# Portable HPC for High-Performance Simulation

Federico Carminati
Future of Many-Core
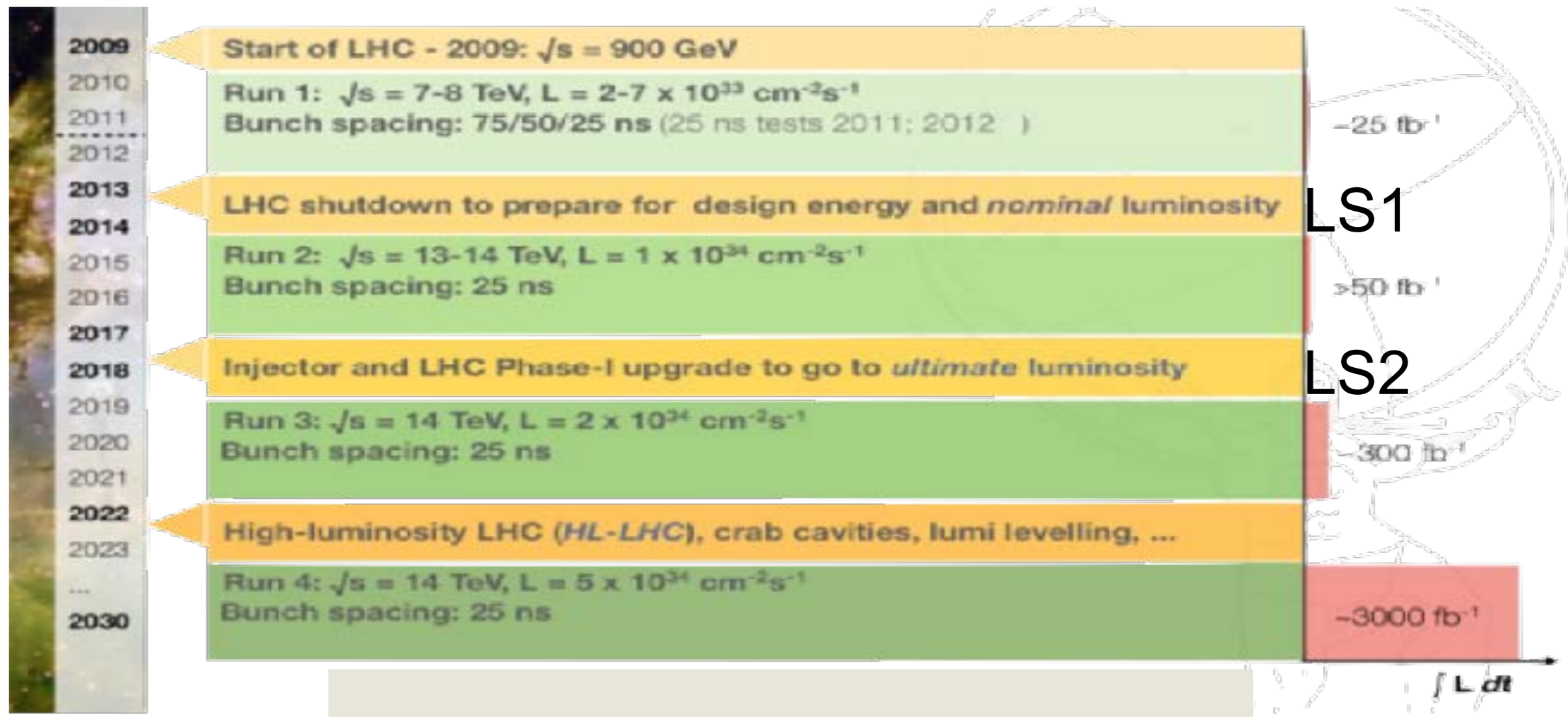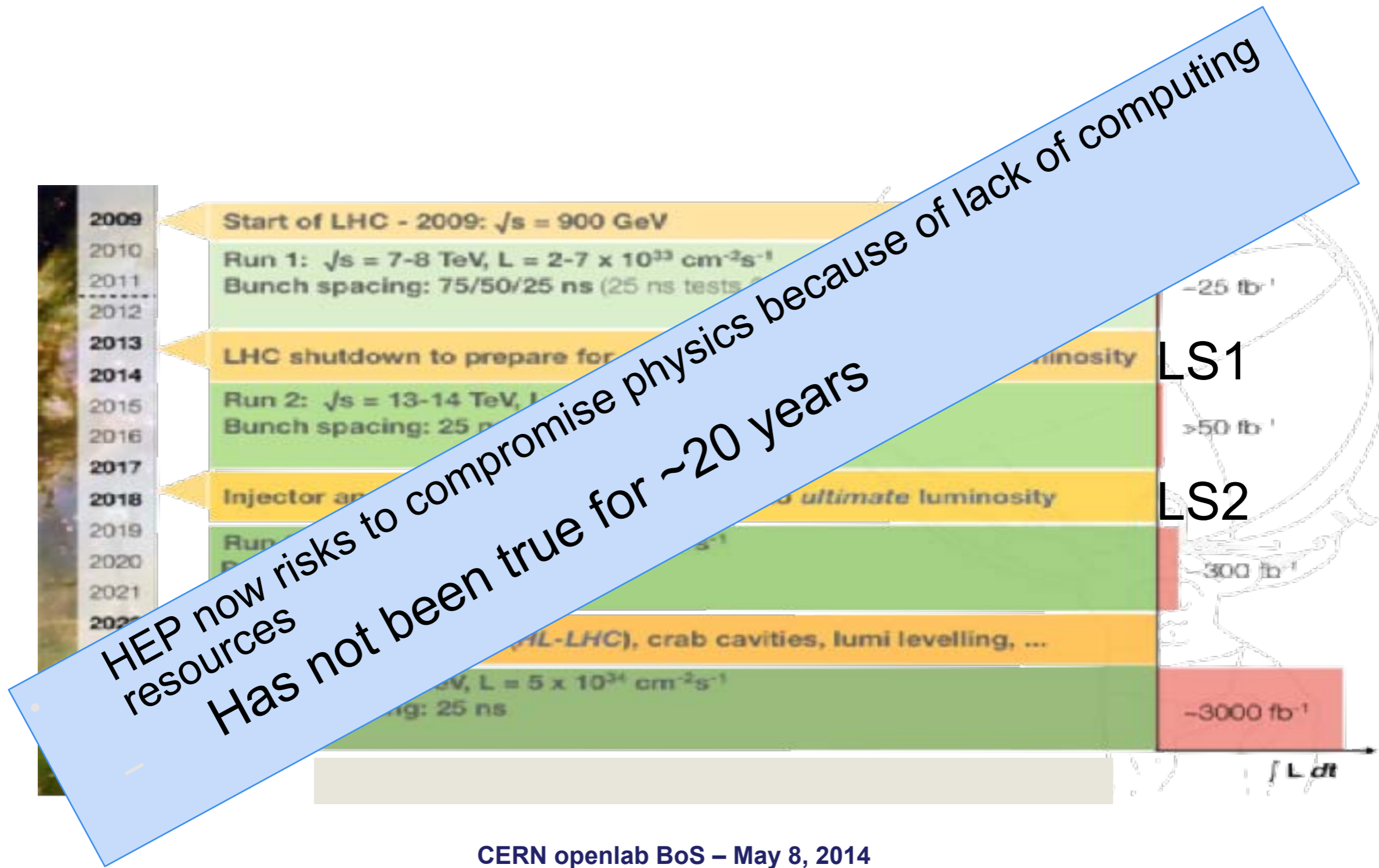Computing in Science
29 May 2014 Wigner Datacenter

# Thanks to

- John Apostolakis (CERN)
- Marilena Bandieramonte (Catania U.)
- Georgios Bitzes (CERN openlab)
- Rene Brun (CERN)
- Philippe Canal (FNAL)
- Federico Carminati (CERN)
- Johannes Christof De Fine Licht (CERN)
- Laurent Duhem (Intel)
- Daniel Elvira (FNAL)
- Andrei Gheata (CERN)
- Soon Yung Jun (FNAL)
- Guilherme Lima (FNAL)
- Mihaly Novak (CERN)
- Raman Sehgal (BARC)
- Oksana Shadura (CERN)
- Sandro Christian Wenzel (CERN)

# A luminous future for HEP...



| Year | Description | | Integrated Luminosity |
|------|-------------|------|----------------------|
| 2009 | Start of LHC - 2009: $\sqrt{s} = 900$ GeV | | |
| 2010–2012 | Run 1: $\sqrt{s} = 7\text{-}8$ TeV, L = $2\text{-}7 \times 10^{33}$ cm$^{-2}$s$^{-1}$. Bunch spacing: 75/50/25 ns (25 ns tests 2011: 2012 ) | | ~25 fb$^{-1}$ |
| 2013–2014 | LHC shutdown to prepare for design energy and *nominal* luminosity | LS1 | |
| 2015–2016 | Run 2: $\sqrt{s} = 13\text{-}14$ TeV, L = $1 \times 10^{34}$ cm$^{-2}$s$^{-1}$. Bunch spacing: 25 ns | | >50 fb$^{-1}$ |
| 2018 | Injector and LHC Phase-I upgrade to go to *ultimate* luminosity | LS2 | |
| 2019–2021 | Run 3: $\sqrt{s} = 14$ TeV, L = $2 \times 10^{34}$ cm$^{-2}$s$^{-1}$. Bunch spacing: 25 ns | | ~300 fb$^{-1}$ |
| 2022–2023 | High-luminosity LHC (*HL-LHC*), crab cavities, lumi levelling, ... | | |
| ...2030 | Run 4: $\sqrt{s} = 14$ TeV, L = $5 \times 10^{34}$ cm$^{-2}$s$^{-1}$. Bunch spacing: 25 ns | | ~3000 fb$^{-1}$ |

$\int L\,dt$

# A luminous future for HEP...



Start of LHC - 2009: $\sqrt{s}$ = 900 GeV

Run 1: $\sqrt{s}$ = 7-8 TeV, L = 2-7 x $10^{33}$ cm$^{-2}$s$^{-1}$
Bunch spacing: 75/50/25 ns (25 ns tests)

LHC shutdown to prepare for ... inosity    LS1

Run 2: $\sqrt{s}$ = 13-14 TeV, ...
Bunch spacing: 25 n...

