

# GPU Computing Coursework: CUDA-Accelerated Prioritized Experience Replay (PER) Sum-Tree Theory, Implementation, Tests, and Expected Results

Aditya Pandey

February 9, 2026

## Contents

<b>1 Overview</b>	<b>2</b>
<b>2 Reference Paper</b>	<b>2</b>
<b>3 Theory</b>	<b>2</b>
3.1 Prioritized Experience Replay (PER) . . . . .	2
3.2 Importance Sampling (IS) Weights . . . . .	2
3.3 Why a Sum-Tree / Segment Tree . . . . .	3
<b>4 Build and Runtime Requirements</b>	<b>3</b>
4.1 Prerequisites . . . . .	3
4.2 How the extension is compiled . . . . .	3
<b>5 Implementation Structure (Files)</b>	<b>3</b>
<b>6 Core Data Structure: Flat Sum-Tree Layout</b>	<b>4</b>
<b>7 CPU Reference: <code>cpu_tree.py</code></b>	<b>4</b>
7.1 Intent . . . . .	4
7.2 Functions (line-by-line intent) . . . . .	4
7.3 Full source . . . . .	4
<b>8 GPU Sum-Tree Wrapper: <code>tree.py</code></b>	<b>5</b>
8.1 Intent . . . . .	5
8.2 Key operations . . . . .	5
8.3 Full source . . . . .	5
<b>9 Baseline Sampler: <code>baseline_cumsum.py</code></b>	<b>7</b>
9.1 Intent . . . . .	7
9.2 Full source . . . . .	7
<b>10 Extension Loader: <code>_ext.py</code></b>	<b>7</b>
10.1 Intent . . . . .	7
10.2 Full source . . . . .	7
<b>11 Standalone Builder: <code>build_ext.py</code></b>	<b>8</b>
11.1 Intent . . . . .	8
11.2 Full source . . . . .	8

<b>12 C++ Binding Layer: per_ext.cpp</b>	<b>8</b>
12.1 Intent . . . . .	8
12.2 Full source . . . . .	9
<b>13 CUDA Kernels: per_kernels.cu</b>	<b>9</b>
13.1 Design intent and correctness . . . . .	9
13.2 Full source . . . . .	9
13.3 Expected results from kernels . . . . .	12
<b>14 Minimal GPU PER Replay Buffer: replay.py</b>	<b>12</b>
14.1 Intent . . . . .	12
14.2 Detailed behavior (function-by-function) . . . . .	12
14.3 Full source . . . . .	13
<b>15 Benchmarks: bench_per.py</b>	<b>15</b>
15.1 Intent . . . . .	15
15.2 Expected benchmark trends . . . . .	15
15.3 Full source . . . . .	16
<b>16 Example Training Script: train_cartpole.py</b>	<b>18</b>
16.1 Intent . . . . .	18
16.2 Expected outcome . . . . .	18
16.3 Full source . . . . .	18
<b>17 Tests: Correctness and Expected Results</b>	<b>21</b>
17.1 How to run . . . . .	21
17.2 test_update.py: Tree invariants . . . . .	21
17.3 test_sample.py: Sampling distribution . . . . .	22
17.4 test_per_replay.py: Replay sampling + updates . . . . .	22
<b>18 Note on per_ext.py Filename</b>	<b>23</b>
<b>19 Conclusion</b>	<b>24</b>

# 1 Overview

This report documents a CUDA-accelerated implementation of a **sum-tree** (flat segment tree) to support **Prioritized Experience Replay (PER)** for reinforcement learning. The project provides:

- A CPU reference sum-tree (`cpu_tree.py`),
- A GPU sum-tree wrapper (`tree.py`) backed by a PyTorch C++/CUDA extension (`per_ext.cpp`, `per_kernels.cu`),
- A baseline GPU sampler using `cumsum + searchsorted` (`baseline_cumsum.py`),
- A minimal GPU PER replay buffer (`replay.py`),
- Benchmarks (`bench_per.py`),
- Unit tests for updates, sampling distribution, and replay (`test_update.py`, `test_sample.py`, `test_per_replay.py`),
- An example CartPole training script using PER (`train_cartpole.py`).

The design goal is to make priority **updates** and **sampling** efficient on GPU without recomputing a full prefix-sum CDF every time.

## 2 Reference Paper

This implementation is based on the PER idea introduced in:

- T. Schaul et al., *Prioritized Experience Replay*, arXiv:1511.05952, 2015. DOI: 10.48550/arXiv.1511.05952.

## 3 Theory

### 3.1 Prioritized Experience Replay (PER)

In PER, transitions are sampled non-uniformly, typically proportional to a priority value derived from TD-error:

$$p_i \leftarrow (|\delta_i| + \epsilon)^\alpha$$

and the sampling probability becomes

$$P(i) = \frac{p_i}{\sum_k p_k}.$$

The exponent  $\alpha \in [0, 1]$  controls prioritization strength;  $\epsilon > 0$  prevents zero probability.

### 3.2 Importance Sampling (IS) Weights

PER biases the training distribution, so learning is corrected using importance sampling weights:

$$w_i = (N \cdot P(i))^{-\beta},$$

where  $N$  is the current buffer size and  $\beta$  anneals from  $\beta_0$  toward 1. Typically weights are normalized:

$$\tilde{w}_i = \frac{w_i}{\max_j w_j} \leq 1.$$

### 3.3 Why a Sum-Tree / Segment Tree

A baseline approach samples via a CDF:

$$\text{cdf}[j] = \sum_{k \leq j} p_k,$$

then binary-search for a uniform  $u \in [0, \sum p]$ . On GPU, `torch.cumsum + torch.searchsorted` is easy but recomputing `cdf` is  $\Theta(N)$ .

