# The Infinite Server Problem[*][†]

## Christian Coester[1], Elias Koutsoupias[2], and Philip Lazos[3]

1    Department of Computer Science, University of Oxford, Oxford, UK
     christian.coester@cs.ox.ac.uk
2    Department of Computer Science, University of Oxford, Oxford, UK
     elias.koutsoupias@cs.ox.ac.uk
3    Department of Computer Science, University of Oxford, Oxford, UK
     filippos.lazos@cs.ox.ac.uk

### Abstract

We study a variant of the $k$-server problem, the infinite server problem, in which infinitely many servers reside initially at a particular point of the metric space and serve a sequence of requests. In the framework of competitive analysis, we show a surprisingly tight connection between this problem and the $(h, k)$-server problem, in which an online algorithm with $k$ servers competes against an offline algorithm with $h$ servers. Specifically, we show that the infinite server problem has bounded competitive ratio if and only if the $(h, k)$-server problem has bounded competitive ratio for some $k = O(h)$. We give a lower bound of 3.146 for the competitive ratio of the infinite server problem, which implies the same lower bound for the $(h, k)$-server problem even when $k/h \to \infty$ and holds also for the line metric; the previous known bounds were 2.4 for general metric spaces and 2 for the line. For weighted trees and layered graphs we obtain upper bounds, although they depend on the depth. Of particular interest is the infinite server problem on the line, which we show to be equivalent to the seemingly easier case in which all requests are in a fixed bounded interval away from the original position of the servers. This is a special case of a more general reduction from arbitrary metric spaces to bounded subspaces. Unfortunately, classical approaches (double coverage and generalizations, work function algorithm, balancing algorithms) fail even for this special case.

## 1    Introduction

The $k$-server problem is a fundamental well-studied online problem [19, 16]. In this problem $k$ servers serve a sequence of requests. The servers reside at $k$ points of a metric space $M$ and requests are simply points of $M$. Serving a request entails moving one of the servers to the request. The objective is to minimize the total distance traveled by the servers. The most interesting variant of the problem is its online version, in which the requests appear one-by-one and the online algorithm must decide how to serve a request without knowing the future requests. It is known that the deterministic $k$-server problem has competitive ratio between $k$ and $2k - 1$ for every metric space with at least $k + 1$ distinct points [19, 17].

In this paper, we study the *infinite server problem*, the variant of the $k$-server problem in which there are infinitely many servers, all of them initially residing at a given point, the

---

*source*[1]. At first glance it may appear that the lower bound of $k$ for the $k$-server problem would imply an unbounded competitive ratio for the infinite server problem. But consider, for example, the version of the $k$-server problem on uniform metric spaces (i.e. the distance between any two points is 1), and observe that the infinite server problem has competitive ratio 1 for this case.

The infinite server problem is closely related to the $(h, k)$-server problem, the resource augmentation version of the $k$-server problem in which the online algorithm has $k$ servers and competes against an offline algorithm for $h \leq k$ servers. This model is also known as *weak adversaries* [2, 15]. One major open problem in competitive analysis is whether the $(h, k)$-server problem has bounded competitive ratio when $k \gg h$. Here we show a, perhaps surprising, tight connection between the infinite server problem and the $(h, k)$-server problem, which also allows us to improve lower bounds for the latter.

The infinite server problem is also a considerable generalization of the ski-rental problem, since the ski-rental problem is essentially a special case of the infinite server problem when the metric space is an isosceles triangle.

Besides its theoretical appeal, our three main reasons for investigating the infinite server problem are the following. First, the competitive ratio of the $k$-server problem goes to infinity as $k \to \infty$, but for $k = \infty$ it goes back to a small constant at least on some metric spaces. This suggests that the high competitive ratio of the $k$-server problem is somewhat artificial. Second, the problem allows to model applications where the number of servers is so high that it is not a limitation in practice, or where more servers can be bought. A price for buying new servers can be modeled easily by appropriate placement of the source in the metric space. Third, the relationship between the infinite server problem and the $(h, k)$-server problem allows for new ways to tackle the latter.

## 1.1 Previous Work

The $k$-server problem was first formulated by Manasse et al. [19], to generalize a variety of online settings whose stepwise cost had a 'metric'-like structure. They built on previous work by Sleator and Tarjan [20], the genesis of competitive analysis, on the paging problem. This problem can be easily recast as a $k$-server instance for the uniform metric and was already known to be $k$-competitive.

Manasse et al. [19] also showed that the competitive ratio of the $k$-server problem is at least $k$ on any metric space with more than $k$ points. They then proposed the renowned $k$-server conjecture, stating that this bound is tight. This has been shown to be true for $k = 2$ [19] and for several special metric spaces [6, 7, 18, 19, 20]. A stream of refinements [12, 3] led to better competitive ratios for general metric spaces until [17] showed that a competitive ratio of $2k - 1$ can be achieved on any metric space. Chasing the competitive ratio for the deterministic (and randomized) $k$-server problem has been pivotal for the development of competitive analysis. For a more in depth view on the history of the $k$-server problem and further related work, we refer to [16].

In the weak adversaries setting, significantly less is known. For the $(h, k)$-server problem, the exact competitive ratio is $\frac{k}{k-h+1}$ on uniform metrics (equivalent to paging) [20] and weighted stars (equivalent to weighted paging) [21]. Bansal et al. [1] showed recently for weighted trees that the competitive ratio as $k/h \to \infty$ is bounded by a constant depending on the depth of the tree. On general metrics, the $(h, k)$-server problem is still poorly understood.

---

[1] We first learned about this problem from Kamal Jain [14].

No algorithm is known for general metrics that performs better than disabling the $k - h$ extra servers and using $h$ servers only. In fact, for the line it was shown [2, 1] that the Double Coverage Algorithm and the Work Function Algorithm – despite achieving the optimal competitive ratio of $h$ if $k = h$ [6, 4] – perform strictly worse in the resource augmentation setting than disabling the $k - h$ extra servers and applying the same algorithm to $h$ servers only. For the case that $h$ is not fixed, the Work Function Algorithm was shown to be $2h$-competitive simultaneously against any number $h \leq k$ of offline servers [15].

