

FX Market Making with Internal Liquidity

Alexander Barzykin¹, Robert Boyce^{1, 2}, and Eyal Neuman²

¹HSBC

²Department of Mathematics, Imperial College London

December 5, 2025

Abstract

As the FX markets continue to evolve, many institutions have started offering passive access to their internal liquidity pools. Market makers act as principal and have the opportunity to fill those orders as part of their risk management, or they may choose to adjust pricing to their external OTC franchise to facilitate the matching flow. It is, a priori, unclear how the strategies managing internal liquidity should depend on market conditions, the market maker’s risk appetite, and the placement algorithms deployed by participating clients. The market maker’s actions in the presence of passive orders are relevant not only for their own objectives, but also for those liquidity providers who have certain expectations of the execution speed. In this work, we investigate the optimal multi-objective strategy of a market maker with an option to take liquidity on an internal exchange, and draw important qualitative insights for real-world trading.

1 Introduction

Internalisation has long been a key component of efficient algorithmic execution in foreign exchange (FX) markets, primarily due to its reduced visibility and consequently minimal market impact. Traditional internalisation typically involves client-to-client matching, where a liquidity provider acts as a neutral intermediary. This mechanism, classified as Internalisation Type A by the Foreign Exchange Professionals Association (FXPA) [[FXPA guidance](#)], offers certain advantages but also inherent limitations, as matching opportunities tend to be scarce.

Since FX trading remains largely over-the-counter (OTC), interaction with OTC liquidity is expected to be the primary driver of internalisation. The FXPA defines Internalisation Type B as the offsetting of commercial flow by a liquidity provider. This naturally involves client algorithms trading against a market maker’s pricing stream, benefitting from ultra-low latency and potentially tighter spreads due to inventory skew.

Recently, several institutions have begun offering passive access to internal liquidity, either through conventional limit or pegged orders. Dynamic orders are typically pegged to an internally maintained fair reference price, allowing client algorithms to communicate in the “high-frequency language” of the market maker, without requiring high-frequency order management. Market makers may fill these orders to meet their risk management objectives; importantly, they can also adjust pricing in their OTC franchise to facilitate such fills. In this way, client algorithms can interact with deep OTC liquidity through the internalisation mechanism.

One might tacitly assume that market makers directly transfer pegged order quantities onto their pricing ladder at the same price and then immediately fill them upon receiving an opposing

trade. However, such a naïve approach would clearly be detrimental to the market maker’s risk management and would reduce potential P&L. In practice, market makers must instead solve an optimal market-making problem in the presence of an additional liquidity source, be it a single limit order or, more generally, a limit order book on an internal exchange.

Understanding the underlying mathematical formulation of this problem can improve the transparency of internalisation, in line with recent FXPA guidance. Orders on the internal exchange cannot be taken for granted: they may be cancelled, be of finite but unknown size, or follow an unpredictable strategy and thus may not always be available. Moreover, the market-making desk may prefer to fill client orders on the internal exchange, as doing so supports the firm’s client algo desk, which manages those orders, and ultimately provides better service to clients. The market-making desk therefore faces a complex and multifaceted trade-off.

The mathematical finance literature on market making originates from the model of Avellaneda and Stoikov [12], later solved in closed form by Guéant et al. [8]. Subsequent extensions include the linear-quadratic framework of Cartea et al [13], influential order effects in Cartea et al. [1], continuous hedging in Barzykin et al. [5], and competition among market makers in Boyce et al. [11]. Interactions between market makers and clients have also been explored, notably through the game-theoretic approach of Cartea and Sánchez-Betancourt [4], while internalisation has been examined in various contexts, including the FX market in Butz and Oomen [6]. Since Avellaneda and Stoikov [12], financial markets and the role of market makers have evolved substantially. A key recent innovation is the emergence of internal exchanges, enabling clients to provide liquidity directly to market makers. This remains largely unexplored, with the exception of Morimoto [3], who study optimal execution under unlimited internal liquidity with price impact. Relatedly, passive impact has attracted renewed attention, most notably in Chahdi et al. [9].

In this work, we present the first quantitative investigation of internal exchange management from the market maker’s perspective. We develop a model in which the market maker continuously streams a price ladder of multiple sizes to external, liquidity-taking clients while optimally timing trades with liquidity-providing clients on an internal exchange. The market maker aims to maximise P&L while remaining averse to large inventory positions and unfilled internal orders. Internal exchange orders are transient, that is they may be cancelled, executed, or reappear later, capturing realistic client execution patterns such as iceberg, TWAP, or full-amount strategies. From a mathematical perspective, our formulation of the market maker’s problem leads to a stochastic control problem over external prices, consistent with classical market-making frameworks [12, 7], combined with repeated optimal stopping decisions for the timing of internal trades, as in the optimal execution framework of [2, 3]. The resulting Hamilton–Jacobi–Bellman quasi-variational inequality (HJBQVI) is then solved numerically.

We derive practical insights with direct relevance to real-world trading. In particular, we demonstrate that there exists an execution threshold, i.e. the inventory level beyond which the market maker will instantaneously fill the limit order. Otherwise, the market maker will adjust pricing to external OTC clients accelerating the move towards the execution threshold, thus facilitating the fill of the limit order. The degree of price skew and the execution threshold level depend on the market maker’s risk aversion, limit order depth and expected placement strategy, as well as on flow facilitation initiatives. Importantly, the optimal strategy significantly outperforms a naïve benchmark strategy that directly incorporates internal exchange orders into the OTC pricing ladder.