A sum-tree stores hierarchical partial sums so that:

- **update** of one leaf is  $\mathcal{O}(\log L)$ ,
- **sample** of one index is  $\mathcal{O}(\log L)$ ,

where  $L$  is the next power-of-two above capacity.

Sampling: draw  $u \sim \mathcal{U}[0, \text{total}]$  and descend:

```
if  $u \leq s_{\text{left}}$  : go left  
else :  $u \leftarrow u - s_{\text{left}}$ , go right
```

until reaching a leaf.

## 4 Build and Runtime Requirements

### 4.1 Prerequisites

- Python + PyTorch with CUDA support
- A CUDA-capable GPU and CUDA toolkit compatible with the PyTorch build
- A C++ compiler toolchain suitable for PyTorch extensions

### 4.2 How the extension is compiled

Two entry points exist:

- `per/_ext.py` compiles/loads at import-time using `torch.utils.cpp_extension.load`.
- `build_ext.py` compiles the extension as a standalone smoke test.

## 5 Implementation Structure (Files)

File	Responsibility
<code>cpu_tree.py</code>	CPU reference flat sum-tree (NumPy): update + sample + total.
<code>tree.py</code>	GPU sum-tree wrapper around the extension: update, sample, priority queries.
<code>baseline_cumsum.py</code> <code>_ext.py</code>	Baseline GPU sampler using <code>cumsum+searchsorted</code> . Loads/compiles CUDA extension at import time.
<code>per_ext.cpp</code>	PyBind11 bindings for CUDA functions.
<code>per_kernels.cu</code>	CUDA kernels: update, coalesced update, sample.
<code>replay.py</code>	Minimal PER replay buffer on GPU using <code>GpuSumTree</code> .
<code>bench_per.py</code>	Benchmarks comparing update/sample latency across implementations.
<code>train_cartpole.py</code>	Example training script on CartPole using <code>GPUReplayPER</code> .
<code>test_update.py</code>	Tests sum-tree update invariants (root sum, internal sums).

---

<code>test_sample.py</code>	Tests sampling distribution matches priorities (approx).
<code>test_per_replay.py</code>	Tests replay sampling shapes and priority update behavior.

---

## 6 Core Data Structure: Flat Sum-Tree Layout

Both CPU and GPU implementations use the same flat representation:

- Let capacity be  $C$  (not necessarily power-of-two).
- Let  $L = 2^{\lceil \log_2 C \rceil}$ .
- Tree array has size  $2L$  (index 0 unused, root at index 1).
- Leaves are stored at indices  $[L, L + C)$ .

Invariants:

$$\text{tree}[i] = \text{tree}[2i] + \text{tree}[2i + 1] \quad \text{for all internal } i,$$

and

$$\text{tree}[1] = \sum_{k=0}^{C-1} p_k.$$

## 7 CPU Reference: `cpu_tree.py`

### 7.1 Intent

`CpuSumTree` is a correctness reference. Its behavior should match the GPU version, making it useful for validation and debugging.

### 7.2 Functions (line-by-line intent)

- `__init__`: compute  $L$ , allocate tree array.
- `update(idx, new_p)`:
  - for each update: compute delta at leaf,
  - apply delta upwards to the root.
- `sample(u)`: descend from root to leaf using the left-sum rule.
- `total`: return root sum.

### 7.3 Full source

Listing 1: `cpu_tree.py`

```

1 import numpy as np
2
3 class CpuSumTree:
4     """
5         Flat sum-tree on CPU (numpy).
6         - capacity = N (not forced to power of 2)
7         - internal nodes stored in tree[1..2*L-1], leaves at [L..L+N-1]
8     """
9     def __init__(self, capacity: int):
10         self.capacity = int(capacity)
11         self.L = 1 << (self.capacity - 1).bit_length()
12         self.tree = np.zeros(2 * self.L, dtype=np.float32)
13

```

```

14     def update(self, idx: np.ndarray, new_p: np.ndarray):
15         # idx: [B] int64, new_p: [B] float32
16         for i, p in zip(idx, new_p):
17             pos = self.L + int(i)
18             delta = float(p) - float(self.tree[pos])
19             self.tree[pos] = float(p)
20
21             pos >= 1
22             while pos >= 1:
23                 self.tree[pos] += delta
24                 pos >= 1
25
26     def sample(self, u: np.ndarray) -> np.ndarray:
27         # u: [B] float32 in [0, tree[1])
28         out = np.empty(u.shape[0], dtype=np.int64)
29         for j, x0 in enumerate(u):
30             x = float(x0)
31             node = 1
32             while node < self.L:
33                 left = node << 1
34                 s_left = float(self.tree[left])
35                 if x <= s_left:
36                     node = left
37                 else:
38                     x -= s_left
39                     node = left + 1
40             leaf = node - self.L
41             if leaf >= self.capacity:
42                 leaf = self.capacity - 1
43             out[j] = leaf
44         return out
45
46     @property
47     def total(self) -> float:
48         return float(self.tree[1])

```

## 8 GPU Sum-Tree Wrapper: `tree.py`

### 8.1 Intent

`GpuSumTree` stores the tree as a CUDA `float32` tensor and delegates update/sample to the extension.

### 8.2 Key operations

- `update`: calls `per_ext.update(tree, idx, new_p, L)`.
- `sample`: draws stratified uniforms and calls `per_ext.sample`.
- `leaf_priorities`: reads leaf values directly from `tree`.
- `update_coalesced`: alternative update kernel for lower contention.

