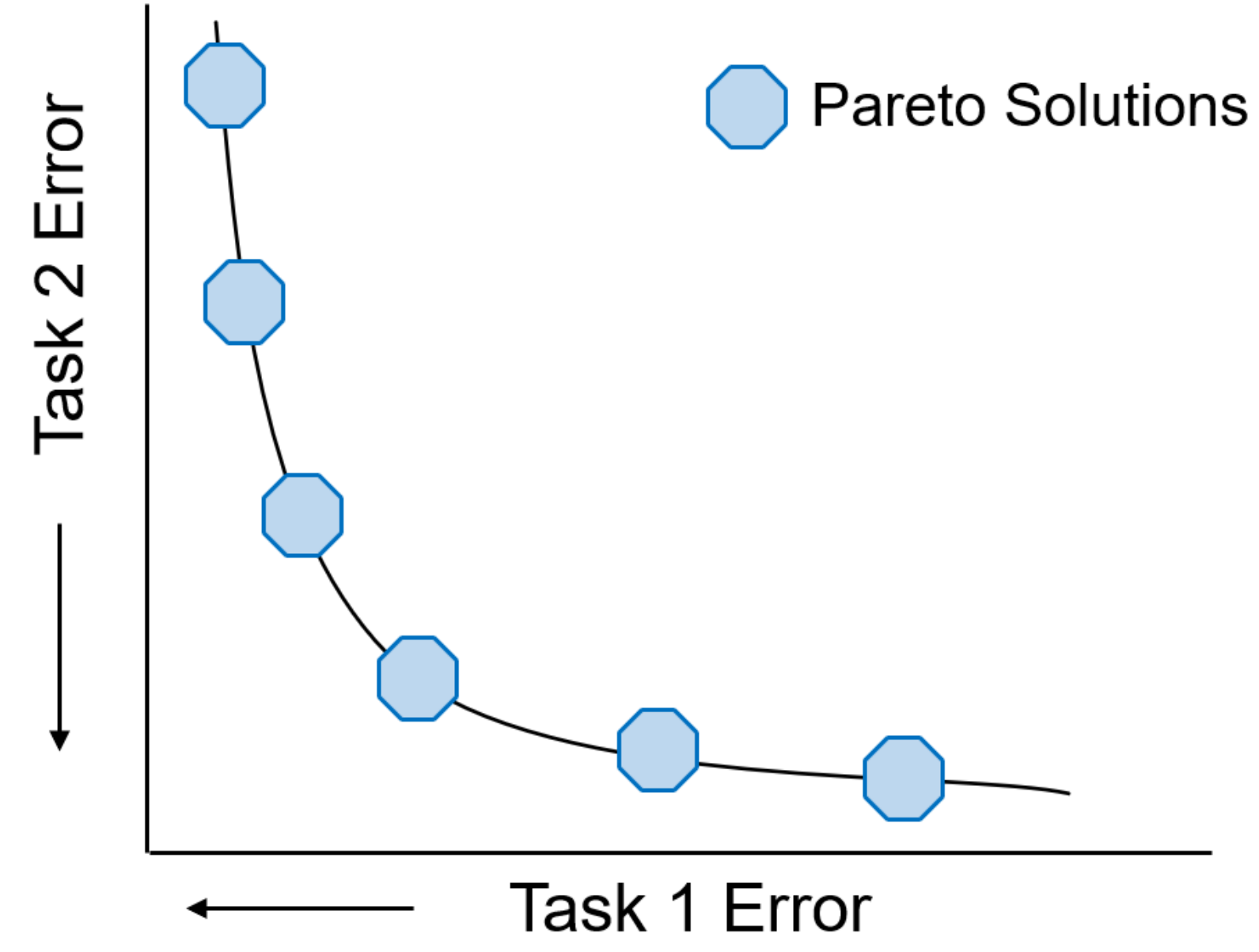


Problem Definition and Contribution

Goal: To generate widely distributed Pareto solutions with different trade-offs for multi-task learning (MTL).



Then MTL practitioners can easily select their preferred solution(s) with different optimal trade-offs.