2 The model

Let $T > 0$ denote the trading period. We fix a filtered probability space $(\Omega, \mathcal{F}, \mathbb{F} = \{\mathcal{F}_t\}_{t \in [0, T]}, \mathbb{P})$ satisfying the usual hypothesis. Let $(S_t)_{t \geq 0}$ be the price process of a risky asset such that $S_t = S_0 + \sigma W_t$, where W is a standard Brownian motion and S_0, σ are positive constants. We consider a market maker who continuously provides liquidity on both sides of the order book, quoting bid and ask prices $(S^{b,z}, S^{a,z})$ that are streamed and adjusted according to the clients' order sizes $\mathcal{Z} = \{z_k, 1 \leq k \leq K\}$:

$$S_t^{b,z} = S_t - \delta_t^{b,z} \quad \text{and} \quad S_t^{a,z} = S_t + \delta_t^{a,z}, \quad 0 \leq t \leq T, \quad z \in \mathcal{Z}.$$

Note that the half-spreads $(\delta^{b,z}, \delta^{a,z})$ are controlled by the market maker and chosen from the set of admissible half-spreads

$$\mathcal{D} = \left\{ \delta : \delta \text{ progressively measurable s.t. } \mathbb{E} \left[\int_0^T \delta_t^2 dt \right] < \infty \right\}. \quad (2.1)$$

Market buy and sell orders of OTC clients are modelled by independent counting processes $N^{a,z} = (N_t^{a,z})_{t \geq 0}$ and $N^{b,z} = (N_t^{b,z})_{t \geq 0}$ for any order size $z \in \mathcal{Z}$, with intensities

$$\Lambda^{b/a,z}(\delta_t^{b/a,z}) = \lambda^{b/a,z} \exp \left(-\kappa^z \delta_t^{b/a,z} \right) \quad 0 \leq t \leq T. \quad (2.2)$$

Here $\lambda^{b/a,z}$ and κ^z are positive constants.

We assume that, in the internal exchange, the client has placed a limit order on one side of the book. Without loss of generality, we take this to be an ask limit order with an initial size of $\bar{\ell}$; that is, the client intends to sell, and the dealer would buy if a trade occurs.

The order size at any time t , denoted by $\{L_t\}_{0,T}$ which is a càdlàg process that determines the liquidity in the internal exchange and it is given by,

$$L_t = \bar{\ell} - \bar{\ell} \int_0^t \mathbb{1}_{\{L_s > 0\}} dC_s + \bar{\ell} \int_0^t \mathbb{1}_{\{L_s \leq 0\}} dA_s + \bar{\ell} R_t - M_t, \quad 0 \leq t \leq T, \quad (2.3)$$

Note that $\bar{\ell} > 0$ corresponds to an order being present at time $t = 0$, whereas $\bar{\ell} \leq 0$ indicates the absence of an order. Moreover, larger magnitudes of $\bar{\ell}$ decrease the likelihood of the order appearing, since liquidity can be consumed by the market maker only when $L_t > 0$ in (2.3). The processes (A, C, M, R) are defined and characterized below.

- $(C_t)_{t \geq 0}$ represents the cancellation of the order by the client and is modelled as an independent Poisson process with intensity ν and unit jump size. An order can be cancelled by the client only if it is present, so jumps of C affect L only when liquidity is available at that time.
- $(A_t)_{t \geq 0}$ represents the arrival of new orders and may also be interpreted as the non-instantaneous replenishment of a previously filled limit order. It is modelled as a Poisson process with intensity μ and unit jump size. While, in practice, a new order could arrive while another is still active, such events are sufficiently rare that we restrict arrivals to occur only when no order is currently present. This assumption is consistent with the behaviour of a TWAP placement algorithm.
- $(M_t)_{t \geq 0}$ represents the cumulative market orders submitted by the dealer. We assume that these orders occur at jump times, which are \mathbb{F} -stopping times $\{\tau_n\}_{n \geq 1}$ controlled by the market

maker, and that each transaction is of unit size. Market orders can only occur at times when $L_t > 0$, and are zero otherwise. Hence, they are chosen from the class

$$\mathbb{T} = \{\tau : \tau \text{ is an } \mathbb{F}\text{-stopping time and } L_{\tau-} \geq 1\}. \quad (2.4)$$

We then define $M_t = \sum_{n \geq 1} \mathbb{1}_{\{\tau_n \leq t\}}$.

- $(R_t)_{t \geq 0}$ represents the replenishment process of an order immediately after the dealer consumes the last unit of liquidity. We assume that $R_t = \sum_{n=1}^{\infty} \chi_n \mathbb{1}_{\{\tau_n \leq t\}} \mathbb{1}_{\{L_{\tau_n}=0\}}$, where $(\chi_n)_{n \geq 1}$ is a sequence of i.i.d. Bernoulli random variables with parameter $p \in [0, 1]$. Replenishment typically occurs when an internal exchange order is part of a larger iceberg order.