In terms of lower bounds, it is known that unlike for (weighted) paging, the competitive ratio does not converge to 1 on general metrics even as $k/h \to \infty$. Prior to this work, the best known lower bounds were 2 on the line, due to Bar-Noy and Schieber (see [5, p. 175]), and 2.4 for general metrics as shown recently by Bansal et al. [1].

The closest publication to this work is by Csirik et al. [10], which studies a problem that is essentially the special case of the infinite server problem on the uniform metric space augmented by a far away source. It is cast as a paging problem where new cache slots can be bought at a fixed price per unit and gives matching upper and lower bounds of $\approx 3.146$ on the competitive ratio.

## 1.2    Our Results

Our main result is an equivalence theorem between the infinite server problem and the $(h, k)$-server problem, presented in Section 2. It states that the infinite server problem is competitive on every metric space if and only if the $(h, k)$-server problem is $O(1)$-competitive on every metric space as $k/h \to \infty$. We show further that it is not even necessary to let $k/h$ tend to infinity because in the positive case, there must also exist some $k = O(h)$. The theorem holds also if "every metric space" is replaced by "the real line".

In Section 3 we present upper and lower bounds on the competitive ratio of the infinite server problem on a variety of metric spaces. Extending the work in [10], we present a tight lower bound for non-discrete spaces, which is then turned into a 3.146 lower bound for the $(h, k)$ setting. To our knowledge, this is the largest bound on the weak adversaries setting for any metric space, as $k/h \to \infty$. We show how recent work by Bansal et al. [1] can be adapted to give an upper bound on the competitive ratio of the infinite server problem on bounded-depth weighted trees. We also consider layered graph metrics, which are equivalent (up to a factor of 2) to general graph metrics. We have not settled the case for their competitive ratio, but we present a natural algorithm with tight analysis and pose challenges for further research. The main open problem is whether there exists a metric space on which the infinite server problem is not competitive.

In Section 4 we show how a variety of known algorithms such as the work function and balancing algorithms fail for the infinite server problem, even on the real line. We focus in particular on a class of speed-adjusted variants of the well-known double coverage algorithm.

Finally, we present a useful reduction from arbitrary metric spaces to bounded subspaces in Section 5. In particular, the infinite server problem on the line is competitive if and only if it is competitive for the special case where requests are restricted to some bounded interval further away from the source.

## 1.3    Preliminaries

Let $M = (M, d)$ be a metric space and let $s$ be a point of $M$. In the *infinite server problem on* $(M, s)$, an unbounded number of servers starts at point $s$ and serves a finite sequence

$\sigma = (\sigma_0 = s, \sigma_1, \sigma_2, \ldots, \sigma_m)$ of requests $\sigma_i \in M$. Serving a request entails moving one of the servers to it. The goal is to minimize the total distance traveled by the servers.

We drop $s$ in the notation if the location of the source is not relevant or understood. We refer to the action of moving a server from the source to another point as *spawning*. Throughout this work we use the letter $d$ for the metric associated with the metric space.

In the online setting, the requests are revealed one by one and need to be served immediately without knowledge of future requests. All algorithms considered in this paper are deterministic. An algorithm is called *lazy* if it moves only one server to serve a request at an unoccupied point and moves no server if the requested point is already covered. An algorithm is called *local* [9] if it moves a server from $a$ to $b$ only if there is no server at some other point $c$ on a shortest path from $a$ to $b$, i.e. with $d(a,b) = d(a,c) + d(c,b)$. It is easy to see that any algorithm can be turned into a lazy and local algorithm without increasing its cost (i.e. the total distance traveled by all servers).

For an algorithm $ALG$, we denote by $ALG(\sigma)$ its cost on the request sequence $\sigma$. Similarly, we write $OPT(\sigma)$ for the optimal (offline) cost.

An online algorithm $ALG$ is $\rho$-*competitive* for $\rho \geq 1$ if $ALG(\sigma) \leq \rho OPT(\sigma) + c$ for all $\sigma$, where $c$ is a constant independent of $\sigma$. The *competitive ratio* of an algorithm is the infimum of all such $\rho$. We say that an algorithm is *competitive* if it is $\rho$-competitive for some $\rho$. We also call an online problem itself ($\rho$-)competitive if it admits such an algorithm. If the additive term $c$ in the definition is 0, then the algorithm is also called *strictly $\rho$-competitive* [11].

The $(h,k)$-*server problem on $M$* is defined like the infinite server problem except that the number of servers is $k$ for the online algorithm and $h$ for the optimal (offline) algorithm against whom it is compared in the definition of competitiveness. For this problem, the servers are not required to start at the same point, although a different initial configuration would only affect the additive term $c$. The problem is interesting only when $k \geq h$. The case $h = k$ is the standard $k$-server problem and the case $k \geq h$ is known as the *weak adversaries* model. One major open problem is to determine the competitive ratio of the $(h,k)$-server problem as $k$ tends to infinity.

We will sometimes write $OPT_h$ and $OPT_\infty$ for the optimal offline algorithm, where the index specifies the number of servers available.

The following two propositions will be useful later in the paper.

▶ **Proposition 1.** *If for every metric space there exists a competitive algorithm for the infinite server problem, then there exists a universal competitive ratio $\rho$ such that the infinite server problem is* strictly *$\rho$-competitive on every metric space.*

**Proof.** We first show the existence of $\rho$ such that the infinite server problem is $\rho$-competitive (strictly or not) on every metric space. Suppose such $\rho$ does not exist, then for every $n \in \mathbb{N}$ we can find a metric space $M_n$ containing some point $s_n$ such that the infinite server problem on $(M_n, s_n)$ is not $n$-competitive. Consider the metric space obtained by taking the disjoint union of all spaces $M_n$ and gluing all the points $s_n$ together. The infinite server problem would not be competitive on this metric space, in contradiction to the assumption.