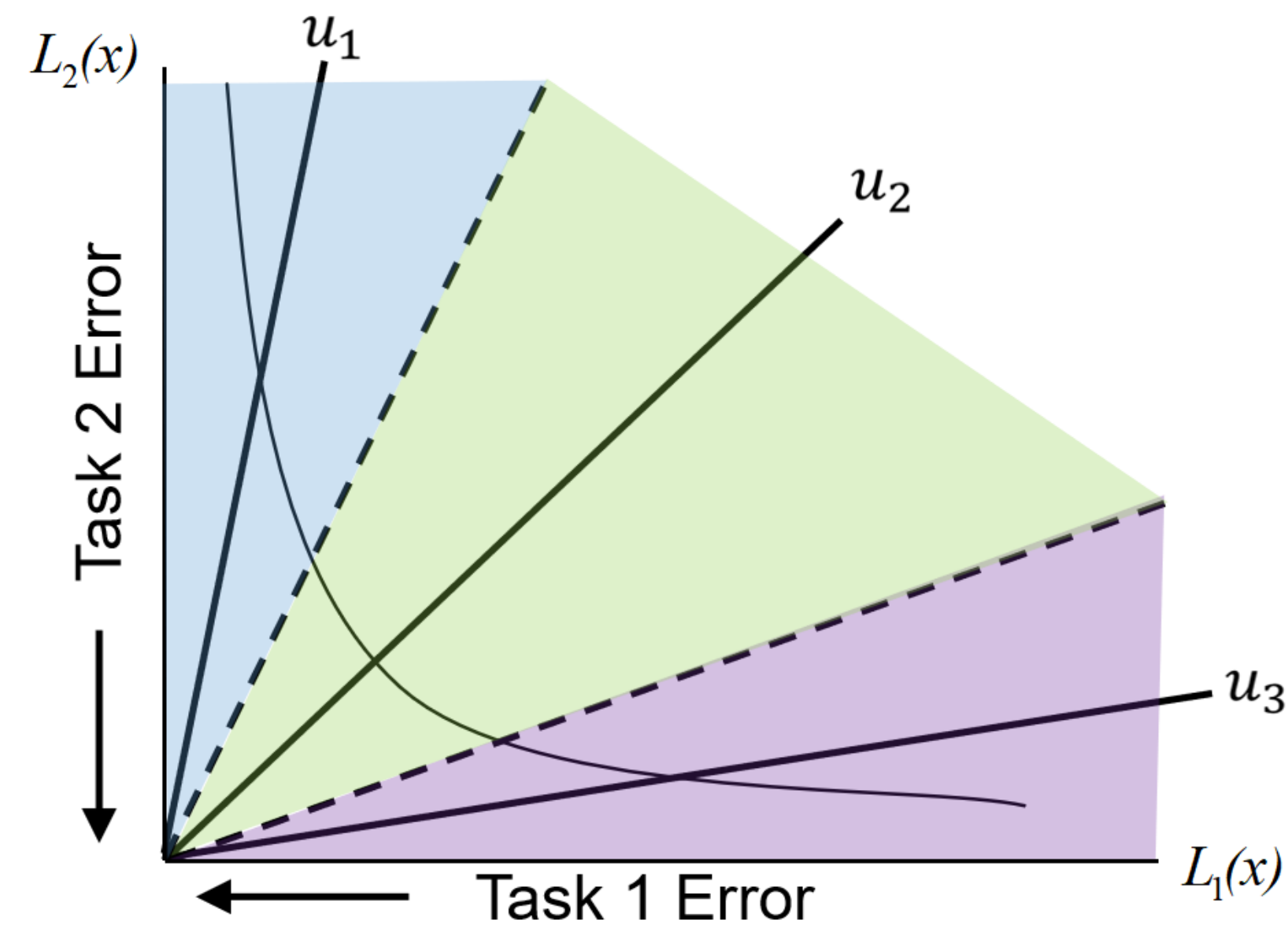
Motivations:

- Current linear scalarization methods need exhaustive weights search which could be inefficient.
- Existing adaptive weight methods can not make different trade-offs among tasks.

Main Contributions:

- A novel method to decompose a MTL problem into multiple subproblems with different trade-offs.
- Show the proposed Pareto MTL can be reformulated as a linear scalarization approach to solve MTL with dynamically adaptive weights.
- A scalable optimization algorithm to solve all constrained subproblems with different preferences.

Key Idea:



Pareto MTL decomposes a given MTL problem into several subproblems with a set of preference vectors. Each MTL subproblem aims at finding one Pareto solution in its restricted preference region.

