

TL;DR: I propose a simple model of competition for earlier transaction execution to make sense of several aspects of transaction ordering policies:

Can we avoid latency competition by using bidding to determine the ordering and inclusion of transactions? (not really)

How does latency competition look like in a batch auction world? (we have zero average profit for bidders in equilibrium when accounting for the cost of latency investment)

What are the performance differences between batch auctions and hybrid formats such as the time boost proposal for Arbitrum? (depends on the parameterization, but in general they look very similar)

The results are particularly relevant for roll-up sequencing but should also inform the broader transaction ordering discourse for L1s or off-chain aggregators.

Many thanks to [Quintus](#), [Akaki](#), [Sergio](#) and [Alejo](#) for discussions and feedback on earlier versions. [Forum](#) for discussion and comments.

Suppose you want to choose a transaction ordering policy for your new roll-up sequencer, L1 blockchain, DEX aggregator, or new financial exchange. The following categories seem to cover your available options pretty well:

First Come First Serve orders transactions by time stamp of arrival at whatever server orders the transactions. Questions of decentralized implementation aside, this policy seems appealing and intuitive to many. It is the go-to policy in traditional finance and hence users are used to interacting with it. FCFS appeals to basic intuitions about fairness. But there is also an efficiency argument to be made: FCFS provides an incentive to incorporate new external information quickly into the state of the system.

Bidding Based Ordering : Orders are processed in batches or blocks. Transactions within a batch interval are ordered according to a function of the attached bids . The function can (hypothetically) be arbitrarily complex: we could order transactions by bid, we could auction individual slots in the batch, or even allow combinatorial bidding where users express preferences over the entire content of the block. Also the “bid” can sometimes be interpreted broadly, for example in Ethereum block building, the bid could contain the amount of MEV the user allows the block builder to extract from him.

Random ordering or other non-conventional policies : Questions of implementation aside, [randomness](#) is a means to achieve ex-ante fairness when ranking transactions. In a different direction, [verifiable sequencing rules](#) are designed to make it detectable if a sequencer deviates from the rule. But, so far, they only have been designed to order AMM transactions for a single trading pair.

Hybrid policies : The above policies can be mixed and matched. For example, FCFS can be implemented with discrete buckets where orders within a bucket are ordered randomly. A [recent proposal](#) by Offchain Labs, orders transactions by a scoring rule called “time boost” that scores transactions by a combination of time stamps and bids.

While these policies look very different from each other, a substantial aspect of all of them is that they organize a contest for earlier transaction execution among those users that care about it (arb traders, liquidators, etc.). The term contest here has the usual meaning from economics: users exert effort (investing in latency, spending money on bidding, spamming your server with transactions) to produce a signal (a timestamp, a bid, a set of transactions IDs) and based on these signals, we decide which transactions to include and in which order (and therefore decide who wins the different contests).

The framing as a contest is helpful, in so far as it gives us an indication of what it means to choose a transaction ordering policy: we organize a contest among users and users maximize whatever signal maximizes their chances of winning the contest. Thus, we need to decide on what dimension we want them to compete: latency investment, expenditure on bidding, the number of transaction requests you receive, increasing entropy, a mixture of all of them, etc?

In the following, I want to focus on the first two categories of transaction ordering policies, time stamp and bidding based policies, and hybrid policies mixing between the first two categories. This is because these policies have an implicit or explicit focus on efficiency (broadly understood) which is desirable.

To motivate the following discussion I would like to start with three somehow obvious observations:

A pure bidding mechanism is impossible: This simply follows from the fact that the bidding phase of an auction cannot run forever. Thus, necessarily some bids end in different batches than others, and time plays a role.

Latency competition necessarily happens, even in a pure batch auction : This is a corollary of the first point. In a batch auction, there still is an advantage for low latency bidders when approaching the end of the batch. This has for example been [well-documented](#) in Ethereum block building, where searchers specialized in CEX-DAX arbitrage need low latency to be competitive in “top of block” MEV.

Illustration of the latency game at the end of a batch auction, source: 0xpandebug on [Arbitrum Research](#)

Bids are more informative signals of users’ preferences than time stamps : If time stamps are used to order transactions, then users are incentivized to create early time stamps. This can lead to investment into colocation and latency reduction. While the negative effects of excessive latency competition has [been recognized](#), investment in latency is not wasteful per

se. It incentivizes quick information incorporation in the state. Moreover, it allows users to express their value for transaction inclusion; users who can generate more value from fast inclusion will invest more in latency reduction. However, everything else being equal, this signal is necessarily less information than a bid because in contrast to bidding, it happens ex-ante before the precise value for transaction inclusion is known to users.