### 8.3 Full source

Listing 2: `tree.py`

```

1 import torch
2 from ._ext import per_ext

```

```

3  class GpuSumTree:
4      """
5          Flat segment tree:
6              - L = next power of two >= capacity
7              - tree is size 2*L, float32 on CUDA
8              - root at index 1
9              - leaves at [L : L+capacity)
10         """
11
12     def __init__(self, capacity: int, device: str = "cuda"):
13         assert capacity > 0
14         self.capacity = int(capacity)
15         self.L = 1 << (self.capacity - 1).bit_length()
16         self.device = torch.device(device)
17         self.tree = torch.zeros(2 * self.L, device=self.device, dtype=
18                               torch.float32)
19
20     @property
21     def total(self) -> torch.Tensor:
22         return self.tree[1]
23
24     def leaf_values(self) -> torch.Tensor:
25         return self.tree[self.L : self.L + self.capacity]
26
27     def update(self, idx: torch.Tensor, new_p: torch.Tensor):
28         """
29             idx: int64 CUDA tensor of shape [n] in [0, capacity)
30             new_p: float32 CUDA tensor of shape [n]
31         """
32         per_ext.update(self.tree, idx, new_p, self.L)
33
34     def sample(self, batch_size: int) -> torch.Tensor:
35         # total priority mass on GPU
36         total = self.tree[1].item()
37         if total <= 0.0:
38             raise RuntimeError("Cannot sample: total priority is 0")
39
40         # stratified sampling: u_i in [i*seg, (i+1)*seg)
41         seg = total / batch_size
42         u = (torch.arange(batch_size, device=self.tree.device, dtype=torch.
43                           float32) +
44               torch.rand(batch_size, device=self.tree.device, dtype=torch.
45                           float32)) * seg
46
47         return per_ext.sample(self.tree, u, self.L, self.capacity)
48
49     def leaf_priorities(self, idx: torch.Tensor) -> torch.Tensor:
50         """
51             idx: int64 CUDA tensor [B] in [0, capacity)
52             returns: float32 CUDA tensor [B] of leaf priorities
53         """
54         return self.tree[self.L + idx]
55
56     def total_priority(self) -> torch.Tensor:
57         """Return total priority mass (CUDA scalar tensor)."""
58         return self.tree[1]
59
60     def update_coalesced(self, idx, new_p):
61         per_ext.update_coalesced(self.tree, idx, new_p, self.L)

```

## 9 Baseline Sampler: baseline\_cumsum.py

### 9.1 Intent

This baseline recomputes a full CDF via `torch.cumsum` and uses `torch.searchsorted`. It is correct but typically slower for frequent updates due to  $\Theta(N)$  prefix sum recomputation.

### 9.2 Full source

Listing 3: baseline\_cumsum.py

```
1 import torch
2
3 class GpuCumsumSampler:
4     """
5         Baseline GPU sampler using cumsum + searchsorted.
6         Priorities stored as a flat CUDA tensor.
7     """
8     def __init__(self, capacity, device="cuda"):
9         self.capacity = capacity
10        self.device = device
11        self.priorities = torch.zeros(capacity, device=device, dtype=torch
12                                      .float32)
13
14    @torch.no_grad()
15    def update(self, idx, new_p):
16        self.priorities[idx] = new_p
17
18    @torch.no_grad()
19    def sample(self, batch_size):
20        cdf = torch.cumsum(self.priorities, dim=0)
21        total = cdf[-1]
22
23        seg = total / batch_size
24        u = (torch.arange(batch_size, device=self.device) +
25             torch.rand(batch_size, device=self.device)) * seg
26
27        idx = torch.searchsorted(cdf, u, right=False)
28        idx.clamp_(0, self.capacity - 1)
29        return idx
```

## 10 Extension Loader: \_ext.py

### 10.1 Intent

Uses PyTorch's `load()` to compile and load a C++/CUDA extension at import-time. It points to:

- `ext/per_ext.cpp`
- `ext/per_kernels.cu`

### 10.2 Full source

Listing 4: \_ext.py

```
1 from pathlib import Path
2 from torch.utils.cpp_extension import load
3
4 _THIS_DIR = Path(__file__).resolve().parent
```

```

5 _ROOT = _THIS_DIR.parent
6
7 # Unique name avoids collisions if you have multiple builds
8 _EXT_NAME = "per_ext"
9
10 per_ext = load(
11     name=_EXT_NAME,
12     sources=[
13         str(_ROOT / "ext" / "per_ext.cpp"),
14         str(_ROOT / "ext" / "per_kernels.cu"),
15     ],
16     extra_cuda_cflags=["-O3", "--use_fast_math"],
17     extra_cflags=["-O3"],
18     verbose=True,
19 )

```

## 11 Standalone Builder: build\_ext.py

### 11.1 Intent

Compiles the extension without importing the package, and performs a minimal CUDA availability check.

### 11.2 Full source

Listing 5: build\_ext.py

```

1 from torch.utils.cpp_extension import load
2 import torch
3
4 per_ext = load(
5     name="per_ext",
6     sources=["ext/per_ext.cpp", "ext/per_kernels.cu"],
7     extra_cuda_cflags=["--use_fast_math"],
8     extra_cflags=["-O3"],
9     verbose=True,
10 )
11
12 if __name__ == "__main__":
13     print("Built per_ext OK")
14     # quick smoke test
15     assert torch.cuda.is_available()
16     print("CUDA available:", torch.cuda.is_available())

```

## 12 C++ Binding Layer: per\_ext.cpp

### 12.1 Intent

Registers the CUDA entry points into a Python module via PyBind11:

- update(tree, idx, new\_p, L)
- update\_coalesced(tree, idx, new\_p, L)
- sample(tree, u, L, capacity)

## 12.2 Full source

Listing 6: per\_ext.cpp

```
1 #include <torch/extension.h>
2
3 // CUDA forward decls
4 void update_cuda(torch::Tensor tree,
5                   torch::Tensor idx,
6                   torch::Tensor new_p,
7                   int64_t L);
8
9 void update_coalesced_cuda(torch::Tensor tree,
10                           torch::Tensor idx,
11                           torch::Tensor new_p,
12                           int64_t L);
13
14 torch::Tensor sample_cuda(torch::Tensor tree,
15                           torch::Tensor u,
16                           int64_t L,
17                           int64_t capacity);
18
19 PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
20     m.def("update", &update_cuda, "Sum-tree update (CUDA)");
21     m.def("update_coalesced", &update_coalesced_cuda, "Sum-tree update (
22         warp-coalesced) (CUDA)");
23     m.def("sample", &sample_cuda, "Sum-tree sample (CUDA)");
24 }
```

## 13 CUDA Kernels: per\_kernels.cu

### 13.1 Design intent and correctness

**Update kernel.** Each thread updates one leaf:

1. Locate leaf `pos = L + idx[i]`.
2. Use `atomicExch` to set leaf to the new value and retrieve the old value.
3. Compute `delta = new - old`.
4. Propagate `delta` upward using `atomicAdd` on parents until reaching the root.

