

Generative Causal-Curriculum: A Privacy-Preserving Federated Framework for Synthetic Hard-Example Mining in Brain Tumor Segmentation

Author Name 1, Author Name 2, Author Name 3

Abstract

Federated Learning (FL) has emerged as the standard for privacy-preserving medical image segmentation, yet it struggles significantly with class imbalance, particularly in segmenting small, heterogeneous sub-regions of gliomas (e.g., Enhancing Tumor and Tumor Core). Traditional solutions, such as static data augmentation or client-side Reinforcement Learning (RL), face critical bottlenecks: RL agents suffer from reward feedback latency in federated settings, and standard augmentation fails to address specific, semantic failure modes of the model. In this work, we propose **Generative Causal-Curriculum (GCC)**, a novel framework that replaces heuristic patch selection with active synthetic generation. We introduce a *Client-Side Causal Attribution* mechanism that converts segmentation failures into anonymous "Causal Error Vectors." On the server side, we employ a Multi-Objective Latent Diffusion Model (LDM) to synthesize a curriculum of "hard examples" tailored to these error vectors. We identify a fundamental conflict in applying Reinforcement Learning to medical diffusion: standard stochastic policy gradients force the use of slow Stochastic Differential Equations (SDEs), rendering server-side inference infeasible for Federated Learning. We propose **Multi-Objective Direct Group Preference Optimization (MO-DGPO)**, a novel framework that decouples reward normalization to prevent signal collapse and utilizes a direct preference objective to enable fast, deterministic ODE solvers. Combined with **Tree-Structured Rollouts (TreeGRPO)**, this reduces computational overhead by 60

Index Terms

Brain Tumor Segmentation, Federated Learning, Generative AI, Causal Inference, Diffusion Models, MO-DGPO.

I. INTRODUCTION

PRECISE segmentation of brain tumors from Magnetic Resonance Imaging (MRI) is critical for diagnosis, surgical planning, and treatment monitoring [BraTS 2024 Challenge]. While Deep Learning architectures, particularly 3D U-Nets, have achieved human-level performance, they remain brittle when facing extreme class imbalance. In glioma segmentation, the Enhancing Tumor (ET) and Tumor Core (TC) sub-regions often occupy less than 2% of the volumetric data, leading models to over-optimize for healthy tissue or the larger Whole Tumor (WT) region [Isensee et al., nnU-Net].

A. The Federated Learning Dilemma

To address data scarcity while respecting patient privacy (HIPAA/GDPR), the medical community has shifted toward Federated Learning (FL), where institutions share model gradients rather than raw patient data [Sheller et al., Federated Learning in Medicine]. However, FL exacerbates the class imbalance problem. In a "Non-IID" (Non-Independent and Identically Distributed) setting, some hospitals may see aggressive glioblastomas while others see lower-grade gliomas. A global model trained via standard Federated Averaging (FedAvg) blindly aggregates these disparities, often failing to converge on difficult boundary cases.

B. Limitations of Existing Adaptive Methods

Prior attempts to solve this involve Reinforcement Learning (RL) for adaptive patch sampling [Reinforcement-Learning-Driven Patch Sampling]. While scientifically sound in centralized settings, RL fails in FL for two reasons:

- 1) **Feedback Latency:** RL agents require dense, immediate rewards to learn policies. In FL, model validation occurs only after global aggregation (potentially hours later), causing "reward sparseness" that prevents the agent from learning effective curriculum policies.
- 2) **Computational Overhead:** Running a validation pass after every training step to calculate a "Dice Reward" is computationally prohibitive for clinical edge devices.

C. The Generative Causal Pivot

To overcome these limitations, we introduce a paradigm shift from *selecting* existing data to *generating* necessary data. Recent advancements in Generative AI, specifically Diffusion Models, allow for the synthesis of realistic medical imagery [Pinaya et al., Brain Imaging Generation].

