

---

# On Distributed Larger-Than-Memory Subset Selection With Pairwise Submodular Functions

---

Maximilian Böther<sup>1</sup> Abraham Sebastian<sup>2</sup> Pranjal Awasthi<sup>2</sup> Ana Klimovic<sup>3</sup> Srikumar Ramalingam<sup>2</sup>

## Abstract

Many learning problems hinge on the fundamental problem of subset selection, i.e., identifying a subset of important and representative points. For example, selecting the most significant samples in ML training cannot only reduce training costs but also enhance model quality. Submodularity, a discrete analogue of convexity, is commonly used for solving subset selection problems. However, existing algorithms for optimizing submodular functions are sequential, and the prior distributed methods require at least one central machine to fit the target subset. In this paper, we relax the requirement of having a central machine for the target subset by proposing a novel distributed bounding algorithm with provable approximation guarantees. The algorithm iteratively bounds the minimum and maximum utility values to select high quality points and discard the unimportant ones. When bounding does not find the complete subset, we use a multi-round, partition-based distributed greedy algorithm to identify the remaining subset. We show that these algorithms find high quality subsets on CIFAR-100 and ImageNet with marginal or no loss in quality compared to centralized methods, and scale to a dataset with 13 billion points.

## 1. Introduction

The increasing volume of collected data requires identifying highly informative subsets of features or datasets to cost-effectively train high-quality models (Bhardwaj et al., 2022; Böther et al., 2023). For instance, vision datasets have scaled from ImageNet’s 1.2 M samples (Deng et al.,

2009) to LAION’s 5 B samples (Schuhmann et al., 2022). A single autonomous vehicle collects terabytes of sensor data daily (Kazhemiaka et al., 2021) and recent language/vision models pretrain on billions of examples from books and big webpage collections (Gao et al., 2020; Shen et al., 2023). The ever-increasing scale of datasets makes it challenging and costly to train ML models on all available data. Selecting high quality samples also improves model quality compared to training with larger, less informative datasets. Subset selection is also crucial in applications beyond ML model training, such as feature selection, dictionary learning, and numerous compressed sensing applications (Krause & Cevher, 2010; Das & Kempe, 2011).

Subset selection is a well studied problem with many competing algorithms that rely on submodularity, coresets, and other clustering-based methods. Various greedy algorithms offer strong approximation guarantees but are inherently sequential and centralized (Nemhauser et al., 1978). Existing distributed methods typically partition the entire dataset, run the greedy algorithm on each partition, and use the greedy algorithm again on the union of the subsets from the different partitions (Mirzasoleiman et al., 2016; Barbosa et al., 2015). This final step, mainly used to achieve strong approximation guarantees, requires a machine that holds the results from all partitions. This is not feasible in practice, and we are not aware of any prior methods that are tested on really massive datasets.

In this paper, we present a novel bounding algorithm that iteratively tightens the maximum and minimum utilities of the individual points. This allows to identify high-utility points and discard less informative ones in a distributed fashion. This algorithm is highly parallelizable and can be implemented in distributed data processing frameworks, such as Apache Beam (Akidau et al., 2015). If the bounding algorithm does not compute the complete subset, we employ a distributed, multi-round, partition-based greedy algorithm to achieve high-quality subsets. With bounding and multi-round partition-based optimization, we are able to select high quality subsets that achieve similar quality to those obtained by centralized algorithms. Neither our bounding method nor the distributed greedy algorithm requires a central machine with sufficient memory to hold the final subset.

---

<sup>1</sup>ETH Zurich, Switzerland. Work done while author was interning at Google. <sup>2</sup>Google, United States of America <sup>3</sup>ETH Zurich, Switzerland. Correspondence to: Maximilian Böther <mboether@inf.ethz.ch>, Srikumar Ramalingam <rsrikumar@google.com>.

Among many competing subset selection methods, we opt for pairwise submodular functions for several reasons. First, submodular functions have been shown to produce state-of-the-art results in a variety of applications (c.f. Section 2). Second, they allow us to use a graph for encoding point similarities without needing to consider more complex interactions, i.e., hyper-edges which cannot be constructed efficiently. Furthermore, many other techniques based on  $k$ -medioids (Park & Jun, 2009) or prototypes (Kim et al., 2016) can also be seen as the maximization of pairwise submodular functions. We make the following contributions.

- We design a bounding algorithm that, by continuously adjusting maximum and minimum utility values, identifies high utility points to expand the subset and discards uninformative ones to reduce the ground set. We also show theoretical approximation guarantees.
- The bounding algorithm does not always find a complete subset. To compute the remaining subset, when necessary, we introduce a distributed multi-round partition-based algorithm that empirically achieves similar quality results as the centralized greedy algorithm.
- We show that our distributed methods lead to similar quality results as centralized methods on the CIFAR-100 (50k samples) and ImageNet (1.2M samples) datasets (Krizhevsky & Hinton, 2009; Deng et al., 2009). We also study scalability on a 13 B Perturbed-ImageNet dataset that we generate from ImageNet.

## 2. Prior Work

**Submodular subset selection.** Several discrete problems, such as cut functions, coverage, and entropy, can be formulated as the minimization or maximization of submodular functions (c.f. Appendix A). Many subset selection tasks can be modeled as submodular maximization problems and applications of submodularity include data selection, feature engineering, sensor placement, and influence maximization (Carbonell & Goldstein, 1998; Simon et al., 2007; Krause & Golovin, 2014; Golovin & Krause, 2011; Wei et al., 2015; Lin & Bilmes, 2011; Kaushal et al., 2019; Prasad et al., 2014; Badanidiyuru et al., 2014; Kim et al., 2016; Prasad et al., 2014; Golovin & Krause, 2011; Kim et al., 2016; Ramalingam et al., 2021; Kothawade et al., 2021; 2022). While submodular functions can be minimized in polynomial time, their maximization is NP-hard. Efficient greedy maximization algorithms with approximation guarantees exist (Nemhauser et al., 1978). Even after four decades, this simple greedy method is the gold standard for centralized submodular maximization.

**Distributed algorithms.** Mirzasoileiman et al. (2016) present the distributed GREEDI algorithm for maximizing monotonic submodular functions under a cardinality constraint. The data is partitioned arbitrarily across machines,

and each machine runs the centralized greedy algorithm on the data assigned to it. Then, the greedy algorithm is run again on the union of the results from the machines. This can be implemented in the distributed MapReduce framework (Dean & Ghemawat, 2008). Barbosa et al. (2015) extend GREEDI to RANDGREEDI by changing the way the data is assigned to each machine. Instead of assigning arbitrarily, the points are assigned to individual machines uniformly at random. This randomization improves the theoretical guarantees and leads to constant-factor worst case approximation guarantees.

Kumar et al. (2015) develop SAMPLE&PRUNE, a MapReduce algorithm for maximizing a monotonic submodular function under a cardinality or matroid constraint. The idea is to sample a set, and prune points from the ground set. SAMPLE&PRUNE has a constant approximation bound but assumes that the memory per machine is  $\mathcal{O}(kn^\delta)$ , where  $k$  is the cardinality of the final subset,  $n$  is the number of input points, and  $\delta > 0$ . Liu & Vondrák (2018) discuss a MapReduce algorithm with constant number of rounds while assuming that there exists a central machine with  $\mathcal{O}(\sqrt{nk} \log(k))$  memory.

Subset selection is related to coresets, which the idea is that the solution to an optimization on this coreset closely matches the solution from the entire dataset. Composable coreset-based algorithms divide the entire dataset into blocks, and compute the coresets on the individual blocks followed by the computation of the coreset on the union (Indyk et al., 2014; Mirrokni & Zadimoghaddam, 2015). In the context of clustering, parallel algorithms have also been studied for the  $k$ -center objective (Ene et al., 2011; Im & Moseley, 2015; Malkomes et al., 2015; McClintock & Wirth, 2016; Ramalingam et al., 2023). Other subset selection techniques that sample during model training exist (Katharopoulos & Fleuret, 2018; Mirzasoileiman et al., 2020; Pooladzandi et al., 2022).

**Limitations of prior methods.** We see two main limitations: (1) almost all methods rely on having a central machine that can hold the subset completely to show strong theoretical guarantees, and (2) actual validation on really large datasets has received limited to no attention. We take a more practical approach and demonstrate results on massive datasets with billions of points, by leveraging a subclass of submodular functions, which is effective on a large class of subset selection problems.

## 3. Centralized Subset Selection

In this section, we describe the problem setting and the centralized greedy algorithm. We refer to Appendix A for the definitions of submodularity and monotonicity. Given a ground set  $V$ , our goal is to find a subset  $S \subseteq V$  of

size  $k \leq |V|$  that maximizes a submodular and monotonically non-decreasing set function  $f : 2^V \rightarrow \mathbb{R}$ . We want to determine  $S = \arg \max_{S' \text{ with } |S'|=k} f(S')$ . Since  $f$  is submodular, we can use an efficient greedy algorithm to compute a subset  $S'$  with constant approximation guarantee  $f(S') \geq (1 - \frac{1}{e}) f(S_{\text{OPT}})$  (Nemhauser et al., 1978). As outlined in Algorithm 1, the greedy algorithm repeatedly chooses the data point  $v \in V \setminus S$  that maximizes the marginal gain upon adding the element to the subset.

---

**Algorithm 1** Centralized greedy algorithm for computing subset of size  $k \in \mathbb{N}$ .

---

```

 $S \leftarrow \emptyset$ 
while  $|S| < k$  do
     $S \leftarrow S \cup \{\arg \max_{s \in V \setminus S} f(S \cup \{s\}) - f(S)\}$ 
end while
    
```

---

Let  $\alpha, \beta$  be balancing parameters. Additionally, let  $E$  denote the edges of a nearest neighbor graph,  $u(v)$  denote the utility of  $v$ , and  $s(v_1, v_2)$  denote the similarity of  $v_1$  and  $v_2$ , i.e., the edge weights in the neighbor graph. We focus on a subclass of pairwise submodular function  $f(S) = \alpha \sum_{v \in S} u(v) - \beta \sum_{(v_1, v_2) \in E; v_1, v_2 \in S} s(v_1, v_2)$ , commonly used in subset selection tasks (Wei et al., 2015; Kim et al., 2016; Ramalingam et al., 2021), as well as graph cuts (Kolmogorov & Zabini, 2004). This class of functions enables implementing the algorithm using a priority queue as shown in Algorithm 2.

---

**Algorithm 2** Centralized implementation of Algorithm 1 using a priority queue.

---

```

 $q \leftarrow$  new priority queue
 $S \leftarrow \emptyset$ 
 $\forall v \in V$  : Insert  $v$  into  $q$  with weight  $u(v)$ 
while  $|S| < k$  do
     $v_1 \leftarrow q.\text{popmax}()$ 
    for each  $v_2 \in p$  with  $s(v_1, v_2) > 0$  do
         $q.\text{decrease\_weight\_by}(v_2, \beta \cdot s(v_1, v_2))$ 
    end for
     $S \leftarrow S \cup \{v_1\}$ 
end while
    
```

---

We initialize all of  $V$  into the queue with their utility scores as priority. Then, we repeatedly pop the element  $v_1$  with highest priority, add it to the subset, and update the priority of all neighboring elements  $v_2$  where  $(v_1, v_2) \in E$ . Since we only update the priorities of the neighbors of  $v_1$ , this enables to efficiently build up  $S$  without the need to repeatedly calculate  $f(S \cup \{v\})$  for all  $v$ . We continue to iterate until we build a subset of cardinality  $k$ .

**Scaling challenges.** We see two major roadblocks. First, the algorithm is inherently sequential and relies on identifying the element that maximizes the marginal gain over

the entire dataset at each step. The computation of the marginal gain with respect to the entire dataset makes it difficult to parallelize. Second, the implementation requires large amounts of DRAM to load the keys and similarity values. For datasets with billions of points, this requires more DRAM than commodity VMs or even high-end servers offer. For example, storing 5 billion 64-bit keys and values in the priority queue, and keeping track of 10 nearest neighbors with 64-bit IDs and distances requires 880 GB of RAM.

## 4. Distributed Submodular Subset Selection

Our algorithm has two components: (1) bounding, and (2) distributed greedy optimization. First, the *bounding* algorithm alternates between (1) removing points from the ground set that are less likely to be in the subset and (2) adding points in the subset that are more likely to be in the optimum set, using their minimum and maximum utility. We can implement the bounding algorithm using the distributed data processing frameworks with a large number of machines without the need to store the final subset in a single machine. We propose both exact (Section 4.1) and approximate bounding (Section 4.2) algorithms, with approximation guarantees (Section 4.3). Second, as the bounding algorithm is not guaranteed to find the complete subset in all cases, we use a distributed greedy algorithm based on partitioning to find the remaining points if necessary (Section 4.4).