Problem Formulation

Problem: Consider a MTL problem with m correlated tasks:

$$\min_{\theta} \mathcal{L}(\theta) = (\mathcal{L}_1(\theta), \mathcal{L}_2(\theta), \dots, \mathcal{L}_m(\theta))^T. \quad (1)$$

Traditional linear scalarization:

$$\min_{\theta} \mathcal{L}(\theta) = \sum_{i=1}^m w_i \mathcal{L}_i(\theta). \quad (2)$$

where w_i is hard to be set.

Pareto MTL: Decompose a MTL problem with a set of unit preference vectors $\{u_1, u_2, \dots, u_K\}$ in R_+^m :

$$\min_{\theta} \mathcal{L}(\theta) = (\mathcal{L}_1(\theta), \dots, \mathcal{L}_m(\theta))^T, s.t. \mathcal{L}(\theta) \in \Omega_k, \quad (3)$$

where $\Omega_k (k = 1, \dots, K)$ is a subregion:

$$\Omega_k = \{v \in R_+^m | u_j^T v \leq u_k^T v, \forall j = 1, \dots, K\}. \quad (4)$$

The constraints can be reformulated as:

$$\mathcal{G}_j(\theta_t) = (u_j - u_k)^T \mathcal{L}(\theta_t) \leq 0, \forall j = 1, \dots, K. \quad (5)$$

Solving the Subproblems: Find a valid direction to minimize all loss functions and activated constraints.

$$(d_t, \alpha_t) = \arg \min_{d \in R^n, \alpha \in R} \alpha + \frac{1}{2} \|d\|^2$$

$$s.t. \quad \nabla \mathcal{L}_i(\theta_t)^T d \leq \alpha, i = 1, \dots, m.$$

$$\nabla \mathcal{G}_j(\theta_t)^T d \leq \alpha, j \in I_{\epsilon}(\theta_t), \quad (6)$$

where $I_{\epsilon}(\theta)$ is the index set of activated constraints.

Pareto MTL as Adaptive Linear Scalarization:

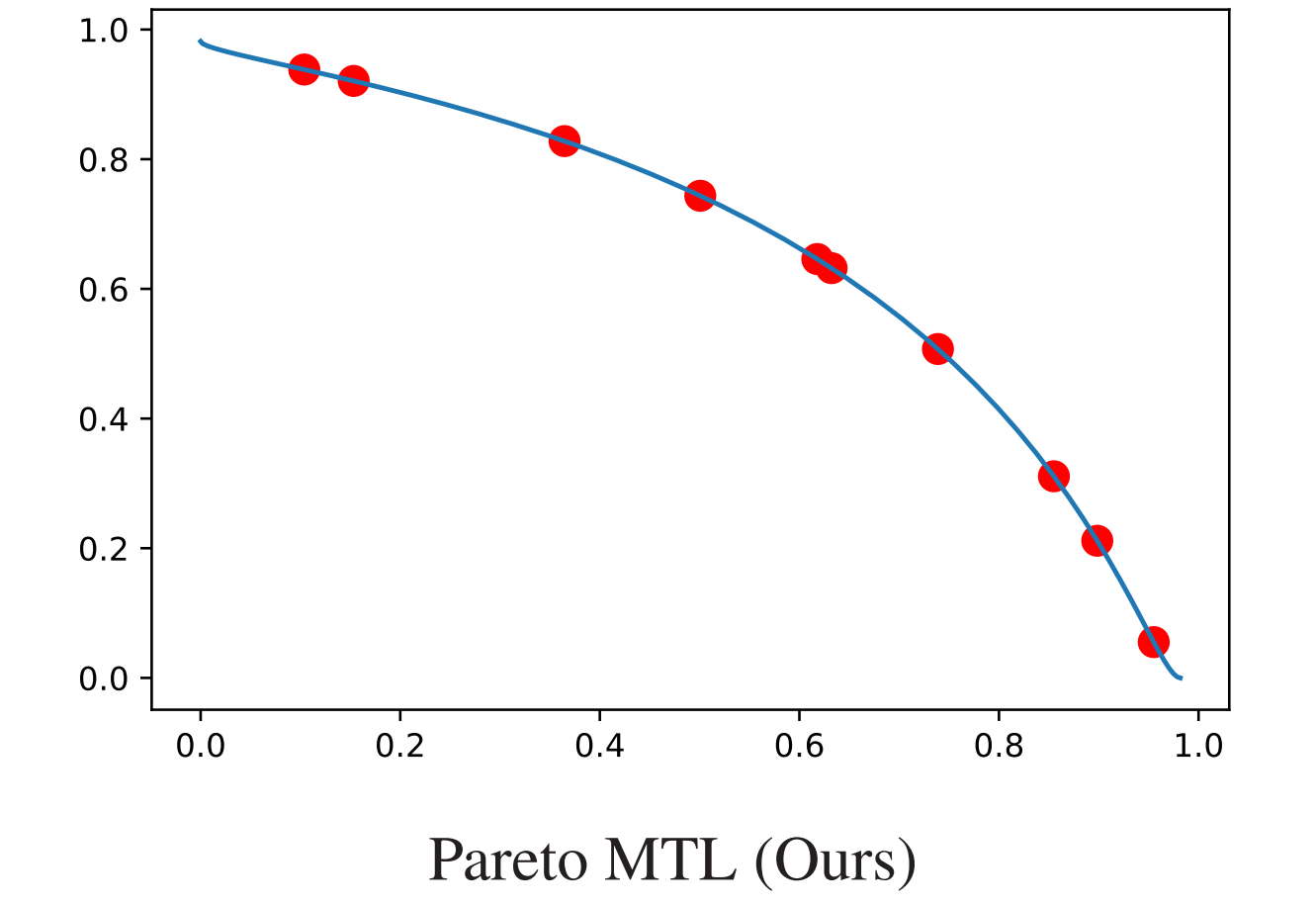
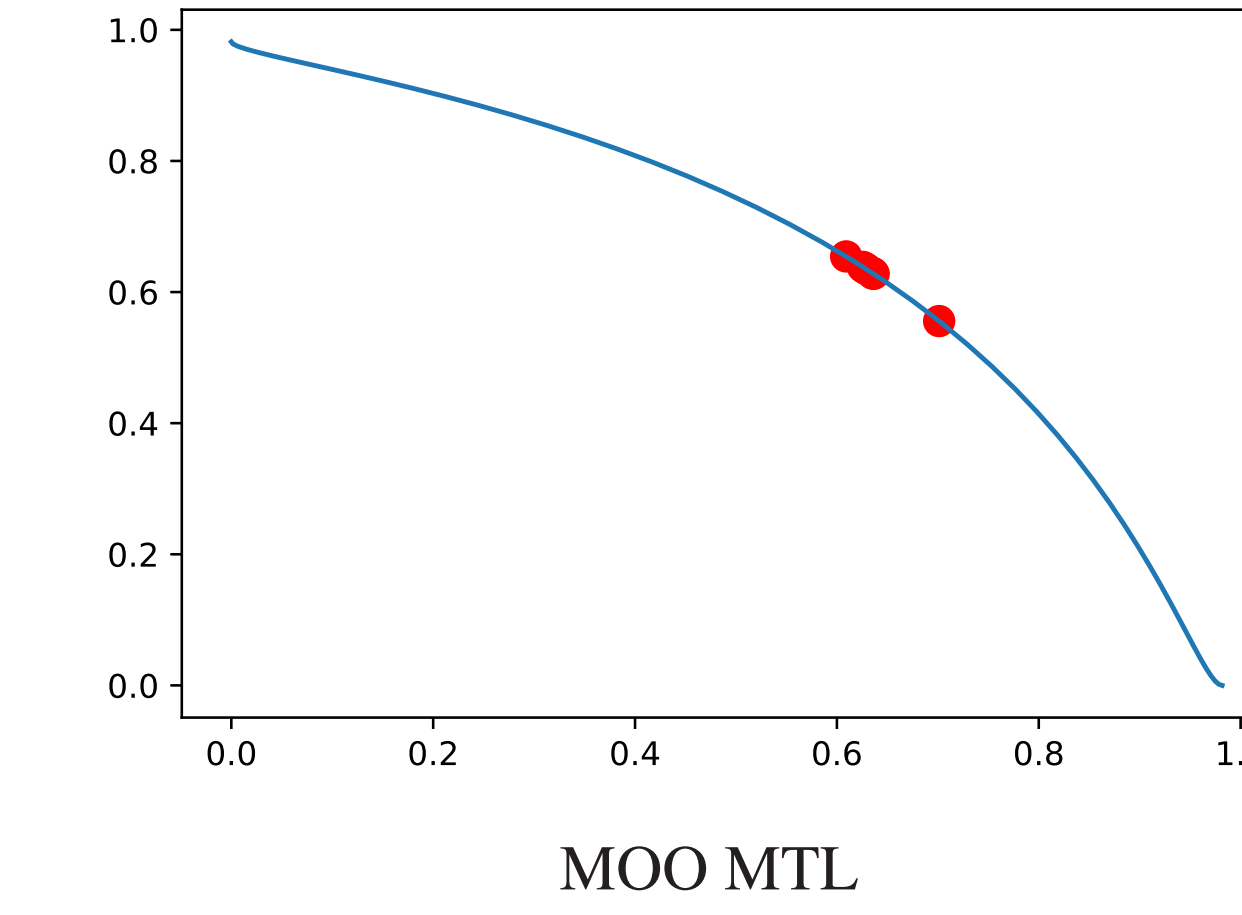