Analogously we can also find a universal constant $c$ that works for all metric spaces as additive constant in the definition of $\rho$-competitiveness. A scaling argument shows that also $c = 0$ works.                                                                                              ◀

With a very similar argument we get:

▶ **Proposition 2.** *Let $k = k(h)$ be a function of $h$. Suppose that for every metric space $M$ and for all $h$ there exists an $O(1)$-competitive algorithm for the $(h,k)$-server problem on $M$. Then there exists a universal competitive ratio $\rho$ such that the $(h,k)$-server problem is strictly $\rho$-competitive on every metric space if all servers start at the same point.*

## 2 Equivalence of Infinite Servers and Weak Adversaries

The main result of this section is the following tight connection between the infinite server problem and the weak adversaries model.

▶ **Theorem 3.** *The following are equivalent:*
1. *The infinite server problem is competitive.*
2. *The $(h, k)$-server problem is $O(1)$-competitive as $k/h \to \infty$.*
3. *For each $h$ there exists $k = O(h)$ so that the $(h, k)$-server problem is $O(1)$-competitive.*
*The three statements above are also equivalent if we fix the metric space to be the real line.*

The implication "$3 \implies 2$" is trivial. The proof of the equivalence theorem consists in its core of two reductions. Theorem 4 contains the easier of the two reductions, which is from the infinite server problem to the $k$-server problem against weak adversaries ("$2 \implies 1$"). By Propositions 1 and 2, it suffices to consider only strictly competitive algorithms. Theorem 5 proves essentially the inverse for general metric spaces, and Theorem 8 specializes it to the line ("$1 \implies 3$").

As a corollary of the theorem we get the non-trivial implication "$2 \implies 3$", a potentially useful statement towards resolving the major open problem about weak adversaries: "Is Statement 2 true?" This highlights the importance of the infinite server problem.

▶ **Theorem 4.** *Fix a metric space $M$ and consider algorithms with all servers starting at some $s \in M$. If for every $h$ there exists $k = k(h)$ such that the $(h, k)$-server problem on $M$ is strictly $\rho$-competitive, for some constant $\rho$, then there exists a strictly $\rho$-competitive online strategy for the infinite server problem on $M$.*

**Proof.** Let $ALG_{k(h)}$ denote an online algorithm with $k(h)$ servers that is strictly $\rho$-competitive against an optimal algorithm $OPT_h$ for $h$ servers, i.e.

$$ALG_{k(h)}(\sigma) \leq \rho OPT_h(\sigma) \tag{1}$$

for every request sequence $\sigma$. Without loss of generality, algorithm $ALG_{k(h)}$ is lazy.

For every request sequence $\sigma$, consider the equivalence relation $\equiv_\sigma$ on natural numbers in which $h \equiv_\sigma h'$ if and only if $ALG_{k(h)}(\sigma)$ and $ALG_{k(h')}(\sigma)$ serve $\sigma$ in exactly the same way (i. e., make exactly the same moves). To every $\sigma$, we associate an equivalence class $H(\sigma)$ of $\equiv_\sigma$ that satisfies
- $H(\sigma)$ is infinite,
- $H(\sigma r) \subseteq H(\sigma)$, for every request $r$.

This is done inductively in the length of $\sigma$ (in a manner reminiscent of König's lemma) as follows: For the base case when $\sigma$ is the empty request sequence, $H(\sigma) = \mathbb{N}$. For the induction step, suppose that we have defined $H(\sigma)$. Consider the equivalence classes of $\equiv_{\sigma r}$, a refinement of the equivalence classes of $\equiv_\sigma$. Since there are only finitely many possible ways to serve $r$, they partition $H(\sigma)$ into finitely many parts. At least one of these parts is infinite and we select it to be $H(\sigma r)$; if there is more than one such sets, we select one arbitrarily, say the lexicographically first.

Given such a mapping $H$, we define the online algorithm $ALG_\infty$ which serves every $\sigma$ in the same way as all the online algorithms $ALG_{k(h)}$ for $h \in H(\sigma)$. The second property of $H$ guarantees that $ALG_\infty$ is a well-defined online algorithm.

By construction, $ALG_\infty(\sigma) = ALG_{k(h)}(\sigma)$ for every $h \in H(\sigma)$. To finish the proof, observe that since $H(\sigma)$ is infinite, it contains some $h$ greater than the length of $\sigma$, and for such an $h$ we have $OPT_\infty(\sigma) = OPT_h(\sigma)$. Substituting these to (1), we see that $ALG_\infty$ is strictly $\rho$-competitive. ◀

We now show the reduction from the $k$-server problem against weak adversaries to the infinite server problem on general metric spaces.

▶ **Theorem 5.** *If the infinite server problem on general metric spaces is strictly $\tilde{\rho}$-competitive, then there exists a constant $\rho$ such that the $(h, k)$-server problem is $\rho$-competitive, for $k = O(h)$. In particular, for every $\epsilon > 0$, we can take $\rho = (3 + \epsilon)\tilde{\rho}$ and any $k \geq (1 + 1/\epsilon)\tilde{\rho}h$.*

**Proof.** Fix some metric space $M$ and a point $s \in M$. We will describe a strictly $\rho$-competitive algorithm for the $(h, k)$-server problem on $M$ for the case that all servers start at $s$. This implies a (not necessarily strictly) $\rho$-competitive algorithm for any initial configuration.

The idea is to simulate a strictly $\tilde{\rho}$-competitive infinite server algorithm, but whenever it would spawn a $(k + 1)$-st server, we bring all servers back to the origin and restart the algorithm. The problem is that the overhead cost for returning the servers to the origin, may be very high. To compensate for this, we assume that every time the servers return to the origin, they pretend to start from a different point further away from the origin. This motivates the following notation:

▶ **Definition 6.** Given a metric $M$, a point $s \in M$, and a value $w \geq 0$, we will use the notation $M_{s \oplus w}$ to denote the metric derived from $M$ when we increase the distance of $s$ from every other point by $w$; we will also denote the relocated point by $s \oplus w$.

