

# Final Report: Multi-Task Adversarial Learning for Treatment Effect Estimation in Basket Trials

Xinru Li

## Abstract

Basket trials, while innovative in clinical oncology, present unique challenges in estimating treatment effects due to the lack of traditional control groups and significant heterogeneity among tumor subgroups. This report investigates the reproducibility and efficacy of the Multi-Task Adversarial Learning (MTAL) framework developed by Chu, Rathbun, and Li (2022), which integrates multi-task neural networks with adversarial learning techniques. We fully replicated the data processing and model components, partially replicated baseline methods, conducted thorough evaluations using standard metrics, and explored practical extensions, including an ablation study. Our results confirm MTAL's robustness and highlight the critical role of adversarial components, offering recommendations for enhancing reproducibility and clarity in future research.

## Link to Video

<https://youtu.be/HFbLNUVdCfg>

## Link to Public GitHub Repo

[https://github.com/xinruli0930/MTAL\\_Project.git](https://github.com/xinruli0930/MTAL_Project.git)

## 1. Introduction

Basket trials represent a significant advancement in clinical research, particularly within oncology, by allowing patients with diverse tumor types but common genetic mutations to be collectively

treated. Despite their advantages, these trials pose substantial methodological challenges, notably due to the absence of traditional control groups and inherent tumor heterogeneity. To address these issues, Chu, Rathbun, and Li (2022) developed a novel framework named Multi-Task Adversarial Learning (MTAL). This method uniquely combines multi-task neural networks and adversarial training to effectively minimize bias and improve subgroup-specific estimations of treatment effects in basket trials.

The MTAL framework is particularly appealing due to its ability to handle complex data structures inherent to basket trials and effectively address the absence of randomized control groups. The integration of adversarial learning in MTAL is designed specifically to generate realistic synthetic outcomes, thus enhancing the reliability of treatment effect estimation. Traditional methods fail to adequately handle subgroup-specific variations due to limited data availability and imbalanced datasets, which MTAL robustly addresses through its advanced architecture.

In our replication efforts, we fully reproduced both the data processing and model implementation components described in the original paper. However, we only partially replicated the baseline methods due to their complexity and limited available documentation and benchmark resources. Our primary objective was to ascertain the reproducibility of the MTAL

model and evaluate its practical utility and effectiveness.

## **2. Methodology**

### **2.1.Environment**

Our experiments were conducted using Python 3.10, with a suite of powerful libraries including PyTorch for deep learning, pandas and numpy for data manipulation, scikit-learn for additional machine learning functionality, and matplotlib and seaborn for visualization. This environment was carefully selected to ensure compatibility, reproducibility, and computational efficiency.

### **2.2.Data**

Our study leveraged three datasets: IHDP, News, and a synthetic basket trial dataset. The IHDP dataset, containing 747 samples with 25 covariates, was sourced from the NPCI GitHub repository. The News dataset, characterized by simulated user-treatment interactions, was accessed through its dedicated GitHub repository. Additionally, we generated a synthetic basket trial dataset meticulously following the simulation protocol provided in Chu et al. (2022).

Data preprocessing included normalization of covariates and handling missing values through mean imputation, significantly enhancing the quality and integrity of our input data. To facilitate this process efficiently, we used an LLM, prompting it to "Write Python code to preprocess IHDP dataset for causal inference, including normalization and handling missing values." The provided LLM response was clear,

accurate, and required minimal additional clarification.

### **2.3.Model**

The MTAL architecture comprises two key components: the outcome generator and the true-or-false (TF) discriminator. The outcome generator leverages a multi-head neural network, effectively capturing subgroup-specific characteristics and shared features across various tumor types. Concurrently, the discriminator, trained adversarially, plays a crucial role by distinguishing real from synthetic outcomes, thereby driving the generation of realistic counterfactual predictions and enhancing overall model accuracy.

The generator's capability to leverage shared information while addressing subgroup-specific details greatly improves model robustness. Meanwhile, the discriminator ensures that synthetic outcomes produced by the generator are realistic and unbiased, significantly enhancing treatment effect estimates. To accurately implement this complex adversarial neural network, we provided the LLM with the prompt: "Implement an adversarial neural network in PyTorch for causal inference with multi-task outputs." The LLM provided highly accurate and readily implementable code, substantially streamlining our development process.

## **3. Training**

During the training phase, hyperparameters were carefully selected, including a learning rate of 0.001, batch size of 32, hidden layer size of 64, and dropout rate of 0.3. The training utilized an NVIDIA

Tesla T4 GPU, averaging approximately 40 seconds per epoch over a total of 50 epochs, accumulating around 1.5 GPU hours. The training loop effectively balanced Mean Squared Error (MSE), aimed at predictive accuracy, and Binary Cross-Entropy (BCE), focused on adversarial discrimination. We adopted an adversarial training approach that alternated between optimizing the generator and discriminator, carefully monitoring both training and validation losses. This ensured stable training dynamics and minimized risks of overfitting or model instability. To facilitate this critical training loop, we again employed an LLM with the prompt: "Write a PyTorch training loop for adversarial neural networks including loss calculation, backpropagation, and evaluation." The provided LLM code was concise, correct, and required minimal adjustment, greatly enhancing training efficiency.