The pricing offset of the internal exchange order relative to the mid-price is denoted by the parameter ρ , which can be positive or negative. The price at which the market maker can trade is given by,

$$P_t = \begin{cases} S_t + \rho & \text{if } L_t > 0 \\ \infty & \text{if } L_t \leq 0. \end{cases}$$

The dealer receives an infinitely unfavorable price when trading is impossible. In practice, internal exchange orders may yield a small fee for the client. From the modeling perspective this can be incorporated by letting $\rho = \tilde{\rho} - \xi$, where $\tilde{\rho}$ is the price offset relative to the mid chosen by the client, and $\xi > 0$ is a constant representing the fee per unit. Throughout the paper, we work directly with ρ .

The dealer's position and cash processes are given by,

$$Q_t = \sum_{z \in \mathcal{Z}} z \left(N_t^{b,z} - N_t^{a,z} \right) + M_t,$$

$$X_t = \sum_{z \in \mathcal{Z}} z \left(\int_0^t S_s^{a,z} dN_s^{a,z} - \int_0^t S_s^{b,z} dN_s^{b,z} \right) - \int_0^t P_s dM_s,$$

respectively. The value function of the maker maker is,

$$v(t, s, q, x, l) = \sup_{\delta^b, \delta^a, (\tau_n)_{n \geq 1}} \mathbb{E} \left[X_T + Q_T S_T - \alpha Q_T^2 - \phi \int_t^T Q_s^2 ds - \psi \int_t^T (L_s)^+ ds \middle| \mathcal{F}_t \right], \quad (2.5)$$

where the supremum is taken over $\delta^b, \delta^a \in \mathcal{D}$ and $\tau_n \in \mathbb{T}$ (see (2.1) and (2.4)). The constants α, ϕ, ψ are nonnegative and $(\cdot)^+$ is the positive part function. The first two terms on the right-hand side of (2.5) represent the terminal value of the market maker's portfolio; that is, the cash position plus the risky asset position valued at mid. The third and fourth terms implement penalties on the terminal and running positions respectively. The fifth term implements a running penalty on unfilled internal exchange orders.

Using the mathematical argument in [14, Chapter 11, Theorem 11.1], the dynamic programming principle yields that the value function v satisfies the following Hamilton-Jacobi-Bellman quasi-variational inequality (HJBQVI),

$$\begin{aligned}
0 = \max & \left\{ \frac{\partial v}{\partial t}(t, s, q, x, l) + \frac{\sigma^2}{2} \frac{\partial^2 v}{\partial s^2}(t, s, q, x, l) - \phi q^2 - \psi \cdot (l)^+ \right. \\
& + \sum_{z \in \mathcal{Z}, i \in \{b, a\}} \left(\sup_{\delta^{i,z}} \left(\lambda^{i,z} e^{-\kappa^z \delta^{i,z}} (v(t, s, q \pm z, x - zs + z\delta^{i,z}, l) - v(t, s, q, x, l)) \right) \right) \\
& + \nu (v(t, s, q, x, l - \bar{\ell}) - v(t, s, q, x, l)) \mathbb{1}_{\{l > 0\}} + \mu (v(t, s, q, x, l + \bar{\ell}) - v(t, s, q, x, l)) \mathbb{1}_{\{l \leq 0\}}, \quad (2.6) \\
& \mathbb{1}_{\{l > 1\}} (v(t, s, q + 1, x - s - \rho, l - 1)) \\
& + \mathbb{1}_{\{l \leq 1\}} (p v(t, s, q + 1, x - s - (\rho \mathbb{1}_{\{l > 0\}} + \infty \mathbb{1}_{\{l \leq 0\}}), l - 1 + \bar{\ell}) \\
& \left. + (1 - p) v(t, s, q + 1, x - s - (\rho \mathbb{1}_{\{l > 0\}} + \infty \mathbb{1}_{\{l \leq 0\}}), l - 1) - v(t, s, q, x, l) \right\},
\end{aligned}$$

with terminal condition

$$v(T, s, q, x, l) = x + qs - \alpha q^2.$$

The terms on the first argument of the maximum in (2.5) arise from Itô's formula for jump diffusions. The second part of the maximum relates to times where an internal exchange order is traded. In particular, when there is no internal exchange standing liquidity, this term evaluates to $-\infty$ and thus the first part of the maximum is larger. The \pm sign in (2.6) relates to terms indexed by b (a) having a plus (minus) sign, respectively. By using the ansatz

$$v(t, s, q, x, l) = x + qs + h(t, q, l),$$

we can reduce the dimension of (2.6). Solving the Hamiltonians in feedback form then yields

$$\begin{aligned}
0 = \max & \left\{ \frac{\partial h}{\partial t}(t, q, l) - \phi q^2 - \psi \cdot (l)^+ \right. \\
& + \sum_{z \in \mathcal{Z}, i \in \{b, a\}} \left(\frac{z \lambda^{i,z} e^{-1}}{\kappa^z} \exp \left(\frac{\kappa^z}{z} (h(t, q \pm z, l) - h(t, q, l)) \right) \right) \\
& + \nu (h(t, q, l - \bar{\ell}) - h(t, q, l)) \mathbb{1}_{\{l > 0\}} + \mu (h(t, q, l + \bar{\ell}) - h(t, q, l)) \mathbb{1}_{\{l \leq 0\}}, \quad (2.7) \\
& (p h(t, q + 1, l - 1 + \bar{\ell}) + (1 - p) h(t, q + 1, l - 1)) \mathbb{1}_{\{l \leq 1\}} \\
& \left. + h(t, q + 1, l - 1) \mathbb{1}_{\{l > 1\}} - h(t, q, l) - (\rho \mathbb{1}_{\{l > 0\}} + \infty \mathbb{1}_{\{l \leq 0\}}) \right\}
\end{aligned}$$

with terminal condition $h(T, q, l) = -\alpha q^2$. The optimal depths are given by,

$$\delta^{i,z}(t, q, l) = \frac{1}{\kappa^z} + \frac{1}{z} (h(t, q, l) - h(t, q \pm z, l)), \quad \text{for } i \in \{b, a\}, \quad (2.8)$$

and the optimal execution times $(\tau_n)_{n \geq 0}$ are the times such that the maximum in (2.7) evaluates as the second term. We refer to the subspace where this holds as the execution region.