`atomicExch` makes duplicate indices in the same update batch safe.

**Coalesced update kernel.** Attempts to reduce atomic contention by grouping updates within a warp when threads share the same parent (uses warp intrinsics).

**Sample kernel.** Each thread performs a tree descent using its own uniform variate `u[i]`. The final leaf index is clamped to `capacity-1` if it lands in the padded region.

## 13.2 Full source

Listing 7: per\_kernels.cu

```
1 #include <torch/extension.h>
2 #include <cuda.h>
3 #include <cuda_runtime.h>
4 #include <ATen/cuda/CUDAContext.h>
```

```

5  static inline void check_cuda(torch::Tensor t, const char* name) {
6      TORCH_CHECK(t.is_cuda(), name, " must be a CUDA tensor");
7      TORCH_CHECK(t.is_contiguous(), name, " must be contiguous");
8  }
9
10 // GPU UPDATE (naive)
11 __global__ void update_kernel(
12     float* tree,                      // size: 2*L (index 0 unused, root at 1)
13     const int64_t* idx,                // leaf indices in [0, capacity)
14     const float* new_p,               // new priorities
15     int64_t n,
16     int64_t L
17 ) {
18     int i = blockIdx.x * blockDim.x + threadIdx.x;
19     if (i >= n) return;
20
21     int64_t pos = L + idx[i];
22
23     // atomic exchange at leaf to handle duplicates safely
24     float old = atomicExch(&tree[pos], new_p[i]);
25     float delta = new_p[i] - old;
26
27     // propagate delta to parents
28     pos >>= 1;
29     while (pos >= 1) {
30         atomicAdd(&tree[pos], delta);
31         pos >>= 1;
32     }
33 }
34
35
36 void update_cuda(torch::Tensor tree, torch::Tensor idx, torch::Tensor
37     new_p, int64_t L) {
38     check_cuda(tree, "tree");
39     check_cuda(idx, "idx");
40     check_cuda(new_p, "new_p");
41
42     const auto n = idx.numel();
43     const int threads = 256;
44     const int blocks = (n + threads - 1) / threads;
45
46     update_kernel<<<blocks, threads, 0, at::cuda::getDefaultCUDASStream()
47         >>>(
48             tree.data_ptr<float>(),
49             idx.data_ptr<int64_t>(),
50             new_p.data_ptr<float>(),
51             n,
52             L
53         );
54 }
55
56 // helper: warp-level shuffle reduction on delta values (optional)
57 static __device__ __forceinline__ float warp_reduce_sum(float v) {
58     for (int offset = 16; offset > 0; offset >>= 1) {
59         v += __shfl_down_sync(0xffffffff, v, offset);
60     }
61     return v;
62 }
63
64 // GPU UPDATE (warp-coalesced atomics)

```

```

63  __global__ void update_kernel_coalesced(
64      float* tree, const int64_t* idx, const float* new_p, int64_t n,
65      int64_t L
66  ) {
67      int i = blockIdx.x * blockDim.x + threadIdx.x;
68      if (i >= n) return;
69
70      int64_t pos = L + idx[i];
71      float old = atomicExch(&tree[pos], new_p[i]);
72      float delta = new_p[i] - old;
73
74      // propagate delta to parents (warp-coalesced atomics)
75      pos >>= 1;
76      while (pos >= 1) {
77          // group threads by parent index to reduce atomics
78          int64_t parent = pos;
79          // find lanes with the same parent
80          unsigned mask = __match_any_sync(0xffffffff, parent);
81          float v = delta;
82          // full warp reduction, then only leader does atomicAdd
83          float all = warp_reduce_sum(v);
84          int leader = __ffs(mask) - 1;
85          if ((threadIdx.x & 31) == leader) {
86              atomicAdd(&tree[parent], all);
87          }
88          pos >>= 1;
89     }
90 }
91
92 void update_coalesced_cuda(torch::Tensor tree, torch::Tensor idx, torch::
93     Tensor new_p, int64_t L) {
94     check_cuda(tree, "tree");
95     check_cuda(idx, "idx");
96     check_cuda(new_p, "new_p");
97
98     const auto n = idx.numel();
99     const int threads = 256;
100    const int blocks = (n + threads - 1) / threads;
101
102    update_kernel_coalesced<<<blocks, threads, 0, at::cuda::
103        getDefaultCUDASStream()>>>(
104            tree.data_ptr<float>(), idx.data_ptr<int64_t>(),
105            new_p.data_ptr<float>(), n, L);
106 }
107
108 // GPU SAMPLE
109 //__global__ void sample_kernel(
110 //    const float* tree, const float* u, int64_t* out, int64_t B, int64_t L,
111 //    int64_t capacity
112 //) {
113 //    int i = blockIdx.x * blockDim.x + threadIdx.x;
114 //    if (i >= B) return;
115 //
116 //    float x = u[i];
117 //    int64_t node = 1;
118 //    while (node < L) {
119 //        int64_t left = node << 1;
120 //        float s_left = tree[left];
121 //        if (x <= s_left) {
122 //            node = left;
123 //        }
124 //    }
125 //}

