

DreamShard: Generalizable Embedding Table Placement for Recommender Systems

Daochen Zha[†], Louis Feng[‡], Qiaoyu Tan^{*}, Zirui Liu[†], Kwei-Herng Lai[†], Bhargav Bhushanam[‡], Yuandong Tian[‡], Arun Kejariwal, Xia Hu[†]

[†]Rice University
[‡]Meta Platforms, Inc.
^{*}Texas A&M University

Background of Embedding Tables

- **What is embedding table?** Embedding learning is a core technique for modeling categorical features in deep recommendation models. It maps sparse categorical features into dense vectors.
- **What is the problem?** Industrial recommendation models demand an extremely large number of parameters for embedding tables, requiring multi-terabyte memory. We have to partition the tables and put them on multiple GPUs.

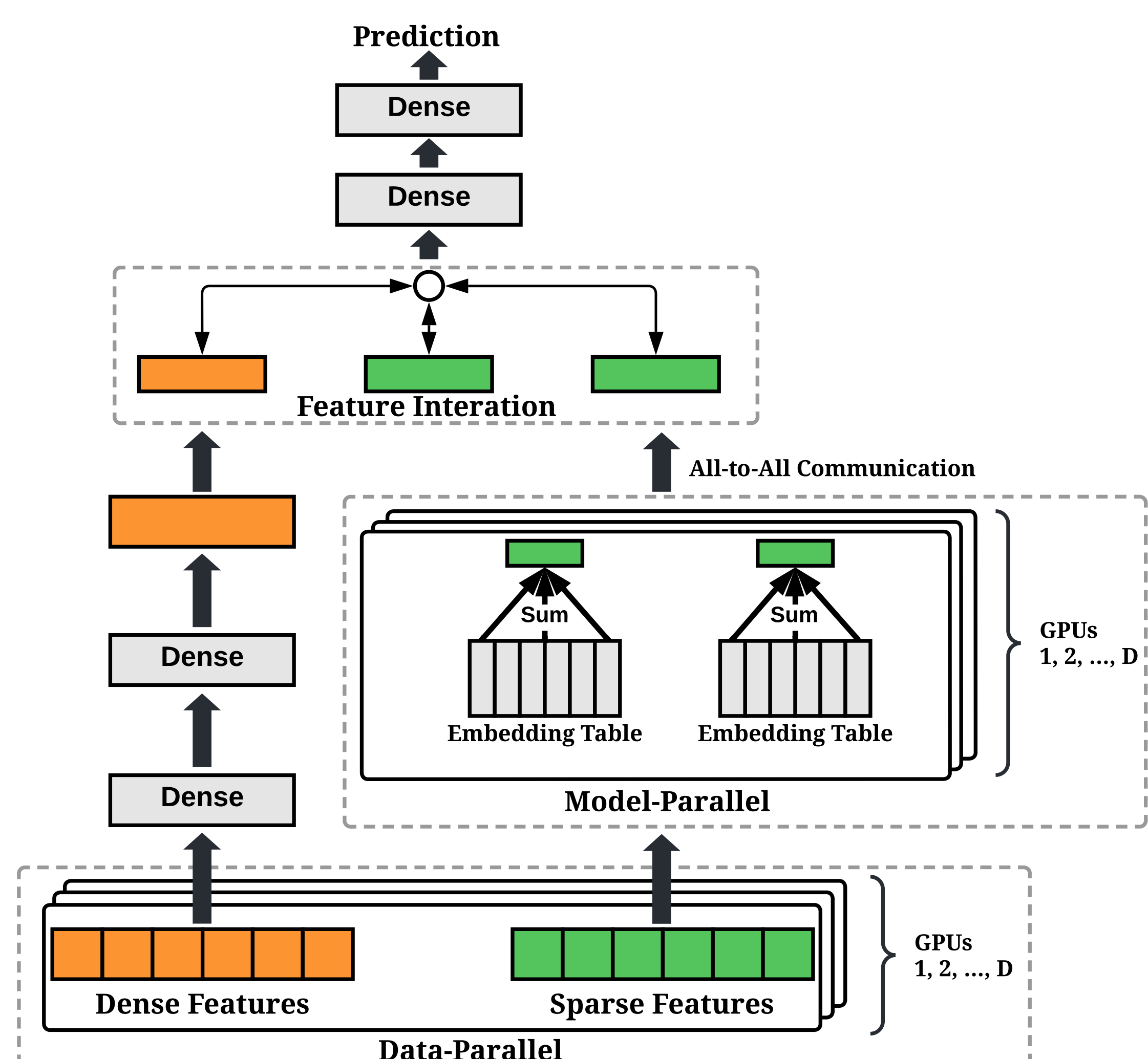


Figure: A typical recommendation model with dense and sparse features. The system exploits a combination of model parallelism (i.e., the embedding tables are partitioned into different devices) and data parallelism (i.e., replicating MLPs on each device and partitioning training data into different devices). The embedding vectors obtained from embedding lookup are appropriately sliced and transferred to the target devices through an all-to-all communication.

