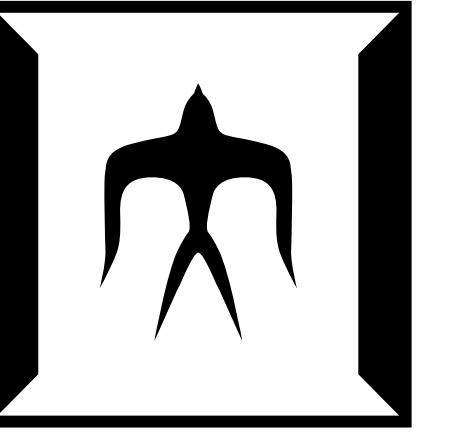


Ikra: Leveraging Object-oriented Abstractions



in a Ruby-to-CUDA JIT Translator



Matthias Springer, Hidehiko Masuhara
(Tokyo Institute of Technology)

Overview

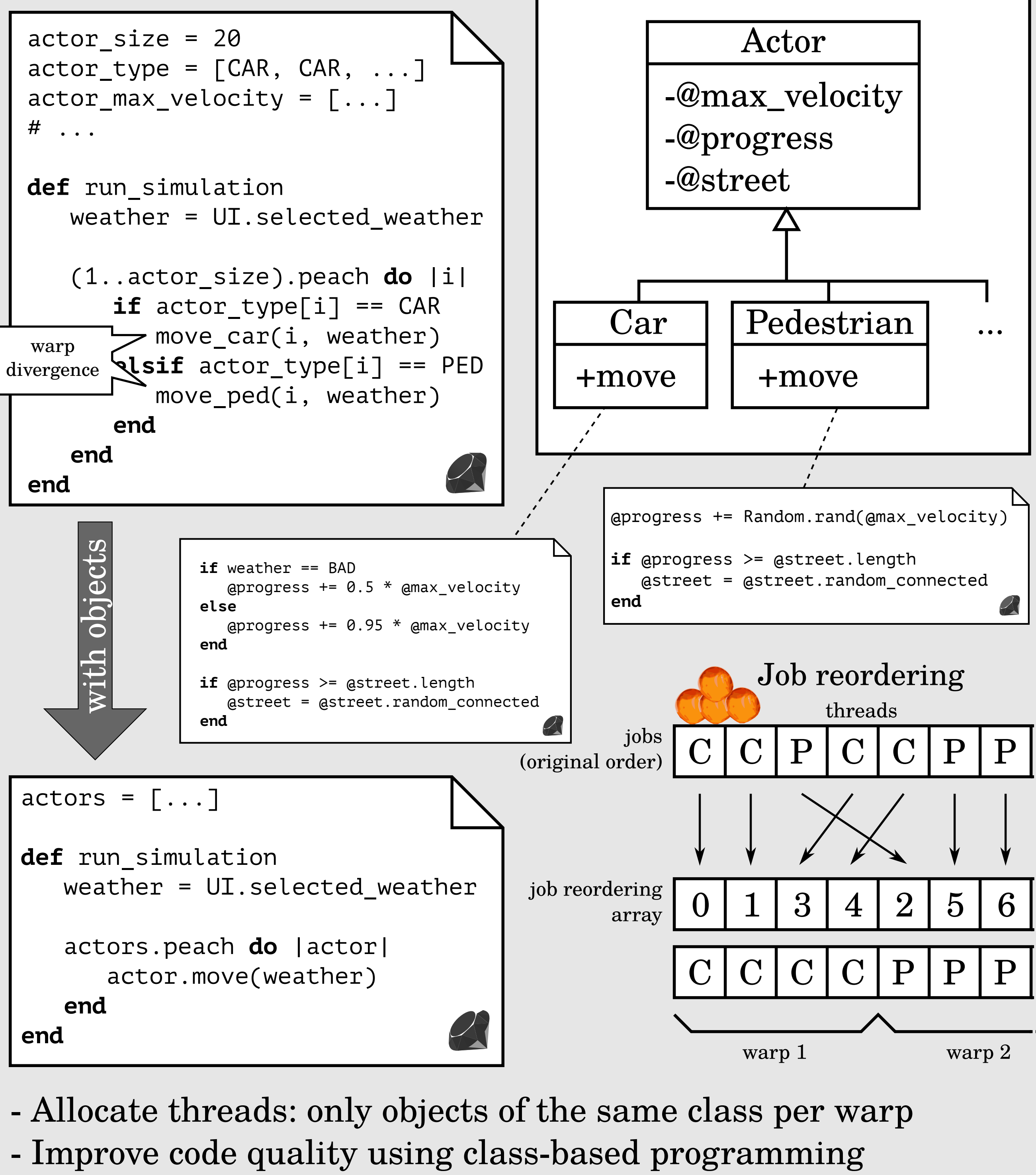
- Acceleration of Ruby programs with GPUs (CUDA)
- High-level Goal: GPGPU for Ruby programmers
- Source code analysis and type inference at runtime