```

```

118     } else {
119         x -= s_left;
120         node = left + 1;
121     }
122 }
123 int64_t leaf = node - L;
124 if (leaf >= capacity) leaf = capacity - 1;
125 out[i] = leaf;
126 }
127
128 torch::Tensor sample_cuda(torch::Tensor tree, torch::Tensor u, int64_t L,
129   int64_t capacity) {
130     check_cuda(tree, "tree");
131     check_cuda(u, "u");
132
133     const auto B = u.numel();
134     auto out = torch::empty({B}, torch::TensorOptions().device(u.device())
135       .dtype(torch::kInt64));
136
137     const int threads = 256;
138     const int blocks = (B + threads - 1) / threads;
139     sample_kernel<<<blocks, threads, 0, at::cuda::getDefautCUDASream()
140     >>>(
141       tree.data_ptr<float>(), u.data_ptr<float>(), out.data_ptr<int64_t
142       >(), B, L, capacity);
143     return out;
144 }
```

### 13.3 Expected results from kernels

- After any sequence of updates, `tree[1]` equals the sum of leaf values (within float tolerance).
- Sampling approximates the intended discrete distribution proportional to leaf priorities.
- Coalesced updates should have similar correctness; performance depends on contention patterns.

## 14 Minimal GPU PER Replay Buffer: `replay.py`

### 14.1 Intent

`GPUReplayPER` stores transitions in GPU tensors and uses `GpuSumTree` for proportional sampling and importance weights.

### 14.2 Detailed behavior (function-by-function)

**PERBatch dataclass.** Holds sampled batch tensors: `obs`, `act`, `rew`, `next_obs`, `done`, `weights`, `indices`.

**`__init__`.** Allocates GPU storage for transitions and initializes:

- ring buffer pointers `pos` and `size` on CPU,
- a `GpuSumTree` of size `capacity`,
- `_max_priority` for initializing new transitions.

**`add`.** Writes one transition at `pos`, sets its tree priority to `_max_priority`, then advances the ring.

sample.

- Computes  $\beta$  by linear anneal from `beta0` to 1.0 over `beta_steps`.
- Samples indices proportional to priorities via `tree.sample`.
- Gathers tensors with advanced indexing `storage[idx]`.
- Computes normalized IS weights on GPU:

$$w = (N \cdot p)^{-\beta}, \quad w \leftarrow \frac{w}{\max w}.$$

`update_priorities`. Computes new leaf priorities:

$$p \leftarrow (|td| + \epsilon)^\alpha$$

and updates the sum-tree. Tracks `_max_priority` for future insertions.

### 14.3 Full source

Listing 8: `replay.py`

```
1 import torch
2 from dataclasses import dataclass
3 from .tree import GpuSumTree
4
5 @dataclass
6 class PERBatch:
7     obs: torch.Tensor
8     act: torch.Tensor
9     rew: torch.Tensor
10    next_obs: torch.Tensor
11    done: torch.Tensor
12    weights: torch.Tensor
13    indices: torch.Tensor
14
15 class GPUReplayPER:
16     """
17         Minimal GPU PER replay:
18             - stores transitions in CUDA tensors
19             - priorities in a GPU sum-tree
20             - proportional sampling + importance weights
21     """
22     def __init__(
23         self,
24         capacity: int,
25         obs_shape,
26         device: str = "cuda",
27         alpha: float = 0.6,
28         beta0: float = 0.4,
29         eps: float = 1e-6,
30         dtype_obs=torch.float32,
31     ):
32         self.device = torch.device(device)
33         self.capacity = int(capacity)
34         self.alpha = float(alpha)
35         self.beta0 = float(beta0)
36         self.eps = float(eps)
37
38         # ring buffer state (keep on CPU for simplicity)
```

```

39     self.pos = 0
40     self.size = 0
41
42     # storage on GPU
43     self.obs = torch.empty((capacity, *obs_shape), device=self.device,
44                           dtype=dtype_obs)
44     self.next_obs = torch.empty((capacity, *obs_shape), device=self.
45                                 device, dtype=dtype_obs)
45     self.act = torch.empty((capacity,), device=self.device, dtype=
46                           torch.int64)
46     self.rew = torch.empty((capacity,), device=self.device, dtype=
47                           torch.float32)
47     self.done = torch.empty((capacity,), device=self.device, dtype=
48                           torch.float32)
48
49     # priorities
50     self.tree = GpuSumTree(capacity, device=device)
51     self._max_priority = 1.0 # CPU scalar is fine
52
53     def __len__(self):
54         return self.size
55
56     @torch.no_grad()
57     def add(self, obs, act, rew, next_obs, done):
58         """
59             obs/next_obs can be numpy or torch; will be copied to GPU.
60             done should be bool or 0/1.
61         """
62         i = self.pos
63
64         # copy to GPU tensors
65         self.obs[i].copy_(torch.as_tensor(obs, device=self.device))
66         self.next_obs[i].copy_(torch.as_tensor(next_obs, device=self.
67                               device))
67         self.act[i] = int(act)
68         self.rew[i] = float(rew)
69         self.done[i] = float(done)
70
71         # set priority for this slot to current max (common PER trick)
72         p = torch.tensor([self._max_priority], device=self.device, dtype=
73                           torch.float32)
73         idx = torch.tensor([i], device=self.device, dtype=torch.int64)
74         self.tree.update(idx, p)
75
76         # advance ring
77         self.pos = (self.pos + 1) % self.capacity
78         self.size = min(self.size + 1, self.capacity)
79
80     @torch.no_grad()
81     def sample(self, batch_size: int, step: int, beta_steps: int = 200_000
82               ) -> PERBatch:
83         """
84             beta anneals from beta0 to 1.0 over beta_steps.
85         """
86         if self.size == 0:
87             raise RuntimeError("Cannot sample from empty buffer")
88
89         # beta schedule
90         frac = min(1.0, step / float(beta_steps))
90         beta = self.beta0 + frac * (1.0 - self.beta0)