Injector an... ultimate luminosity    LS2

Run...

...HL-LHC), crab cavities, lumi levelling, ...

...V, L = 5 x $10^{34}$ cm$^{-2}$s$^{-1}$
...g: 25 ns

~25 fb$^{-1}$
>50 fb$^{-1}$
~300 fb$^{-1}$
~3000 fb$^{-1}$

$\int L\, dt$

2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
202...

**HEP now risks to compromise physics because of lack of computing resources**
**Has not been true for ~20 years**

# The Eight dimensions

- The "dimensions of performance"
  - Vectors
  - Instruction Pipelining
  - Instruction Level Parallelism (ILP)
  - Hardware threading
  - Clock frequency
  - Multi-core
  - Multi-socket
  - Multi-node

Micro-parallelism: gain in throughput and in time-to-solution

Very little gain to be expected and no action to be taken

Gain in memory footprint and time-to-solution but not in throughput

Possibly running different jobs as we do now is the best solution

# The Eight dimensions

- The "dimensions of performance"
  - Vectors
  - Instruction Pipelining
  - Instruction Level Parallelism (ILP)
  - Hardware threading
  - Clock frequency
  - Multi-core
  - Multi-socket
  - Multi-node

Micro-parallelism: gain in throughput and in time-to-solution

Very little gain to be expected and no action to be taken

Gain in memory footprint and time-to-solution but not in throughput

Possibly running different jobs as we do now is the best solution

Expected limits on performance scaling

|  | SIMD | ILP | HW |
|---|---|---|---|
| THEORY | 8 | 4 | 1.35 |
| OPTIMISED | 6 | 1.57 | 1.25 |
| HEP | 1 | 0.8 | 1.25 |

Expected limits on performance scaling (multiplied)

|  | SIMD | ILP | HW |
|---|---|---|---|
| THEORY | 8 | 32 | 43.2 |
| OPTIMISED | 6 | 9.43 | 11.79 |
| HEP | 1 | 0.8 | 1 |

OpenLab@CHEP12

# Classical HEP transport is mostly local

- Event- or event track-level parallelism will better use resources but won't improve these points

> - Geometry navigation (local)
> - Material – X-section tables
> - Particle type - physics processes

> 50 per cent of the time spent in 0.7% volumes

ATLAS volumes sorted by transport time. The same behavior is observed for most HEP geometries.

- Navigating very large data structures
- No locality
- OO abused: very deep instruction stack
- Cache misses

# Introduced basketized transport

Deal with particles in parallel

Output buffer(s)

# Introduced basketized transport

Deal with particles in parallel

Output buffer(s)

# Introduced basketized transport

Deal with particles in parallel

Output buffer(s)

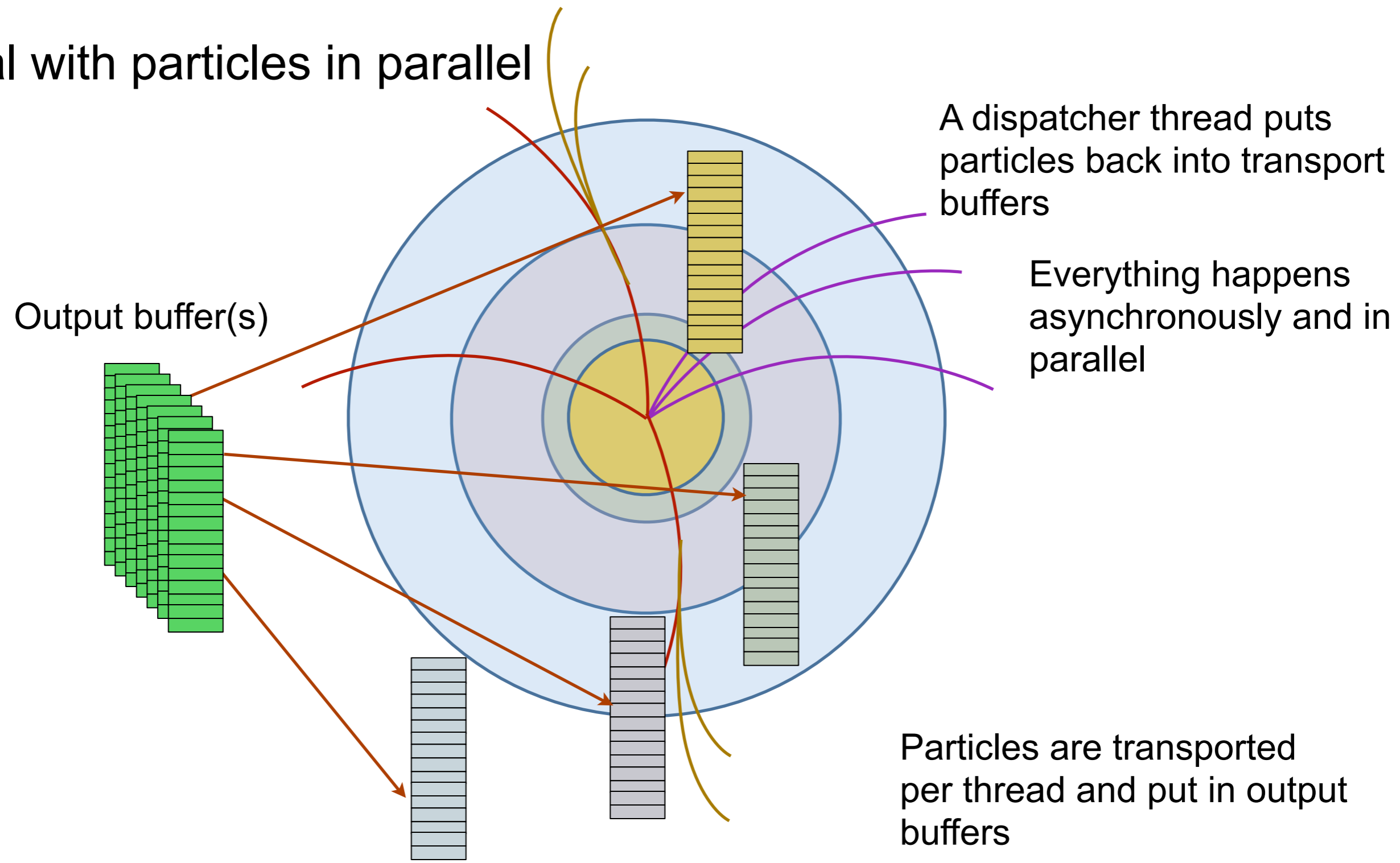Particles are transported per thread and put in output buffers