Restrictions in parallel sections: No restrictions
 - No meta programming/reflection outside of
 - Dynamic typing should be avoided parallel sections

Related Work:

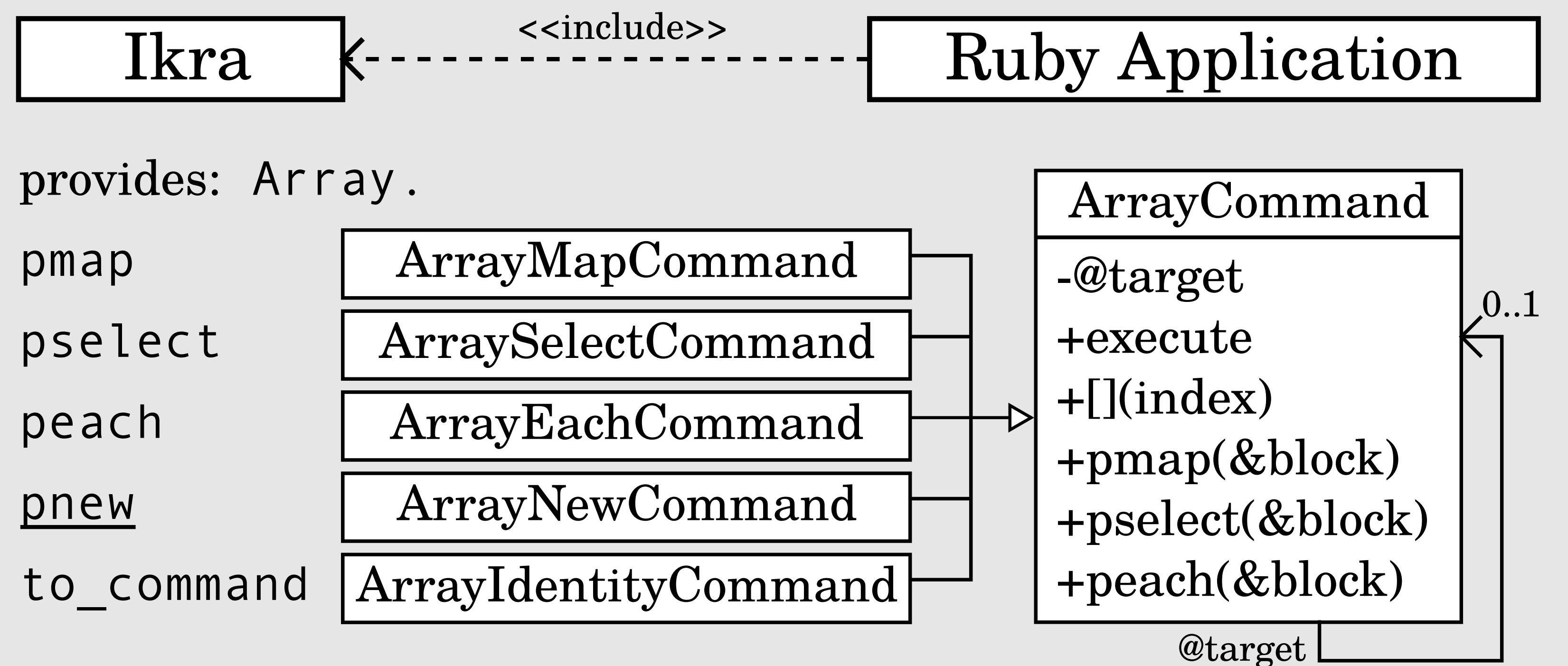
- Job Reordering
Zhang et al. On-the-fly elimination of dynamic irregularities for GPU computing. ASPLOS XVI.
- Kernel Fusion
Wahib et al. Scalable kernel fusion for memory-bound GPU applications. SC '14.
- Columnar Object Layout
Mattis et al. Columnar objects: improving the performance of analytical applications. ONWARD! 2015

