

Finding Mobile Data Under Delay Constraints With Searching Costs *

Amotz Bar-Noy^{1,3}, Panagiotis Cheilaris², and Yi Feng¹

¹Department of Computer Science, The Graduate Center of the City University of New York

²Alfréd Rényi Institute of Mathematics, Hungarian Academy of Sciences, Budapest, Hungary.

³Department of Computer and Information Science, Brooklyn College of the City University of New York

Abstract

A token is hidden in one of several boxes and then the boxes are locked. The probability of placing the token in each of the boxes is known. A searcher is looking for the token by unlocking boxes where each box is associated with an unlocking cost. The searcher conducts its search in rounds and must find the token in a predetermined number of rounds. In each round, the searcher may unlock any set of locked boxes concurrently. The optimization goal is to minimize the expected cost of unlocking boxes until the token is found. The motivation and main application of this game is the task of paging a mobile user (token) who is roaming in a zone of cells (boxes) in a Cellular Network system. Here, the unlocking costs reflect cell congestions and the placing probabilities represent the likelihood of the user residing in particular cells. Another application is the task of finding some data (token) that may be known to one of the sensors (boxes) of a Sensor Network. Here, the unlocking costs reflect the energy consumption of querying sensors and the placing probabilities represent the likelihood of the data being found in particular sensors. In general, we call mobile data any entity that has to be searched for.

The special case, in which all the boxes have equal unlocking costs has been well studied in recent years and several optimal polynomial time solutions exist. To the best of our knowledge, this paper is the first to study the general problem in which each box may be associated with a different unlocking cost. We first present three special interesting and important cases for which optimal polynomial time algorithms exist: (i) There is no a priori knowledge about the location of the token and therefore all the placing probabilities are the same. (ii) There are no delay constraints so in each round only one box is unlocked. (iii) The token is atypical in the sense that it is more likely to be placed in boxes whose unlocking cost is low. Next, we consider the case of a typical token for which the unlocking cost of any box is proportional to the probability of placing the token in this box. We show that computing the optimal strategy is strongly NP-Hard for any number of unlocking rounds, we provide a PTAS algorithm, and analyze a greedy solution. We propose a natural dynamic programming heuristic that unlocks the boxes in a non-increasing order of the ratio probability over cost. For two rounds, we prove that this strategy is an $8/7 \approx 1.143$ -approximation solution for an arbitrary token and a $\frac{7-2\sqrt{7}}{28-10\sqrt{7}} \approx 1.108$ -approximation for a typical token and that both bounds are tight. We conduct simulations in the application of Cellular Networks, which show that for any number of rounds the approximation factor is even less than $8/7$. We test our algorithms on random data (that either follows the *Zipf* distribution or the uniform distribution) and on “real” data that includes 171929 appearances of 996 users in 5625 cells. The results indicate that our algorithms perform remarkably better than their guaranteed theoretical bounds.

*The work in this paper was supported by the NSF program award CNS-0626606.

1 Introduction

Consider the following combinatorial game. A token is hidden in one out of N boxes following some probability distribution. The boxes are then locked and the only known information about the location of the token is the probability distribution. Each box is associated with an unlocking cost. A searcher needs to find the token as fast as possible by unlocking boxes while minimizing the expected unlocking cost. The searcher is given D ($1 \leq D \leq N$) rounds to find the token where in each round the searcher may unlock any set of locked boxes.

Let the boxes be $\{C_1, \dots, C_N\}$, let w_1, \dots, w_N be the unlocking costs, and let p_1, \dots, p_N be the placing probabilities: with probability p_i the token is placed in box C_i and all the probabilities are independent. The fastest but the most expensive search strategy would unlock all the N boxes in one round, (*a blanket search*). The other extreme is to unlock one box per round terminating once the token is found (*a sequential search*). In general, a search strategy for D rounds is an ordered D -partition $\mathcal{A} = \langle A_1, \dots, A_D \rangle$ of the boxes, such that in the i th round, all the boxes in the set A_i are unlocked if the token was not found during the previous $(i - 1)$ rounds. The search process terminates in round d if the token is found in one of the boxes of the set A_d . Then the cost for the searcher is the total cost of unlocking all the boxes from the sets A_1, \dots, A_d .

The ultimate goal is to minimize both the number of rounds and the expected unlocking cost until the token is found. These are the two main criteria in evaluating the efficiency of a specific search strategy. The problem studied in this paper is a common way to attack bi-criteria optimization problems by constraining one criterion and optimizing the other: *Given the delay constraint of finding the token in at most D search rounds, design a search strategy with minimum expected unlocking cost.*

Example: Let $N = 3$, the placing probabilities be 0.5, 0.2, 0.3, and the unlocking costs be 0.1, 0.2, 0.7 for boxes C_1, C_2, C_3 , respectively. The cost of unlocking all the boxes in one round is 1. Suppose now that the token must be found in $D = 2$ rounds. One possible search strategy is $\{\{C_1\}, \{C_2, C_3\}\}$. For this strategy, with probability 0.5 the token is found in C_1 for a cost of 0.1. Otherwise, with probability $(0.2 + 0.3)$ all the boxes are unlocked for a cost of 1. Thus, the total expected cost is $0.5 \cdot 0.1 + 0.5 \cdot 1 = 0.55$. Another possible search strategy is $\{\{C_1, C_2\}, \{C_3\}\}$. For this strategy, with probability $(0.5 + 0.2)$ the token is found in the first round for a cost of $(0.1 + 0.2)$. Otherwise, with probability 0.3, all the boxes are unlocked for a cost of 1. Thus, the total expected cost is $0.7 \cdot 0.3 + 0.3 \cdot 1 = 0.51$. The above two search strategies follow the non-increasing order p_i/w_i . The expected cost of the search strategy $\{\{C_2\}, \{C_1, C_3\}\}$ that “violates” this order is $0.2 \cdot 0.2 + (0.5 + 0.3) \cdot 1 = 0.84$. Finally, it is not hard to see that with three rounds the best strategy is to unlock the boxes following the order C_1, C_2, C_3 . With probability 0.5, only C_1 is unlocked for a cost of 0.1, with probability 0.2, both C_1 and C_2 are unlocked for the cost of $(0.1 + 0.2)$, and with probability 0.3, all boxes are unlocked for the cost of 1. The total expected cost is therefore $0.5 \cdot 0.1 + 0.2 \cdot (0.1 + 0.2) + 0.3 \cdot 1 = 0.41$ which is the best possible for this example.

Motivation: The main application to the above game is the task of paging a mobile user (token) that is roaming among the cells (boxes) of a Cellular Network (e.g., [15]). When a call to a user arrives, the system must locate the exact cell in which the user resides to establish a connection. If the user reports its new location whenever it crosses boundaries of cells, then the system “knows” its exact location at any time and the task of finding this user becomes trivial. Since Cellular Networks are expected to have many cells (mini-cells or micro-cells) and mobile users are expected to move very fast, the user might cross boundaries of cells very frequently. This would make it infeasible for the user to report its new location each time it enters a different cell because the reporting

process consumes the “expensive” resources: time, energy, and uplink bandwidth. Indeed, many existing location management schemes instruct mobile users to report less often. A common location management framework partitions the cells into location areas (zones), each with possibly many cells. A user reports its new location to the system only when it crosses zone boundaries (e.g., [21]). When a call to a user arrives, the system may page some or all the N cells (boxes) in some zone to find the user. Although the choice of a location management scheme to minimize the overall use of system resources depends on many parameters, such a paging step is common to most of the schemes. Frequently, the system is looking for a mobile user without knowing the exact location of this user. However, in many cases, some a priori knowledge about the whereabouts of the user is known. This knowledge can be modeled with N probability values, one value associated with each of the N cells: with probability p_i the mobile user resides in cell C_i and all the probabilities are independent. This a priori knowledge could either be supplied by the user itself, be extracted from history logs maintained by the system, be based on recent reports and calls involving this user, or be based on some mobility patterns. A paramount task in any location management scheme is to design, analyze, implement, and evaluate efficient paging (search) strategies for a mobile user while taking advantage of any partial knowledge of its whereabouts. The optimization goals are to find the user as fast as possible while “paying” as little as possible for paging the cells.

A special case of this problem in which $w_i = 1$ for all $1 \leq i \leq N$ was studied thoroughly. This special case corresponds to searching for a mobile user in D rounds while minimizing the expected number of cells paged. Efficient polynomial time dynamic programming solutions are known for this case. The scope of this paper is the general case for which the paging costs are not the same for all cells. This mainly reflects the fact that Cellular Networks are highly correlated with user density and cell congestion. As a result, paging a cell with more users could be more expensive than paging a cell with less users. In addition, there are a lot of other factors affecting the cost of paging a cell, which include the maintenance cost of the base stations, the different regulations on radiation emission, etc.