```

```

91     # sample indices proportional to priorities
92     idx = self.tree.sample(batch_size) # int64 CUDA [B]
93
94     # gather transitions
95     obs = self.obs[idx]
96     next_obs = self.next_obs[idx]
97     act = self.act[idx]
98     rew = self.rew[idx]
99     done = self.done[idx]
100
101
102     # importance weights
103     total = self.tree.total_priority() # CUDA scalar
104     prio = self.tree.leaf_priorities(idx).clamp_min(self.eps)
105     p = prio / total.clamp_min(self.eps)
106
107     N = float(self.size)
108     weights = (N * p).pow(-beta)
109     weights = weights / weights.max().clamp_min(self.eps)
110
111     return PERBatch(obs=obs, act=act, rew=rew, next_obs=next_obs, done
112                     =done, weights=weights, indices=idx)
113
114     @torch.no_grad()
115     def update_priorities(self, indices: torch.Tensor, td_error: torch.
116                           Tensor):
117         """
118             indices: int64 CUDA [B]
119             td_error: float tensor [B] on CUDA
120             priority = (|td| + eps)^alpha
121         """
122
123         p = (td_error.abs() + self.eps).pow(self.alpha).to(torch.float32)
124         self.tree.update(indices, p)
125
126         # track max priority on CPU to init new entries
127         max_p = float(p.max().item())
128         if max_p > self._max_priority:
129             self._max_priority = max_p

```

## 15 Benchmarks: `bench_per.py`

### 15.1 Intent

Measures update and sample latency for:

- `GpuSumTree` (extension),
- baseline `GpuCumsumSampler`,
- CPU reference `CpuSumTree`.

### 15.2 Expected benchmark trends

- With frequent updates, sum-tree should outperform `cumsum` baseline by avoiding full prefix-sum recomputation.
- Sampling performance depends on capacity and batch size, typically favoring GPU for large workloads.

### 15.3 Full source

Listing 9: bench\_per.py

```
1 import time
2 import torch
3 import numpy as np
4
5 from per.tree import GpuSumTree
6 from per.baseline_cumsum import GpuCumsumSampler
7 from per.cpu_tree import CpuSumTree
8
9
10 def time_cuda(fn, iters=100):
11     # warmup
12     for _ in range(5):
13         fn()
14     torch.cuda.synchronize()
15
16     t0 = time.time()
17     for _ in range(iters):
18         fn()
19     torch.cuda.synchronize()
20     return (time.time() - t0) / iters
21
22
23 def time_cpu(fn, iters=200):
24     for _ in range(10):
25         fn()
26     t0 = time.time()
27     for _ in range(iters):
28         fn()
29     return (time.time() - t0) / iters
30
31
32 def benchmark(capacity, batch_size, iters=200):
33     print(f"\n==== capacity={capacity} batch={batch_size} ===")
34     device = "cuda" if torch.cuda.is_available() else "cpu"
35
36     # Data
37     idx_gpu = torch.randint(0, capacity, (batch_size,), device=device,
38                           dtype=torch.int64)
39     p_gpu = torch.rand(batch_size, device=device, dtype=torch.float32)
40
41     # GPU sum-tree
42     if device == "cuda":
43         t = GpuSumTree(capacity)
44
45         def upd_tree():
46             t.update(idx_gpu, p_gpu)
47
48         def samp_tree():
49             _ = t.sample(batch_size)
50
51         t_upd = time_cuda(upd_tree, iters=iters)
52         # Need non-zero total before sample timing
53         t.update(torch.arange(capacity, device="cuda", dtype=torch.int64),
54                 torch.ones(capacity, device="cuda", dtype=torch.float32))
55         t_samp = time_cuda(samp_tree, iters=iters)
```

```

56     print(f"GpuSumTree update: {t_upd*1e6:.2f} us")
57     print(f"GpuSumTree sample: {t_samp*1e6:.2f} us")
58
59     # Baseline cumsum sampler
60     b = GpuCumsumSampler(capacity)
61     b.update(idx_gpu, p_gpu)
62
63     def upd_base():
64         b.update(idx_gpu, p_gpu)
65
66     def samp_base():
67         _ = b.sample(batch_size)
68
69     t_upd2 = time_cuda(upd_base, iters=iters)
70     # make priorities non-zero
71     b.priorities[:] = 1.0
72     t_samp2 = time_cuda(samp_base, iters=iters)
73
74     print(f"Baseline update: {t_upd2*1e6:.2f} us")
75     print(f"Baseline sample: {t_samp2*1e6:.2f} us")
76
77     # CPU sum-tree
78     c = CpuSumTree(capacity)
79     idx_cpu = np.random.randint(0, capacity, size=(batch_size,), dtype=np.
80         int64)
80     p_cpu = np.random.rand(batch_size).astype(np.float32)
81
82     def upd_cpu():
83         c.update(idx_cpu, p_cpu)
84
85     def samp_cpu():
86         total = c.total
87         if total <= 0:
88             return
89         u = (np.arange(batch_size, dtype=np.float32) + np.random.rand(
90             batch_size).astype(np.float32)) * (total / batch_size)
90         _ = c.sample(u)
91
92     t_upd_cpu = time_cpu(upd_cpu, iters=iters)
93     # make non-zero
94     c.update(np.arange(capacity, dtype=np.int64), np.ones(capacity, dtype=
95         np.float32))
95     t_samp_cpu = time_cpu(samp_cpu, iters=iters)
96
97     print(f"CpuSumTree update: {t_upd_cpu*1e6:.2f} us")
98     print(f"CpuSumTree sample: {t_samp_cpu*1e6:.2f} us")
99
100
101 def main():
102     caps = [2**12, 2**14, 2**16, 100_000]
103     batches = [256, 1024, 4096]
104     for cap in caps:
105         for b in batches:
106             benchmark(cap, b, iters=200)
107
108
109 if __name__ == "__main__":
110     main()

```

## 16 Example Training Script: `train_cartpole.py`

### 16.1 Intent

Demonstrates how `GPUReplayPER` can be used in a DQN-like loop for `CartPole-v1`, including:

- epsilon-greedy action selection,
- replay insertions,
- PER sampling with IS weights,
- TD-error-based priority updates,
- periodic target network updates.