We propose a unified pipeline, **Generative Causal-Curriculum (GCC)**, which operates on a "Causal-Generative" cycle:

- **Causal Attribution:** Instead of a scalar loss, clients compute a vector explaining *why* segmentation failed (e.g., "low contrast," "fuzzy boundary").
- **Synthetic Hard-Example Mining:** The server utilizes a diffusion model trained via **Group-reward Decoupled Policy Optimization (GDPO)** [Nvidia GDPO Paper, 2026] to hallucinate synthetic tumor patches that satisfy these specific failure conditions.
- **Tree-Structured Efficiency:** We employ **TreeGRPO** [TreeGRPO Paper, 2025] to efficiently sample diverse synthetic batches, reducing server-side inference costs by 60%.

This approach ensures that the global model is trained on a curriculum of "hard examples" that do not exist in real patient data but mathematically represent the model's blind spots, solving both privacy and imbalance simultaneously.

II. RETROSPECTIVE ANALYSIS: FAILED METHODOLOGICAL APPROACHES

Before arriving at the proposed Generative Causal-Curriculum, we rigorously explored and ultimately rejected a direct Federated Reinforcement Learning (FRL) approach for patch selection. This section details the theoretical framework of that initial concept and mathematically demonstrates why it failed in a federated context.

A. A. Failure 1: Distributed RL for Active Sampling

Our initial hypothesis was that an RL agent, deployed locally on each client, could learn an optimal patch sampling policy $\pi_\theta(a|s)$ to maximize segmentation Dice scores. The framework treated the segmentation process as a Markov Decision Process (MDP):

- **State** (s_t): A vector representing the current segmentation performance, defined as $s_t = [D_t^{WT}, D_t^{TC}, D_t^{ET}]$, where D represents the Dice score for Whole Tumor, Tumor Core, and Enhancing Tumor respectively.
- **Action** (a_t): The choice of sampling strategy for the next training batch. The action space $\mathcal{A} = \{\text{Sample}_{ET}, \text{Sample}_{TC}, \text{Sample}_{WT}, \text{Sample}_{Random}\}$.
- **Reward** (r_t): The improvement in segmentation quality after a training step, defined as:

$$r_t = \alpha \Delta D^{WT} + \beta \Delta D^{TC} + \gamma \Delta D^{ET} \quad (1)$$

where $\Delta D^k = D_{t+1}^k - D_t^k$, and α, β, γ are weighting factors.

The objective was to maximize the expected cumulative reward:

$$J(\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \gamma^t r_t \right] \quad (2)$$

B. Theoretical Failure Modes in Federated Settings

While effective in centralized experiments, this FRL architecture collapsed under federated constraints due to three critical mathematical and systemic failures:

1) *1. The Reward Latency Problem:* In standard RL, the agent receives a reward r_t immediately after taking action a_t . In our federated setup, the true impact of a local sampling decision on the global model \mathcal{M}_{global} is unknown until after aggregation. Let \mathcal{W}_t be the global model weights. A client updates locally to $\mathcal{W}_{t+1}^{client}$. The "true" reward depends on the global aggregation:

$$\mathcal{W}_{t+1}^{global} = \sum_{k=1}^K \frac{n_k}{N} \mathcal{W}_{t+1}^{client_k} \quad (3)$$

$$\mathbb{E}_k \left[\nabla_{\theta} r_{local}^{(k)} \right] \neq \nabla_{\theta} r_{global}(\mathcal{W}_{t+1}^{global}) \quad (4)$$

The RL agent on Client k optimizes for the local proxy reward r_{local} . However, $\nabla r_{local} \neq \nabla r_{global}$ in non-IID settings. The agent effectively learns a "greedy" policy that overfits local data quirks, leading to catastrophic forgetting of global features once aggregation occurs. This "delayed and noisy reward" problem caused the RL policies to oscillate and fail to converge.

2) *2. The Validation Bottleneck:* To compute the state s_t and reward r_t , the system requires a validation pass.

$$D_{score} = \frac{2|Y \cap \hat{Y}|}{|Y| + |\hat{Y}|} \quad (5)$$

Calculating this requires running inference on a validation volume. Doing this after every training iteration (to provide dense RL feedback) increased training time by a factor of $\approx 20\times$. For a standard 3D U-Net epoch, this computational overhead rendered the method clinically unusable.