3 Optimal strategy

In this section, we investigate the behaviour of the optimal strategy obtained by numerically solving (2.7) using a backward Euler scheme. An anonymised subsample of HSBC GBPUSD trade data was

used to calibrate the intensities $\lambda^{a,z}$ and $\lambda^{b,z}$, as well as the sensitivity of fill probabilities to quotes, κ^z , for $z \in \mathcal{Z} = \{1, 5, 10\}$. Specifically, we set $\kappa^1 = 1.5$, $\kappa^5 = 1.0$, $\kappa^{10} = 0.5$, $\lambda^{b,1} = \lambda^{a,1} = 0.2$, $\lambda^{b,5} = \lambda^{a,5} = 0.005$, and $\lambda^{b,10} = \lambda^{a,10} = 0.001$. Each κ is measured in bps^{-1} , and each $\lambda^{a,z}$, $\lambda^{b,z}$ in seconds^{-1} . The inventory penalties are given by $\phi = \alpha = 0.001$, and the penalty for unfilled internal exchange orders is $\psi = 0.01$. The time horizon is $T = 300$ seconds, which represents a reasonable high-frequency risk management period. We consider a range of client prices \tilde{p} and three parameter configurations for (2.3) corresponding to different client algos.

- **Iceberg.** The client executes without interruption, but the total size of the order is never visible. As such, when the order is filled by the market maker, it is immediately replenished. At some point, the client finishes executing, and the order is no longer renewed. To represent this, the replenishment probability is set to $p = 0.9$, meaning there is a 10% chance that, after liquidity is taken, it is not renewed. Additionally, there is a small probability that the order is spontaneously cancelled by the client, given by $\nu = 0.001 \text{ s}^{-1}$. For simplicity, we do not allow new orders after the iceberg order finishes executing, so $\mu = 0$. We let the order size be one million notional, and therefore set $\bar{\ell} = 1$.
- **TWAP.** The client executes at a constant pace, but orders arrive with pauses of random lengths from the perspective of the market maker. There is no instantaneous replenishment, so we set $p = 0$. The arrival process A represents the renewal of the order some time after it has been consumed, and we set its intensity to $\mu = 0.05 \text{ s}^{-1}$. As in the case of the iceberg strategy, we let $\nu = 0.001 \text{ s}^{-1}$ and $\bar{\ell} = 1$.
- **Full Amount.** The client places their entire order at once and never updates it, except for the possibility of cancellation. Therefore, there is no instantaneous replenishment ($p = 0$) and no arrivals ($\mu = 0$). However, the initial order size is larger than one, and we set $\bar{\ell} = 10$. As before, the cancellation rate is $\nu = 0.001 \text{ s}^{-1}$.

3.1 Optimal quotes

In Figure 1 we plot the optimal bid-depths $\delta^{b,z}$ and ask-depths $\delta^{a,z}$ at $t = 0$ for $z \in \mathcal{Z}$ for different inventory positions of the market maker, in the case where client orders are placed at mid price ($\tilde{p} = 0$). We chose the snapshot at $t = 0$ as the strategy is approximately stationary away from the terminal time T . We observe in Figure 1 (left panel) that, in the iceberg scenario, the ask-side quotes are lower near the execution region (i.e., when the market maker has a short position greater than one unit, $q < -1$) and the limit order is present (dashed lines), compared with the classical Avellaneda-Stoikov benchmark (solid lines). This occurs because the market maker is more willing to accumulate inventory, given the possibility of closing the position through the limit order if necessary. The effect becomes more pronounced with larger ask sizes and weaker with larger bid sizes, forming a substantial passive impact on both bid and ask prices. In contrast, this effect is marginal in the presence of a TWAP order due to the time intervals between successive renewals of the client's orders (see right panel).

In Figure 2, we show the full (one-time) client order scenario for an outstanding order of 10 units (left) and one unit remaining (right). Prices adjust much more relative to the no-internal-order case when the order is full at $L = 10$, since greater available liquidity allows larger positions to be closed internally. Note that the trading region here is now $q < 5$ due to the fact that the order will not repeat itself and the market maker can mitigate the internal execution urgency term (the fifth term on the right-hand side of (2.5)) by taking the liquidity. This happens both when $L = 10$ and $L = 1$ due to the linearity of the penalty for positive L .

