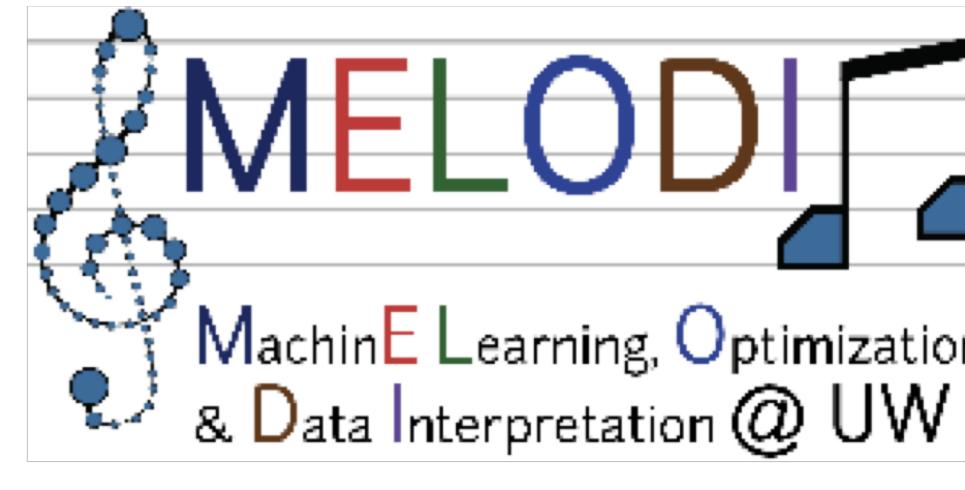


# Diverse Ensemble Evolution: Curriculum based Data-Model Marriage

Tianyi Zhou, Shengjie Wang, Jeff Bilmes, University of Washington, Seattle



- Background:** An ensemble of light-weight models can outperform a large model, and can be trained at multiple edge devices with limited computational power.
- Problem:** How to train an ensemble of models jointly, so they all have different and diverse expertise, but are complementary to each other.
- Basic Idea:** teacher adaptively assigns training samples to multiple learners for each epoch of their training processes.
- Previous works on ensemble training:** only encourage diversity before training begins, not adaptive, always train each model on the whole training set.
- Our Method (DivE<sup>2</sup>):** can evolve diverse and complementary expertise on different models faster, achieve better performance with less training time.
- Applications:** distributed/federated machine learning; edge computing.

## Combinatorial-Continuous Hybrid Optimization

$$\max_{W} \max_{A \subseteq E, A \in \mathcal{I}_v \cap \mathcal{I}_u} G(A, W) \triangleq \sum_{(v_j, u_i) \in A} (\beta - \ell(v_j; w_i)) + \gamma F_{\text{inter}}(A) + \lambda F_{\text{intra}}(A)$$

$$\mathcal{I}_u = \{A | A \subseteq E, |A \cap \delta(u)| \leq p, \forall u \in U\}$$

Model Selecting (at most  $p$ ) Sample

$$\mathcal{I}_v = \{A | A \subseteq E, |A \cap \delta(v)| \leq k, \forall v \in V\}$$

