ignores the volume on each side and often oscillates (bid-ask bounce) even if true value hasn't changed [5] [6] . To incorporate order book information, practitioners use **volume-weighted prices**. A first improvement is the **weighted mid-price**: for example,

$$P_w = \frac{Q_b}{Q_b + Q_a} P_a \; + \; \frac{Q_a}{Q_b + Q_a} P_b,$$

which leans the price toward the side with greater volume [7] . If more volume is on the bid (high $Q_b$), $P_w$ shifts upward toward the ask price, reflecting that buying pressure. Conversely, heavy ask volume shifts it toward the bid. This weighted-mid is closer to the eventual "efficient" price than the raw mid-price and responds to imbalance intuitively [7] .

**Stoikov's Micro-Price.** Building on this idea, Sasha Stoikov formally defined the **micro-price** as a more rigorous fair value estimator that **adjusts the mid-price for both the spread and the current imbalance** [8] . In essence, the micro-price $P_{\micro}$ can be expressed as:

$$P_{\micro} = \text{Mid} \; + \; f(\text{spread}, Q_b, Q_a),$$

where $f$ is an adjustment that raises or lowers the mid based on the bid-ask spread and the volume imbalance [9] . Intuitively, if the bid side is stronger (larger $Q_b$), $P_{\micro}$ will sit closer to the ask (above the mid), and if the ask side dominates, $P_{\micro}$ will be below the mid. In Stoikov's 2017 paper, the micro-price is defined as the **limit of a sequence of conditional expected mid-prices**, given the order book state [8] . By construction, this makes the micro-price a **martingale** – its future expected value (given current information) equals its current value [10] . In other words, $E[P_{\micro}(t+\Delta t) \mid \text{book at }t] = P_{\micro}(t)$, implying no predictable drift under the information in the LOB. This property holds because the micro-price is essentially the *fair price* of the asset *conditional on the order book state* [10] . It encapsulates the best estimate of true value at that moment, so any subsequent price change is, by definition, unpredictable (a random martingale increment). In practice, Stoikov's micro-price can be estimated in real time and has been shown empirically to **predict short-term price moves better than either the mid-price or simple weighted mid-price** [8] . It's also less noisy, since it filters out the transient bid-ask bouncing by incorporating imbalance information [11] .

*Why is it considered a martingale?* – Because it's derived as an expectation. Stoikov models the order book as a Markov process where the mid-price can move up, down, or stay, and imbalance influences those transition probabilities [10] . The micro-price comes out as the **expected mid-price in the "distant future" given the current state** (spread and imbalance), essentially an efficient price. Since an expectation of future price (under a fair model) is by definition equal to the current micro-price, it has no systematic bias in its evolution (no alpha left). This aligns with the efficient market idea: after adjusting for known information (spread and volumes), the resulting price should follow a martingale. Any deviation between the micro-price and the current market price suggests a temporary imbalance that will likely resolve as the market moves.

## "Fade vs. Chase": When to Provide Liquidity vs. Take It

The **Micro-Structure Sniper** strategy uses a game-theoretic logic to decide when to **provide liquidity** (place passive limit orders, a.k.a. **fade** the current move) and when to **take liquidity** (send marketable orders, a.k.a. **chase** the price). The core signal for this decision is the divergence between the micro-price and the actual mid-price. Think of the micro-price as the *true* price pressure indicator and the mid-price as the current market consensus. Any gap between them implies a potential trading edge:

- If **micro-price is significantly higher than the mid-price**, the book is imbalanced to the buy side and suggests *bullish* pressure. Here the strategy might "**chase**" – i.e. take liquidity by *buying* aggressively (e.g. hit the ask) because the current mid is arguably too low. In this case, crossing the spread to buy ensures you get in before the price rises to where the micro-price indicates it should be. After your buy, as the mid-price catches up to the micro-price (rising toward fair value), you can hopefully sell later at a profit. Chasing is essentially momentum trading on microstructure signals: you go with the pressure. The upside of chasing is capturing an imminent move (ensuring execution), but the downside is paying the spread/fees and risking that the signal was false or already decayed.

- If **micro-price is notably below the mid-price**, it indicates *bearish* pressure (excess sell volume), so the strategy might chase by *selling* (take the bid) before a drop. This is the symmetric case on the short side.