Mobile Data: The scope of this paper is very general. Let *mobile data* be an abstraction of any entity in a network whose exact location is not known to the system at the time when a specific query is looking for this data. Instead, the system knows that the mobile data may be found in one out of N locations. The system has a *profile* for the data which is represented as a vector of probabilities: with probability p_i the data is in location C_i and all the probabilities are independent. Whenever the system queries location C_i to see if it has the data it pays a cost of w_i . Paging mobile users in cellular networks is one application to this general setting but there are more applications. Consider a wireless sensor network that accumulates some information (e.g., weather or traffic). Mobile data may be any information that can be found in this sensor network. In order to save battery energy, the sensors do not push the information but only reply to queries. As a result, the system needs to pull the data by probing the sensors. The above framework models the pull task where the objectives are to minimize the time it takes to get the data and to minimize the expected cost incurred by the sensors that are probed. The above two applications are for wireless networks, but one could think of similar applications in any kind of network, for example, the task of looking for some data in the Internet or in a peer-to-peer network.

Prior art and related work: The general framework of mobile user location management has been studied a lot in the recent fifteen years. See the survey [1]. Modeling uncertainty of locations of mobiles as a probability distribution vector was first studied in [18, 20]. The paper [19] introduced the user profile based paging scheme, under which the problem solved in this paper is discussed. The papers [3, 15–17, 19] described optimal solutions based on dynamic programming for

the special version of this problem, in which all cells are of equal cost. The papers [16, 19] studied how to minimize the expected number of paged cells given an *average* (as opposed to worst-case) delay constraint using relaxation to a continuous model [19] or with a weakly polynomial dynamic programming solution [16]. The problem of paging more than one user for a conference call was studied in [5, 7, 12–14]. The problem of online paging a mobile user (in contrast to predetermined offline paging) was studied in [6]. The paper [10] explored a similar problem in which the order of cells is dictated in the context of TTL flood searching in sensor networks. The problem of paging mobile users with an inaccurate information of the user location probabilities was studied in [4].

Contributions: To the best of our knowledge, this generalized version of the problem, *i.e.*, w_i being an arbitrary number for each C_i , was not studied prior to our work. Indeed, the algorithms that generate optimal search strategies when $w_i = 1$ may provide very bad performance because they ignore the different costs. Thus, our goal is to explore different solutions and approaches for the general case of arbitrary cost values.

We start with three interesting and important special cases for which polynomial time algorithms that produce optimal search strategies exist. In the first, the searcher has no a priori knowledge about the location of the token and therefore all the probabilities are the same. We show that this case is “dual” to the traditional case in which all the costs are the same. Therefore, the known optimal dynamic programming solutions can be applied to this case as well. In the second, there are no delay constraints and the searcher may search the token in N rounds. We show that ordering the boxes by a non-increasing order of the ratio p_i/w_i implies an optimal sequential search. In the third, the token is *atypical* in the sense that it is more likely to be placed in boxes with low unlocking costs. We show that applying the known dynamic programming solutions on the boxes ordered by the non-increasing order of the ratio p_i/w_i yields an optimal search strategy. We denote this algorithm by FRO (Follow Ratio Order).

Next, we consider the case of a token that has a higher probability of being placed in “expensive” boxes (in Cellular Networks, this corresponds to a mobile user who follows the massive behavior of other users). Let a *typical* token be a token for which p_i is proportional to w_i . We show that for a typical token the problem has a Polynomial Time Approximation Scheme (PTAS). This is best possible because we also show that the problem is strongly NP-hard already for $D = 2$ and therefore a Fully PTAS (FPTAS) is impossible. We also show that the problem for a typical token is similar to a known version of the load balancing problem. As a result, we analyze the natural greedy solution borrowed from the load balancing problem and show similar results to those that were known for the load balancing problem.

Next we address the case of $D = 2$ rounds. For an arbitrary token that may be placed in any box with any probability, we show that the FRO strategy has a tight guaranteed $8/7 \approx 1.143$ approximation ratio. For a typical token, we design another solution that implies a slightly better tight guaranteed $\frac{7-2\sqrt{7}}{28-10\sqrt{7}} \approx 1.108$ approximation ratio.

We complement our theoretical results with simulations for Cellular Networks, given in the appendix. First, we analyze real data from Shenzhen, China, provided by China Unicom that includes 171929 appearances of 996 users in 5625 cells. We show, by cross validation, that the cost vector (determined by the congestion) and the users’ probability vectors (determined by frequencies of appearances in cells) follow the *Zipf* distribution. We mainly explore the cases for which we do not have concrete theoretical results. We cover the ranges of D between 2 and N and types of users between typical and atypical users. We implement the FRO algorithm and an adaptation of the greedy load balancing algorithm denoted by \mathcal{S} [9]. We compare these polynomial time algorithms

with the optimal exponential time solution for small values of N and with the sequential optimal strategy for larger values of N . We test the algorithms on random data (that either follows the *Zipf* distribution or the uniform distribution) and real data. We report that for all values of D and for all types of users, FRO performs well and that the theoretical bound $8/7$ is a pessimistic bound. On the other hand, \mathcal{S} performs competitively only for typical or almost typical users since it “ignores” the paging costs.

Paper organization: In most of the paper we will use the token-boxes terminology. We will use the users-cells terminology to describe our simulations. Section 2 provides formal definitions and some preliminaries. Section 3 presents three cases for which optimal polynomial time algorithms exist. Section 4 analyzes the case of typical tokens. Section 5 analyzes the performance FRO for two rounds; some of the relevant technical proofs appear in Appendices A and B. Section 6 concludes with some open problems. Appendix C reports the results of the simulation work on real data and simulated data.

2 Preliminaries

Denote the N boxes by C_1, C_2, \dots, C_N . Let $\mathbf{p} = \langle p_1, p_2, \dots, p_N \rangle$ be the vector of independent probabilities of the token being placed in these boxes respectively. Let $\mathbf{w} = \langle w_1, w_2, \dots, w_N \rangle$ be the vector of costs of unlocking these boxes respectively. Denote the delay constraint for finding the token by D , $1 \leq D \leq N$. An *instance* to the problem is the quadruple $I = (N, D, \mathbf{w}, \mathbf{p})$. An instance I is a *uniform cost instance* if $w_1 = \dots = w_N = 1/N$ and a *uniform probability instance* if $p_1 = \dots = p_N = 1/N$.

A *search strategy* $\mathcal{A} = \langle A_1, \dots, A_D \rangle$ is an ordered D -partition of the boxes, such that in the i th round, all the boxes in the set A_i are unlocked. The search process terminates in round d if the token is found in one of the boxes of the set A_d . For a given search strategy \mathcal{A} and a round $1 \leq d \leq D$, the *round probability* is $P_d = \sum_{C_i \in A_d} p_i$ and the *round cost* is $W_d = \sum_{C_i \in A_d} w_i$. The cost of search strategy \mathcal{A} on an instance $I = (N, D, \mathbf{w}, \mathbf{p})$ is denoted by $\text{cost}(\mathcal{A}, I)$ and when the definition of I is clear it is denoted by $\text{cost}(\mathcal{A})$.

Proposition 1. *The following are two different but equivalent ways to compute the search cost of a search strategy $\mathcal{A} = \langle A_1, \dots, A_D \rangle$ on the instance $I = (N, D, \mathbf{w}, \mathbf{p})$:*

$$\text{cost}(\mathcal{A}, I) = \sum_{d=1}^D \left(P_d \sum_{i=1}^d W_i \right) \qquad \text{cost}(\mathcal{A}, I) = \sum_{d=1}^D \left(W_d \sum_{i=d}^D P_i \right) \quad (1)$$

Proof. The first equation follows since with probability P_d the token is found during the d th round and the strategy pays the cost of the first d sets A_1, \dots, A_d of the partition. The second equation follows since the strategy pays the cost of the d th round only if the token is in a box belonging to the last $(D - d + 1)$ sets A_d, \dots, A_D of the partition. \square

An *optimal search strategy* is a search strategy \mathcal{O} such that $\text{cost}(\mathcal{O})$ is the minimum among all possible search strategies. An *optimal algorithm* is an algorithm that generates optimal search strategies for all possible instances. Let OPT be an optimal algorithm and let algorithm ALG be another search algorithm. ALG is a $(1 + \varepsilon)$ -approximation, if $\text{cost}(\text{ALG})/\text{cost}(\text{OPT}) \leq (1 + \varepsilon)$ for any instance.

Normalizing the cost and the probability vectors: We observe the following basic fact that allows us to assume without loss of generality that $\sum_{i=1}^N w_i = 1$ and that $\sum_{i=1}^N p_i = 1$.

Proposition 2. *Let \mathcal{O} be an optimal search strategy on boxes with probabilities p_i and costs w_i , $1 \leq i \leq N$. \mathcal{O} is also an optimal search strategy on boxes with probabilities p'_i and costs w'_i ,*

if $w'_i = c_w \cdot w_i$ and $p'_i = c_p \cdot p_i$, where $1 \leq i \leq N$ and c_w and c_p are positive constants. The approximation ratio of any non-optimal solution is retained, too.

Types of tokens: We classify tokens by their *typicality*. In one extreme, the probability vector of a *typical token* is proportional to its cost vector. By Proposition 2, without loss of generality, for a typical token, $p_i = w_i$ for any box C_i . In the other extreme, an *atypical token* is more likely to be located in lower cost boxes. Formally, the costs and the probabilities are in opposite order. A token with no typicality association is called an *arbitrary* token. Such a token may be located in any box C_i with arbitrary cost w_i and arbitrary probability p_i .