Sample Selecting (at most  $k$ ) Model

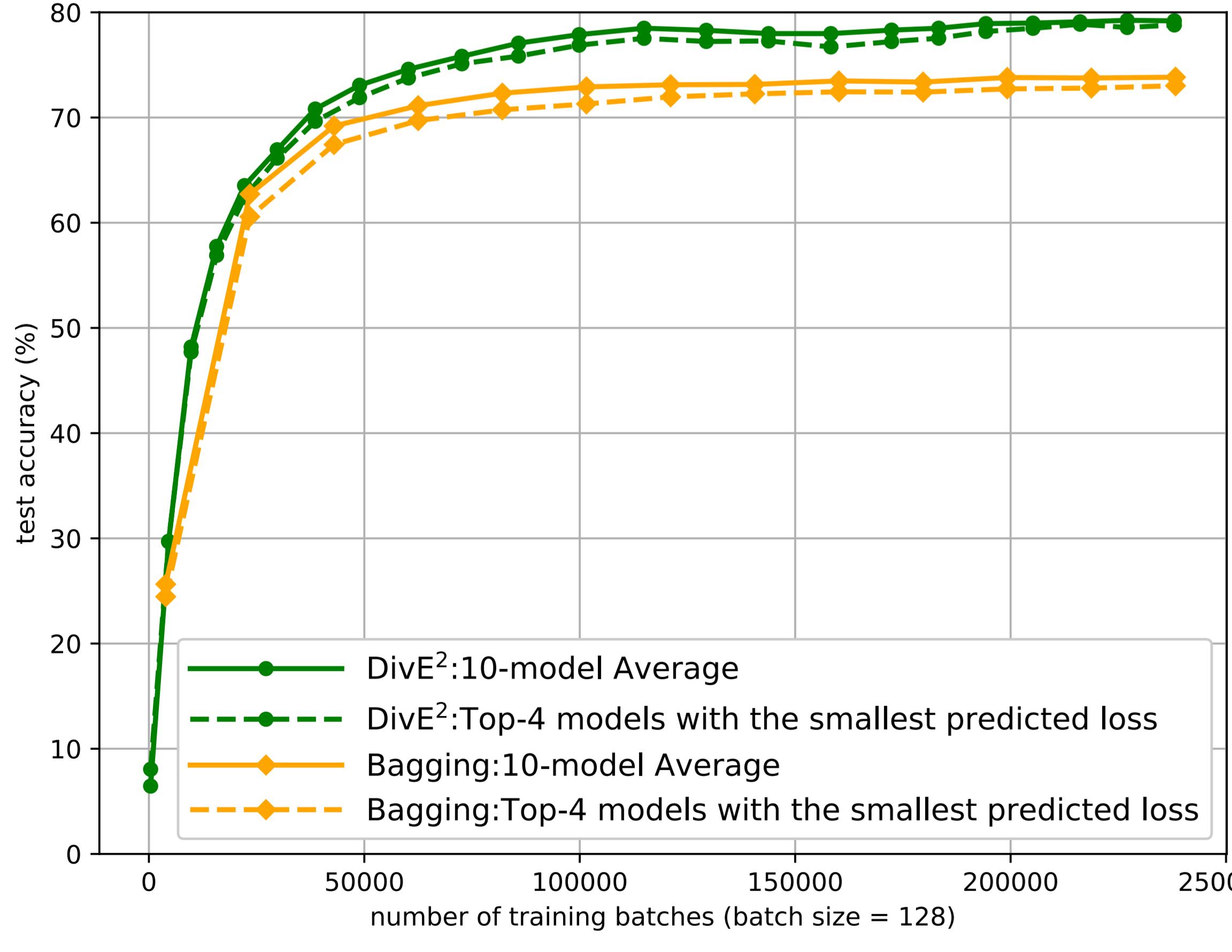
$$F_{\text{inter}}(A) \triangleq \sum_{i,j \in [m], i < j} F(\delta(u_i) \cup \delta(u_j) \cap A)$$

Inter-model Diversity (Submodular)