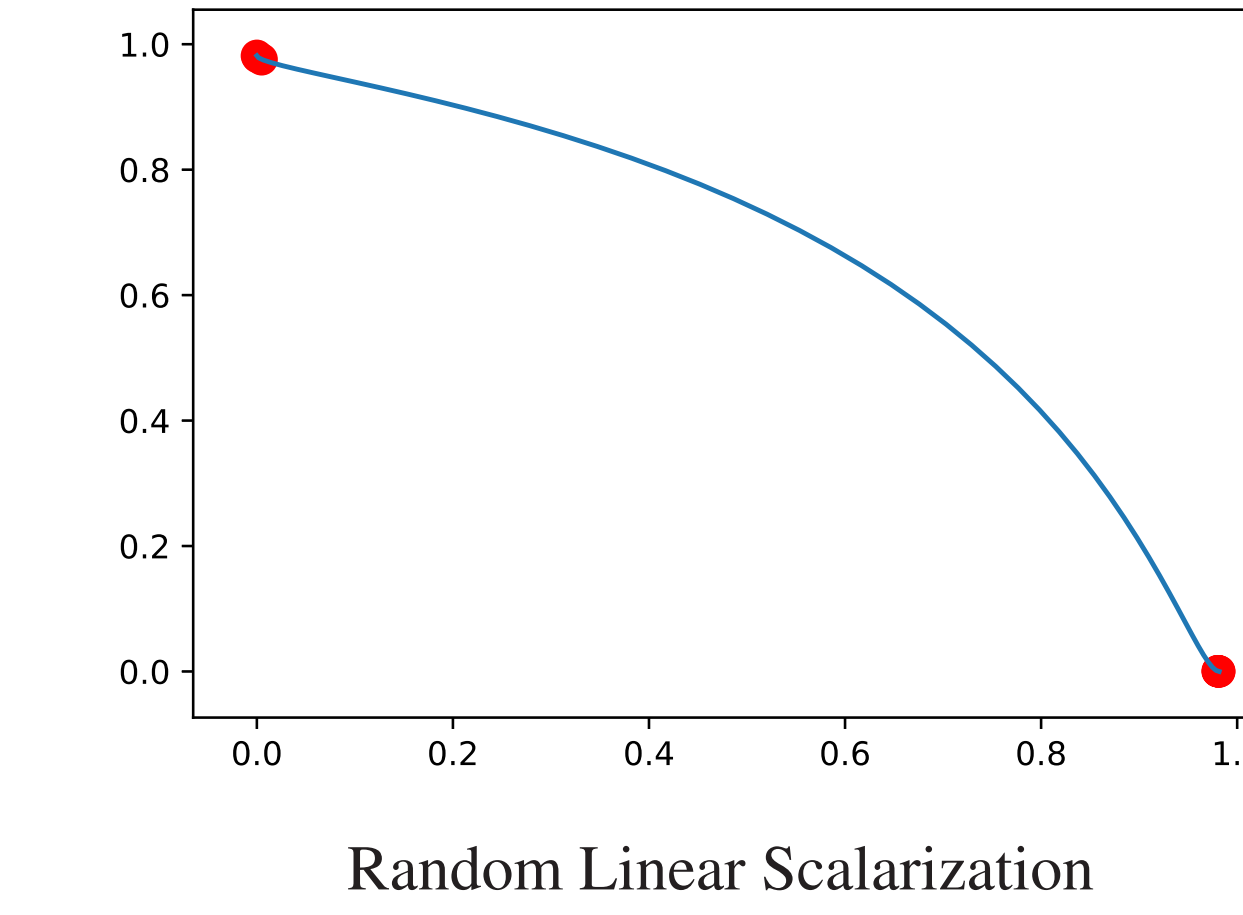
$$\mathcal{L}(\theta_t) = \sum_{i=1}^m \alpha_i \mathcal{L}_i(\theta_t),$$