Optimal polynomial time algorithms: We say that an ordered partition $\mathcal{A} = \langle A_1, \dots, A_D \rangle$ respects the order of boxes $\langle C_1, C_2, \dots, C_N \rangle$ if there are no A_i, A_j with $i < j$, such that $C_{i'} \in A_i$ and $C_{j'} \in A_j$ with $i' > j'$. Given an order of the boxes, one can find a minimum cost partition that respects that order in polynomial time by slightly modifying the dynamic programming methods described in [3, 15–17, 19] to include costs of boxes. A naive implementation implies an $O(N^2D)$ algorithm. We next show a more efficiently implementation. Proof is omitted.

Theorem 1. *The dynamic programming scheme from [3] can be implemented in $\Theta(ND)$ time to find a minimum cost partition that respects a given order on the boxes.*

Unfortunately, for the general problem, as we will prove later, it is impossible to find in polynomial time the order of boxes in an optimal search strategy unless $P = NP$.

Algorithm FRO: Consider the N boxes ordered by the non-increasing order of the p_i/w_i ratio. That is, $p_1/w_1 \geq p_2/w_2 \geq \dots \geq p_N/w_N$. We call the algorithm that computes the minimum cost partition that respects the above order Algorithm FRO (follow ratio order).

3 Special Cases with Optimal Search Strategies

The traditional version of the problem with uniform cost instances can be solved by polynomial time algorithms. In this section, we present three additional special cases for which optimal search strategies can be found with polynomial time algorithms. The first case is when the instances are uniform probability instances; the second case is when there are no delay constraints ($D = N$); and the third case is when the token is atypical.

3.1 Uniform probabilities

Assume that the cost of unlocking each box is $1/N$. Several optimal polynomial time dynamic programming schemes solve this case [3, 15, 16, 19]. In all of these solutions, the boxes are ordered following a non-increasing order of p_i . Then, recursively the minimum cost of searching in the first (or last) n boxes in d rounds is computed based on the optimal solutions for $d - 1$ rounds and any number of boxes that were computed in earlier stages. The other uniform case is when all the probabilities are $1/N$. This is the case when the searcher has no clue for the whereabouts of the token but still knows the unlocking costs associated with the boxes. In the next lemma, we show the duality of this case to the uniform costs case.

Let $I = (\mathbf{w} = \langle w_1, \dots, w_N \rangle, \mathbf{p} = \langle 1/N, \dots, 1/N \rangle)$ be a uniform probability instance such that $w_1 \leq \dots \leq w_N$. Define its *dual* instance, as $I' = (\mathbf{w}' = \langle \frac{1}{N}, \dots, \frac{1}{N} \rangle, \mathbf{p}' = \langle p'_1 = w_N, \dots, p'_N = w_1 \rangle)$. I' is a uniform cost instance whose probability vector is ordered in a non-decreasing order $p'_1 \geq \dots \geq p'_N$.

Lemma 1. *Let $\mathcal{A} = \langle A_1, \dots, A_D \rangle$ be a partition and let $\mathcal{A}' = \langle A'_1 = A_D, \dots, A'_D = A_1 \rangle$ be its reverse partition. Then $\text{cost}(\mathcal{A}, I) = \text{cost}(\mathcal{A}', I')$.*

Proof. By the left equation in (1), the cost of \mathcal{A} for the uniform probability instance is

$$\text{cost}(\mathcal{A}, I) = \frac{|A_1|}{N}W_1 + \frac{|A_2|}{N}(W_1 + W_2) + \cdots + \frac{|A_D|}{N}(W_1 + \cdots + W_D).$$

By the right equation in (1), the cost of \mathcal{A}' for the uniform cost instance is

$$\text{cost}(\mathcal{A}', I') = \frac{|A'_1|}{N}(P'_1 + \cdots + P'_D) + \cdots + \frac{|A'_{D-1}|}{N}(P'_{D-1} + P'_D) + \frac{|A'_D|}{N}P'_D.$$

The lemma follows since A'_d is the set A_{D+1-d} and since $P'_d = W_{D+1-d}$ for all $1 \leq d \leq D$. \square

The above lemma proves that there exists a “cost preserving” 1 – 1 mapping between the set of all uniform probability instances and the set of all uniform cost instances. Therefore,

Corollary 1. *Let I be a uniform probability instance and let I' be its dual uniform cost instance. Let \mathcal{O}' be the optimal search strategy for I' . Then the reverse partition \mathcal{O} of the partition \mathcal{O}' is an optimal search strategy for I .*

The above corollary gives a polynomial time reduction from the uniform probabilities case to the uniform costs case. Therefore, the known optimal polynomial time solutions to the uniform costs case apply also to the uniform probabilities case.

3.2 $D = N$ rounds

The following is a necessary condition for any optimal search strategy.

Lemma 2. *Let $\mathcal{A} = \langle A_1, \dots, A_D \rangle$ be an optimal solution for an arbitrary instance I . Then $P_d/W_d \geq P_{d+1}/W_{d+1}$, for $1 \leq d < D$.*

Proof. Fix d , $1 \leq d < D$. Define another search strategy $\mathcal{A}' = \langle A_1, \dots, A_{d-1}, A_{d+1}, A_d, \dots, A_D \rangle$ that is obtained from \mathcal{A} by swapping A_d and A_{d+1} . By Proposition 1, it follows that $\text{cost}(\mathcal{A}') - \text{cost}(\mathcal{A}) = P_d W_{d+1} - P_{d+1} W_d$. This is because in both strategies all the costs that incurred by W_i for $i < d$ and $i > d + 1$ are the same and in both strategies the terms $P_d W_d$ and $P_{d+1} W_{d+1}$ are part of the cost. Now, since \mathcal{A} is optimal, it follows that $P_d W_{d+1} - P_{d+1} W_d \geq 0$, which is equivalent to $P_d/W_d \geq P_{d+1}/W_{d+1}$. \square

Corollary 2. *For $D = N$, the optimal search strategy unlocks the boxes in a non-increasing order of p_i/w_i , one box per round.*

The above corollary implies that the polynomial time algorithm FRO generates an optimal search strategy for the case $D = N$.

3.3 Atypical tokens

For an atypical token, the non-increasing order of the probability vector corresponds to the non-decreasing order of the cost vector. The following lemma shows that an optimal search strategy can be found in polynomial time for such tokens.

Lemma 3. *If there is an order of the boxes such that $p_1 \geq \dots \geq p_N$ while $w_1 \leq \dots \leq w_N$, then algorithm FRO follows this order and generates an optimal search strategy.*

Proof. Assume towards a contradiction that there exists an optimal solution $\mathcal{A} = \langle A_1, \dots, A_D \rangle$ that does not respect the non-decreasing order of the ratios p_i/w_i . This implies the existence of indices $j < i$ (and therefore $p_i < p_j$ and $w_i > w_j$) such that $C_i \in A_d$ and $C_j \in A_{d+1}$ for some $1 \leq d < D$. Define a new partition \mathcal{A}' which is almost identical to \mathcal{A} except that C_i and C_j are swapped. To get a contradiction, we show that the cost of \mathcal{A}' is smaller than the cost of \mathcal{A} .

Let $P_d = P'_d + p_i$, $W_d = W'_d + w_i$, $P_{d+1} = P'_{d+1} + p_j$, and $W_{d+1} = W'_{d+1} + w_j$. A careful examination of the terms in Eq. (1) from Proposition 1 reveals that both $\text{cost}(\mathcal{A})$ and $\text{cost}(\mathcal{A}')$ share some identical terms while $\text{cost}(\mathcal{A})$ has unique terms $P'_d w_i + p_j w_i + p_j W'_{d+1}$ and $\text{cost}(\mathcal{A}')$ has unique terms $P'_d w_j + p_i w_j + p_i W'_{d+1}$. Now, since $p_i < p_j$ and $w_i > w_j$, it follows that $\text{cost}(\mathcal{A}') < \text{cost}(\mathcal{A})$. \square

4 Typical Tokens

In this section, we discuss and analyze search strategies for the special case of typical tokens. Recall that a typical token is more likely to be found in high cost boxes and therefore, after normalizing the cost and probability values, we assume that $p_i = w_i$ for all $1 \leq i \leq N$. We prove that the problem is strongly NP-Hard even in this special case. In particular, we show that the problem is essentially a variation of a known load balancing problem. This enables us to apply a known *polynomial time approximation scheme* (PTAS) solution for the load balancing problem to our problem of searching for typical tokens. In addition, we analyze the performance of a natural load balancing greedy algorithm applied to the search problem. We first show how to compute the cost for a typical token by simplifying the equations from Proposition 1.

