# Fast Distributed $k$-Means with a Small Number of Rounds

Tom Hess[1]     Ron Visbord[2]     Sivan Sabato[3]

[1,3]Department of Computer Science, Ben Gurion University, Beer-Sheva, Israel
[2]Independent researcher
[1]`tomhe@bgu.ac.il`, [2]`ronvisbord@gmail.com`, [3]`sabatos@cs.bgu.ac.il`

## Abstract

We propose a new algorithm for $k$-means clustering in a distributed setting, where the data is distributed across many machines, and a coordinator communicates with these machines to calculate the output clustering. Our algorithm guarantees a cost approximation factor and a number of communication rounds that depend only on the computational capacity of the coordinator. Moreover, the algorithm includes a built-in stopping mechanism, which allows it to use fewer communication rounds whenever possible. We show both theoretically and empirically that in many natural cases, indeed $1-4$ rounds suffice. In comparison with the popular $k$-`means||` algorithm, our approach allows exploiting a larger coordinator capacity to obtain a smaller number of rounds. Our experiments show that the $k$-means cost obtained by the proposed algorithm is usually better than the cost obtained by $k$-`means||`, even when the latter is allowed a larger number of rounds. Moreover, the machine running time in our approach is considerably smaller than that of $k$-`means||`. Code for running the algorithm and experiments is available at `https://github.com/selotape/distributed_k_means`.

## 1   Introduction

Modern datasets can be very large, requiring algorithms that can handle massive amounts of data. This need drives the development of distributed algorithms, which use many machines that work in parallel to solve the given problem faster. In some cases, the data is already split among separate machines, again calling for a distributed solution. In this work, we study the classical problem of $k$-means clustering Sebestyen (1962) in the distributed setting. The goal of a $k$-means clustering algorithm is to select cluster centers from the input dataset that induce a $k$-means cost as close as possible to the smallest cost that can be obtained for the dataset. In a distributed framework, a main bottleneck in many practical settings is communication. Most distributed algorithms run in communication rounds, where in each round each machine performs an individual task, and the machines synchronize and communicate after each round. The number of rounds is a crucial factor in the practical performance of distributed algorithms, since each such round requires synchronization and communication between the machines, which are costly and can cause time delays. Therefore, reducing the number of rounds as much as possible is a key goal for distributed algorithms.

We focus on a common practical distributed computation model (Ene et al., 2011; Guha et al., 2019), in which one machine, called the *coordinator*, communicates with all other machines, while the data to cluster is distributed among the machines. We consider the case where the coordinator is

capable of running heavier computations, while the machines are more limited in their computation power and do not communicate among themselves. This model is suitable, for instance, when the dataset to cluster is partitioned between low-end mobile devices, and the coordinator is a stronger machine. The data may be split among the machines for the purpose of performing the distributed computation, or it may be partitioned among the devices to begin with, for instance if each device has independently collected data points (e.g., by taking pictures using the device's camera). We do not make any assumption on the partition of the data, thus we support also non-i.i.d. data. We note that the model that we consider is different from the Federated Learning model **??**, which emphasizes other requirements, such as privacy.

One of the most popular distributed clustering algorithms is $k$-`means||` (Bahmani et al., 2012). This algorithm approximates the optimal $k$-means cost on the dataset up to a constant approximation factor, assuming that this cost is bounded away from zero (see the example and discussion in Bachem et al. 2017a), and that a sufficient number of communication rounds is performed. However, $k$-`means||` does not have an adaptive mechanism to decide how many communication rounds to run. Therefore, the number of rounds is usually set heuristically, in which case the guarantee for a constant approximation factor might not hold. Other distributed algorithms (e.g., Balcan et al., 2013) use only a single round of communication by definition, but do not scale well when the number of machines is large.

In this work, we propose the new distributed $k$-means clustering algorithm, `SOCCER` (Sampling, Optimal Clustering Cost Estimation, Removal), which guarantees a constant approximation factor that depends only on the computational capacity of the coordinator, without requiring the optimal clustering cost to be bounded away from zero. Moreover, the algorithm automatically stops once a sufficient number of rounds has been completed, which can be much earlier than the worst-case number of rounds. We demonstrate that in many natural datasets, the number of rounds required by `SOCCER` is much smaller than the worst-case upper bound. In particular, we prove that `SOCCER` stops after a single round if the dataset is drawn from a high-dimensional Gaussian mixture. In addition, we prove that there are datasets such that `SOCCER` stops after one round and obtains a constant approximation factor, while $k$-`means||` requires $k - 1$ rounds for the same result.

We empirically compare `SOCCER` to $k$-`means||` on synthetic and real datasets, showing that indeed in practical scenarios, `SOCCER` stops after $1 - 4$ rounds. In contrast, $k$-`means||` does not have a stopping condition, and when stopped after the same or a similar number of rounds, it usually obtains a worse clustering cost. Moreover, the machine run time of `SOCCER` is almost always significantly smaller than that of $k$-`means||` for a comparable final cost.

Our technique is based on letting the coordinator run a (centralized) clustering algorithm on a limited number of points, and using this clustering to calculate an estimate of a truncated version of the optimal $k$-means cost on the dataset. This provides information to the machines that allows them to progressively remove points from their part of the dataset. When all the points are removed, the algorithm stops and calculates the final clustering from the centers selected in the centralized clustering runs. `SOCCER` combines clustering approaches designed for two different settings: The distributed setting (Ene et al., 2011) and the online setting (Hess et al., 2021). Ene et al. (2011) iteratively samples points from the machines and then removes points that are close to them from consideration. We show that calculating a clustering on the point sample, along with a technique adapted from Hess et al. (2021), lead to a more accurate removal of points. This provides a practical and successful algorithm with approximation guarantees that depend only on the number of points that the coordinator can cluster.

**Our contribution**    To summarize, `SOCCER` is a new distributed $k$-means algorithm that is

equipped with theoretical guarantees on its cost approximation factor and number of communication rounds, and requires an even smaller number of rounds in practice. Our experiments demonstrate its practical advantages in comparison with $k$-means||, in a distributed model which allows the coordinator to calculate a clustering on a limited number of points. Some of the proofs and experiment results are deferred to the appendices.

## 2   Related work

A naive approach to distributed clustering would be to implement a centralized algorithm in a straightforward manner under the distributed model. However, this tends to be impractical, since it requires a large number of communication rounds (see, e.g., the discussion in Bahmani et al., 2012). Therefore, algorithms that are specifically tailored to the distributed setting have been suggested. Many of the algorithms that we mention below select more than $k$ centers. It is then standard to use a weighted centralized $k$-means algorithm to reduce the number of centers to exactly $k$. It is known (e.g., Guha et al., 2003, Theorem 4) that this preserves approximation guarantees up to constants.

One common technique used in many distributed algorithms has the following structure: Each machine calculates a set of representatives of its own data (sometimes called *coresets*). these are then sent to the coordinator, which uses them to calculate a set of centers (Ailon et al., 2009; Balcan et al., 2013; Feldman et al., 2020; Bachem et al., 2017b). These algorithms require a small constant number of communication rounds. However, the technique has the drawback that the run time and the memory size of the coordinator increase with the number of machines after suppressing a certain threshold, while in SOCCER the running time of the coordinator and the machines improves linearly with the number of the machines (see experiments in Bahmani et al. 2012). Some works address the setting of distributed $k$-means with outliers (Guo and Li, 2018; Guha et al., 2019; Chen et al., 2018). These algorithms also require coordinator resources that increase with the number of machines. Other distributed algorithms obtain superior guarantees, but under strong structural assumptions on the data, such as a small aspect ratio or perturbation-resilient instances (Voevodski, 2021), or on the partition of the data into machines (Bhaskara and Wijewardena, 2018).

As mentioned above, one of the most successful distributed $k$-means algorithms to date is $k$-means||(Bahmani et al., 2012), which is widely used in practice (e.g., in the MLLib library of Apache Spark, Meng et al., 2016) and also has theoretical guarantees. $k$-means|| proposes a distributed seeding algorithm that selects a small number of potential centers. The worst-case number of rounds of $k$-means|| is $O(\log(n/\mathsf{opt}))$, where $\mathsf{opt}$ is the optimal $k$-means cost of the dataset. This guarantee requires the optimal $k$-means cost to be bounded away from zero (see also the example in Bachem et al., 2017a). Bachem et al. (2017a) show that if the variance of the dataset is bounded and $\mathsf{opt}$ is bounded away from zero, then $k$-means|| can be stopped after a constant number of rounds. However, this requires additional information about the dataset.

Ene et al. (2011) proposed a distributed $k$-median algorithm (which can easily be adapted to $k$-means) with a number of communication rounds that depends on the memory size of the coordinator. In each round, each machine draws two random sub-samples from its data, and sends them to the coordinator. The coordinator adds the first sample from each machine to the output clustering, and uses the second sample to calculate a threshold using a simple quantile statistic. Then, the threshold and most of the points received by the coordinator are sent to all the machines. Each machine then removes from its dataset the points whose distance to the current clustering does not exceed the threshold. The total number of removed points is by definition a fixed fraction

of the dataset. The final round occurs when the remaining points in the machines fit entirely in the coordinator memory. Chen et al. (2016) proposed a variation on this idea that reduces the total communication, while increasing the number of rounds. Kumar et al. (2015) generalized this technique to other related problems. Despite its theoretical guarantees, the algorithm of Ene et al. (2011) has significant disadvantages in practice. First, it always uses the worst-case number of rounds. In addition, in practice, on reasonable dataset sizes, the sub-sampling does not significantly reduce the number of points relative to the original dataset. This means that calculating the final clustering is not much easier than calculating a clustering on the original dataset. In addition, the number of points sent from the coordinator to the machines is large, leading to both to a heavy communication requirement and a heavy computation in each machine. We show below how the approach of SOCCER avoids these issues.

Most existing algorithms for distributed clustering can be applied to both the $k$-means and $k$-medoids formulations with slight adaptations, and so the same body of work is generally relevant for both formulations. The difference in formulation between $k$-means and $k$-medoids is manifested in the constant approximation factor when using a black box offline clustering algorithm.

# 3   Setting and notation

For an integer $l$, denote $[l] := \{1, \ldots, l\}$. Let $(X, \rho)$ be a finite metric space, where $X$ is a set of size $n$ and $\rho : X \times X \to \mathbb{R}_+$ is a metric. For a point $x \in X$ and a set $T \subseteq X$, let $\rho(x, T) := \min_{y \in T} \rho(x, y)$. For simplicity, we use set notations for datasets, although they can include duplicates. For an integer $k \geq 2$, a $k$-clustering of $X$ is a set of (at most) $k$ points from $X$ which represent cluster centers. Given a set $S \subseteq X$, the $k$-means cost of $T$ on $S$ is $\mathsf{cost}(S, T) := \sum_{x \in S} \rho(x, T)^2$. The goal when clustering $X$ is to find a clustering $T$ with a low cost $\mathsf{cost}(X, T)$. We denote by OPT an optimal $k$-means clustering: $\text{OPT} \in \operatorname{argmin}_{T \subseteq X, |T| \leq k} \mathsf{cost}(X, T)$.

A (centralized) $k$-means algorithm $\mathcal{A}$ takes as input a finite set of points $S$ and the parameter $k$, and outputs a $k$-clustering of $S$, denoted $\mathcal{A}(S, k)$. For $\beta \geq 1$, $\mathcal{A}$ is a $\beta$-approximation $k$-means algorithm on $(X, \rho)$, if for all input sets $S \subseteq X$, $\mathsf{cost}(S, \mathcal{A}(S, k)) \leq \beta \cdot \mathsf{cost}(S, \text{OPT}_S)$, where $\text{OPT}_S$ is an optimal solution on $S$ with centers from $S$: $\text{OPT}_S \in \operatorname{argmin}_{T \subseteq S, |T| \leq k} \mathsf{cost}(S, T)$. In the centralized setting, the best known approximation constant for an efficient $k$-means algorithm is 9 for a general metric space and 6.357 for Euclidean spaces (Ahmadian et al., 2019).

In the coordinator model (Guha et al., 2019) which we study, the data $X$ is arbitrarily partitioned among $m$ machines, where $X_j$ denotes the set of points in machine $j \in [m]$. The machines communicate directly only with the coordinator. Broadcasts from the coordinator to the machines are counted as a single transmission. The computation is conducted in rounds, where in each round the machines perform an individual task and then communicate with the coordinator.

# 4   The guarantees of SOCCER

In this section, we present the guarantees of SOCCER, our new distributed $k$-means algorithm, which is described in detail in Section 5. Similarly to Ene et al. (2011), we assume a bound of $\tilde{\Theta}(kn^\epsilon)$ on the number of points that can be stored in the memory of the coordinator, where $\epsilon \in (0, 1)$ is a parameter linking the coordinator size with the dataset size. Specifically, we assume that the coordinator can calculate a (centralized) clustering over a dataset of size $\eta(\epsilon) = 36kn^\epsilon \log(\frac{1.1k}{\delta\epsilon})$, and can store the same order of magnitude of data points. We further assume it has access to a

centralized black-box $k$-means algorithm $\mathcal{A}$ that can be used for this purpose. The clustering is used by SOCCER to calculate an estimate of a truncated version of the optimal attainable $k$-means cost for the dataset $X$. As mentioned above, most distributed algorithms select more than $k$ centers. This number can then reduced to $k$ using a standard weighted clustering technique. SOCCER selects only slightly more than $k$ centers, making the final reduction step easier. The following theorem gives the guarantees of SOCCER. The proof is provided in Section 6 and the appendices referenced there.

**Theorem 4.1.** *Suppose that the size of the dataset $X$ is a sufficiently large $n$. Suppose that SOCCER runs with a confidence parameter $\delta \in (0,1)$, a coordinator parameter $\epsilon \in (0,1)$, and number of centers $k \geq 5$, and suppose that the black-box algorithm $\mathcal{A}$ is a $\beta$-approximation $k$-means algorithm. Denote the total number of communication rounds until SOCCER stops by $I$, and denote the set of cluster centers it selects by $C_{\mathrm{out}}$. Then, with probability at least $1 - \delta$,*

- *$I < \frac{1}{\epsilon} - 1$;*

- *$|C_{\mathrm{out}}| \leq I \cdot (k + 9 \log \frac{1.1k}{\delta\epsilon})$;*

- *$\mathsf{cost}(X, C_{\mathrm{out}}) \leq I \cdot (80\beta + 44) \cdot \mathsf{cost}(X, \mathrm{OPT})$;*

- *The total number of points transmitted to the coordinator is at most $I \cdot \eta(\epsilon) = 72Ikn^\epsilon \log(\frac{1.1k}{\delta\epsilon})$.*

- *The total number of points broadcasted from the coordinator is at most $I \cdot (k + 9 \log \frac{1.1k}{\delta\epsilon})$.*

We note that while the theorem above lists specific constants, these are in fact interdependent, so that, for instance, one can allow a larger coordinator memory constant, to obtain a significantly smaller cost approximation constant; see also the discussion in Section 6.

