

Competitive Clustering of Stochastic Communication Patterns on a Ring^{*}

Chen Avin¹, Louis Cohen², and Stefan Schmid³

¹ Ben-Gurion University of the Negev, Israel

² Ecole Normale Supérieure Paris Saclay, France

³ Aalborg University, Denmark & TU Berlin, Berlin, Germany

Abstract. This paper studies a fundamental dynamic clustering problem. The input is an online sequence of pairwise communication requests between n nodes (e.g., tasks or virtual machines). Our goal is to minimize the communication cost by partitioning the communicating nodes into ℓ clusters (e.g., physical servers) of size k (e.g., number of virtual machine slots). We assume that if the communicating nodes are located in the same cluster, the communication request costs 0; if the nodes are located in different clusters, the request is served remotely using inter-cluster communication, at cost 1. Additionally, we can migrate: a node from one cluster to another at cost $\alpha \geq 1$.

We initiate the study of a stochastic problem variant where the communication pattern follows a fixed distribution, set by an adversary. Thus, the online algorithm needs to find a good tradeoff between benefitting from quickly moving to a seemingly good configuration (of low inter-cluster communication costs), and the risk of prematurely ending up in a configuration which later turns out to be bad, entailing high migration costs.

Our main technical contribution is a deterministic online algorithm which is $O(\log n)$ -competitive with high probability (w.h.p.), for a specific but fundamental class of problems: namely on ring graphs.

1 Introduction

Modern distributed systems are often highly virtualized and feature unprecedented resource allocation flexibilities. For example, these flexibilities can be exploited to improve resource utilization, making it possible to multiplex more applications over the same shared physical infrastructure, reducing operational costs and increasing profits. However, exploiting these resource allocation flexibilities is non-trivial, especially since workloads and resource requirements are time-varying.

This paper studies a fundamental dynamic resource allocation problem underlying many network-intensive distributed applications, e.g., batch processing or streaming applications, or scale-out databases. To minimize the resource footprint (in terms of bandwidth) of such applications as well as latency, we want to colocate frequently communicating tasks or virtual machines on the same physical server, saving communication across the network. The

^{*} Research supported by the German-Israeli Foundation for Scientific Research and Development, G.I.F. No I-1245-407.6/2014). Part of the research was done while the second author was visiting Ben Gurion University and TU Berlin.

underlying problem can be seen as a clustering problem [3]: nodes (the tasks or virtual machines) need to be partitioned into different clusters (the physical servers), minimizing inter-cluster communications.

The clustering problem is challenging as the detailed communication patterns are often stochastic and the specific distribution unknown ahead of time. In other words, a clustering algorithm must *deal with uncertainties*: although two nodes may have communicated frequently in the past, it can turn out later that it is better to collocate different node pairs. Accordingly, clustering decisions may have to be reconsidered, which entails migrations.

Our Contributions. This paper initiates the study of a natural dynamic clustering problem where communication patterns follow an unknown distribution, chosen by an adversary: the distribution represents the worst-case for the given online algorithm, and communication requests are drawn i.i.d. from this distribution. Our goal is to devise online algorithms which perform well against an optimal offline algorithm which has perfect knowledge of the distribution. Our main technical contribution is a deterministic online algorithm which, for a special but fundamental request pattern family, namely the ring, achieves a competitive ratio of $O(\log n)$, with high probability (w.h.p.), i.e., with probability at least $1 - 1/n^c$, where n is the total number of nodes and c is a constant.

Novelty and Challenges. Our work presents an interesting new perspective on several classic problems. For example, our problem is related to the fundamental statistical problem of guessing the most likely distribution (and its parameters) from which a small set of samples is drawn. Indeed, one natural strategy of the online algorithm could be to first simply sample requests, and once a good estimation of the actual distribution emerges, directly move to the optimal clustering configuration. However, as we will show in this paper, the competitive ratio of this strategy can be very bad: the communication cost paid by the online algorithm during sampling can be high. Accordingly, the online algorithm is forced to eliminate distributions early on, i.e., it needs to migrate to seemingly low-cost configurations. And here lies another difference to classic distribution learning problems: in our model, an online algorithm needs to pay for changing configurations, i.e., when revising the “guessed distribution”. In other words, our problem features an interesting *combination of distribution learning and efficient searching*. It turns out that amortizing the migration costs with the expected benefits (i.e., the reduced communication costs) at the new configuration however is not easy. For example, if the request distribution is uniform, i.e., if all clustering configurations have the same probability, the best strategy is not to move: the migration costs cannot be amortized. However, if the distribution is “almost uniform”, migrations are required and “pay off”. Clearly, distinguishing between uniform and almost uniform distributions is difficult from an online perspective.

Organization. The remainder of this paper is organized as follows. In Section 2, we introduce our formal model. In Section 3, we provide intuition about our problem and highlight the challenges. In Section 4, we present our deterministic online algorithm, and we analyze it formally in Section 5. After reviewing related work in Section 6, we conclude our contribution in Section 7.

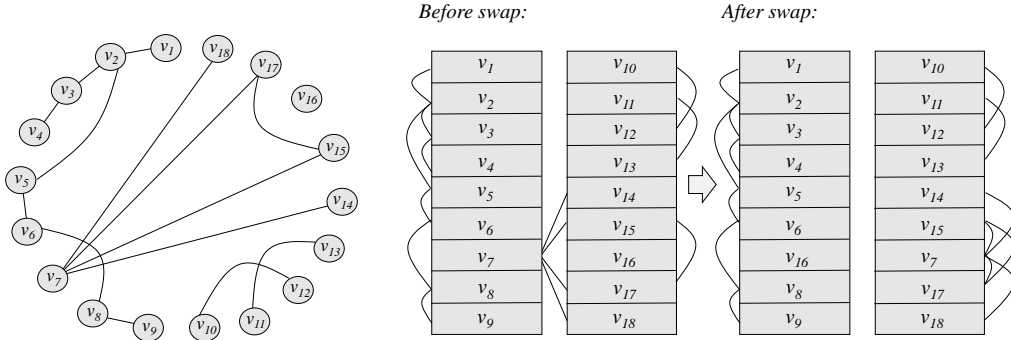


Fig. 1.1. Example: Communication patterns drawn from a certain distribution (on the left, represented as a communication graph) need to be learned and clustered. In this example, we have $\ell = 2$ clusters of size $k = 9$. In the middle, a bad clustering is shown: there are four inter-cluster edges (“before swap”). However, by swapping nodes v_7 and v_{16} , all inter-cluster edges can be removed (on the right in the figure). Note that different edges can have different frequencies, which however are not depicted in this example.