Proposition 3. *Let $I = (N, D, \mathbf{p}, \mathbf{p})$ be an instance of the typical token problem and $\mathcal{A} = \langle A_1, \dots, A_D \rangle$ be a search strategy. Then, $\text{cost}(\mathcal{A}, I) = 1/2 + \sum_{d=1}^D (P_d)^2/2$*

Proof. For all $1 \leq d \leq D$, denote by $P_d = W_d$ the probability as well as the cost of the set A_d . By Proposition 1,

$$\text{cost}(\mathcal{A}, I) = \sum_{d=1}^D \left(P_d \sum_{i=1}^d W_i \right) = \sum_{1 \leq i \leq j \leq D} P_i P_j = \frac{1}{2} \left(\left(\sum_{d=1}^D P_d \right)^2 + \sum_{d=1}^D (P_d)^2 \right)$$

The proposition follows since $\sum_{d=1}^D P_d = 1$. \square

The above proposition implies that for a typical token and a given search strategy \mathcal{A} , the cost is the same regardless of the order of the D sets in \mathcal{A} . That is, the task of finding an efficient D -round search strategy becomes the task of partitioning the N boxes into D sets while minimizing a particular cost function. One can view the sets as machines and the probabilities as the processing times of tasks. Then the problem is almost identical to the known *load balancing* problem from [9,11] whose goal is to minimize the L_2 norm of the tasks' completion time on all the machines. Formally,

Definition 1. *Let $\{T_1, \dots, T_N\}$ be N tasks and $\{M_1, \dots, M_D\}$ be D machines. The processing time for T_i on any machine is p_i . The load balancing problem is to find an allocation of tasks to machines that minimizes the L_2 norm of the completion time on all machines $\sum_{d=1}^D (\sum_{T_i \in M_d} p_i)^2$.*

If X is the optimization goal of the load balancing problem and Y is the optimization goal of the search problem, then $X = (1 + Y)/2$ by Proposition 3 and Definition 1. As a result, we now show that any approximation algorithm to the load balancing problem is also an approximation algorithm to the searching for typical tokens problem with an even smaller approximation factor.

Lemma 4. *Let ALG be a $(1 + \varepsilon)$ -approximation algorithm to the load balancing problem where $(1 + \varepsilon)$ is a tight bound. Then ALG is a $(1 + \varepsilon')$ -approximation algorithm to the searching for typical tokens problem, where $\varepsilon' < \frac{\varepsilon}{2}$ for all instances and $\varepsilon' \geq \frac{\varepsilon}{D+1}$ for some instance.*

Proof. Let OPT be an optimal solution to the load balancing problem. Let \mathcal{O} and \mathcal{A} be the L_2 norm of solutions OPT and ALG , respectively. Let OPT' be an optimal solution to the search typical tokens problem. Let $ALG' = ALG$ be a $(1 + \varepsilon')$ -approximation solution to the search problem. Let $\mathcal{O}' = \text{cost}(OPT')$ and $\mathcal{A}' = \text{cost}(ALG')$. By definition, we have $\mathcal{A} \leq (1 + \varepsilon)\mathcal{O}$ for all instances in

the load balancing problem. Therefore,

$$1 + \varepsilon' = \frac{\mathcal{A}'}{\mathcal{O}'} = \frac{0.5 + 0.5\mathcal{A}}{0.5 + 0.5\mathcal{O}} \leq \frac{1 + (1 + \varepsilon)\mathcal{O}}{1 + \mathcal{O}} = 1 + \frac{\mathcal{O}}{1 + \mathcal{O}}\varepsilon \text{ for all instances.}$$

Since $(1 + \varepsilon)$ is a tight bound to the load balancing problem, we have $\mathcal{A} = (1 + \varepsilon)\mathcal{O}$ for some instance. Therefore,

$$1 + \varepsilon' = \frac{\mathcal{A}'}{\mathcal{O}'} = \frac{0.5 + 0.5\mathcal{A}}{0.5 + 0.5\mathcal{O}} = \frac{1 + (1 + \varepsilon)\mathcal{O}}{1 + \mathcal{O}} = 1 + \frac{\mathcal{O}}{1 + \mathcal{O}}\varepsilon \text{ for some instance.}$$

Since all $P_i \geq 0$ and $\sum_{i=1}^D P_i = 1$, it follows that the cost $\sum_{i=1}^D P_i^2$ (Definition 1) of any solution to the load balancing problem is greater than $1/D$ and smaller than 1. In particular $1/D \leq \mathcal{O} \leq 1$ and therefore $\mathcal{O}/(1 + \mathcal{O}) \leq 1/2$ for all instances and $\mathcal{O}/(1 + \mathcal{O}) \geq 1/(D + 1)$ for some instance. \square

The paper [2] pointed out that the load balancing problem is strongly NP-Hard. The next theorem follows since Lemma 4 gives a lower bound to the approximation ratio to the searching for typical tokens problem without N as a parameter.

Theorem 2. *For any $N > D \geq 2$, the search problem is strongly NP-Hard even for a typical token.*

Constant approximation ratio algorithms to the load balancing problem have been introduced in [9, 11] and a PTAS to this problem has been presented in [2]. The following theorem is another corollary of Lemma 4.

Theorem 3. *For any $N > D \geq 2$, there exists a constant approximation ratio algorithm and a PTAS to the searching for typical tokens problem.*

Algorithm \mathcal{S} in [9] is a natural greedy algorithm. Essentially, it allocates boxes one by one in a non-increasing order of p_i to the set A_d currently with the smallest P_d . In some cases, like when $D = 2$, one can compute exactly the approximation ratio of \mathcal{S} for both the load balancing and the searching for typical tokens problems. We omit the technical details. The specific approximation

Case	Load Balancing	Searching Typical Tokens
$N \leq 4$, any D	1	1
$D = 2$, $N \geq 5$	tight ≈ 1.0285	tight ≈ 1.0143
$D = 3$, $N \geq 5$	$[83/81, 25/24]$	$[326/324, 49/48]$
even $D \geq 4$, $N \geq 5$	$[37/36, 25/24]$	$[1 + 1/36(D + 1), 49/48]$
odd $D \geq 5$, $N > 5$	$[1 + (D - 1)/36D, 25/24]$	$[1 + (D - 1)/36D(D + 1), 49/48]$

Table 1: Approximation ratio of \mathcal{S} -algorithm

ratios for different values of N and D of Algorithm \mathcal{S} for the load balancing problem and for the searching for typical tokens problem (by lemma 4) are shown in Table 1. In the table, the meaning of $[x, y]$ is that there exists a y approximation factor and there cannot be an approximation factor smaller than x .

5 Algorithm FRO for Two Rounds

The main result of this section is the analysis of the performance of Algorithm FRO for two rounds on arbitrary tokens whose cost vector has no correlation to the probability vector. Note that the problem for an arbitrary token is strongly NP hard when $2 \leq D < N$ because it is strongly NP-hard

already for a typical token. Therefore, our goal is to prove a guaranteed approximation factor. We show that the approximation ratio of FRO is $8/7 \approx 1.143$. For typical tokens, we slightly improve the ratio to $\frac{7-2\sqrt{7}}{28-10\sqrt{7}} \approx 1.108$. For both cases, we provide instances that show that these bounds are tight.

5.1 Arbitrary tokens

Theorem 4. *Algorithm FRO has an approximation ratio $8/7$ when $D = 2$ and $N > 2$, and the ratio $8/7$ is attainable.*

Proof. First, we show that an upper bound on the approximation ratio of FRO for $N = 3, 4$ is also an upper bound for all $N > 2$ (Lemma 5). Next, we prove that the upper bound on the approximation ratio of FRO for $N = 3, 4$ is $8/7$ (Lemmas 6, 7). Finally, for every N , we construct an instance with approximation ratio $8/7$ (Lemma 8). \square

Lemma 5. *For each instance with $N \geq 4$ for which the approximation ratio of FRO is ρ , there exists an instance, with either $N = 3, 4$, for which the approximation ratio of FRO is at least ρ .*

Proof. Let the instance with $N \geq 4$ consist of boxes $\{C_1, \dots, C_N\}$. Let the OPT partition be $\langle X, Y \rangle$ and let the FRO partition be $\langle X', Y' \rangle$. Define the following four subsets of the N boxes: $A = X \cap X'$, $B = X' \setminus X$, $C = Y' \setminus Y$, $D = Y \cap Y'$. Hence, by definition $\text{OPT} = \langle A \cup C, B \cup D \rangle$ and $\text{FRO} = \langle A \cup B, C \cup D \rangle$.

If all four sets A, B, C, D are not empty, we construct an instance with $N = 4$ boxes as follows. The boxes are $\{C_A, C_B, C_C, C_D\}$, such that $p_I = \sum_{c_i \in I} p_i$ and $w_I = \sum_{c_i \in I} w_i$, where $i = 1 \dots N$ and $I \in \{A, B, C, D\}$. Let OPT_4 be the optimal solution for $N = 4$. It follows that $\text{cost}(\text{OPT}_4) = \text{cost}(\text{OPT}_N)$. This is because; (i) $\text{cost}(\text{OPT}_4)$ cannot be less than $\text{cost}(\text{OPT}_N)$, otherwise OPT_N would not be optimal (taking the corresponding OPT_4 partition); (ii) $\text{cost}(\text{OPT}_4)$ can reach $\text{cost}(\text{OPT}_N)$ by taking $\text{OPT}_4 = \langle \{C_A C_C\}, \{C_B C_D\} \rangle$ (by definition of partition $\{A, B, C, D\}$). It also follows that $\text{FRO}_4 = \langle \{C_A C_B\}, \{C_C C_D\} \rangle$ because otherwise, FRO_N will not be $\langle A \cup B, C \cup D \rangle$ by definition. Thus, $\text{cost}(\text{FRO}_4) \geq \text{cost}(\text{FRO}_N)$. Therefore, $\text{cost}(\text{FRO}_N)/\text{cost}(\text{OPT}_N) \leq \text{cost}(\text{FRO}_4)/\text{cost}(\text{OPT}_4)$.