Before presenting the algorithm, we compare the guarantees above to the closest relevant results. In comparison with the algorithm of Ene et al. (2011) (henceforth EIM11), SOCCER uses the same number of communication rounds in the worst case. However, as will be made evident below, unlike EIM11, it can use considerably fewer rounds on many natural datasets. SOCCER selects $\tilde{O}(k)$ centers, and these are all the points that the coordinator ever broadcasts to the machines. In contrast, EIM11 selects $\Omega(kn^\epsilon \log(n))$ centers and broadcasts all of them, thus its total communication to the machines is significantly larger. This also affects the computation resources required from the machines, as discussed in more detail in Section 5. Like SOCCER, EIM11 also obtains a constant approximation factor. While its approximation constant is smaller, the issues mentioned above make the algorithm impractical, as we observe in Section 8.

To compare these guarantees to $k$-means||, note that the worst-case number of rounds of $k$-means|| is $O(\log(n/\mathsf{opt}))$, while in the theorem above (as in Ene et al., 2011) it is $1/\epsilon = \tilde{O}(\log(n)/\log(L/k))$, where $L$ is the limitation on the coordinator. If $L$ is set to $\tilde{\Theta}(k)$ then the worst-case number of rounds of SOCCER is similar to that of $k$-means||, except that it does not require $\mathsf{opt}$ to be bounded away from zero. In addition, and unlike $k$-means||, in this approach a larger $L$ can be used to reduce the worst-case number of iterations. Moreover, as seen below, SOCCER stops on its own when the number of rounds is sufficient for the dataset. In contrast, the actual number of rounds of $k$-means|| is a hyper-parameter. In the next section, we give the full description of SOCCER.

---

**Algorithm 1** SOCCER

---

**input** $\delta \in (0,1)$ (confidence), $k \in \mathbb{N}$, $n \in \mathbb{N}$ (data size), $\mathcal{A}$ (a centralized $k$-means algorithm), $\epsilon \in (0,1)$ (coordinator parameter), $m \in [n]$ (number of machines).

At the beginning of the run, machine $j$ holds data $X_j$, where $X := \cup_{j \in [m]} X_j$.

1: $C_{\text{out}} \leftarrow \emptyset$, $N \leftarrow n$.
2: **while** $N > \eta(\epsilon)$ **do**
3:     $\alpha \leftarrow \eta(\epsilon)/N$.
4:     For $l \in \{1,2\}$, each machine $j$ adds each point in $X_j$ to a set $P_j^l$ with independent probability $\alpha$.
5:     Each machine $j$ sends $P_j^1$ and $P_j^2$ to the coordinator.
6:     In the coordinator:
7:         $P_1 := \bigcup\limits_{j \in [m]} P_j^1, P_2 := \bigcup\limits_{j \in [m]} P_j^2$.
8:         $C_{\text{iter}} \leftarrow \mathcal{A}(P_1, k_+)$.
9:         $v := 2\mathsf{cost}_{\frac{3}{2}(k+1)d_k}(P_2, C_{\text{iter}})/(3kd_k)$.    △ *See definition in Section 5*
10:       $C_{\text{out}} \leftarrow C_{\text{out}} \cup C_{\text{iter}}$.
11:       Broadcast $(v, C_{\text{iter}})$ to each of the machines.
12:     Removal: Each machine $j$ updates:
       $X_j \leftarrow \{x \in X_j \mid \rho(x, C_{\text{iter}})^2 > v\}$.
13:     Each machine sends $N_j := |X_j|$ to the coordinator, which sets $N \leftarrow \sum_{j \in [m]} N_j$.
14:     $V \leftarrow \sum\limits_{j \in [m]} X_j$.    △ *The machines sends $|X_j|$ to the coordinator.*
15: **end while**
16: All the machines send $X_j$ to the coordinator, which sets $V \leftarrow \cup_{j \in [m]} X_j$.
17: The coordinator calculates $C_{\text{out}} \leftarrow C_{\text{out}} \cup \mathcal{A}(V, k)$.
18: **return** $C_{\text{out}}$.

---

## 5   The SOCCER algorithm

SOCCER is listed in Alg. 1. It uses the notations $k_+ := k + 9\log(1.1k/(\delta\epsilon))$, $d_k := 6.5\log(1.1k/(\delta\epsilon))$. The underlying structure of SOCCER is superficially similar to that of EIM11, which was described in Section 2. It runs a loop, where in each iteration, each machine sends the coordinator a sub-sample of its points. The coordinator then sends data points and a threshold to all machines. Then, each machine removes from its data the points that are closer to the sent points than the threshold. This is repeated in rounds, until the number of remaining points is small enough so they can be stored in full in the coordinator.

Despite the external similarity in structure, SOCCER is crucially different from EIM11 and its variants, which send most of the points received by the coordinator to the machines, and remove a fixed fraction of the dataset in each round. The coordinator in SOCCER uses the sub-samples received from the machines as input to the centralized black-box $k$-means clustering algorithm $\mathcal{A}$. Then, it calculates an estimate of the truncated $k$-means cost of the centers selected by $\mathcal{A}$ on the entire dataset. This estimate is then used to calculate the threshold that the machines use to remove points from their dataset. The method for estimating the cost and calculating the threshold is based on a technique first proposed in Hess et al. (2021), which addresses a different setting of (centralized) online no-substitution clustering. In that work, the estimate is used for the purpose

of on-the-fly center selection, when clustering a stream of points. Our analysis shows how this type of estimate can be used to improve performance in the distributed setting, despite its original use for a completely different purpose.

In SOCCER, in each iteration (corresponding to a communication round), each machine $j$ creates two sub-samples from its dataset, $P_j^1$ and $P_j^2$, which are then sent to the coordinator. These sub-samples are drawn independently at random from the machine's current set of points, where their sizes are set so that the total number of points sent to the coordinator by all machines is $\eta(\epsilon)$. The coordinator merges these sub-sample pairs into the respective sets $P_1$ and $P_2$. It then calculates a $k_+$-means clustering on $P_1$ using $\mathcal{A}$, denoted $C_{\text{iter}}$, and calculates a threshold using the *truncated* cost of $C_{\text{iter}}$ on $P_2$: For two sets $S, T \subseteq X$ and an integer $l$, the $l$-truncated cost of $T$ on $S$, denoted $\text{cost}_l(S, T)$, is the total cost of the clustering after removing the $l$ points in $S$ that incur the most cost.

The coordinator adds $C_{\text{iter}}$ to the output set $C_{\text{out}}$, and sends $v$ and $C_{\text{iter}}$ to each of the machines. Then, each machine removes from its dataset all the points whose distance from $C_{\text{iter}}$ is at most $\sqrt{v}$. Our analysis below shows that the truncated cost can be used to lower-bound the cost of points that belong to large clusters in the optimal $k$-means clustering of $X$. As a result, points that are $\sqrt{v}$-close to some center in $C_{\text{iter}}$ are sufficiently close to an optimal center to guarantee the final approximation factor.

Lastly, when sufficiently many points have been removed in each machine so that the entire remaining data can be handled by the coordinator, the loop terminates and the remaining points are sent to the coordinator, which calculates a $k$-clustering on them and adds the output centers to $C_{\text{out}}$.

We note that the main computational burden in the machines is to calculate the distances of the data points they store from the points broadcasted by the coordinator. Therefore, the number of broadcasted points needs to be small for this burden to be reasonable. Indeed, in SOCCER this number is only $k_+ = k + 9\log(1.1k/(\epsilon\delta))$. In contrast, in EIM11 this number is $9kn^\epsilon \log(n/\delta)$. Thus, for large datasets, the computational requirements from the machines in SOCCER are lighter by orders of magnitude than those of EIM11.

$C_{\text{iter}}$ and $v$ are calculated similarly to the centralized online clustering algorithm of Hess et al. (2021) mentioned above. However, our constants are significantly smaller, as a result of a tighter analysis (see Appendix A.1). The improvement of the constants is of significant practical importance: These constants are used by SOCCER. If they were too large, as in Hess et al. (2021), then SOCCER would be impractical. For instance, in Hess et al. (2021), the number of outliers removed when calculating the truncated cost is very large. Using the same number in SOCCER would have caused the fraction of removed points in each round to be too small, leading to a large number of rounds. Moreover, these constants cannot be easily changed without careful analysis, since they are inter-dependent. Finding an appropriate assignment of constants that makes the algorithm practical while guaranteeing the desired behaviour requires a delicate balance of many competing quantities.

In the next section, we state the main lemma that we prove to derive the guarantees of SOCCER.

## 6   Main lemma

We now give the main lemma that allows us to prove Theorem 4.1. First, we define necessary notation. Consider the contents of the machine datasets $\{X_j\}_{j \in [m]}$ at the beginning of iteration $i$ in line 12 of Alg. 1, and let $V_i := \cup_{j \in [m]} X_j$. Denote the points removed at iteration $i$ by $R_i := V_i \setminus V_{i+1}$.

Let $C_{\text{iter}}^i$ and $\alpha_i$ be the values of $C_{\text{iter}}$ and $\alpha$, respectively, as calculated by `SOCCER` at iteration $i$. To prove Theorem 4.1, we provide the following lemma, which is proved in Appendix A.1.

**Lemma 6.1.** *Assume that `SOCCER` runs with the parameters given in Theorem 4.1. Let $i \leq I$. With probability at least $1 - \delta\epsilon$,*

- $\text{cost}(R_i, C_{\text{iter}}^i) \leq (80\beta + 44) \cdot \text{cost}(V_i, \text{OPT})$;

- $|V_{i+1}| \leq 5.5kd_k/\alpha_i$.

The first part of this lemma shows that in round $i$, the calculated cluster $C_{\text{iter}}^i$ obtains a constant approximation on all the removed points in this round. This is later used to prove the overall approximation guarantee. The second part bounds the number of remaining points in each round, which is used to upper bound the number of rounds. Theorem 4.1 can now be proved using the lemma. The proof is provided in Appendix A.2

We note that the constants in Lemma 6.1 and, consequently, in Theorem 4.1, are interdependent. In particular, increasing the coordinator's capacity by a constant factor would decrease the cost approximation constant, since a larger memory constant would allow $P_1$ and $P_2$ to be larger, making them more representative of the full data, and leading to a smaller cost approximation factor. In addition, it would allow reducing the threshold for removal, again improving the accuracy at the expense of a larger coordinator capacity.

# 7 Beyond worst-case: Why `SOCCER` can stop after fewer rounds

As discussed above, a main desideratum of the distributed algorithm is to use a small number of communication rounds. While the worst-case number of rounds for `SOCCER` is $\Theta(1/\epsilon)$, it stops earlier if sufficiently many data points are removed from the machine datasets, so that the current total data size can be handled by the coordinator. If this is the case, then also the approximation factor and the number of selected centers are smaller, as can be see in Theorem 4.1.

We now show that indeed, `SOCCER` is likely to require fewer rounds on many natural datasets. `SOCCER` calculates in each round the clustering $C_{\text{iter}}$ based on the sub-sample sent from each machine. Our analysis shows that $C_{\text{iter}}$ obtains a near-optimal clustering cost on points that in the optimal solution belong to clusters that are larger than $d_k/\alpha = \tilde{O}(n^{1-\epsilon}/k)$. Such points will typically be sufficiently close to $C_{\text{iter}}$ to be removed from the machine dataset in the removal step. The number of points in small optimal clusters can be at most $kd_k/\alpha = O(n^{1-\epsilon})$. In many natural cases, and in particular when $n$ is sufficiently large, the optimal solution will have even fewer points, perhaps none, in such small clusters. Thus, almost all points will be removed in the first round. As a simple example, consider a dataset drawn from a $k$-Gaussian mixture. The following result shows that `SOCCER` requires a single round to cluster such a dataset. The proof is provided in Appendix A.3.

**Theorem 7.1.** *Let $X$ be a dataset drawn from a $k$-spherical Gaussian mixture. For sufficiently large $d$ and $n$, if $\epsilon \geq \log\log(n/\delta)/\log n$, then with high probability, `SOCCER` when running on $X$ will stop after one round, and output a clustering with a constant cost approximation factor.*

This property of `SOCCER` is contrasted with EIM11, which removes the same fraction of points in each round, regardless of the structure of the data, and so never stops early. To compare to `k-means||`, recall that it has no stopping mechanism and its number of rounds is set heuristically. Moreover, the following theorem, proved in Appendix A.4, shows that there are cases in which

$k$-means|| requires $k-1$ rounds to get any finite approximation factor, while SOCCER stops after a single round and finds the optimal clustering.

**Theorem 7.2.** *Let $k \in \mathbb{N}$. For any $n_0 \in \mathbb{N}$, there exists a dataset $X$ of size $n \geq n_0$, such that if $k$-means|| runs on $X$ for fewer than $k-1$ rounds, then it does not obtain a finite multiplicative approximation factor, while with probability at least $1 - \delta$, SOCCER stops after a single round and returns the optimal clustering.*

The experiments below demonstrate that also in practice, in many cases SOCCER requires few rounds.

# 8 Experiments

We report experiments on synthetic and real datasets. The code is provided at `https://github.com/selotape/distribu`
The experiments were performed on a single multi-core machine with a standard Intel processor, which ran the code of the coordinator and of all the machines. We could not run EIM11 (Ene et al., 2011) on these datasets, since, as explained in Section 5, in this algorithm the coordinator broadcasts a very large number of points to the machines. Since each machine is required to calculate the distance from each of its data points to the broadcasted set of points, this leads to a very large machine running time. For instance, for $k = 100$, $n = 10^7$, and $\epsilon = 0.1$, the coordinator broadcasts 72,000 points to the machines in each round, compared to about 200 points sent by SOCCER and $k$-means||. As a result, the machine running time of EIM11 is more than a hundred-fold larger, making this algorithm far from competitive in terms of machine run time, and impractical to run in our environment.