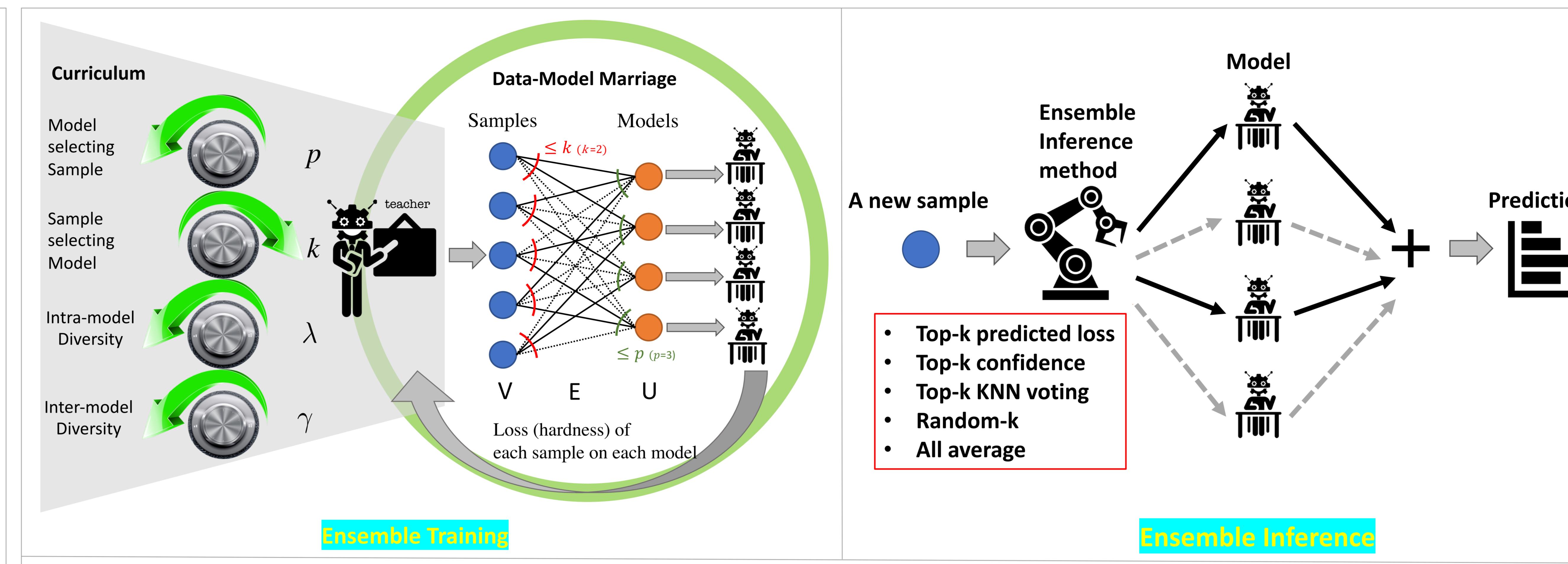
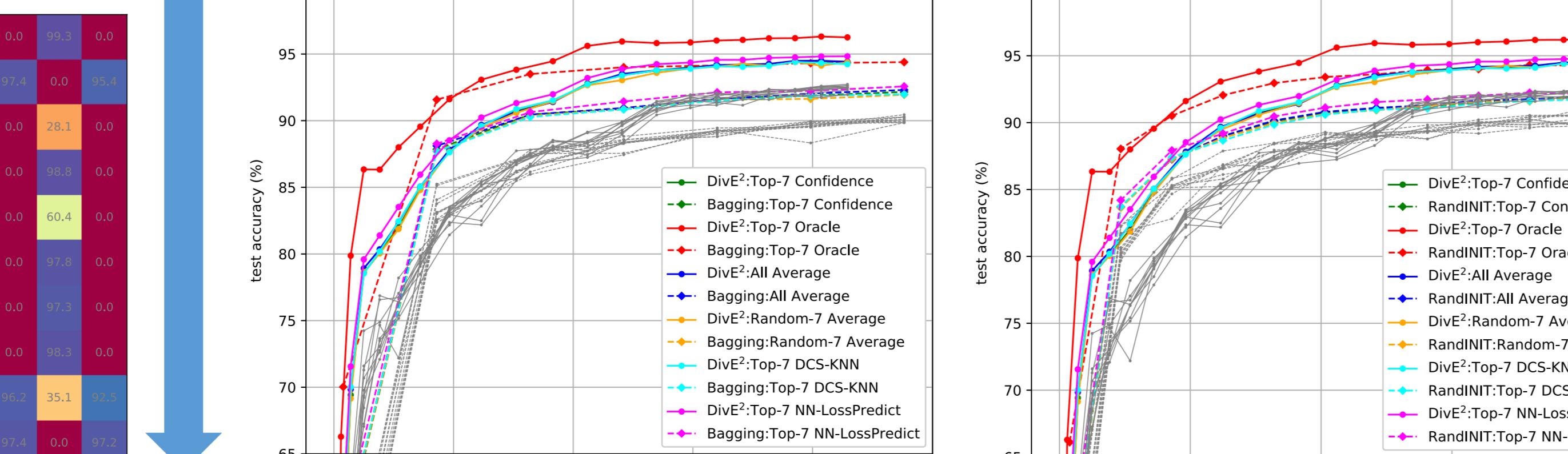