### 4.1. Exact Bounding

In our exact bounding algorithm, the addition of points to the subset and removal of points from the ground set are guaranteed to preserve the subset quality. The bounding algorithm uses two metrics — the minimum and maximum utilities of a point — that depends on the current partial subset  $S'$  and the current ground set  $V$ .

**Definition 4.1** (Minimum Utility). For a point  $v_1 \in V$ , its *minimum utility* is its utility considering all its neighbors that have not been discarded:  $U_{\min}(v_1) = u(v_1) - \frac{\beta}{\alpha} \sum_{v_2 \in V \cup S' \wedge (v_1, v_2) \in E} s(v_1, v_2)$ .

**Definition 4.2** (Maximum Utility). For a point  $v_1 \in V$ , its *maximum utility* is its utility when only considering the neighbors in the partial subset  $S'$ :  $U_{\max}(v_1) = u(v_1) - \frac{\beta}{\alpha} \sum_{v_2 \in S' \wedge (v_1, v_2) \in E} s(v_1, v_2)$ .

Before we present the bounding algorithm, we introduce two basic blocks for growing the current subset  $S'$ , and shrinking the ground set  $V$  based on the minimum and maximum utilities of the points. Let  $U_{\min}^i$  and  $U_{\max}^i$  denote the  $i$ -th largest minimum and maximum utility points, respectively.

**Grow and Shrink Steps.** Let  $S^*$  denote the optimal solution. The idea behind growth is to identify points with minimum utility higher than the  $k$ -th largest maximum util-

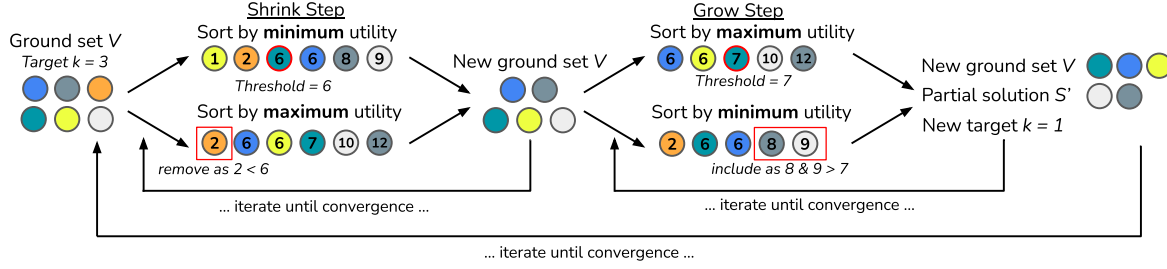


Figure 1. Visualization of the distributed bounding procedure of finding a 50% subset for 6 data points.

ity and include these in the subset, where  $k$  is the number of points we still need to find. The following lemma and the algorithm below summarizes this.

**Lemma 4.3.** For  $v \in V$ , if  $U_{\min}(v) > U_{\max}^k$ , then  $v \in S^*$ .

---

**Algorithm 3**  $\text{Grow}(V, S, k)$

---

$\forall v \in V$ : Compute  $U_{\min}(v)$  and  $U_{\max}(v)$ .  
 Threshold  $\leftarrow U_{\max}^k$   
 $S' \leftarrow \{v \in V \mid U_{\min}(v) > \text{Threshold}\}$   
**return**  $V \setminus S', S' \cup S, k - |S'|$

---

Analogously, as per the following lemma and the algorithm below, we eliminate points with very low maximum utility, i.e., points whose maximum utility is lower than the  $k$ -th largest minimum utility.

**Lemma 4.4.** For  $v \in V$ , if  $U_{\max}(v) < U_{\min}^k$ , then  $v \notin S^*$ .

---

**Algorithm 4**  $\text{Shrink}(V, S, k)$

---

$\forall v \in V$ : Compute  $U_{\min}(v)$  and  $U_{\max}(v)$ .  
 Threshold  $\leftarrow U_{\min}^k$   
**return**  $\{v \in V \mid U_{\max}(v) \geq \text{Threshold}\}$

---

As described in Algorithm 5, we use grow and shrink blocks, and repeat each block until convergence before we alternate to the other one. We may be able to grow more points after repeated shrink steps, or shrink more points after repeated grow steps since the grow step will decrease the maximum utility of some points, and the shrink step will increase the minimum utility of some points. We give a visual intuition in Figure 1.

---

**Algorithm 5**  $\text{Bounding}(V, k)$

---

$S' \leftarrow \emptyset$   
**repeat**  
     **repeat**  $V \leftarrow \text{Reduce}(V, S', k)$  **until**  $V$  **converges**  
     **repeat**  $V, S', k \leftarrow \text{Grow}(V, S', k)$  **until**  $S'$  **converges**  
**until**  $S'$  and  $V$  **converge**  
**return**  $V, S', k$

---

## 4.2. Approximate Bounding

While the exact bounding algorithm can provide optimal subsets, even better than the centralized greedy algorithm, it is not guaranteed to find the complete subset. In practice, we find that the algorithm converges very quickly yielding very incomplete subsets (Section 5.2). This is due to the strict condition to only discard and add points when we are absolutely certain. In this subsection, we relax this constraint, and present an approximate bounding algorithm using the notion of expected utility as defined below.

**Definition 4.5** (Expected Utility). For any point  $v_1 \in V$ , its *expected utility* considers only a subset of neighbors:  $U_{\text{exp}}(v_1) = u(v_1) - \frac{\beta}{\alpha} \sum_{v_2 \in U(N_{v_1}) \cup S'} s(v_1, v_2)$ , where  $U(N_{v_1})$  denotes the set obtained by sampling the neighbors of  $v_1$  from the ground set, either uniformly or weighted based on the similarity values.

Note that we always consider all neighbors in the current partial solution  $S'$ , and, in our implementation (Section 5), only sample if the number of neighbors in  $S'$  is too small. With the expected utility, the approximate bounding algorithm works the same as the exact bounding where we replace the minimum utility with the expected utility. Replacing the minimum utility with the expected utility intuitively leads to removal of a larger number of points from the ground set and the addition of a larger number of points in the subset. While approximate bounding is not guaranteed to find subsets of optimal quality as in the case of exact bounding, we show provable guarantees on the quality of the subset with respect to the centralized greedy solution in the next subsection. Note that this algorithm might grow  $S'$  larger than needed, in which case we sample a subset of the correct size uniformly at random.

## 4.3. Theoretical Analysis of Approximate Bounding

In this section we provide theoretical justification for the approximate bounding procedure. First notice the by performing *exact* bounding followed by running the centralized greedy submodular procedure at the end at least get a  $\frac{1}{2}$ -approximation guarantee since the exact bounding procedure never discards elements from the optimal set  $S^*$ . For this theoretical analysis, we sample the neighborhood uniformly at random. This means for each  $v \in V$ , each vertex

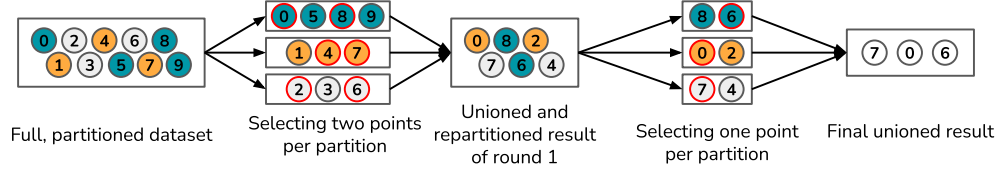


Figure 2. Visualization of the distributed submodular algorithm finding a subset of size 3 out of 10 points using 2 rounds with 3 partitions. The partitioning is given by color; the selected points per partition are marked with a red border, and the numbers represent IDs.

in the neighborhood  $N_v$  is chosen independently at random with probability  $p$  when calculating  $U_{\min}(v)$ .

**Theorem 4.6.** *Let the given input instance be such that each non-zero similarity value lies in the range  $[a, b]$ . Furthermore, let  $\gamma > 1$  be such that at the beginning of the algorithm, for each  $v \in V$  it holds that  $\frac{U_{\max}(v)}{U_{\min}(v)} \leq \gamma$ . If random sampling with probability  $p$  is used for approximate bounding then with probability at least  $1 - |V|e^{-p^2(1-p)^2k\frac{a^2}{(b-a)^2}}$  the algorithm outputs a set  $S$  of size  $k$  such that*

$$f(S) \geq \frac{1}{2(1 + \gamma(1 - p^2))} f(S^*). \quad (1)$$

The approximation guarantee becomes better as  $p$  increases, which is expected. For  $p = 1$  we recover the  $\frac{1}{2}$ -guarantee of exact bounding. The guarantee also depends on the ratio of the maximum and the minimum utilities which in turn is related to the balancing parameters  $\alpha, \beta$ . In the worst case, when this ratio becomes infinity, we get a vacuous bound which is expected since we use approximate bounding. We refer to [Appendix B](#) for the proof.

#### 4.4. Distributed Greedy Algorithm

The bounding algorithm is not always guaranteed to provide the complete subset, and in such cases, we use a partition-based distributed algorithm to compute the rest of the points in the subset. Prior distributed algorithms first partition the entire dataset, run centralized greedy in each of these partitions, and typically use a final step where a centralized greedy is run on the union of the individual subsets computed in each of the partitions (c.f. [Section 2](#)). However, this final step of running the centralized greedy on the union is infeasible when the size of the subset is large. We skip this final subset selection on the union, and instead produce smaller subsets in each of the partitions whose union directly lead to the desired subset. A union can be implemented without materializing all data in memory by data processing frameworks such as Spark ([Zaharia et al., 2016](#)), Flume ([Chambers et al., 2010](#)), or Beam ([Akidau et al., 2015](#)). To produce high quality subset, we employ several rounds where we iterate over finding subsets in different partitions, and partitioning the union of the computed subsets for the next round.

We show an example in [Figure 2](#) for the distributed greedy

**Algorithm 6** *Distributed greedy using adaptive partitioning over  $m$  machines for  $r$  rounds to find subset  $S$  of size  $k$ .*

---

```

partition_cap  $\leftarrow \lceil \frac{|V|}{m} \rceil$ 
 $V_0 \leftarrow V$ 
for round = 1  $\dots$   $r$  do
   $n_{\text{round}} \leftarrow \Delta(|V|, r, \text{round}, k)$ 
   $m_{\text{round}} \leftarrow \lceil \frac{n_{\text{round}}}{\text{partition\_cap}} \rceil$ 
   $P_1, \dots, P_{m_{\text{round}}} \leftarrow$  random partition of  $V_{\text{round}-1}$ 
  for each partition  $i = 1 \dots m_{\text{round}}$  in parallel do
     $S_i \leftarrow$  Centralized Greedy (Alg. 2) on  $P_i$ 
  end for
   $V_{\text{round}} \leftarrow \bigcup_{i=1 \dots m_{\text{round}}} S_i$ 
end for
 $S \leftarrow$  subsample( $V_r, k$ ) if  $|V_r| > k$  else  $V_r$ 
    
```

---

algorithm outlined in [Algorithm 6](#). Given the initial dataset size  $|V|$ , the overall number of rounds  $r$ , the current round, and the target size of the last round  $k$ , the function  $\Delta$  gives us the number of data points to keep in this round. This can, for example, be a linear interpolation between  $|V|$  at the beginning and  $k$  in the last round. However, many choices of  $\Delta$  are possible and the only constraint is that  $\Delta(|V|, r, r, k) = k$  to ensure we output a subset of the correct size.

After determining the target size  $n_r$  for the current round, we determine how many partitions, and therefore nodes, to use this round. We differentiate two variants of the algorithm. In [Algorithm 6](#), we show our adaptive partitioning algorithm where we scale the number of partitions based on how many nodes we actually need to fit the current dataset  $V_{\text{round}}$ . The intuition behind adaptive partitioning is that using more partitions will perform worse due to less global information, and hence we want to utilize the minimum number of partitions needed. We can also disable adaptive partitioning, in which the number of partitions each round  $m_r$  is always equal to the number of machines  $m$  available at the beginning. In this paper, we limit ourselves to partitioning uniform at random. Then, in parallel on each node, we start the centralized algorithm on its assigned partition, with the target to find a dataset of size  $\lceil n_r/m_r \rceil$ . We discard any neighborhood relation across partitions while computing the submodular objective. After all partitions have been processed, we union the results of each partition.

We repeat this procedure for  $r$  rounds. At the end, the resulting dataset  $V_r$  might contain up to  $m_r$  additional points due to rounding. Hence, we obtain  $S$  by subsampling  $k$  points from  $V_r$ . In Section 5, we show that sufficient rounds yield high quality subsets in practice, close to centralized greedy methods, despite not using centralized greedy on the final union of computed subsets, as done in prior methods. Implementation details can be found in Appendix C.

## 5. Evaluation

We demonstrate that our distributed algorithms achieve similar quality subsets as the centralized greedy algorithm. We evaluate this by comparing them to objective obtained using the centralized algorithm.