Both SOCCER and $k$-means|| output more than $k$ centers. The output $k$ clustering was calculated using the standard weighted $k$-means approach described in Section 2, using the $k$-means algorithm of python's `scikit-learn` (Pedregosa et al., 2011), which was also used as our centralized black-box $k$-means algorithm for the intermediate clustering calculations of the coordinator in SOCCER. To reduce variance, we fixed the sample sizes $P_1$ and $P_2$ to be exactly an $\alpha$ fraction of the current data. The parameter $l$ of $k$-means||, which determines the number of points to select in each round, was set to $2k$, as in Bahmani et al. (2012) and in the default setting of `MLLib` (Meng et al., 2016). We calculated $k$-means clusterings using each of the two algorithms, for several values of $k$, on both synthetic and real datasets. The properties of the tested datasets are listed in Table 1.

Table 1: Properties of datasets

| Dataset | # points | Dim. |
|---|---|---|
| $k$-Gaussian Mixture | 10M | 15 |
| Higgs | 11M | 28 |
| Census1990 | 2.45M | 68 |
| KDDCup1999 | 4.8M | 42 |
| BigCross | 11.6M | 57 |

For SOCCER, we set $\delta = 0.1$ in all the experiments, and tested several values of $\epsilon$. For $k$-means||, we tested stopping after each round between 1 and 5. Each experiment was repeated 10 times; we

9

Table 2: Some of the experiment results (See Appendix A.5 for full results). Top: Comparing `SOCCER` and $k$-`means||` when each is running a single round. Bottom: $k$-`means||` results for 2 and 5 rounds. The factors in parenthesis for $k$-`means||` results provide the ratio between the $k$-`means||` cost or time to the corresponding values of `SOCCER`.

| Dataset | $k$ | $\epsilon$ | $|P_1|$ | Cost | T (seconds) | Cost | T (seconds) |
|---|---|---|---|---|---|---|---|
| | | | | SOCCER, one round | | $k$-means\|\|, one round | |
| Gau | 25 | 0.05 | 11,316 | 150 | 0.37 | $168 \cdot 10^3$ (x6,340) | 0.05 (x0.14) |
| | 100 | 0.05 | 56,440 | 150 | 0.68 | $1,079 \cdot 10^3$ (x1,773) | 0.05 (x0.07) |
| Hig | 25 | 0.1 | 25,335 | $144 \cdot 10^6$ | 0.32 | $171 \cdot 10^6$ (x1.19) | 0.05 (x0.16) |
| | 100 | 0.05 | 56,440 | $122 \cdot 10^6$ | 0.48 | $137 \cdot 10^6$ (x1.12) | 0.06 (x0.12) |
| Cen | 25 | 0.1 | 22,018 | $188 \cdot 10^6$ | 0.09 | $418 \cdot 10^6$ (x2.22) | 0.05 (x0.56) |
| | 100 | 0.1 | 109,813 | $132 \cdot 10^6$ | 0.13 | $264 \cdot 10^6$ (x2) | 0.05 (x0.38) |
| KDD | 25 | 0.2 | 110,088 | $112 \cdot 10^{12}$ | 0.15 | $254 \cdot 10^{12}$ (x2.08) | 0.06 (x0.4) |
| | 100 | 0.2 | 549,037 | $743 \cdot 10^{10}$ | 0.26 | $5,175 \cdot 10^{10}$ (x6.97) | 0.06 (x0.23) |
| Big | 25 | 0.1 | 25,335 | $332 \cdot 10^{10}$ | 0.38 | $519 \cdot 10^{10}$ (x1.56) | 0.18 (x0.47) |
| | 100 | 0.1 | 126,354 | $152 \cdot 10^{10}$ | 0.53 | $241 \cdot 10^{10}$ (x1.86) | 0.18 (x0.34) |

| Dataset | $k$ | Cost | T (seconds) | Cost | T (seconds) |
|---|---|---|---|---|---|
| | | $k$-means\|\|, 2 rounds | | $k$-means\|\|, 5 rounds | |
| Gau | 25 | 37,350 (x246) | 0.33 (x0.89) | 164 (x1.1) | 1.98 (x5.35) |
| | 100 | 25,866 (x172) | 1.09 (x1.6) | 167 (x1.1) | 7.09 (x10.4) |
| Hig | 25 | $153 \cdot 10^6$ (x1.06) | 0.31 (x0.96) | $139 \cdot 10^6$ (x1.06) | 1.59 (x4.96) |
| | 100 | $125 \cdot 10^6$ (x1.06) | 0.85 (x1.77) | $115 \cdot 10^6$ (x0.94) | 5.62 (x11.7) |
| Cen | 25 | $218 \cdot 10^6$ (x1.15) | 0.15 (x1.66) | $185 \cdot 10^6$ (x0.98) | 0.6 (x6.66) |
| | 100 | $133 \cdot 10^6$ (x1) | 0.31 (x2.38) | $109 \cdot 10^6$ (x0.82) | 1.66 (x12.76) |
| KDD | 25 | $157 \cdot 10^{12}$ (x1.4) | 0.23 (x1.53) | $126 \cdot 10^{12}$ (x1.12) | 1.03 (x6.86) |
| | 100 | $649 \cdot 10^{10}$ (x0.87) | 0.54 (x2.07) | $795 \cdot 10^{10}$ (x1.07) | 3.05 (x11.73) |
| Big | 25 | $519 \cdot 10^{10}$ (x1.66) | 0.18 (x0.47) | $330 \cdot 10^{10}$ (x0.99) | 2.13 (x5.6) |
| | 100 | $169 \cdot 10^{10}$ (x1.11) | 1.09 (x2.06) | $150 \cdot 10^{10}$ (x0.99) | 6.17 (x11.64) |

report the average of each result. Standard deviations (reported in Appendix A.5) were usually smaller than 2% of the reported mean.

First, we generated for each tested $k$ a synthetic dataset drawn from a $k$-Gaussian mixture in $\mathbb{R}^{15}$. The mean of each Gaussian was randomly drawn from the unit cube in $\mathbb{R}^{15}$, and all Gaussian were all set to be spherical with isotropic variance $\sigma = 0.001$. The weight distribution of the Gaussians in the mixture was set according to the Zipf distribution, proportionally to $i^\gamma$, where $\gamma = 1.5$. Each dataset consisted of ten million points drawn from this distribution. We provide the code and seed for generating these datasets at `https://github.com/selotape/distributed_k_means`. We then tested the algorithms on four real-world datasets with millions of points, which were used in previous papers studying similar settings: `HIGGS`, `KDDCup1999` (Baldi et al., 2014) and `Census1990`, all from the UCI repository (Dua and Graff, 2017), and `Bigcross` (Ackermann et al., 2012).

Table 2 provides some of the results of running the algorithms on the each of the datasets. Results of `SOCCER` for all values of $\epsilon$ and for $k$-`means||` after all rounds between 1 and 5 are reported in full in Appendix A.6. In Table 2 (Top), we report the value of $\epsilon$ and the induced coordinator clustering size $|P_1|$ that resulted in `SOCCER` stopping after a single round, and provide

the obtained cost and the machine running time of SOCCER and of $k$-means|| after a single round. This provides a direct comparison with the same number of rounds. In Table 2 (Bottom), we report the results of $k$-means|| for the same experiments after two and five rounds, for comparison to SOCCER after a single round. The reported machine running time was calculated by taking the sum, over all rounds, of the maximal machine running time in each round based on 50 machines. The communication complexity of SOCCER per round is $2|P_1|$. The communication complexity of $k$-means|| per round is $l = 2k$. While for large $\epsilon$, the total communication is much larger in SOCCER, the average communication complexity per machine in SOCCER is smaller, since it is $2|P_1|$ divided by the number of machines.

For the $k$-Gaussian mixtures, the first two rows of Table 2 (Top) show that when the coordinator is allowed to cluster $|P_1| \approx 500 \cdot k$ points, SOCCER stops after a single round. In comparison, when stopping $k$-means|| after one round, its resulting clustering cost is three orders of magnitude larger than that of SOCCER. As can be seen in Table 2 (Bottom), even after five rounds, the cost obtained by $k$-means|| is still somewhat larger than the one obtained in one round by SOCCER, at which point the machine running time is also larger than that of SOCCER. For all the coordinator sizes that we tested (see Appendix A.6), the output cost of SOCCER for the Gaussian mixtures was almost identical (and approximately optimal) regardless of coordinator sizes, which only affected the number of rounds.

For the other datasets, it can be seen that the cost obtained by SOCCER after its single round is lower than that obtained by $k$-means|| after one round, and almost always also after two rounds. In addition, the machine running time of SOCCER after one round is usually significantly smaller than that of $k$-means|| after running the number of rounds necessary to obtain a comparable cost.

In all of our experiments, SOCCER stopped after a smaller number of rounds than the worst-case guarantee of $1/\epsilon - 1$. In particular, Table 3 reports experiments in which $\epsilon = 0.01$ and so the coordinator size was very small (see Appendix A.6 for other values of $\epsilon$). In this case, the worst-case number of rounds is 99, while the true number of rounds was usually between 2 and 4. Even with this small coordinator size, the number of rounds required by $k$-means|| to obtain a comparable cost was usually much larger, as can be seen by comparing to the two rightmost columns in Table 3. The machine running time in $k$-means|| was also usually significantly larger. Note that unlike $k$-means||, in which each round requires the same running time, in SOCCER each additional round is considerably faster, due to the removal of points. Regarding the dependence of the cost upper bound on the total number of rounds of SOCCER in Theorem 4.1, it can be seen in Appendix A.6 that in practice the cost is similar for the same dataset for different coordinator sizes, although they each lead to a different total number of rounds.

We conclude that overall, if the coordinator is allowed to calculate a clustering for a moderate number of points, SOCCER usually stops after a small number of rounds, and obtains a comparable or better cost than $k$-means||, even if the latter runs for a larger number of rounds. In addition, SOCCER requires significantly less machine running time to achieve a comparable cost.

Our run time comparison focuses on machine running times, showing that in this respect SOCCER is considerably faster. Our premise is that the coordinator is significantly stronger, and its computation time is not a bottleneck. However, one may still be interested in reducing the coordinator running time as well. To speed up the coordinator, a faster clustering implementation can be used. However, existing fast implementations are typically less successful on more difficult datasets. We demonstrate this approach by replacing the black-box k-means implementation used for the experiments above with the faster `MiniBatchKMeans` implementation from `scikit-learn`. The results, reported in Appendix A.7, show that in almost all experiments, SOCCER obtains a similar cost to

Table 3: Results of experiments with $\epsilon = 0.01$. 'R' of SOCCER gives the number of rounds it required. $k$-means|| was run until a cost that is up to 2% from that of SOCCER.

| Data | $k$ | $|P_1|$ | SOCCER, $\epsilon = 0.01$ | | | $k$-means|| | |
|------|-----|---------|------|------|------|------|------|
| | | | R | Cost | T | R | T |
| Gau | 25 | 6,000 | 3 | 150 | 0.72 | 15 | 12.3 |
| | 100 | 30,000 | 2 | 150 | 0.95 | 15 | 47 |
| Hig | 25 | 6,000 | 3 | $134 \cdot 10^6$ | 0.6 | 8 | 3.8 |
| | 100 | 30,000 | 2 | $120 \cdot 10^6$ | 0.68 | 3 | 2.1 |
| Cen | 25 | 6,000 | 4 | $176 \cdot 10^6$ | 0.2 | 8 | 1.3 |
| | 100 | 30,000 | 3 | $110 \cdot 10^6$ | 0.3 | 5 | 1.7 |
| KDD | 25 | 6,000 | 11 | $114 \cdot 10^{12}$ | 1 | 10 | 3.3 |
| | 100 | 30,000 | 7 | $597 \cdot 10^{10}$ | 1.1 | 10 | 10.9 |
| Big | 25 | 6,000 | 3 | $319 \cdot 10^{10}$ | 0.9 | 8 | 5.1 |
| | 100 | 31,000 | 2-3 | $154 \cdot 10^{10}$ | 0.93 | 4 | 3.9 |

$k$-means|| with a comparable total running time and fewer rounds. A notable exception is the KDDCup1999 dataset. In this dataset, MiniBatchKMeans fails to find a clustering with a reasonable cost, even when running on the entire set of points. We believe this is because this dataset includes many outliers ?, which are not well handled by this implementation. This highlights the importance of using a black box that is suitable for the task at hand.

# 9 Conclusion

In this work, we presented a new distributed $k$-means clustering algorithm that can require as little as one or two communication rounds, and stops on its own without having to specify the number of rounds as a parameter. Given a restriction on the maximal number of points that can be clustered by the coordinator using a centralized $k$-means algorithm, our algorithm obtains a constant approximation factor, as well as a constant upper bound on the number of rounds. Our experiments demonstrate its effectiveness on various datasets, where it usually obtains a smaller cost than $k$-means|| using fewer rounds. We believe that the techniques used in SOCCER can further be used to support robustness against outliers and machine failures, and we intend to study these challenges in future work.

**References**

Marcel R Ackermann, Marcus Märtens, Christoph Raupach, Kamil Swierkot, Christiane Lammersen, and Christian Sohler. StreamKM++: A clustering algorithm for data streams. *Journal of Experimental Algorithmics (JEA)*, 17:2–4, 2012.

Sara Ahmadian, Ashkan Norouzi-Fard, Ola Svensson, and Justin Ward. Better guarantees for k-

means and euclidean k-median by primal-dual algorithms. *SIAM Journal on Computing*, 49(4): FOCS17–97, 2019.

Nir Ailon, Ragesh Jaiswal, and Claire Monteleoni. Streaming k-means approximation. In *Advances in neural information processing systems*, pages 10–18, 2009.

Olivier Bachem, Mario Lucic, and Andreas Krause. Distributed and provably good seedings for k-means in constant rounds. In *International Conference on Machine Learning*, pages 292–300. PMLR, 2017a.

Olivier Bachem, Mario Lucic, and Andreas Krause. Practical coreset constructions for machine learning. *arXiv preprint arXiv:1703.06476*, 2017b.

Bahman Bahmani, Benjamin Moseley, Andrea Vattani, Ravi Kumar, and Sergei Vassilvitskii. Scalable k-means++. *arXiv preprint arXiv:1203.6402*, 2012.