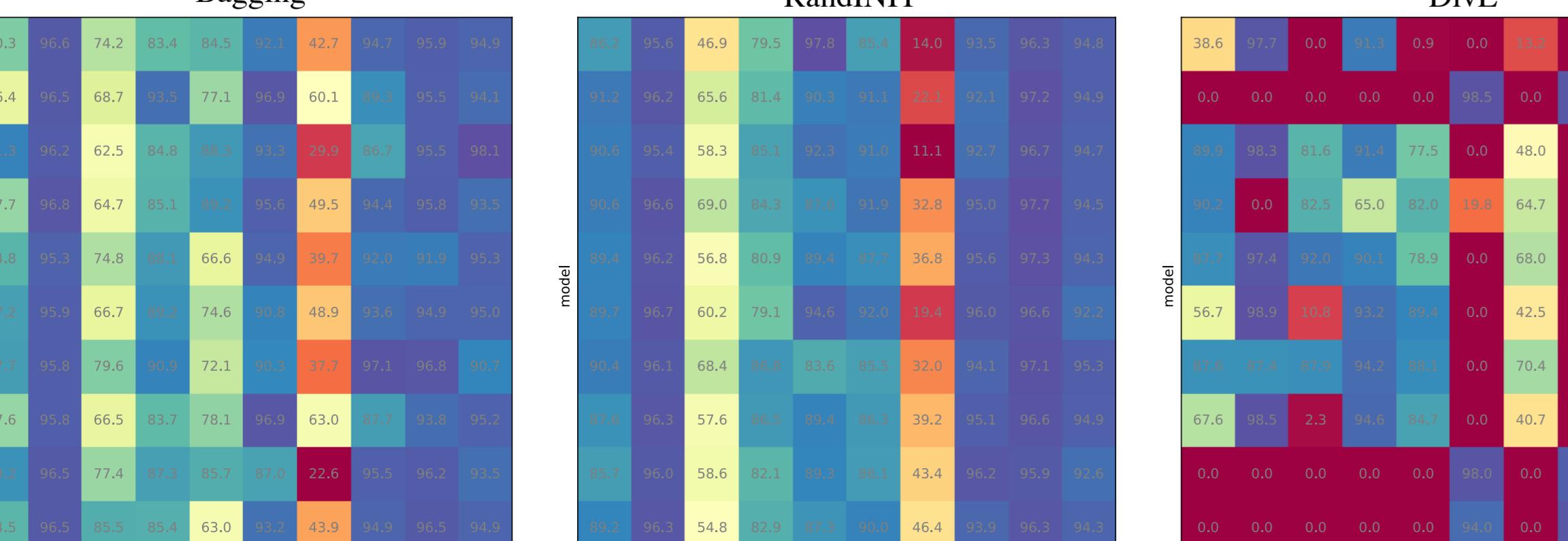
