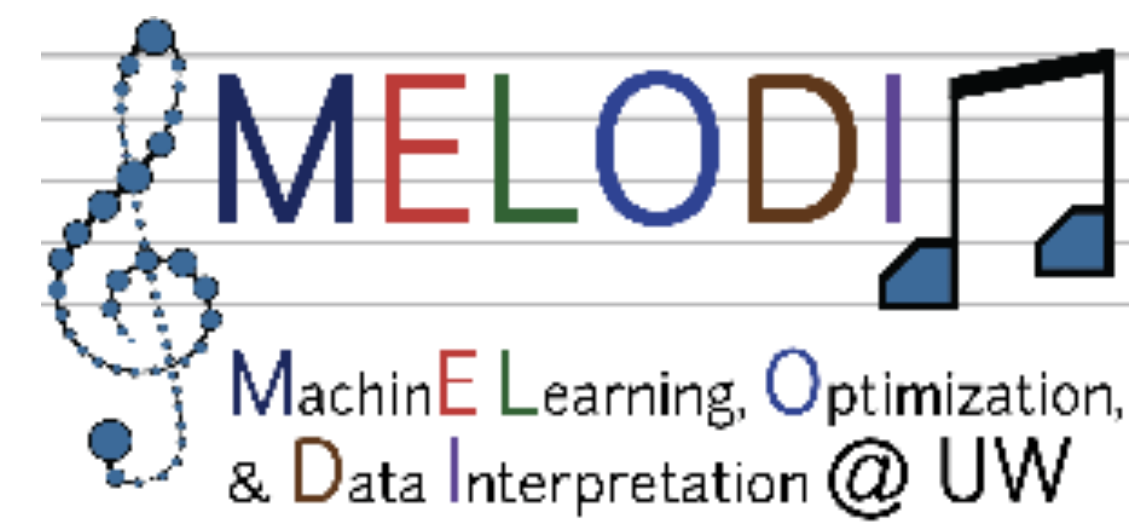


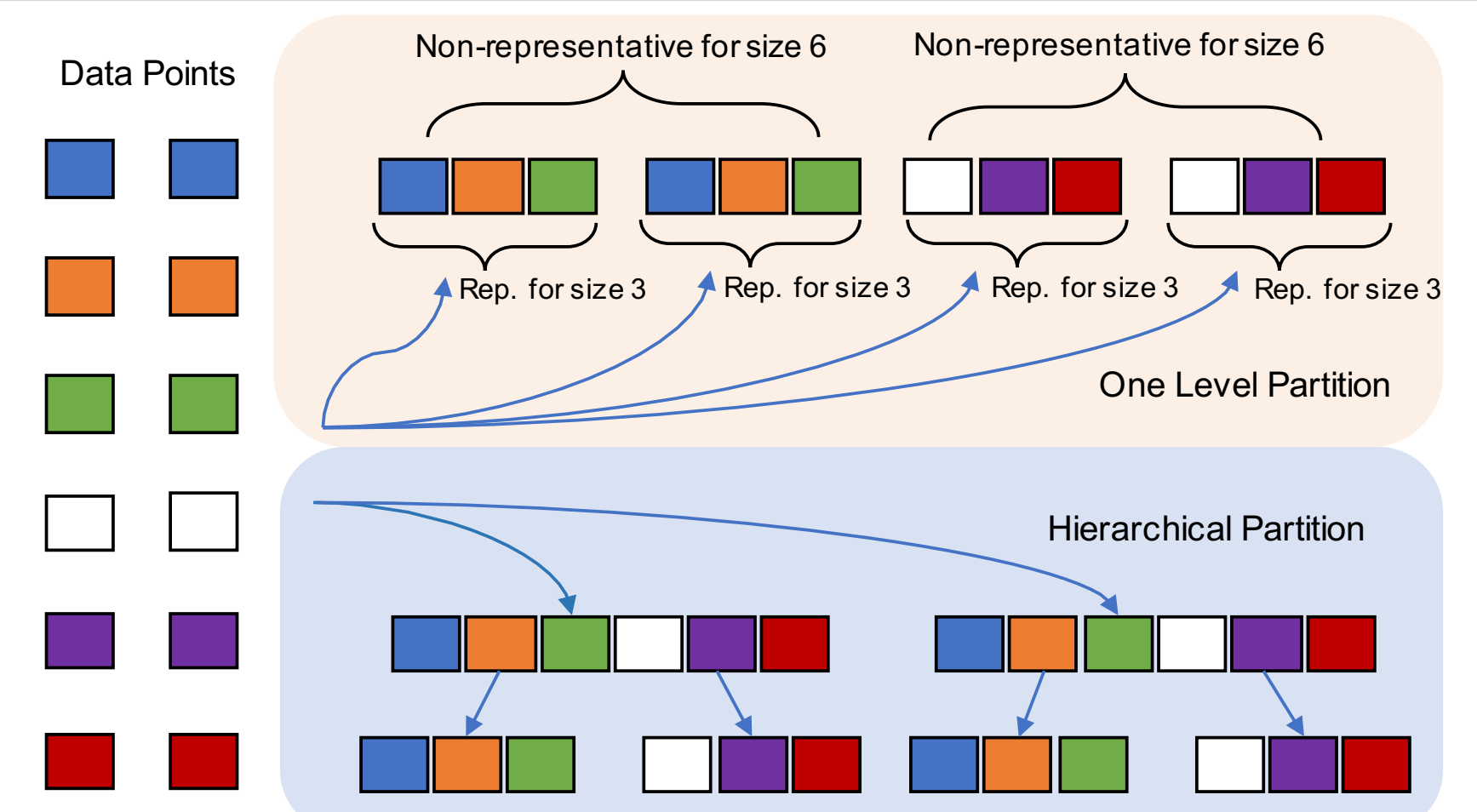


Fixing Mini-batch Sequences with Hierarchical Robust Partitioning

Shengjie Wang, Wenruo Bai, Chandrashekar Lavanaia,
Jeff A. Bilmes **University of Washington, Seattle**



- Problem:** Stochastic gradient methods require random access to data points which can be costly, especially for resource limited devices such as mobile devices and computing at the edge. Also, the randomly sampled mini-batches may contain redundant information.
- Approach:** We propose a general and efficient hierarchical robust partitioning framework to generate a deterministic sequence of mini-batches, one that offers assurances of being high quality, unlike a randomly drawn sequence.
- Experiments:** We test our fixed sequences to train neural networks on CIFAR-100 and Imagenet datasets and achieve significantly improved performance compared to randomly sampled sequences.