$$\text{where } \alpha_i = \lambda_i + \sum_{j \in I_{\epsilon}(\theta)} \beta_j (u_{ji} - u_{ki}), \quad (7)$$

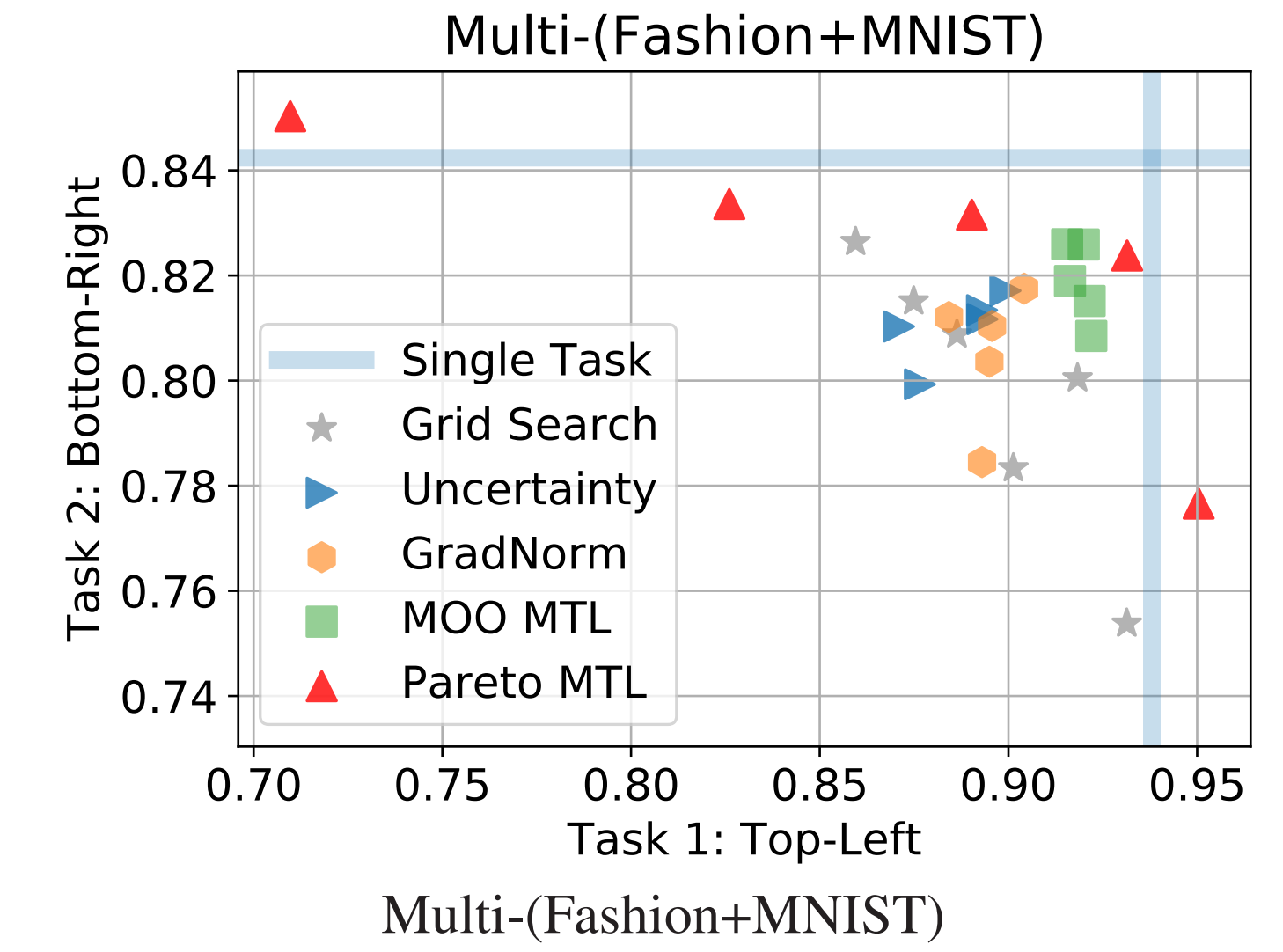