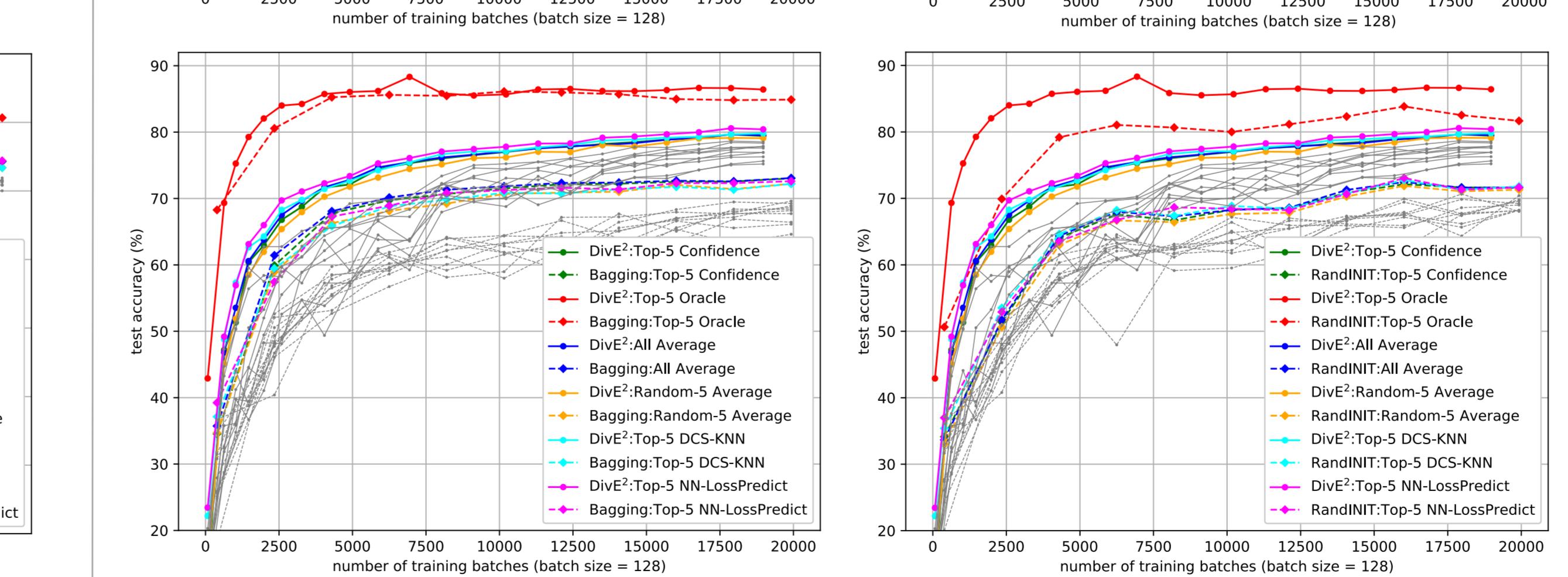
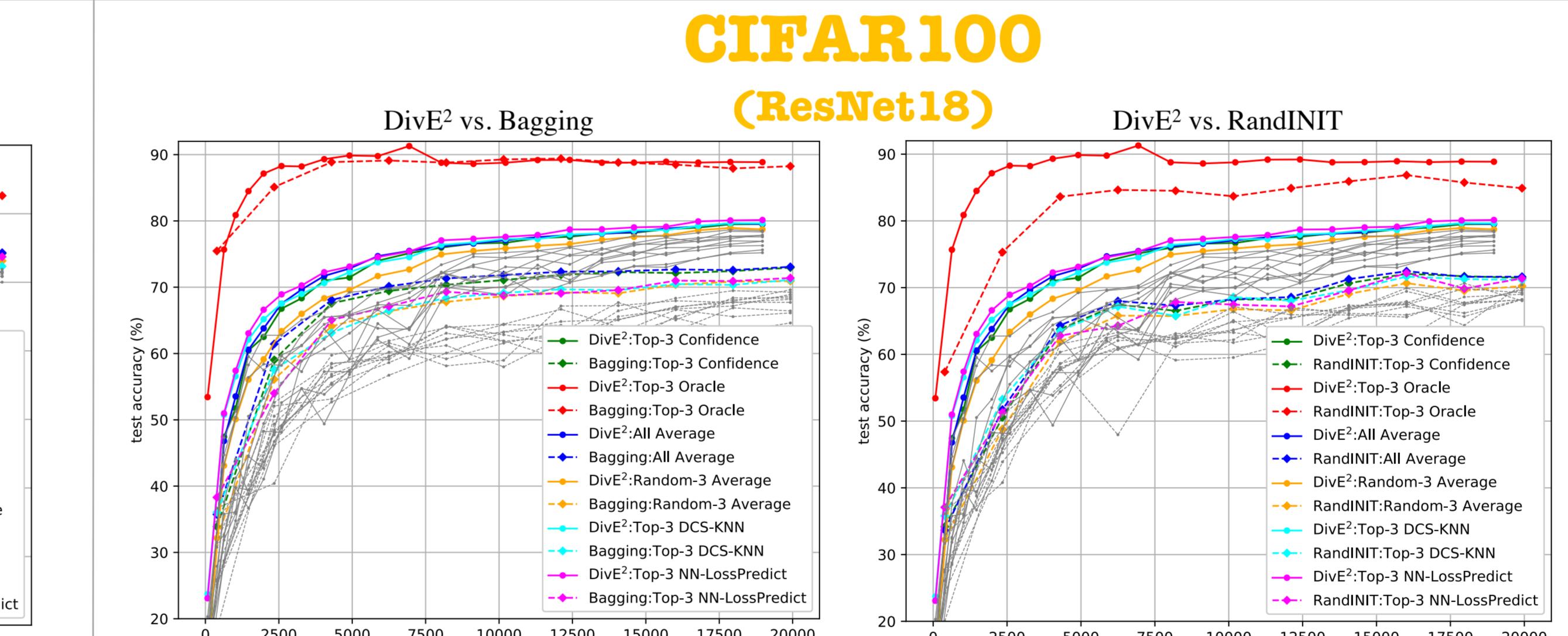