### 16.2 Expected outcome

On a correct setup (and with reasonable hyperparameters), training should learn a policy that achieves high CartPole return. Exact learning curves vary by random seed and system.

### 16.3 Full source

Listing 10: `train_cartpole.py`

```
1 import time
2 import random
3 from dataclasses import dataclass
4
5 import numpy as np
6 import torch
7 import torch.nn as nn
8 import torch.nn.functional as F
9
10 import gymnasium as gym
11
12 from per.replay import GPUReplayPER
13
14
15 @dataclass
16 class Config:
17     env_id: str = "CartPole-v1"
18     seed: int = 0
19
20     total_steps: int = 300_000
21     warmup_steps: int = 10_000
22
23     capacity: int = 100_000
24     batch_size: int = 256
25     gamma: float = 0.99
26
27     lr: float = 1e-3
28     train_freq: int = 1
29     target_update_freq: int = 1000
30
31     # epsilon-greedy
32     eps_start: float = 1.0
33     eps_end: float = 0.05
34     eps_decay_steps: int = 100_000
35
36     # PER
```

```

37     alpha: float = 0.6
38     beta0: float = 0.4
39     beta_anneal_steps: int = 200_000
40
41
42 class QNet(nn.Module):
43     def __init__(self, obs_dim, act_dim):
44         super().__init__()
45         self.net = nn.Sequential(
46             nn.Linear(obs_dim, 128),
47             nn.ReLU(),
48             nn.Linear(128, 128),
49             nn.ReLU(),
50             nn.Linear(128, act_dim),
51         )
52
53     def forward(self, x):
54         return self.net(x)
55
56
57 def linear_schedule(step, start, end, duration):
58     t = min(step / duration, 1.0)
59     return start + t * (end - start)
60
61
62 def main(cfg: Config):
63     random.seed(cfg.seed)
64     np.random.seed(cfg.seed)
65     torch.manual_seed(cfg.seed)
66
67     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
68
69     env = gym.make(cfg.env_id)
70     obs, _ = env.reset(seed=cfg.seed)
71     obs_dim = obs.shape[0]
72     act_dim = env.action_space.n
73
74     q = QNet(obs_dim, act_dim).to(device)
75     q_targ = QNet(obs_dim, act_dim).to(device)
76     q_targ.load_state_dict(q.state_dict())
77
78     opt = torch.optim.Adam(q.parameters(), lr=cfg.lr)
79
80     buf = GPUReplayPER(
81         capacity=cfg.capacity,
82         obs_shape=(obs_dim,),
83         alpha=cfg.alpha,
84         beta0=cfg.beta0,
85         beta_anneal_steps=cfg.beta_anneal_steps,
86         device=device,
87     )
88
89     ep_ret = 0.0
90     ep_len = 0
91     t0 = time.time()
92
93     for step in range(cfg.total_steps):
94         eps = linear_schedule(step, cfg.eps_start, cfg.eps_end, cfg.
95             eps_decay_steps)
96         if random.random() < eps:

```

```

96         act = env.action_space.sample()
97     else:
98         with torch.no_grad():
99             o = torch.tensor(obs, dtype=torch.float32, device=device).
100                unsqueeze(0)
101            act = torch.argmax(q(o), dim=1).item()
102
103    next_obs, rew, terminated, truncated, _ = env.step(act)
104    done = terminated or truncated
105
106    buf.add(
107        torch.tensor(obs, dtype=torch.float32),
108        act,
109        float(rew),
110        torch.tensor(next_obs, dtype=torch.float32),
111        done,
112    )
113
114    obs = next_obs
115    ep_ret += rew
116    ep_len += 1
117
118    if done:
119        obs, _ = env.reset()
120        ep_ret = 0.0
121        ep_len = 0
122
123    if step < cfg.warmup_steps:
124        continue
125
126    if step % cfg.train_freq == 0:
127        batch = buf.sample(batch_size=cfg.batch_size, step=step)
128
129        with torch.no_grad():
130            q_next = q_targ(batch.next_obs).max(dim=1).values
131            target = batch.rew + cfg.gamma * (1.0 - batch.done.float())
132                * q_next
133
134            q_pred = q(batch.obs).gather(1, batch.act.long().unsqueeze(1))
135                .squeeze(1)
136            td = target - q_pred
137
138            loss = (batch.weights * td.pow(2)).mean()
139
140            opt.zero_grad()
141            loss.backward()
142            opt.step()
143
144            buf.update_priorities(batch.indices, td.detach())
145
146    env.close()
147
148
149 if __name__ == "__main__":
150     main(Config())

```

## 17 Tests: Correctness and Expected Results

### 17.1 How to run

Typical command:

- `pytest -q`

These tests are designed primarily for CUDA; if CUDA is unavailable they return early.

### 17.2 `test_update.py`: Tree invariants

**Expected results.** All assertions pass:

- Root equals sum of leaves after random updates (within tolerance).
- Random internal nodes satisfy `tree[i] = tree[2i] + tree[2i+1]`.