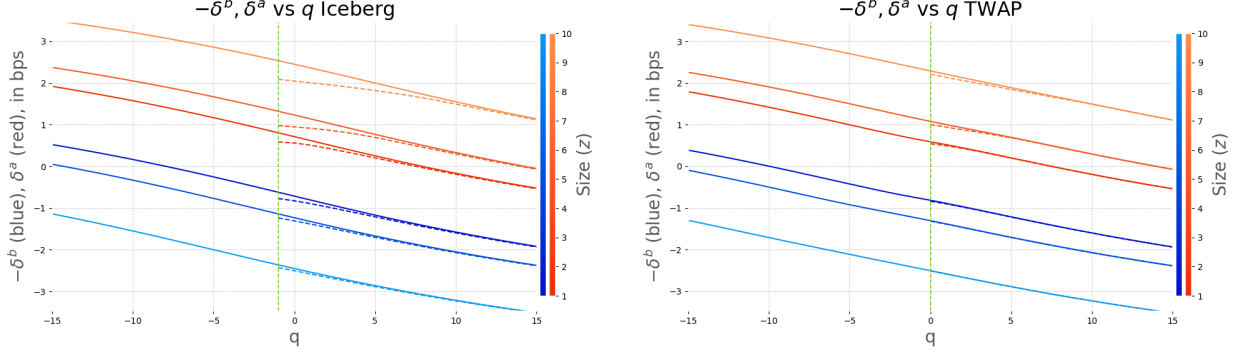


Figure 1: Ask and bid depths, $\delta^{b,z}(0, q, 1)$ and $\delta^{a,z}(0, q, 1)$, when an internal exchange order is present (dashed lines), and when it is not (solid lines), for $z \in \mathcal{Z}$, in the presence of an iceberg (left) and TWAP (right) client algorithm in the internal exchange. Bid (ask) depths are shown in blue (red), with lighter shades corresponding to larger values of z . Dashed lines corresponding to cases where the internal exchange order is available are shown only for values of q outside the execution region, where the internal exchange order is taken. The area to the left of the green line indicates the region where the market maker trades with the internal exchange.

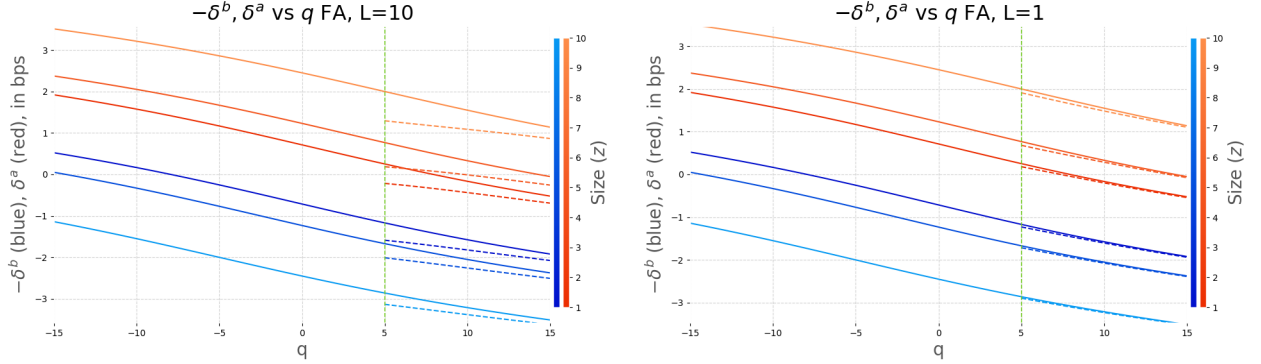


Figure 2: Ask and bid depths when the internal exchange order is present as in Figure 1, now for when the internal exchange order is of full amount type with 10 (left) and 1 (right) units remaining.

We therefore conclude that the type and magnitude of order determines the extent of price adjustment due to internal exchange orders, which can be perceived externally as passive price impact. The pricing adjustments are more pronounced when the expected remaining liquidity on the internal exchange is high. Since clients rarely disclose their trading strategies, a data-driven approach may be needed to infer internal order properties and decide on external price adjustments, initially assuming a standalone order and updating quotes after the first fill.

3.2 When to trade with the internal exchange

In Figure 1, we observed that, for the chosen parameters, it was optimal to trade with the internal exchange whenever the position was negative in the TWAP case, and when the position was sufficiently negative in the iceberg case. In contrast, in Figure 2, it was optimal to trade when the position was already positive. This indicates that the optimal execution boundary depends on the model parameters. Figure 3 illustrates how this boundary changes with the client's price ρ in the

iceberg and TWAP cases. Recall that, in this section, there is no fee ($\xi = 0$), so $\rho = \tilde{\rho}$.

We observe that when the client posts aggressively, with a price below the mid ($\tilde{\rho} < 0$), the market maker may be willing to trade with the internal exchange order even if doing so worsens their position (a behaviour known as risk increasing). This pattern reflects a trade-off between optimising P&L, managing inventory costs, and accommodating the urgency of executing internal orders (as described below (2.5)). Owing to the quadratic inventory risk term, holding a small position is penalised less per unit than holding a large one; hence, the market maker may prefer to hold a small negative position rather than always remain non-negative when $\tilde{\rho}$ is large.

One of the main conclusions from Figure 3 is that clients can expect the time required to fill their orders to be more sensitive to their price offset from the mid when using a TWAP strategy than when using an iceberg order. We demonstrate this result quantitatively in Table 2 in Section 4. The main reason for the difference in trading regions between the iceberg and TWAP scenarios for aggressive client orders ($\tilde{\rho} < 0$) is that, at these prices, iceberg orders effectively represent a batch of existing orders that the market maker can use to make an immediate profit and to close a short position when needed. Therefore, the market maker tends to postpone consuming them until holding a sufficiently negative inventory. In contrast, TWAP orders arrive at an exponential rate and are only placed when there are no outstanding client orders. Hence, it is more profitable to consume them immediately regardless of inventory, as otherwise new orders will not arrive.