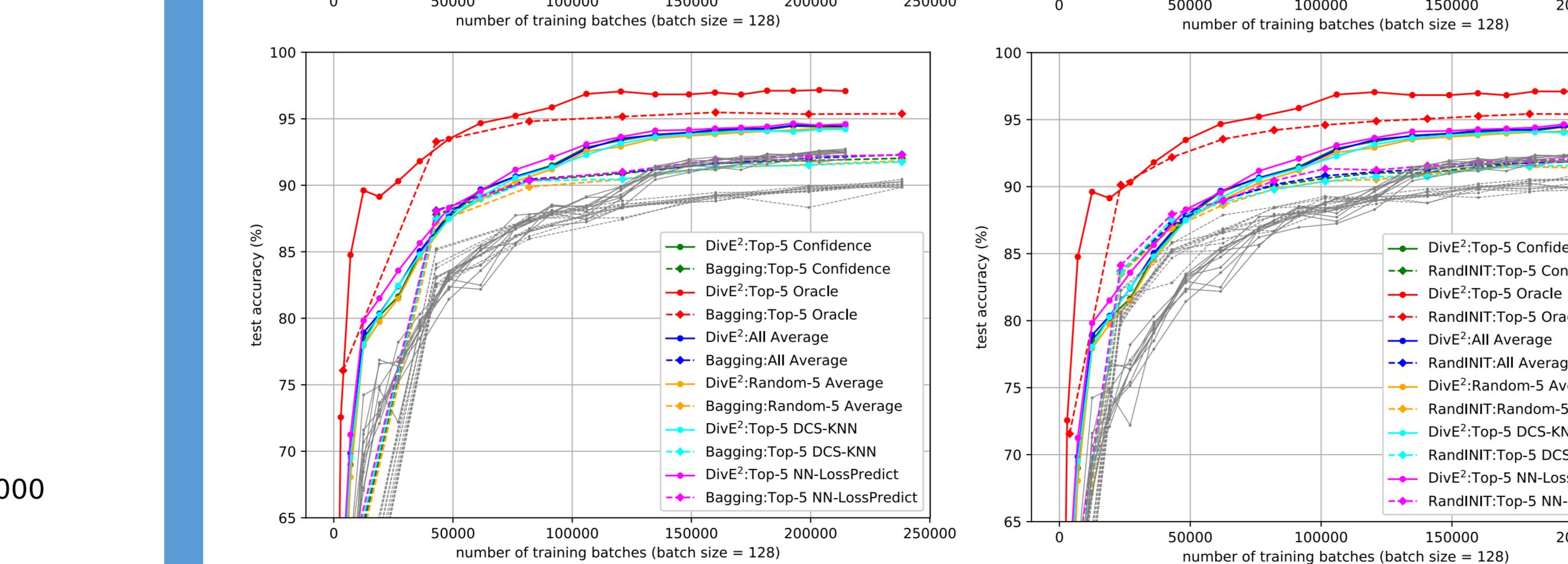
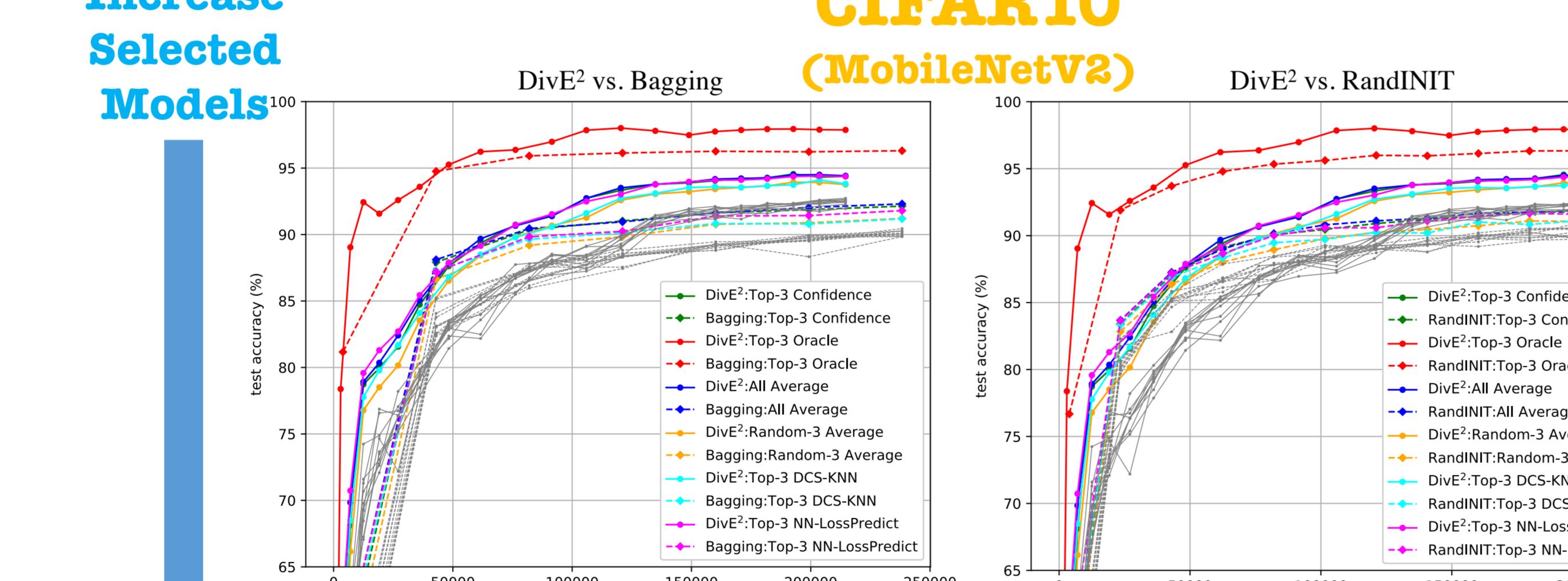
$$F_{\text{intra}}(A) = \sum_{i \in [m]} F'(\delta(u_i) \cap A)$$