- Robust Submodular Partition with Cardinality Constraint:**
$$\max_{\pi \in \Pi(V, k)} \min_{i=1:m} f(\pi_i(V)).$$
 Where f is a submodular function, $\Pi(V, k)$ is all partitions of ground set V with size k . The max-min formulation enforces the worst mini-batch to be representative of V .

- Algorithm 1 - Greedy:** Pick the worst mini-batch and add the item with highest gain. To generate the sequence, we order the mini-batches by decreasing function values.

Theoretical Guarantees:

Theorem 1. For submodular function f on ground set V and mini-batch size k , suppose $m = |V|/k$, Algorithm 1 gives an approximation ratio of $\frac{e-1}{(e-1)m+1}$.

The bound almost matches the best known bound for the un-constrained robust submodular partition problem within a factor of $m/(m+1)$.

Algorithm 1: Cardinality Constrained Submodular Robust Partition (RobustPartitionK(f, V, k))

```

input : f, V, k
1 m := |V|/k R := V Let A1 = A2 = ... = Am = ∅
  while R ≠ ∅ do
2   j* ∈ argminj, |Aj| < k f(Aj) ; // least valued block.
3   v* ∈ argmaxv ∈ R f(v|Aj*) ; // best for block.
4   Aj* := Aj* ∪ {v*} ; // add to block.
5   R = R \ {v*}
6 end
7 Sort Aj's by f(Aj) so that f(Aj1) ≥ f(Aj2) ≥ ... ≥ f(Ajm)
8 return (Aj1, Aj2, ..., Ajm)

```

- Limitations of Algorithm 1:** 1) When using lazy greedy for the greedy step (line 3), the memory cost is proportional to $m|V| = |V|^2/k$. 2) if mini-batch size k is small, every mini-batch is not capable of representing the ground set V , we may get redundancies from combination of consecutive mini-batches.

Algorithm 2 – Run Algorithm 1 in hierarchy.

- Memory Efficiency and Group Representativeness:** The peak memory cost of Algorithm 2 is $\max_{i=1:r} m_i k_{i-1}$. Even for $r = 2$, $k_1 = |V|/2$, the peak memory is halved. In addition, mini-batches grouped by the hierarchical structure are also representative.

Theoretical Guarantees:

Definition 1. We run Algorithm 1 with ground set size V' , constraint size k' and $m' = |V'|/k'$, the greedy step gets executed $T = |V'|$ times, and we get a sequence of sets $Q = (A_1^T, A_2^T, \dots, A_{m'}^T)$ as the output, with $A_{j'}^T$ having the minimal evaluation, i.e., $j' \in \text{argmin}_{i=1:m'} f(A_i^T)$. There exists an earliest greedy step $1 \leq t \leq T$ such that $|A_{j'}^t| = k'$, and $j' \in \text{argmin}_{i=1:m'} f(A_i^t)$, we define $\tau := \min_{i=1:m'} |A_i^t|$.

Theorem 2. If we have $\tau \geq 2$ for every call to Algorithm 1 from Algorithm 2, then we achieve an approximation ratio of $\left(\frac{\tau-1}{2\tau-1}\right)^r \frac{k_r}{|V|}$.

Algorithm 2: Hierarchical Submodular Robust Partitioning

```

input : f, V, k1, ..., kr
1 k0 := |V|; Q1 := (V) ; // Qi's store sequence of sets to
  further partition
2 for i := 1; i ≤ r; i := i + 1 do
3   mi := ki-1/ki ; // mi: number of blocks for the next
  partition
4   Qi+1 = (); // Qi+1 initialized with an empty sequence
5   for j := 1; j ≤ |Qi|; j := j + 1 do
6     A1, ..., Ami = RobustPartitionK(f, Qi[j], ki) ;
7     // Qi[j]: jth set in the sequence
8     Append A1, ..., Ami to Qi+1 ; // Add partitioned
  blocks of Qi[j] to solution
9 end
10 return Qr+1

```

Empirical Results

- Choice of Submodular Function:** Nearest Neighbor Submodular Function (a special case of a facility location function):

$$f_{NN}(S) = \sum_{v \in V} \max_{v' \in S} \text{sim}(v, v'), \quad \text{sim}(v_1, v_2) = e^{-\frac{\|x(v_1) - x(v_2)\|_2}{\sigma}}$$

f_{NN} naturally captures the maximum likelihood estimates over the given data set for a nearest-neighbor classifier. As f_{NN} only requires similarities between data points within the same class, for large dataset such as Imagenet, we can afford to compute such sparse similarity graph.

Experiments on CIFAR-100 and Imagenet:

For CIFAR-100, we use the WRN-28-8 network, and for Imagenet, we use the Resnet-18. Compared to 30 randomly generated sequences on CIFAR-100, and the p-value for on test set accuracy is 0.0009. For Imagenet, we compare to 15 randomly generated sequences and the p-value for top-1 accuracy is 0.0151.