Let $ALG_\infty$ denote a strictly $\tilde{\rho}$-competitive online algorithm for the infinite server problem. We now define an online algorithm $ALG_k$ for $k$ servers (all starting at $s$). We will make use of the notation $A(\sigma; s)$ to denote the cost of algorithm $A$ to serve the request sequence $\sigma$ when all servers start at $s$.

▶ **Definition 7** ($ALG_k$ derived from $ALG_\infty$)**.** Algorithm $ALG_k$ runs in phases with the initial phase being the 0th phase. At the beginning of every phase, all servers of $ALG_k$ are at $s$. In every phase $i$, the algorithm simulates the infinite server algorithm $ALG_\infty$, whose servers start at $s \oplus w_i$ for some $w_i \geq 0$. The parameters $w_i$ are determined online, and initially $w_0 = 0$. Whenever $ALG_\infty$ spawns a server from $s \oplus w_i$, algorithm $ALG_k$ spawns a server from $s$.

The phase ends just before $ALG_\infty$ spawns its $(k + 1)$-st server or when the request sequence ends. In the former case, all servers of $ALG_k$ return to $s$ to start the $(i + 1)$-st phase. To determine the starting point of the simulated algorithm of the next phase, we set

$$w_{i+1} = \epsilon \, \frac{OPT_h(\sigma_i; s)}{h} \,, \tag{2}$$

where $\sigma_i$ is the sequence of requests during phase $i$.

Let $n$ be the number of phases. The cost of $ALG_k$ for the requests in phase $i < n$ is $ALG_\infty(\sigma_i; s \oplus w_i) - kw_i$; the last term is subtracted because the $k$ servers do not have to actually travel the distance between $s \oplus w_i$ and $s$. However for the last phase no such term can be subtracted since we do not know how many servers are spawned during the phase, and we can only bound the cost from above by $ALG_\infty(\sigma_n; s \oplus w_n)$. The cost of returning the servers to $s$ at the end of a phase can at most double the cost during the phase.

From this, we see that the total cost of $ALG_k$ in phase $i$ is

$$cost_i \leq \begin{cases} 2\left(ALG_\infty(\sigma_i; s \oplus w_i) - kw_i\right) & \text{for } i < n \\ ALG_\infty(\sigma_n; s \oplus w_n) & \text{for } i = n \,. \end{cases}$$

Since $ALG_\infty$ is strictly $\tilde{\rho}$-competitive, we have

$$ALG_\infty(\sigma_i; s \oplus w_i) \leq \tilde{\rho} \, OPT_\infty(\sigma_i; s \oplus w_i)$$
$$\leq \tilde{\rho} \, OPT_h(\sigma_i; s \oplus w_i)$$
$$\leq \tilde{\rho} \, (OPT_h(\sigma_i; s) + hw_i)$$

and substituting this in the expression for the cost, we can bound the total cost by

$$ALG_k(\sigma; s) = \sum_{i=0}^{n} cost_i \leq 2 \sum_{i=0}^{n-1} \left( \tilde{\rho}(OPT_h(\sigma_i; s) + hw_i) - kw_i \right) + \tilde{\rho}(OPT_h(\sigma_n; s) + hw_n)$$
$$= 2 \sum_{i=0}^{n-1} \left( \tilde{\rho}OPT_h(\sigma_i; s) - (k - \tilde{\rho}h)w_i \right) + \tilde{\rho}OPT_h(\sigma_n; s) + \tilde{\rho}hw_n \, .$$

The parameters $w_i$ and $k$ were selected so that the summation telescopes, and we are left with

$$ALG_k(\sigma; s) \leq 2\,\tilde{\rho}\,OPT_h(\sigma_{n-1}; s) + \tilde{\rho}\,OPT_h(\sigma_n; s) + \tilde{\rho}\,\epsilon\,OPT_h(\sigma_{n-1}; s)$$
$$\leq (3 + \epsilon)\,\tilde{\rho}\,OPT_h(\sigma; s) \, . \qquad\qquad \blacktriangleleft$$

The previous reduction requires the infinite server problem to be competitive on *every* metric space. The following variant only requires the infinite server problem to be competitive on the line.

▶ **Theorem 8.** *If the infinite server problem on the line is $\rho$-competitive, then for every $h \in \mathbb{N}$ and $\epsilon > 0$, the $(h, k)$-server problem on the line is $(3 + \epsilon)\rho$-competitive, when $k \geq 2\lceil(1 + 1/\epsilon)\rho h\rceil$.*

**Proof.** A straightforward adaptation of the proof of the previous lemma, shows the existence of a $(3 + \epsilon)\rho$-competitive algorithm for the interval $[0, \infty)$, when $k \geq 2(1 + 1/\epsilon)\rho h$. By doubling the number of online servers so that half of them are used in each half-line, we get a $(3 + \epsilon)\rho$-competitive algorithm for the entire line, when $k \geq 2\lceil(1 + 1/\epsilon)\rho h\rceil$.

Note that the proof assumes strictly competitive algorithms. But, by a straightforward scaling argument, if the infinite server problem on the line is $\rho$-competitive, then it is also strictly $\rho$-competitive. This in turn implies a strictly $\rho$-competitive online algorithm for $M_{0\oplus w}$, since this space is isometric to the subspace $\{-w\} \cup (0, \infty)$ of the line. ◀

In the next section we look at some particular metric spaces and give upper and lower bounds on the competitive ratio.