Intra-model Diversity (Submodular)

## Experiments



### Increase Selected Models



- Curriculum for Data-Model Marriage:** gradually tune four knobs controlling the data assignment process (formulated as a bipartite graph matching), which includes two matroid constraints and two diversity terms.
- Model selecting Sample & Sample Selecting Model:** in earlier stages, every model selects  $p$  samples easier to learn; in later stages, every sample selects  $k$  models which provide better prediction; use them together to avoid imbalance.
- Inter-model & Intra-model Diversity:** intra-model diversity selects diverse training samples for each model; inter-model diversity encourages less overlap between training samples assigned to different models; we start from large diversity and gradually reduce it, so the learners expand their expertise regions and gradually turn focus on data they are good at.

## Algorithm

### Algorithm 1 SELECTLEARN( $k, p, \lambda, \gamma, \{w_i^0\}_{i=1}^m$ )

```

1: Input:  $\{v_j\}_{j=1}^n, \{\ell(\cdot; w_i)\}_{i=1}^m, \pi(\cdot; \eta)$ 
2: Output:  $\{w_i^t\}_{i=1}^m$ 
3: Initialize:  $w_i \leftarrow w_i^0 \quad \forall i \in [m], t = 0$ 
4: while not "converged" do
5:    $W \leftarrow \{w_i^t\}_{i=1}^m$ , define  $G(\cdot, W)$  by  $W$ ;
6:    $\hat{A} \leftarrow \text{SUBMODULARMAX}(G(\cdot, W), k, p)$ ;
7:   if  $G(\hat{A}, W) > G(A, W)$  then
8:      $A \leftarrow \hat{A}$ ;
9:   end if
10:  for  $i \in \{1, \dots, m\}$  do
11:     $-\nabla \hat{H}(w_i^t) \leftarrow \frac{\partial}{\partial w_i^t} \sum_{v_j \in V(A \cap \delta(u_i))} \ell(v_j; w_i^t)$ ;
12:     $w_i^{t+1} \leftarrow w_i^t + \pi(w_i^t, -\nabla \hat{H}(w_i^t)^{1:t}; \eta^t)$ ;
13:  end for
14:   $t \leftarrow t + 1$ ;
15: end while
```

### Algorithm 2 Diverse Ensemble Evolution (DivE<sup>2</sup>)

```

1: Input:  $\{(x_j, y_j)\}_{j=1}^n, \{w_i^0\}_{i=1}^m, \pi(\cdot; \eta), \mu, \Delta_k, \Delta_p, T$ 
2: Output:  $\{w_i^t\}_{i=1}^m$ 
3: Initialize:  $k \leq m, p \geq 1$  s.t.  $mp \leq nk$ ,
    $\lambda \in [0, 1], \gamma \in [0, 1]$ 
4: for  $t \in \{1, \dots, T\}$  do
5:    $\{w_i^t\}_{i=1}^m \leftarrow \text{SELECTLEARN}(k, p, \lambda, \gamma, \{w_i^{t-1}\}_{i=1}^m)$ ;
6:    $\lambda \leftarrow (1 - \mu) \cdot \lambda, \gamma \leftarrow (1 - \mu) \cdot \gamma$ ;
7:    $k \leftarrow \max\{\lceil k - \Delta_k \rceil, 1\}, p \leftarrow \min\{p + \Delta_p, n\}$ ;
8: end for
```