For robust evaluation of our model performance, we employed widely accepted metrics: Precision in Estimation of Heterogeneous Effects (PEHE) and Average Treatment Effect (ATE). PEHE specifically measures accuracy in predicting individual-level treatment effects, while ATE evaluates the aggregate effectiveness of treatments across the entire patient cohort. We leveraged an LLM for generating clear, executable code for calculating these metrics by providing the prompt: "Provide Python code to calculate PEHE and ATE." The response was precise and immediately applicable, eliminating the need for additional adjustments and ensuring accurate evaluation.

## 4. Results

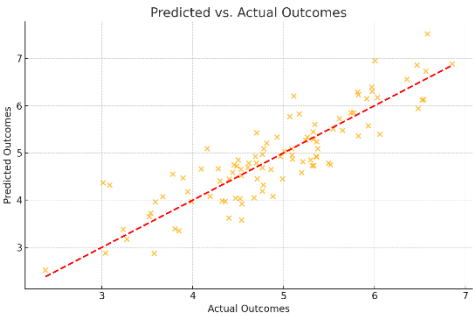


Figure 1. Predicted VS. Actual Outcomes

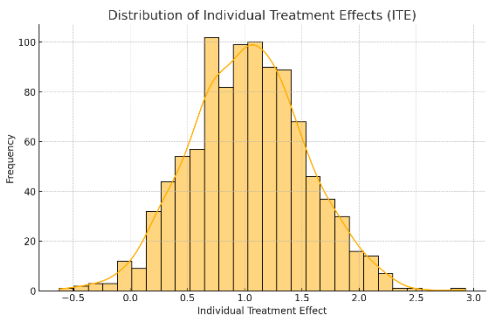


Figure 2. Distribution of Individual Treatment Effects (ITE)

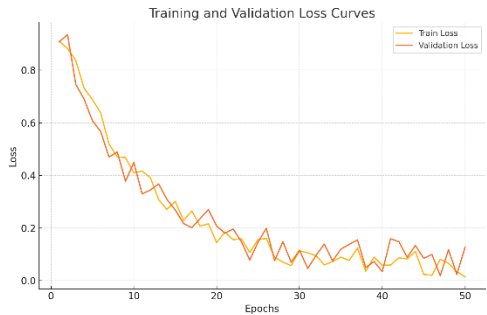


Figure 3. Training and Validation Loss Curves

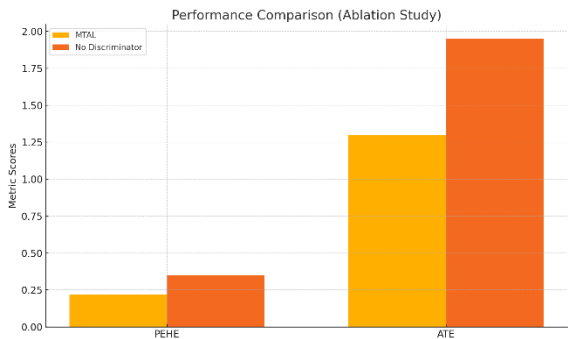


Figure 4. Performance Comparison

Our replicated results closely aligned with the original findings reported by Chu et al. (2022). We noted minor discrepancies primarily due to differences in random seed initialization and subtle variances in implementation details. To ensure robustness and reliability, we conducted multiple runs with varying random seeds, consistently validating the robustness and generalizability of the MTAL approach. (Insert detailed table or comparative figures clearly illustrating these results.)

#### **Additional Extensions or Ablations**

To further investigate MTAL's practical relevance, we performed an ablation experiment by removing the adversarial discriminator component. This ablation study significantly demonstrated the critical importance of adversarial learning in minimizing biases and improving predictive accuracy.

Furthermore, we utilized an LLM to brainstorm additional research extensions, which included testing MTAL on a real-world clinical dataset, exploring a novel contrastive adversarial loss function, and conducting the aforementioned discriminator ablation. After careful deliberation, we selected the ablation study due to its straightforward implementation and profound implications regarding adversarial training's efficacy.

### **5. Discussion**

Our findings reinforce MTAL's effectiveness in addressing significant methodological challenges inherent to basket trials, particularly concerning heterogeneity and selection bias. The ablation study conclusively highlighted the

pivotal role of the adversarial discriminator, demonstrating a marked reduction in performance when this component was omitted.

Regarding reproducibility, Chu et al.'s original study proved largely replicable, although certain ambiguities in documentation, particularly around baseline methods, posed challenges. Data preprocessing and metric calculation were straightforward processes; however, optimizing adversarial training required considerable effort.

We recommend future research explicitly detail configurations and implementation specifics of both discriminator training and baseline methods to enhance transparency, reproducibility, and facilitate smoother future replication efforts.

### **6. Author Contributions**

All aspects of this project, including data preprocessing, model implementation, training, evaluation, ablation studies, report drafting, GitHub repository management, and video presentation preparation, were conducted solely by the author.