## 3 Upper and Lower Bounds

Unlike the $k$-server problem, which is 1-competitive if and only if the metric spaces has at most $k$ points and conjectured $k$-competitive otherwise, the situation is more diverse for the infinite server problem. For example, on uniform metric spaces (where all distances are the same) the problem is trivially 1-competitive even if the metric space consists of uncountably many points. This is because an optimal strategy in this case is to spawn a server to every requested point. More generally, this strategy achieves a finite competitive ratio on any metric space where distances are bounded from below and above by positive constants. This suggests that statements about the competitive ratio for the infinite server problem cannot be as simple as the (conjectured) dichotomy for the $k$-server problem, which depends only on the number of points of the metric space. In this section we derive bounds on the competitive ratio for particular classes of metric spaces.

## 3.1 Weighted Trees

We consider the infinite server problem on metric spaces that can be modeled by edge-weighted trees. The points of the metric space are the nodes of the tree, and the distance between two nodes is the sum of edge weights along their connecting path. We choose the source of the metric space as the root of the tree, and define the depth of the tree as the maximal number of edges from the root to a leaf. The number of nodes can be infinite (otherwise the infinite server problem is trivially 1-competitive), but we assume the depth to be finite.

An upper bound on the competitive ratio of such trees follows easily from an upper bound for the $(h, k)$-server on such trees [1] and the equivalence theorem:

▶ **Theorem 9.** *The competitive ratio of the infinite server problem on trees of depth $d$ is at most $O(2^d \cdot d)$.*

**Proof.** Bansal et al. [1, Theorem 1.3] showed that the competitive ratio of the $(h, k)$-server problem on trees of depth $d$ is at most $O(2^d \cdot d)$ provided that $k/h$ is large enough. Inspection of the proof in [1] shows that if all servers start at the root, it is in fact strictly $O(2^d \cdot d)$-competitive. Thus, Theorem 4 implies the result for the infinite server problem.     ◀

## 3.2 Non-Discrete Spaces and Spaces with Small Infinite Subspaces

The following theorem gives a lower bound of 3.146 on the competitive ratio of the infinite server problem on any metric space containing an infinite subspace of a diameter that is small compared to the subspace's distance from the source. For example, every non-discrete metric space has this property (unless the source is the only non-discrete point), since non-discrete metric spaces contain infinite subspaces of arbitrarily small diameter. The theorem is a generalization of such a lower bound established in [10] for a variant of the paging problem where cache cells can be bought. Crucial parts of the subsequent proof are as in [10].

▶ **Theorem 10.** *Let $M$ be a metric space containing an infinite subspace $M_0 \subset M$ of finite diameter $\delta$ and a point $s \in M \setminus M_0$ such that the infimum $\Delta$ of distances between $s$ and points in $M_0$ is positive. Let $\lambda > 3.146$ be the largest real solution to*

$$\lambda = 2 + \ln \lambda \,. \tag{3}$$

*The competitive ratio of any deterministic online algorithm for the infinite server problem on $(M, s)$ is bounded from below by a value that converges to $\lambda$ as $\Delta/\delta \to \infty$. In particular, the competitive ratio is at least $\lambda$ if $M \setminus \{s\}$ contains a non-discrete part.*

**Proof.** By scaling the metric, we can assume that $\delta = 1$. Let $p_1, p_2, p_3, \ldots$ be infinitely many distinct points in $M_0$.

Fix some lazy deterministic online algorithm $ALG$. We consider the request sequence that always requests the point $p_i$ with $i$ minimal such that $p_i$ is not occupied by a server of $ALG$. We call a move of a server between two points in $M_0$ *local* (i.e. every move that does not spawn is local). Let $f_j$ be the cumulative cost of local moves incurred to $ALG$ until it spawns its $j$th server. Let $\sigma_k$ be this request sequence that is stopped right after $ALG$ spawns its $k$th server, for some large $k$. The total online cost is

$$ALG(\sigma_k) \geq k\Delta + f_k \,. \tag{4}$$

Let $h = \lceil k/\lambda \rceil$. We consider several offline algorithms that start behaving the same way, so we think of it as one algorithm initially that is forked into several algorithms later. The

offline algorithms make use of only $h$ servers and they begin by spawning them to the points $p_1, \ldots, p_h$. They do not need to move any servers until $ALG$ spawns its $h$th server. Whenever $ALG$ spawns its $j$th server for some $j \geq h$, every offline algorithm is forked to $h$ distinct algorithms: Each of them moves a different server to $p_{j+1}$ (to prepare for the next request, which will be at $p_{j+1}$). We will keep the invariant that each offline algorithm already has a server at the next request. To this end, whenever $ALG$ does a local move from $p$ to $p'$, every offline algorithm that does not have a server at $p$ moves a server from $p'$ to $p$; note that the algorithm had a server at $p'$ by the invariant, and the next request will be at $p$.

When $ALG$ has $j$ spawned servers ($j \geq h$), the offline algorithms are in $\binom{j}{h-1}$ different configurations, each of which occurs equally often among them. If $ALG$ does a local move from $p$ to $p'$, there are $\binom{j-1}{h-1}$ different offline configurations for which a local move is made in the opposite direction. Thus, for each local move by $ALG$ while having $j$ servers in total, a portion $\binom{j-1}{h-1}/\binom{j}{h-1} = \frac{j-h+1}{j}$ of the offline algorithms move a server in the opposite direction for the same cost.