6

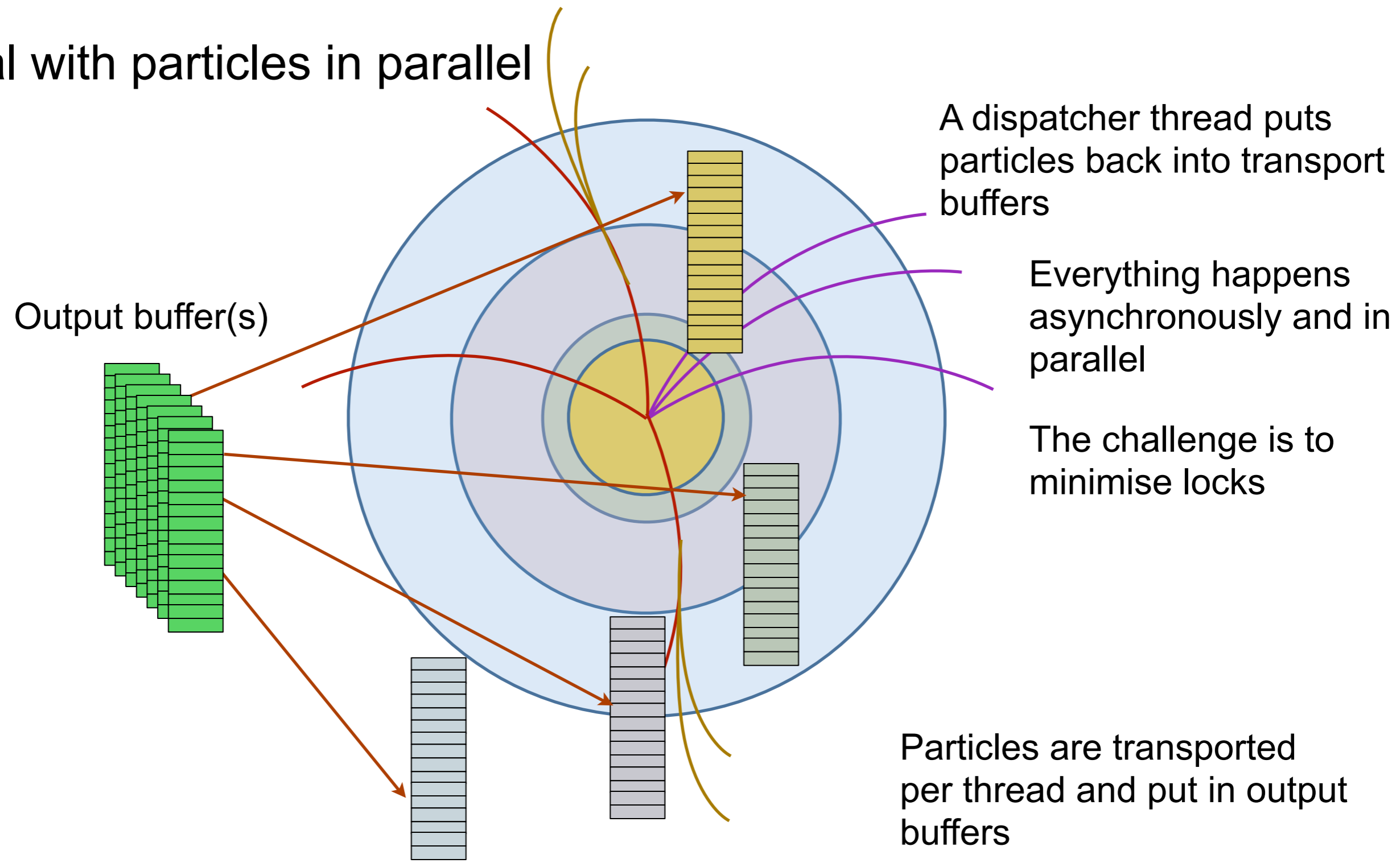# Introduced basketized transport

Deal with particles in parallel

A dispatcher thread puts particles back into transport buffers

Output buffer(s)
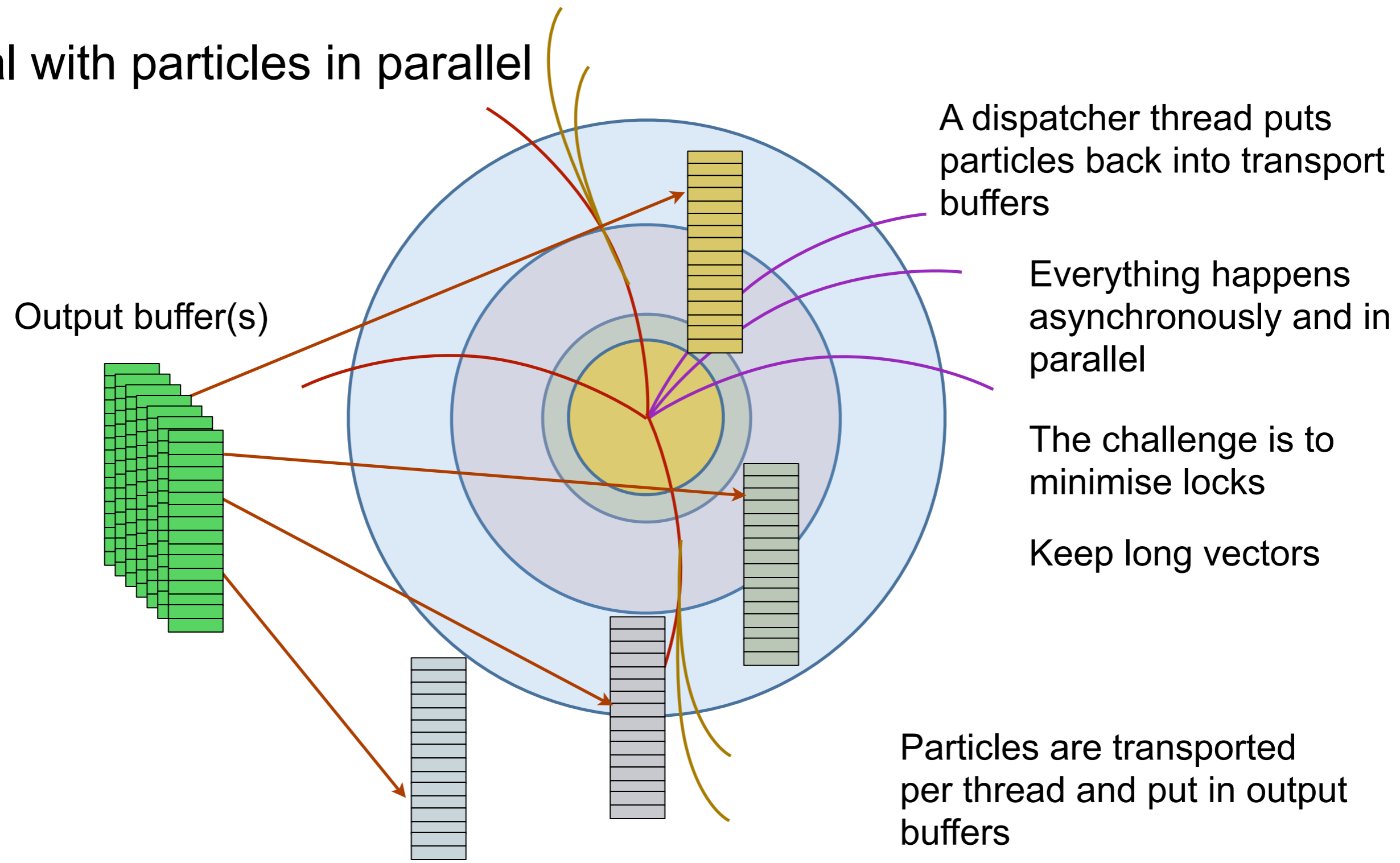
Particles are transported per thread and put in output buffers

# Introduced basketized transport

**Deal with particles in parallel**

Output buffer(s)

A dispatcher thread puts particles back into transport buffers
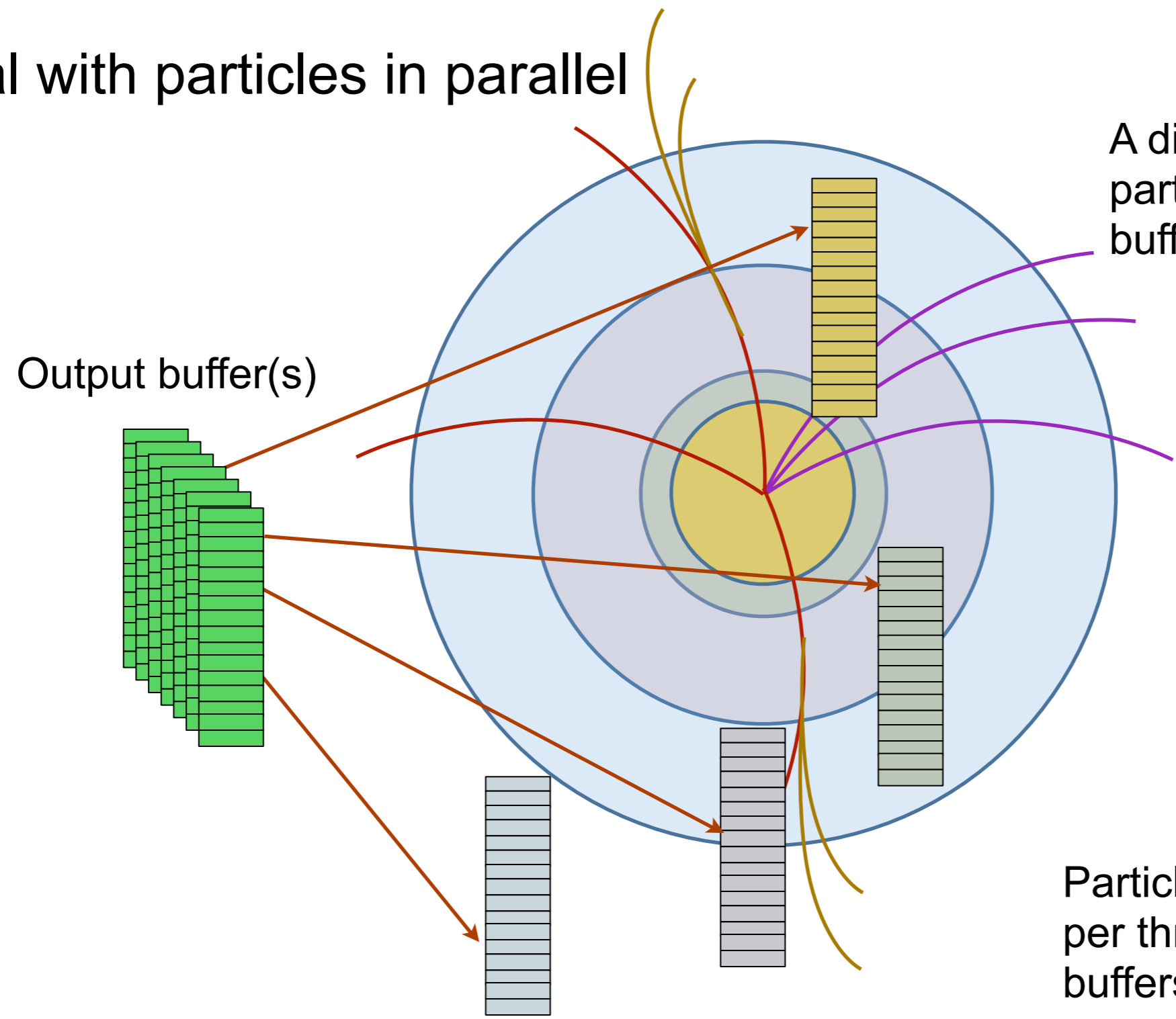
Everything happens asynchronously and in parallel

Particles are transported per thread and put in output buffers

# Introduced basketized transport

Deal with particles in parallel

Output buffer(s)