Object Support



- Allocate threads: only objects of the same class per warp
- Improve code quality using class-based programming

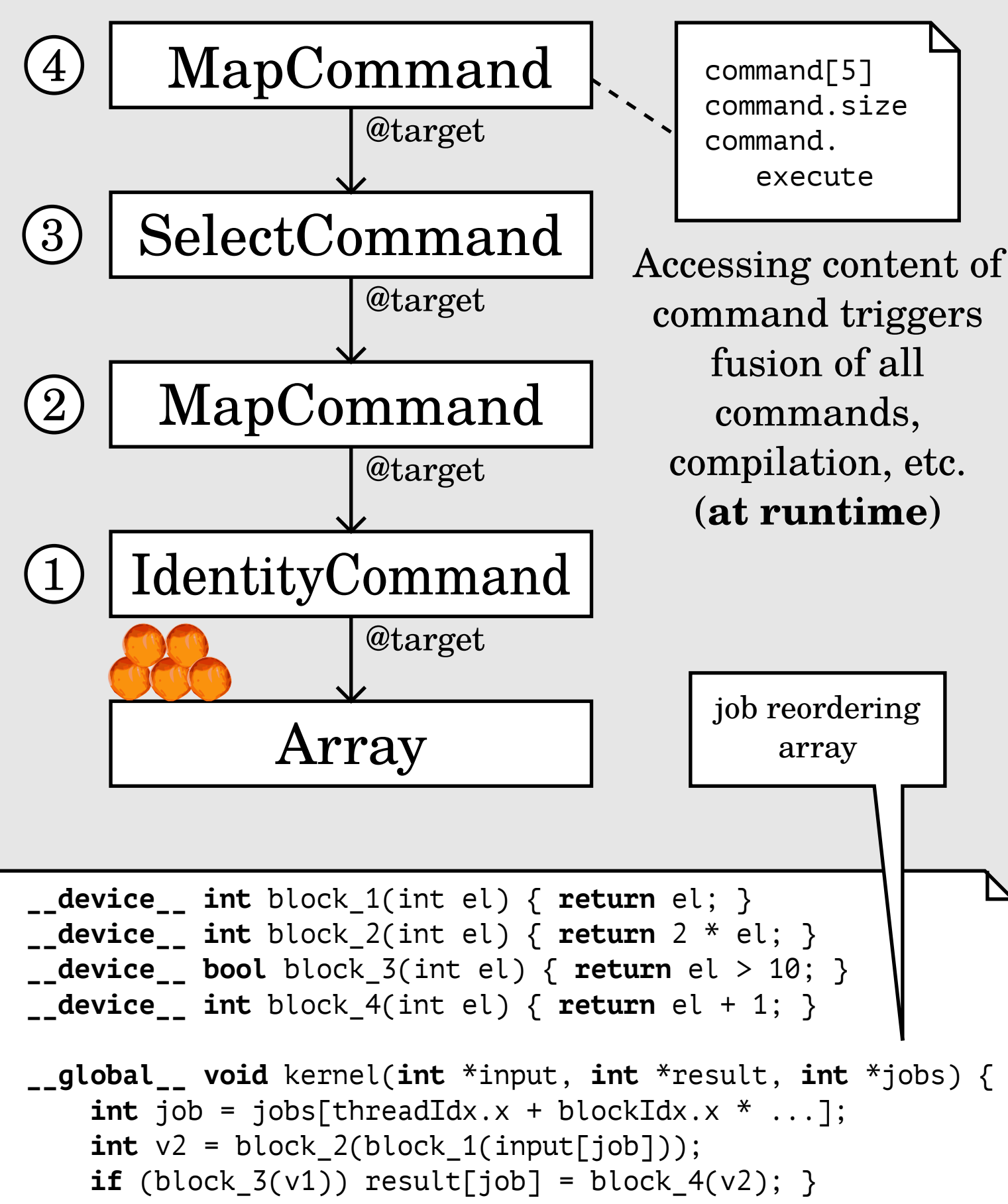
Architecture



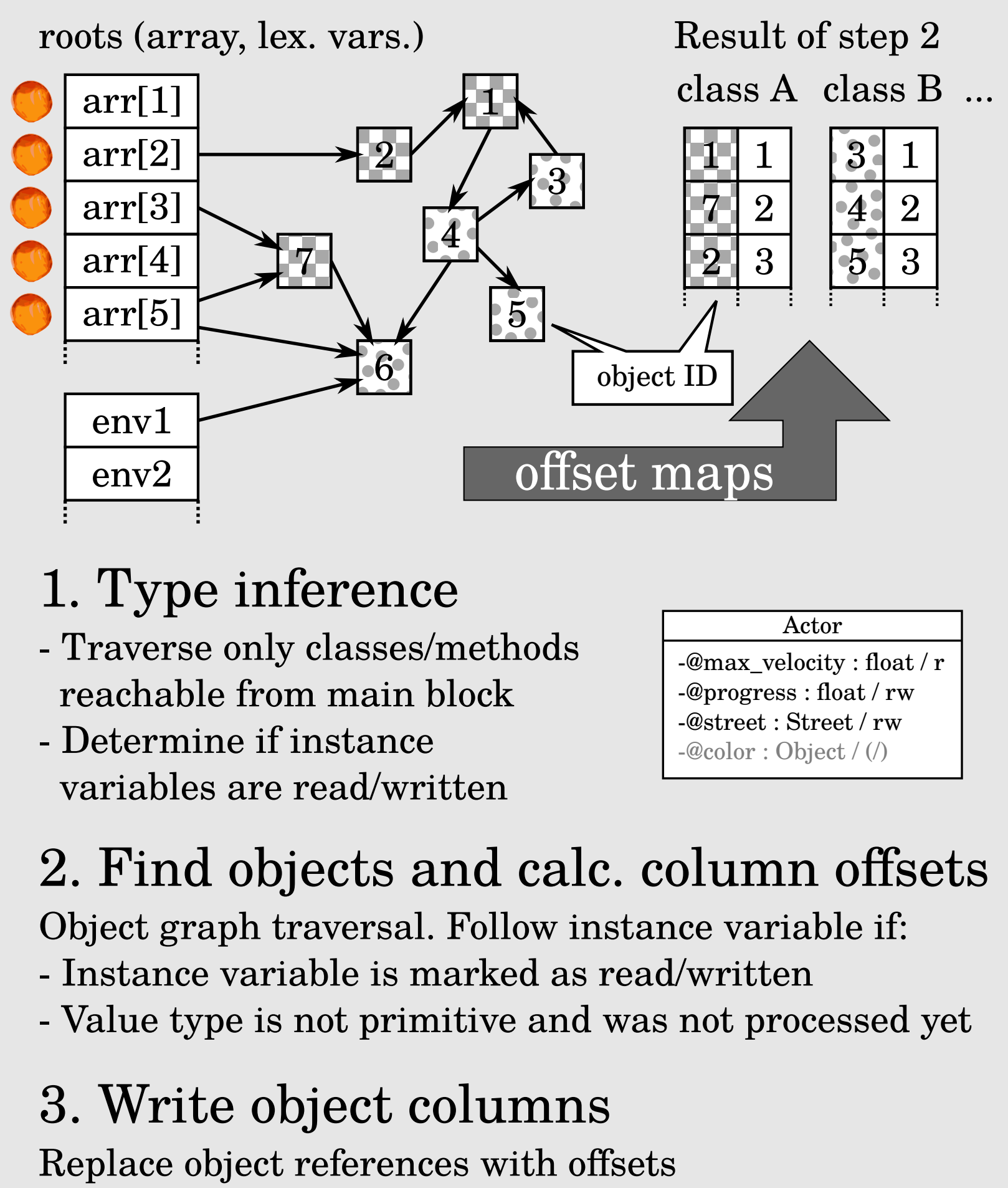
Overview:

1. Merge cascaded commands (**kernel fusion**)
2. Infer types and whether instance variables are read or written
3. Generate CUDA code
4. Compile shared library (nvcc)
5. Reorder jobs (avoiding **warp divergence**)
6. Trace reachable objects, allocate and transfer objects (Ruby FFI)
7. Invoke kernel
8. Write back written columns

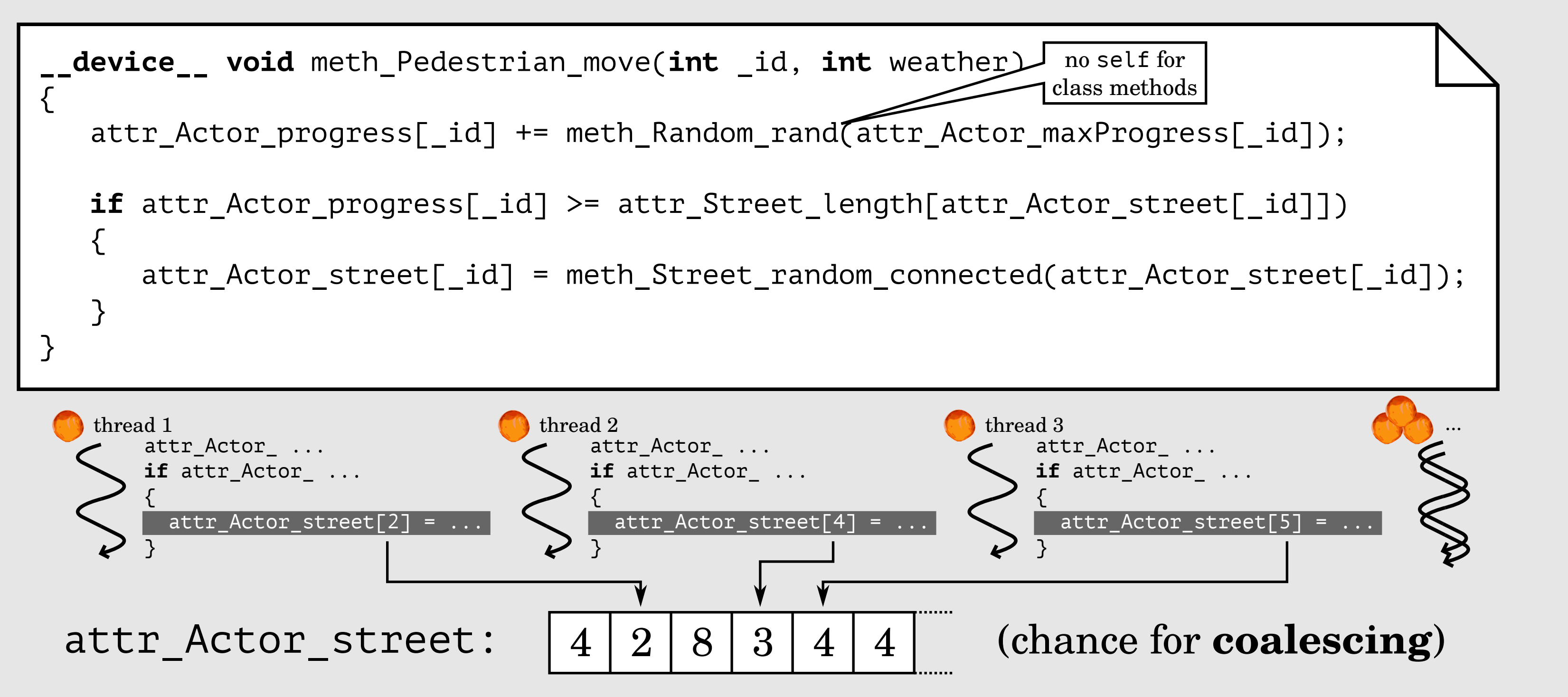