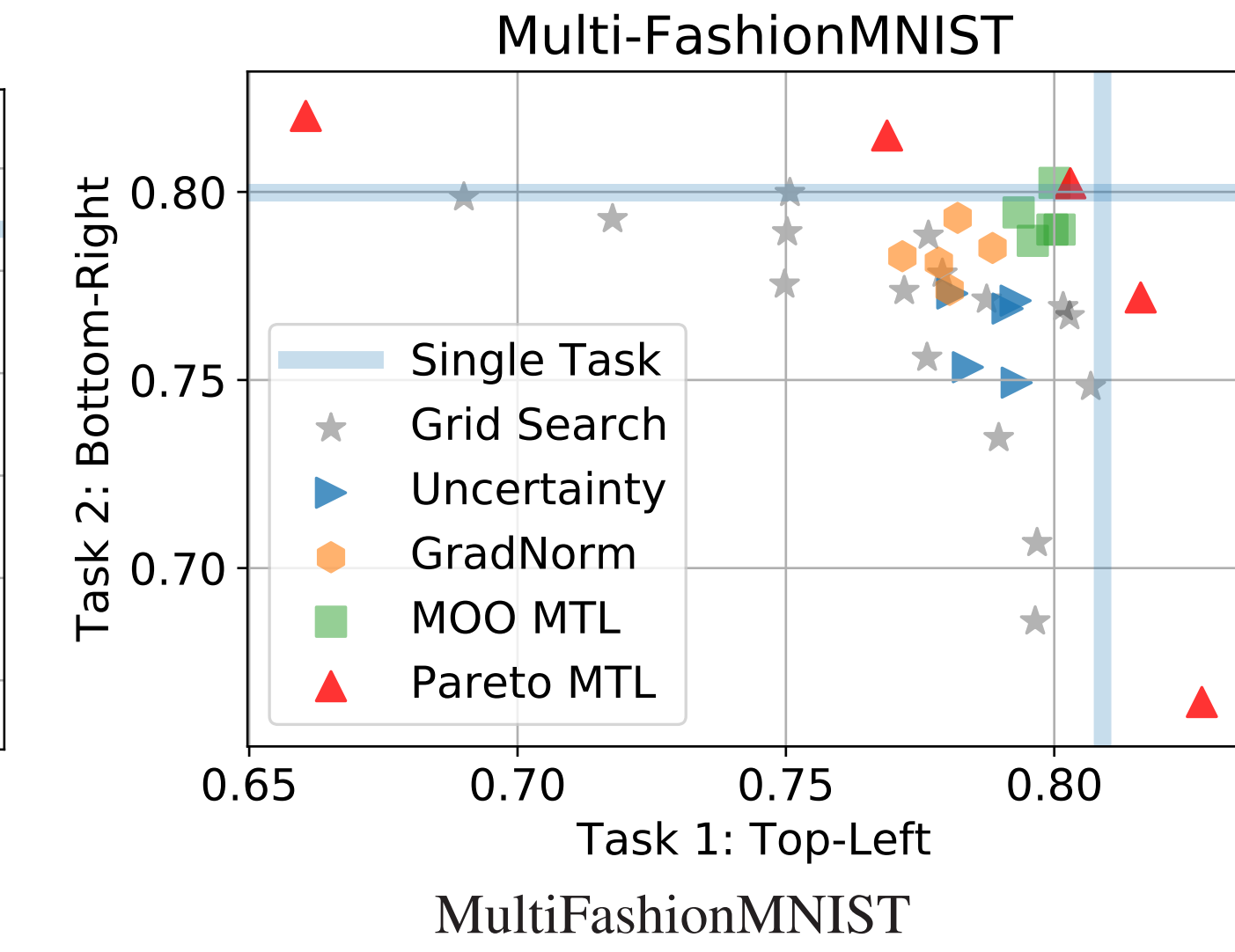
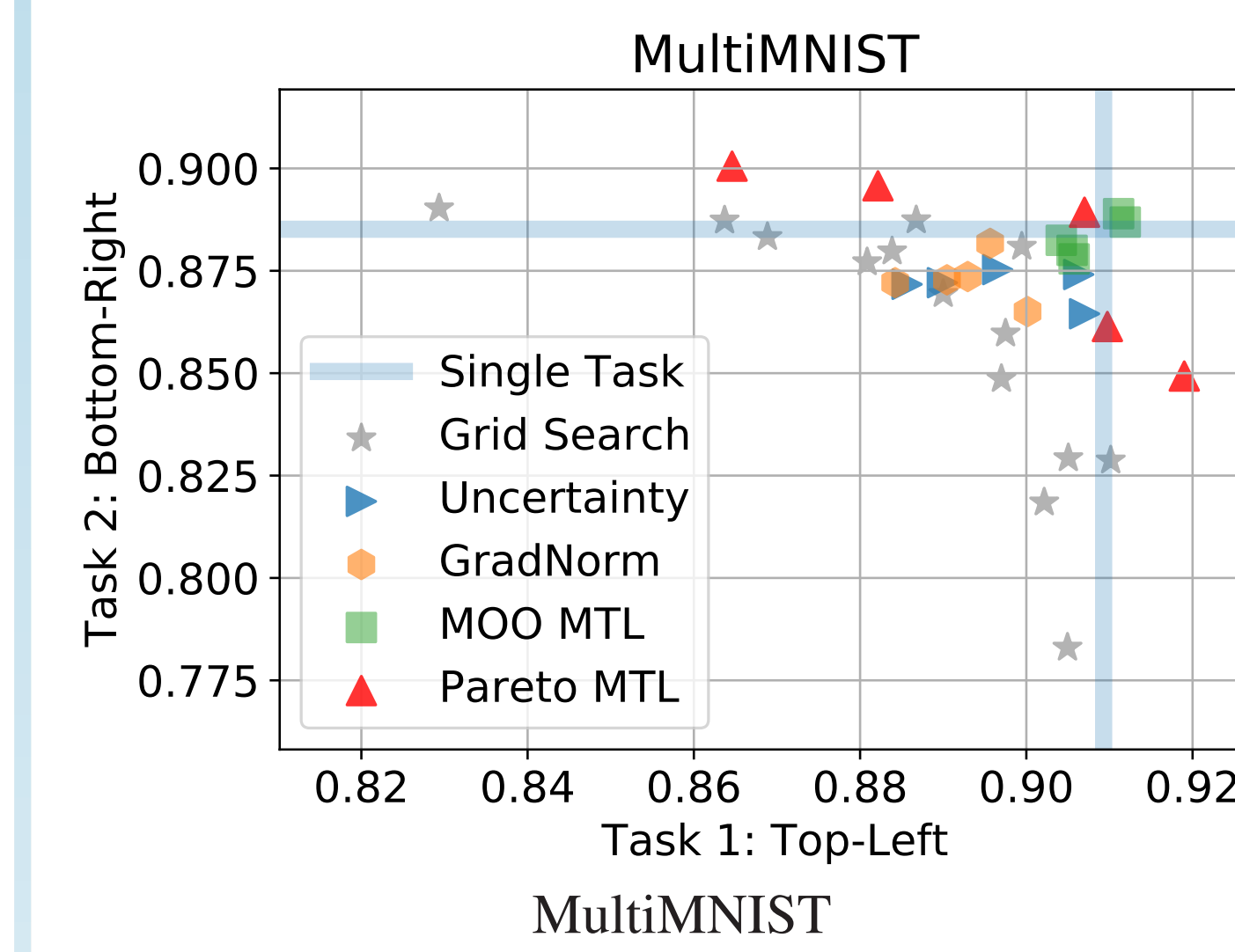
where λ_i and β_j are obtained by solving the dual form of problem (6) with assigned reference vector u_k .

