

# CompactGpu: Massively Parallel Memory Defragmentation on GPUs



東京工業大学  
Tokyo Institute of Technology

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

## 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

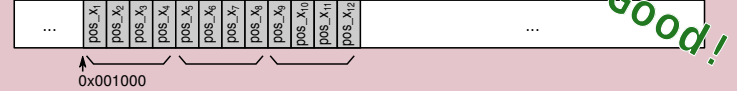
## Related Work

- 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?

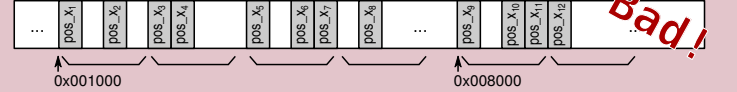
## Background: GPU Architecture and Dyn. Memory Allocation

- 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

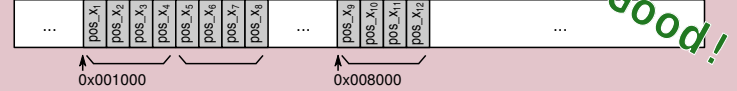
(a) Compact SOA Layout: 3 memory transactions required



(b) Fragmented SOA Layout: 6 memory transactions required

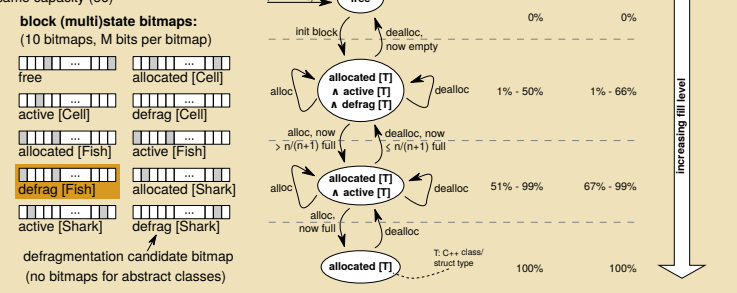
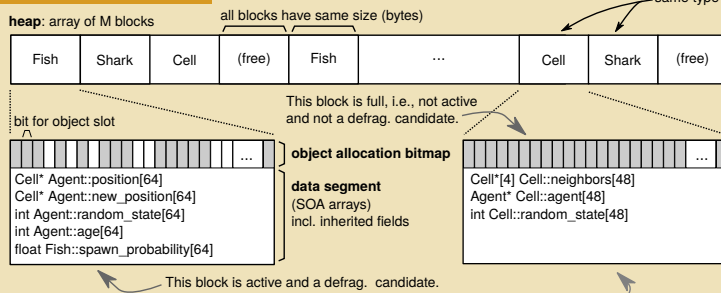


(c) Clustered SOA Layout: 3 memory transactions required



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

## DynaSOAr Heap Layout



## Defragmentation by Block Merging: parallel\_defrag<Fish>()

Running example:  
Fish-and-Sharks simulation

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

**Definition of Defragmentation Candidates:** Depends on defrag. factor  $n$  (problem-spec., compile-time parameter)

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