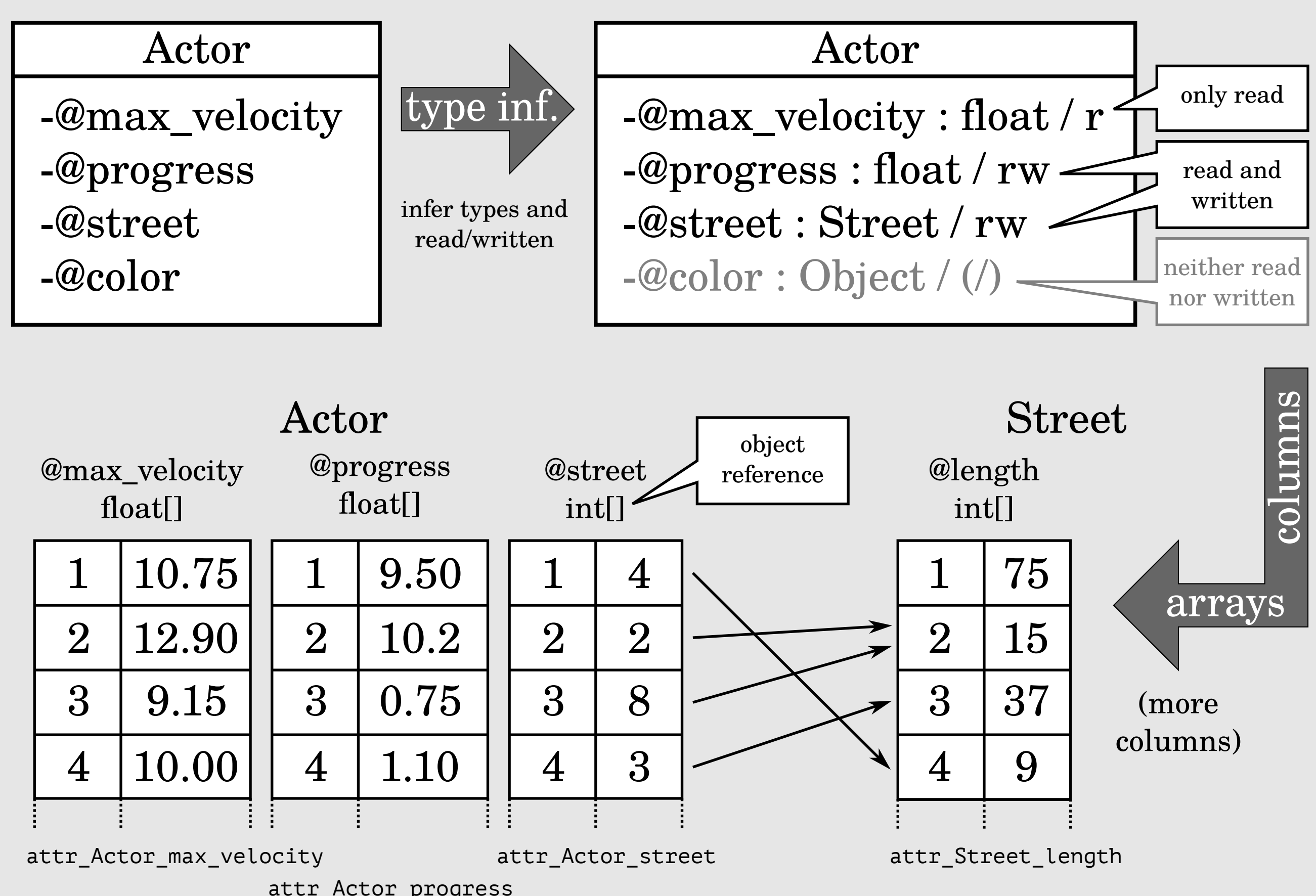
Kernel Fusion



Heap Object Tracer



Columnar Object Layout



- Idea: Represent all objects of a class as fields of arrays (columnar layout)
- Benefit: Chance for coalescing when accessing the same column in parallel
- Implementation: Heap Object Tracer converts object graph to columnar layout