Maria Florina Balcan, Steven Ehrlich, and Yingyu Liang. Distributed k-means and k-median clustering on general topologies. *arXiv preprint arXiv:1306.0604*, 2013.

Pierre Baldi, Peter Sadowski, and Daniel Whiteson. Searching for exotic particles in high-energy physics with deep learning. *Nature communications*, 5(1):1–9, 2014.

Aditya Bhaskara and Maheshakya Wijewardena. Distributed clustering via lsh based data partitioning. In *International Conference on Machine Learning*, pages 570–579. PMLR, 2018.

Jiecao Chen, He Sun, David Woodruff, and Qin Zhang. Communication-optimal distributed clustering. *Advances in Neural Information Processing Systems*, 29:3727–3735, 2016.

Jiecao Chen, Erfan Sadeqi Azer, and Qin Zhang. A practical algorithm for distributed clustering and outlier detection. *arXiv preprint arXiv:1805.09495*, 2018.

Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL http://archive.ics.uci.edu/ml.

Alina Ene, Sungjin Im, and Benjamin Moseley. Fast clustering using mapreduce. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 681–689, 2011.

Dan Feldman, Melanie Schmidt, and Christian Sohler. Turning big data into tiny data: Constant-size coresets for k-means, pca, and projective clustering. *SIAM Journal on Computing*, 49(3): 601–657, 2020.

Sudipto Guha, Adam Meyerson, Nina Mishra, Rajeev Motwani, and Liadan O'Callaghan. Clustering data streams: Theory and practice. *IEEE transactions on knowledge and data engineering*, 15(3):515–528, 2003.

Sudipto Guha, Yi Li, and Qin Zhang. Distributed partial clustering. *ACM Transactions on Parallel Computing (TOPC)*, 6(3):1–20, 2019.

Xiangyu Guo and Shi Li. Distributed *k*-clustering for data with heavy noise. *arXiv preprint arXiv:1810.07852*, 2018.

Tom Hess, Michal Moshkovitz, and Sivan Sabato. A constant approximation algorithm for sequential no-substitution k-median clustering under a random arrival order. *arXiv preprint arXiv:2102.04050*, 2021.

Ravi Kumar, Benjamin Moseley, Sergei Vassilvitskii, and Andrea Vattani. Fast greedy algorithms in mapreduce and streaming. *ACM Transactions on Parallel Computing (TOPC)*, 2(3):1–22, 2015.

Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338, 2000.

Xiangrui Meng, Joseph Bradley, Burak Yavuz, Evan Sparks, Shivaram Venkataraman, Davies Liu, Jeremy Freeman, DB Tsai, Manish Amde, Sean Owen, et al. Mllib: Machine learning in apache spark. *The Journal of Machine Learning Research*, 17(1):1235–1241, 2016.

Rajeev Motwani and Prabhakar Raghavan. Randomized algorithms. *ACM Computing Surveys (CSUR)*, 28(1):33–37, 1996.

F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

George S Sebestyen. *Decision-making processes in pattern recognition (ACM monograph series)*. Macmillan Publishing Co., Inc., 1962.

Konstantin Voevodski. Large scale k-median clustering for stable clustering instances. In *International Conference on Artificial Intelligence and Statistics*, pages 2890–2898. PMLR, 2021.

# A  Appendix

## A.1  Proof of Lemma 6.1

In this section, we give the proof of Lemma 6.1. We first give an auxiliary lemma. This lemma is a variation on results that were proved in Hess et al. (2021) (henceforth abbreviated to HMS21), where the latter have significantly larger constants. An additional difference is that HMS21 proved the results for $k$-median. The adaptation to $k$-means is straightforward, but affects some constants.

Denote a $k$-means solution from $X$ which is optimal for some subset $Y \subseteq X$ by the notation $\mathrm{OPT}_Y^X := \min_{T \subseteq X, |T| \le k} \mathsf{cost}(Y, T)$. Consider the optimal clusters induced by $\mathrm{OPT}_Y^X$ on $Y$. HMS21 defines small optimal clusters as those optimal clusters which include at most $150 \log(32k/(\delta))/\alpha$ points. In order to obtain guarantees with smaller constants, we use a variant of this definition. Recall the notation $d_k := 6.5 \log(1.1k/(\delta\epsilon))$, $k_+ := k + 9 \log(1.1k/(\delta\epsilon))$. Denote $d_k' := 6.5 \log(1.1k/\delta), k_+' := k + 9 \log(1.1k/\delta)$, which are equal to $d_k, k_+$ with $\epsilon = 1$. Let $F_\alpha(Y)$ be the $d_k'/\alpha$ points in $Y$ that are furthest from $\mathrm{OPT}_Y^X$. Define small optimal clusters to be those that of size at most $d_k'/\alpha$ after removing the points in $F_\alpha(Y)$. Any larger cluster is called a *large optimal cluster*. Denote by $C_{\mathrm{small}}^\alpha(Y)$ the set of points in $Y$ that belong to small optimal clusters in $Y$, and its complement by $C_{\mathrm{large}}^\alpha(Y) := Y \setminus C_{\mathrm{small}}^\alpha$. The following lemma provides results that are adaptations of results from HMS21, where latter have larger constants and hold for the original definition of small optimal clusters.

**Lemma A.1** (Adaptation of results from HMS21). *Let $Y \subseteq X$. Let $\alpha, \delta \in (0, 1)$ and set $k_+', d_k'$ as defined in Alg. 1. Let $P_1, P_2 \subseteq Y$ be two independent samples of size $\alpha|Y|$, selected uniformly at random from $Y$. Let $\mathcal{A}$ be a $\beta$-approximation $k$-means algorithm, and define $T := \mathcal{A}(P_1, k_+')$ . With probability at least $1 - \delta$,*

1.  $\mathsf{cost}(C_{\mathrm{large}}^\alpha(Y) \setminus F_\alpha(Y), T) \le (36\beta + 20)\mathsf{cost}(Y, \mathrm{OPT})$;

2.  $\psi := \frac{2}{3\alpha}\mathsf{cost}_{\frac{3}{2}(k+1)d_k'}(P_2, T) \le \mathsf{cost}_{(k+1)d_k'/\alpha}(Y, T)$;

3.  $|\{x \in Y \mid \rho(x, T)^2 > \psi\alpha/(kd_k')\}| \le 5.5kd_k'/\alpha$.

*Proof Sketch.* The lemma is derived by adapting results from HMS21 to our setting. The adaptation is consists of following the same proofs with minor technical differences; we give a sketch of the differences below.

The three parts of the lemma are derived by adapting lemmas 5.7, 5.8, and 5.9 of HMS21 to our setting. The original claim in Lemma 5.7 is proved for $\mathrm{OPT}_Y^X := \min_{T \subseteq X, |T| \le k} \mathsf{cost}(Y, T)$, however it is easy to see that $\mathsf{cost}(Y, \mathrm{OPT}_Y^X) \le \mathsf{cost}(Y, \mathrm{OPT})$. In addition, the original claim does not subtract $F_\alpha(Y)$ on the LHS. This subtraction allows us to get improved final constants. Lemma 5.9 gives the claim in part 3 for $Y \setminus (P_1 \cup P_2)$, while in our case it holds for $Y$ (with a different constant). This is because in our case, $P_1$ and $P_2$ are independent samples, while in HMS21 they are non-overlapping.

The main differences between the original lemmas and the version we give here are in the definition of large clusters and in the resulting constants. In particular, HMS21 provided guarantees for $d_k' = 150 \log(\frac{32k}{\delta})$ and $k_+' = k + 38 \log(\frac{32k}{\delta})$. In the current work, we define $d_k' = 6.5 \log(\frac{1.1k}{\delta})$ and $k_+' = k + 9 \log(\frac{1.1k}{\delta})$. In addition, as described above, our definition of small optimal clusters ignores the points in $F_\alpha(Y)$. In addition to the new definition of small clusters and a tightening of the constants in the analysis, the improved constants are also due to the fact that unlike HMS21,

15

we require fewer events to hold. For instance, we do not require the optimal points to be outside of $P_1$ and $P_2$. This allows reducing the factor in the log in the definition of $d_k'$.

The constant factor is further improved by using the tighter version of the multiplicative Chernoff bound (Motwani and Raghavan, 1996) to tighten the constants in Lemma 5.4 of HMS21. Reducing the constants in $d_k'$ leads to a reduction in other constants as well, including those in $k_+'$. An additional improvement in constants stems by assuming that $n$ is not too small, which allows avoiding certain edge cases. In particular, this allows improving the constants in the guarantees provided in HMS21 for linear bin divisions, and these affect the final result.

The approximation factor of Lemma 5.7 in HMS21 is $18\beta + 10$. This factor is reduced to $9\beta + 5.5$ for $k$-medians using the techniques above. For $k$-means, the triangle inequality used in several places in that proof needs to be replaced by the weak triangle inequality, leading to a final approximation factor of $36\beta + 20$. $\square$

The following corollary is immediate, by applying the lemma above to the intermediate calculations in `SOCCER`, and replacing $\delta, d_k, k_+$ by $\delta\epsilon, d_k', k_+'$, respectively. For simplicity, we take the sizes of $P_1$ and $P_2$ in `SOCCER` to be exactly an $\alpha$ fraction of $V_i$. For the independent sampling mechanism of $\{P_j^l\}$ used in Alg. 1, this holds in expectation, and with a high probability for large data sizes, up to a negligible correction. It can also be enforced exactly and for all dataset sizes, by letting the coordinator set the number of sample points that each machine should send, based on a draw from the relevant multinomial distribution. However, since this would have a negligible effect in most cases, and makes the algorithm unnecessarily more complicated, we chose to present the simpler mechanism in Alg. 1.

**Corollary A.2.** *Assume that* `SOCCER` *runs with the parameters given in Theorem 4.1. Let* $\alpha_i, C_{\mathrm{iter}}^i, v_i$ *be the respective values of* $\alpha, C_{\mathrm{iter}}, \psi, v$ *calculated at iteration* $i$ *of* `SOCCER`. *Let* $\psi_i = v_i k d_k / \alpha_i$. *Let* $V_i$ *be the remaining dataset at the beginning of round* $i$ *of* `SOCCER`.

*With probability at least* $1 - \delta\epsilon$,

- $\mathsf{cost}(C_{\mathrm{large}}^\alpha(V_i) \setminus F_{\alpha_i}(V_i), C_{\mathrm{iter}}^i) \leq (36\beta + 20)\mathsf{cost}(V_i, \mathrm{OPT})$;

- $\psi_i \leq \mathsf{cost}_{(k+1)d_k/\alpha_i}(V_i, C_{\mathrm{iter}}^i)$;

- $|V_{i+1}| \leq 5.5 k d_k / \alpha_i$.

We now use this corollary to prove Lemma 6.1.

*Proof of Lemma 6.1.* To prove the first part of the lemma, we separately bound the cost of $R_i \cap (C_{\mathrm{large}}^\alpha(V_i) \setminus F_{\alpha_i}(V_i))$ and $R_i \cap (C_{\mathrm{small}}^\alpha(V_i) \cup F_{\alpha_i}(V_i))$ with respect to $C_{\mathrm{iter}}^i$. For the first part, the bound follows from part 1 of Cor. A.2, since

$$\mathsf{cost}(R_i \cap (C_{\mathrm{large}}^\alpha(V_i) \setminus F_{\alpha_i}(V_i)), C_{\mathrm{iter}}^i)$$
$$\leq \mathsf{cost}(C_{\mathrm{large}}^\alpha(V_i) \setminus F_{\alpha_i}(V_i), C_{\mathrm{iter}}^i)$$
$$\leq (36\beta + 20)\mathsf{cost}(V_i, \mathrm{OPT}). \tag{1}$$

Next, we consider the second part. Note that by the definition of small clusters, $|C_{\mathrm{small}}^\alpha(V_i) \cup F_{\alpha_i}(V_i)| \leq (k+1)d_k/\alpha_i$. Hence, we get

$$\mathsf{cost}_{(k+1)d_k/\alpha_i}(V_i, C_{\mathrm{iter}}^i)$$
$$\leq \mathsf{cost}(V_i \setminus (C_{\mathrm{small}}^\alpha(V_i) \cup F_{\alpha_i}(V_i)), C_{\mathrm{iter}}^i)$$
$$= \mathsf{cost}(C_{\mathrm{large}}^\alpha(V_i) \setminus F_{\alpha_i}(V_i), C_{\mathrm{iter}}^i).$$

16

Hence, by combing the above equation with the first and second parts of Cor. A.2, we get that

$$\psi_i \leq (36\beta + 20)\mathsf{cost}(V_i, \mathrm{OPT}).$$

Note that by the definition of $R_i$ and the equation above,

$$\forall x \in R_i, \quad \rho(x, C_{\mathrm{iter}}^i) \leq v^i = \frac{\psi_i}{kd_k/\alpha}$$
$$\leq \frac{(36\beta + 20)\mathsf{cost}(V_i, \mathrm{OPT})}{kd_k/\alpha}.$$

Hence,

$$\mathsf{cost}(R_i \cap (C_{\mathrm{small}}^\alpha(V_i) \cup F_{\alpha_i}(V_i)), C_{\mathrm{iter}}^i) \tag{2}$$
$$\leq |R_i \cap (C_{\mathrm{small}}^\alpha(V_i) \cup F_{\alpha_i}(V_i))| \frac{(36\beta + 20)\mathsf{cost}(V_i, \mathrm{OPT})}{kd_k/\alpha}$$
$$\leq \frac{k+1}{k}(36\beta + 20)\mathsf{cost}(V_i, \mathrm{OPT})$$
$$\leq (44\beta + 24)\mathsf{cost}(V_i, \mathrm{OPT}), \tag{3}$$

Where the second inequality follows since,

$$|(C_{\mathrm{small}}^\alpha(V_i) \cup F_{\alpha_i}(V_i))| < (k+1)d_k/\alpha,$$

and the third inequality follows since $k \geq 5$. By combining Eq. (1) and Eq. (3), we get that

$$\mathsf{cost}(R_i, C_{\mathrm{iter}}^i) =$$
$$\mathsf{cost}(R_i \cap (C_{\mathrm{small}}^\alpha(V_i) \cup F_{\alpha_i}(V_i)), C_{\mathrm{iter}}^i)$$
$$+ \mathsf{cost}(R_i \cap (C_{\mathrm{large}}^\alpha(V_i) \setminus F_{\alpha_i}(V_i)), C_{\mathrm{iter}}^i)$$
$$\leq (80\beta + 44)\mathsf{cost}(V_i, \mathrm{OPT}).$$

This completes the proof of the first part of the lemma. The second part of the lemma is the same as the third part of Cor. A.2. □

## A.2 Proof of Theorem 4.1

*Proof of Theorem 4.1.* By a union bound, the event in Lemma 6.1 holds in all of the first $\min(I, 1/\epsilon)$ rounds with probability at least $1 - \delta$.