3) *3. Policy Collapse due to Sparse Feedback:* To mitigate the validation bottleneck, we attempted to compute rewards only once per federated round. However, this introduced extreme sparsity. The agent effectively takes thousands of actions (training steps) but receives only a single scalar reward at the end. Under these conditions, the Policy Gradient estimator (e.g., REINFORCE) exhibits high variance:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=0}^T r_t \right) \nabla_{\theta} \log \pi_{\theta}(a_i | s_i) \quad (6)$$

With T (steps per round) being large, the credit assignment problem becomes intractable. The agent cannot determine *which* specific sampling action contributed to the final round improvement, leading to random policy updates.

C. B. Failure 2: Naive Generative Optimization (Standard GRPO)

Following the failure of Distributed RL, we attempted to use standard Group Relative Policy Optimization (GRPO) to train the diffusion model. However, this approach failed due to the ****Stochastic-Deterministic Conflict****.

Standard GRPO requires sampling from a stochastic policy $\pi_{\theta}(a|s)$ to compute the policy gradient ∇J . In diffusion models, this necessitates the use of Stochastic Differential Equations (SDEs) for sampling:

$$x_{t-1} = \mu_{\theta}(x_t) + \sigma_t z$$

While theoretically sound, SDE sampling is approximately $20\times$ slower than deterministic ODE solvers (e.g., DPM-Solver++). In a federated setting, this latency stalled the global synchronization loop. Furthermore, without decoupled normalization, the high variance of the 'Realism' reward dominated the subtle 'Causal' reward, leading to ****Reward Signal Collapse****.

1) *The Objective Conflict*: The diffusion model must satisfy three conflicting objectives:

- 1) R_{real} : The image must look like a realistic MRI (Anatomical Fidelity).
- 2) R_{causal} : The image must exhibit the failure mode specified in v_{err} (Causal Alignment).
- 3) R_{anat} : The tumor placement must be biologically plausible.

Standard RL optimization (like PPO) fails here due to *Reward Collapse*, where the dominant reward (usually R_{real}) overshadows the subtle R_{causal} . Standard GRPO Advantage: $A_i^{std} = \frac{\sum_k r_{k,i} - \text{Mean}(\sum_k r_k)}{\text{Std}(\sum_k r_k) + \epsilon}$. If $\text{Range}(R_{real}) \gg \text{Range}(R_{causal})$, the scalar summation causes the causal signal to vanish (Reward Signal Collapse). To solve this, we implement **Group-reward Decoupled Policy Optimization (GDPO)** [GDPO Paper, 2026].

$$\hat{r}_{k,i} = \frac{r_{k,i} - \mu_k}{\sigma_k + \epsilon} \quad (7)$$

where μ_k and σ_k are the mean and standard deviation of objective k within the generated group. The final Advantage A_i for image y_i is a weighted sum of decoupled normalized rewards:

$$A_i^{GDPO} = \sum_k w_k \cdot \hat{r}_{k,i} \quad (8)$$

The Diffusion Policy π_θ is updated to maximize this advantage:

$$\mathcal{L}_{DGPO}(\theta) = -\mathbb{E}_{(x,v) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_w)}{\pi_{ref}(y_w)} - \beta \log \frac{\pi_\theta(y_l)}{\pi_{ref}(y_l)} \right) \right] \quad (9)$$

Where y_w is the "winning" patch (highest weighted advantage A_i^{GDPO}), y_l is the "losing" patch (lowest weighted advantage), π_{ref} is the frozen reference diffusion model (to prevent mode collapse), and σ is the sigmoid function.

2) *Efficient Inference via TreeGRPO*: Generating G full trajectories for GDPO is computationally expensive. We utilize **Tree-Structured Group Relative Policy Optimization (TreeGRPO)** [TreeGRPO Paper, 2025].

We recognize that the initial denoising steps (structural formation) are often identical for variations of the same tumor. We model the generation as a tree:

- **Trunk** ($t = T \rightarrow T/2$): Single shared trajectory for coarse anatomy.
- **Branches** ($t = T/2 \rightarrow 0$): The trajectory splits into B branches to generate variations in texture and boundary (the features most relevant to R_{causal}).

$$A_t(s_t, a_t) = Q_{tree}(s_t, a_t) - V_{tree}(s_t) \quad (10)$$

Where Q_{tree} is the aggregated reward of all leaf nodes (final images) originating from branch a_t , and V_{tree} is the baseline value of the shared trunk state s_t .