A dispatcher thread puts particles back into transport buffers

Everything happens asynchronously and in parallel

The challenge is to minimise locks

Particles are transported per thread and put in output buffers

# Introduced basketized transport

**Deal with particles in parallel**

Output buffer(s)
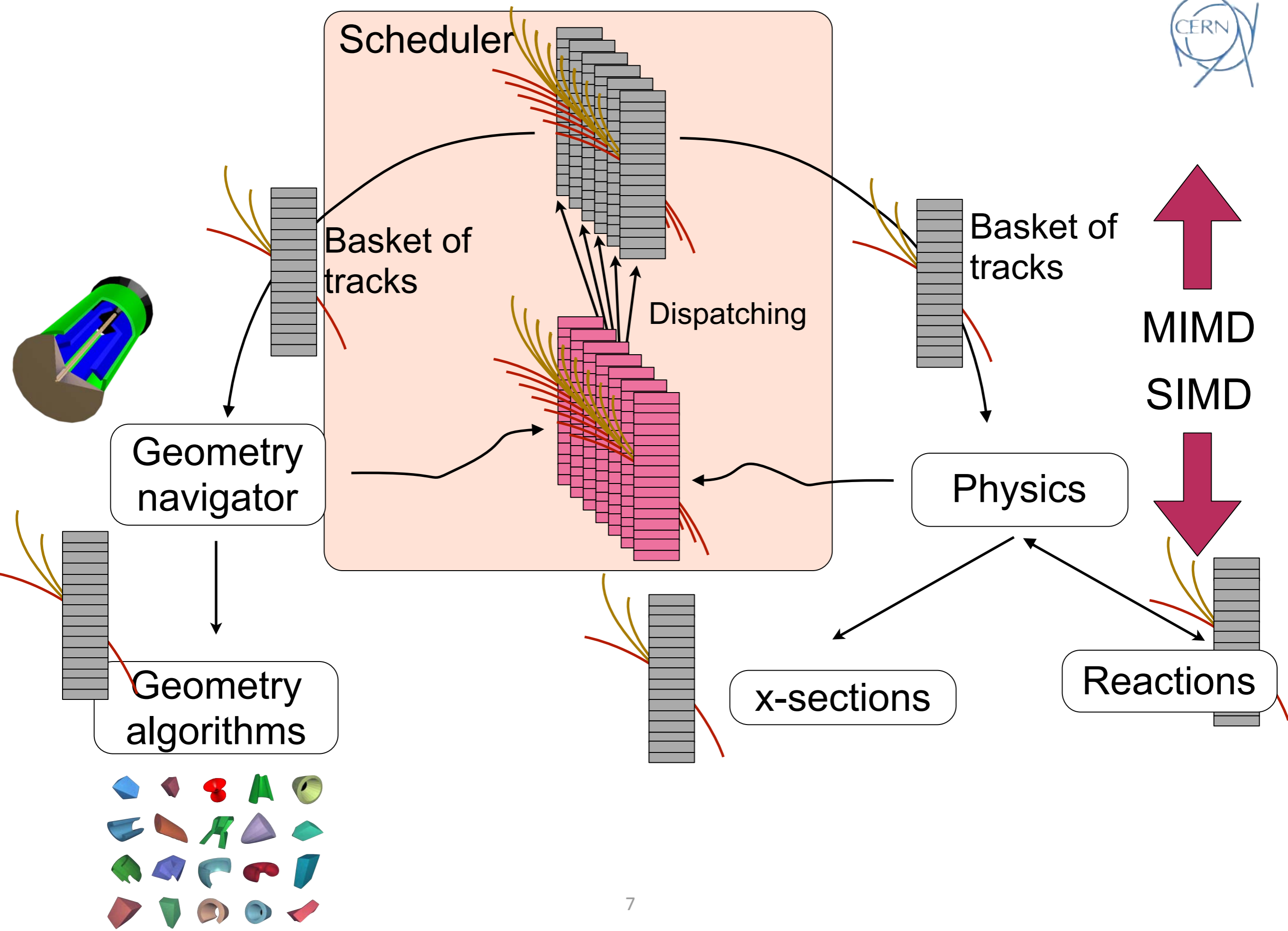
A dispatcher thread puts particles back into transport buffers
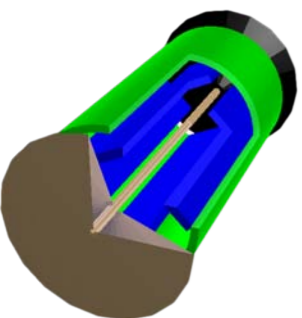
Everything happens asynchronously and in parallel

The challenge is to minimise locks

Keep long vectors

Particles are transported per thread and put in output buffers

# Introduced basketized transport

**Deal with particles in parallel**



Output buffer(s)

A dispatcher thread puts particles back into transport buffers

Everything happens asynchronously and in parallel

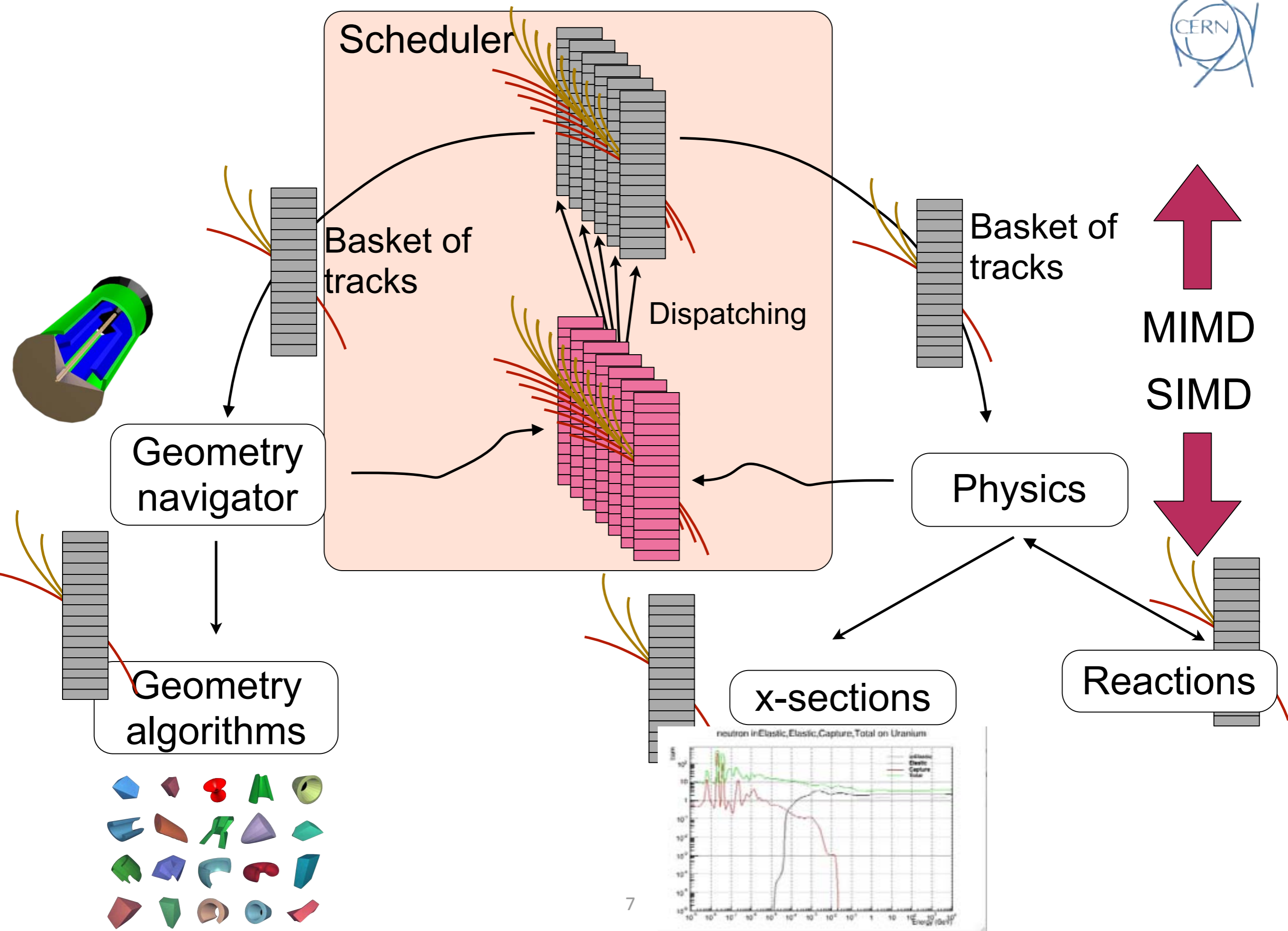The challenge is to minimise locks

Keep long vectors

Avoid memory explosion

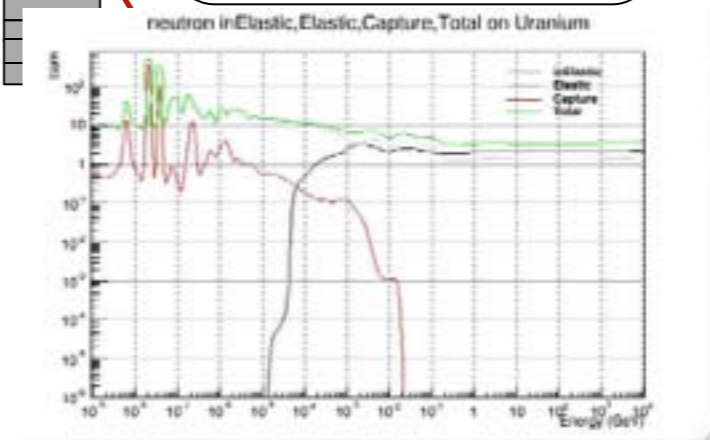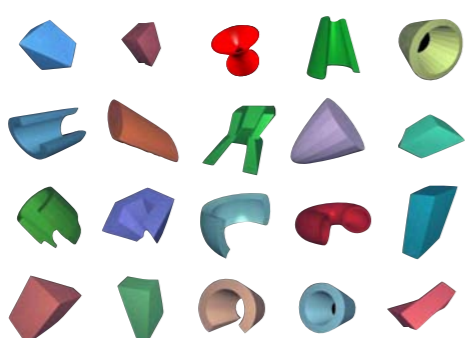Particles are transported per thread and put in output buffers