If we agree with the three statements, it seems necessary to focus on hybrid policies that take both bidding and time into consideration: batch auctions and other bidding based procedures are in reality a particular instances of hybrid policies and FCFS is, in most cases, not appealing. But how should we choose, on economic grounds, among different policies that use combinations of time stamps and bids to order transactions?

A first step is to analyze the equilibrium bidding and equilibrium latency investment of users interacting with the policy and to derive the equilibrium welfare and revenue achieved in different designs.

Some disclaimers: The following analysis is purely economic. I abstract away from questions of consensus and implementation and assume that the policies considered can be implemented, because they are run by a trusted centralized sequencer or because we know how to decentralize these policies in a satisfying way. I abstract away from incentive compatibility problems (MEV extraction, censoring etc.) on the side of the party that implements the transaction ordering policy and assume that the policies are implemented as stated.

Equilibrium Analysis of Bidding and Latency Investment

The starting point of my analysis is a bidding and/or latency race between two bidders, who each want their transaction to be executed before the other bidder's transaction (I would expect similar results to hold for more than two bidders). A typical situation that triggers such race could, for example, be an arbitrage opportunity arising through a price discrepancy between an off-chain CEX and an on-chain DEX. Another typical example would be a competition for executing a liquidation. While there are other MEV games played in reality, where bidders have more complicated preferences over transactions orderings than just about how two transactions are ordered relative to each other, these atomic contest for earlier inclusion constitute a large fraction of trading activity on most platforms. Moreover, many other strategies contain an element of it, as it might be a necessary part of the execution of a more complicated trade.

I assume that the value of earlier execution can be different for the two bidders, for example they could have different amounts of liquidity deployed on different platforms so that they can extract different amounts of value from an arb. However, both of the two bidders should have a non-negative value for having their transaction be executed first. The situation is a race between the two bidders in the sense that the bidder, whose transaction is included later, cannot extract any value.

I assume that there are two sources of uncertainty for the two bidders:

The bidders are uncertain about the competitor's value of winning.

The bidders are uncertain about each others' latency. Timestamps are uncertain and random but correlated with (ex-ante) latency investment decisions.

Batch auctions

In a batch auction, all transactions arriving within a pre-specified time window (according to some time stamping scheme) are ordered according to their attached bid, with the highest bid transaction being executed first (if feasible), the second highest bid transaction second (if feasible) etc. A bid cut-off (or reserve price) can be used to bound the total number of transactions being included. A variety of payment rules can be used. In the following, I assume that the competition between the two bidders follow a first price auction format with two possible interpretations: 1) the payment rule is pay-as-you bid, but the capacity is bounded so that the lower bid transaction is not executed 2) the lower bid transaction is reverted if it does not land in the higher slot.

Qualitatively very similar results would also hold for an all-pay batch auction.

For the analysis I normalize time so that bidding happens during a unit time interval $[0, 1]$. The timing is as follows:

An arbitrage opportunity arrives uniformly at random in the unit interval. When the arbitrage opportunity occurs, bidders learn their valuations. A player i has a valuation v_i to have his transaction included first, where v_i is distributed according to $F_i(x)$.

Bidders send a bid for inclusion. Depending on their latency and time of observing the arb, their bid gets included in the current or in the next batch.

At the end of the batch, bids are evaluated according to a first price auction. The higher bid transaction is executed and pays the attached bid.

Assume that there are two bidders with valuations v_1 and v_2 , distributed, for simplicity, i.i.d. uniform on the unit interval. Qualitatively, the analysis carries over to non-uniform valuations. The assumption of independent valuations models the case where there is heterogeneity between the two bidders, while the common component in valuations is known with certainty.

with parameters g, q, q' and c, c' , and transactions are ordered by the score

$$\pi - t \pi - t$$

where t is the timestamp of arrival of a bid at the sequencer and π is the time boost. The parameters have the following interpretation: Parameter g gives the maximal time boost, a user can get from bidding. In particular, transactions finalize after waiting g . Parameter c is the marginal cost per unit of time (normalized by g). The auction is all-pay: Bidders need to pay the time boost fee no matter how transactions are ordered. However, the analysis for a standard winner-pay auction looks very similar (I comment on this below).