- If **micro-price and mid-price are roughly in line** (or the divergence is small), the market is near equilibrium. Here an HFT market-making approach leans toward **fading**, i.e. **providing liquidity** on the side of the spread that anticipates reversion. For example, if the mid and micro-price are equal or only slightly off, any small price flickers might mean-revert. A sniper might place passive orders on both sides (earning the spread) or on the side suggested by micro-price. **Fading** means you *fade the move* – e.g. if the price just inched up but micro-price doesn't confirm a strong imbalance, you might place sell limits (providing liquidity at the ask), expecting the price uptick to revert back down. By providing a limit order, you earn the bid-ask spread if filled, essentially betting that the apparent move will **"fade"** away rather than continue.

- In **game-theoretic terms**, each HFT must anticipate others' actions. If an imbalance is extreme, many traders will try to "sniper" the stale price – meaning if you don't chase, someone else will, and the price will move quickly. Thus, when the micro-price divergence is large (clear imbalance), the equilibrium tends to be *taking* liquidity (everyone races to grab the mispricing). Conversely, when imbalance is mild or ambiguous, aggressive action by others is less certain; providing liquidity (fading) might be optimal to earn spread income while the market stabilizes. The strategy continuously weighs: **Is the expected short-term move (indicated by micro-price) larger than the spread and fees?** If yes, *take* the liquidity (because paying the spread is justified by the move). If not, it's better to *provide* liquidity and earn the spread, assuming mean reversion.

In practice, a Micro-Structure Sniper will rapidly flip between these modes. For instance, suppose best bid is \$100 (size 500) and best ask is \$100.05 (size 100). If a large buy order consumes the 100 shares at \$100.05, the mid moves up, but now the order book might show \$100.05 bid vs \$100.10 ask (new mid \$100.075). If the micro-price (accounting for the new imbalance) is still below \$100.075, it signals that the move might have overshot – a fade opportunity. The sniper could place a sell limit at \$100.10 (providing liquidity at the new ask) expecting the price to dip back. On the other hand, if that initial state had shown \$100 bid vs \$100.05 ask with *5000* on the bid and only 100 on the ask (huge imbalance), the micro-price would have been very close to \$100.05 (or beyond). A sniper seeing that would **not** wait – it would immediately lift the ask (buy at \$100.05) to capture the likely imminent jump, effectively **chasing** the microprice signal. In summary, the strategy **provides liquidity (passive)** when the order book suggests mean reversion or only slight pressure, and **takes liquidity (active)** when a strong imbalance implies a price will move before a passive order can profit. This **make-or-take decision** is critical in HFT market making: providing liquidity earns the spread but risks adverse selection if the market moves against you, while taking liquidity secures entry ahead of a move but incurs cost. The Micro-Structure Sniper continuously plays this game-theoretic balancing act, aiming to **"snipe" tiny inefficiencies**: fading when it thinks others will hesitate, and pouncing (chasing) when it predicts a rapid price shift.

## 2. Developer Implementation Guide

### Low-Latency Execution Engine Architecture

Designing a high-frequency trading engine requires a focus on **low latency** and high throughput. The system should be event-driven, reacting to each **tick-level L2 update** (Level 2 market data, i.e. order book changes) as quickly as possible. A typical architecture includes the following components:

1. **Market Data Feed Handler** – This module connects to the exchange's data feed (using a binary protocol or API) and receives order book updates in real time. It should parse messages (new orders,

cancels, trades) and update an in-memory order book. To handle thousands of updates per second, implement this as a non-blocking, asynchronous loop. For example, in Python one could use `asyncio` to read from a socket or message queue without stalling other tasks. In C++ you might use epoll or a lock-free queue mechanism to ingest data. The key is that each tick (e.g., a change in best bid/ask or depth) triggers an event in the engine's event loop. **Minimize any I/O or compute blocking** in this loop – e.g., avoid heavy computations or disk access here. Many HFT engines run a single-threaded event loop for market data to preserve order and avoid context-switch overhead (each event is processed in sequence). If using multiple threads, careful coordination is needed to avoid race conditions in order book state.