We use the average cost of all offline algorithms we considered as an upper bound on the optimal cost. The cost of spawning $h$ servers is at most $h(\Delta + 1)$, and the average cost while $ALG$ has $j$ spawned servers (for $j = h, \ldots, k-1$) is at most $\frac{j-h+1}{j}(f_{j+1} - f_j) + 1$ (with the "+1" coming from the move when offline algorithms fork). Hence,

$$OPT(\sigma_k) \leq h(\Delta + 1) + k - h + \sum_{j=h}^{k-1} \frac{j-h+1}{j}(f_{j+1} - f_j),$$

$$\leq h\Delta + k + \frac{k-h}{k-1}f_k - \frac{f_h}{h} - \sum_{j=h+1}^{k-1} \frac{h-1}{j(j-1)}f_j,$$

Note that $\frac{f_k}{k}$ is bounded from above because otherwise $ALG$ would not be competitive, and it is bounded from below by 0. Thus, $L = \liminf_{k \to \infty} \frac{f_k}{k}$ exists. In the following we use the asymptotic notation $o(1)$ for terms that disappear as $k \to \infty$. We can choose arbitrarily large values of $k$ such that $\frac{f_k}{k} = L + o(1)$. Since $h = \lceil k/\lambda \rceil$, we have $\frac{f_j}{j} \geq L + o(1)$ for all $j \geq h$. Moreover, $\sum_{j=h+1}^{k-1} \frac{1}{j-1} = \ln(\lambda) + o(1)$. This allows us to simplify the previous bound to

$$OPT(\sigma_k) \leq \frac{k}{\lambda}\Big(\Delta + \lambda + \big(\lambda - 1 - \ln(\lambda)\big)L + o(1)\Big)$$

$$= \frac{k}{\lambda}\big(\Delta + L + \lambda + o(1)\big),$$

where the last step uses equation (3).

The competitive ratio is at least

$$\frac{ALG(\sigma_k) + O(1)}{OPT(\sigma_k)} \geq \frac{k\Delta + f_k + O(1)}{\frac{k}{\lambda}\big(\Delta + L + \lambda + o(1)\big)}$$

$$= \lambda \cdot \frac{\Delta + L}{\Delta + L + \lambda} + o(1).$$

The fraction in the last term tends to 1 as $\Delta \to \infty$.                                          ◄

This bound is tight due to a matching upper bound in [10] that shows (translated to the terminology of the infinite server problem) that a competitive ratio of $\lambda$ can be achieved on metric spaces where all pairwise distances are 1 except that the source is at some larger distance $\Delta$ from the other points.

The previous theorem together with the equivalence theorem also allows us to obtain a new lower bound for the $k$-server problem against weak adversaries.

▶ **Corollary 11.** *For sufficiently large h, there is no* 3.146-*competitive algorithm for the* $(h, k)$-*server problem on the line, even if* $k \to \infty$.

**Proof.** By a scaling argument it is easy to see that if the infinite server problem on the line is $\rho$-competitive, then it is also *strictly* $\rho$-competitive. Thus, the statement follows from Theorems 4 and 10.                                                                                   ◀

This improves upon both the previous best known lower bounds of 2 for this problem on the line [5, p. 175] and 2.4 on general metric spaces [1].

## 3.3   Layered Graphs

A *layered graph of depth D* is a graph whose (potentially infinitely many) nodes can be arranged in layers $0, 1, \ldots, D$ so that all edges run between adjacent layers and each node – except for a single node in layer 0 – is connected to at least one node of the previous layer. The induced metric space is the set of nodes with the distance being the minimal number of edges of a connecting path. For the purposes of the infinite server problem, the single node in layer 0 is the source. We assume $D \geq 2$ to avoid trivial cases.

Note that a connected graph is layered if and only if it is bipartite. Moreover, any graph can be embedded into a bipartite graph by adding a new node in the middle of each edge. So essentially, layered graphs capture *all* graph metrics.

Let *Move Only Outwards (MOO)* be some lazy and local algorithm for the infinite server problem on layered graphs that moves servers along edges only in the direction away from the source. Not surprisingly, the competitive ratio of this simple algorithm is quite bad and we show that it is exactly $D - 1/2$. Nonetheless, at least for $D \leq 3$ this is actually the optimal competitive ratio.

▶ **Theorem 12.** *The competitive ratio of MOO is exactly* $D - \frac{1}{2}$.

▶ **Theorem 13.** *The competitive ratio of the infinite server problem on layered graphs of depth D is exactly* 1.5 *for* $D = 2$, *exactly* 2.5 *for* $D = 3$ *and at least* 3 *for* $D \geq 4$.

Both theorems are proved in the full version of this paper. It remains an open problem to close the gap between the lower bound of 3 and the upper bound of 3.5 for $D = 4$. More importantly, we are interested in the question whether an algorithm better than MOO exists for large $D$, achieving a competitive ratio of less than $D - 1/2$ on any layered graph of depth $D$. Note that if no algorithm with a competitive ratio of $O(1)$ as $D \to \infty$ exists, then the infinite server problem on general metric spaces would not be competitive.

For large $D$, the lower bound of 3 is certainly not tight: Consider a layered graph where each layer contains one node except that the bottom layer contains infinitely many nodes. By Theorem 10 (and a matching upper bound shown in [10]), the competitive ratio on this graph converges to $\lambda \approx 3.146$ as $D \to \infty$.

## 4   Algorithms with Unbounded Competitive Ratio

We examine the performance of classical algorithms known for the $k$-server problem when applied to the infinite server problem, focusing on the line as a particularly appealing metric space. We also consider a generalization of the Double Coverage algorithm for the line with adjusted server speeds. This idea has proved successful for the $(h, k)$-server problem (and hence the infinite server problem) on weighted trees [1]. However, neither of these algorithms is competitive for the infinite server problem on the line. Proofs of the results of this section

as well as the definitions of the algorithms of the following theorem can be found in the full version of this paper.

▶ **Theorem 14.** *The Work Function Algorithm [8, 17], Balance [19] and Balance2 [13] are not competitive for the infinite server problem on the line.*

Perhaps more surprising than for the above three algorithms is that a class of algorithms extending the Double Coverage (DC) algorithm [6] is also not competitive for the infinite server problem. The basic DC algorithm on the line serves each request by an adjacent server. If the request lies between two servers, both servers move towards it at equal speed until one of them reaches the request. A sensible extension of this algorithm seems to be to give different speeds to servers, so that they move away from the source faster than towards it.

