# CompactGpu: Massively Parallel Memory Defragmentation on GPUs

Matthias Springer (Tokyo Institute of Technology)     https://github.com/prg-titech/dynasoar

東京工業大学
Tokyo Institute of Technology

## Why Defragment GPU Memory?

- Space efficiency: **Reduce memory usage**
- Improve runtime performance:
  Accessing compact data requires fewer vector accesses
  → **Better memory coalescing**

## Design Requirements

- Extension to the DynaSOAr dynamic GPU memory allocator
- **Parallel, in-place, stop-the-world** defragmentation approach
- To reduce defragmentation overhead: Uniform control flow, little synchronization, efficient memory access
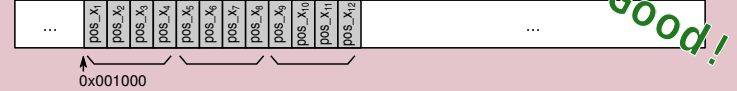
## Related Work

- R. Veldema, M. Philippsen. Parallel Memory Defragmentation on GPUs. MSPC '12
  Assumes many **different allocation sizes**, not in-place, large runtime overhead
- M. Springer, H. Masuhara. DynaSOAr: A Parallel Memory Allocator for Object-oriented Programming on GPUs with Efficient Memory Access. ECOOP '19
- H. Boehm. Space Efficient Conservative Garbage Collection. PLDI '93
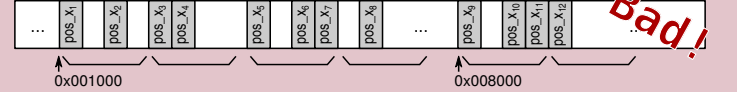  Similar problem: How to **find all pointers to moved objects** that must be rewritten?

## Background: GPU Architecture and Dyn. Memory Allocation

- Pattern: Many small allocations, mostly same size
- For good mem. access performance: **Structure of Arrays (SOA)** data layout
- Recent NVIDIA GPUs have 128-byte vector registers
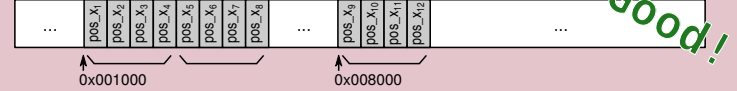  → Memory access in aligned, 128-byte transactions



**(a) Compact SOA Layout:** 3 memory transactions required   *Good!*

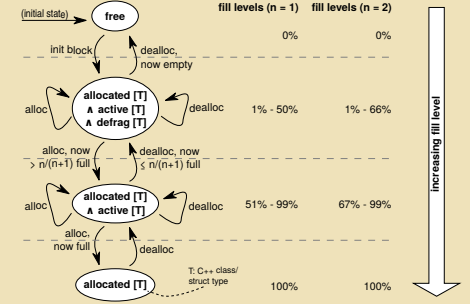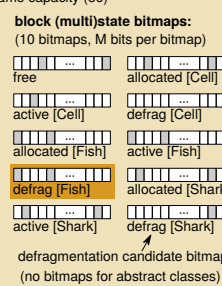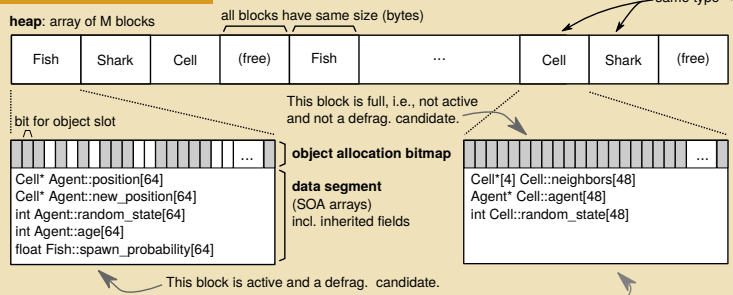**(b) Fragmented SOA Layout:** 6 memory transactions required   *Bad!*

**(c) Clustered SOA Layout:** 3 memory transactions required   *Good!*

*For illustration purposes:* Vector length 32 byte (4 scalars) instead of 128 byte (32 scalars). N-body sim.
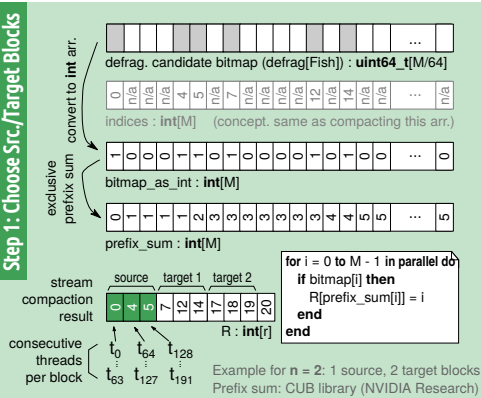
### DynaSOAr Heap Layout



**heap**: array of M blocks

all blocks have same size (bytes)

same type ⇒ same capacity (56)

This block is full, i.e., not active and not a defrag. candidate.

bit for object slot

**object allocation bitmap**

Cell* Agent::position[64]
Cell* Agent::new_position[64]
int Agent::random_state[64]
int Agent::age[64]
float Fish::spawn_probability[64]

**data segment** (SOA arrays) incl. inherited fields

Cell*[4] Cell::neighbors[48]
Agent* Cell::agent[48]
int Cell::random_state[48]

This block is active and a defrag. candidate.