To prove the first part of the theorem, we show that under this joint event, SOCCER stops after at most $1/\epsilon$ rounds. By the definition in line 3, $\alpha_i = \eta(\epsilon)/|V_i|$. Also, $\eta(\epsilon) = 36kn^\epsilon \log(\frac{1.1k}{\delta\epsilon})$ and $d_k = 6.5 \log(\frac{1.1k}{\delta\epsilon})$. Hence, by the second part of Lemma 6.1,

$$|V_{i+1}| \leq 5.5kd_k/\alpha_i = 5.5|V_i|kd_k/\eta(\epsilon) < |V_i|/n^\epsilon.$$

Since $V_1 = n$, it follows by induction that $|V_{i+1}| \leq n^{1-i\epsilon}$. Recall that the stopping condition of the main loop of SOCCER is $|V_{i+1}| \leq \eta(\epsilon)$. Clearly, we have $n^{1-i\epsilon} \leq \eta(\epsilon)$ once $i \geq (1 - \frac{\log(36k\log(\frac{1.1k}{\delta\epsilon}))}{\log n})\frac{1}{\epsilon} -$

17

1. Therefore, the total number of communication rounds is at most $(1 - \frac{\log(36k \log(\frac{1.1k}{\delta\epsilon}))}{\log n})\frac{1}{\epsilon} < 1/\epsilon - 1$. This proves the first part of the theorem.

Next, we prove the cost approximation bound (the third part of the theorem).

Since $\{R_i\}$ and $V_I$ are a partition of $X$ and $C_{\text{iter}}^i \subseteq C_{\text{out}}$ for all iterations $i$, we have

$$\mathsf{cost}(X, C_{\text{out}}) \leq \sum_{i \in [I-1]} \mathsf{cost}(R_i, C_{\text{iter}}^i) + \mathsf{cost}(V_I, C_{\text{iter}}^I),$$

where $C_{\text{iter}}^I$ is the result of the $k$-means clustering performed in line 17 of Alg. 1.

Using the first part of Lemma 6.1 and recalling that $C_{\text{iter}}^I$ is a $\beta$ approximation $k$-means solution for $V_I$, we get

$$\mathsf{cost}(X, C_{\text{out}}) \leq \sum_{i \in [I]} (80\beta + 44) \cdot \mathsf{cost}(V_i, \text{OPT}) \leq \sum_{i \in [I]} (80\beta + 44) \cdot \mathsf{cost}(X, \text{OPT})$$
$$= I \cdot (80\beta + 44) \cdot \mathsf{cost}(X, \text{OPT}).$$

This completes the proof of the third part of the theorem. The second, fourth, and fifth parts follow directly from the definition of SOCCER. This completes the proof. □

## A.3   Proof of Theorem 7.1

*Proof of Theorem 7.1.* Consider a $k$-Gaussian mixture in dimension $d$. Suppose that the Gaussians are all spherical with covariance matrix $\sigma^2 I$. For $Z \sim N(\mu, \sigma^2 I)$, it is known (see, e.g. Laurent and Massart, 2000) that for $\gamma > 0$,

$$\mathbb{P}[||Z - \mu||_2^2 \leq \sigma^2 d \cdot (1 + 2\sqrt{\frac{\log(\frac{1}{\gamma})}{d}} + \frac{2\log(\frac{1}{\gamma})}{d})] \geq 1 - \gamma. \tag{4}$$

In other words, for large values of $d$, almost all the points drawn from each Gaussian are about $\sigma\sqrt{d}$-far from the mean of the Gaussian. Thus, with high probability, the optimal $k$-clustering cost for a dataset of size $n$ is $\Theta(n\sigma^2 d)$.

Suppose that SOCCER runs on a dataset drawn from this $k$-mixture. In the first iteration, SOCCER calculates a $k_+$-clustering over a random sample of points from the dataset, using the $\beta$ approximation algorithm $\mathcal{A}$. Since $k_+$ does not depend on $d$, it is easy to see that for a large enough $d$, the average distance of the dataset points from any $k_+$ centers cannot be significantly smaller than the average distance of these points from their Gaussian centers. Therefore, the cost of the calculated $k_+$ clustering on the dataset is $\Theta(n\sigma^2 d)$.

We now show that SOCCER stops after one round. First, we consider the value of $\mathsf{cost}_{\frac{3}{2}(k+1)d_k}(P_2, C_{\text{iter}})$, which is used to calculate $v$ in line 9 of Alg. 1. This is the cost of $C_{\text{iter}}$ on $P_2$ after removing the $\frac{3}{2}(k+1)d_k$ points that are furthest from $C_{\text{iter}}$. Note that $|P_2| = \alpha n$. Therefore, the fraction of points from $P_2$ disregarded in the calculation of the truncated cost is $kd_k/(\alpha n) = 6.5k \log(\frac{1.1k}{\delta})/(\alpha n) = 6.5k \log(\frac{1.1k}{\delta})/\eta(\epsilon) = O(n^{-\epsilon})$. Therefore, this fraction goes to zero for large $n$. It follows that $\mathsf{cost}_{\frac{3}{2}(k+1)d_k}(P_2, C_{\text{iter}}) = |P_2|\Theta(\sigma^2 d) = \Theta(\alpha n \sigma^2 d)$. Since $\alpha = \eta(\epsilon)/n = \Theta(n^{\epsilon-1})$, we get $\mathsf{cost}_{\frac{3}{2}(k+1)d_k}(P_2, C_{\text{iter}}) = \Theta(n^\epsilon \sigma^2 d)$. The threshold is thus $v = \Theta(\mathsf{cost}_{\frac{3}{2}(k+1)d_k}(P_2, C_{\text{iter}})) = \Theta(n^\epsilon \sigma^2 d)$.

We now show that SOCCER stops after a single round, by showing that with high probability, all the points in $X$ are closer to $C_{\text{iter}}$ than $v$. From Eq. (4) with $\gamma = \log(1/\delta)/n$, we get that with a probability at least $1 - \log(1/\delta)/n$, a point drawn from a Gaussian has a square distance of $O(\sigma^2(d + \log(n/\delta)))$ to the center of the Gaussian. Hence, with probability at least $1 - \delta$, this holds for all the points in the dataset. Clearly, this implies that the centers $C_{\text{iter}}$ selected by $\mathcal{A}$ must include centers with a square distance of $O(\sigma^2(d + \log(n/\delta)))$ from the Gaussian mean, otherwise the $\beta$-approximation guarantee would not hold. By the assumption of the theorem, $\epsilon > \log\log(n/\delta)/\log n$. Therefore, $n^\epsilon > \log(n/\delta)$. It follows that $v = \Omega(\sigma^2 d \log(n/\delta))$. Thus, for a large enough $d$, $v$ is larger than the distance of all points from $C_{\text{iter}}$. As a result, all the dataset points are removed in the first round of SOCCER, and the algorithm completes after one round. By Theorem 4.1, this implies also that the cost of the output clustering is a constant approximation of the optimal cost. □

## A.4 Proof of Theorem 7.2

*Proof of Theorem 7.2.* The proof is based on an example of Bachem et al. (2017a) of a hard instance for $k$-means||. Bachem et al. (2017a, Theorem 2) describes a dataset of size $2k - 2$ such that for any value of the $k$-means|| parameter $l$, $k$-means|| requires at least $k - 1$ rounds to obtain a constant approximation. The dataset in the example includes $k$ distinct points $\{x_i\}_{i\in[k]}$, where $x_1$ has $k - 1$ copies in the dataset and each of $x_2, \ldots, x_k$ appear a single time in the dataset.

To prove the claim in the theorem, we construct a dataset of size $n > n_0$ based on this example, by duplicating the above dataset $z = \lceil n_0/(2k-2) \rceil$ times. The number of rounds required by $k$-means|| remains the same, as can be verified by following the proof of Theorem 2 in Bachem et al. (2017a).

In contrast, we now show that SOCCER stops after one round on this dataset. Consider the sub-sample $P_1$, which is calculated in the first round of SOCCER. For any $i \in [k]$,

$$\begin{aligned}
\mathbb{P}[x_i \notin P_1] &\leq (1-\alpha)^z \leq \exp(-\alpha \cdot z) \\
&\leq \exp\left(-\frac{\eta(\epsilon)}{2z(k-1)} \cdot z\right) \leq \exp\left(-\frac{\eta(\epsilon)}{2(k-1)}\right) \\
&\leq \delta/k.
\end{aligned}$$

The last inequality follows since $\eta(\epsilon) = 36kn^\epsilon \log(\frac{k}{\delta})$. Therefore, with probability at least $1 - \delta$, an instance of each of $\{x_i\}_{i\in[k]}$ is found in $P_1$. Hence, the optimal clustering for $P_1$ includes all the distinct points from $X$, and has a cost of zero. As a result, also $C_{\text{iter}}$ must have a cost of zero and so it also includes all the distinct points from $X$, leading to the removal of all the dataset points from each of the machines in line 12. Therefore, SOCCER stops after one round and returns an optimal clustering. □

## A.5   Full experiment results

In this section, we provide the full results of all the experiments described in Section 8. Each table reports experiments on one of the datasets in Table 1. The results are divided to two subsections. Appendix A.6 shows the results of `SOCCER` and $k$-`means||` when the standard Kmeans implementation is used as a black-box algorithm for `SOCCER`, and Appendix A.7 shows the results when `MiniBatchKMeans` is used as the black box.

## A.6   Results for standard Kmeans as black-box for `SOCCER`

The results for `SOCCER` and $k$-`means||`, when standard Kmeans algorithm used as black-box for `SOCCER` are provided below in Table 4, Table 5, Table 6, Table 7, and Table 8.

Table 4: `k-GaussiansMixture` dataset experiments with Standard KMeans as black-box. 'T' stands for time in seconds.

| $k$ | ALG | $\epsilon$ | $P_1$ | Output size | Rounds | Cost | T (machine) | T (Total) |
|---|---|---|---|---|---|---|---|---|
| | | 0.2 | 126,978 | 90 | 1 | 150.1±0 | 0.32±0.07 | 6.56±0.22 |
| | SOCCER | 0.1 | 25,335 | 96 | 1 | 150.2±0 | 0.44±0.08 | 2.51±0.12 |
| | | 0.05 | 11,316 | 127±1 | 1 | 150.3±0 | 0.37±0.09 | 1.75±0.12 |
| | | 0.01 | 5,939 | 348 | 3 | 150.1±0 | 0.73±0.07 | 4.11±0.21 |
| 25 | | - | - | 51 | 1 | 1,688,270.3±951992.1 | 0.05±0 | 0.15±0.03 |
| | | - | - | 101 | 2 | 37,530.5±46409 | 0.33±0.01 | 0.43±0.02 |
| | $k$-means\|\| | - | - | 151 | 3 | 196.9±18.6 | 0.76±0.03 | 0.87±0.03 |
| | | - | - | 201 | 4 | 171.2±4.7 | 1.32±0.06 | 1.42±0.06 |
| | | - | - | 251 | 5 | 164.4±2.1 | 1.98±0.07 | 2.1±0.07 |
| | | 0.2 | 285,296 | 121 | 1 | 150.1±0 | 0.39±0.08 | 16.39±0.26 |
| | SOCCER | 0.1 | 56,924 | 127 | 1 | 150.2±0 | 0.48±0.08 | 5.03±0.2 |
| | | 0.05 | 25,427 | 137±1 | 1 | 150.3±0 | 0.57±0.08 | 3.43±0.12 |
| | | 0.01 | 13,344 | 346 | 2 | 150.2±0 | 0.79±0.09 | 5.02±0.15 |
| 50 | | - | - | 101 | 1 | 1,283,640.5±558248.2 | 0.05±0 | 0.24±0.04 |
| | | - | - | 201 | 2 | 13,399.3±11108.8 | 0.62±0.03 | 0.81±0.04 |
| | $k$-means\|\| | - | - | 301 | 3 | 211.5±6.8 | 1.45±0.04 | 1.66±0.04 |
| | | - | - | 401 | 4 | 174.9±2.3 | 2.52±0.05 | 2.71±0.07 |
| | | - | - | 501 | 5 | 166±1 | 3.83±0.08 | 4.04±0.12 |
| | | 0.2 | 633,271 | 177 | 1 | 150.1±0 | 0.53±0.05 | 73.11±0.73 |
| | SOCCER | 0.1 | 126,354 | 183 | 1 | 150.1±0 | 0.67±0.1 | 13.77±0.3 |
| | | 0.05 | 56,440 | 212±41 | 1 | 150.3±0 | 0.68±0.11 | 8.22±0.35 |
| | | 0.01 | 29,620 | 428±42 | 2 | 150.2±0 | 0.95±0.13 | 10.78±0.23 |
| 100 | | - | - | 201 | 1 | 1,079,458.8±266,814.6 | 0.05±0.01 | 0.45±0.13 |
| | | - | - | 401 | 2 | 25,866.5±16072.3 | 1.09±0.03 | 1.51±0.14 |
| | $k$-means\|\| | - | - | 601 | 3 | 226.9±61.2 | 2.67±0.06 | 3.09±0.1 |
| | | - | - | 801 | 4 | 176.6±3.1 | 4.75±0.05 | 5.21±0.08 |
| | | - | - | 1001 | 5 | 167.2±1.7 | 7.09±0.1 | 7.52±0.11 |
| | | 0.1 | 277,721 | 297±18 | 1 | 150.1±0 | 0.87±0.09 | 42.15±0.26 |
| | SOCCER | 0.05 | 124,053 | 371±81 | 1 | 150.3±0 | 0.93±0.12 | 21.97±0.55 |
| | | 0.01 | 65,104 | 648±61 | 2 | 150.2±0 | 1.26±0.1 | 27.1±0.52 |
| 200 | | - | - | 401 | 1 | 1,104,954±201,686.7 | 0.05±0.01 | 0.82±0.06 |
| | | - | - | 801 | 2 | 26,593.9±9,916.1 | 2.08±0.04 | 3±0.14 |
| | $k$-means\|\| | - | - | 1201 | 3 | 218.8±14.9 | 4.97±0.09 | 5.89±0.1 |
| | | - | - | 1601 | 4 | 175.7±1.5 | 8.72±0.1 | 9.74±0.17 |
| | | - | - | 2001 | 5 | 167±1.8 | 13.16±0.12 | 14.18±0.17 |