**Datasets.** We use three datasets: (1) CIFAR-100 (Krizhevsky & Hinton, 2009) with 100 classes and 50k points, (2) ImageNet (Deng et al., 2009) with 1k classes and ca. 1.2 M points, and (3) Perturbed-ImageNet with 1k classes and ca. 13 B images. We obtain Perturbed-ImageNet by perturbing each point of ImageNet in embedding space into 10k vectors, leading to 13 B embedding vectors for stress-testing our distributed algorithm.

**Embeddings, similarities, and utilities.** For both CIFAR-100 and Imagenet, we generate predictions and embeddings for all points using a coarsely-trained ResNet-56 (He et al., 2016) trained on a random 10% subset. We employ SGD with Nesterov momentum 0.9, using 450/90 epochs for CIFAR/ImageNet. The base learning rate is 1.0 for CIFAR and 0.1 for ImageNet, and is reduced by a tenth at epochs 15, 200, 300, 400 (CIFAR) and 5, 30, 69, 80 (ImageNet). We extract the penultimate layer features to generate 64-dimensional embeddings for CIFAR, and 2048-dimensional embeddings for ImageNet.

We build a 10-nearest neighbor graph  $(G, E)$  in the embedding space, using the fast similarity search by Guo et al. (2020). Note that this nearest neighbor graph is not symmetric. As bounding and scoring require a symmetric graph, we symmetrize the graph, such that data-points have a varying amount of, but at least 10 neighbors. This leads to an average number of 15/16 neighbors for CIFAR-100/ImageNet, respectively. Based on this graph, we set  $s(v_1, v_2)$  in the submodular objective  $f(S) = \alpha \sum_{v \in S} u(v) - \beta \sum_{(v_1, v_2) \in E; v_1, v_2 \in S} s(v_1, v_2)$  to the cosine similarities of neighboring points. We set  $\beta = 1 - \alpha$  and only mention the value of  $\alpha$  subsequently.

We use margin-based uncertainty (Scheffer et al., 2001) as the utility. In multi-class settings, margin provides high utility for points that are hard to classify, i.e., close to the decision boundaries:  $u(x_i) = 1 - (P(\text{top} | x_i) - P(\text{sec} | x_i))$  where  $P(c | x_i)$  denotes the probability for class  $c$  predicted by the model for the example  $x_i$ , and *best* and *sec* denote

the best and the second best classes as per the predictions. Intuitively, a coarse model has already learned to classify easy points, while uncertain points are harder to learn and, therefore, more important. We center the utilities by subtracting the minimum utility from all values.

**Normalization.** For each dataset, to compare the algorithms, for the same parameter group (dataset,  $\alpha/\beta$ , and target subset size  $k$ ), we map the objective from the centralized greedy to 100%, and the lowest observed score to 0%. This normalization enables us to uniformly interpret a percent point as a gain over the worst case, and emphasize instances where the basic centralized score is exceeded.

### 5.1. Distributed Greedy Algorithm

We fix the delta function to linear interpolation with a factor of 0.75, i.e.,  $\Delta(|V|, r, \text{round}, k) = \lceil 0.75 \cdot (r - \text{round}) \cdot \frac{|V| - k}{r} \rceil + k$ . For an ablation on the interpolation, we refer to Appendix E.

**Multiple partitions and multiple rounds.** Figure 3 shows the normalized scores without adaptive partitioning depending on the subset size and  $\alpha$  resp.  $\beta$  on the CIFAR-100 dataset. In all settings, (i) fewer partitions lead to a higher score and (ii) more rounds lead to a higher score. The reason is that when we partition the data into many sets, the neighborhood relation is lost across partitions. We refer to Appendix D for a visualization. Nevertheless, we still obtain high quality subsets by using multiple rounds. For CIFAR, 2 partitions / 1 round have a 80% score, but with 32 rounds, we obtain 98% score. Similarly, on ImageNet, the score increases from 86% to 98%. We note that 32 partitions / 1 round does not have a score of 0 because the random partitioning in the adaptive partitioning experiment (same parameter group) leads to a slightly lower score.

For all settings, we find that the use of multiple rounds is more effective when the budget  $k$  is smaller. For example, for CIFAR and  $\alpha = 0.9$ , when finding a 10% subset, for 16 partitions, going from 1 to 32 rounds increases the relative score from 15% to 74%, while for a 50% subset, it only increases the relative score from 8% to 18%. For the other values of  $\alpha$ , we find a similar trend. For larger target sizes, the steps per round are smaller, such that there is less potential for improvement from repartitioning.

*We observe that more partitions generally decrease and more rounds increase the submodular objective, especially for smaller subsets.*

**Adaptive partitioning.** In Figure 4, we show that adaptive partitioning increases the score drastically since each round, we get closer to the centralized version and lose fewer edges in the neighbor graph. The benefit of adaptivity is higher the smaller the target subset size is, since for smaller subsets, fewer partitions are necessary. For example, for 10% sub-

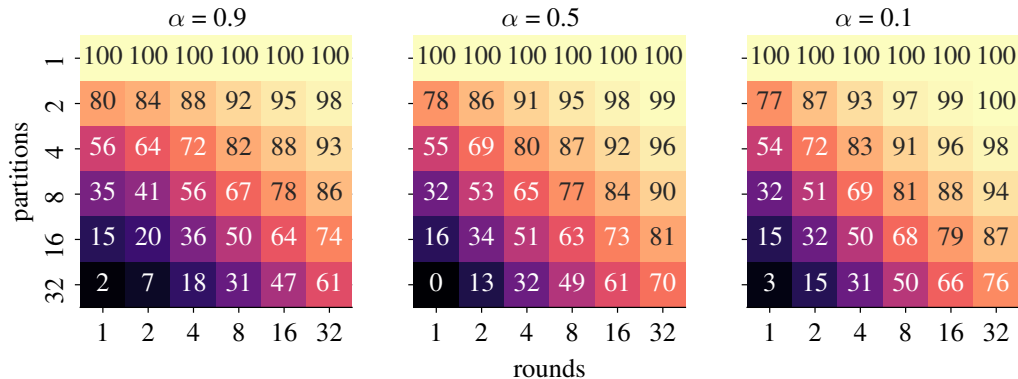


Figure 3. Normalized scores for finding a 10% subset on CIFAR-100, depending on the number of partitions, rounds, and  $\alpha$ . The full version can be found in Figure 12 (CIFAR) and Figure 13 (ImageNet). Here, 100 denotes the quality of centralized greedy algorithm.

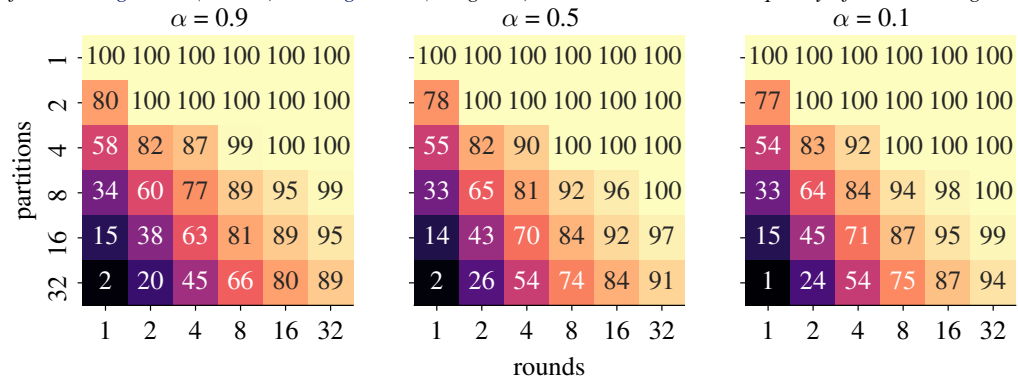


Figure 4. Normalized scores for finding a 10% subset on CIFAR-100, depending on the number of partitions, rounds, and  $\alpha$ , using adaptive partitioning. The full version can be found in Figure 14 (CIFAR) and Figure 15 (ImageNet).

sets, for both CIFAR and ImageNet, even with 32 partitions, we obtain a score of around 90%, whereas without adaptive partitioning, the highest scores were 60 to 75%. Furthermore, using adaptive partitioning is less resource-intensive, since it requires fewer parallel machines.

*Adaptive partitioning often reaches the quality of centralized method, and uses fewer parallel resources.*

## 5.2. Bounding

Table 1 shows the results for  $\alpha = 0.9$ . To evaluate each configuration, we investigate the number of included/excluded datapoints, and the centralized score (1 partition/1 round) when processing the bounding result with the greedy algorithm. For other configurations, see Appendix G.

**Exact bounding.** For CIFAR (50k points), in the 10% subset case, bounding excludes ca. 22% of the points in 16 rounds. In the 50% subset case, it does not make any decision, and for the 80% subset, it includes 4% in 9 rounds. For ImageNet (1.2M points), it does not make a decision for 10% and 50% subsets, and includes less than 1% of the dataset in 5 rounds. For both datasets, in the 80% subset case, the final score is slightly higher than the centralized non-bounding score, just by including a few data points.

Overall, exact bounding only includes or excludes a few points for very small or very large subsets, respectively. We observe this behavior because a smaller target size leads to easier requirements for shrinking, which results in more excluded points (Algorithm 3). The same argument holds for large target sizes and inclusion. Intuitively, when the target subset size is extreme, it is easier for the algorithm to make a decision, since overall more points are included/excluded.

**Approximate bounding.** We test approximate bounding with uniform sampling, i.e., all neighbors of a data point have the same chance of being sampled, and weighted sampling, i.e., the sampling probability is uniform to the pairwise interaction between the neighbors. For both sampling types, we test sampling a 30% and 70% neighborhood. Generally, for both CIFAR and ImageNet, the 30% neighborhood is able to include and exclude many points. However, the 70% neighborhood still struggles to make decisions, especially for the 50% subset setting. We also find that the larger the neighborhood, i.e., the more information we use, the higher the score is.

In the 10% subset setting, for both CIFAR and ImageNet, when sampling a 30% neighborhood, the algorithm excludes around 50% of the dataset (22% and < 1% for exact bounding on CIFAR and ImageNet), reducing the load

Table 1. Bounding results for  $\alpha = 0.9$ . Each cell gives the number of included / excluded points, the number of grow / reduce rounds, and the centralized score after running the centralized greedy algorithm on the output. We achieve high quality results for most settings with uniform sampling, and occasionally even outperforming centralized methods.

Algorithm	CIFAR-100			ImageNet		
	10 % Subset	50 % Subset	80 % Subset	10 % Subset	50 % Subset	80 % Subset
Regular	0 / 10769	0 / 0	2002 / 0	0 / 0	0 / 0	8043 / 0
	1 / 16	1 / 1	9 / 2	1 / 1	1 / 1	5 / 2
30 % Uniform Sampling	100.01 %	100.0 %	100.55 %	100.0 %	100.0 %	100.1 %
	3 / 25743	24972 / 14896	39999 / 0	16 / 756625	633027 / 274635	1024932 / 0
70 % Uniform Sampling	2 / 10	9 / 7	4 / 2	3 / 10	34 / 12	5 / 2
	100.0 %	97.39 %	85.95 %	100.01 %	99.77 %	92.72 %
30 % Weighted Sampling	0 / 17785	0 / 0	30724 / 0	0 / 470133	0 / 0	192019 / 0
	1 / 9	1 / 1	165 / 2	1 / 14	1 / 1	189 / 2
70 % Weighted Sampling	100.06 %	100.0 %	103.51 %	100.0 %	100.0 %	102.13 %
	3 / 25729	24996 / 14853	39999 / 0	16 / 756511	639492 / 273628	1024934 / 0
30 % Weighted Sampling	2 / 8	6 / 9	4 / 2	2 / 10	9 / 11	4 / 1
	100.0 %	75.59 %	81.79 %	100.01 %	81.95 %	78.52 %
70 % Weighted Sampling	0 / 17748	0 / 0	40000 / 0	0 / 469860	0 / 0	1024904 / 0
	1 / 9	1 / 1	6 / 1	1 / 14	1 / 1	12 / 2
	100.06 %	100.0 %	88.17 %	100.0 %	100.0 %	95.73 %

on the centralized algorithm and requiring fewer partitions at the start. In the 80 % subset setting, the algorithm often finds (almost) the entire subset without the need for the greedy algorithm. Interestingly, it struggles to do so with 70 % uniform sampling, where for ImageNet, it includes ten times fewer points and takes 189 rounds, compared to five to twelve rounds in the other settings. Weighted sampling helps the algorithm to converge faster, but sometimes has a worse score than uniform sampling, which we attribute to the bias in the sampling strategy that not always is optimal. Last, for 30 % neighborhoods, the algorithm is even able to both include and exclude data points in the 50 % subset setting, which is something that exact bounding and even the 70 % neighborhood struggle with. Approximate sampling empirically performs well, keeping scores of over 90 %, and is implementable in a parallel setting.