If either of the four sets A, B, C, D is empty, we can similarly construct an instance of $N = 3$ with at least the same approximation ratio. Finally, if there are at least two empty sets among A, B, C, D , then OPT is FRO , and thus the approximation ratio is 1. \square

Corollary 3. *An upper bound ρ on the approximation ratio of FRO for all instances with $N = 3, 4$ is also an upper bound on the approximation ratio for all instances with $N > 4$.*

It remains to prove the following two lemmas to complete the proof of Theorem 4. The full proofs of Lemma 6 and Lemma 7 are long involving tedious case analysis. We present case analysis and the proofs of some of the cases in Appendices A and B respectively.

Lemma 6. *For all instances of $N = 3$, the approximation ratio $\rho = \text{cost}(\text{FRO})/\text{cost}(\text{OPT}) \leq 8/7$.*

Lemma 7. *For all instances of $N = 4$, the approximation ratio $\rho = \text{cost}(\text{FRO})/\text{cost}(\text{OPT}) \leq 8/7$.*

In the next lemma we show that the $8/7$ upper bound on the approximation ratio is tight.

Lemma 8. *The approximation ratio $8/7$ is attainable for any $N > 2$.*

Proof. Consider the instance with $p_1 = 1/4$, $p_2 = 3/4$, $p_3 = \dots = p_N = 0$ and $w_1 = 1/5$, $w_2 = 3/5$, $w_3 = \dots = w_N = 1/(5(N-2))$. FRO outputs the partition $(1|23 \dots N)$ the cost of which is $4/5$ while the optimal partition is $(2|13 \dots N)$ the cost of which is $7/10$. The ratio is $8/7$. \square

5.2 Typical tokens

For a typical token, FRO does not distinguish among the boxes since all the ratios are 1. Therefore, we assume that an adversary picks the worst permutation on the boxes that is respected by FRO. Still, we can prove a smaller guaranteed approximation factor for this case.

Theorem 5. *Algorithm FRO has an approximation ratio $\frac{7-2\sqrt{7}}{28-10\sqrt{7}} \approx 1.108$ when $D = 2$ and $N > D$, and this ratio is attainable.*

Proof. (outline) The proof is very similar to the proof of Theorem 4. First, the reduction lemma, Lemma 5, is correct for any type of token. Then the proofs of the equivalent Lemmas (but with a different ratio) to Lemma 6 and Lemma 7 are easier since there are fewer variables. Finally, the next lemma demonstrates the instance for which the ratio is attainable. \square

Lemma 9. *The approximation ratio $\frac{7-2\sqrt{7}}{28-10\sqrt{7}}$ is attainable for any $N > 2$.*

Proof. For simplicity, assume that $0/0 = 1$ since one can replace each zero value with a small ε . Consider the instance with $p_1 = w_1 = x$, $p_2 = w_2 = 1 - 2x$, $p_3 = w_3 = x$, and $p_4 = w_4 = \dots = p_N = w_N = 0$. FRO, that respects this order, outputs the partition $(1|23\dots N)$ the cost of which is $1 - x + x^2$ while the optimal partition is $(2|13\dots N)$ the cost of which is $1 - 2x + 4x^2$. The maximum of the ratio is achieved for $x = (3 - \sqrt{7})/2$ and is $\frac{7-2\sqrt{7}}{28-10\sqrt{7}}$. \square

6 Open Problems

We conjecture that the $8/7$ bound for the FRO algorithm holds also for $D > 2$. Similarly to the $D = 2$ case, we can reduce the problem instances with any N to instances with $D \leq N \leq D^2$. However, we do not know how to handle all these cases, even for $D = 3$. We only know how to resolve an instance with $N = 4$ and $D = 3$ showing a $8/7$ lower bound for algorithm FRO.

Are there better strategies than FRO? By examining the instance that forces FRO to have an $8/7$ approximation factor, one can conclude that a possible improvement is to give more priority to boxes with high probability or boxes with low cost even if their ratio is not that large. Preliminary experiments indicate that such modifications gain no significant improvement. Further experimental study is desired for similar strategies or for other strategies.

Is there a way to “measure” the typicality of a token? If so, can we design more competitive algorithms to find tokens with certain level of typicality? Recall, that we have an optimal solution for an atypical token and a PTAS solution for a typical token.

The uniform cost case was investigated also in many settings for paging multiple users in Cellular Networks. These settings of finding more than one hidden token could be addressed in the non-uniform case as well.

References

- [1] AKYILDIZ, I. F., MCNAIR, J., HO, J., UZUNALIOGLU, H., AND WANG, W. Mobility management in next-generation wireless systems. In *Proc. IEEE* (1999), pp. 1347–1384.
- [2] ALON, N., AZAR, Y., WOEGINGER, G. J., AND YADID, T. Approximation schemes for scheduling on parallel machines. *J. Scheduling* 1, 1 (1998), 55–66.
- [3] BAR-NOY, A., FENG, Y., AND GOLIN, M. J. Paging mobile users efficiently and optimally. In *Proc. IEEE Conference on Computer Communications* (2007), pp. 1910–1918.
- [4] BAR-NOY, A., AND KLUKOWSKA, J. Finding mobile data: efficiency vs. location inaccuracy. In *Proc. Annual European Symposium on Algorithms (ESA)* (2007), pp. 111–122.
- [5] BAR-NOY, A., AND MALEWICZ, G. Establishing wireless conference calls under delay constraints. *J. Algorithms* 51, 2 (2004), 145–169.
- [6] BAR-NOY, A., AND MANSOUR, Y. Competitive on-line paging strategies for mobile users under delay constraints. In *Proc. ACM Symposium on Principles of Distributed Computing (PODC)* (2004), pp. 256–265.
- [7] BAR-NOY, A., AND NAOR, Z. Efficient multicast search under delay and bandwidth constraints. *Wireless Networks* 12, 6 (2006), 747–757.
- [8] BUCHANAN, M. Ecological modelling: the mathematical mirror to animal nature. *Nature* 453 (2008), 714–716.
- [9] CHANDRA, A. K., AND WONG, C. K. Worst-case analysis of a placement algorithm related to storage allocation. *SIAM J. Comput.* 4, 3 (1975), 249–263.
- [10] CHANG, N. B., AND LIU, M. Revisiting the TTL-based controlled flooding search: optimality and randomization. In *Proc. 10th Annual International Conference on Mobile Computing and Networking (MOBICOM)* (2004), pp. 85–99.
- [11] CODY, R. A., AND COFFMAN, E. G. Record allocation for minimizing expected retrieval costs on drum-like storage devices. *J. ACM* 23, 1 (1976), 103–115.
- [12] EPSTEIN, L., AND LEVIN, A. The conference call search problem in wireless networks. *Theor. Comput. Sci.* 359, 1-3 (2006), 418–429.
- [13] EPSTEIN, L., AND LEVIN, A. A PTAS for delay minimization in establishing wireless conference calls. *Discrete Optimization* 5, 1 (2008), 88–96.
- [14] GAU, R.-H., AND HAAS, Z. J. Concurrent search of mobile users in cellular networks. *IEEE/ACM Trans. Netw.* 12, 1 (2004), 117–130.
- [15] GOODMAN, D. J., KRISHNAN, P., AND SUGLA, B. Minimizing queuing delays and number of messages in mobile phone location. *Mobile Netw. and Appl.* 1, 1 (1996), 39–48.
- [16] KRISHNAMACHARI, B., GAU, R.-H., WICKER, S. B., AND HAAS, Z. J. Optimal sequential paging in cellular wireless networks. *Wireless Netw.* 10, 2 (2004), 121–131.
- [17] MADHAVAPEDDY, S., BASU, K., AND ROBERTS, A. *Adaptive paging algorithms for cellular systems*, vol. 1. Kluwer Academic Publishers, Norwell, MA, USA, 1996, pp. 83–101.
- [18] ROSE, C. State-based paging/registration: a greedy technique. *IEEE Trans. Veh. Tech.* 48, 1 (January 1999), 166–173.
- [19] ROSE, C., AND YATES, R. D. Minimizing the average cost of paging under delay constraints. *Wireless Netw.* 1, 2 (1995), 211–219.
- [20] ROSE, C., AND YATES, R. D. Ensemble polling strategies for increased paging capacity in mobile communication networks. *Wireless Netw.* 3, 2 (1997), 159–167.
- [21] TSG/WG. *Section 7.6.5.18, 3GPP TS 09.02 Mobile Application Part (MAP) Specification, ver. 8.8.1*, 2008.

A Proof of Lemma 6

For $N = 3$, $D = 2$, assume boxes C_1, C_2, C_3 , are sorted in decreasing order of p_i/w_i . The only possible FRO ordered partitions are $\langle \{C_1, C_2\}, \{C_3\} \rangle$ and $\langle \{C_1\}, \{C_2, C_3\} \rangle$. We use the shorthand notation 12|3 and 1|23, respectively, for the above ordered partitions. The only possible OPT partitions that are not of the FRO form are 2|13 and 13|2, because 3|12 has cost at least as much as 12|3, and 23|1 has cost at least as much as 1|23.