We consider here only the half-line $[0, \infty)$ with the source at the left border 0. Let $x_i$ be the position of the $i$th server from the right. We use the notation $x_i$ both for its position and for the server itself. As servers do not overtake each other, $x_i$ is the $i$th spawned server. Let $\mathcal{S} = \{s_i \geq 1 \mid i \in \mathbb{N} \text{ and } i \geq 2\}$ for a monotonic (non-decreasing or non-increasing) sequence of speeds $s_i$. The algorithm $\mathcal{S}$-DC is defined as follows:

- If there exist servers $x_{i+1}$ and $x_i$ to the left and right of the request, move them towards it with speeds $s_{i+1}$ and 1 respectively until one of the two reaches it.
- If a request does not have a server to its right, move the rightmost server to the request.

If $s_i = 1$ for all $i$, this is precisely the original DC algorithm.

▶ **Theorem 15.** *Algorithm $\mathcal{S}$-DC is not competitive for any $\mathcal{S}$.*

The intuitive reason is that servers move to the right either too slowly or too fast: Imagine repeatedly requesting the same $n$ points in some small interval away from the source, until $\mathcal{S}$-DC covers all $n$ points. One case is that $\mathcal{S}$-DC spawns too slowly and is therefore defeated by an adversary covering these $n$ positions immediately with $n$ servers. In the other case, the adversary will also use $n$ servers to cover the initial group of requests and then shift its group of servers slowly towards the source, always making requests at the new positions of these offline servers. As $\mathcal{S}$-DC tries to cover the new requests, it is tricked into spawning too many servers. Both cases lead to an unbounded competitive ratio.

## 5   Reduction to Bounded Spaces

In this section we show a reduction from the infinite server problem on general metric spaces to bounded subspaces. Specifically, a metric space can be partitioned into "rings" of points whose distance from the source is between $r^n$ and $r^{n+1}$, where $r > 1$ is fixed and $n \in \mathbb{Z}$. We show that if the infinite server problem is strictly $\rho$-competitive on each ring, then it is competitive on the entire metric space.

▶ **Theorem 16.** *Let $M$ be a metric space and $s \in M$ and let $r > 1$. For $n \in \mathbb{Z}$ let $M_n = \{s\} \cup \{p \in M \mid d(s, p) \in [r^n, r^{n+1})\}$. If for each $n$ the infinite server problem on $(M_n, s)$ is strictly $\rho$-competitive, then on $(M, s)$ it is strictly $\frac{4r-1}{r-1}\rho$-competitive.*

**Proof.** Let $ALG_n$ be a $\rho$-competitive algorithm for the infinite server problem on $(M_n, s)$.

For a request sequence $\sigma$, let $\sigma_n$ be the subsequence of requests in $M_n$. Let $ALG$ be the algorithm for $(M, s)$ that uses different servers for each of the subsequences $\sigma_n$ and serves them independently according to $ALG_n$.

The total online cost is $ALG(\sigma) = \sum_n ALG_n(\sigma_n) \leq \rho \sum_n OPT(\sigma_n)$. To finish the proof, it suffices to show that

$$\sum_n OPT(\sigma_n) \leq \frac{4r-1}{r-1} OPT(\sigma) \, . \tag{5}$$

Thus, we only need to analyze the offline cost. We do this for each offline server separately. Fix some offline server $x$. Let $N_0$ and $N_1$ be the minimal and maximal values of $n$ such that $x$ visits $M_n$. We can assume without loss of generality (by adding virtual points to the metric space) that whenever $x$ moves from $M_n$ to $M_{n'}$ for some $n < n'$, it travels across points $p_{n+1}, p_{n+2}, \ldots, p_{n'}$ with $d(s, p_i) = r^i$, and similarly for $n > n'$.

The movements of server $x$ can be tracked by many servers, one server $x_n$ in every set $M_n$ for $N_0 \leq n \leq N_1$. When server $x$ is in $M_n$, server $x_n$ is exactly at the same position tracking the movement of $x$. When server $x$ exits $M_n$ at some point $p$ at the boundary to $M_{n-1}$ or $M_{n+1}$, server $x_n$ freezes at $p$. The movement cost of $x_n$ can be partitioned into the cost of deploying $x_n$ at the first point visited in $M_n$, the tracking cost within $M_n$, and the cost of of relocating $x_n$ whenever $x$ re-enters $M_n$ at a location different from the last exiting location.

The total tracking cost of all servers $x_n$ is bounded by the distance traveled by $x$. The cost of deploying all servers $x_n$ is $\sum_{n=N_0}^{N_1} r^n \leq \sum_{n=-\infty}^{N_1} r^n = r^{N_1+1}/(r-1)$, which is at most $\frac{r}{r-1}$ times the total movement of server $x$, because the latter is at least $r^{N_1}$.

To bound the relocating cost, say $x$ exits $M_n$ at $p$ and re-enters it at $p'$. Then $p$ and $p'$ are at the boundary of $M_n$ and $M_{n+u}$ for $u \in \{-1, +1\}$. Let $b$ be the distance traveled by $x$ in $M_{n+u}$ between the times when it is entered at $p$ and when it is next exited. If this exiting is at $p'$, then the relocating cost $d(p, p')$ is at most $b$ by the triangle inequality. Otherwise, $x$ exits $M_{n+u}$ at a point $p''$ at the boundary of $M_{n+u}$ and $M_{n+2u}$. If $u = 1$, then $d(p, p') \leq d(s, p) + d(s, p') = 2r^{n+1}$ and $b \geq d(p, p'') \geq d(s, p'') - d(s, p) = r^{n+2} - r^{n+1} = (r-1)r^{n+1}$. If $u = -1$, then $d(p, p') \leq d(s, p) + d(s, p') = 2r^n$ and $b \geq d(p, p'') \geq d(s, p) - d(s, p'') = r^n - r^{n-1} = \frac{r-1}{r}r^n$. In both cases, the relocating cost $d(p, p')$ is at most $\frac{2r}{r-1}b$. Thus, the total relocating cost of all servers $x_n$ is at most $\frac{2r}{r-1}$ times the total distance traveled by $x$.