Table 5: `Higgs` dataset experiments with Standard KMeans as black-box. 'T' stands for time in seconds.

| $k$ | ALG | $\epsilon$ | $P_1$ | Output size | Rounds | Cost ($\cdot 10^6$) | T (Machine) | T (Total) |
|---|---|---|---|---|---|---|---|---|
| | | 0.2 | 126,978 | 92 | 1 | 129±0.38 | 0.3±0.02 | 112.01±4.91 |
| | SOCCER | 0.1 | 25,335 | 121 | 1 | 144±2.76 | 0.32±0.05 | 14.59±1.39 |
| | | 0.05 | 11,316 | 204 | 2 | 144±1.53 | 0.31±0.03 | 9.23±0.33 |
| | | 0.01 | 5,939 | 348 | 3 | 134±1.09 | 0.63±0.08 | 7.88±0.2 |
| 25 | | - | - | 51 | 1 | 171±3.99 | 0.05±0 | 0.16±0.02 |
| | | - | - | 101 | 2 | 153±1.47 | 0.31±0.02 | 0.43±0.03 |
| | $k$-means\|\| | - | - | 151 | 3 | 148±1.41 | 0.68±0.09 | 0.81±0.09 |
| | | - | - | 201 | 4 | 143±0.98 | 1.06±0.05 | 1.19±0.06 |
| | | - | - | 251 | 5 | 139±0.58 | 1.59±0.04 | 1.72±0.03 |
| | | 0.2 | 285,296 | 122.8±0.4 | 1 | 117±0.19 | 0.35±0.02 | 318.99±5.72 |
| | SOCCER | 0.1 | 56,924 | 177 | 1 | 134±1.75 | 0.38±0.04 | 54.53±3.34 |
| | | 0.05 | 25,427 | 183 | 1 | 128±0.79 | 0.39±0.05 | 20.6±1.46 |
| | | 0.01 | 13,344 | 346 | 2 | 124±0.51 | 0.56±0.04 | 17.16±0.76 |
| 50 | | - | - | 101 | 1 | 153±1.66 | 0.05±0 | 0.25±0.04 |
| | | - | - | 201 | 2 | 139±1.14 | 0.5±0.03 | 0.72±0.05 |
| | $k$-means\|\| | - | - | 301 | 3 | 133±0.63 | 1.13±0.05 | 1.38±0.08 |
| | | - | - | 401 | 4 | 129±0.65 | 1.95±0.03 | 2.2±0.05 |
| | | - | - | 501 | 5 | 127±0.45 | 2.96±0.1 | 3.24±0.11 |
| | | 0.2 | 633,272 | 178±1 | 1 | 106±0.09 | 0.45±0.03 | 908.56±9.99 |
| | SOCCER | 0.1 | 126,354 | 283 | 1 | 131±1.86 | 0.44±0.02 | 191.25±12.93 |
| | | 0.05 | 56,440 | 289 | 1 | 122±0.55 | 0.48±0.04 | 69.38±5.04 |
| | | 0.01 | 29,620 | 508 | 2 | 120±0.59 | 0.68±0.07 | 55.61±2.24 |
| 100 | | - | - | 201 | 1 | 137±0.85 | 0.06±0.01 | 0.44±0.04 |
| | | - | - | 401 | 2 | 125±0.92 | 0.85±0.03 | 1.29±0.06 |
| | $k$-means\|\| | - | - | 601 | 3 | 120±0.66 | 2.08±0.05 | 2.58±0.07 |
| | | - | - | 801 | 4 | 117±0.5 | 3.75±0.05 | 4.29±0.1 |
| | | - | - | 1001 | 5 | 115±0.61 | 5.62±0.08 | 6.15±0.08 |
| | | 0.1 | 277,721 | 470±20 | 1 | 119±2.93 | 0.65±0.02 | 671.95±18.65 |
| | SOCCER | 0.05 | 124,053 | 496 | 1 | 115±0.51 | 0.67±0.05 | 251.17±16.4 |
| | | 0.01 | 65,104 | 820 | 2 | 119±0.79 | 0.84±0.02 | 192.33±5.38 |
| 200 | | - | - | 401 | 1 | 122±1.22 | 0.06±0 | 0.84±0.04 |
| | | - | - | 801 | 2 | 112±0.27 | 1.69±0.07 | 2.6±0.08 |
| | $k$-means\|\| | - | - | 1201 | 3 | 108±0.26 | 4.1±0.1 | 5.13±0.1 |
| | | - | - | 1601 | 4 | 106±0.31 | 7.13±0.15 | 8.33±0.17 |
| | | - | - | 2001 | 5 | 104±0.21 | 11.08±0.33 | 12.42±0.31 |

Table 6: `Census1990` dataset experiments with Standard KMeans as black-box. 'T' stands for time in seconds.

| $k$ | ALG | $\epsilon$ | $P_1$ | Output size | Rounds | Cost $(\cdot 10^6)$ | T. Machine | T. Total |
|---|---|---|---|---|---|---|---|---|
|  | SOCCER | 0.2 | 95,908 | 90 | 1 | 172±1.34 | 0.12±0.03 | 17.67±1.15 |
|  |  | 0.1 | 22,018 | 121 | 1 | 188±4.89 | 0.1±0.02 | 5.4±0.3 |
|  |  | 0.05 | 10,550 | 204 | 2 | 179±2.68 | 0.11±0.01 | 5.69±0.24 |
|  |  | 0.01 | 5,856 | 489 | 4 | 176±1.07 | 0.23±0.02 | 8.69±0.17 |
| 25 |  | - | - | 51 | 1 | 418±133.76 | 0.05±0 | 0.16±0.04 |
|  |  | - | - | 101 | 2 | 218±11.09 | 0.15±0 | 0.27±0.03 |
|  | $k$-means\|\| | - | - | 151 | 3 | 199±3.39 | 0.28±0 | 0.4±0.03 |
|  |  | - | - | 201 | 4 | 188±2.8 | 0.44±0.01 | 0.56±0.03 |
|  |  | - | - | 251 | 5 | 185±3.03 | 0.61±0.01 | 0.76±0.06 |
|  | SOCCER | 0.2 | 215,487 | 121 | 1 | 131±1.47 | 0.11±0.01 | 49.62±3.63 |
|  |  | 0.1 | 49,471 | 177 | 1 | 156±2.77 | 0.11±0.01 | 14.57±1.24 |
|  |  | 0.05 | 23,704 | 266 | 2 | 140±2.39 | 0.14±0.03 | 15.61±0.65 |
|  |  | 0.01 | 13,158 | 592 | 4 | 138±1.55 | 0.24±0.03 | 20.17±0.53 |
| 50 |  | - | - | 101 | 1 | 318±79.46 | 0.05±0 | 0.27±0.05 |
|  |  | - | - | 201 | 2 | 169±6.17 | 0.2±0 | 0.41±0.02 |
|  | $k$-means\|\| | - | - | 301 | 3 | 153±3.33 | 0.42±0.01 | 0.67±0.05 |
|  |  | - | - | 401 | 4 | 144±1.97 | 0.66±0.01 | 0.94±0.07 |
|  |  | - | - | 501 | 5 | 140±1.43 | 0.97±0.02 | 1.3±0.07 |
|  | SOCCER | 0.2 | 478,318 | 177 | 1 | 100±0.81 | 0.15±0.02 | 172.23±11.82 |
|  |  | 0.1 | 109,813 | 283 | 1 | 132±2.7 | 0.14±0.01 | 43.41±2 |
|  |  | 0.05 | 52,616 | 378 | 2 | 110±0.87 | 0.17±0.03 | 45.09±2.1 |
|  |  | 0.01 | 29,207 | 712 | 3 | 110±0.88 | 0.29±0.04 | 51.47±1.52 |
| 100 |  | - | - | 201 | 1 | 264±67.1 | 0.05±0 | 0.44±0.02 |
|  |  | - | - | 401 | 2 | 133±2.3 | 0.31±0.01 | 0.77±0.08 |
|  | $k$-means\|\| | - | - | 601 | 3 | 119±1.05 | 0.65±0.01 | 1.15±0.06 |
|  |  | - | - | 801 | 4 | 112±0.97 | 1.1±0.01 | 1.69±0.13 |
|  |  | - | - | 1001 | 5 | 109±0.67 | 1.67±0.03 | 2.21±0.08 |
|  | SOCCER | 0.1 | 241,364 | 489 | 1 | 111±1.9 | 0.19±0.01 | 152.95±9.98 |
|  |  | 0.05 | 115,648 | 563.2 | 1,2 | 89.7±1.13 | 0.22±0.03 | 143.42±10.41 |
|  |  | 0.01 | 64,197 | 1130 | 3 | 87.3±0.54 | 0.41±0.03 | 147.67±6.72 |
| 200 |  | - | - | 401 | 1 | 224±49.18 | 0.05±0 | 0.92±0.14 |
|  |  | - | - | 801 | 2 | 104±2.75 | 0.5±0.01 | 1.45±0.07 |
|  | $k$-means\|\| | - | - | 1201 | 3 | 93.8±0.97 | 1.14±0.02 | 2.27±0.09 |
|  |  | - | - | 1601 | 4 | 88.7±0.45 | 1.96±0.03 | 3.19±0.06 |
|  |  | - | - | 2001 | 5 | 87±0.56 | 3.03±0.06 | 4.35±0.1 |

Table 7: `KDDCup1999` dataset experiments with Standard KMeans as black-box. 'T' stands for time in seconds.

| $k$ | ALG | $\epsilon$ | $P_1$ | Output size | Rounds | Cost $(\cdot 10^{12})$ | T. Machine | T. Total |
|---|---|---|---|---|---|---|---|---|
| | | 0.2 | 110,088 | 115 | 1 | 112.79±10.71 | 0.15±0.02 | 9.25±1.84 |
| | SOCCER | 0.1 | 23,590 | 236±40 | 2.2±0.4 | 118.21±18.54 | 0.24±0.02 | 5.62±0.96 |
| | | 0.05 | 10,920 | 433 | 4 | 130.33±12.46 | 0.35±0.03 | 6.04±0.32 |
| | | 0.01 | 5,896 | 1324±49 | 11.2±0.42 | 113.55±10.09 | 1.01±0.09 | 13.91±0.73 |
| 25 | | - | - | 51 | 1 | 253.76±34.98 | 0.07±0 | 0.18±0.03 |
| | | - | - | 101 | 2 | 157.12±12.26 | 0.23±0 | 0.34±0.03 |
| | $k$-means‖ | - | - | 151 | 3 | 148.23±22.31 | 0.44±0.01 | 0.55±0.03 |
| | | - | - | 201 | 4 | 124.1±3.3 | 0.71±0.01 | 0.82±0.04 |
| | | - | - | 251 | 5 | 126.4±11.82 | 1.03±0.01 | 1.15±0.02 |
| | | 0.2 | 247,347 | 171 | 1 | 21.77±1.33 | 0.18±0.03 | 23.34±3.9 |
| | SOCCER | 0.1 | 53,003 | 304 | 2 | 23.71±4.96 | 0.3±0.02 | 11.84±0.83 |
| | | 0.05 | 24,535 | 515±70 | 3.5±0.5 | 23.95±3.82 | 0.41±0.02 | 11.1±1.21 |
| | | 0.01 | 13,249 | 1352±62 | 8.8±0.4 | 22.98±2.19 | 0.96±0.06 | 19.23±1.05 |
| 50 | | - | - | 101 | 1 | 108.72±29.12 | 0.07±0 | 0.25±0.02 |
| | | - | - | 201 | 2 | 37.17±10.96 | 0.33±0.01 | 0.56±0.08 |
| | $k$-means‖ | - | - | 301 | 3 | 35.18±4.36 | 0.69±0.01 | 0.91±0.05 |
| | | - | - | 401 | 4 | 35.19±5.25 | 1.13±0.01 | 1.36±0.04 |
| | | - | - | 501 | 5 | 32.8±5.59 | 1.69±0.02 | 1.97±0.08 |
| | | 0.2 | 549,037 | 277 | 1 | 7.43±0.66 | 0.27±0.04 | 71.35±7.87 |
| | SOCCER | 0.1 | 117,651 | 466 | 2 | 8.07±0.76 | 0.39±0.02 | 29.1±2.11 |
| | | 0.05 | 54,461 | 667 | 3 | 7.13±0.52 | 0.55±0.03 | 24.85±1.15 |
| | | 0.01 | 29,409 | 1528 | 7 | 5.97±0.36 | 1.15±0.07 | 37.81±1.3 |
| 100 | | - | - | 201 | 1 | 51.75±15.63 | 0.06±0 | 0.47±0.05 |
| | | - | - | 401 | 2 | 6.49±0.74 | 0.54±0.01 | 0.92±0.04 |
| | $k$-means‖ | - | - | 601 | 3 | 8.41±0.8 | 1.16±0.02 | 1.62±0.14 |
| | | - | - | 801 | 4 | 8.54±1.34 | 1.99±0.03 | 2.4±0.04 |
| | | - | - | 1001 | 5 | 7.95±0.49 | 3.05±0.03 | 3.55±0.13 |
| | | 0.1 | 258,592 | 778 | 2 | 3.06±0.08 | 0.56±0.04 | 93.74±6.77 |
| | SOCCER | 0.05 | 119,705 | 1088 | 3 | 2.89±0.33 | 0.73±0.05 | 64.74±4.33 |
| | | 0.01 | 64,641 | 2060 | 6 | 2.46±0.21 | 1.46±0.07 | 82.2±2.99 |
| 200 | | - | - | 401 | 1 | 10.85±2.14 | 0.06±0 | 0.89±0.18 |
| | | - | - | 801 | 2 | 1.71±0.5 | 0.92±0.02 | 1.81±0.23 |
| | $k$-means‖ | - | - | 1201 | 3 | 2.62±0.35 | 2.11±0.05 | 3.06±0.25 |
| | | - | - | 1601 | 4 | 2.76±0.11 | 3.64±0.09 | 4.6±0.12 |
| | | - | - | 2001 | 5 | 2.41±0.35 | 5.69±0.08 | 6.73±0.11 |