Embedding Table Placement Problem

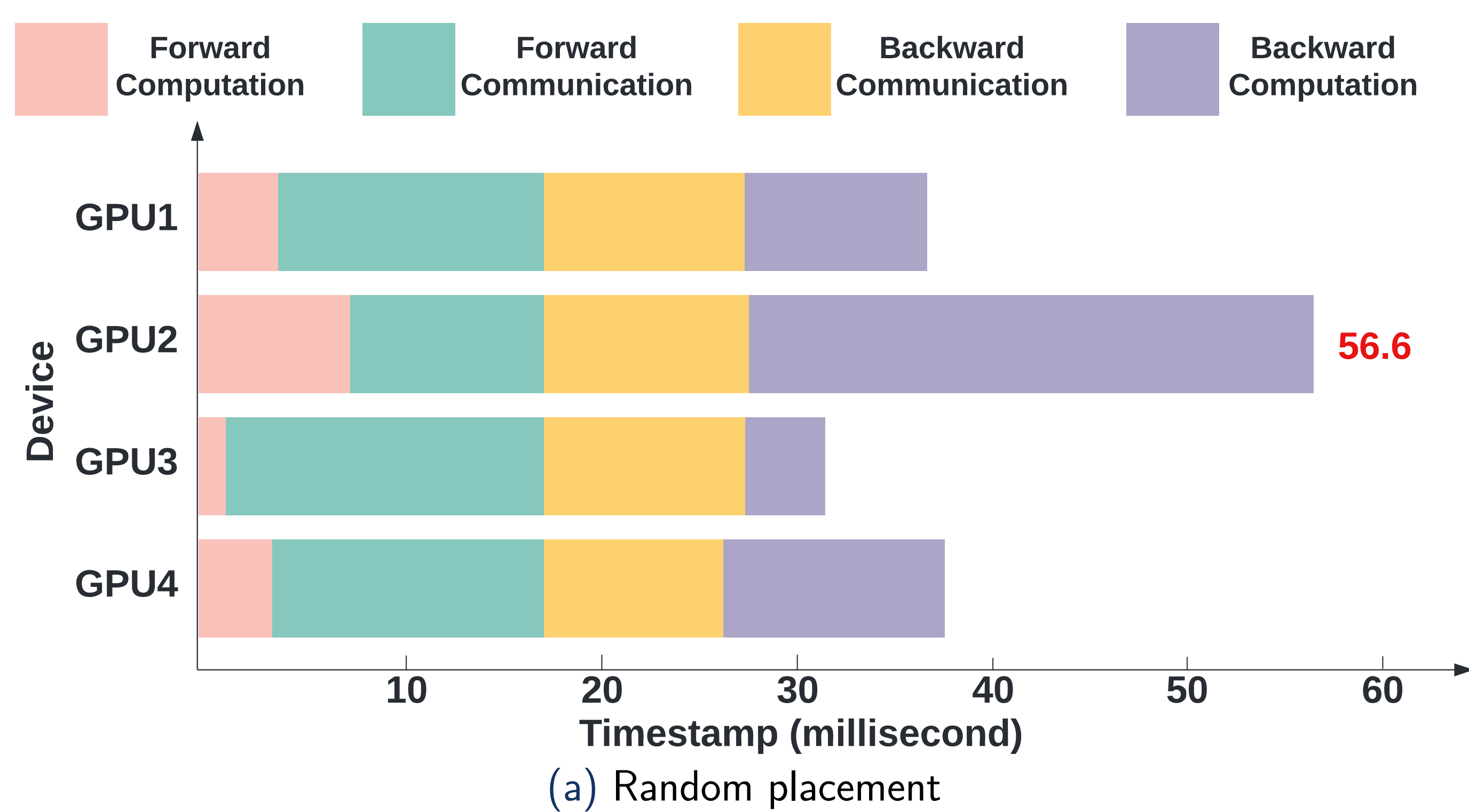


Figure: Naively placing tables will easily lead to imbalances.