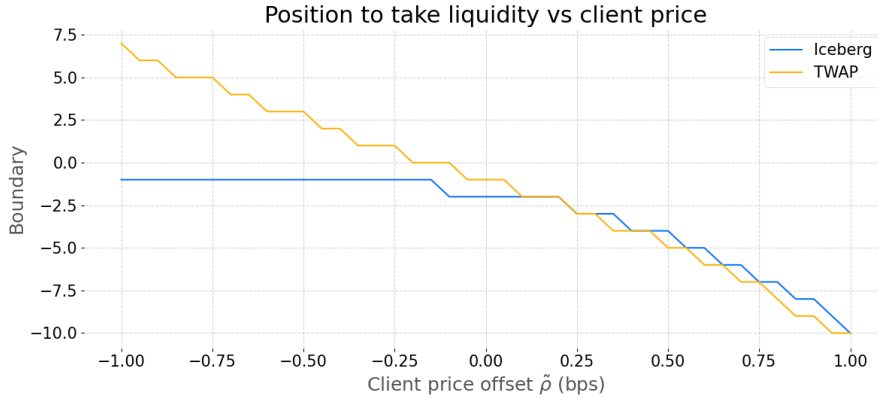


Figure 3: The largest position at which it is optimal to fill the internal exchange order vs. as the client’s price offset $\tilde{\rho}$. The blue line illustrates the iceberg order scenario, while the orange line shows TWAP case.

4 Naïve benchmark and performance comparison

In this section, we compare the performance and behaviour of the optimal strategy obtained by solving the HJBQVI (2.7) and using the optimal depths (2.8), with common heuristic benchmark strategies. We use the same parameters as in Section 3. Moreover, due to the 24-hour nature of the FX market, we evaluate the optimal strategy at time $t = 0$ to neglect terminal inventory constraints in our comparison. This assumption is justified, as the strategy becomes time-independent when sufficiently far from the end of the time horizon. The following algorithm describes the naïve benchmark strategy.

In Algorithm 1 when the market maker’s position is negative and an internal exchange order

Algorithm 1 Naïve benchmark strategy

Initiate Set $t = 0$ and $L_0 = \bar{\ell} \geq 0$ in (2.3).
while $t \leq T$ **do**
 if the available liquidity $L_t \leq 0$ **then**
 Use Avellaneda–Stoikov quotes (i.e. (2.2) at time zero with $L \equiv 0$),
 else
 if $Q_t < 0$ **then**
 Purchase the internal exchange outstanding order and update $L_{t+\Delta} \leftarrow L_t - 1$
 else
 Use Avellaneda–Stoikov bid quotes and **adjusted** ask quotes in (4.1).
 Set $t \leftarrow t + \Delta$.

is present, it is filled¹. If no such order exists, prices are set using the Avellaneda–Stoikov model, that is the case when there is no internal exchange ($L \equiv 0$). Prices on the ask side are adjusted according to a hubristic rule which is described below.

We let $\delta^{a,z_k,AS}(0, q)$ denote the optimal half-spreads at $t = 0$ in the Avellaneda–Stoikov model, i.e., those corresponding to our model with $L \equiv 0$ and allowing for multiple order sizes (see [7]). In the presence of a client’s sell order of size l at a price $\tilde{\rho}$ above the mid-price, the market maker inserts an order on the ask side of their pricing ladder with size l at a price $\tilde{\rho} + \iota$ above mid, where $\iota > 0$ is the margin introduced by the dealer. We then compute a new volume-weighted average price (VWAP) to obtain the naïve strategy’s half-spreads.

Let the order sizes $z_1 < z_2, \dots$ and let $z_i = \min\{z \in \mathcal{Z} : \tilde{\rho} + \iota < \delta^{a,z,AS}(0, q)\}$. The naïve strategy’s half-spreads when there is an order on the internal exchange are then given by $\tilde{\delta}^{a,z_j,BM}(0, q, l) = \delta^{a,z_j,AS}(0, q)$ if $\delta^{a,z_j,AS}(0, q) < \tilde{\rho} + \iota$, which is when $j < i$, and for $j \geq i$,

$$\begin{aligned} \tilde{\delta}^{a,z_j,BM}(0, q, l) = & \frac{1}{z_j} \left(l(\tilde{\rho} + \iota) + z_{i-1} \delta^{a,z_{i-1},AS}(t, q) + \right. \\ & + l \sum_{r=i+1}^j \left(\left(\frac{z_{r-1} \delta^{a,z_{r-1},AS}(t, q) - z_{r-2} \delta^{a,z_{r-2},AS}(t, q)}{z_{r-1} - z_{r-2}} \right) \right) \\ & \left. + \sum_{r=i}^j \left((z_r - z_{r-1} - l) \left(\frac{z_r \delta^{a,z_r,AS}(t, q) - z_{r-1} \delta^{a,z_{r-1},AS}(t, q)}{z_r - z_{r-1}} \right) \right) \right) \end{aligned}$$