2. **Internal Order Book & Signal Calculator** – The engine maintains a live snapshot of the **limit order book (LOB)** (at least the top levels needed for the strategy). On each tick event, it **updates the LOB state** (e.g., update best bid price/size when a level-1 quote changes). Immediately after updating, it **calculates signals** like OBI and micro-price. This calculation is very fast – just arithmetic on a few numbers – but it must be done carefully for streaming data. For example, if using Python, you could predefine functions to compute OBI = $(Q_b - Q_a)/(Q_b+Q_a)$ and micro-price = $P_b + \frac{Q_a}{Q_a+Q_b}(P_a - P_b)$ (which simplifies to the weighted average formula) on each update. With Python's `asyncio`, you might implement this as a callback or coroutine that runs on every book change event. If performance becomes an issue in Python, you can use C++ extensions or numpy for vectorized operations (however, since these are simple scalar calculations per tick, the bottleneck is more likely the Python interpreter overhead per message). Some practitioners interface Python with C++ for critical paths – for instance, using C++ to handle the feed parsing and signal computation, then passing the results to Python for decision logic. This hybrid approach leverages C++ speed while keeping high-level strategy logic in Python for flexibility. **Latency considerations:** Aim to process each tick-to-signal in microseconds. In modern HFT, "tick-to-trade" latency (time from incoming tick to outbound order) can be a few microseconds in optimized C++ systems [12]. A pure Python event loop will be slower (likely milliseconds), which may be acceptable for intraday strategies that aren't latency-arbitrage, but it's still crucial to optimize (e.g., avoid unnecessary allocations, use asyncio's event loop efficiently). The system should handle thousands of events per second without backlog.

3. **Decision Logic (Strategy Core)** – This is the brain that implements the **'Fade vs. Chase' logic** described above. After each tick's signals are updated, the strategy module decides whether to place an order, cancel an order, or do nothing. This logic can run as part of the event loop (synchronously after each update) or on a separate thread that continuously checks the latest state. In a low-latency system, it's common to integrate it directly in the event loop so there's virtually no delay between data update and decision. The logic will compare the micro-price vs. mid-price, evaluate current inventory and outstanding orders, and then formulate an action: e.g., "Microprice is $100.05 vs mid $100 – significant divergence; send a market buy to chase." Risk checks should be embedded here (don't overshoot size limits, etc.). It might help to maintain a small state machine: e.g., *State:* currently have a passive order resting or not; *Event:* new imbalance signal says price likely to move up; *Action:* if no order resting, send a buy; if an order was resting on the wrong side, cancel it, etc.

4. **Order Execution Handler** – This module interfaces with the exchange's trading API/protocol (e.g. FIX, OUCH, REST, etc.) to send orders and receive order acknowledgments and fills. Key responsibilities include:

5. **Low-latency order submission:** Use the fastest possible method (for many exchanges, a binary protocol or colocated direct connection is necessary – e.g., FIX might be too slow for true HFT). The handler should format the order message (buy/sell, price, size, etc.) and send it immediately when instructed by the strategy logic. Non-blocking I/O is important here as well (in Python, use asynchronous sockets or a separate thread to avoid delaying the market data loop).

6. **Order Queue Management:** When the strategy places passive limit orders (to provide liquidity), the execution handler should track them in an **order queue** structure. You'll want to assign each order an ID and maintain its status (open, partially filled, filled, or canceled). If multiple orders are outstanding (perhaps one on the bid, one on the ask), keep them in a list or dictionary. The handler processes exchange **acknowledgments** – e.g., an order accepted message, fills, cancel confirmations – updating the status. This is critical for knowing your position and remaining order sizes in real time.

7. **Cancellations & Amendments:** The strategy often needs to cancel orders very quickly when market conditions change. For instance, if you were fading (providing liquidity) on the ask but suddenly a large buy imbalance appears (microprice jumps up), you risk being adversely selected (your sell order is now very likely to get picked off at an unfavorable time). The strategy should immediately send a cancel for that resting order. The execution handler should prioritize cancel messages and handle the logic: e.g., throttle if too many cancel attempts (to avoid exchange rejection or undue bandwidth). It should also handle order replacements (modify price) if the strategy chooses to reprice instead of outright cancel. In a low-latency engine, **cancel/replace speed** is as important as order send speed – stale orders left in the book for even a few milliseconds can be "latency arbitraged" by faster traders.

8. **Latency and Risk Management:** *Latency arbitrage risk* refers to the danger that due to latency, your view of the market is slightly outdated and faster players can trade against you if you post a stale quote. To mitigate this, the entire loop – from tick to decision to order/cancel – must be as fast as possible. Additionally, the strategy might include safeguards like **quote timeout**: e.g., cancel any resting order that sits more than X milliseconds without fill, since old quotes become risky. The execution handler can enforce these time-based cancels. It's also wise to implement basic risk controls: e.g., a kill-switch if too many orders execute in a short time or if inventory goes out of bounds (to avoid runaway losses due to a bug or rare event).

**Event Loop Considerations:** The engine's event loop should likely be single-threaded for core logic to avoid synchronization delays, but you can employ a multi-thread or multi-process design for non-critical tasks. For example, logging and analytics can be offloaded so they don't slow the main loop (e.g., log to memory and have another thread periodically write to disk). In C++ you might use lock-free ring buffers to pass data between threads. In Python, one could use `asyncio` with `create_task` for things like sending orders asynchronously while the main loop continues processing new ticks. The overall flow is: **market data event -> update state -> compute signals -> decide -> send order/cancel -> back to waiting for next event**, all within a few milliseconds or less.

**Technology stack:** Favor technologies that minimize latency: e.g., use kernel-bypass network drivers if possible (Solarflare/OpenOnload or DPDK) to cut down network IO latency. If deploying on a cloud VPS, ensure it's in the same region as the exchange's servers to reduce round-trip time. Many HFT systems eventually deploy on **co-located servers** or VPS instances in data centers near the exchange matching engine. A recommended development path is to **first build and test locally**, then deploy to a low-latency cloud/VPS for production. Using free/public data during development is helpful – for example, one could use the IEX exchange's public TOPS feed (which provides top-of-book quotes for U.S. stocks) to test the strategy logic in real-time without fees, before moving to a paid data feed [13] . After local validation, migrate

the code to a VPS close to the exchange (for instance, a New York-region server for U.S. equities) to improve network latency from, say, tens of milliseconds down to 1–2 ms or less. This two-step approach (local then cloud) ensures you catch bugs and refine logic in a safe environment, then achieve the needed speed for live trading.

### Real-Time Calculation: OBI & Micro-Price in Streaming Data

Calculating **Order Book Imbalance and Micro-Price in real time** is straightforward mathematically, but the implementation must be highly optimized. Here's how to do it:

- **Order Book Data Structure:** Maintain variables for `best_bid_price`, `best_bid_size`, `best_ask_price`, `best_ask_size`. Each L2 update from the feed will indicate if any of these changed. For example, if an update says a new order at the best ask or a cancellation at the bid, update the respective price or size. Often, using an array or dictionary for the full depth can allow computing deeper imbalance (e.g., including more levels), but for the micro-price we typically focus on level-1 data. In Python, a simple approach is to update these four variables on each tick event. In C++, you might keep a struct for best quotes and update fields.

- **OBI Computation:** Once the best bid/ask sizes are updated, compute `OBI = (best_bid_size - best_ask_size) / (best_bid_size + best_ask_size)` (taking care to handle division by zero if order book is empty, which in normal operation won't happen at top of book). This gives a number between -1 and 1. In code, this is a few CPU operations – negligible in cost. The main consideration is doing it in-line with the feed handler without introducing Python-level looping over many entries (since we only use top-of-book, it's constant time). If using asyncio, you might have something like:

```python
async def on_order_book_update(msg):
    best_bid_price = msg.best_bid_price
    best_bid_size = msg.best_bid_size
    best_ask_price = msg.best_ask_price
    best_ask_size = msg.best_ask_size
    # Compute imbalance
    if best_bid_size + best_ask_size > 0:
        obi = (best_bid_size - best_ask_size) / (best_bid_size + best_ask_size)
    else:
        obi = 0  # (should not happen for non-empty book)
    # Compute micro-price
    micro_price = best_bid_price + (best_ask_price - best_bid_price) *
 (best_bid_size / float(best_bid_size + best_ask_size))
    # Now use these values in strategy logic...
```

This pseudocode shows updating state and computing signals in one async callback. In practice, you might separate the concerns (update state first, then compute signals).

- **Micro-Price Computation:** Using the formula above, `micro_price = (Q_b/(Q_a+Q_b)) * P_a + (Q_a/(Q_a+Q_b)) * P_b`. This is equivalent to `mid_price + (spread * (Q_b - Q_a) / (2*(Q_a+Q_b)))` if using the -1 to +1 imbalance definition. Any of these forms yield the same

result [7] . The code should use fast arithmetic; avoid any unnecessary conversions or Python loops. If using libraries, note that simple operations are usually fine without NumPy, but you could use NumPy if you were computing many instruments in parallel. Given this strategy likely focuses on a single instrument at a time (or a small handful), plain arithmetic is sufficient. In C++, these computations would be done in a few nanoseconds; in Python, a few microseconds.

- **Asynchronous Processing:** While the above runs synchronously on each tick, one could also pipeline it. For instance, if the feed is very heavy, you might accumulate a batch of updates and process the latest one, skipping intermediate ones if they come faster than you can handle (this is a technique if you can tolerate some misses – often HFT strategies cannot, but if your trading logic only needs the latest state, you might drop some updates in extremely bursty scenarios). Python's asyncio event loop will queue callbacks if it receives another tick while still processing the previous; ensuring your processing is super fast prevents long queues. If you find Python is not keeping up, consider using a **C++ extension or Cython** for the hot path of updating the order book and computing OBI/micro-price. That way, the heavy tick processing runs at C++ speed, and only the decision (which might be slightly slower) runs in Python. This kind of binding is common – for example, using C++ for a matching engine or feed handler and exposing it via Python bindings (pybind11, CFFI, etc.) to a Python strategy.

- **Latency Metrics:** It's good practice to monitor how long each tick processing takes ("tick-to-trade latency"). You can timestamp when a market data message arrives and when your order is sent out. For a Python-based system, you might see, say, 1ms from tick to order on average (depending on hardware), whereas a fully optimized C++ could be under 10μs [12] . If your strategy is truly high-frequency, you'll want to continuously optimize this. If it's slightly slower (still intraday but not competing on the microsecond arms race), a few milliseconds may be acceptable. Still, the faster your loop, the better your **queue position** in reacting to new information.

## Execution Handler: Order Management & Latency Considerations

The **Execution Handler** is where your strategy's intentions meet the real market. A robust implementation will manage orders lifecycle and guard against latency pitfalls:

- **Order Queues and IDs:** Each order sent should be stored with a unique client ID. The exchange will typically echo an order ID or you may generate one. Track attributes: price, size, side, time sent, etc. Use a data structure that allows quick lookup by ID (for updates/cancels). In Python, a dict mapping order_id -> order_object is handy. In C++, perhaps an `unordered_map` or arrays if using sequential IDs.

- **Matching Engine Interaction:** Because this strategy often places passive orders to capture spread, it is sensitive to **queue position** on the exchange's order book. The execution handler can't control matching (that's on exchange side), but it **must manage the outcome**: when your order gets filled partially or fully, you'll get a fill report. Upon fill, update your inventory and P&L, and consider removing or reducing the order from your tracking. If only partial, update remaining size. All this should happen very fast – ideally handle the exchange messages (fills, cancels, etc.) with the same event-driven approach as market data.

- **Latency Arbitrage Risk Management:** One risk is that your passive orders get "picked off" when you're at a disadvantage. For example, if you posted a buy at the bid and an adverse event (like a large sell order) comes, you might get filled just as the price is about to drop further. This is known as being adversely selected. To mitigate this, the strategy might employ **tight cancellation rules**:

- If the micro-price signal flips or the market moves against your order, send a cancel immediately. The code should handle a cancel similarly to a new order – with low latency. Many HFT strategies use **"cancel/replace"** frequently, updating quotes several times a second (or much more) as needed.
- Use **short time-to-live** on orders: for instance, cancel any order that's been resting for more than, say, 50 milliseconds without fill. The rationale is that if it hasn't filled in 50ms, conditions might have changed or faster traders might know something you don't, so better to remove it rather than sit exposed. Implementing this requires a timing mechanism – e.g., record the timestamp when order was placed, and on each new tick or a periodic timer, check open orders for expiry. You could integrate this into the event loop by scheduling a callback (with `asyncio.call_later` in Python, for example) to cancel the order after X milliseconds if it's still open.

- Be aware of **exchange-specific latency**: There's latency from you to the exchange and back. In a cloud VPS scenario, maybe ~1ms each way, plus matching engine internal latency. This means after you decide to cancel, there's a window where you might still get filled before the cancel is processed. Your software should be prepared to handle out-of-order events like a fill confirmation for an order you just sent a cancel for. Typically, the exchange will indicate if a fill happened before the cancel was effective. Your handler logic should thus process **fills before cancels** if they come in that order, to avoid erroneously thinking an order was gone while it actually executed.