In summary, approximate bounding includes and excludes many more points than exact bounding at the cost of a sometimes slightly lower objective. For large subsets, bounding often finds the entire subset on its own. However, as it is massively parallelizable, it still performs exceptionally well from a runtime-objective tradeoff standpoint.

**Lowering the utility coefficient  $\alpha$ .** For  $\alpha \in \{0.1, 0.5\}$  we find that no configuration includes/excludes any points. The smaller values of  $\alpha$  increase the role of the pairwise terms. This makes the minimum and maximum utilities drift apart, leading to less growing/shrinking of points.

*Bounding often performs better, and occasionally even better than the centralized method, when the utility dominates the objective function.*

### 5.3. Dataset with 13 Billion Points

We demonstrate the scalability of our approach on a dataset with around 13 billion datapoints. For the distributed sub-

modular algorithm, we use 16 partitions with 350 GB of memory per partition. We test 1, 2, and 8 rounds for  $\alpha = 0.9$ . We cannot normalize the scores since we are unable to determine the centralized score.

**10 % subset.** The score increases from 1 058 841 312 (1 round) to 1 092 474 410 (2 rounds), up to 1 145 682 717 (8 rounds). We also test the bounding procedure. Exact bounding includes 0.007 % of the 13 billion points and excludes 10 %. Approximate bounding with a 30 % neighborhood includes 0.7 % of the points and excludes 60 % for both uniform and weighted sampling. All bounding approaches reach a score of slightly above 100 % of the score without bounding after 8 rounds.

**50 % subset.** For the distributed greedy algorithm, the score increases from 4 168 989 874 (1 round), to 4 200 071 672 (2 rounds), up to 4 250 047 523 (8 rounds).

## 6. Discussion

We present a distributed bounding algorithm for optimizing pairwise submodular functions for subset selection, which does not assume the subset fits in the main memory of a single machine. We also propose a multi-round, partitioning-based, distributed algorithm that yields high-quality subsets with minimal to no quality loss compared to centralized algorithms. We conduct a comprehensive analysis of both the bounding and the distributed greedy algorithm in various settings and providing insights on achieving high quality subsets. We find that using more rounds help to reach near-centralized scores, especially when adaptively adjusting the number of parallel partitions. In some cases, approximate bounding is able to exclude over 50 % of the dataset, and sometimes solves the entire problem on its own.

Despite proving guarantees for the bounding algorithm, showing approximation guarantees for the distributed



greedy algorithm remains future work. Occasionally, we observe that the approximate bounding can even outperform centralized methods, and we plan to explore this further.

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## Acknowledgments

We would like to thank Sameer Agarwal, Sanjiv Kumar, Christian Tjandraatmadja, Ayan Chakrabarti, Daniel Glasner, and Morteza Zadimoghaddam for their valuable help.

## References

- Akidau, T., Bradshaw, R., Chambers, C., Chernyak, S., Fernández-Moctezuma, R. J., Lax, R., McVeety, S., Mills, D., Perry, F., Schmidt, E., and Whittle, S. The dataflow model: A practical approach to balancing correctness, latency, and cost in massive-scale, unbounded, out-of-order data processing. *Proceedings of the VLDB Endowment*, 8 (12):1792–1803, 2015. doi: 10.14778/2824032.2824076.
- Badanidiyuru, A., Mirzasoleiman, B., Karbasi, A., and Krause, A. Streaming submodular maximization: massive data summarization on the fly. In *Proceedings of the Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2014. doi: 10.1145/2623330.2623637.
- Barbosa, R., Ene, A., Nguyen, H., and Ward, J. The power of randomization: Distributed submodular maximization on massive datasets. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2015.
- Bhardwaj, R., Xia, Z., Ananthanarayanan, G., Jiang, J., Shu, Y., Karianakis, N., Hsieh, K., Bahl, P., and Stoica, I. Ekya: Continuous learning of video analytics models on edge compute servers. In *Proceedings of the Symposium on Networked Systems Design and Implementation (NSDI)*, 2022.
- Böther, M., Gsteiger, V., Robroek, T., and Klimovic, A. Modyn: A platform for model training on dynamic datasets with sample-level data selection. 2023. doi: 10.48550/arXiv.2312.06254.
- Carbonell, J. and Goldstein, J. The use of mmr, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of the Conference on Research and Development in Information Retrieval (SIGIR)*, 1998. doi: 10.1145/290941.291025.
- Chambers, C., Raniwala, A., Perry, F., Adams, S., Henry, R. R., Bradshaw, R., and Weizenbaum, N. FlumeJava: Easy, efficient data-parallel pipelines. In *Proceedings of the Conference on Programming Language Design and Implementation (PLDI)*, 2010. doi: 10.1145/1806596.1806638.
- Das, A. and Kempe, D. Submodular meets spectral: greedy algorithms for subset selection, sparse approximation and dictionary selection. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2011.
- Dean, J. and Ghemawat, S. Mapreduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1), 2008. doi: 10.1145/1327452.1327492.
- Deng, J., Dong, W., Socher, R., Li, L., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- Ene, A., Im, S., and Moseley, B. Fast clustering using mapreduce. In *Proceedings of the International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2011. doi: 10.1145/2020408.2020515.
- Gao, L., Biderman, S., Black, S., Golding, L., Hoppe, T., Foster, C., Phang, J., He, H., Thite, A., Nabeshima, N., Presser, S., and Leahy, C. The pile: An 800gb dataset of diverse text for language modeling, 2020.
- Golovin, D. and Krause, A. Adaptive submodularity: Theory and applications in active learning and stochastic optimization. *Journal of Artificial Intelligence Research*, 42(1), 2011.
- Guo, R., Sun, P., Lindgren, E., Geng, Q., Simcha, D., Chern, F., and Kumar, S. Accelerating large-scale inference with anisotropic vector quantization. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2020.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. doi: 10.1109/CVPR.2016.90.
- Im, S. and Moseley, B. Fast and better distributed mapreduce algorithms for k-center clustering. In *Proceedings of the Symposium on Parallelism in Algorithms and Architectures (SPAA)*, 2015. doi: 10.1145/2755573.2755607.
- Indyk, P., Mahabadi, S., Mahdian, M., and Mirrokni, V. S. Composable core-sets for diversity and coverage maximization. In *Proceedings of the Symposium on Principles of Database Systems (PODS)*, 2014. doi: 10.1145/2594538.2594560.

- Katharopoulos, A. and Fleuret, F. Not all samples are created equal: Deep learning with importance sampling. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.
- Kaushal, V., Iyer, R., Kothawade, S., Mahadev, R., Doctor, K., and Ramakrishnan, G. Learning from less data: A unified data subset selection and active learning framework for computer vision. In *Proceedings of the Winter Conference on Applications of Computer Vision (WACV)*, 2019. doi: 10.1109/wacv.2019.00142.
- Kazhamiaka, F., Zaharia, M., and Bailis, P. Challenges and opportunities for autonomous vehicle query systems. In *Proceedings of the Conference on Innovative Data Systems Research (CIDR)*, 2021.
- Kim, B., Koyejo, O., and Khanna, R. Examples are not enough, learn to criticize! criticism for interpretability. In *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS)*, 2016.
- Kolmogorov, V. and Zabini, R. What energy functions can be minimized via graph cuts? *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(2), 2004. doi: 10.1109/tpami.2004.1262177.
- Kothawade, S., Beck, N., Killamsetty, K., and Iyer, R. K. SIMILAR: submodular information measures based active learning in realistic scenarios. In *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- Kothawade, S., Kaushal, V., Ramakrishnan, G., Bilmes, J., and Iyer, R. PRISM: a rich class of parameterized submodular information measures for guided data subset selection. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(9), 2022. doi: 10.1609/aaai.v36i9.21264.
- Krause, A. and Cevher, V. Submodular dictionary selection for sparse representation. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2010.
- Krause, A. and Golovin, D. Submodular function maximization, 2014.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. Technical report, University of Toronto, Toronto, Ontario, 2009.
- Kumar, R., Moseley, B., Vassilvitskii, S., and Vattani, A. Fast greedy algorithms in mapreduce and streaming. *ACM Transactions on Parallel Computing*, 2(3), 2015. doi: 10.1145/2809814.
- Lin, H. and Bilmes, J. A. A class of submodular functions for document summarization. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics: Human Language Technologies (HLT)*, 2011.
- Liu, P. and Vondrák, J. Submodular optimization in the mapreduce model. In *Proceedings of the Symposium on Simplicity in Algorithms (SOSA)*, 2018. doi: 10.4230/OASICS.SOSA.2019.18.
- Malkomes, G., Kusner, M. J., Chen, W., Weinberger, K. Q., and Moseley, B. Fast distributed k-center clustering with outliers on massive data. In *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS)*, 2015.
- McClintock, J. and Wirth, A. Efficient parallel algorithms for k-center clustering. In *Proceedings of the International Conference on Parallel Processing (ICPP)*, 2016. doi: 10.1109/icpp.2016.22.
- Mirroknj, V. S. and Zadimoghaddam, M. Randomized composable core-sets for distributed submodular maximization. In *Proceedings of the Symposium on Theory of Computing (STOC)*, 2015. doi: 10.1145/2746539.2746624.
- Mirzasoleiman, B., Karbasi, A., Sarkar, R., and Krause, A. Distributed submodular maximization. *Journal of Machine Learning Research*, 17(235), 2016.
- Mirzasoleiman, B., Bilmes, J. A., and Leskovec, J. Coresets for data-efficient training of machine learning models. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2020.
- Nemhauser, G. L., Wolsey, L. A., and Fisher, M. L. An analysis of approximations for maximizing submodular set functions—i. *Mathematical Programming*, 14(1):265–294, 1978. doi: 10.1007/bf01588971.
- Park, H.-S. and Jun, C.-H. A simple and fast algorithm for k-medoids clustering. *Expert Systems with Applications*, 36(2), 2009. doi: 10.1016/j.eswa.2008.01.039.
- Pooladzandi, O., Davini, D., and Mirzasoleiman, B. Adaptive second order coresets for data-efficient machine learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2022.
- Prasad, A., Jegelka, S., and Batra, D. Submodular meets structured: Finding diverse subsets in exponentially-large structured item sets. In *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS)*, 2014.
- Ramalingam, S., Glasner, D., Patel, K., Vemulapalli, R., Jayasumana, S., and Kumar, S. Less is more: Selecting informative and diverse subsets with balancing constraints. 2021. doi: 10.48550/arXiv.2104.12835.
- Ramalingam, S., Awasthi, P., and Kumar, S. A weighted k-center algorithm for data subset selection. 2023. doi: /10.48550/arXiv.2312.10602.

- Scheffer, T., Decomain, C., and Wrobel, S. Active hidden markov models for information extraction. In *Advances in Intelligent Data Analysis (IDA)*, 2001. doi: 10.1007/3-540-44816-0\\_31.
- Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., Schramowski, P., Kundurthy, S., Crowson, K., Schmidt, L., Kaczmarczyk, R., and Jitsev, J. LAION-5B: an open large-scale dataset for training next generation image-text models. In *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- Shen, Z., Tao, T., Ma, L., Neiswanger, W., Liu, Z., Wang, H., Tan, B., Hestness, J., Vassilieva, N., Soboleva, D., and Xing, E. SlimPajama-DC: understanding data combinations for llm training. 2023. doi: 10.48550/arXiv.2309.10818.
- Simon, I., Snavely, N., and Seitz, S. M. Scene summarization for online image collections. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2007. doi: 10.1109/iccv.2007.4408863.
- Wei, K., Iyer, R. K., and Bilmes, J. A. Submodularity in data subset selection and active learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2015.
- Zaharia, M., Xin, R. S., Wendell, P., Das, T., Armbrust, M., Dave, A., Meng, X., Rosen, J., Venkataraman, S., Franklin, M. J., Ghodsi, A., Gonzalez, J., Shenker, S., and Stoica, I. Apache spark: a unified engine for big data processing. *Communications of the ACM*, 59(11):56–65, 2016. doi: 10.1145/2934664.

## A. Submodularity and Monotonicity

**Definition A.1** (Submodularity). Let  $\Omega$  be a finite set. A set function  $F : 2^\Omega \rightarrow \mathbb{R}$  is *submodular* if for all  $A, B \subseteq \Omega$  with  $B \subseteq A$  and  $e \in \Omega \setminus A$ ,  $F(A \cup \{e\}) - F(A) \leq F(B \cup \{e\}) - F(B)$ .

This property is also referred to as *diminishing returns* since the gain diminishes as we add elements.

**Definition A.2** (Monotonicity). A set function  $F$  is *monotonically non-decreasing* if for  $B \subseteq A$ ,  $F(B) \leq F(A)$ .

## B. Theoretical Analysis of Approximate Bounding

In this section, we provide the proof for [Theorem 4.6](#).