Experiments & Results

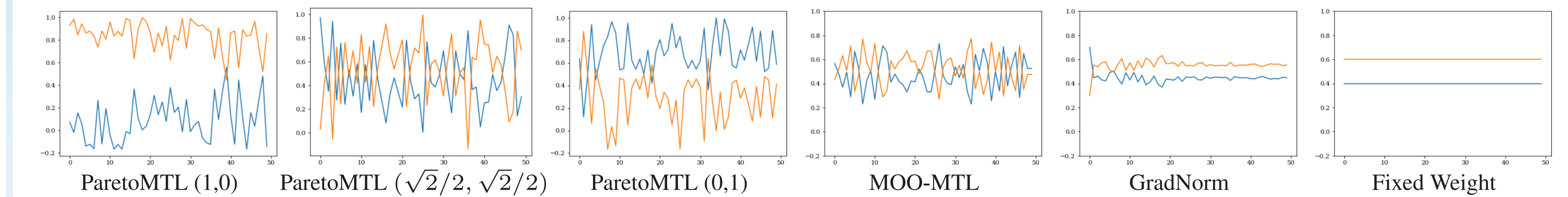
Synthetic Example with Concave Pareto Front:



MultiMNIST, MultiFashionMNIST and Multi-(Fashion+MNIST):



Adaptive Weights for Different Algorithms:



Pareto MTL with Randomly Generated Preference Vectors:

