

# DynaSOAr: A Parallel Memory Allocator for Object-oriented Programming on GPUs



東京工業大学  
Tokyo Institute of Technology

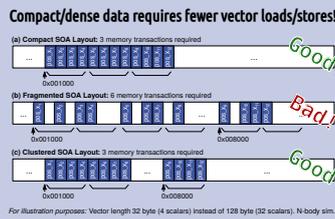
Matthias Springer and Hidehiko Masuhara (Tokyo Institute of Technology)

See our presentation @ ECOOP on Thursday, 15:40 - 16:00!

## Motivation

- **Object-oriented programming (OOP)** has many applications in high-performance computing (HPC), but it is **too slow**.
- **Esp. dynamic memory allocation**, which is very useful in OOP.
- **Main performance problem on GPUs: Memory access.**
- **DynaSOAr: A framework for efficient OOP on GPUs (3 parts)**
- 1. **Data layout DSL: Store objects in a Structure of Arrays (SOA).**
- 2. **Dyn. mem. allocator: Allocate objects with low fragmentation.**
- 3. **Parallel do-all operation: Expressing parallelism over an obj. set.**

## Memory Coalescing



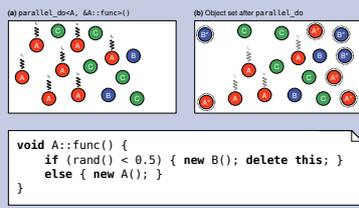
## Related Work

- **Default CUDA allocator: Slow, unoptimized, unreliable**
- **Other state-of-the-art GPU allocators:**
  - **Halloc [1]** and **ScatterAlloc/mallocMC [2]**
  - **Lack optimizations for structured data: Fast (de)allocation but inefficient usage of allocated memory.**
  - **Based on hashing techniques: Leads to high fragmentation.**
- **DynaSOAr has detailed knowledge about the structure of its allocations (fields of classes/structs).** This allows us to apply more optimizations (e.g., SOA layout). Other allocators just allocate raw bytes.
- **Data layout DSL based on Ikra-Cpp [3].**

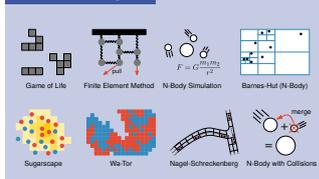
- [1] A. Adineez, D. Pleter. Halloc: A High-Throughput Dynamic Memory Allocator for GPU Architectures. GPU Technology Conference, 2004.
- [2] M. Steinberger, M. Kenzel, B. Kalnz, D. Schmalstieg. ScatterAlloc: Massively Parallel Dynamic Memory Allocation for the GPU. In Par 2012.
- [3] M. Springer, H. Masuhara. Ikra-Cpp: A C++/CUDA DSL for Object-Oriented Programming with Structure-of-Arrays Layout. WPMVP 2018.

## SMMO: Single-Method Multiple-Objects

- **Obj.-orient. programming model** for SIMD architectures.
- **OOP-speech for Single-Instruction Multiple-Data (SIMD)**
- **Run method for all objects of a type T:** parallel\_do<T, &T::func>
- **Method can create/delete objects.**
- **In parallel: Spawns a CUDA kernel.**
- **Deleting an object of type T other than this (bound object) is forbidden.**



## SMMO Examples



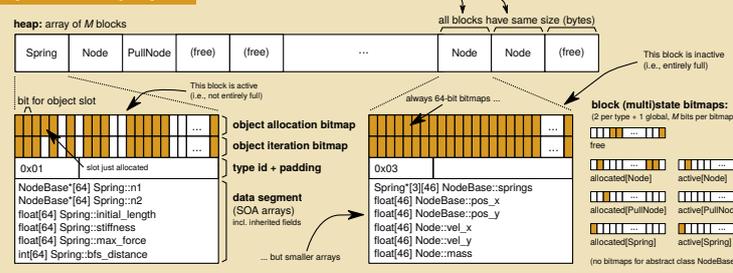
```

#include "dynasoar.h"
class Body {
public:
    Body(float pos_x, float pos_y, float pos_z, float vel_x, float vel_y, float vel_z, float mass) : pos_x(pos_x), pos_y(pos_y), pos_z(pos_z), vel_x(vel_x), vel_y(vel_y), vel_z(vel_z), mass(mass) {}
};

class BodyAllocator : public Allocator<Body> {
public:
    BodyAllocator(int num_bodies, float max_force) : num_bodies(num_bodies), max_force(max_force) {}
};

int main() {
    auto h_allocator = new BodyAllocator(10000, 10000); // Initialize DynaSOAr heap.
    auto h_allocator = new BodyAllocator(10000, 10000); // Initialize DynaSOAr heap.
    for (int i = 0; i < 10000; ++i) {
        h_allocator->allocate<Body>(i);
        delete h_allocator; // Release all GPU memory.
    }
}
    
```

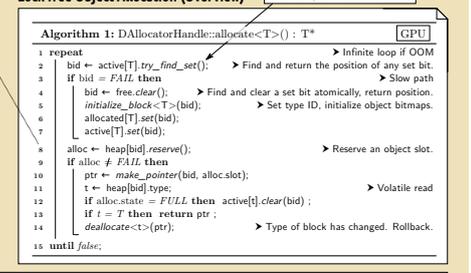
## DynaSOAr Heap Layout



```

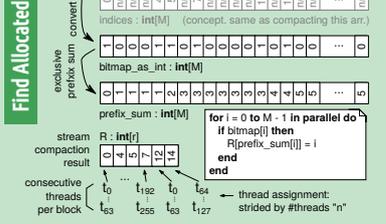
Algorithm 4: Block reservation | (int, state) | reserving capacity 64 | GPU
1 repeat
2   if pos < 2048 then return FAIL;
3   if pos < 1024 then return FAIL;
4   if pos < 512 then return FAIL;
5   if pos < 256 then return FAIL;
6   if pos < 128 then return FAIL;
7   if pos < 64 then return FAIL;
8   if pos < 32 then return FAIL;
9   if pos < 16 then return FAIL;
10  if pos < 8 then return FAIL;
11  if pos < 4 then return FAIL;
12  if pos < 2 then return FAIL;
13  if pos < 1 then return FAIL;
14  if pos < 0 then return FAIL;
15  return FAIL;
    
```

## Lock-free Object Allocation (Overview)

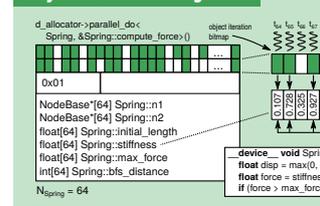


## Main Challenges

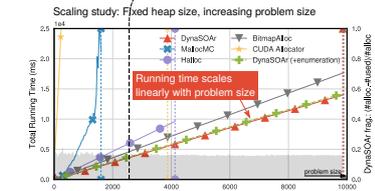
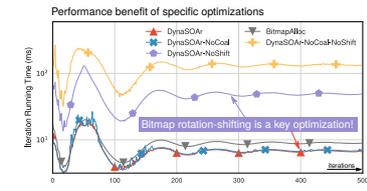
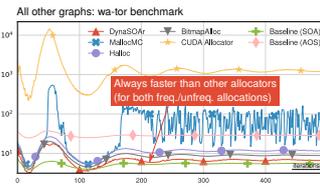
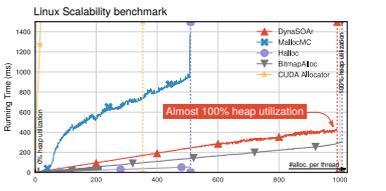
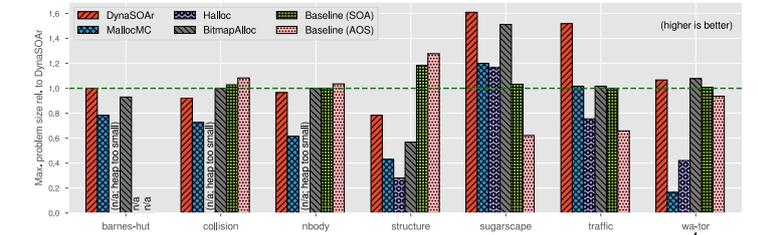
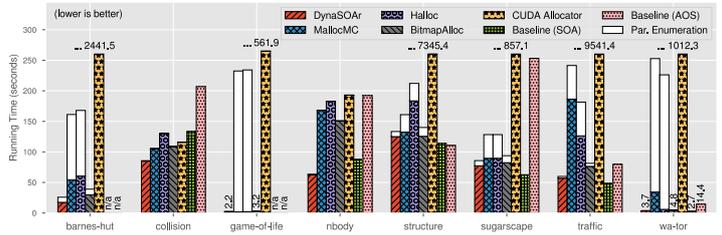
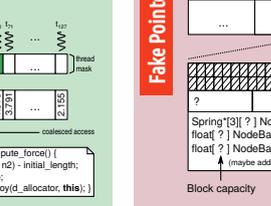
- **Lock-free implementation: Based on atomic operations and retry loops.**
- **Safe memory reclamations: All blocks have the same structure regardless of state, fill level and type. I.e., bitmaps and type ID are always at the same offset.**



## Object-to-Thread Assignment



## Fake Pointers



## Other Optimizations

- To find active/free blocks with a logarithmic number of accesses, block states are indexed with **hierarchical, lock-free bitmaps**.
- To reduce allocation contention, **simultaneous allocation requests are coalesced/combined on a per-warp basis [4]**.
- To reduce #retries, **bitmaps are rotating-shifted before ffs**.
- To further reduce fragmentation, try Alg. 1 Line 2 up to r times.

## Memory Fragmentation

- **Diminishes the benefit of SOA.**
  - **No internal/external fragmentation by design.**
- $$F = \frac{1}{\#Blocks} \sum_{b \in Blocks} \frac{\#free\ slots(b)}{\#slots(b)}$$
- (only allocated blocks)