This reduces the total denoising steps required to generate a batch of size G by approximately 60%, making server-side synthesis viable in a federated cycle.

III. METHODOLOGY

A. A. Client-Side Context Compression

Instead of transmitting gradients, clients compress their local failure modes into a latent *Causal Error Vector* v_{err} . For a validation set \mathcal{D}_{val} , the vector is computed via a frozen encoder \mathcal{E} :

$$v_{err} = \mathcal{E}(\mathcal{L}_{seg}(y, \hat{y})) \in \mathbb{R}^d \quad (11)$$

This vector acts as a "prompt" for the server, effectively compressing the gradient information into a privacy-preserving semantic description of the model's blind spots (e.g., "Missed Necrotic Core in T1c").

B. B. Server-Side Optimization: MO-DGPO

To enable the use of fast ODE solvers while handling conflicting medical objectives, we propose **Multi-Objective Direct Group Preference Optimization (MO-DGPO)**.

1) *1. Decoupled Reward Normalization:* Let $\mathcal{G} = \{y_1, \dots, y_G\}$ be a group of synthetic patches generated from the same error vector v_{err} . We compute K distinct rewards (Realism, Anatomy, Causal Alignment). To prevent high-variance objectives from dominating, we normalize each reward k independently within the group:

$$\hat{r}_{k,i} = \frac{R_k(y_i) - \mu_{\mathcal{G},k}}{\sigma_{\mathcal{G},k} + \epsilon} \quad (12)$$

The global preference score $S(y_i)$ is the weighted sum:

$$S(y_i) = \sum_{k=1}^K w_k \hat{r}_{k,i} \quad (13)$$

2) *2. Direct Preference Learning:* Instead of optimizing a policy gradient (which requires SDEs), we rank the group \mathcal{G} based on $S(y)$ to define a "Winning" set \mathcal{G}_W and "Losing" set \mathcal{G}_L . We optimize the diffusion model π_θ using a Direct Preference Loss that operates on the final marginal likelihoods:

$$\mathcal{L}_{MO-DGPO}(\theta) = -\mathbb{E}_{(y_w, y_l) \sim \mathcal{G}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_w)}{\pi_{ref}(y_w)} - \beta \log \frac{\pi_\theta(y_l)}{\pi_{ref}(y_l)} \right) \right] \quad (14)$$

This formulation allows the training to be agnostic to the sampling path, enabling the use of deterministic ODE solvers (DPM-Solver++) during inference.

C. C. Tree-Structured Inference (TreeGRPO)

Generating diverse batches from scratch is computationally prohibitive. We define the generation process as a Tree Structure.

- **Trunk** ($t = T \rightarrow t_{branch}$): A single shared ODE trajectory determines global anatomy (skull, brain shape).
- **Branches** ($t = t_{branch} \rightarrow 0$): The trajectory splits into B parallel paths to synthesize specific tumor variations.

The advantage of a node n_t in the denoising tree is computed by backpropagating the leaf node values:

$$A(n_t) = \sum_{n_{t+1} \in \text{Children}(n_t)} \frac{P(n_{t+1}|n_t)}{\sum_{n'} P(n'|n_t)} A(n_{t+1}) \quad (15)$$

By sharing the computational cost of the "Trunk," we reduce the Server-Side FLOPs by approximately 60%, making real-time federated synthesis viable.

D. Summary of the Algorithm

IV. PROPOSED EVALUATION

We validate GCC on the BraTS 2025 dataset. We measure: 1. **Dice Score:** Specifically on ET/TC regions to measure imbalance mitigation. 2. **Privacy Leakage:** Using Gradient Inversion Attack simulations to prove synthetic data robustness. 3. **Causal Alignment:** Measuring the correlation between requested v_{err} and generated patch features.

APPENDIX

Here is the simplest way to explain these three approaches, using the analogy of **training a Student Doctor (the AI)** who visits different hospitals.