Thus, the sum of deployment, tracking and relocating cost of the servers $x_n$ is at most $\frac{4r-1}{r-1}$ times the distance traveled by $x$. This shows (5), giving the statement of the theorem. ◀

The last theorem can also be slightly generalized to the case where instead of *strict* $\rho$-competitiveness, an additive term proportional to $r^n$ is allowed. It is not difficult to show the following specialization for the line, where the premise can be weakened to require competitiveness only on a single interval:

▶ **Corollary 17.** *Let $0 < a < b$. The infinite server problem is competitive on the line if and only if it is competitive on $(\{0\} \cup [a, b], 0)$.*

## 6 Open Problems

The most obvious open problem is whether the infinite server problem is competitive on general metric spaces. A challenging special case is to resolve the question for the real line. Similarly, improving the MOO algorithm and settling the question for layered graphs remains open. It would also be interesting to find a metric space with a competitive ratio greater than 3.146 for the infinite server problem or the $(h, k)$-server problem when $k \gg h$. Another possible line of research is to consider randomized algorithms.

## References

1 Nikhil Bansal, Marek Eliáš, Łukasz Jeż, and Grigorios Koumoutsos. The $(h, k)$-server problem on bounded depth trees. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms, 2017*, pages 1022–1037. SIAM, 2017. `doi:10.1137/1.9781611974782.65`.

2 Nikhil Bansal, Marek Eliáš, Łukasz Jeż, Grigorios Koumoutsos, and Kirk Pruhs. Tight bounds for double coverage against weak adversaries. In *International Workshop on Approximation and Online Algorithms*, pages 47–58. Springer, 2015. `doi:10.1007/978-3-319-28684-6_5`.

3 Yair Bartal and Eddie Grove. The harmonic k-server algorithm is competitive. *J. ACM*, 47(1):1–15, January 2000. `doi:10.1145/331605.331606`.

4 Yair Bartal and Elias Koutsoupias. On the competitive ratio of the work function algorithm for the k-server problem. *Theoretical Computer Science*, 324(2):337–345, 2004. `doi:10.1016/j.tcs.2004.06.001`.

5 Allan Borodin and Ran El-Yaniv. *Online Computation and Competitive Analysis*. Cambridge University Press, New York, NY, USA, 1998.

6 Marek Chrobak, Howard Karloff, Tom Payne, and Sundar Vishwanathan. New results on server problems. *SIAM J. Discret. Math.*, 4(2):172–181, March 1991. `doi:10.1137/0404017`.

7 Marek Chrobak and Lawrence L. Larmore. An optimal on-line algorithm for k-servers on trees. *SIAM J. Comput.*, 20(1):144–148, February 1991. `doi:10.1137/0220008`.

8 Marek Chrobak and Lawrence L. Larmore. The server problem and on-line games. In *On-line Algorithms, volume 7 of DIMACS Series in Discrete Mathematics and Theoretical Computer Science*. Citeseer, 1992.

9 Ilan Reuven Cohen, Alon Eden, Amos Fiat, and Łukasz Jeż. Pricing online decisions: Beyond auctions. In *Proceedings of the Twenty-sixth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA'15, pages 73–91, Philadelphia, PA, USA, 2015. Society for Industrial and Applied Mathematics. `doi:10.1137/1.9781611973730.7`.

10 János Csirik, Csanád Imreh, John Noga, Steven S. Seiden, and Gerhard J. Woeginger. Buying a constant competitive ratio for paging. In *Proceedings of the 9th Annual European Symposium on Algorithms*, ESA'01, pages 98–108, London, UK, 2001. Springer-Verlag. `doi:10.1007/3-540-44676-1_8`.

11 Yuval Emek, Pierre Fraigniaud, Amos Korman, and Adi Rosén. On the additive constant of the k-server work function algorithm. In *International Workshop on Approximation and Online Algorithms*, pages 128–134. Springer, 2009. `doi:10.1007/978-3-642-12450-1_12`.

12 Amos Fiat, Yuval Rabani, and Yiftach Ravid. Competitive k-server algorithms. In *Proceedings of the Thirty First Annual Symposium on Foundations of Computer Science*, pages 454–463 vol.2, Oct 1990. `doi:10.1109/FSCS.1990.89566`.

13 Sandy Irani and Ronitt Rubinfeld. A competitive 2-server algorithm. *Information Processing Letters*, 39(2):85–91, 1991. `doi:10.1016/0020-0190(91)90160-J`.

14 Kamal Jain. Personal Communication.

15 Elias Koutsoupias. Weak adversaries for the k-server problem. In *Proceedings of the 40th Annual Symposium on Foundations of Computer Science*, FOCS'99, Washington, DC, USA, 1999. IEEE Computer Society. `doi:10.1109/SFFCS.1999.814616`.

16 Elias Koutsoupias. The k-server problem. *Computer Science Review*, 3(2):105–118, May 2009. `doi:10.1016/j.cosrev.2009.04.002`.

17 Elias Koutsoupias and Christos H. Papadimitriou. On the k-server conjecture. *J. ACM*, 42(5):971–983, September 1995. `doi:10.1145/210118.210128`.

18 Elias Koutsoupias and Christos H. Papadimitriou. The 2-evader problem. *Information Processing Letters*, 57(5):249–252, March 1996. `doi:10.1016/0020-0190(96)00010-5`.

**19** Mark Manasse, Lyle McGeoch, and Daniel Sleator. Competitive algorithms for on-line problems. In *Proceedings of the Twentieth Annual ACM Symposium on Theory of Computing*, STOC'88, pages 322–333, New York, NY, USA, 1988. ACM. `doi:10.1145/62212.62243`.

**20** Daniel D. Sleator and Robert E. Tarjan. Amortized efficiency of list update and paging rules. *Communications of the ACM*, 28(2):202–208, February 1985. `doi:10.1145/2786.2793`.

**21** Neal Young. The k-server dual and loose competitiveness for paging. *Algorithmica*, 11(6):525–541, 1994. `doi:10.1007/BF01189992`.