Table 8: `BigCross` dataset experiments with Standard KMeans as black-box. 'T' stands for time in seconds.

| $k$ | ALG | $\epsilon$ | $P_1$ | Output size | Rounds | Cost ($\cdot 10^{10}$) | T. Machine | T. Total |
|---|---|---|---|---|---|---|---|---|
| | | 0.2 | 126,978 | 90 | 1 | 328± 5 | 0.43±0.06 | 60.32±3.7 |
| | SOCCER | 0.1 | 25,335 | 106±13 | 1 | 332± 7 | 0.39±0.03 | 11.95±0.9 |
| | | 0.05 | 11,316 | 204 | 2 | 345± 5 | 0.4±0.05 | 6.82±0.24 |
| | | 0.01 | 5,939 | 358±13 | 3 | 319± 2 | 0.87±0.08 | 7.78±0.26 |
| 25 | | - | - | 51 | 1 | 519± 40 | 0.18±0.01 | 0.27±0.03 |
| | | - | - | 101 | 2 | 367± 9 | 0.49±0.01 | 0.6±0.02 |
| | $k$-means\|\| | - | - | 151 | 3 | 350± 5 | 0.93±0.04 | 1.05±0.06 |
| | | - | - | 201 | 4 | 339± 6 | 1.59±0.05 | 1.71±0.06 |
| | | - | - | 251 | 5 | 330± 6 | 2.14±0.09 | 2.28±0.1 |
| | | 0.2 | 285,296 | 121 | 1 | 224± 4 | 0.41±0.01 | 164.48±16.65 |
| | SOCCER | 0.1 | 56,924 | 127 | 1 | 221± 3 | 0.47±0.05 | 35.69±2.03 |
| | | 0.05 | 25,427 | 266 | 2 | 242± 3 | 0.5±0.05 | 19.17±1.26 |
| | | 0.01 | 13,344 | 444 | 3 | 215± 1 | 0.83±0.06 | 17.42±0.61 |
| 50 | | - | - | 101 | 1 | 365± 28 | 0.18±0.02 | 0.39±0.03 |
| | | - | - | 201 | 2 | 244± 6 | 0.78±0.07 | 1.02±0.09 |
| | $k$-means\|\| | - | - | 301 | 3 | 230± 2 | 1.33±0.01 | 1.58±0.03 |
| | | - | - | 401 | 4 | 223± 2 | 2.27±0.07 | 2.55±0.1 |
| | | - | - | 501 | 5 | 217± 1 | 3.54±0.19 | 3.84±0.12 |
| | | 0.2 | 633,272 | 177 | 1 | 151± 1 | 0.54±0.06 | 510.31±41.78 |
| | SOCCER | 0.1 | 126,354 | 183 | 1 | 152± 1 | 0.53±0.03 | 105.06±9.78 |
| | | 0.05 | 56,440 | 289 | 1 | 170± 2 | 0.61±0.05 | 48.38±3.31 |
| | | 0.01 | 29,620 | 580± 50 | 2,3 | 154± 2 | 0.94±0.07 | 48.47±3.34 |
| 100 | | - | - | 201 | 1 | 242± 20 | 0.18±0.03 | 0.56±0.06 |
| | | - | - | 401 | 2 | 169± 2 | 1.1±0.04 | 1.54±0.05 |
| | $k$-means\|\| | - | - | 601 | 3 | 157± 2 | 2.43±0.18 | 2.94±0.18 |
| | | - | - | 801 | 4 | 153± 1 | 3.94±0.25 | 4.49±0.33 |
| | | - | - | 1001 | 5 | 150± 1 | 6.17±0.31 | 6.71±0.32 |
| | | 0.1 | 277,721 | 289 | 1 | 103± 0 | 0.73±0.04 | 336.87±27.02 |
| | SOCCER | 0.05 | 124,053 | 496 | 1 | 117± 2 | 0.74±0.04 | 142.19±9.06 |
| | | 0.01 | 65,104 | 820 | 2 | 109± 1 | 1.13±0.07 | 121.35±4.43 |
| 200 | | - | - | 401 | 1 | 166± 9 | 0.21±0.06 | 1.04±0.07 |
| | | - | - | 801 | 2 | 119± 1 | 1.8±0.03 | 2.72±0.03 |
| | $k$-means\|\| | - | - | 1201 | 3 | 111± 1 | 4.08±0.13 | 5.31±0.28 |
| | | - | - | 1601 | 4 | 107± 1 | 7.31±0.09 | 8.56±0.09 |
| | | - | - | 2001 | 5 | 106± 0 | 10.86±0.32 | 12.24±0.32 |

## A.7 Results for `MiniBatchKMeans` as black-box for SOCCER

The results for SOCCER and $k$-`means||`, when `MiniBatchKMeans` used as black-box for SOCCER are provided below in Table 9, Table 10 , Table 11, , Table 12, and Table 13.

Table 9: $k$-`GaussianMixture` dataset experiments with `MiniBatchKMeans` used as black-box. 'T' stands for time in seconds.

| $k$ | ALG | $\epsilon$ | $P_1$ | Output size | Rounds | Cost | T (machine) | T (Total) |
|---|---|---|---|---|---|---|---|---|
|  | SOCCER | 0.2 | 126,978 | 105±13 | 1 | 150.2±0.2 | 0.32±0.06 | 1.03±0.2 |
|  |  | 0.1 | 25,335 | 161±50 | 1.6±0.5 | 150.3±0.1 | 0.49±0.12 | 1.14±0.25 |
|  |  | 0.05 | 11,316 | 178±53 | 1.5±0.5 | 150.5±0.1 | 0.49±0.14 | 1.05±0.28 |
|  |  | 0.01 | 5,939 | 348 | 3 | 150.1±0 | 0.74±0.12 | 1.67±0.18 |
| 25 | $k$-means\|\| | - | - | 51 | 1 | 1,688,270.3±951992.1 | 0.05±0 | 0.15±0.03 |
|  |  | - | - | 101 | 2 | 37,530.5±46409 | 0.33±0.01 | 0.43±0.02 |
|  |  | - | - | 151 | 3 | 196.9±18.6 | 0.76±0.03 | 0.87±0.03 |
|  |  | - | - | 201 | 4 | 171.2±4.7 | 1.32±0.06 | 1.42±0.06 |
|  |  | - | - | 251 | 5 | 164.4±2.1 | 1.98±0.07 | 2.1±0.07 |
|  | SOCCER | 0.2 | 285,296 | 171 | 1 | 150.4±0.3 | 0.52±0.1 | 1.91±0.28 |
|  |  | 0.1 | 56,924 | 216±41 | 1.5±0.5 | 152.1±5.2 | 0.51±0.07 | 1.5±0.31 |
|  |  | 0.05 | 25,427 | 244±42 | 1.7±0.5 | 150.6±0.2 | 0.65±0.14 | 1.48±0.23 |
|  |  | 0.01 | 13,344 | 361±47 | 2.1±0.3 | 150.3±0.1 | 0.83±0.09 | 1.86±0.18 |
| 50 | $k$-means\|\| | - | - | 101 | 1 | 1,283,640.5±558248.2 | 0.05±0 | 0.24±0.04 |
|  |  | - | - | 201 | 2 | 13,399.3±11108.8 | 0.62±0.03 | 0.81±0.04 |
|  |  | - | - | 301 | 3 | 211.5±6.8 | 1.45±0.04 | 1.66±0.04 |
|  |  | - | - | 401 | 4 | 174.9±2.3 | 2.52±0.05 | 2.71±0.07 |
|  |  | - | - | 501 | 5 | 166±1 | 3.83±0.08 | 4.04±0.12 |
|  | SOCCER | 0.2 | 633,271 | 277 | 1 | 628.2±612.3 | 0.67±0.13 | 4.19±0.76 |
|  |  | 0.1 | 126,354 | 406±52 | 2 | 154±6.1 | 0.91±0.11 | 3.22±0.3 |
|  |  | 0.05 | 56,440 | 390±54 | 1.9±0.3 | 150.7±0.5 | 0.94±0.1 | 2.7±0.37 |
|  |  | 0.01 | 29,620 | 550±54 | 2.4±0.5 | 150.3±0 | 1.05±0.13 | 2.91±0.21 |
| 100 | $k$-means\|\| | - | - | 201 | 1 | 1,079,458.8±266814.6 | 0.05±0.01 | 0.45±0.13 |
|  |  | - | - | 401 | 2 | 25,866.5±16072.3 | 1.09±0.03 | 1.51±0.14 |
|  |  | - | - | 601 | 3 | 226.9±61.2 | 2.67±0.06 | 3.09±0.1 |
|  |  | - | - | 801 | 4 | 176.6±3.1 | 4.75±0.05 | 5.21±0.08 |
|  |  | - | - | 1001 | 5 | 167.2±1.7 | 7.09±0.1 | 7.52±0.11 |
|  | SOCCER | 0.1 | 277,721 | 733±96 | 1.9±0.3 | 156.9±4.8 | 1.42±0.17 | 6.94±0.31 |
|  |  | 0.05 | 124,053 | 757±74 | 2 | 150.9±0.6 | 1.47±0.12 | 5.13±0.29 |
|  |  | 0.01 | 65,104 | 1008±86 | 3 | 855.6±1139.9 | 1.71±0.11 | 5.53±0.26 |
| 200 | $k$-means\|\| | - | - | 401 | 1 | 1,104,954±201686.7 | 0.05±0.01 | 0.82±0.06 |
|  |  | - | - | 801 | 2 | 26,593.9±9916.1 | 2.08±0.04 | 3±0.14 |
|  |  | - | - | 1201 | 3 | 218.8±14.9 | 4.97±0.09 | 5.89±0.1 |
|  |  | - | - | 1601 | 4 | 175.7±1.5 | 8.72±0.1 | 9.74±0.17 |
|  |  | - | - | 2001 | 5 | 167±1.8 | 13.16±0.12 | 14.18±0.17 |

Table 10: `Higgs` dataset experiments with MiniBatchKMeans used as black-box. 'T' stands for time in seconds.

| $k$ | Alg | $\epsilon$ | $P_1$ | Output size | Rounds | Cost ($\cdot 10^6$) | T (machine) | T (Total) |
|---|---|---|---|---|---|---|---|---|
| | | 0.2 | 126,978 | 92±2 | 1 | 129.5±0.9 | 0.27±0.03 | 1.04±0.12 |
| | SOCCER | 0.1 | 25,335 | 121 | 1 | 141.5±2 | 0.29±0.02 | 0.83±0.09 |
| | | 0.05 | 11,316 | 204 | 2 | 144.5±1.8 | 0.31±0.04 | 0.99±0.11 |
| | | 0.01 | 5,939 | 348 | 3 | 135.2±0.7 | 0.49±0.04 | 1.52±0.19 |
| 25 | | - | - | 51 | 1 | 171±4 | 0.05±0 | 0.16±0.02 |
| | | - | - | 101 | 2 | 153±1.5 | 0.31±0.02 | 0.43±0.03 |
| | $k$-means\|\| | - | - | 151 | 3 | 148±1.4 | 0.68±0.09 | 0.81±0.09 |
| | | - | - | 201 | 4 | 143±1 | 1.06±0.05 | 1.19±0.06 |
| | | - | - | 251 | 5 | 139±0.6 | 1.59±0.04 | 1.72±0.03 |
| | | 0.2 | 285,296 | 123±0.4 | 1 | 118.1±0.3 | 0.35±0.04 | 1.98±0.22 |
| | SOCCER | 0.1 | 56,924 | 177 | 1 | 133.4±1.6 | 0.37±0.04 | 1.26±0.18 |
| | | 0.05 | 25,427 | 183 | 1 | 128.2±0.8 | 0.38±0.05 | 1.16±0.11 |
| | | 0.01 | 13,344 | 346 | 2 | 124.9±0.4 | 0.54±0.06 | 1.69±0.16 |
| 50 | | - | - | 101 | 1 | 153±1.7 | 0.05±0 | 0.25±0.04 |
| | | - | - | 201 | 2 | 139±1.1 | 0.5±0.03 | 0.72±0.05 |
| | $k$-means\|\| | - | - | 301 | 3 | 133±0.6 | 1.13±0.05 | 1.38±0.08 |
| | | - | - | 401 | 4 | 129±0.7 | 1.95±0.03 | 2.2±0.05 |
| | | - | - | 501 | 5 | 127±0.4 | 2.96±0.1 | 3.24±0.11 |
| | | 0.2 | 633,272 | 179±0.4 | 1 | 106.7±0.2 | 0.55±0.03 | 4.03±0.2 |
| | SOCCER | 0.1 | 126,354 | 283 | 1 | 127.5±1.4 | 0.49±0.06 | 2.06±0.2 |
| | | 0.05 | 56,440 | 289 | 1 | 121.5±0.7 | 0.51±0.06 | 1.89±0.18 |
| | | 0.01 | 29,620 | 508 | 2 | 121.1±0.7 | 0.73±0.06 | 2.62±0.14 |
| 100 | | - | - | 201 | 1 | 137±0.8 | 0.06±0.01 | 0.44±0.04 |
| | | - | - | 401 | 2 | 125±0.9 | 0.85±0.03 | 1.29±0.06 |
| | $k$-means\|\| | - | - | 601 | 3 | 120±0.7 | 2.08±0.05 | 2.58±0.07 |
| | | - | - | 801 | 4 | 117±0.5 | 3.75±0.05 | 4.29±0.1 |
| | | - | - | 1001 | 5 | 115±0.6 | 5.62±0.08 | 6.15±0.08 |
| | | 0.1 | 277,721 | 480±17 | 1 | 117.3±1.3 | 0.79±0.06 | 4.23±0.31 |
| | SOCCER | 0.05 | 124,053 | 496 | 1 | 115.3±0.4 | 0.77±0.07 | 3.54±0.21 |
| | | 0.01 | 65,104 | 820 | 2 | 116.7±0.7 | 0.94±0.08 | 4.45±0.34 |
| 200 | | - | - | 401 | 1 | 122±1.2 | 0.06±0 | 0.84±0.04 |
| | | - | - | 801 | 2 | 112±0.3 | 1.69±0.07 | 2.6±0.08 |
| | $k$-means\|\| | - | - | 1201 | 3 | 108±0.3 | 4.1±0.1 | 5.13±0.1 |
| | | - | - | 1601 | 4 | 106±0.3 | 7.13±0.15 | 8.33±0.17 |
| | | - | - | 2001 | 5 | 104±0.2 | 11.08±0.33 | 12.42±0.31 |