The costs for the above mentioned partitions are:

$$\begin{aligned} \text{cost}(12|3) &= 1 - p_1w_3 - p_2w_3 & \text{cost}(1|23) &= 1 - p_1w_2 - p_1w_3 \\ \text{cost}(2|13) &= 1 - p_2w_1 - p_2w_3 & \text{cost}(13|2) &= 1 - p_1w_2 - p_3w_2 \end{aligned}$$

Therefore we have to compute the worst ratio for four possible cases of FRO and OPT ordered partitions:

$$\begin{aligned} \text{FRO : } 12|3 \text{ and OPT : } 2|13 & & \text{FRO : } 12|3 \text{ and OPT : } 13|2 \\ \text{FRO : } 1|23 \text{ and OPT : } 2|13 & & \text{FRO : } 1|23 \text{ and OPT : } 13|2 \end{aligned}$$

Because of the space limitation, we prove the first of the four cases. The rest can be proved by similar methods.

Proof. Assume FRO: 12|3 and OPT: 2|13. Since FRO outputs 12|3, the cost of 12|3 is not worse than the cost of the other FRO ordered partition 1|23, which implies

$$p_1w_2 \leq p_2w_3. \quad (2)$$

The ratio $\text{cost}(\text{FRO})/\text{cost}(\text{OPT})$ is:

$$\rho = \frac{1 - p_1w_3 - p_2w_3}{1 - p_2w_1 - p_2w_3}.$$

Let $d = (p_1w_2 - p_2w_1)/(w_1 + w_2) \geq 0$. Consider a new input, by substituting p_1 with $p_1 - d$ and p_2 with $p_2 + d$. In the new input, equation (2) still holds and therefore FRO outputs the same partition. Moreover, in the new input, OPT outputs the same partition too. The ratio ρ' of the new input is at least the ratio of the original input, because:

$$\rho' \geq \rho \iff (p_1w_2 - p_2w_1)(w_1 + w_3)(1 - p_1w_3 - p_2w_3) \geq 0$$

and each factor of the product $(p_1w_2 - p_2w_1)(w_1 + w_3)(1 - p_1w_3 - p_2w_3)$ is non-negative.

The above transformation of p_1 to $p_1 - d$ and p_2 to $p_2 + d$ has the effect of creating a new input in which the p/w ratio of boxes 1 and 2 is the same. It is possible to solve exactly the optimization problem of maximizing the ratio of FRO over OPT under this additional constraint. The solution of the optimization problem gives a maximum ratio of $8/7$ for input: $\mathbf{p} = \langle 1/4, 3/4, 0 \rangle$, $\mathbf{w} = \langle 1/5, 3/5, 1/5 \rangle$. In the following we sketch a proof of $\rho \leq 8/7$. Since boxes 1 and 2 have the same p/w ratio, we write $p_1 = rw_1$ and $p_2 = rw_2$. We want to prove that the approximation ratio is:

$$\frac{1 - rw_1w_3 - rw_2w_3}{1 - rw_1w_2 - rw_2w_3} \leq \frac{8}{7}$$

or equivalently $r(8w_1w_2 - 7w_1w_3 + w_2w_3) \leq 1$. Therefore, it is enough to study the maximization problem:

$$\begin{aligned} & \text{maximize } r(8w_1w_2 - 7w_1w_3 + w_2w_3) \text{ under:} \\ & w_1 + w_2 + w_3 = 1 \wedge w_3 \geq w_1 \wedge w_2 \geq w_3 \wedge rw_1 + rw_2 \leq 1. \end{aligned}$$

The constraint $w_3 \geq w_1$ is implied by $p_1w_2 \leq p_2w_3$, the constraint $w_2 \geq w_3$ is implied by $p_2w_1 \leq p_1w_3$ (which is implied by $\text{cost}(2|13) \leq \text{cost}(12|3)$), the constraint $rw_1 + rw_2 \leq 1$ is implied by $p_1 + p_2 \leq 1$. Writing

$w_3 = 1 - w_1 + w_2$ and by the fact that $r \leq 1/(w_1 + w_2)$ we solve the following maximization problem (which might have a bigger solution):

$$\begin{aligned} & \text{maximize } \frac{8w_1w_2 - 7w_1(1 - w_1 - w_2) + w_2(1 - w_1 - w_2)}{w_1 + w_2} \\ & \text{under: } w_1 \geq 0 \wedge 1 - w_2 \geq 2w_1 \wedge 2w_2 \geq 1 - w_1. \end{aligned}$$

The value of the maximization function depends on the values of w_1, w_2 . The range of w_1, w_2 is a triangle in the plane \mathbb{R}^2 . We will compute the maximum in each of the segments of the form $w_2 = w_2^0 - \frac{1}{2}w_1$, with $w_1 \in [0, \frac{2}{3}(w_2^0 - \frac{1}{2})]$, and $w_2^0 \in [\frac{1}{2}, 1]$ (each value of w_2^0 defines a segment in the triangle, parallel to the $1 - 2w_2 = 2w_1$ side of the triangle). The value of the maximization function is:

$$\frac{w_2^0(1 - w_2^0) + (18w_2^0 - 9)w_1 - 25w_1^2}{w_2^0 - w_1}.$$

For $w_2^0 \leq 25/47$, the maximum is achieved at $w_1 = \frac{2}{3}(w_2^0 - \frac{1}{2})$ and it is $(2 - w_2^0)/3 < 1$. For $w_2^0 \geq 25/47$, the maximum is achieved at $w_1 = \frac{1}{5}(5w_2^0 - 2\sqrt{2}\sqrt{w_2^0 + (w_2^0)^2})$ and it is:

$$-\frac{20\sqrt{2}(w_2^0)^2}{\sqrt{w_2^0(w_2^0 + 1)}} + \left(32 - \frac{20\sqrt{2}}{\sqrt{w_2^0(w_2^0 + 1)}}\right)w_2^0 + 9.$$

The above as a function of $w_2^0 \in [1/2, 1]$ is convex and therefore its maximum is attained at one of the extremes of the range $[1/2, 1]$, namely $w_2^0 = 1$ for a maximum value of 1, as expected. Therefore $w_1 = 1/5$, which implies $w_2 = 3/5$. We choose the maximum possible r to satisfy the constraint $r \leq 1/(w_1 + w_2)$, which is $r = 5/4$. \square

B Proof of Lemma 7

For $N = 4$ and $D = 2$, if there is one round with three boxes in one of the \mathcal{FRO} and \mathcal{OPT} solutions, then there are two boxes in this round that are in the same round in the other solution (for example, boxes 1 and 3 when \mathcal{FRO} : 123|4 and \mathcal{OPT} : 14|23). By combining the above two boxes to one (summing the probabilities and weights), we get an instance of $N = 3, D = 2$, with the same approximation ratio, i.e., we have reduced the problem of finding the worst approximation ratio to the case $N = 3$ that we studied before. Therefore, we only have to consider the following cases:

$$\begin{array}{ll} \mathcal{FRO} : 12|34 \text{ and } \mathcal{OPT} : 13|24 & \mathcal{FRO} : 12|34 \text{ and } \mathcal{OPT} : 24|13 \\ \mathcal{FRO} : 12|34 \text{ and } \mathcal{OPT} : 23|14 & \mathcal{FRO} : 12|34 \text{ and } \mathcal{OPT} : 14|23 \end{array}$$

We will consider inputs in which the probabilities or weights do not sum up to 1 and therefore need the following observation, whose proof is based on the fact that a ratio of two costs is oblivious to a scaling of p_i s or w_i s and we omit the details. We denote the p/w -ratio of box i with r_i .

Observation 1. *Scaling p_i s and w_i s by any positive factor will not affect the approximation ratio.*

Again, because of space considerations, we only analyze the first case, where \mathcal{FRO} is 12|34 and \mathcal{OPT} is 13|24. The other cases are similar and we will give a sketch later.

Lemma 10. *For instances of $N = 4$ in which $\mathcal{OPT} = (13|24)$ and $\mathcal{FRO} = (12|34)$, the approximation ratio is at most $8/7$.*

In order to prove the above, we transform gradually any instance of the above form to instances that have worse and worse approximation ratio. For the final instance, it is possible to solve an optimization problem and show that the worst approximation ratio is $8/7$, like the $N = 3$ case.