- **Objective.** Given a number of embedding tables, we aim to find the optimal strategy to place them to minimize the max cost.

Challenges

- It is challenging to estimate the costs. The total cost of multiple tables in a shard is not the sum of the single table costs within the shard due to parallelism and operator fusion.
- The partitioning problem is known to be NP-hard.

DreamShard Framework

We present DreamShard, a reinforcement learning framework based on estimated MDP, for embedding table placement.

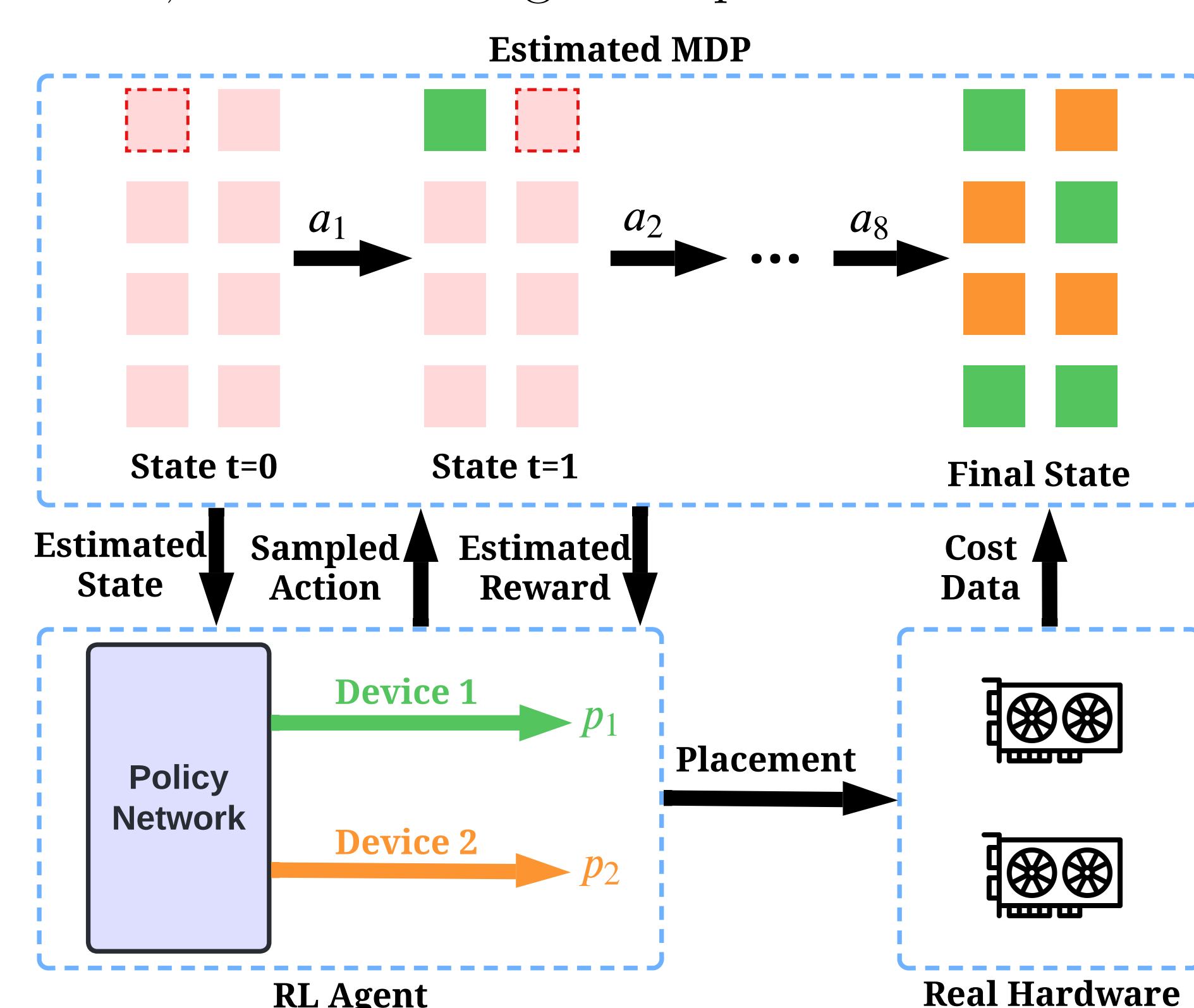


Figure: DreamShard framework. The agent interacts with the estimated MDP, which is trained with the cost data collected from GPUs.

Results