2 Model

We consider the problem of partitioning n nodes $V = \{v_1, v_2, \dots, v_n\}$ into ℓ clusters of capacity k each. We assume that $n = \ell \cdot k$, i.e., nodes perfectly fit into the available clusters, and there is no slack. We call a specific node-cluster assignment a *configuration* c . We assume that the communication request is generated from a fixed distribution \mathcal{D} , chosen in a worst-case manner by the adversary. The sequence of actual requests $\sigma(\mathcal{D}) = (\sigma_1, \sigma_2, \dots, \sigma_T)$, is sampled i.i.d. from this distribution: the communication event at time t is a (directed) node pair $\sigma_t = (v_i, v_j)$. Alternatively, we represent the distribution \mathcal{D} as a weighted graph $G = (V, E)$. For an edge $(v_i, v_j) \in E(G)$, let the weight of the edge $p(v_i, v_j)$ denote the probability of a communication request from between v_i and v_j : each edge $e \in E$ has a certain probability $p(e)$ and $\sum_{e \in E} p(e) = 1$. A request (i.e., edge in G) $\sigma_t = (v_i, v_j)$ is called *internal* if v_i and v_j belong to the same cluster at the current configuration (i.e., at the time of the request); otherwise, the request (edge) is called *external*. We will assume that the communication cost of an external request is 1 and the cost of an internal request is 0.

Note that each configuration uniquely defines external edges that form a “cut”, interconnecting ℓ clusters in G . Therefore in the following, we will treat the terms “configuration” and “cut” as synonyms and use them interchangeably; we will refer to them by c . Moreover, we define the probability of a cut (or identically a configuration) c as the sum of the probabilities of its *external edges*: $p(c) = \sum_{e \in c} p(e)$. We also note that there are many configurations which are symmetric, i.e., they are equivalent up to cluster renaming. Accordingly, in the following, we will only focus on the actually different (i.e., non-isomorphic) configurations.

To reduce external communication costs, an algorithm can change the current configuration by using *node swaps*. Swapping a node pair costs 2α (two node migrations of cost α each). Since the request probability of different configurations/cuts differs, the goal of the al-

gorithm will be to quickly guess and move toward a good cut, a configuration that reduces its future cost. Figure 1.1 shows an example.

In particular, we are interested in the *online problem variant*: we assume that the distribution \mathcal{D} of the communication pattern (and hence the σ we observe is generated from) is initially *unknown* to the online algorithm. Nevertheless, we want the performance of an online clustering algorithm, *ON*, to be similar to the one of a hypothetical offline algorithm, *OFF*, which knows the request distribution as well as the number of requests σ , henceforth denoted by $|\sigma|$, ahead of time. In particular, *OFF* can move before any request occurs or σ is generated.

We aim to minimize the competitive ratio, the worst ratio of the online algorithm cost divided by the offline algorithm cost (for a given distribution \mathcal{D} and the same starting configuration c_0):

$$\rho = \max_{\sigma(\mathcal{D})} \frac{ON(\sigma(\mathcal{D}))}{OFF(\sigma(\mathcal{D}))}$$

Here, the cost $ON(\sigma(\mathcal{D}))$ of any algorithm *ON* for a sequence $\sigma(\mathcal{D})$ is the sum of the overall communication costs and the migration costs. We consider bounds on ρ with high probability.

As a first step, we focus on partitioning problems where $\ell = 2$ and consider fundamental ring communication patterns. That is, the communication graph G is the cycle graph and the event space is defined over the edges $E = \{(v_1, v_2), (v_2, v_3), \dots, (v_{n-1}, v_n), (v_n, v_1)\}$. Moreover, we assume configurations that minimize the cut, that is nodes are partitioned according to contiguous subsequences of the identifier space. Each cluster is (up to modulo) of the form, $\{(v_i, v_{i+1}, \dots, v_{i+k-1})\}$. This communication pattern is not only fundamental but also captures the aspects and inherent tradeoffs rendering the problem non-trivial. In this model, an algorithm changes configurations using *rotations* (either clockwise or counter-clockwise). See Figure 3.1.

3 The Challenge of Dynamic Clustering

In order to acquaint ourselves with the problem and understand the fundamental challenges involved in dynamic clustering, we first provide some examples and discuss naive strategies. Let us consider an example with $n = 2k$ nodes divided into $\ell = 2$ clusters of size k . There are k possible configurations/cuts: $\{c_0, c_1, \dots, c_{k-1}\}$. At one end of the algorithmic spectrum lies a lazy algorithm which never moves, let's call it *LAZY*. At the other end of the spectrum lies a very proactive algorithm which greedily moves to the configuration which so far received the least external requests, let's call it *GREEDY*. Both *LAZY* and *GREEDY* are doomed to fail, i.e., they have a large competitive ratio: *LAZY* fails under a request distribution where the initial external cut has probability 1, i.e., $p(c_0) = 1$ and for any $i > 0$, $p(c_i) = 0$: *LAZY* pays for all requests, while after a simple node swap all communication costs would be 0. *GREEDY* fails in uniform distributions, i.e., if $p(c_i) = 1/k$ for all i : the best configuration is continuously changing, and in particular, the best cut is likely to be at distance $\Omega(k)$ from the initial configuration c_0 : *GREEDY* quickly incurs migration costs in the order of $\Omega(\alpha \cdot k)$, while staying at the same location would cost $1/k$ per request. Thus, the competitive ratios grow super-linearly in the number of requests and in the number of nodes.