To summarise, the prices quoted by the market maker are,

$$\delta^{a,z_j,BM}(0, q, l) = \begin{cases} \tilde{\delta}^{a,z_j,BM}(0, q, l) & \text{if } l > 0, \\ \delta^{a,z_j,AS}(0, q) & \text{if } l \leq 0. \end{cases} \quad (4.1)$$

The prices on the bid-side remain unchanged in the presence of the client’s order, that is $\delta^{a,z_j,BM}(0, q, l) = \delta^{a,z_j,AS}(0, q)$ (see Algorithm 1).²

In Table 1, we compare the mean and standard deviation of the simulated P&L, defined as $P\&L = X_T + Q_T S_T$ (see (2.5)), across different types of client orders and client pricing levels:

¹In reality, this may depend on the client’s posted price $\tilde{\rho}$; however, this is reasonable for aggressive and mildly passive pricing.

²Alternatively, one can consider removing equivalent size from the bid ladder.

aggressive ($\rho = -0.2$), mid ($\rho = 0$), and passive ($\rho = 0.2$). In all cases, the market maker’s initial position is zero, fees are set to $\xi = 0$, the time step is $\Delta = 0.3$ s, the time horizon is $T = 300$ seconds, and the number of simulated trajectories is 5,000. We compare the performance of the optimal strategy with that of the naïve benchmark strategy using a margin of $\iota = 0.1$. We observe that the optimal strategy consistently outperforms the naïve benchmark, and that the market maker’s profits are slightly more sensitive to the client’s pricing offset $\tilde{\rho}$ in the TWAP and full-amount cases than in the iceberg case.

In Table 2, we observe that the expected first fill times of internal client orders under the naïve benchmark strategy are generally shorter than those under the optimal strategy. This results from the naïve strategy’s more aggressive pricing adjustments and its execution rule of always taking when the inventory is negative. In particular, the market maker readily goes short and therefore fills internal exchange orders more frequently. We also observe that the change in the mean time to first fill under the optimal strategy is more sensitive to the client’s pricing offset $\tilde{\rho}$ when a TWAP order is used than when an iceberg order is used. Finally, when a full-amount order is used, the time to first fill is always zero under the optimal strategy, as the initial position of zero lies within the region where the market maker trades with the internal exchange (see Figure 2).

		P&L (K\$)		
		Aggressive: $\rho = -0.2$	At mid: $\rho = 0$	Passive: $\rho = 0.2$
Iceberg	Optimal	3.788 (1.62)	3.677 (1.6)	3.603 (1.6)
	Naïve	3.384 (1.69)	3.36 (1.69)	3.347 (1.69)
TWAP	Optimal	3.84 (1.65)	3.708 (1.61)	3.673 (1.59)
	Naïve	3.381 (1.5)	3.399 (1.53)	3.405 (1.56)
FA	Optimal	3.744 (2.0)	3.547 (2.0)	3.354 (2.0)
	Naïve	3.452 (1.6)	3.433 (1.58)	3.404 (1.57)

Table 1: Mean (standard deviation) of P&L in thousand \$, in the presence of different client algos with passive, at mid, and aggressive placement for iceberg, TWAP and full amount (FA) scenarios.

		Time to First Fill (seconds)		
		Aggressive: $\rho = -0.2$	At mid: $\rho = 0$	Passive: $\rho = 0.2$
Iceberg	Optimal	21.133 (31.57)	47.965 (53.67)	56.681 (63.16)
	Naïve	5.147 (6.5)	7.59 (10.36)	12.145 (18.64)
TWAP	Optimal	Always zero	24.518 (34.69)	68.221 (68.3)
	Naïve	5.571 (7.7)	7.797 (10.24)	12.507 (18.12)
FA	Optimal	Always zero	Always zero	Always zero
	Naïve	5.309 (6.7)	7.402 (9.82)	12.063 (17.34)

Table 2: Mean (standard deviation) of time to fill for the first internal exchange order in seconds, in the presence of different client algos with passive, at mid, and aggressive placement for iceberg, TWAP and full amount (FA) scenarios.

Although execution times under the optimal strategy are slower than those of the naïve strategy, the latter substantially reduces the market maker’s P&L. To sustain internal flow without harming performance, an additional compensation mechanism is required. Figure 4 shows the impact of fees ξ and margins ι on P&L (left) and fill rate (right). Margins are the least effective option: while they improve the naïve strategy’s P&L similarly to fees, they also reduce the volume of internal exchange

executions. Fully compensating the naïve strategy would require a large fee, which could in practice discourage internal exchange participation, an effect not captured by the model. In contrast, the optimal strategy benefits from fees, improving both P&L and turnover.