Running example: Fish-and-Sharks simulation

**block (multi)state bitmaps:** (10 bitmaps, M bits per bitmap)

free / allocated [Cell]
active [Cell] / defrag [Cell]
allocated [Fish] / active [Fish]
defrag [Fish] / allocated [Shark]
active [Shark] / defrag [Shark]

defragmentation candidate bitmap (no bitmaps for abstract classes)

fill levels (n = 1)   fill levels (n = 2)

T: C++ class / struct type

## Defragmentation by Block Merging:
### parallel_defrag<Fish>()

$$F = \frac{1}{\#Blocks} \sum_{b \in Blocks} \frac{\#free\ slots(b)}{\#slots(b)}$$

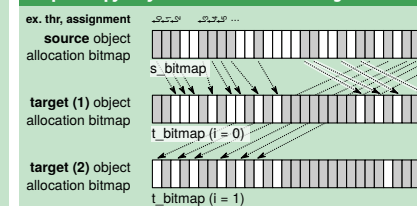**Definition of Defragmentation Candidates:**

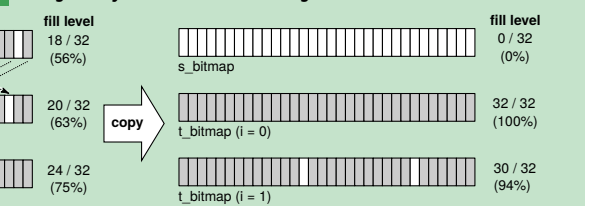**Depends on defrag. factor n (problem-spec., compile-time parameter)**

n = 1: ≤ 50% full
n = 2: ≤ 66.6% full
Arbitrary n: ≤ n/(n+1) full
Guaranteed target frag.: 1/(n+1)

**Step 1: Choose Src./Target Blocks**

defrag. candidate bitmap (defrag[Fish]) : **uint64_t**[M/64]

convert to **int** arr.

indices : **int**[M] (concept. same as compacting this arr.)

exclusive prefix sum

bitmap_as_int : **int**[M]

prefix_sum : **int**[M]

stream compaction result

consecutive threads per block

R : **int**[r]

```
for i = 0 to M - 1 in parallel do
    if bitmap[i] then
        R[prefix_sum[i]] = i
    end
end
```

Example for **n = 2**: 1 source, 2 target blocks
Prefix sum: CUB library (NVIDIA Research)

**Step 2: Copy Objects from Source to Target Blocks:** Merge every source block into n target blocks.

ex. thr. assignment

**source** object allocation bitmap   s_bitmap   18 / 32 (56%)   →   0 / 32 (0%)

**target (1)** object allocation bitmap   t_bitmap (i = 0)   20 / 32 (63%)   **copy**   32 / 32 (100%)

**target (2)** object allocation bitmap   t_bitmap (i = 1)   24 / 32 (75%)   →   30 / 32 (94%)

**Step 3: Store Forwarding Ptrs. in Source**

Fish* Fish::forwarding_ptr[64]

Overwrite data segment with pointers.

**Step 4: Rewrite Pointers to Relocated Objects with Bitmap**

```
for all Fish*& ptr in parallel do
    s_bid = extract_block_id(ptr)
    if defrag[Fish][s_bid] then
        s_oid = extract_object_id(ptr)
        ptr = heap[s_bid].forwarding_ptr[s_oid]
    end
end
```

How to find all Fish*/Agent* values on the heap?
- Option 1: Scan heap, look for anything that could be a pointer. *Slow!*
- Option 2: Utilize DynaSOAr's data layout DSL. **Scan only mem. locations of SOA arrays with base type Fish*/Agent*.** *Fast!*

— Fast: defrag[T] bitmap largely cached in L2

Memory transactions: 2 memory reads + 1 write for relocated objects, 1 memory read for all others

**Step 5: Update Block State Bitmaps**

Blocks may now be empty, full and/or no longer defrag. candidates.

**Step 6: If there are > n defrag. candidates left, go to Step 1.**

**Generalization: Other Allocators?**

Many other GPU allocators (Halloc, ScatterAlloc) use hashing (**very high frag.**) and **do not utilize SOA.**