Scheduler

Basket of tracks

Dispatching

Basket of tracks

MIMD

SIMD

Geometry navigator

Physics

Geometry algorithms

x-sections

Reactions

Scheduler

Basket of tracks

Dispatching

Basket of tracks

Geometry navigator

Physics

MIMD

SIMD

Geometry algorithms

x-sections
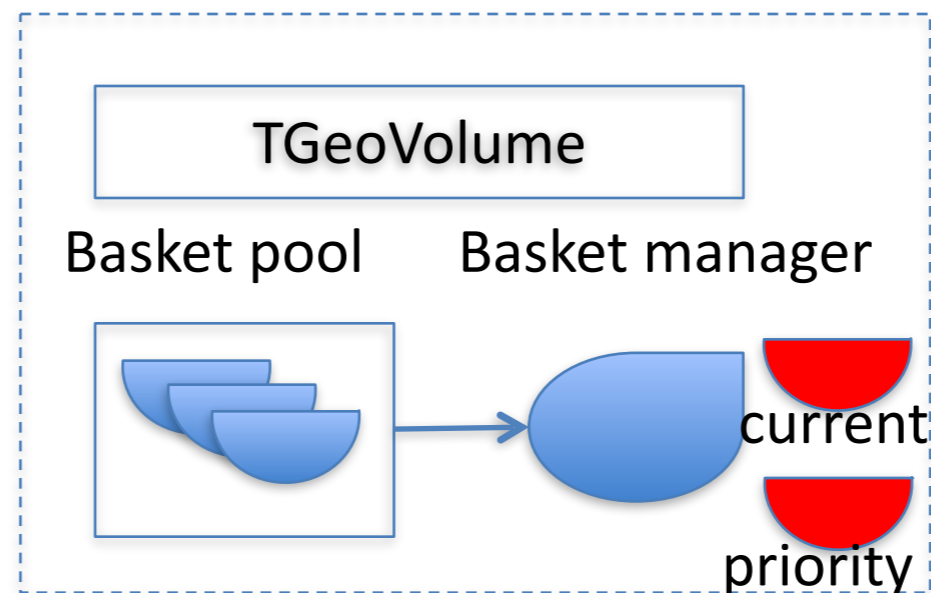
Reactions

neutron inElastic,Elastic,Capture,Total on Uranium
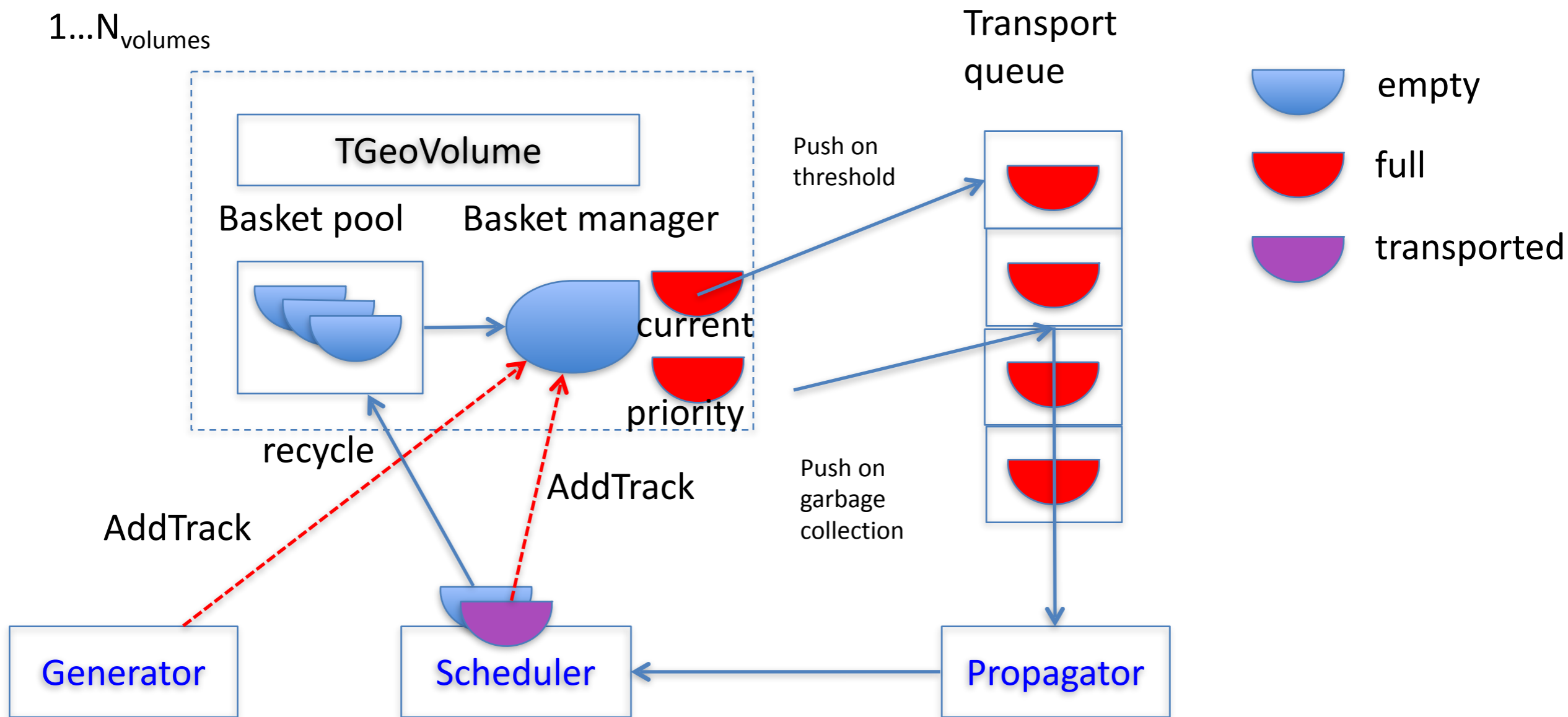
7

# Basket managers

- One basket manager per volume
  - Receiving tracks entering the volume from generator or scheduler
  - Accessed by scheduler only
- Pool of empty baskets, one current basket + one basket for prioritized tracks
- Lock-free access for unique scheduler (only one thread can add tracks)
- Transportability threshold per manager
  - If threshold reached when adding tracks, the current basket is pushed in the work queue and replaced from the pool. Tracks added with the priority flag go to the priority basket which gets pushed to the priority side of the queue
  - Threshold(vol) = Ntracks_in_flight(vol)/2N_threads rounded to %4 (min 4, max 256)

$1...N_{volumes}$
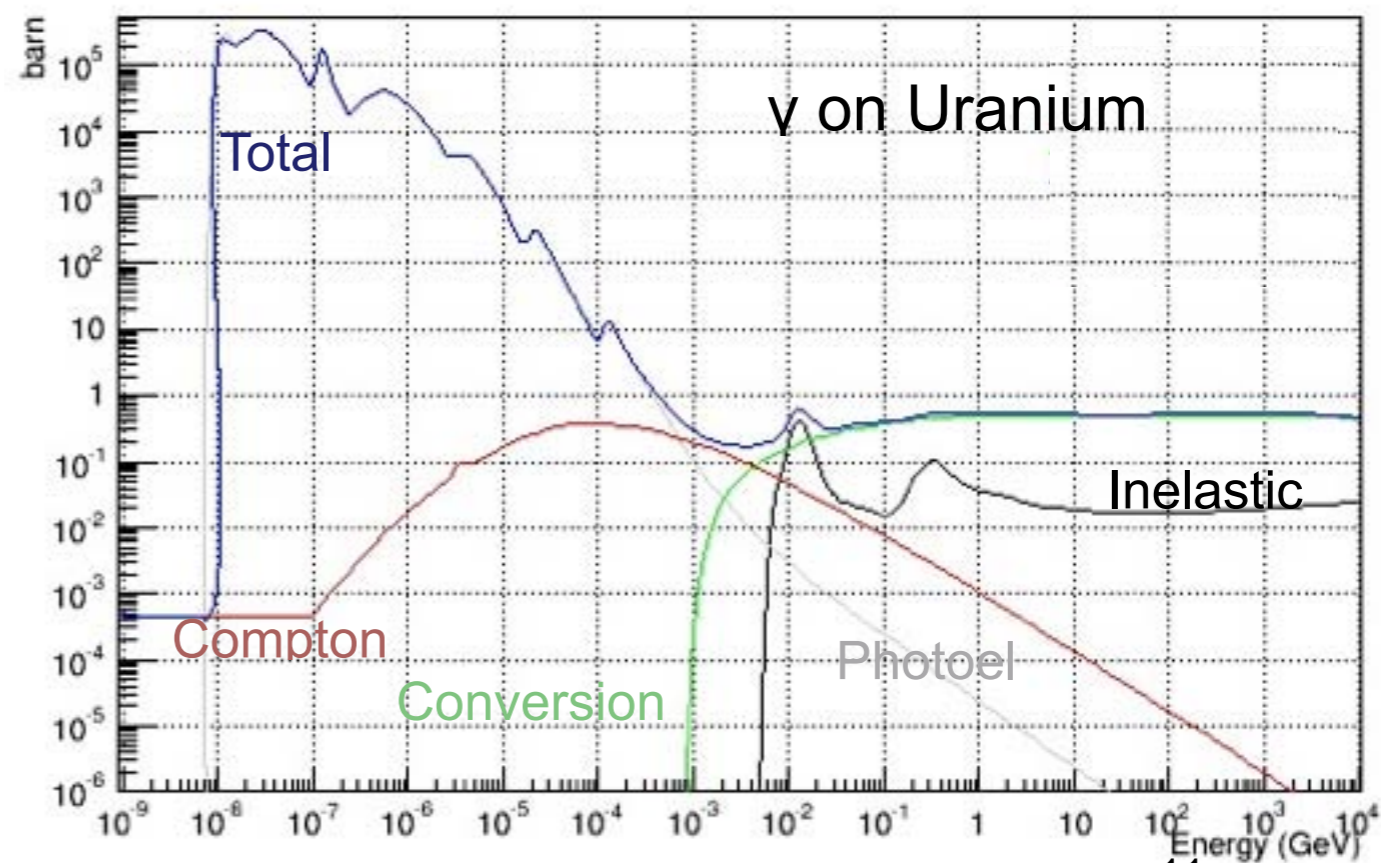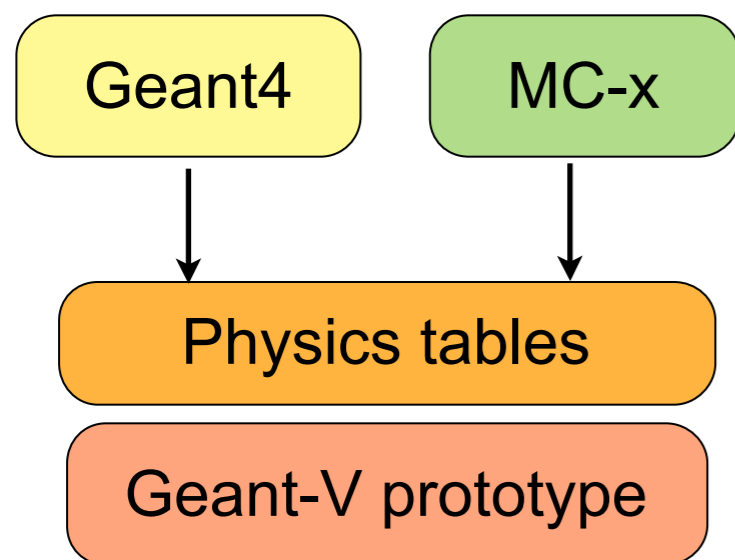
# Basket lifecycle