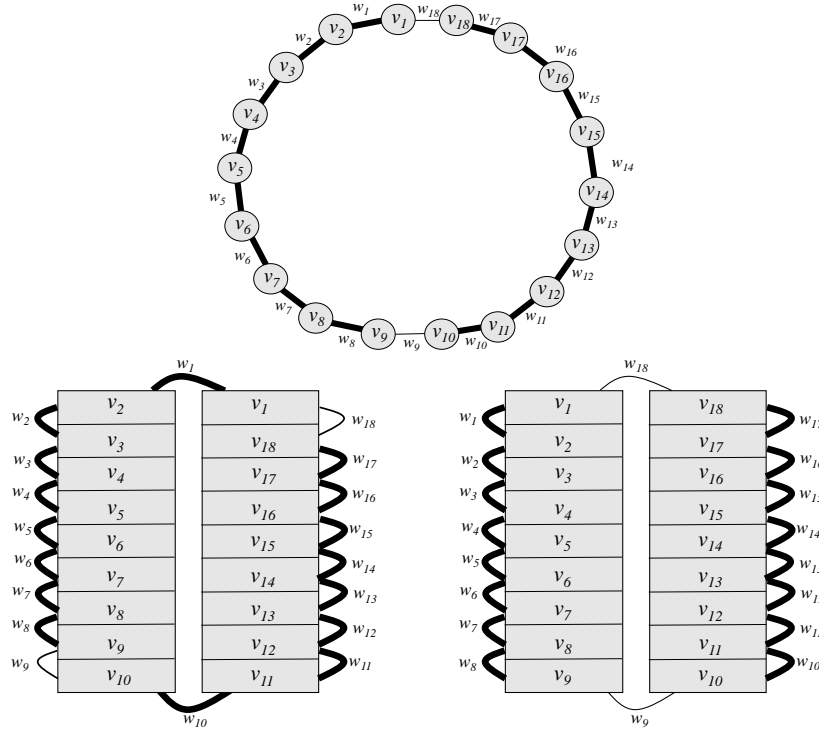


Fig. 3.1. Weighted ring communication pattern: frequently used edges (in *bold*) should not be part of the cut. The cut can be changed using rotations: in the figure, a counter-clockwise rotation leads from the middle to the right configuration.

Another intuitive strategy could be to wait in the initial configuration c_0 for some time, simply observing and sampling the actual distribution, until a “sufficiently accurate” estimation of the distribution is obtained. Then, we move directly to the (hopefully) optimal configuration. Thus, the problem boils down to the classic statistical problem of estimating the distribution (and its parameters) from samples. However, it is easy to see that waiting for the optimal distribution to emerge is costly. Imagine for example a scenario where the initial configuration/cut c_0 has a high probability, and there are two additional cuts c_1 and c_2 which have almost the same low probability (for example polynomially low probability). Clearly, waiting at c_0 to learn whether c_1 or c_2 is better is not only very costly, but it may also be pointless: even if the online algorithm ended up at c_1 although c_2 was a little bit better, the resulting competitive ratio could be still small.

Thus, the key challenge of our problem lies in its required joint optimization of *learning and searching*: while learning the distribution, an efficient search algorithm must be employed to minimize reconfiguration costs. In particular, the following criteria need to be met:

1. *Migrate early...*: An online algorithm should migrate away from a suboptimal configuration early, possibly long before the optimal configuration can be guessed.

2. ... *but not too early*...: An online algorithm should avoid frequent migrations, e.g., due to a wrong or poor estimate of the actual request distribution.
3. ... *and locally*: Especially if the length of σ is small (small number of requests), it may not make sense to migrate to an optimal but faraway location, even if the distribution is known: even *OFF* would not move there.

4 Deterministic and Competitive Clustering

With these intuitions and challenges in mind, we present our solution. Let us first start with the offline algorithm. It is easy to see that *OFF*, knowing the distribution as well as the number of requests, only moves *once* in time (i.e., one move consisting of multiple migrations or node swaps): namely *in the beginning* and *to the configuration providing an optimal cost-benefit tradeoff*. Concretely, *OFF* computes for each configuration c_i , its expected cost-benefit tradeoff: the communication cost of configuration c_i is $|\sigma| \cdot p(c_i)$ and the cost of moving there is $2\alpha \cdot d(c_0, c_i)$, where $d(\cdot, \cdot)$ is the rotation distance between the two configurations (the smallest number of rotation moves to reach the other configuration). Thus, *OFF* will move to $c_{OFF} := \operatorname{argmin}_{c_i} p(c_i) + (2\alpha \cdot d(c_0, c_i))/|\sigma|$ (note that this configuration is not necessarily unique). In the following, we will use the short form $d_i = d(c_0, c_i)$ to denote distances relative to c_0 , the initial configuration.

The online algorithm is more interesting. The competitive and deterministic online algorithm presented in this paper relies on three key ideas:

- *Eliminating bad configurations*: We define conditions for configurations which, if met, allow us to eliminate the corresponding configurations once and for all. In particular, we will guarantee (w.h.p.) that an online algorithm be competitive (even) if it never moves back to such a configuration anymore in the future. In other words, our online algorithm will only move between configurations for which this condition is not true yet.
- *Local migrations and growing-radius search strategy*: In order to avoid high migration costs, our online algorithm is local in the sense that it only moves to nearby cuts/configurations once the condition of the current configuration is met and it needs to be eliminated. Concretely, our online algorithm is based on a growing-radius search strategy: we only migrate to valid configurations lying within the given radius. Only if no such configurations exist, the search radius is increased.
- *Amortization*: The radius growth strategy alone is not sufficient to provide the necessary amortization for being competitive. Two additions are required:
 1. *Directed search*: An online algorithm may still incur a high migration cost when frequently moving back-and-forth within a given radius, chasing the next best configuration. Therefore, our proposed online algorithm first moves in one direction only (clockwise), and then in the other direction, bounding the number of times the c_0 configuration is crossed.
 2. *Lazy expansion*: Even once all configurations within this radius have been eliminated, the online algorithm should not immediately move to configurations in the next larger interval. Rather, the algorithm waits until a certain amount of requests have been accumulated, allowing to amortize the migrations (an “insurance”).