Time boost formula (source: [Arbitrum research](#))

Let us again consider the scenario where arb opportunities appear uniformly at random on a time interval and assume that the arb is only good for the earlier transaction. If a bidder has a lower score he will get a payout of zero, but still needs to pay.

The analysis of time boost is more complex than the previous analysis for the batch auction, since now the precise time stamps of the two bidders matter and not only whether the time stamp is before the batch cut-off or not. Thus, making the same analysis as previously, where the values and time stamps of the bidder follows some distribution is generally intractable. However, we can analyze the two extreme cases, where 1) the value of the arb is commonly known by the bidders, but the time stamps are uncertain, 2) the value of the arb to the other bidder is uncertain, but the time stamps are certain.

Known value, uncertain time stamp

Let us first analyze the case where time stamps are uncertain. They are generated according to a probability distribution, which is a function of the latency of the two bidders. The value v of winning is commonly known and the same to both bidders. Bidder 1 wins if and only if his score is higher

$$\pi_1 - t_1 \geq \pi_2 - t_2 \Leftrightarrow \pi_1 + \Delta \geq \pi_2, \pi_1 - t_1 \geq \pi_2 - t_2 \Leftrightarrow \pi_1 + \Delta \geq \pi_2,$$

where $\Delta := t_2 - t_1$ is the difference in time stamps. Thus, it suffices to know the bids π_1, π_2 and the difference Δ in time stamps to determine the winner. In particular, it suffices to know the distribution F_{Δ} of Δ (and not necessarily the individual distributions of the time stamps t_1 and t_2). I analyze the symmetric case where F_{Δ} is symmetric and unimodal with mode 0, which captures the case of equal latency distributions. Inverting the time boost formula, the fee F for buying a time increment of π is

$$F = c\pi g - \pi. F = \frac{c\pi}{g} (g - \pi).$$

If latency is noisy enough the model has a pure strategy symmetric equilibrium:

Equilibrium derivation

In equilibrium,

for low valuations neither bidder will make a time boost bid,

for high valuations bidders produce the same score on average,

the value threshold at which bidders start bidding is increasing in c and decreasing in g and increasing in the variance of time stamps.

To maximize bidding revenue, the parameter c should be chosen small and the maximal time boost g large.

Unknown value, known time stamp

In the previous section I assumed that the value v is commonly known. Now I consider the private value setting, analogous to our analysis of the batch auction. For tractability, I assume that the latency of the two bidders is common knowledge among the two bidders. This simplification is not innocent, but it is a reasonable approximation of reality when bidders interact repeatedly and can estimate the expected time stamp of the competitor with very good accuracy. Second, for simplicity, I consider a linear approximation for the boost formula:

$$F = c\pi g - \pi \approx c\pi g. F = \frac{c\pi}{g} (g - \pi) \approx c\pi g.$$

For a reasonably large boost parameter, e.g. $g = 10$ when $c = 1$ the marginal cost $c(g - \pi)$ (which is relevant for deriving optimal bidding policies) is approximated very well by c/g and we can expect to get qualitatively very similar results for the true boost formula, as long as g is sufficiently large. In the following, I assume $g \geq c$. The equilibrium analysis now follows an all-pay auction with a head start for the lower latency bidder:

Equilibrium derivation

In equilibrium:

For low valuations neither bidder will make a time boost bid.

For high valuations, the bidders produce the same score for the same valuation.

The value threshold at which bidders start bidding is increasing in c , decreasing in g and increasing in the latency difference.

A lower latency bidder underbids relative to a standard all-pay auction, a high latency bidder overbids relative to a standard all-pay auction.

To maximize bidding revenue, the parameter c should be chosen small and the maximal time boost g large.

Both versions of the model, the one with uncertainty about time stamps and the one with uncertainty about valuations make qualitatively similar predictions. An immediate implication is that the parameter c in the time boost formula should be selected small to increase participation and revenue. In a similar way, the maximal time boost g should be selected large. However, the latter comes with trade-offs, as finalization of bids is slower. I comment on these parameter choices and compare the performance in the next section. The equilibrium would look qualitatively similar (pooling for small valuations and separation for high valuations) for a winner-pay instead of an all-pay-auction: The threshold at which bidders start bidding is smaller in that case (the square root disappears).