# Scheduling policies

- Workload balancing
  - Divide the work evenly to scale with number of workers
  - Queue control: garbage collection on work queue depletion
  - Improvement: schedule physics as separate task (process selection and discrete processes post-step)
- Memory management
  - Not active currently, the idea it to trigger hit/digits collection and memory cleanup on thresholds
- Keep large vectors
  - Raise transportability thresholds per volume
  - Postpone sparse tracks when not in garbage collection mode
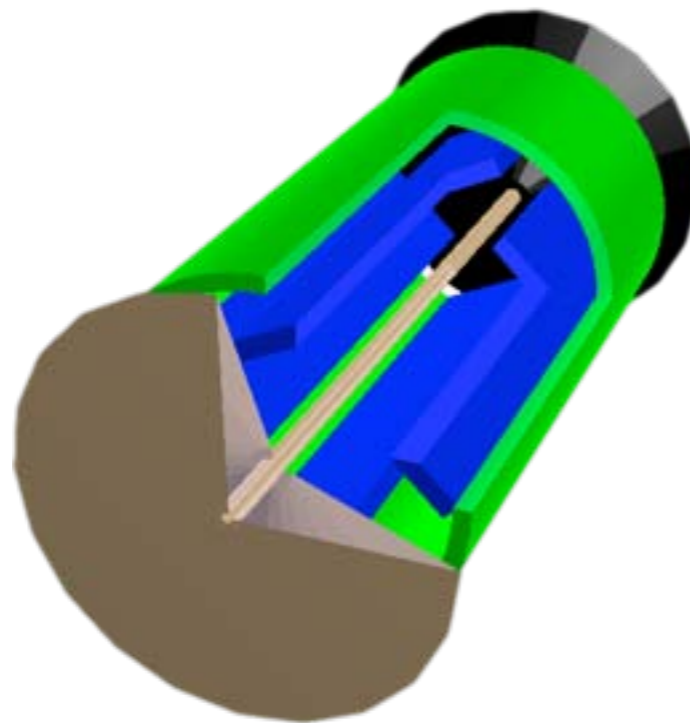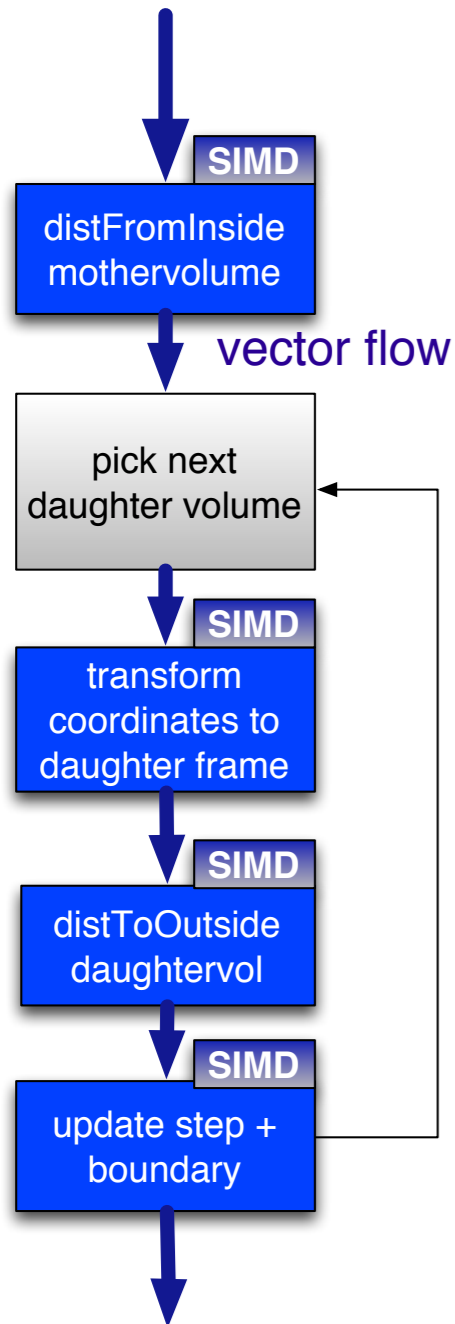- Trigger single track mode when vectorization gives just overhead

# Physics

- A lightweight physics for realistic shower development
- Select the major mechanisms
  - Bremsstrahlung, e+ annihilation, Compton, Decay, Delta ray, Elastic hadron, Inelastic hadron, Pair production, Photoelectric, Capture + dE/dx & MS
- Tabulate all x-secs (100 bins -> 90MB)
- Generate (10-50) final states (300kB per final state & element)
- Not good as Geant4, but it could be the seed of a fast simulation option
- Independent from the MonteCarlo that actually generates the tables



Geant4 → Physics tables

MC-x → Physics tables

Geant-V prototype

11

11

# Recap of performance status

* provided new optimized vector interfaces for some elementary solids and geometric base classes ( implemented important functions for particle navigation )

* overall performance gain in a standard navigation benchmark ( in toy detector with 4 boxes, 3 tubes, 2 cones ) - comparison to ROOT/ 5.34.17

**SIMD**
distFromInside mothervolume

vector flow

pick next daughter volume

**SIMD**
transform coordinates to daughter frame

**SIMD**
distToOutside daughtervol

**SIMD**
update step + boundary

|  | 16 particles | 1024 particles |
|---|---|---|
| Intel IvyBridge (AVX) | ~2.8x | ~4.0x |
| Intel Haswell (AVX2) | ~3.0x | ~5.0x |
| Intel Xeon-Phi (AVX512) | ~4.1x | ~4.8x |

**Xeon-Phi and Haswell benchmarks by CERN Openlab (Georgios Bitzes)**

**CHEP13 paper: http://arxiv.org/pdf/1312.0816.pdf**
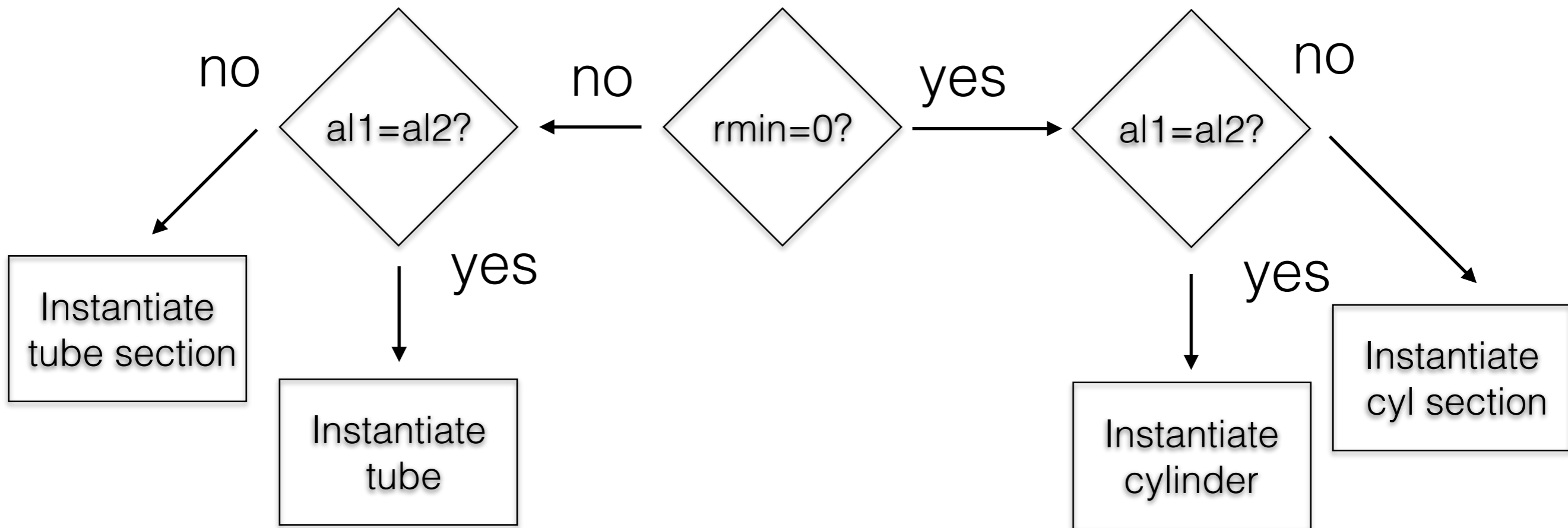
# Portable HPC?

- Straight "vectorisation" of existing code is difficult to impossible
- Resulting code is hard to read and maintain
  - And it is largely compiler-dependency
- Porting to different high end devices is very difficult
- Explored solution is to use template specification for solid placement, specialisation and code generation
- Highly optimised modular "codelets" à la STL are used to construct algorithms

# Solid specialisation

CreateTube(rmin, rmax, al1, al2)

no

al1=al2?

no

rmin=0?

yes

al1=al2?

no

Instantiate tube section

yes

Instantiate tube

yes

Instantiate cylinder

Instantiate cyl section

# Illustrating scalar/SIMD abstraction and kernels
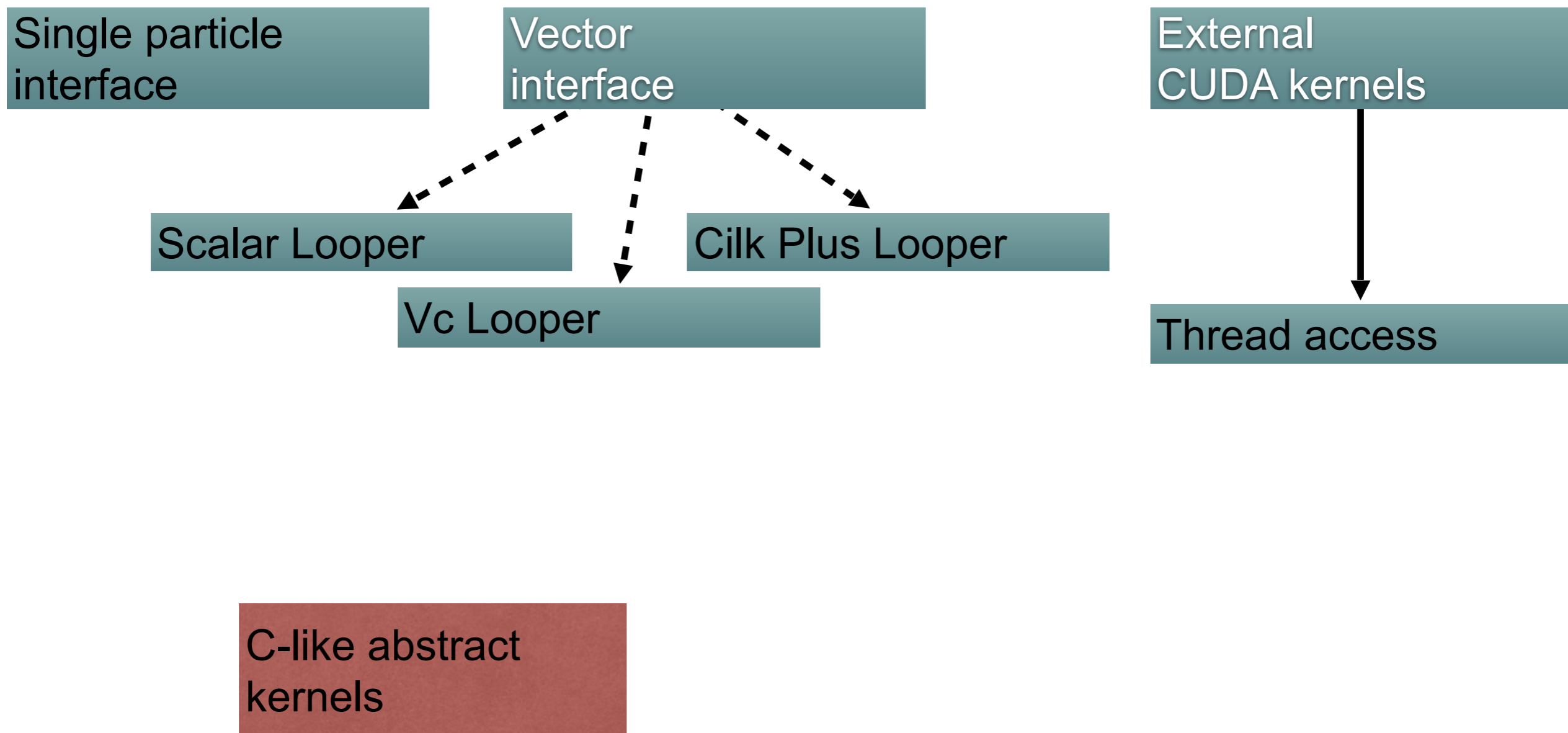
Single particle
interface

Vector
interface

External
CUDA kernels

C-like abstract
kernels

# Illustrating scalar/SIMD abstraction and kernels

Single particle interface

Vector interface

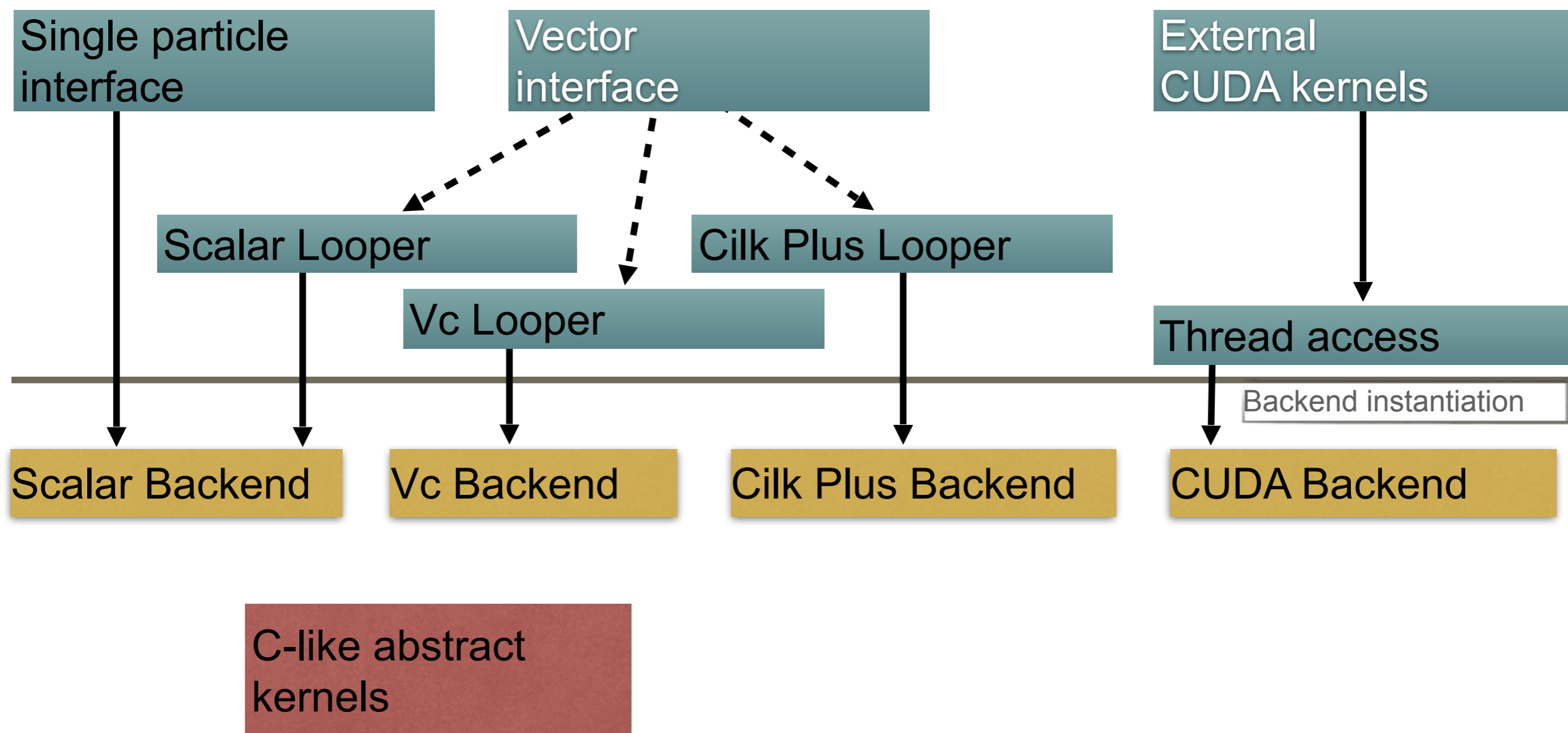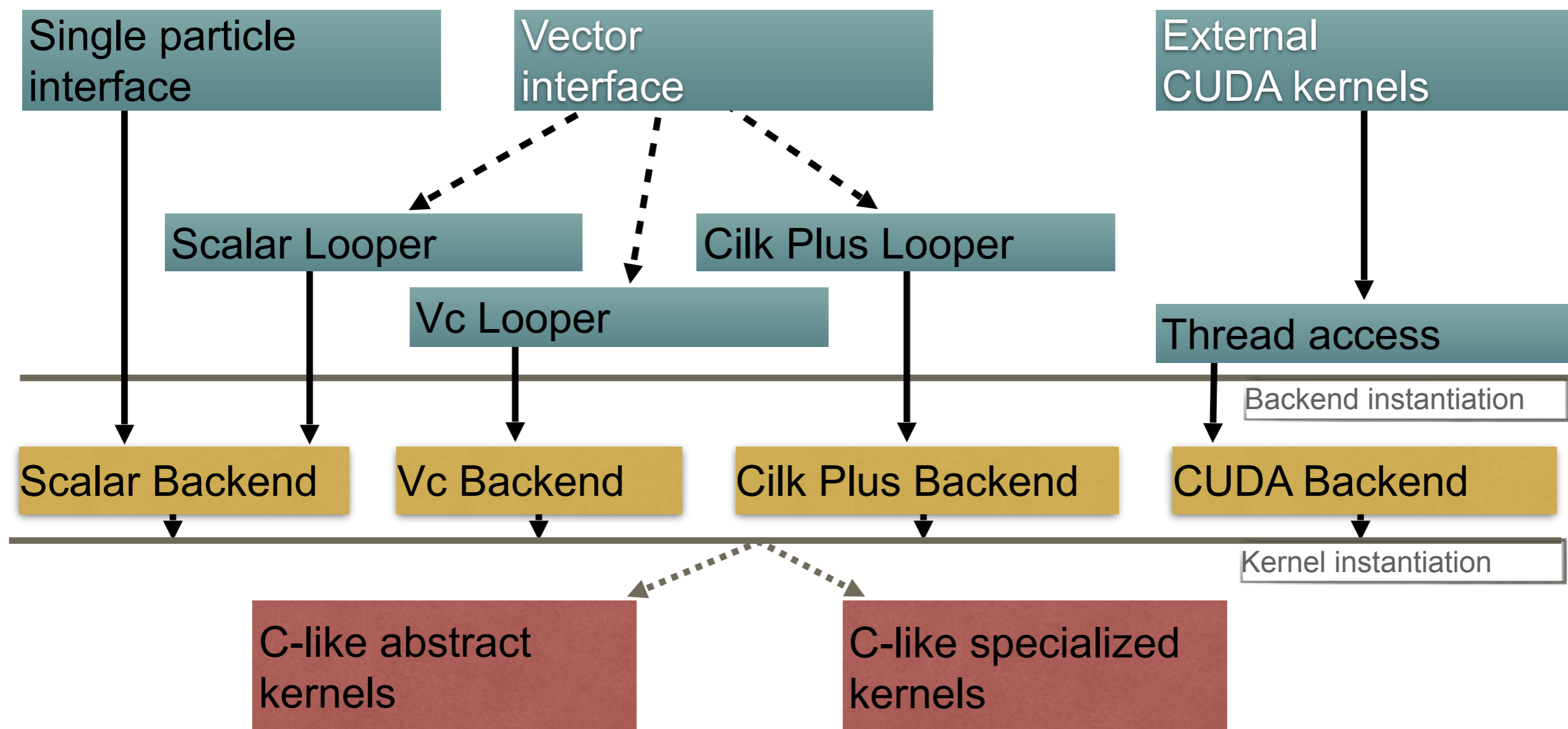External CUDA kernels