*Proof.* We first bound the performance of the algorithm with respect to a generic approximate algorithm for computing  $U_{\min}(v)$  and then instantiate the selection procedure to be random i.i.d. sampling to get the desired final bound. For any set  $S$ , the objective value at  $S$  is

$$f(S) = \sum_{v \in S} u(v) - \frac{\beta}{\alpha} \sum_{v_1, v_2 \in S} s(v_1, v_2). \quad (2)$$

Alternately, we can write the above as

$$f(S) = \sum_{v \in S} \tilde{u}(v, S), \quad (3)$$

where  $\tilde{u}(v, S)$  is defined as

$$\tilde{u}(v, S) = u(v) - \frac{\beta}{2\alpha} \sum_{v_2 \in S} s(v, v_2). \quad (4)$$

Next consider a generic procedure for computing  $U_{\min}$  values that outputs values  $\tilde{U}_{\min}(v)$  such that  $U_{\min}(v) \leq \tilde{U}_{\min}(v) \leq rU_{\min}(v)$  holds for all  $v \in V$ , for some  $r > 1$ . Next consider a surrogate objective  $\hat{f}(S)$  defined as

$$\hat{f}(S) = \sum_{v \in S} \hat{u}(v, S) \quad (5)$$

$$=: \sum_{v \in S} \max(\tilde{u}(v, S), \tilde{U}_{\min}(v)). \quad (6)$$

Notice that the surrogate objective has the property that for any set  $S$ ,  $f(S) \leq \hat{f}(S) \leq rf(S)$ . Furthermore, with respect to the surrogate objective the grow and reduce operations are performing exact bounding as no element of the optimal solution will be discarded. As a result, the algorithm will output a set  $S$  such that

$$f(S) \geq \frac{1}{2r} f(S^*). \quad (7)$$

Finally, it remains to bound the value of  $r$ . For this, we use the fact that we sample the neighbors uniformly at random with probability  $p$ . Consider a vertex  $v$  and let  $\mu(v) = \sum_{v_1 \in N_v} s(v, v_1)$ . If we select each element with probability  $p$ , from Chernoff bounds (and a union bound over all the elements), we get that for each  $v \in V$ , the sum included in the selected elements will be at least  $(1 - \delta)\mu(v)$  except with probability at most  $e^{-\frac{\delta^2 p^2 k a^2}{(b-a)^2}}$ . Finally, choosing  $\delta = 1 - p$  and using the fact that the ratio of the maximum to minimum utility is bounded by  $\gamma$ , we get that, with high probability,  $r \leq (1 + \gamma(1 - p)^2)$ .  $\square$

## C. Implementation of bounding and scoring

The assumption in this paper is that we cannot hold the entire target subset in memory. In order to implement (i) bounding and (ii) subset scoring under this assumption, we use an Apache Beam-style programming pattern. We assume that the number of neighbor interactions is limited, i.e., there is a (small)  $k \in \mathbb{N}$  s.t. for all  $v \in V$ ,  $|\{v' \in V \mid s(v, v') > 0\}| \leq k$ . In our evaluation (Section 5) this holds true, since we only consider the distance to the  $k$  nearest neighbors. For both bounding and scoring, the difficulty is that when iterating over the interacting neighbors of a data points, we cannot easily perform a  $\mathcal{O}(1)$  check whether the neighbor is in the subset, as the subset is not in memory. Instead, we need to cleverly use distributed joins. We assume that we start with a `PCollection`<sup>1</sup> of (node id, list of neighbors) tuples and a `PCollection` of (node id, utility) values.

**Bounding.** We generate the *fanned-out neighbor graph*, i.e., we iterate over the (node id, neighbor list) tuples, and for each neighbor in the list, emit a triple (neighbor id, node id,  $s(\text{node}, \text{neighbor})$ ). Note that the neighbor id becomes the triple key. We keep track of a `PCollection` for the current partial result with (node, utility) tuples. To get the minimum utility of all yet-unassigned data points, we need to find all neighbors of a point that are either in the partial solution (always considered) or not yet assigned.

To this end, we perform a distributed three-way join of the fanned neighbor graph, the current solution, and the currently unassigned points. Within this join, we filter and invert the neighbor graph: For a node  $a$ , we know whether it is in the partial result by checking whether there is a join partner in that collection. If not, we check if there is a join partner from the currently unassigned points. If not, we discard the edge, since  $a$  is neither in the partial solution nor in the unassigned points, i.e., it has been removed in a shrink step. Now, we iterate over all join partners from the neighbor graph, i.e., nodes  $b$  that were neighbors of  $a$ . We emit 4-tuples of the form  $(b, a, s(a, b), \text{boolean indicating whether } a \text{ is in the partial solution})$ , recovering the original edges before generating the fanned-out neighbor graph, where  $a$  was a neighbor. In these 4-tuples,  $a$  is guaranteed to be not discarded or in the partial solution, i.e., it is relevant for the minimum utility. However,  $b$  might be already discarded or in the partial solution. Hence, we perform another join of the 4-tuples with the unassigned points collection, and discard tuples where there is no join partner. If there is a join partner, we can sample the neighborhood (when using approximate bounding), and use the boolean flag to know whether a neighbor should always be included to determine the minimum utility. We then emit (node, min\_utility) tuples. The maximum utility can be implemented analogously, but only needs to consider the points that are part of the partial solution. The remaining parts of the bounding algorithm can be straightforwardly implemented using map operations.

**Scoring.** Distributed scoring is implemented similarly to the minimum utility calculation. We generate the fanned-out neighbor graph, and then filter it by joining the solution, giving us all neighbors that are part of the solution. We can then invert the result again, and reduce it to have a score per-datapoint, and then reduce this to an overall score by summing up the individual scores, as our function is decomposable.

## D. Subset Visualization

In this section, we visualize the chosen subset on the example of finding a 10% subset of CIFAR-100 with  $\alpha = 0.9$ , depending on the number of partitions for one round. To this end, we take the embeddings of the data points and reduce them to the 2-dimensional plane using t-SNE. We use scikit-learn’s implementation with default settings, i.e., PCA initialization and automatic learning rate. The results are shown in Figure 5. The centralized version chooses data points more uniformly across the plane, while using more partitions creates local clusters. This is because the random partitioning loses information about local edges. The algorithm per partition focuses more on the utility of data points, since it cannot reason about their influence on the diversity score. Overall, this leads to local utility clusters in case of many partitions, whereas the single partition algorithm would have distributed the points more evenly across the plane.

## E. Ablation on $\Delta$

In this section, we perform an ablation study on the delta function (Section 4.4). We use the linear delta function from Section 5, but vary the interpolation factor which we set to  $\gamma = 0.75$  in the evaluation. This means for  $\gamma \in \{0.25, 0.5, 1\}$ , we test  $\Delta(|V|, r, r_{\text{curr}}, k) = \lceil \gamma \cdot (r - r_{\text{curr}}) \cdot \frac{|V| - k}{r} \rceil + k$ . Note that other choices of  $\Delta$  are possible, but we limit ourselves to linear interpolation. We only evaluate the influence of  $\gamma$  on 10% and 50% subsets since the influence of

<sup>1</sup>See the Apache Beam Programming Guide for an introduction to its programming model.

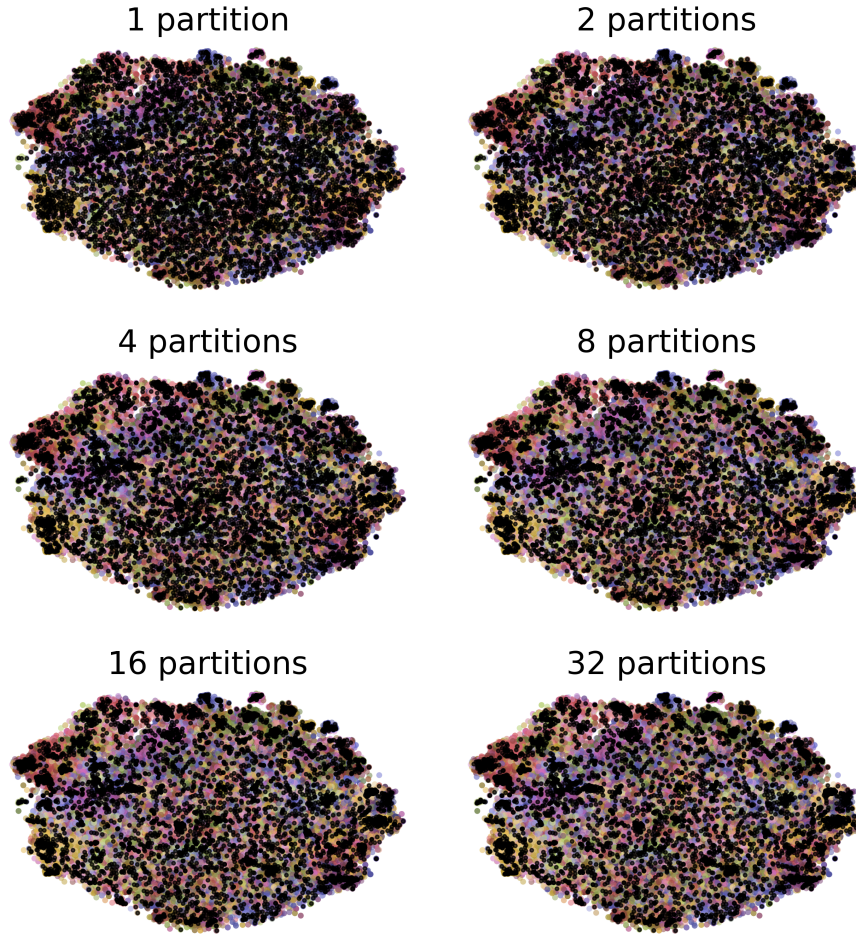


Figure 5. A rasterized visualization of the chosen 5 000 points out of the 50 000 points in CIFAR-100. The points are colored by label, and chosen data points are depicted as black.

the delta function for large subsets is limited.

We first evaluate  $\gamma$  without adaptive partitioning since adaptive partitioning is biased towards smaller values, as they allow for fewer partitions. To ease the comparison, we investigate the difference of the normalized score to the base case of  $\gamma = 0.75$ , i.e., positive values indicate a higher normalized score, and negative values indicate a lower negative score. For CIFAR-100, the results can be found in Figures 6 to 8 and for ImageNet, the results can be found in Figures 9 to 11. Note that decimal places are truncated in the plots.

For  $\gamma = 1$ , i.e., when increasing the intermediate partition sizes, for both CIFAR and ImageNet, the scores are mostly similar or worse than  $\gamma = 0.75$ , across all configurations. There are just very few instances where scores are marginally higher. At the same time, the compute and storage costs are higher due to the larger intermediate partition sizes.

For  $\gamma = 0.5$ , i.e., when decreasing the intermediate partition sizes, we find an increase in scores in many scenarios. This holds especially for  $\alpha = 0.9$ , i.e., when utility is very important. Intuitively, having smaller intermediate partitions forces the algorithm to make inclusion/exclusion decisions earlier. For settings where utility is important, this is beneficial since there is less noise by diverse data points in the selection process. The benefit is higher the more partitions are used. For  $\alpha = 0.1$  the inverse effect happens, as the more partitions are used,  $\gamma = 0.75$  performs better than  $\gamma = 0.5$ . Further lowering  $\gamma$  to 0.25, these effects are amplified.

*For larger values of  $\alpha$ , smaller values of  $\gamma$  are preferred, while smaller values of  $\alpha$  benefit from larger intermediate partitions.*

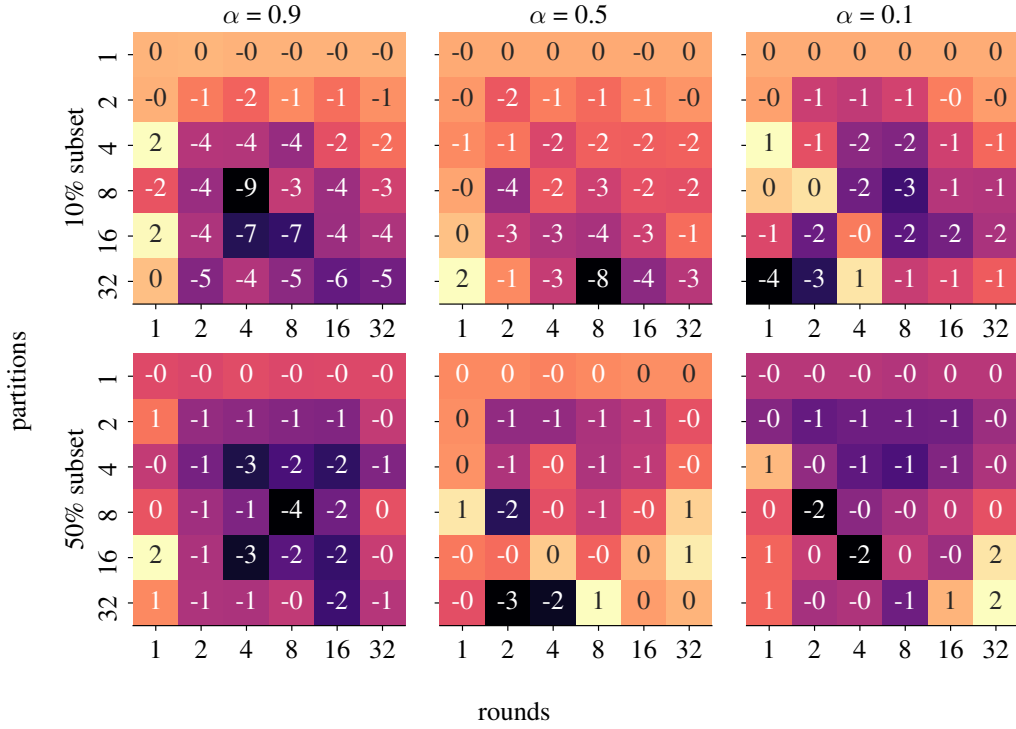


Figure 6. Difference in normalized score of  $\gamma = 1$  to  $\gamma = 0.75$  on CIFAR-100, depending on the subset size, the number of partitions, rounds, and  $\alpha$ .