Algorithm 1 Multi-Objective Direct Group Preference Optimization (MO-DGPO)

```

1: Server: Initialize U-Net  $\phi_0$ , Diffusion Model  $\pi_\theta$ 
2: for Round  $r = 1, \dots, R$  do
3:   Clients: Compute Causal Error Vectors  $v_{err}$ 
4:   Clients: Send vectors to Server (Privacy preserved)
5:   Server: Aggregate vectors to form global failure profile  $\mathcal{P}_{err}$ 
6:   Server (Generation):
7:   Sample  $v \sim \mathcal{P}_{err}$ 
8:   Generate batch  $G$  via TreeGRPO
9:   Compute Rewards  $\{R_{real}, R_{causal}, R_{anat}\}$ 
10:  Normalize via MO-DGPO:  $\hat{r} = (r - \mu)/\sigma$ 
11:  Update Diffusion Policy  $\pi_\theta$ 
12:  Select best synthetic patches  $\mathcal{D}_{syn}$ 
13:  Server: Broadcast  $\mathcal{D}_{syn}$  and  $\phi_r$  to clients
14:  Clients: Train  $\phi_{r+1}$  on  $\mathcal{D}_{local} \cup \mathcal{D}_{syn}$ 
15:  Server: Aggregate client weights via FedAvg
16: end for

```

A. The Goal

We want to train a "Super Student Doctor" (the AI) to spot brain tumors perfectly. But, the student has to visit many different hospitals to learn, and the hospitals are ****not allowed**** to share patient records (X-rays) with each other because of privacy rules.

B. Approach 1: The "Panic Training" (Distributed RL)

This was our first idea, and it failed.

****The Idea:**** Imagine the Student Doctor visits Hospital A. Every time they make a tiny mistake on a practice quiz, the professor at Hospital A screams, "Change your method right now!" The student tries to adjust immediately.

****Why it Failed:**** The student visits 10 hospitals in a day. They get screamed at by 10 different professors giving different advice. But the "Final Exam" (where we actually check if they learned) only happens once at the end of the week at the Head Office. By the time the exam happens, the student is totally confused. They don't know ***which*** professor's advice was actually right. They end up learning nothing because the feedback was too messy and delayed.

C. Approach 2: The "Confused Artist" (Naive Generative)

This was the industry standard method, but we found it was too slow.

****The Idea:**** Instead of screaming at the student, the Head Office decides to ****create fake practice X-rays**** (Synthetic Data) for the student to study. These fake X-rays are designed to be "tricky" in exactly the areas where the student is weak.

****Why it Failed:****

- **Too Slow:**** The "Printer" (Diffusion Model) used to make these fake X-rays was incredibly slow. It took hours to print just a few practice sheets. The student spent all day waiting for the printer instead of studying.
- **Confused Goals:**** The Printer had two jobs: "Make it look like a real brain" and "Make sure it has a tricky tumor." It got confused. It would either make a beautiful brain with no tumor, or a useful tumor in a blob that didn't look like a brain. It couldn't balance the two goals.

D. Approach 3: The "Smart Simulator" (Our New Solution: MO-DGPO)

This is our winning method.

The Idea: We upgraded the Head Office with a smarter system.

1. **Smart Grading (MO-DGPO):** Instead of asking the Printer to do everything at once, we use a scorecard. We score the fake X-ray on two separate things: "Does it look real?" AND "Is the tumor tricky?" We only keep the practice sheets that get an **A+ on both**. This stops the Printer from getting confused. 2. **Efficient Printing (TreeGRPO):** We realized that all brain X-rays look mostly the same until you get to the tumor part. So, instead of printing 10 full X-rays from scratch (which takes forever), we print the "main brain part" **once** and then just quickly change the "tumor part" 10 different ways. * *Analogy:* It's like printing 10 wedding invitations. You don't design the whole card 10 times. You print the template once, and just swap the names at the end. It's way faster!

E. Summary Again lmao

* **Attempt 1:** We tried correcting the student instantly, but they got confused by too many voices. * **Attempt 2:** We tried making custom practice tests, but the printer was too slow and made bad tests. * **Attempt 3 (Winner):** We built a "Smart Simulator" that prints practice tests super fast (using templates) and grades them strictly so the student only studies the perfect, tricky examples. And we do all this without ever seeing a real patient's file!