- **Concurrency and Atomicity:** If you run multi-threaded (e.g., a separate thread for sending orders), be careful to avoid race conditions updating the same order structures. Using thread-safe queues for communication is one approach (have the strategy thread put an "order request" in a queue that the network thread picks up and sends). With `asyncio`, you might not need threads at all – you can await the network send without blocking the main loop. Just ensure not to mix up messages; use locks or atomic swaps if needed when updating shared state like position.

- **Testing the Execution Handler:** It's wise to test the order logic on a **simulated exchange** or a paper trading environment first, to ensure all the message handling is correct. During development, you can connect to a test exchange (many exchanges have a simulated environment) or use a very simple market (low volume stock) to send test orders. This will flush out issues in order tracking, cancellations, etc., without risking significant capital.

- **Favoring Free/Public APIs:** For initial development, you can use free APIs to practice order placement logic. For example, some exchanges (like certain crypto venues or IEX for equities) offer sandbox or testnet APIs where you can send orders without real money. Although we focus on non-crypto here, the concept stands: use any publicly available *trading simulator* or API to refine your execution handler before going live. This ties back to the earlier point of first doing it local – perhaps simulate on your machine with a fake matching engine – then move to a real environment.

In summary, the developer's goal is to **streamline the event loop** such that from the moment a market data tick arrives, all computations and decisions happen quickly and an order or cancel (if needed) is sent out almost immediately. By structuring the system into modular components (feed handler, calc, decision,

execution) and using asynchronous, non-blocking patterns, one can achieve the low latency required for intraday millisecond-level strategies. Modern HFT engines even measure latency in microseconds, and while a Python-based approach might be a bit slower, careful optimization and eventually using C++ for the hottest paths can get the loop close to the "wire-to-wire" speeds needed [12] . Always remember: every millisecond (or microsecond) of delay in reacting could mean someone else beats you to the trade!

## 3. Validation and Backtesting

### Simulating an Exchange Environment for Backtesting

Before deploying a microstructure strategy live, it's crucial to **backtest it in a realistic simulated environment**. Traditional backtesting (with say, 1-minute bars) is insufficient for HFT – we need a **tick-by-tick, limit order book simulation**. There are two main approaches to simulate the exchange:

- **Historical Market Data Replay:** Obtain historical Level-2 (order book) data or at least top-of-book quote data with timestamps. In a replay, you step through each historical update and "play" it into your strategy as if it were happening live. When your strategy decides to place an order, you insert that order into the simulated order book and then continue processing historical events to see if/ when your order would have filled. This requires writing a **matching engine simulator** that sits between your strategy and the data. Essentially, the simulator maintains an order book and matches your orders against the historical tape. Since your orders didn't actually exist in history, we must simulate what *could* have happened. For example, if your strategy places a buy limit at a price that *was* traded through in history, you likely would have been filled – but it depends on where you were in the queue.

- **Stochastic Order Book Simulation:** Alternatively, one can simulate a synthetic order flow using statistical models. Academic research has many **agent-based models** or Monte Carlo frameworks for LOB dynamics [14] [15] . You could simulate random order arrivals (using Poisson processes for trades, etc.) that approximate real market patterns, and embed your algorithm into that artificial market. This can be useful to test "what-if" scenarios or robustness. However, it's quite complex to get right, and most practitioners prefer to use actual historical data to ensure realism.

For practicality, **replaying real data** with a simulated matching engine is a common choice. Here's how to implement it:

1. **Data Collection:** Gather historical tick data. Ideally, you have full order-by-order data (often called **Level-3 data**). If not, Level-2 data (aggregate sizes at top N levels) can suffice. Free public data can be limited here; however, some exchanges (or academic sources) provide sample LOB data. One free option: use **public trades and quotes (TAQ)** data to reconstruct a top-of-book history. For example, the IEX exchange provides daily history of quotes/trades that you can use to replay one day's order book evolution. Always ensure the data's time resolution is high (millisecond or microsecond timestamps) for accuracy in sequencing events.

2. **Matching Engine Simulation:** In your backtest harness, implement a simplified matching engine:

3. Maintain an order book structure (buy and sell queues at each price). Initially populate it with the state from data (e.g., at 9:30:00, take the first snapshot).

4. Step through events in chronological order. Events might be: *limit order added, order canceled, market order (trade)*, etc. Apply each event to the book. This will update best bid/ask and volumes, just like in real life.

5. Allow your **strategy to interact**: At each event time (or at each simulation time tick), feed the updated market state to your strategy code. The strategy will output order actions (if any) – e.g., "place a buy limit at \$100" or "cancel my previous sell order". When the strategy sends an order, inject it into the simulated order book *with an appropriate position in the queue*. For example, if your order is a buy at \$100 and \$100 is currently the best bid price, you need to decide where in the bid queue your order would sit. If you have Level-3 data, you know how many shares are ahead of you at that price at that moment. If you only have Level-2 (aggregate), you might assume you join behind the existing volume.

6. **Order Execution Logic:** As you continue replaying events, whenever a trade event happens at a price where your order is resting, you need to determine if your order gets filled. This is the tricky part: it depends on your queue position and the trade size. For instance, if you were 5th in line at \$100 with 100 shares posted, and a sell market order comes that hits 200 shares at \$100, will you be filled? Possibly yes, if the 200 shares were enough to eat through at least 4 other orders ahead of you. To model this, track your **queue position and remaining volume at that price**. Each time there's an execution (market order) that hits that price, subtract the executed volume from the *front* of the queue. If your order is reached (your queue position goes to 0 or negative), that means your order was executed (partially or fully). Fill your order accordingly (and perhaps remove it if fully filled). If an order book event is a cancellation at your price level, that could move your order up in queue. You can model that as well – e.g., if 100 shares cancel in front of you, your queue position decreases by 100.

7. **Queue Position Modeling:** Because historical data often doesn't tell us exactly where our order *would* be, we use models. A **conservative approach** is to assume you're always at the **back of the queue** when you join. This means you only get filled if the entire displayed volume (at the time you placed the order) plus any additions before the price moves, were executed. This is a worst-case assumption (other cancellations don't help you; you're last in line). It avoids overly optimistic fills. Another approach is a **probabilistic queue model** [16] [17] : for example, assume that when some volume cancels at your level, there's a certain probability it was in front of you vs. behind you. Advanced simulators like **HftBacktest** offer models like *risk-averse* (assume all cancels are behind you – i.e., none of them move you up, which is very conservative) or *proportional/random* (assume cancels are equally likely to be anywhere in the queue) [17] [18] . These models essentially randomize your fill outcomes in a statistically reasonable way. For instance, a Probabilistic Queue Model might say: if 100 shares cancel at this price and you're currently, say, 300 shares from the front, maybe there's a 1/3 chance those cancels were ahead of you (so your position improves) [19] [20] . By calibrating such models (even using the formulas from Erik Rigtorp's method or others [13] [21] ), you can get realistic fill rates.

8. **Avoiding Unrealistic Results:** The goal of the above is to avoid the common backtest fallacy of assuming *any* touch of your price equals a fill. In a real order book, your limit order might not get executed even if the market trades through your price, if there was not enough volume to reach your spot in queue. By modeling queue depth and **fill probability**, you ensure the backtest doesn't overestimate profits. For example, if your strategy is to always sit on the bid and capture one tick, a

naive backtest might fill you on every upward tick and show huge profits. A realistic backtest will often show that many times your order would only be partially filled or not filled at all because other orders were ahead (especially in a crowded trade). Incorporating this can drastically reduce the hypothetical P&L – which is good, because it prevents false confidence.

9. **Performance of Backtest:** Simulating every tick and order can be slow, but since this is *offline*, you can afford to do a detailed simulation. Still, efficient data structures (like maintaining heaps or sorted lists for book levels) and vectorized operations where possible (maybe using pandas or numpy to handle time series) will help when running through millions of events. You might also consider event-driven backtest frameworks or libraries. For instance, **NautilusTrader** or **HftBacktest** are libraries that handle LOB simulation and allow you to plugin strategies [22] [18] . These can save time, as they come with built-in queue models and matching engine logic. Always test your backtest engine with known scenarios (for example, send a limit order in a simple up-move scenario and verify the fill logic does what you expect).

## Modeling Queue Position & Fill Probability

To emphasize, modeling the execution is just as important as modeling the signal. **Queue position** can be thought of as an additional state variable in your strategy. When backtesting, keep track of how many shares are ahead of your order at a given price. Some advanced points:

- **Order Acknowledgment Delay:** In live trading, when you send an order, there's a tiny delay before it actually appears in the order book (exchange processing time). In a fast market, by the time your order is live, the queue ahead could be different. A rigorous backtest might simulate that by **inserting your order with a slight delay** (e.g., if replaying with timestamps, maybe assume your order hits the book 1–2 ms after you decided to place it, which could affect your position if another order came in meanwhile).