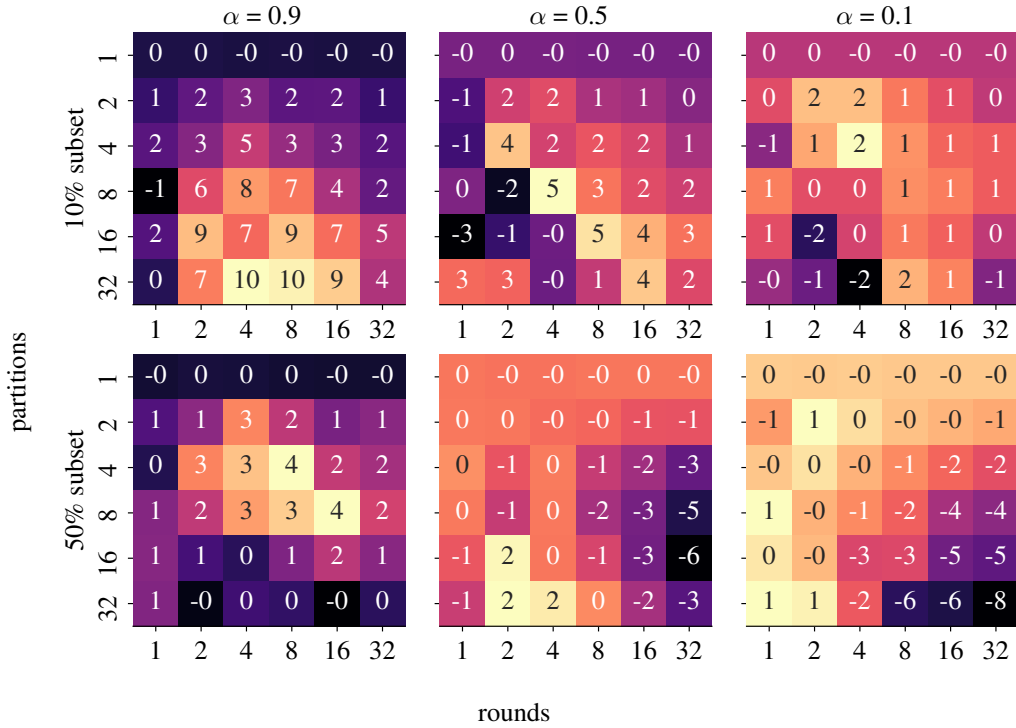


Figure 7. Difference in normalized score of  $\gamma = 0.5$  to  $\gamma = 0.75$  on CIFAR-100, depending on the subset size, the number of partitions, rounds, and  $\alpha$ .

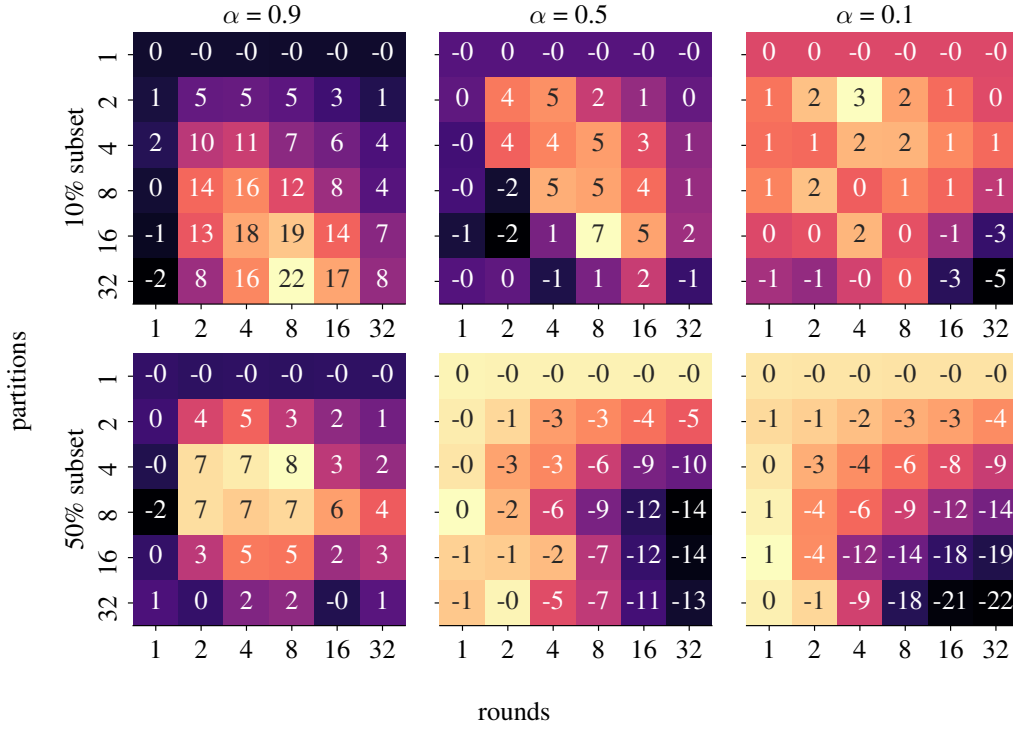


Figure 8. Difference in normalized score of  $\gamma = 0.25$  to  $\gamma = 0.75$  on CIFAR-100, depending on the subset size, the number of partitions, rounds, and  $\alpha$ .

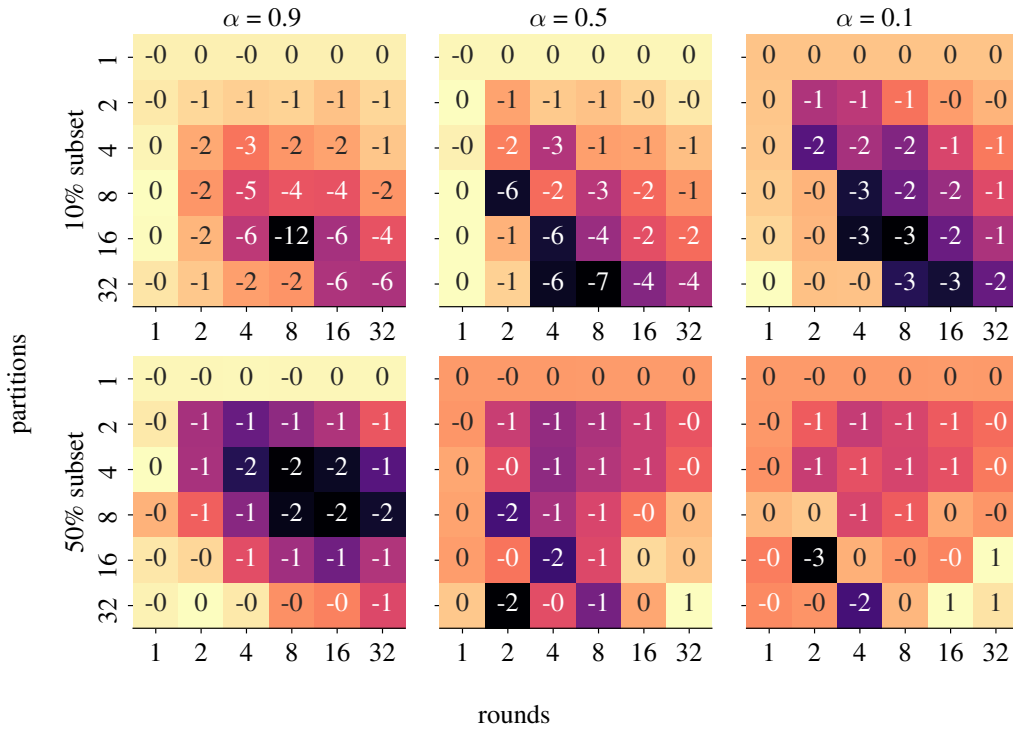


Figure 9. Difference in normalized score of  $\gamma = 1$  to  $\gamma = 0.75$  on ImageNet, depending on the subset size, the number of partitions, rounds, and  $\alpha$ .



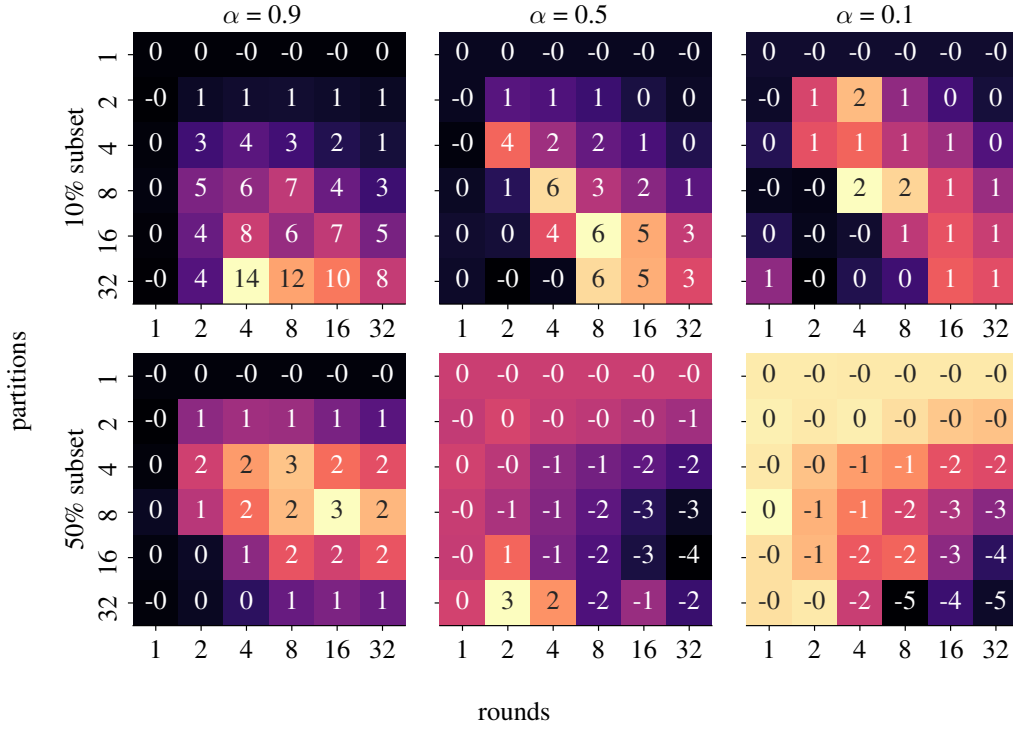


Figure 10. Difference in normalized score of  $\gamma = 0.5$  to  $\gamma = 0.75$  on ImageNet, depending on the subset size, the number of partitions, rounds, and  $\alpha$ .