Observation 2. *ρ is maximized if $w_1 = 0$.*

Proof. Consider the input $\mathbf{p} = \langle p_1, p_2, p_3, p_4 \rangle$ and $\mathbf{w} = \langle w_1, w_2, w_3, w_4 \rangle$ and assume it is ordered according to p/w ratio and that FRO:12|34 and OPT:13|24. Set $\rho = \text{cost}(\text{FRO})/\text{cost}(\text{OPT})$. We have the following costs:

$$\begin{aligned}\text{cost}(\text{FRO}) &= \text{cost}(12|34) = (p_1 + p_2)(w_1 + w_2) + (p_3 + p_4)(w_1 + w_2 + w_3 + w_4), \\ \text{cost}(1|234) &= p_1 w_1 + (p_2 + p_3 + p_4)(w_1 + w_2 + w_3 + w_4), \\ \text{cost}(123|4) &= (p_1 + p_2 + p_3)(w_1 + w_2 + w_3) + p_4(w_1 + w_2 + w_3 + w_4), \\ \text{cost}(\text{OPT}) &= \text{cost}(13|24) = (p_1 + p_3)(w_1 + w_3) + (p_2 + p_4)(w_1 + w_2 + w_3 + w_4).\end{aligned}$$

Moreover,

$$\text{cost}(12|34) \leq \text{cost}(1|234), \quad \text{cost}(12|34) \leq \text{cost}(123|4).$$

Now change the input by setting $w_1 = 0$. Then, for that input, the order of boxes does not change, and it is still the case that FRO: 12|34. Moreover the new approximation ratio is:

$$\begin{aligned}\rho' &= \text{cost}'(\text{FRO})/\text{cost}'(\text{OPT}) \geq (\text{cost}(\text{FRO}) - w_1)/(\text{cost}(\text{OPT}) - w_1) \\ &= (\rho - w_1/\text{cost}(\text{OPT}))/(1 - w_1/\text{cost}(\text{OPT})) \geq \rho\end{aligned}$$

This implies setting $w_1 = 0$ will not decrease the approximation ratio ρ . □

Similarly, we can prove the following.

Observation 3. ρ is maximized if $p_4 = 0$.

Given $w_1 = 0$ and $p_4 = 0$, we now prove:

Observation 4. For $\frac{p_2}{w_2} = \frac{p_3}{w_3}$, ρ is maximized.

Proof. Let $r_2 = p_2/w_2$, and $r_3 = p_3/w_3$. Initially $r_2 \geq r_3$. Since FRO is 12|34, we have $\text{cost}(12|34) \leq \text{cost}(1|234)$ and $\text{cost}(12|34) \leq \text{cost}(123|4)$, which implies $w_2 \leq p_2(w_3 + w_4)/p_1$ and $p_3 \leq (p_1 + p_2)w_3/w_4$. Observe that for the extreme values $w_2 = p_2(w_3 + w_4)/p_1$ and $p_3 = (p_1 + p_2)w_3/w_4$, we have $r_2 \leq r_3$ and therefore there are values of p_3 and w_2 both greater or equal to the initial respective values, for which the ratios of the two boxes become the same. Since for

$$\rho = \frac{\text{cost}(\text{FRO})}{\text{cost}(\text{OPT})} = \frac{(p_1 + p_2)w_2 + p_3(w_2 + w_3 + w_4)}{(p_1 + p_3)w_3 + p_2(w_2 + w_3 + w_4)}$$

we have $\frac{\partial \rho}{\partial p_3} \geq 0$ and $\frac{\partial \rho}{\partial w_2} \geq 0$, the approximation ratio when $r_2 = r_3$ is greater or equal than the initial approximation ratio. □

Now, with a similar optimization analysis like the one done for $N = 3$, it can be proven that the worst possible ratio is $8/7$ (we omit the tedious details). Otherwise, one can use a symbolic optimization program like Mathematica and get the same result.

When OPT outputs 23|14, we transform the input as follows: We set p_4 to 0, decrease w_2 until the first two boxes have the same ratio ($r_1 = r_2$), and increase p_3 until the first of the following events happens: a) $r_2 = r_3$ or b) $\text{cost}(12|34) = \text{cost}(123|4)$. It can be proven by taking derivatives of the approximation ratio ρ that the above transformations never make ρ smaller. Now what remains is to study the two cases with a) $p_4 = 0$, $r_2 = r_3$ and b) $p_4 = 0$, $\text{cost}(12|34) = \text{cost}(123|4)$. With an optimization analysis like above, it is possible to prove that $8/7$ is an upper bound for ρ .

When OPT outputs 14|23, we transform the input as follows: We set w_1 to 0, decrease p_3 until $r_3 = r_4$, and increase w_2 until the first of the following events happens: a) $r_2 = r_3$ or b) $\text{cost}(12|34) = \text{cost}(1|234)$. Again it can be shown by taking derivatives of ρ that the above transformations never make ρ smaller and what remains is to study the two cases with a) $w_1 = 0$, $r_2 = r_3$ and b) $w_1 = 0$, $\text{cost}(12|34) = \text{cost}(1|234)$. Again for both of the above cases, it can be proven that $8/7$ is an upper bound for ρ .

Finally, when OPT outputs 14|23, we do the following transformations: We decrease w_2 until $r_1 = r_2$ and decrease p_3 until $r_3 = r_4$. Again it can be shown by taking derivatives of ρ that the above transformations never make ρ smaller, and (by an optimization analysis) that the maximum possible ρ for an input of the last form is $8/7$.

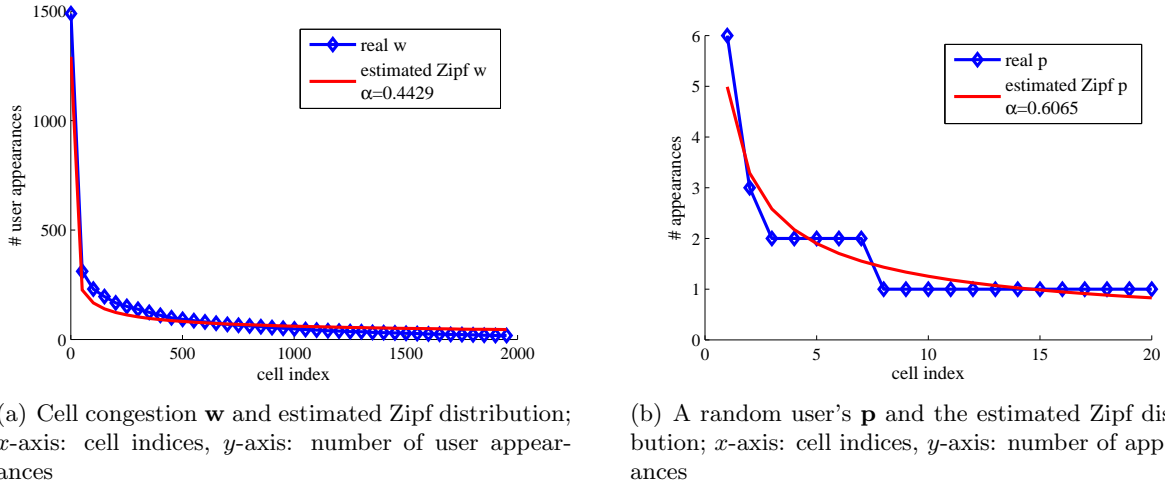


Figure 1: Data Analysis

C Simulation

We conduct simulation on the major application of this problem: paging a mobile user (token) in cells (boxes) in a cellular network system. Parameters of the problem, N , D , \mathbf{p} and \mathbf{w} , retain the same semantics. We mainly explore some aspects of the problem that are not covered by theoretical results. We acquire real user data from Shenzhen, China, provided by China Unicom. First, we analyze and model the user probability and cell congestion. Second, we provide our implemented algorithms and metrics to evaluate the performance of algorithms. Next, we test the performance of algorithms \mathcal{S} and FRO on Zipf data in order to evaluate their actual performance, and on random uniform data in order to test their worst case and average performance without any assumption about user locations. Finally, we evaluate the approximation ratio of \mathcal{S} and FRO on real user data.

Data Analysis: In [8], the authors pointed out that human beings are more likely to appear in N different places with a frequency that follows the *Zipf* distribution, which is defined as $\mathbf{p} = (p_1, \dots, p_N)$ with $p_i = i^{-\alpha} / \sum_{i=1}^N i^{-\alpha}$ where $\alpha \geq 0$ is the Zipf parameter. Zipf distribution obeys a power law. When $\alpha = 0$, it is a uniform distribution; as α grows, the distribution becomes more and more uneven.

We obtain 171929 appearances of 996 users in 5625 cells on 31 consecutive days from the boundary of a metro and a suburban region. For every user appearance, we record pair $\langle \text{user ID}, \text{cell ID} \rangle$ that denotes the user and the cell in which it appears.

Cellwise, we model the cell congestion vector \mathbf{w} of the 5625 cells. We estimate parameter α using the least square method with cross validation and we obtain $\alpha = 0.4429$. The cross validation indicates that the estimates for α differ less than 1% when using the odd entries and even entries of the data separately. This further validates our assumption of a Zipf distribution. The cell congestion \mathbf{w} and the estimated Zipf \mathbf{w} is shown in Figure 1(a).

Userwise, we model the user probability vector \mathbf{p} of all 996 users. Since each user only appears in a limited number of cells, techniques like cross validation cannot be used because we do not have enough samples. We randomly plot many user probability vectors and their corresponding estimated Zipf distribution, and almost all of the plots appear like in Figure 1(b). This indicates that \mathbf{p} follows a Zipf distribution.

Algorithms and Benchmarks: We implement four algorithms: \mathcal{S} and FRO are being evaluated; OPT and OPT- N are metrics to evaluate algorithm performance.