With these high-level ideas in mind, we now describe the algorithm in detail (cf. Algorithm 1). We consider a time t , and assume that the online algorithm is at configuration c_t . The algorithm maintains an array $r[]$ where it counts, for each possible configuration c_0, \dots, c_{k-1} , the number of samples that hit an external edge of the corresponding cut; in other words, $r[]$ is used to estimate the distribution of the communication pattern. Let \mathcal{E} be the set of the eliminated configurations, and let $\bar{\mathcal{E}}$ be the complement of \mathcal{E} : the set of configurations not eliminated yet. R is the search radius, initially $R = 1$. Upon each request, σ_t , we first increment the value of the corresponding configuration in the sampling array $r[]$ (only one configuration is affected by a given external request). We then compare all configurations not eliminated yet to the “seemingly best configuration”: the configuration which received the least (external) requests so far (i.e., $\operatorname{argmin}_{c_i} r[c_i]$). Let $r_{\min} := \min_{c_i} r[c_i]$ be the minimum value. We now eliminate any configuration c_j for which the condition $\operatorname{Cond}(r[c_j], r_{\min})$ is fulfilled: c_j is too far from the optimum. Concretely, w.l.o.g. assume that $r[c_j] > r[c_i]$ and let $\gamma = r[c_i]/r[c_j] < 1$. Then for $\epsilon > 0$ (a parameter for the error probability), we use the following condition:

$$\operatorname{Cond}(j, i) := \begin{cases} \text{True} & r[c_j] \geq \frac{\ln(\frac{1}{\epsilon})}{\ln(\frac{2}{1+\gamma}) - (\frac{1-\gamma}{2})} \\ \text{False} & \text{otherwise} \end{cases} \quad (1)$$

If on this occasion, we eliminated our own current configuration $c(t)$, we then have to decide where we want to move next, using the function $\operatorname{next}(\bar{\mathcal{E}})$ (unless all configurations have been eliminated). The distance from the suggested next configuration c_{next} to c_0 (the initial configuration) may be greater than the current radius R , in which case we double R until $R \geq d_{\operatorname{next}}$. However, before moving, we also test whether $\min_{\{d_{\operatorname{next}} < R\}} (r[c_{\operatorname{next}}]) \geq \alpha \cdot R$. Only if this is fulfilled, we can move to the new configuration c_{next} ; otherwise, we lazily stay on the current configuration.

Let us now elaborate more on the moving strategy. Before going into the details however, let us note that for ease of presentation, we will use two different but equivalent numbering schemes to refer to configurations: depending on what is more useful in the current context. In particular, while when talking about the number of requests, $r[]$, we often enumerate configurations globally, $0, 1, 2, \dots, k$. When discussing moving strategies, we often enumerate configurations relative to c_0 , i.e., $-1, 1, -2, 2, \dots, c_{k/2}$, depending on whether they are located clock- or counter-clock wise from c_0 .

Given this remark, let us consider a simple migration strategy: we could always move to the closest not eliminated configuration next. However, we can show that this strategy is flawed. To see this, consider the following distribution:

$$\forall i \in [1; \frac{k}{2}] : p(c_i) = \frac{1}{k^i}, \quad p(c_0) = \left(1 - \sum_{i \in [1; \frac{k}{2}]} p(c_i)\right), \quad \forall i \in [-\frac{k}{2}, -1] : p(c_i) = 0$$

In such a situation, we have to move away from the configuration c_0 as soon as possible: we pay a cost close to 1 on this configuration, for each request. In particular, we cannot wait until we even observe the first request on c_1 : we would incur high communication costs. Now, however, the algorithm may move in the wrong direction: e.g., to c_1 , and then to the closer configuration not eliminated, c_2 . Thus, eventually all configurations in $[c_0, c_{k/2}[$ may be visited before reaching the minimal configurations.

Algorithm 1 Online Algorithm *ON* (upon receiving request $\sigma(t)$ and current configuration $c(t)$)

Initialize: $r := [0; \dots; 0]$, $\mathcal{E} := \{\}$, $\overline{\mathcal{E}} := [-\frac{k}{2}, \frac{k}{2}]$, $R := 1$ $\epsilon := \frac{1}{n^2}$

- 1: $c_j = c(\sigma(t))$ (* configuration to which $\sigma(t)$ is external *)
- 2: $r[c_j] ++$
- 3: $r_{\min} := \min\{r[i] \mid i \in [1, k]\}$
- 4: **if** $c_j \in \overline{\mathcal{E}}$ **then**
- 5: **if** $\text{Cond}(r[c_j], r_{\min})$ **then**
- 6: remove c_j from $\overline{\mathcal{E}}$
- 7: add c_j to \mathcal{E}
- 8: **end if**
- 9: **end if**
- 10: **if** $c(t) \in \mathcal{E}$ **then**
- 11: $c_{\text{next}} :=$ The next configuration $c_i \in \overline{\mathcal{E}}$ on the *searching path*
- 12: **while** $d_{\text{next}} > R$ **do**
- 13: $R = 2R$
- 14: **end while**
- 15: **if** $r[c(t)] \geq \alpha \cdot d_{\text{next}}$ **then**
- 16: move from $c(t)$ to c_{next}
- 17: $c(t) := c_{\text{next}}$
- 18: **end if**
- 19: **end if**

This is reminiscent of classic line searching [12] type problems like “the goat searches the hole in the fence”-escape problems: moving in one direction only, the goat may risk missing a nearby hole in the other direction. That is, moving greedily in one direction is $\Omega(F)$ competitive only, where F is the circumference of the fence, which in our case means that the competitive ratio is $\Omega(k)$. Accordingly, some combination of search-left and search-right is required. Our search radius R is centered around c_0 at any time during the execution of the algorithm, and we always first explore all remaining non-eliminated configurations in one direction, and then explore the remaining configurations in the other direction. In other words, starting from c_0 , we alternate the search between the positive and negative configurations following the sequence: $(1, -1, 2, 3, -2, -3, \dots, 2^{2i-2} + 1, \dots, 2^{2i}, -2^{2i-1} - 1, \dots, -2^{2i+1}, \dots)$. Thus, configuration c_0 is crossed only a constant number of times per given radius R . We call this sequence the *searching path*.

Given a moving strategy, we next note that we should not move too fast: we introduce a second condition for when it is safe to move. When in a configuration 2^{2i} and before we want to explore configurations in $[-2^{2i+1}, -2^{2i-1}]$, we wait in the configuration c_{\min} between configurations -2^{2i-1} and 2^{2i} , until this configuration fulfills $r[c_{\min}] \geq \alpha \cdot 2^{2i+1}$. Similarly, when moving from the configuration -2^{2i+1} to explore the configurations in $[2^{2i}, 2^{2i+2}]$, we will wait at c_{\min} between $[-2^{2i+1}, 2^{2i}]$, until $r[c_{\min}] \geq \alpha \cdot 2^{2i+2}$.

5 Analysis

We first make some general observations on our elimination condition. Subsequently, we will present a cost-breakdown which will be helpful to analyze the competitive ratio of ON : we will show that each cost component is competitive with respect to the optimal offline algorithm.

The following lemma provides an intuition of our algorithm and its condition.

Lemma 1. *Let $\epsilon > 0$, then if $\text{Cond}(j, i) = \text{True}$,*

$$\Pr(p(c_j) > p(c_i)) \geq 1 - \epsilon$$

Proof. We first prove the following helper claim.

Claim. Assume c_i and c_j occur with the same (unknown) probability, let $b > a$ and assume w.l.o.g. that $r[c_j] > r[c_i]$, then,

$$\Pr(r[c_i] \leq a \text{ and } r[c_j] \geq b \mid p(c_i) = p(c_j)) \leq \Pr(r[c_j] \geq (1 + \delta)X) \leq \left(\frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^X$$

where $\delta = \frac{b-a}{b+a}$ and $X = \frac{b}{1+\delta} = \frac{b+a}{2}$

Proof. The proof idea is to consider two probabilities using known Chernoff Bounds [16]:

$$P_1[\delta_i] := \Pr(r[c_i] \leq (1 - \delta_i)E[r[c_i]]) \leq \left(\frac{e^{-\delta_i}}{(1 - \delta_i)^{1-\delta_i}} \right)^{E[r[c_i]]} \quad (2)$$

and

$$P_2[\delta_j] := \Pr(r[c_j] \geq (1 + \delta_j)E[r[c_j]]) \leq \left(\frac{e^{\delta_j}}{(1 + \delta_j)^{1+\delta_j}} \right)^{E[r[c_j]]} \quad (3)$$

The two events are not independent, but we can bound the probability that both events occur by the maximum of the two probabilities when we assume $p(c_j) = p(c_i)$ and $E[r[c_j]] = E[r[c_i]] = (a + b)/2 = X$. In this case, we have that $\delta = (b - a)/(b + a)$. We now want to bound the maximum of these two probabilities. Towards this objective, we study which one of our bounds is greater and bound the maximum of the probability by the maximum of the bounds. Let $B_1[\delta]$ (resp $B_2[\delta]$) the bound on $P_1[\delta]$ (resp. $P_2[\delta]$).

$$B_1[\delta] = \left(\frac{e^{-\delta}}{(1 - \delta)^{1-\delta}} \right)^X \quad B_2[\delta] = \left(\frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^X$$

To determine which one is greater than the other, we now study the function:

$$F(\delta) = \frac{B_1[\delta]}{B_2[\delta]} = \left(e^{-2\delta} \frac{(1 + \delta)^{1+\delta}}{(1 - \delta)^{1-\delta}} \right)^X$$

We obtain that $\forall \delta \geq 0 \quad F(\delta) \leq 1$, so $\Pr(r[c_i] = a \text{ and } r[c_j] = b \mid p(c_i) = p(c_j)) \leq B_2[\delta] = \left(\frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^X$.

We can now prove Lemma 1. Specifically, we want to prove that $\Pr(\text{Cond}(j, i) \mid p(c_j) \leq p(c_i)) \leq \epsilon$. First note that for $x \leq y$:

$$\Pr(\text{Cond}(j, i) \mid p(c_j) = x, p(c_i) = y) \leq \Pr(\text{Cond}(j, i) \mid p(c_j) = y, p(c_i) = y)$$

Next we bound $\Pr(\text{Cond}(j, i) \mid p(c_j) = p(c_i))$ using Claim 5.

$$\Pr(\text{Cond}(j, i) \mid p(c_j) = p(c_i)) \leq \Pr(r[c_j] \geq (1 + \delta)X) = P_2$$

$$P_2 \leq \left(\frac{e^\delta}{(1 + \delta)^{1 + \delta}} \right)^{\frac{r[c_j]}{1 + \delta}}$$

We want that $P_2 \leq \epsilon$:

$$\begin{aligned} \left(\frac{e^\delta}{(1 + \delta)^{1 + \delta}} \right)^{\frac{r[c_j]}{1 + \delta}} \leq \epsilon &\iff \left(\frac{e^{\frac{\delta}{1 + \delta}}}{(1 + \delta)} \right)^{r[c_j]} \leq \epsilon \iff \left(\frac{\delta}{1 + \delta} - \ln(1 + \delta) \right) r[c_j] \leq \ln(\epsilon) \iff \\ r[c_j] &\geq \frac{\ln(\epsilon)}{\left(\frac{\delta}{1 + \delta} - \ln(1 + \delta) \right)} \iff r[c_j] \geq \frac{\ln(\frac{1}{\epsilon})}{\ln(1 + \delta) - \left(\frac{\delta}{1 + \delta} \right)} \end{aligned}$$

Now let $\gamma = \frac{r[c_i]}{r[c_j]} < 1$, so $\delta = \frac{1 - \gamma}{1 + \gamma}$, and we have:

$$r[c_j] \geq \frac{\ln(\frac{1}{\epsilon})}{\ln\left(1 + \frac{1 - \gamma}{1 + \gamma}\right) - \left(\frac{1 - \gamma}{1 + \frac{1 - \gamma}{1 + \gamma}}\right)} \iff r[c_j] \geq \frac{\ln(\frac{1}{\epsilon})}{\ln\left(\frac{2}{1 + \gamma}\right) - \left(\frac{1 - \gamma}{\frac{2}{1 + \gamma}}\right)} \iff r[c_j] \geq \frac{\ln(\frac{1}{\epsilon})}{\ln\left(\frac{2}{1 + \gamma}\right) - \left(\frac{1 - \gamma}{2}\right)}$$

which concludes the proof of the lemma. \square

5.1 A Cost Breakdown

It is convenient to break down the algorithm costs into different components. In case of *OFF*, the situation is fairly easy: *OFF* simply incurs a migration cost, henceforth denoted by OFF_{mig} , of $\text{OFF}_{mig} = 2\alpha \cdot d_{OFF}$ to move to the optimal location c_{OFF} , where d_{OFF} is the rotation distance between c_0 and c_{OFF} , plus an expected communication cost OFF_{comm} of $|\sigma| \cdot p(c_{OFF})$.

In case of *ON*, the situation is more complicated. In particular, while we do not distinguish between different migration costs for *ON* either, we consider three types of communication costs for *ON*: ON_{elim} is the elimination cost, i.e., the total communication cost incurred while *ON* is waiting on every configuration that has not been eliminated yet, until the condition $\text{Cond}(j, i)$ is fulfilled for the current configuration. ON_{ins} is the ‘‘insurance’’ cost paid by *ON* when waiting in an already eliminated configuration, until being allowed to actually move beyond the current radius to a non-eliminated configuration. Finally, ON_{final} is the communication cost paid by *ON* once it reached its final configuration and all other configurations have been eliminated. (Note that the cost incurred at the final configuration while there are still other, non-eliminated configurations, is counted toward elimination costs.)

The total communication cost ON_{comm} is the sum of these three costs. In the following, we will prove that all these cost components are competitive compared to *OFF*'s overall costs, from which the bound on the competitive ratio is obtained.

5.2 Competitive Ratio

We now prove that our online algorithm ON performs well with high probability (w.h.p.). That is, we derive a competitive ratio of $O(\log k)$ which holds with probability at least $1 - 1/n^c$ for some constant c .

Theorem 1. *The competitive ratio achieved by ON is $\rho \in O(\log n)$ with high probability.*

The remainder of this section is devoted to the proof of this theorem. In particular, we will use our cost breakdown, and express the competitive ratio as (where $\sigma = \sigma(\mathcal{D})$):

$$\rho = \max_{\sigma} \left(\frac{ON(\sigma)}{OFF(\sigma)} \right) = \max_{\sigma} \left(\frac{ON_{mig}(\sigma) + ON_{elim}(\sigma) + ON_{ins}(\sigma) + ON_{final}(\sigma)}{OFF_{comm}(\sigma) + OFF_{mig}(\sigma)} \right)$$

We will prove that each cost component in ON is competitive to OFF 's overall cost, therefore resulting in an $O(\log n \cdot OFF(\sigma))$ bound.

Elimination Costs To calculate the elimination cost (the total cost resulting from waiting at different configurations until $Cond()$ holds for the current configuration), we divide all configurations into two sets: configurations c for which $p(c) \leq 20p_{\min}$ and configurations c' for which $p(c') > 20p_{\min}$. We consider the elimination cost for these two sets in turn.

- All configurations c for which $p(c) \leq 20p_{\min}$. We will consider again two cases. Let $e[c]$ the cost of elimination on a position c (number of requests served until the condition of elimination of c is fulfilled). Either $e[c] \leq 20 \log n$ or $e[c] > 20 \log n$. In the first case we can just say that the number of configuration we have to eliminate is in $O(ON_{migr})$ and so $\sum_{e(c_i) \leq \log n} e(c_i) \leq O(\log n \cdot ON_{migr}) = O(\log n \cdot OFF)$.

For the other case, where $e(c_i) > 20 \cdot \log n$, we use the following claim:

Claim. Let $\Delta = [t_1, t_2]$ be a time interval. We note $r[c](\Delta) = r[c](t_2) - r[c](t_1)$, where $r[c](t)$ is the number of requests on the configuration c at the time t . Then :

If $p(c_j) \leq 20p(c_i)$ and $r[c_j](\Delta) \geq 20 \log n$ then w.h.p. $r[c_j](\Delta) \leq 40r[c_i](\Delta)$.

Proof. First note that from the bound of Eq. (3) w.h.p. $r[c_j](\Delta) \leq 2E[r[c_j](\Delta)]$. Similarly since $E[r[c_i]] \geq \frac{1}{20}E[r[c_j]]$ we have that w.h.p. $r[c_i](\Delta) \geq \frac{1}{2}E[r[c_i](\Delta)] \geq \frac{1}{40}E[r[c_j](\Delta)]$. So w.h.p. $r[c_j](\Delta) \leq 40r[c_i](\Delta)$. \square

From the above lemma and union bound over at most n states we get that w.h.p. $r[c_j](\Delta_j) \leq 40r_{c_{\min}}(\Delta_j)$ for all such configurations, with Δ_j denoting the time interval where we stayed on the configuration c_j , and c_j was not eliminated.

So

$$\begin{aligned} \sum_{e(c_i) \leq \log n} e(c_i) &= \sum_{e(c_i) \leq \log n} r[c_i](\Delta_i) \leq \sum_{e(c_i) \leq \log n} 20r[c_{\min}](\Delta_i) \\ &\leq 20r[c_{\min}]([0, |\sigma|]) = 20r[c_{\min}] \leq O(OFF_{comm}) \end{aligned}$$

In conclusion as $\text{ON}_{\text{elim} \leq 20} = \sum_{e(c_i) \leq \log n} e(c_i) + \sum_{e(c_i) > \log n} e(c_i)$ we have w.h.p.:

$$\frac{\text{ON}_{\text{elim} \leq 20}(\sigma)}{\text{OFF}(\sigma)} = O(1)$$

– All configurations c' for which $p(c') > 20p_{\min}$. For this we claim:

Claim. If $p(c_j) \geq 20p(c_i)$ and $r[c_j] \geq 20 \log n$ then w.h.p. $r[c_j] > 5r[c_i]$ and $\text{Cond}(j, i)$ is True for $\epsilon = \frac{1}{n^2}$.

Proof. Since $r[c_j] \geq 20 \log n$ w.h.p. $E[r[c_j]] \leq 2r[c_j]$. If $r[c_i] > \frac{1}{5}r[c_j]$ then w.h.p. $E[r[c_i]] > \frac{1}{10}r[c_j]$, but this contradicts the assumption that $E[r[c_i]] \leq \frac{1}{20}E[r[c_j]]$. So we have $\frac{r[c_i]}{r[c_j]} \leq \frac{1}{5}$ and $\text{Cond}(j, i)$ holds for $\epsilon = \frac{1}{n^2}$. \square

Now since the number of configurations ON needs to eliminate is lower than $\text{ON}_{\text{mig}}/\alpha \leq \text{ON}_{\text{mig}}$, the total cost ON paid is $O(\text{ON}_{\text{mig}} \cdot \log n)$. But since $\frac{\text{ON}_{\text{mig}}(\sigma)}{\text{OFF}(\sigma)} = O(1)$ (as we show next) we have :

$$\frac{\text{ON}_{\text{elim} > 20}(\sigma)}{\text{OFF}(\sigma)} = O(\log n)$$

To conclude $\text{ON}_{\text{elim}} = \text{ON}_{\text{elim} \leq 20} + \text{ON}_{\text{elim} > 20}$, and: $\text{ON}_{\text{elim}}(\sigma)/\text{OFF}(\sigma) = O(\log n)$.

Migration Cost We distinguish two cases. Let c_{far} be the farthest configuration reached by our online algorithm. Either d_{far} (the distance between c_{far} and c_0) is lower than d_{OFF} , or it is greater than d_{OFF} .

– In the first case, $d_{OFF} \geq d_{far}$, we can prove

Lemma 2. *if $d_{OFF} \geq d_{far}$ then $\text{ON}_{\text{mig}} \leq 6 \cdot \text{OFF}_{\text{mig}}(\sigma)$.*

Proof. $\exists x \in \mathbb{N} \ 2^{2x} \leq d_{far} < 2^{2x+2}$. Then, in the worst case, we have to go to 2^{2x+2} . So

$$\text{ON}_{\text{mig}}(\sigma) \leq \sum_{i=0}^{2x+1} 3 \cdot 2^i \cdot \alpha \leq 6 \cdot 2^{2x+1} \cdot \alpha \leq 6d_{far} \alpha \leq 6 \cdot d_{OFF} \cdot \alpha \leq 6 \cdot \text{OFF}_{\text{mig}}(\sigma)$$

\square

– If $d_{OFF} < d_{far}$, then from Claims 5.2 and claim 5.2 with $\Delta = [0, |\sigma|]$ it follows that w.h.p. $r[c_{OFF}] \geq \Omega(\alpha \cdot d_{far})$: Recall that in our algorithm (line 15) we only move beyond the current radius if the corresponding costs have been amortized. Hence $\text{ON}_{\text{mig}} \leq \text{OFF}_{\text{comm}}$.

In conclusion, in both cases: $\text{ON}_{\text{mig}}(\sigma)/\text{OFF}(\sigma) = O(1)$.

Insurance Costs For the insurance cost we also consider several cases. Let c_{far} be the farthest configuration reached by our online algorithm. Let c_{OFF} denote the location of the offline algorithm. We split ON_{ins} into two parts: $ON_{ins < far}$ and $ON_{ins = far}$. $ON_{ins < far}$ is the insurance cost up to (not including) c_{far} while $ON_{ins = far}$ is the insurance cost paid on c_{far} . The last insurance cost, paid before the last migration to c_{far} , is αd_{far} , so we have $ON_{ins < far} \leq O(ON_{mig}) = O(OFF)$ (see the migration cost analysis).

The only possible problem is therefore $ON_{ins = far}$. Now we consider two cases:

- c_{OFF} is in \mathcal{E} (eliminated configuration). Since c_{OFF} was eliminated before c_{far} it follows from Claims 5.2 and 5.2 that w.h.p. $r[c_{OFF}] > \Omega(r[c_{far}])$ so $ON_{ins = far} < O(OFF_{comm})$.
- c_{OFF} is in $\bar{\mathcal{E}}$. In this case because of our *searching path* and the selection of c_{next} , we have $d_{OFF} \geq d_{next}/2$. Therefore $ON_{ins = far} \leq O(OFF_{mig})$.

Overall we have: $ON_{ins}(\sigma)/OFF(\sigma) = O(1)$.

Final Costs By definition, in the final configuration, all other configurations have been eliminated. Thus, our condition, $\text{Cond}(j, i)$, has been fulfilled at some point for any c_j , with respect to some c_i . The probability that we eliminate a minimum configuration and end up at a suboptimal configuration is small. This follows from Lemma 1, when setting $\epsilon := \frac{1}{n^2}$: once we stopped in a configuration, it is, with high probability, a (not necessarily unique) minimal configuration. Since OFF directly moves to a minimum configuration (which may not be unique), ON cannot incur a higher cost than OFF on a specific minimum configuration, i.e., not more than $r[c_{\min}]$. As the offline algorithm moved from the start to a configuration c_{OFF} and $r[c_{\min}]$ is the configuration with the lowest number of requests, $r[c_{OFF}] \geq r[c_{\min}]$. Thus, $ON_{final}(\sigma) \leq OFF(\sigma)$, and also $ON_{final}(\sigma)/OFF(\sigma) = O(1)$.

Overall Costs In conclusion, with high probability:

$$\rho \leq \max_{\sigma} \left(\frac{ON_{mig}(\sigma) + ON_{elim}(\sigma) + ON_{ins}(\sigma) + ON_{final}(\sigma)}{OFF_{comm}(\sigma) + OFF_{mig}(\sigma)} \right) = O(\log n)$$

6 Related Work

Our paper takes a novel perspective on a range of classic problems. First, clustering and graph partitioning problems as well as repartitioning problems [21] have been studied for many years and in many contexts. These problems are usually NP-complete and even hard to approximate [2]. Especially partitioning problems for two clusters ($\ell = 2$ in our case), known as minimum bisection problems [9], have been studied intensively. Minimum bisection problems are known to allow for good, $O(\log^{1.5} n)$ -factor approximations [13]. Problem variants with $k = 2$ correspond to maximum matching problems, which are polynomial-time solvable. In contrast to our work however, these models assume an offline perspective where the problem input is given ahead of time. In the online world, our problem is related to page (resp. file) migration [4,6] and server migration [5] problems: in these problems, a server needs to be migrated close to requests occurring on a graph, trading off access and migration costs. In the former problem variant, migration costs relate to distance; in the latter, migration costs

relate to the available bandwidth along migration paths. Moreover, in our problem, a ski-rental resp. rent-or-buy like tradeoff between migration and communication costs needs to be found. However, migrations do not occur along a graph but between clusters, and multiple nodes can be migrated simultaneously. The large configuration space also renders solutions based on metrical task system approaches [7] inefficient. Another interesting connection exists to k -server problems [11], where multiple servers can “collaboratively” serve requests. In some sense, our problem can be seen as the opposite problem, where rather than aiming to move servers to the locations where the requests occur, we aim to move away and *avoid* configurations (i.e., cuts) where requests occur. More importantly, compared to classic online migration problems where requests define a unique optimal location from which they can be served at minimal cost (namely at the corresponding graph vertex), in our case, a request only reveals very limited information about the optimal (minimal cost) configuration. In other words, a single request only contains very limited information about how good a current clustering is, and how far (in terms of migrations) we are from an optimal offline location.

Our model can be seen as a generalization of online paging [10,14,15,20,22], and especially its variants *with bypassing* [1,8]. However, in general, in our model, the “cache” is *distributed*: requests occur *between* nodes and *not to* nodes, and costs can be saved by collocation.

Our problem also has connections to online packing problems, where items of different sizes arriving over time need to be packed into a minimal number of bins [18,19]. In contrast to these problems, however, in our case the objective is not to minimize the number of bins but rather the number of “links” between bins, given a fixed number of bins.

The paper closest to ours is [3] which studies online partitioning problems from a deterministic perspective, i.e., σ is generated in a deterministic manner. In this setting, it has been shown that the competitive ratio is inherently high, at least linear in k , and even if the online algorithm is allowed to use larger clusters than the offline algorithm (scenario with augmentation). We in this paper initiate the study of stochastic models where request patterns are drawn from an unknown but fixed distribution, and show that polylogarithmic bounds can be achieved under ring patterns, even without augmentation.

In general, we believe that a key conceptual contribution of our model itself regards the underlying combination of learning and searching. Indeed, while the fundamental problem of how to efficiently learn a distribution has been explored for many decades [17], our perspective comes with an additional locality requirement, namely that searching induces costs (i.e., migrations).

7 Conclusion

This paper initiated the study of a natural cluster learning problem where the search procedure entails costs: communication costs occur in “suboptimal” clustering configurations and migration costs occur when switching between configurations. In particular, we presented an efficient online clustering algorithm which performs well even if compared to an offline algorithm which knows the distribution of the communication pattern ahead of time. Indeed, the $O(\log k)$ competitive ratio is interesting as k is likely to be small in the applications considered in this paper: k corresponds to the number of virtual machines that can be hosted on the same server, e.g., the number of cores. Moreover, we believe that our online approach is interesting in practice as it does not rely on any assumptions on the communication distribution, which may turn out to be wrong.

We believe that our work sheds an interesting new light on multiple classic problems, and opens an interesting field for future research. In particular, it would be interesting to know whether similar competitive ratios can be achieved even for more general communication patterns. Moreover, so far we have only focused on deterministic algorithms, and the exploration of randomized algorithms constitutes another interesting avenue for future research.

References

1. Adamaszek, A., Czumaj, A., Englert, M., Räcke, H.: An $O(\log k)$ -competitive algorithm for generalized caching. In: Proc. 23rd SODA. pp. 1681–1689 (2012)
2. Andreev, K., Räcke, H.: Balanced graph partitioning. *Theory of Computing Systems* 39(6), 929–939 (2006)
3. Avin, C., Loukas, A., Pacut, M., Schmid, S.: Online balanced repartitioning. In: Proc. 30th International Symposium on Distributed Computing (DISC) (2016)
4. Bartal, Y., Charikar, M., Indyk, P.: On page migration and other relaxed task systems. *Theoretical Computer Science* 268(1), 43–66 (2001), also appeared in *Proc. of the 8th SODA*, pages 43–52, 1997
5. Bienkowski, M., Feldmann, A., Grassler, J., Schaffrath, G., Schmid, S.: The wide-area virtual service migration problem: A competitive analysis approach. *IEEE/ACM Transactions on Networking (ToN)* (2014)
6. Black, D.L., Sleator, D.D.: Competitive algorithms for replication and migration problems (1989)
7. Borodin, A., Linial, N., Saks, M.E.: An optimal on-line algorithm for metrical task system. *Journal of the ACM* 39(4), 745–763 (1992), also appeared in *Proc. of the 19th STOC*, pages 373–382, 1987
8. Epstein, L., Imreh, C., Levin, A., Nagy-György, J.: Online file caching with rejection penalties. *Algorithmica* 71(2), 279–306 (2015)
9. Feige, U., Krauthgamer, R.: A polylogarithmic approximation of the minimum bisection. *SIAM Journal on Computing* 31(4), 1090–1118 (2002)
10. Fiat, A., Karp, R.M., Luby, M., McGeoch, L.A., Sleator, D.D., Young, N.E.: Competitive paging algorithms. *Journal of Algorithms* 12(4), 685–699 (1991)
11. Fiat, A., Rabani, Y., Ravid, Y.: Competitive k -server algorithms. *J. Comput. Syst. Sci.* 48(3), 410–428 (1994)
12. Franck, W.: An optimal search problem. *SIAM review* 7(4), 503–512 (1965)
13. Krauthgamer, R., Feige, U.: A polylogarithmic approximation of the minimum bisection. *SIAM Review* 48(1), 99–130 (2006)
14. McGeoch, L.A., Sleator, D.D.: A strongly competitive randomized paging algorithm. *Algorithmica* 6(6), 816–825 (1991)
15. Mendel, M., Seiden, S.S.: Online companion caching. *Theoretical Computer Science* 324(2–3), 183–200 (2004)
16. Mitzenmacher, M., Upfal, E.: *Probability and Computing: Randomized Algorithms and Probabilistic Analysis*. Cambridge University Press, New York, NY, USA (2005)
17. Pöschel, T., Ebeling, W., Rosé, H.: Guessing probability distributions from small samples. *Journal of statistical physics* 80(5–6), 1443–1452 (1995)
18. Ramanan, P.V., Brown, D.J., Lee, C.C., Lee, D.T.: On-line bin packing in linear time. *Journal of Algorithms* 10(3), 305–326 (1989)
19. Seiden, S.S.: On the online bin packing problem. *Journal of the ACM* 49(5), 640–671 (2002)
20. Sleator, D.D., Tarjan, R.E.: Amortized efficiency of list update and paging rules. *Communications of the ACM* 28(2), 202–208 (1985)
21. Vaquero, L., Cuadrado, F., Logothetis, D., Martella, C.: Adaptive partitioning for large-scale dynamic graphs. In: Proc. 4th Annual Symposium on Cloud Computing (SOCC). pp. 35:1–35:2 (2013)
22. Young, N.E.: On-line caching as cache size varies. In: Proc. of the 2nd ACM-SIAM Symp. on Discrete Algorithms (SODA). pp. 241–250 (1991)