An analysis of the latency investment game in the same vein as above for the time boost proposal seems intractable. However, the same high level logic should apply: latency competition will lead to a race to the bottom where all profits are competed away. I therefore conjecture:

Conjecture: For the time-boost proposal and many other hybrid policies, ex-ante investment into latency leads to zero average profit in equilibrium: The ex-ante expected gains from competing in the time boost auction are equal to the cost of latency investment of the bidder.

Comparing batch auctions and time boost

It is instructive to compare the relative performance of the two auction formats. To make the batch auction comparable to the time boost proposal, I assume that if neither of the two bidders make it into the batch because they have too high latency, then their orders are both processed in the next batch. We can then look into the equilibrium in either model as a function of the realized difference in latency $\Delta := |\Delta_2 - \Delta_1|$ between the two bidders where Δ_i is the delay of bidder i when sending a bid to the sequencer.

First, let us look at allocative efficiency: how likely is it that the higher valuation bidder wins the race? In the batch auction this can only happen if the two bidders end in different batches and the faster bidder has the lower valuation which happens with probability $\Delta/2$. Under time boost, this can only happen, if the slower bidder refrains from bidding (which in our equilibrium happens if his valuation is below the threshold $u := c g \Delta$) while having a higher valuation. This leads to the following observation:

Under the batch auction design, the likelihood that the high valuation bidder loses is half the latency difference $\Delta/2$.

Under the time boost design, the likelihood is half the cost of compensating for the latency difference $\frac{c}{g} \Delta/2$.

The likelihood of the higher valuation bidder winning is higher under time boost if and only if the marginal cost of bidding is small $c g \leq 1$.

Next let us look at the welfare loss relative to the first best where the item is always allocated to the higher valuation bidder: under which design is this welfare gap larger? In the first best, the surplus is the expectation of the maximum of the two valuations which is $2/3$ for the i.i.d. uniform case.

Derivation of welfare loss

Under the batch auction design, the welfare gap to the optimum is $\Delta/6$.

Under the time boost design, the welfare gap is $\frac{1}{6} (c g \Delta)^{3/2}$.

The welfare gap is smaller with time boost, if and only if the marginal cost of bidding is small $c g \leq 1/\sqrt{3}$.

Next let's look at bidding revenue for the auctioneer. The ex-ante revenue from bidding is the sum of payments received.

Derivation of revenue

Under the batch auction design, the revenue is $(1 - \Delta)/3$.

Under the time boost design, the revenue is $(1 - (c/g\Delta)^{3/2})/3 - (\frac{c}{g}\Delta)^{3/2}/3$.

The revenue is higher under time boost if and only if the marginal cost of bidding is small $c/g \leq 1/\Delta^3 \frac{c}{g} \leq 1/3\Delta$.

The comparisons indicate that the time boost proposal outperforms the batch auction among several dimensions, if the marginal cost of bidding in time boost is small. There is, however, a catch to this: Decreasing the marginal cost by increasing the boost parameter g , leads to longer finalization times of transactions which has other downsides, some of which I have previously discussed. Moreover, it is a somehow unfair comparison, as I have normalized the size of batches to 1, whereas the g parameter, which plays a similar role as the batch size, is allowed to be larger than 1.

To get an intuition what happens with variable batch size, note that if a unit time interval is subdivided into two batches while there is still one arb per unit of time, then the likelihood that the lower valuation bidder wins the race doubles. Similarly, the welfare gap grows and the revenue decreases in the number of batches. Thus, choosing larger batches has a qualitatively very similar effect as choosing a larger maximal time boost g .

The previous analysis is stylized but has immediate implications for economic design:

If the time boost design is implemented, special attention should be put on the choice of parameters. Choosing the marginal cost of bidding too high is detectable in equilibrium. In that case we can predict little use of time boost bidding (no bidding rather than producing low signals) and the parameters should be adjusted.

While latency competition is much less severe for the batch auction than in a FCFS design, there are still advantages of low latency bidding. Towards the end of a batch, bidders with a faster connection can underbid relative to optimal bidding in a standard first price auction, since there is a substantial likelihood that slower competitors do not make it into the current batch. This has an adverse effect on efficiency and revenue. Moreover, the equilibrium analysis predicts that all profits of the bidders are competed away through latency investment by the bidders in the ex-ante stage before the actual bidding.