FRO: Algorithm FRO uses the dynamic programming scheme in [3, 15–17, 19]. The most efficient implementation in [3] takes $\Theta(ND)$ time.

S: \mathcal{S} is the greedy algorithm introduced in [9]. It is implemented based on \mathbf{p} without using information on \mathbf{w} . It sorts cells by a non-increasing order of p_i and allocate cells in this order, one at a time, to the partition d with smallest sum of allocated cells P_d . Its complexity is $\Theta(ND)$. It is a 49/48-competitive algorithm for typical users.

OPT: Since the problem is NP-Hard, obtaining optimal solution takes exponential time. Our most efficient implementation of OPT computes all permutations of D -partition on the N cells and takes $\Theta(D^N)$ time. On current computers, this implementation allows us to run instances up to $N = 20$ and $D = 2$, $N = 15$ and $D = 3$, or in general $D^N \leq 2^{20}$.

OPT- N : For larger values of D and N , computing an OPT solution takes too long. OPT- N is the polynomial time algorithm that optimally pages cells in $D = N$ rounds. According to Corollary 2, this can be implemented in $\Theta(N \log N)$ time by sorting cells by p_i/w_i . OPT- N gives a lower bound of the cost of OPT.

Randomly generated data: We test the performance of our algorithms on randomly generated data. We run two algorithms, \mathcal{S} and FRO, on Zipf distributed data and uniform random data. Our purpose of using Zipf data is to study how the algorithms work in real world applications; our purpose of using uniform random data is to study how the algorithms behave if we have no knowledge about users' whereabouts. We evaluate small data instances by comparing the results with OPT, and large instances with OPT- N .

First test: We test the performance of \mathcal{S} and FRO on a typical user, that is, $\mathbf{p} = \mathbf{w} \sim \text{Zipf}(\alpha)$. The results are shown in figure 2. We can see that on typical users, both algorithms perform very well, and \mathcal{S} slightly outperforms FRO because it is smarter in breaking ties (of p_i/w_i). We run the test on different N s and D s and the results are very similar. We also run the FRO algorithm with different initial order of cells and the results are similar.

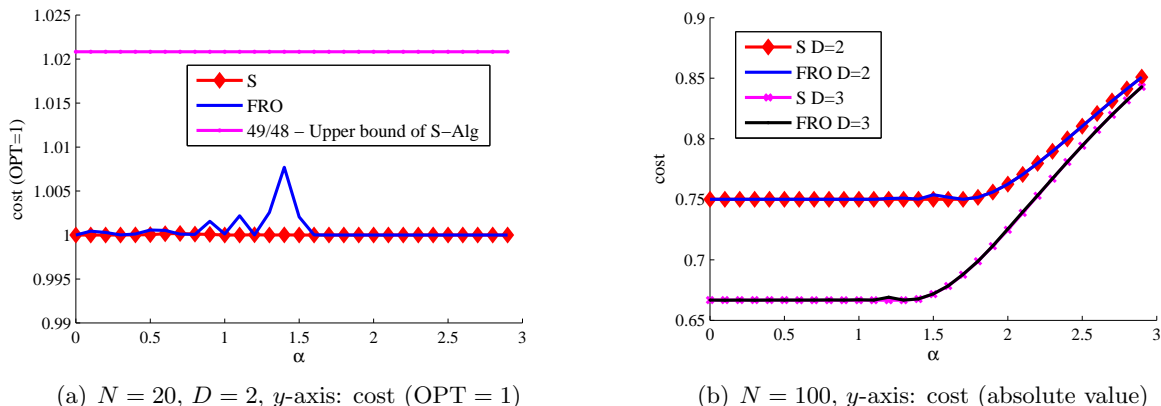
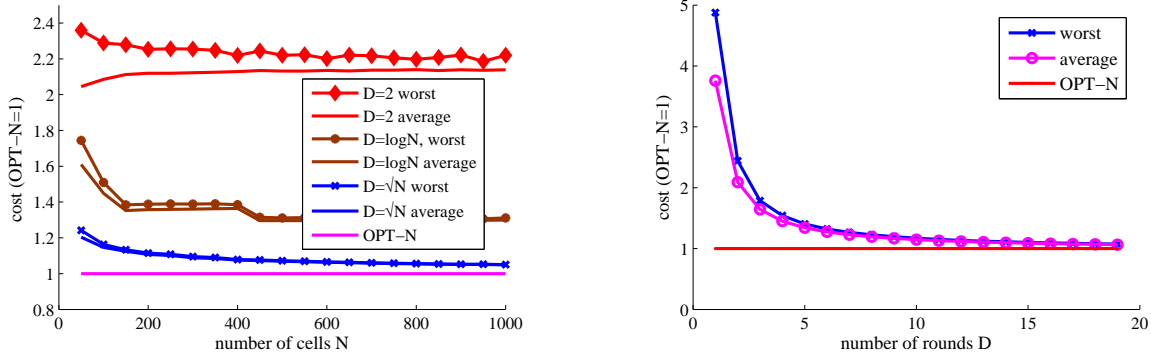


Figure 2: Performance of \mathcal{S} and FRO on a typical user, $\mathbf{p} = \mathbf{w} \sim \text{Zipf}(\alpha)$, x -axis: α

Second Test: We test the performance of algorithm FRO on uniform random \mathbf{p} and \mathbf{w} for larger number of rounds D to complement our theoretical results. Figure 3(a) shows the first result. We can see that FRO performs stably for larger D because the worst-case cost is very close to average cost. We also see that when $D \geq \sqrt{N}$, FRO is almost optimal. This is because when D is larger, there is less flexibility for the optimal solution to not obey the cellwise order. To evaluate the performance of FRO on not very large D , we run the test shown in Figure 3(b). It indicates that FRO performs reasonably well when compared with OPT- N (recall that when $D = 2$, FRO is a 8/7-approximation). FRO also performs very stably because the worst-case and average costs are very close and the curve is very smooth. We also run this test on random Zipf data. We create random Zipf data by shuffling \mathbf{w} on a typical Zipf user. To do such a shuffle we take a random permutation of the elements of the vector \mathbf{w} . The results are very similar.

Third Test: We test FRO on atypical users. We generate atypical users by randomly generating a sorted \mathbf{p} , reversing its order and randomly generating another sorted \mathbf{w} . Both \mathbf{w} and \mathbf{p} are normalized. We test



(a) Average and worst case cost ($\text{OPT}-N = 1$) of algorithm FRO, $N = 50 \dots 1000$, $D = 2, \log N, \sqrt{N}$, uniform random \mathbf{p} and \mathbf{w} , 10,000 runs.

(b) Average and worst case cost ($\text{OPT}-N = 1$) of algorithm FRO, $N = 100$, $D = 1 \dots 20$, uniform random \mathbf{p} and \mathbf{w} , 10,000 runs.

Figure 3: Tests on uniform data (similar results for Zipf \mathbf{p} and \mathbf{w} after random shuffle)

1,000,000 instances for which ($N = 20, D = 2$), and ($N = 15, D = 3$) and the solutions are always optimal. This coincides with our theoretical results.

Fourth Test: Finally, we test FRO for a general user. A general user is a user that follows a Zipf location distribution. However, it may or may not follow the massive behavior. Therefore, we define a general user to be one with $\mathbf{w} \sim \text{Zipf}(\alpha = 0.4429)$ and $\mathbf{p} \sim \text{Zipf}(\alpha' = 0.5)$, while \mathbf{p} and \mathbf{w} are independently shuffled. Lemma 8 shows that the FRO solution can be $8/7$ times worse than the optimal. We measure the average and worst approximation ratio in 1,000,000 runs and the results are shown in Table 2(a). We conclude that for the above random user the FRO algorithm has almost optimal performance on average and worst case performance much better than $8/7$.

(a) Approximation ratio FRO, general user, $D = 2$

metric	$N = 4$	$N = 5$
worst	1.00868	1.00194
average	1.00199	1.00009
metric	$N = 8$	$N = 10$
worst	1.00489	1.00557
average	1.00021	1.00016

(b) Approximation ratio, real user

Setting	Algo.	Average	Worst
$N = 20, D = 2$	\mathcal{S}	1.13475	9.60367
	FRO	1.00000	1.00001
$N = 15, D = 3$	\mathcal{S}	1.07308	5.94067
	FRO	1.00000	1.00036
$N = 10, D = 4$	\mathcal{S}	1.01648	2.52683
	FRO	1.00007	1.00080

Table 2: Approximation Ratios

We also try to tune up algorithm FRO not only taking consideration the p_i/w_i ratio but also p_i and w_i values. Our preliminary results indicate the tuned algorithm does not provide noticeable improvement.

Real User Data: We test the performance of FRO and \mathcal{S} on real user data. In order to obtain the optimal solution in reasonable time, for each user, we only take their top 20 or 15 most frequent cells. This is justifiable because none of the 996 users appears more than once beyond their top 20 favorite cells. Vector \mathbf{p} is acquired directly from the user and \mathbf{w} is a universal vector. Table 2(b) shows the results. We can see that algorithm FRO performs almost optimally with respect to both average and worst-case cost. Algorithm \mathcal{S} performs fine with respect to average cost but poorly for worst-case cost. This indicates there is a certain portion of not-very-typical users, of which some are very atypical.