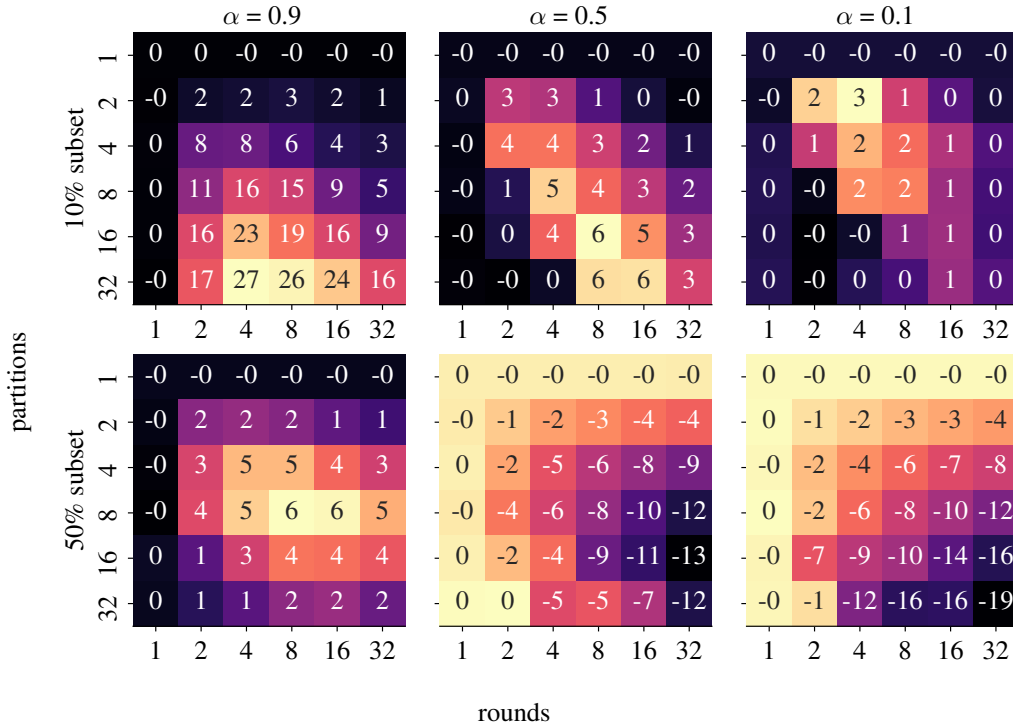


Figure 11. Difference in normalized score of  $\gamma = 0.25$  to  $\gamma = 0.75$  on ImageNet, depending on the subset size, the number of partitions, rounds, and  $\alpha$ .

### F. Additional Figures (without bounding)

In this section, we give the full versions and additional data omitted in the main body of the paper due to space constraints.

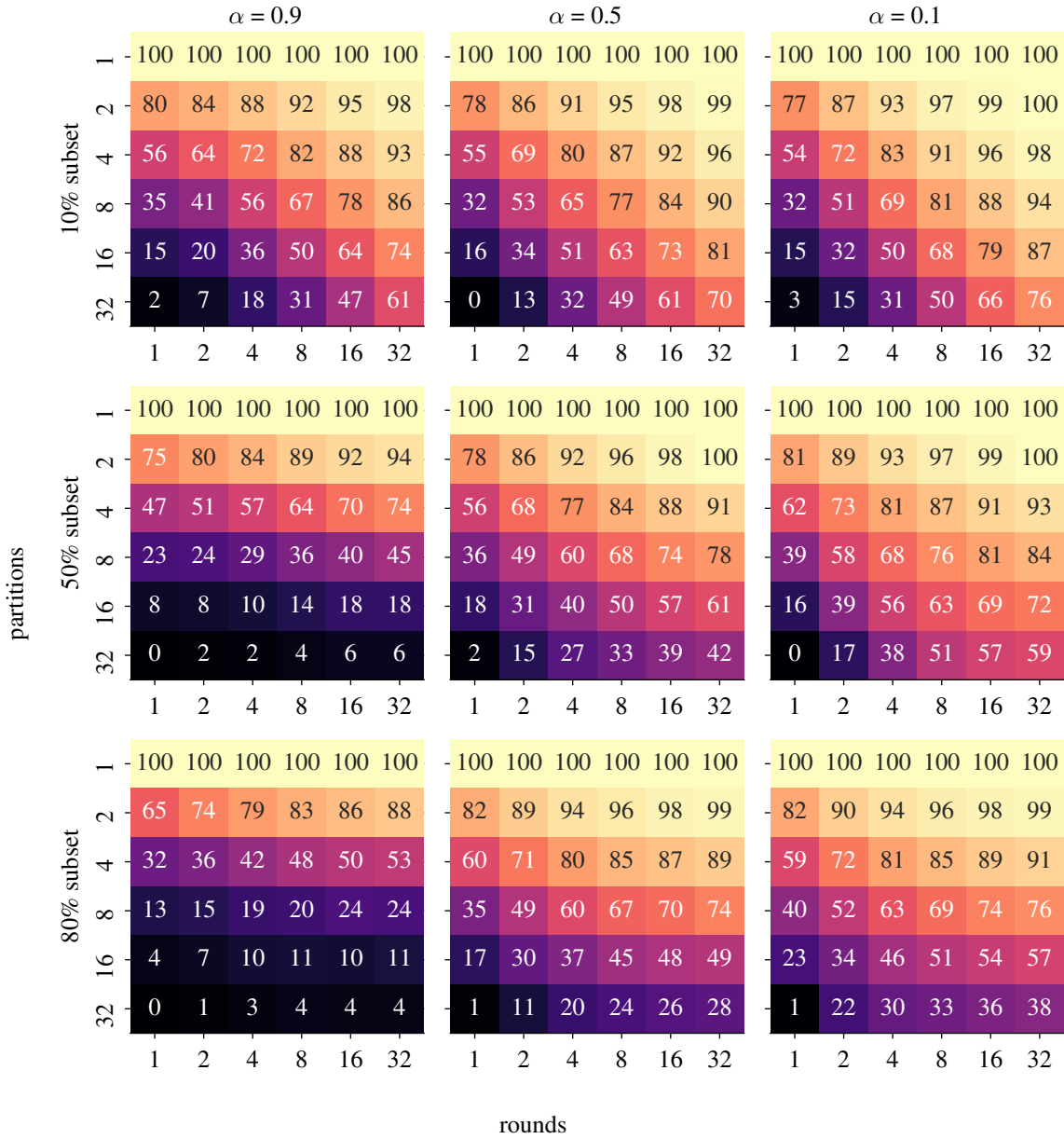


Figure 12. Full version of Figure 3. Normalized scores for finding subsets of CIFAR-100, depending on the subset size, the number of partitions, rounds, and  $\alpha$ .

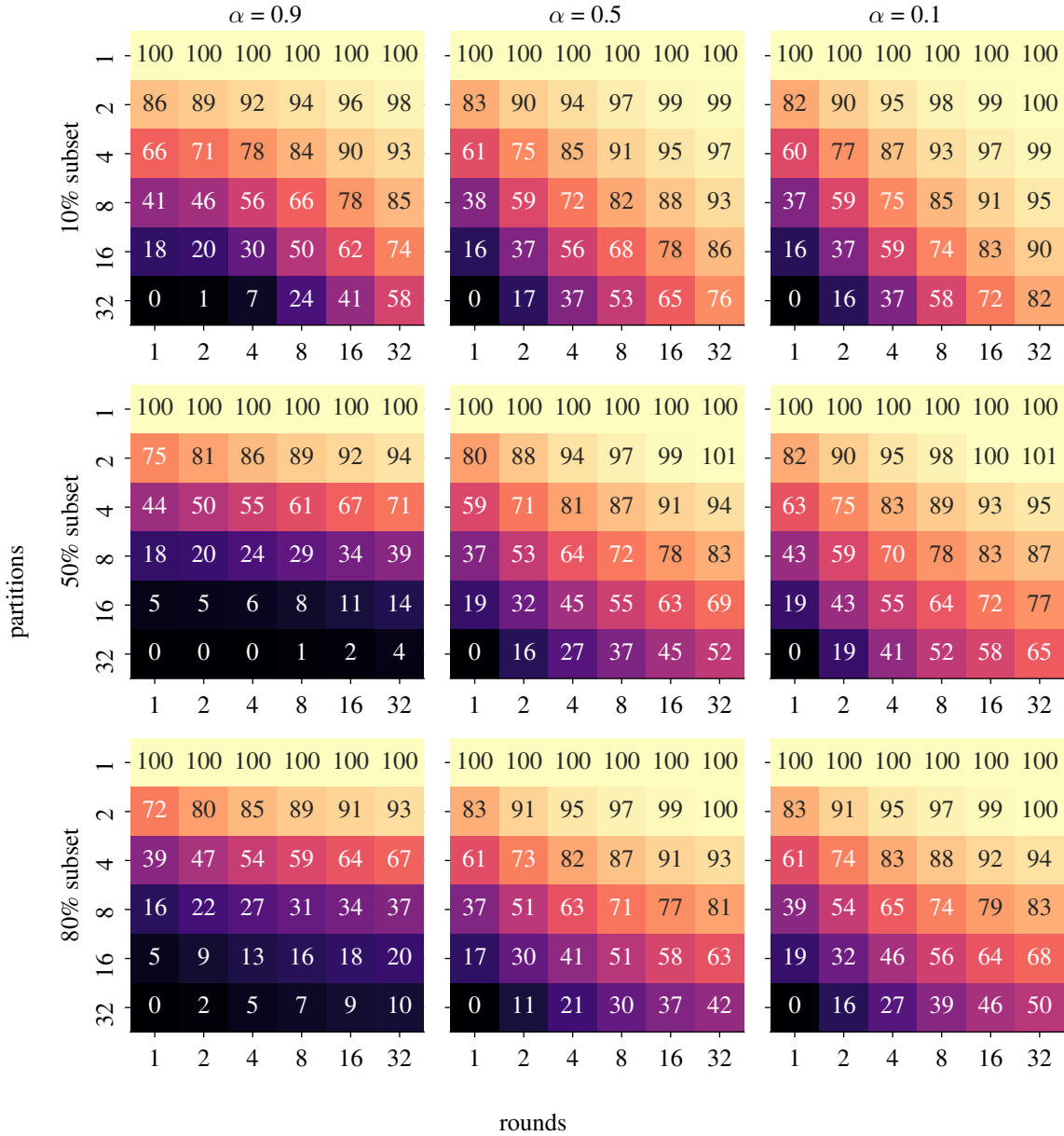


Figure 13. Normalized scores for finding subsets of ImageNet, depending on the subset size, the number of partitions, rounds, and  $\alpha$ .

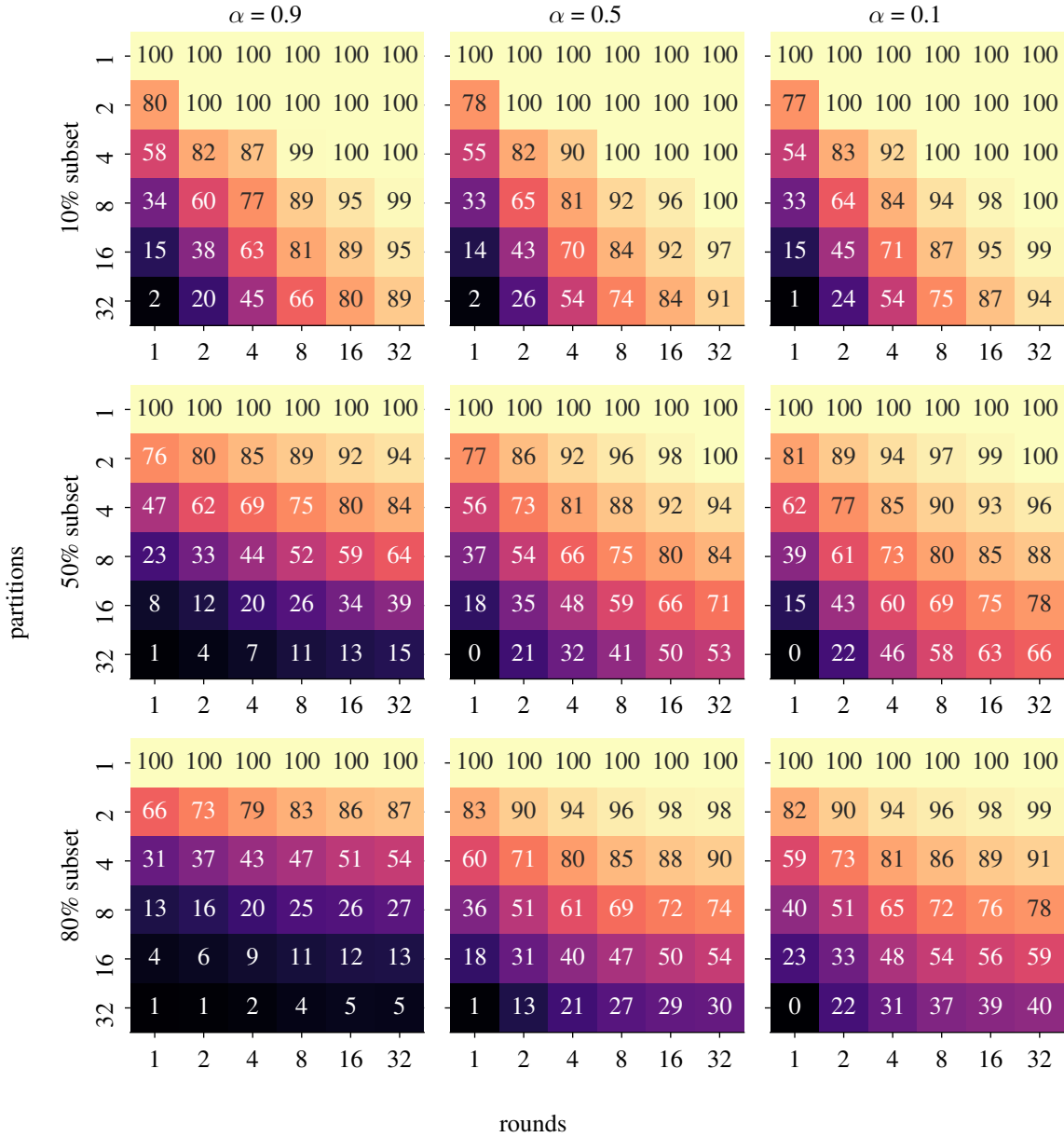


Figure 14. Full version of Figure 4. Normalized scores for finding subsets of CIFAR-100, depending on the subset size, the number of partitions, rounds, and  $\alpha$ , using adaptive partitioning.