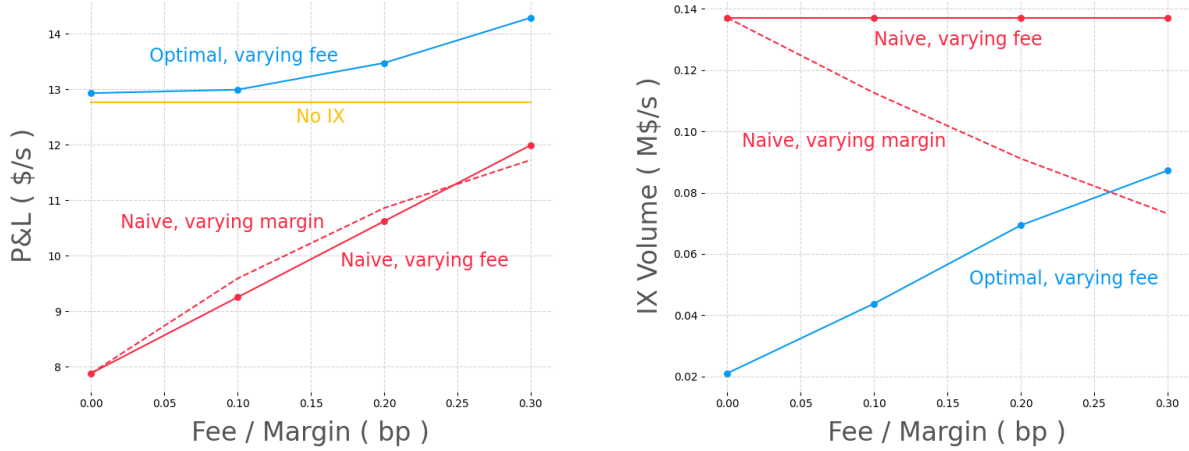


Figure 4: P&L (left) and internal exchange volume per second $\frac{M_T}{T}$ (right) for optimal market maker's strategy (blue) and naïve benchmark strategy (red) in the presence of an iceberg as functions of fee (solid lines) or margin (dashed lines). Reference P&L using Avellaneda-Stoikov pricing without an internal exchange is shown for comparison (orange).

5 Conclusion

We have introduced a model in which a market maker continuously streams prices to clients while also having the option to take liquidity from dynamic passive orders on an internal exchange. We solve for the optimal strategy using an HJBQVI and demonstrate its superior performance compared to a heuristic benchmark strategy that directly incorporates internal exchange orders into the OTC pricing ladder.

The optimal strategy defines an inventory-dependent execution threshold beyond which the market maker is willing to immediately take the internal exchange order. Otherwise, the market maker skews prices to facilitate the opposing flow. The degree of skew and the positioning of the execution threshold depend on the market maker's risk aversion, client order depth and placement strategy, as well as on flow facilitation initiatives.

Notably, the skew mechanism revealed in our analysis highlights a potential origin of passive market impact: clients place passive orders on the internal exchange, prompting the market maker to skew prices for their OTC clients. As described in [10], these clients may then extract alpha from the resulting skew, leading the market to move.

Acknowledgments

The views expressed are those of the authors and do not necessarily reflect the views and practices at HSBC. The authors are grateful to Richard Anthony (HSBC) for helpful discussions and support throughout the project.

References

- [1] Á. Cartea, S. Jaimungal, and J. Ricci. Buy Low Sell High: A High Frequency Trading Perspective. *SIAM J. Financ. Math.*, 5(1):415–444, 2014
- [2] Á. Cartea, and S. Jaimungal. Optimal Execution with Limit and Market Orders. *Quant Finance.*, 15(8):1279–1291, 2015.
- [3] Y. Morimoto. Optimal Execution Strategies Incorporating Internal Liquidity Through Market Making. Preprint, available at SSRN:5074405, 2024.
- [4] Á. Cartea, and L. Sánchez-Betancourt. Brokers and informed traders: dealing with toxic flow and extracting trading signals. *SIAM J. Financ. Math.*, 16(2):243–270 2025
- [5] A. Barzykin, P. Bergault, and O. Guéant. Algorithmic market making in dealer markets with hedging and market impact. *Mathematical Finance*, 33(1):41–79, 2024/06/11 2023.
- [6] M. Butz and R. Oomen. Internalisation by electronic fx spot dealers. *Quant Finance.*, 19(1): 35–56, 01 2019.
- [7] P. Bergault, and O. Guéant. Size matters for OTC market makers: General results and dimensionality reduction techniques *Math Financ.*, 31 (1): 279-322., 2021.
- [8] O. Guéant, C-A. Lehalle, and J. Fernandez-Tapia. Dealing with the inventory risk: a solution to the market making problem *Math. Financ. Econ.*, 7 (4): 477–507., 2013.
- [9] Y.O. Chahdi, M. Rosenbaum, G. Szymanski. A theory of passive market impact Preprint, available at arXiv:2412.07461, 2024.
- [10] A. Barzykin, P. Bergault, O. Guéant, and M. Lemmel. Optimal Quoting under Adverse Selection and Price Reading Preprint, available at arXiv:2508.20225, 2025.
- [11] R. Boyce, M. Herdegen, and L. Sánchez-Betancourt. Market making with exogenous competition. *SIAM J. Financ. Math.*, 16(2):692–706, 2025.
- [12] M. Avellaneda and S. Stoikov. High-frequency trading in a limit order book *Quantitative Finance*, 8(3):217–224, 2008.
- [13] Á. Cartea, S. Jamingual, and J. Penalva. Algorithmic and High Frequency Trading *Cambridge University Press, Cambridge, United Kingdom, 2015* ISBN 9781107091146
- [14] B. Øksendal, and A. Sulem. Applied Stochastic Control of Jump Diffusions *Springer, Berlin, Germany, 2019* ISBN 3540140239
- [FXPA guidance] FXPA Guidance: Definitions & Best Practices for FX Internalisation in Algo Execution Published by the Foreign Exchange Professionals Association (FXPA), July 2025.