Table 11: `Census1990` dataset experiments with MiniBatchKMeans used as black-box. 'T' stands for time in seconds.

| $k$ | ALG | $\epsilon$ | $P_1$ | Output size | Rounds | Cost ($\cdot 10^6$) | T (machine) | T (Total) |
|---|---|---|---|---|---|---|---|---|
| 25 | SOCCER | 0.2 | 95,908 | 90 | 1 | 171.3±1.7 | 0.11±0.05 | 0.96±0.14 |
| | | 0.1 | 22,018 | 121 | 1 | 187.8±4 | 0.11±0.04 | 0.65±0.16 |
| | | 0.05 | 10,550 | 204 | 2 | 179.6±3.6 | 0.14±0.05 | 0.85±0.13 |
| | | 0.01 | 5,856 | 489 | 4 | 175.8±1.6 | 0.3±0.07 | 1.65±0.17 |
| | $k$-means\|\| | - | - | 51 | 1 | 418±133.8 | 0.05±0 | 0.16±0.04 |
| | | - | - | 101 | 2 | 218±11.1 | 0.15±0 | 0.27±0.03 |
| | | - | - | 151 | 3 | 199±3.4 | 0.28±0 | 0.4±0.03 |
| | | - | - | 201 | 4 | 188±2.8 | 0.44±0.01 | 0.56±0.03 |
| | | - | - | 251 | 5 | 185±3 | 0.61±0.01 | 0.76±0.06 |
| 50 | SOCCER | 0.2 | 215,487 | 121 | 1 | 129.1±1.4 | 0.1±0.03 | 1.79±0.19 |
| | | 0.1 | 49,471 | 177 | 1 | 155.2±2.8 | 0.15±0.05 | 1.08±0.14 |
| | | 0.05 | 23,704 | 266 | 2 | 139.8±1.4 | 0.15±0.06 | 1.2±0.16 |
| | | 0.01 | 13,158 | 592 | 4 | 135.6±1.1 | 0.31±0.08 | 2.02±0.18 |
| | $k$-means\|\| | - | - | 101 | 1 | 318±79.5 | 0.05±0 | 0.27±0.05 |
| | | - | - | 201 | 2 | 169±6.2 | 0.2±0 | 0.41±0.02 |
| | | - | - | 301 | 3 | 153±3.3 | 0.42±0.01 | 0.67±0.05 |
| | | - | - | 401 | 4 | 144±2 | 0.66±0.01 | 0.94±0.07 |
| | | - | - | 501 | 5 | 140±1.4 | 0.97±0.02 | 1.3±0.07 |
| 100 | SOCCER | 0.2 | 478,318 | 177 | 1 | 99.6±0.8 | 0.22±0.08 | 4.27±0.34 |
| | | 0.1 | 109,813 | 283 | 1 | 126.8±3.1 | 0.18±0.06 | 1.9±0.18 |
| | | 0.05 | 52,616 | 378 | 2 | 110.5±1.3 | 0.17±0.06 | 2.22±0.16 |
| | | 0.01 | 29,207 | 722±32 | 3.1±0.3 | 105.3±1.4 | 0.39±0.09 | 2.98±0.18 |
| | $k$-means\|\| | - | - | 201 | 1 | 264±67.1 | 0.05±0 | 0.44±0.02 |
| | | - | - | 401 | 2 | 133±2.3 | 0.31±0.01 | 0.77±0.08 |
| | | - | - | 601 | 3 | 119±1 | 0.65±0.01 | 1.15±0.06 |
| | | - | - | 801 | 4 | 112±1 | 1.1±0.01 | 1.69±0.13 |
| | | - | - | 1001 | 5 | 109±0.7 | 1.67±0.03 | 2.21±0.08 |
| 200 | SOCCER | 0.1 | 241,364 | 489 | 1 | 107.8±2.7 | 0.31±0.06 | 4.33±0.43 |
| | | 0.05 | 115,648 | 592 | 2 | 89.4±0.6 | 0.29±0.07 | 4.48±0.26 |
| | | 0.01 | 64,197 | 1130 | 3 | 87.3±0.2 | 0.47±0.08 | 5.49±0.34 |
| | $k$-means\|\| | - | - | 401 | 1 | 224±49.2 | 0.05±0 | 0.92±0.14 |
| | | - | - | 801 | 2 | 104±2.7 | 0.5±0.01 | 1.45±0.07 |
| | | - | - | 1201 | 3 | 93.8±1 | 1.14±0.02 | 2.27±0.09 |
| | | - | - | 1601 | 4 | 88.7±0.4 | 1.96±0.03 | 3.19±0.06 |
| | | - | - | 2001 | 5 | 87±0.6 | 3.03±0.06 | 4.35±0.1 |

Table 12: `KDDCup1999` dataset experiments with MiniBatchKMeans used as black-box. 'T' stands for time in seconds.

| $k$ | Alg | $\epsilon$ | $P_1$ | Output size | Rounds | Cost ($\cdot 10^{10}$) | T (machine) | T (Total) |
|---|---|---|---|---|---|---|---|---|
| | | 0.2 | 110,088 | 115 | 1 | 6,273,229±102,539 | 0.14±0.01 | 0.94±0.15 |
| | SOCCER | 0.1 | 23,590 | 313 | 3 | 3,034,310±1,064,972 | 0.25±0.01 | 1.33±0.13 |
| | | 0.05 | 10,920 | 433 | 4 | 3,982,745±1124135 | 0.37±0.04 | 1.6±0.2 |
| | | 0.01 | 5,896 | 1243±61 | 10.5±0.5 | 2,344,966±946,061 | 0.88±0.06 | 3.88±0.21 |
| 25 | | - | - | 51 | 1 | 254±35 | 0.07±0 | 0.18±0.03 |
| | | - | - | 101 | 2 | 157±12.3 | 0.23±0 | 0.34±0.03 |
| | $k$-means\|\| | - | - | 151 | 3 | 148±22.3 | 0.44±0.01 | 0.55±0.03 |
| | | - | - | 201 | 4 | 124±3.3 | 0.71±0.01 | 0.82±0.04 |
| | | - | - | 251 | 5 | 126±11.8 | 1.03±0.01 | 1.15±0.02 |
| | | 0.2 | 247,347 | 171 | 1 | 5,971,765±351183 | 0.17±0.01 | 1.59±0.18 |
| | SOCCER | 0.1 | 53,003 | 304 | 2 | 5,716,218±507,248 | 0.29±0.01 | 1.5±0.13 |
| | | 0.05 | 24,535 | 555±56 | 3.8±0.4 | 2,995,765±1,379,745 | 0.43±0.03 | 2.05±0.21 |
| | | 0.01 | 13,249 | 1248±47 | 8.1±0.3 | 3,946,555±954,957 | 0.89±0.02 | 3.86±0.23 |
| 50 | | - | - | 101 | 1 | 109±29.1 | 0.07±0 | 0.25±0.02 |
| | | - | - | 201 | 2 | 37.2±11 | 0.33±0.01 | 0.56±0.08 |
| | $k$-means\|\| | - | - | 301 | 3 | 35.2±4.4 | 0.69±0.01 | 0.91±0.05 |
| | | - | - | 401 | 4 | 35.2±5.3 | 1.13±0.01 | 1.36±0.04 |
| | | - | - | 501 | 5 | 32.8±5.6 | 1.69±0.02 | 1.97±0.08 |
| | | 0.2 | 549,037 | 277 | 1 | 5,969,280±355,576 | 0.26±0.02 | 3.22±0.22 |
| | SOCCER | 0.1 | 117,651 | 466 | 2 | 5,587,433±940,351 | 0.38±0.02 | 2.66±0.27 |
| | | 0.05 | 54,461 | 667 | 3 | 3,924,526±674,470 | 0.52±0.02 | 2.89±0.18 |
| | | 0.01 | 29,409 | 1528 | 7 | 3,530,365±841,054 | 1.11±0.03 | 5.36±0.25 |
| 100 | | - | - | 201 | 1 | 51.7±15.6 | 0.06±0 | 0.47±0.05 |
| | | - | - | 401 | 2 | 6.5±0.7 | 0.54±0.01 | 0.92±0.04 |
| | $k$-means\|\| | - | - | 601 | 3 | 8.4±0.8 | 1.16±0.02 | 1.62±0.14 |
| | | - | - | 801 | 4 | 8.5±1.3 | 1.99±0.03 | 2.4±0.04 |
| | | - | - | 1001 | 5 | 8±0.5 | 3.05±0.03 | 3.55±0.13 |
| | | 0.1 | 258,592 | 778 | 2 | 4,561,690±1,036,461 | 0.61±0.04 | 5.3±0.25 |
| | SOCCER | 0.05 | 119,705 | 1088 | 3 | 2,859,826±1,177,854 | 0.78±0.07 | 5.2±0.21 |
| | | 0.01 | 64,641 | 2060 | 6 | 3,814,931±1,066,684 | 1.46±0.08 | 8.22±0.23 |
| 200 | | - | - | 401 | 1 | 10.8±2.1 | 0.06±0 | 0.89±0.18 |
| | | - | - | 801 | 2 | 1.7±0.5 | 0.92±0.02 | 1.81±0.23 |
| | $k$-means\|\| | - | - | 1201 | 3 | 2.6±0.4 | 2.11±0.05 | 3.06±0.25 |
| | | - | - | 1601 | 4 | 2.8±0.1 | 3.64±0.09 | 4.6±0.12 |
| | | - | - | 2001 | 5 | 2.4±0.4 | 5.69±0.08 | 6.73±0.11 |

Table 13: `BigCross` dataset experiments with MiniBatchKMeans used as black-box. 'T' stands for time in seconds.

| $k$ | ALG | epsilon | $P_1$ | Output size | Rounds | Cost ($10^{10}$) | T (machine) | T (Total) |
|---|---|---|---|---|---|---|---|---|
|  |  | 0.2 | 126,978 | 90 | 1 | 328±2.6 | 0.38±0.03 | 1.39±0.13 |
|  | SOCCER | 0.1 | 25,335 | 99±8 | 1 | 327±6.3 | 0.37±0.05 | 0.86±0.12 |
|  |  | 0.05 | 11,316 | 204 | 2 | 345±8 | 0.4±0.03 | 1.1±0.15 |
|  |  | 0.01 | 5,939 | 365±12 | 3 | 318±3.2 | 0.79±0.02 | 1.81±0.11 |
| 25 |  | - | - | 51 | 1 | 519±39.6 | 0.18±0.01 | 0.27±0.03 |
|  |  | - | - | 101 | 2 | 367±8.7 | 0.49±0.01 | 0.6±0.02 |
|  | $k$-means\|\| | - | - | 151 | 3 | 350±4.8 | 0.93±0.04 | 1.05±0.06 |
|  |  | - | - | 201 | 4 | 339±6.3 | 1.59±0.05 | 1.71±0.06 |
|  |  | - | - | 251 | 5 | 330±5.5 | 2.14±0.09 | 2.28±0.1 |
|  |  | 0.2 | 285,296 | 121 | 1 | 222±3.7 | 0.47±0.05 | 2.44±0.23 |
|  | SOCCER | 0.1 | 56,924 | 142±24 | 1 | 221±2.1 | 0.45±0.02 | 1.3±0.17 |
|  |  | 0.05 | 25,427 | 266 | 2 | 242±4.2 | 0.48±0.04 | 1.55±0.19 |
|  |  | 0.01 | 13,344 | 444 | 3 | 217±1.8 | 0.85±0.13 | 2.26±0.2 |
| 50 |  | - | - | 101 | 1 | 365±28 | 0.18±0.02 | 0.39±0.03 |
|  |  | - | - | 201 | 2 | 244±5.9 | 0.78±0.07 | 1.02±0.09 |
|  | $k$-means\|\| | - | - | 301 | 3 | 230±2.1 | 1.33±0.01 | 1.58±0.03 |
|  |  | - | - | 401 | 4 | 223±2.3 | 2.27±0.07 | 2.55±0.1 |
|  |  | - | - | 501 | 5 | 217±1.4 | 3.54±0.19 | 3.84±0.12 |
|  |  | 0.2 | 633,272 | 177 | 1 | 150±1.3 | 0.59±0.05 | 5.17±0.23 |
|  | SOCCER | 0.1 | 126,354 | 193 | 1 | 148±1.2 | 0.59±0.07 | 2.3±0.15 |
|  |  | 0.05 | 56,440 | 298±28 | 1.1±0.3 | 169±2.7 | 0.59±0.05 | 2.15±0.18 |
|  |  | 0.01 | 29,620 | 612 | 3 | 154±1.3 | 0.87±0.03 | 3.1±0.17 |
| 100 |  | - | - | 201 | 1 | 242±20 | 0.18±0.03 | 0.56±0.06 |
|  |  | - | - | 401 | 2 | 169±2.3 | 1.1±0.04 | 1.54±0.05 |
|  | $k$-means\|\| | - | - | 601 | 3 | 157±2 | 2.43±0.18 | 2.94±0.18 |
|  |  | - | - | 801 | 4 | 153±1.5 | 3.94±0.25 | 4.49±0.33 |
|  |  | - | - | 1001 | 5 | 150±1 | 6.17±0.31 | 6.71±0.32 |
|  |  | 0.1 | 277,721 | 289 | 1 | 102±0.3 | 0.74±0.02 | 4.33±0.25 |
|  | SOCCER | 0.05 | 124,053 | 496 | 1 | 116±2.1 | 0.79±0.07 | 3.72±0.16 |
|  |  | 0.01 | 65,104 | 820 | 2 | 109±0.7 | 1.13±0.05 | 4.74±0.18 |
| 200 |  | - | - | 401 | 1 | 166±8.6 | 0.21±0.06 | 1.04±0.07 |
|  |  | - | - | 801 | 2 | 119±1.4 | 1.8±0.03 | 2.72±0.03 |
|  | $k$-means\|\| | - | - | 1201 | 3 | 111±1 | 4.08±0.13 | 5.31±0.28 |
|  |  | - | - | 1601 | 4 | 107±0.5 | 7.31±0.09 | 8.56±0.09 |
|  |  | - | - | 2001 | 5 | 106±0.4 | 10.86±0.32 | 12.24±0.32 |