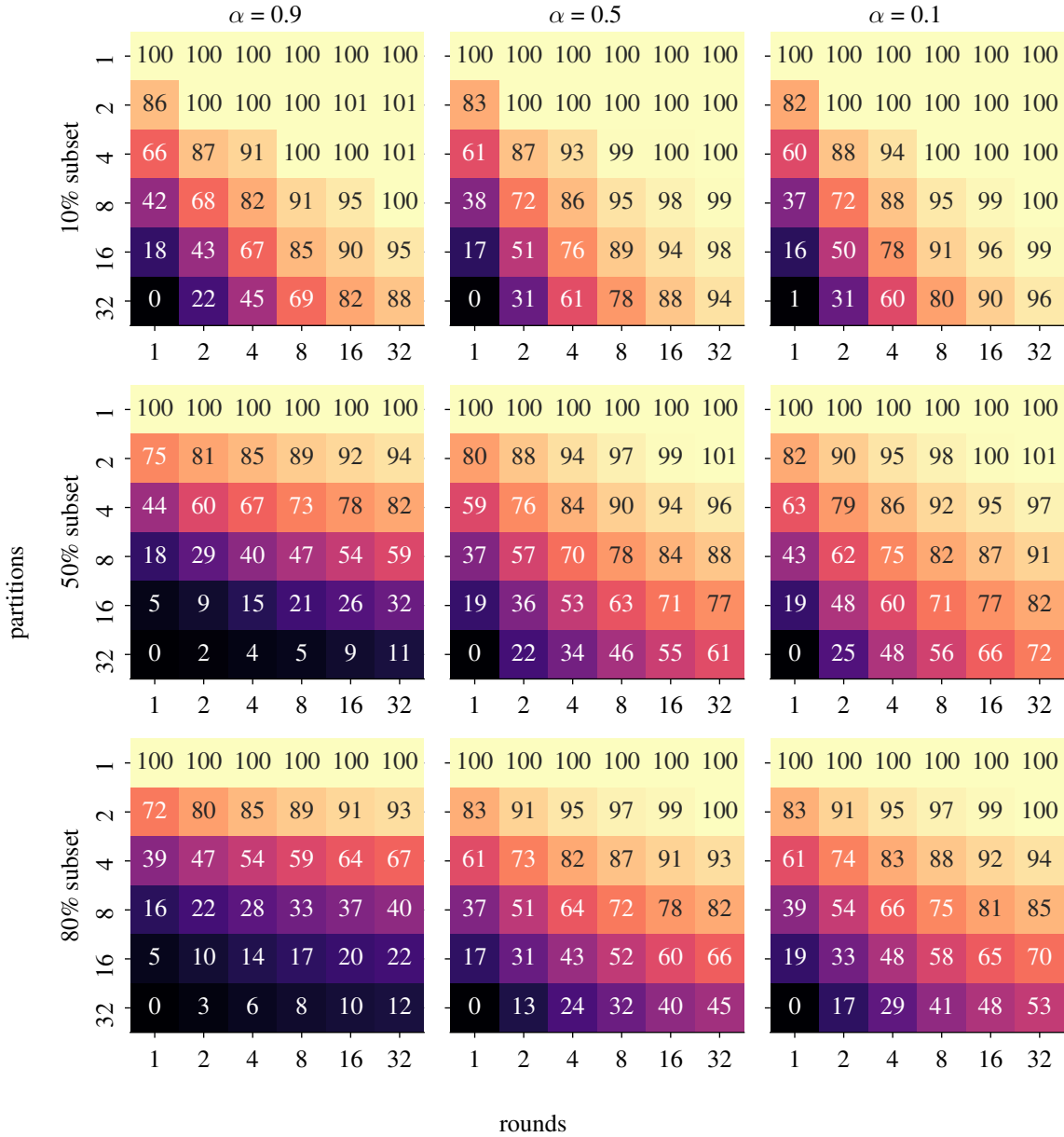


Figure 15. Normalized scores for finding subsets of ImageNet, depending on the subset size, the number of partitions, rounds, and  $\alpha$ , using adaptive partitioning.

## G. Additional Figures (with bounding)

In this section, we give the full heatmaps of the five bounding configurations on ImageNet and CIFAR. For statistics on the number of included points, refer to [Table 1](#). We only show the results for  $\alpha = 0.9$  since, as discussed in [Section 5.2](#), for other values of  $\alpha$ , bounding does not include or exclude points.

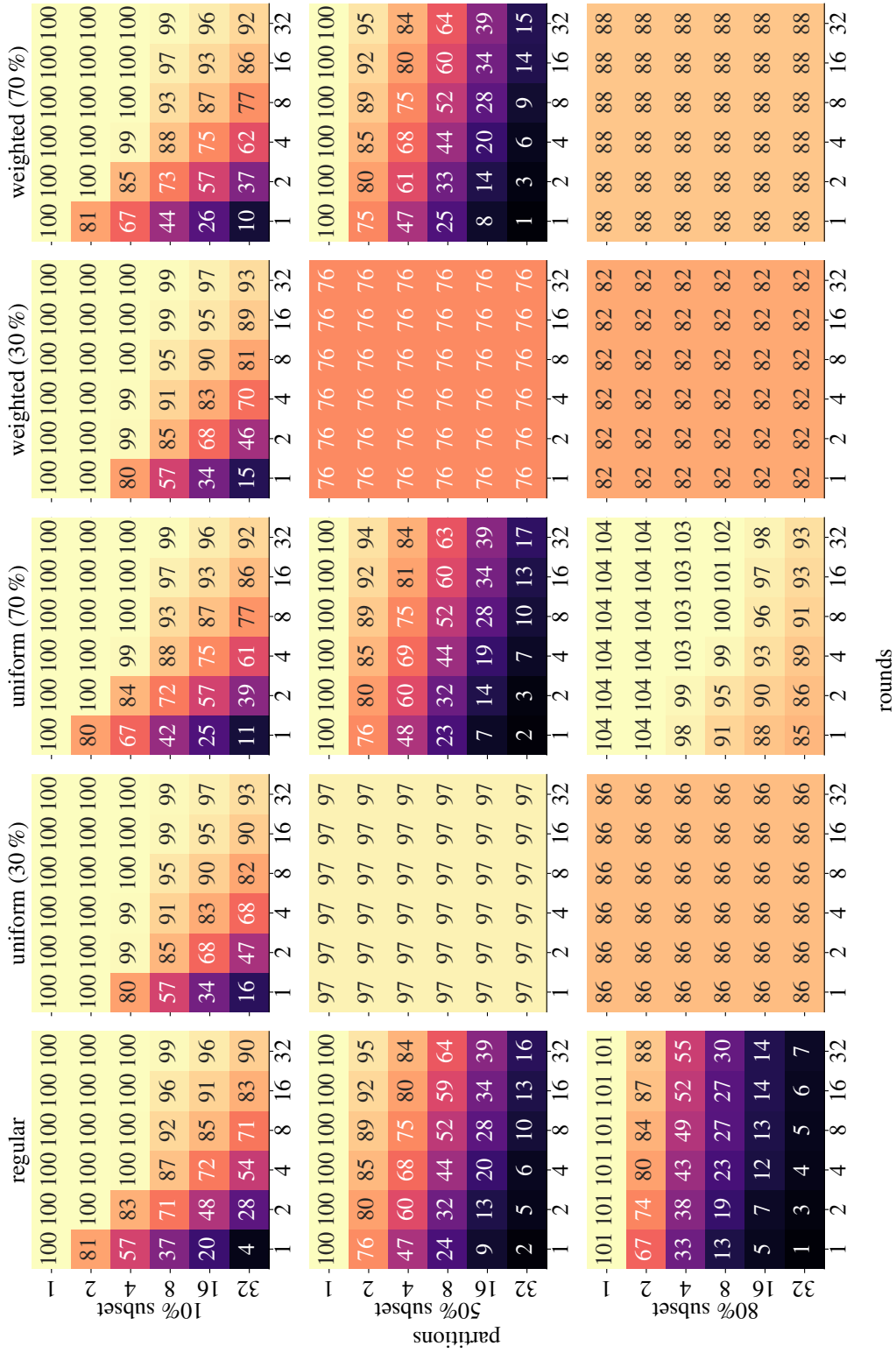


Figure 16. Normalized scores for finding subsets of CIFAR-100, depending on the subset size, the number of partitions, rounds, and type of bounding, using adaptive partitioning.

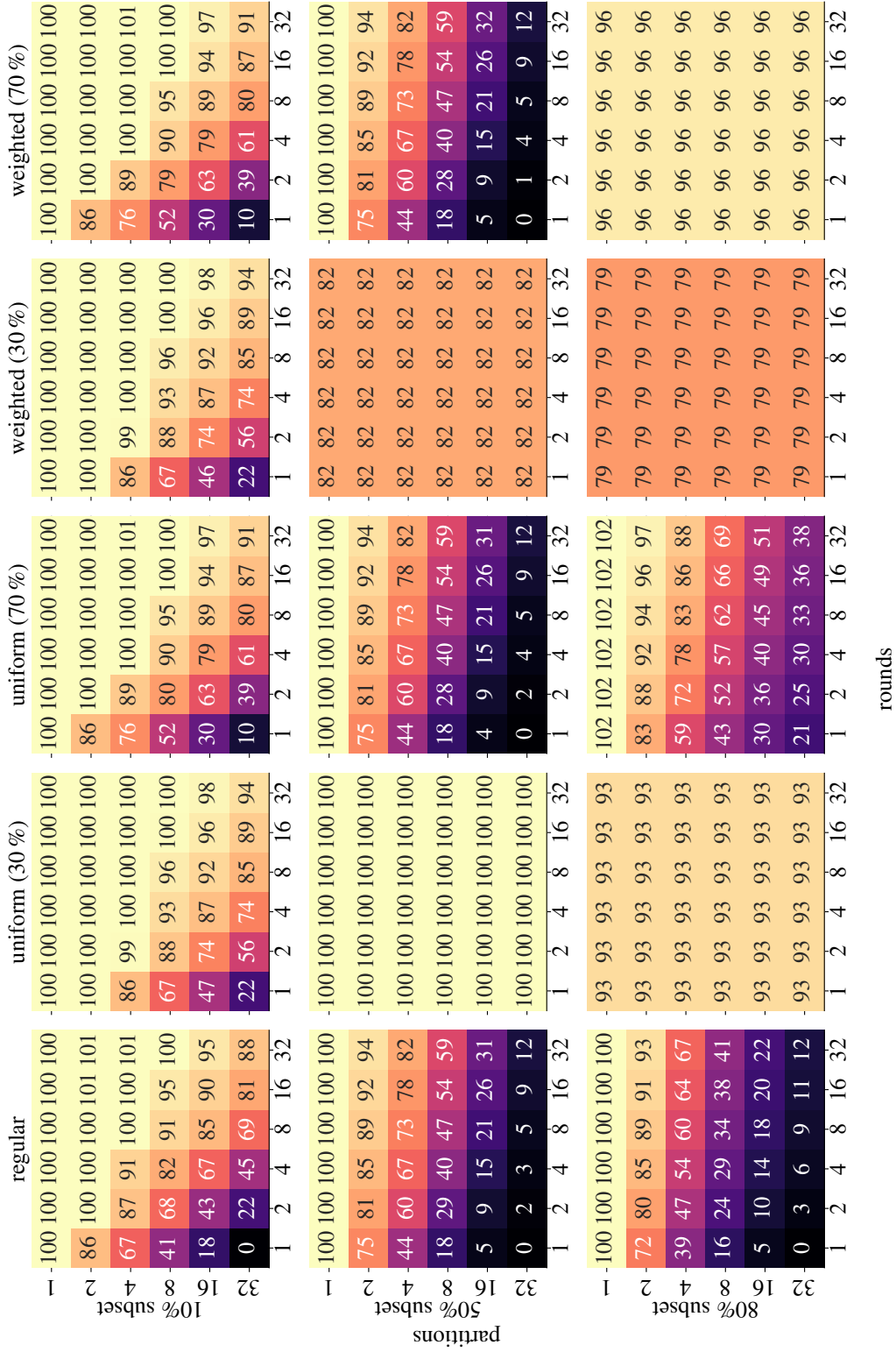


Figure 17. Normalized scores for finding subsets of ImageNet, depending on the subset size, the number of partitions, rounds, and type of bounding, using adaptive partitioning.