Listing 11: `test_update.py`

```
1 import torch
2 from per.tree import GpuSumTree
3
4 def test_sum_invariant_after_random_updates():
5     if not torch.cuda.is_available():
6         return
7
8     cap = 1024
9     t = GpuSumTree(cap)
10
11    # random updates (with duplicates)
12    n = 5000
13    idx = torch.randint(0, cap, (n,), device="cuda", dtype=torch.int64)
14    new_p = torch.rand(n, device="cuda", dtype=torch.float32)
15
16    t.update(idx, new_p)
17
18    leaves_sum = t.leaf_values().sum()
19    root = t.total
20
21    assert torch.allclose(root, leaves_sum, rtol=1e-4, atol=1e-5), (root.
22        item(), leaves_sum.item())
23
24 def test_internal_node_invariant_spotcheck():
25     if not torch.cuda.is_available():
26         return
27
28     cap = 256
29     t = GpuSumTree(cap)
30
31    # set all leaves explicitly (unique update)
32    idx = torch.arange(cap, device="cuda", dtype=torch.int64)
33    new_p = torch.rand(cap, device="cuda", dtype=torch.float32)
34    t.update(idx, new_p)
35
36    # spot check random internal nodes
37    # valid internal indices: [1, L-1]
38    L = t.L
39    internal = torch.randint(1, L, (50,), device="cuda", dtype=torch.int64
40        ).tolist()
```

```

40     for i in internal:
41         left = 2 * i
42         right = left + 1
43         # right child might be out of bounds only if i >= L, but we
44         # sampled < L
45         s = t.tree[left] + t.tree[right]
46         assert torch.allclose(t.tree[i], s, rtol=1e-4, atol=1e-5), f"node
47         {i} mismatch"

```

### 17.3 test\_sample.py: Sampling distribution

**Expected results.** The L1 error between the empirical histogram and the target distribution stays below a forgiving threshold (e.g.  $< 0.10$ ). This is a statistical test, so the tolerance is intentionally loose to avoid flakiness.

Listing 12: test\_sample.py

```

1 import torch
2 from per.tree import GpuSumTree
3
4 def test_sampling_distribution_reasonable():
5     if not torch.cuda.is_available():
6         return
7
8     cap = 128
9     t = GpuSumTree(cap)
10
11    idx = torch.arange(cap, device="cuda", dtype=torch.int64)
12    p = (idx.float() + 1.0) # priorities 1..cap
13    t.update(idx, p)
14
15    # sample many times
16    B = 200_000
17    samp = t.sample(B)
18
19    hist = torch.bincount(samp, minlength=cap).float()
20    hist /= hist.sum()
21
22    target = p / p.sum()
23    l1 = torch.abs(hist - target).sum().item()
24
25    # forgiving threshold to avoid flakiness
26    assert l1 < 0.10, f'L1 error too large: {l1}'

```

### 17.4 test\_per\_replay.py: Replay sampling + updates

**Expected results.**

- Sampled batch shapes match the requested batch size and observation dimensions.
- Indices lie within  $[0, \text{capacity})$ .
- IS weights are positive and normalized to  $\leq 1$ .
- Priority updates run without error.

Listing 13: test\_per\_replay.py

```

1 import torch
2 from per.replay import GPUReplayPER

```

```

3
4 def test_per_replay_sample_shapes():
5     if not torch.cuda.is_available():
6         return
7
8     buf = GPUReplayPER(capacity=1024, obs_shape=(4,))
9     # fill some transitions
10    for i in range(200):
11        obs = torch.randn(4)
12        next_obs = torch.randn(4)
13        buf.add(obs, i % 2, 1.0, next_obs, False)
14
15    batch = buf.sample(batch_size=64, step=0)
16
17    assert batch.obs.shape == (64, 4)
18    assert batch.next_obs.shape == (64, 4)
19    assert batch.act.shape == (64,)
20    assert batch.rew.shape == (64,)
21    assert batch.done.shape == (64,)
22    assert batch.weights.shape == (64,)
23    assert batch.indices.shape == (64,)
24    assert batch.indices.min().item() >= 0
25    assert batch.indices.max().item() < 1024
26    assert batch.weights.min().item() > 0.0
27    assert batch.weights.max().item() <= 1.0 + 1e-6
28
29 def test_update_priorities_runs():
30     if not torch.cuda.is_available():
31         return
32
33     buf = GPUReplayPER(capacity=256, obs_shape=(4,))
34     for i in range(100):
35         buf.add(torch.zeros(4), 0, 0.0, torch.zeros(4), False)
36
37     batch = buf.sample(batch_size=32, step=10)
38     td = torch.randn(32, device="cuda")
39     buf.update_priorities(batch.indices, td) # should not crash

```

## 18 Note on `per_ext.py` Filename

The repository includes a file named `per_ext.py` whose contents are C++ code (not Python). This is likely an artifact or alternate binding stub and should be renamed (e.g., `per_ext_alt.cpp`) to avoid confusion. It is *not* used by `_ext.py`.

Listing 14: `per_ext.py` (contains C++ code; consider renaming)

```

1 #include <torch/extension.h>
2
3 void update_cuda(torch::Tensor tree, torch::Tensor idx, torch::Tensor
4     new_p, int64_t L);
5 void update_opt_cuda(torch::Tensor tree, torch::Tensor idx, torch::Tensor
6     new_p, int64_t L);
7
8 PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
9     m.def("update", &update_cuda, "Sum-tree update (CUDA)");
10    m.def("update_opt", &update_opt_cuda, "Sum-tree update optimized with CUB
11          (CUDA)");
12    // keep sample binding too if you have it
13}

```

---

## 19 Conclusion

This coursework implements a GPU-friendly PER sum-tree with:

- correctness invariants validated by tests,
- GPU kernels designed for safe concurrent updates (via atomics),
- a replay buffer that computes IS weights and supports TD-error priority updates,
- a baseline method for performance comparison.

The design aligns with the original PER paper and demonstrates practical GPU programming considerations (atomics, warp-level primitives, memory layout, and integration via PyTorch extensions).

## References

- [1] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. *Prioritized Experience Replay*. arXiv:1511.05952, 2015. DOI: 10.48550/arXiv.1511.05952.