- **Partial Fills:** Your order might get partially filled over time. The simulation should allow an order to remain in the book after a partial fill (with reduced size) and continue to possibly get more fills later if the price level is hit again. This mirrors reality where you often get filled in pieces.

- **Metric Tracking:** As you simulate, collect metrics like *fill ratio* (what percentage of your posted orders actually got filled), *hold time* (how long orders rest before fill or cancel), slippage vs. theoretical (did you capture the spread or did you have to cross it more often). These help validate if the strategy works as expected. For example, if you see in backtest that every time you placed a passive order you almost never got filled (maybe because the market moved away), you might be too slow or your model might need to be more aggressive (or it might reveal that capturing spread on that instrument is extremely competitive).

- **Free Public Data for Backtesting:** While high-quality tick data is often expensive, you can start with whatever is available. Some resources provide sample LOB data (e.g., the **LOBSTER dataset** for Nasdaq stocks offers detailed order book feeds for academic use, albeit not free beyond small samples). You could use a small sample to test your simulator. Alternatively, use **crypto exchange data** (even though we ignore crypto for strategy logic, their data is freely available and useful to test a simulator since the mechanics of order books are similar). For instance, you can get a day's worth of full order book updates from Binance or Coinbase for free and run your backtest on it – just

treating it as a generic order book. This can be done locally to validate your Micro-Structure Sniper implementation in a controlled way.

In conclusion, validation requires **level-2 precision**: simulate the matching as closely as possible, model your queue position either pessimistically or with realistic randomness, and ensure your strategy's profit in backtest comes from true edge (predicting moves or earning spread) and not from an artifact of assuming perfect fills. By doing this due diligence, you'll avoid the trap of strategies that look great on paper but fail in production.

**Real-Time Monitoring:** Even after backtesting, when you go live, treat the live trading initially as an extension of testing. Monitor the strategy's behavior: Is it getting the fill rates you expected? Is slippage within bounds? You might run the strategy at minimal size and compare live results to backtest predictions. Any large deviation might indicate that your simulation was optimistic (or some live factor wasn't accounted for). Then you can iteratively refine the model. High-frequency trading is an arms race, but with a solid theoretical foundation (OBI and microprice giving you an edge) and a carefully engineered execution system, the Micro-Structure Sniper strategy can systematically capture the bid-ask spread by smartly fading or chasing the microstructure signals, all while managing the subtle realities of order execution in a limit order book environment.

**Sources:** The concept of micro-price as a fair, volume-weighted price comes from Stoikov's research [8] [10]. The predictive power of order book imbalance for short-term price movement is well-documented in market microstructure literature [1] [2]. In implementing such strategies, practitioners emphasize minimizing tick-to-trade latency and accurately simulating matching engine mechanics for backtesting [13] [21]. By focusing on these aspects (signal edge and execution efficiency), one can build a robust HFT sniper that thrives on intraday microstructure inefficiencies.

---

[1] [2] [7] [11] arxiv.org
https://arxiv.org/pdf/2307.15599

[3] [4] Thoughts on High-Frequency Trading Strategies (4) | by FMZQuant | Medium
https://medium.com/@FMZQuant/thoughts-on-high-frequency-trading-strategies-4-21f817837c26

[5] [6] The Micro-Price
https://www.ma.imperial.ac.uk/~ajacquie/Gatheral60/Slides/Gatheral60%20-%20Stoikov.pdf

[8] [9] The Micro-Price: A High Frequency Estimator of Future Prices by Sasha Stoikov :: SSRN
https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2970694

[10] High Frequency Estimator of Future Prices — Micro-price paper & code walkthrough | by Mhfizt | Medium
https://medium.com/@mhfizt/high-frequency-estimator-of-future-prices-micro-price-paper-code-walkthrough-475adb98e91d?responsesOpen=true&sortBy=REVERSE_CHRON

[12] What is tick-to-trade latency? | Databento Microstructure Guide
https://databento.com/microstructure/tick-to-trade

[13] [19] [20] [21] Estimating Order Queue Position | Erik Rigtorp
https://rigtorp.se/2013/06/08/estimating-order-queue-position.html

14 15 16 17 18 22 Order Fill — hftbacktest

https://hftbacktest.readthedocs.io/en/py-v2.1.0/order_fill.html