Scalar Looper

Cilk Plus Looper

Vc Looper

Thread access

C-like abstract kernels

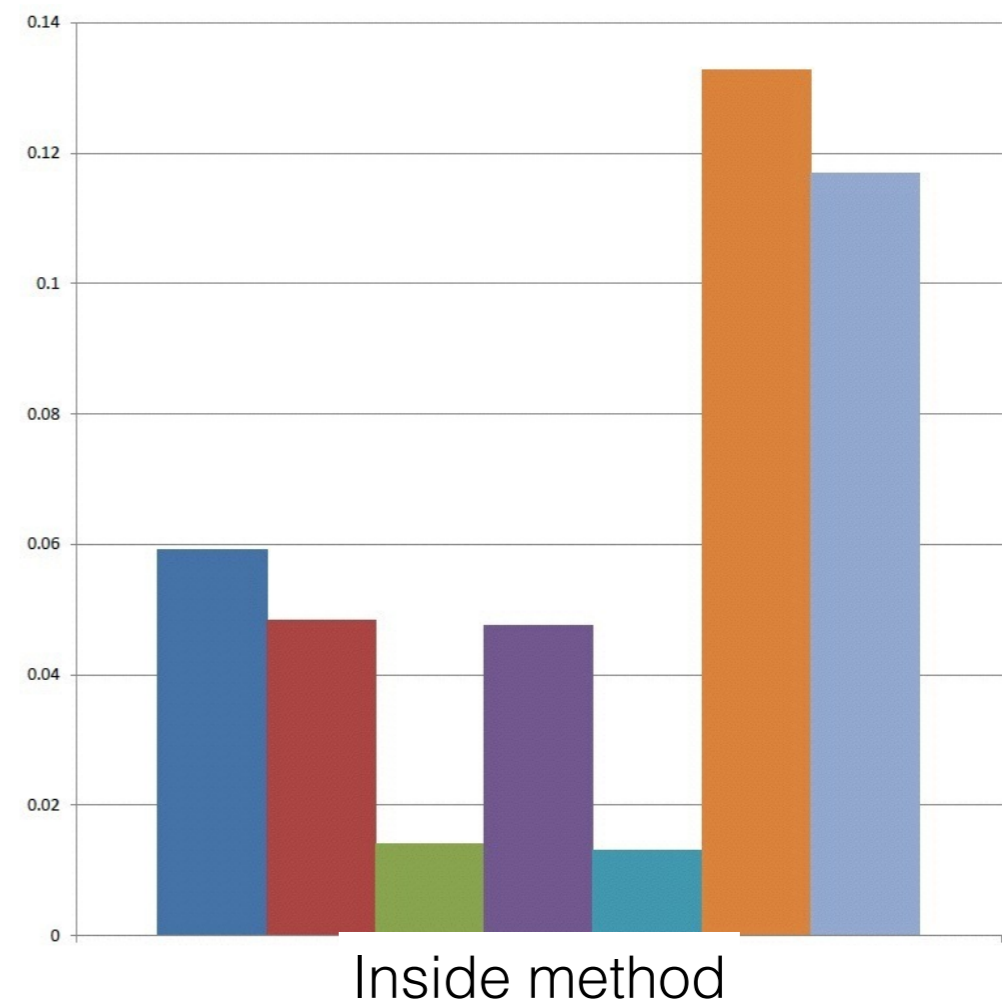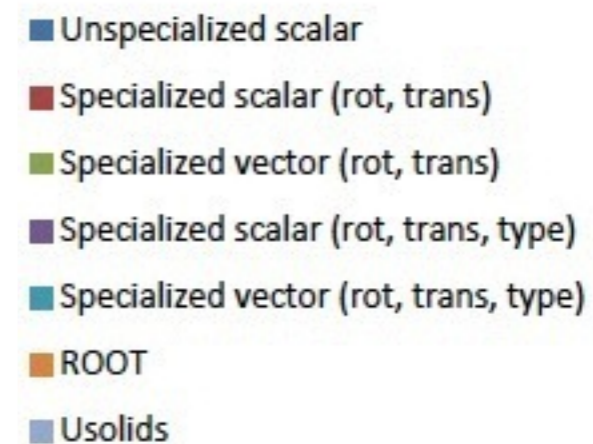# Illustrating scalar/SIMD abstraction and kernels

# Illustrating scalar/SIMD abstraction and kernels
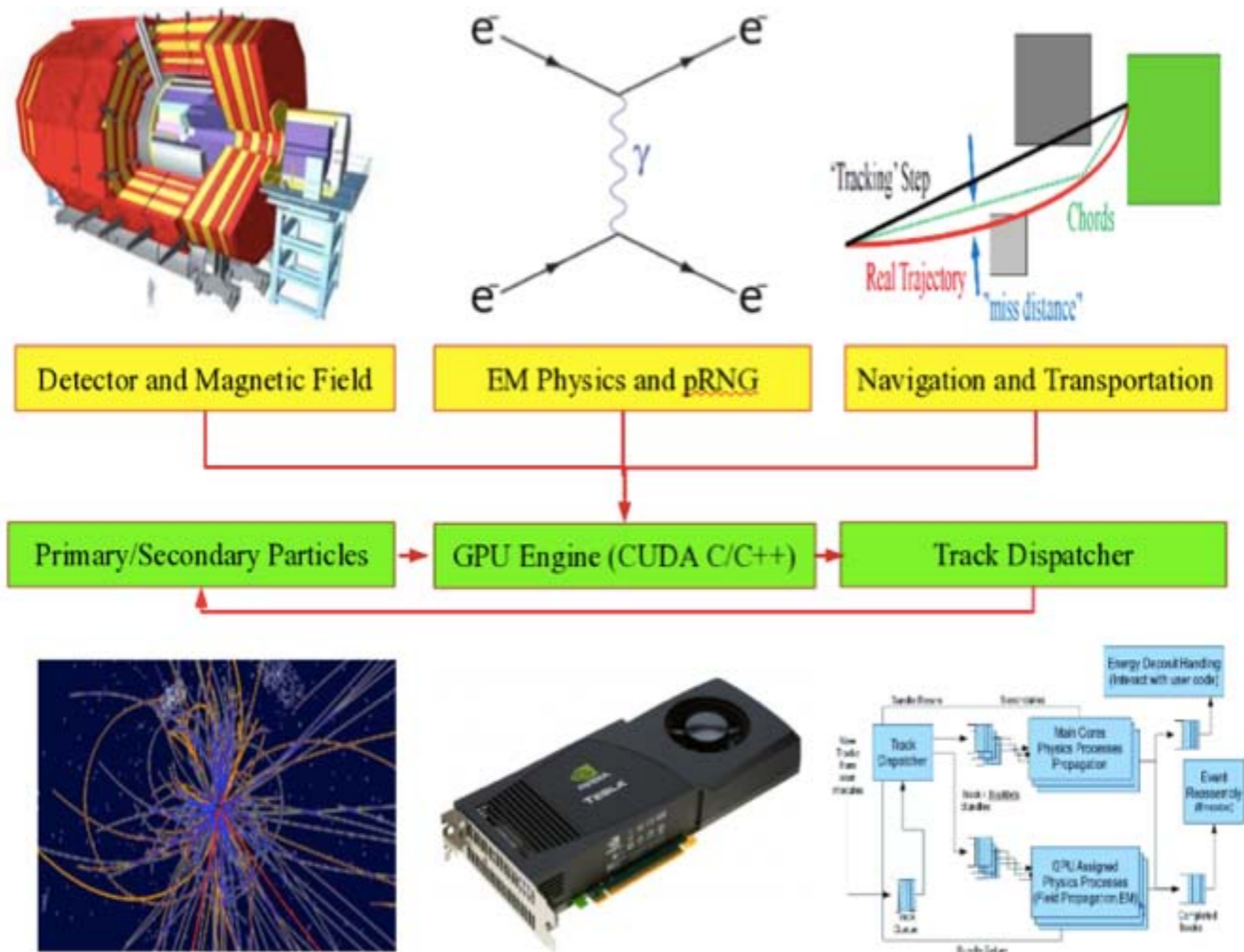
# Performance examples

- Performance of the inside method for a tube on an Intel CPU

- As you can see the performance is improved even for the scalar version

Unspecialized scalar
Specialized scalar (rot, trans)
Specialized vector (rot, trans)
Specialized scalar (rot, trans, type)
Specialized vector (rot, trans, type)
ROOT
Usolids

Inside method

# Particle Transport on the GPU

**Fermilab**

Detector and Magnetic Field | EM Physics and pRNG | Navigation and Transportation

Primary/Secondary Particles → GPU Engine (CUDA C/C++) → Track Dispatcher

**Fermilab**

# Physics Validation of GPU Physics

- Compare simulated physics outputs
  RED                    BLUE



EM Particle Transportation on the GPU - S.Y. Jun @Annual Corcurrent Forum Meeting          April 2, 2014

**Fermilab**

# Performance Evaluation

- Hardware (host + device)

| | Host (CPU) | Device (GPU) |
|---|---|---|
| **M2090** | AMD Opteron™ 6134 32 cores @ 2.4 GHz | Nvidia M2090 (Fermi) 512 cores @ 1.3 GHz |
| **K20** | Intel® Xeon® E5-2620 24 cores @ 2.0 GHz | NVidia K20 (Kepler) 2496 cores @ 0.7GHz |

- Performance measurement

**‡‡ Fermilab**

# Performance: Realistic Simulation

- A simple calorimeter (CMS Ecal) with the CMS b-field map
- Tracking for one step: split kernels (GPIL+sorting+DoIt)

|  | CPU [ms] | GPU [ms] | CPU/GPU |
|---|---|---|---|
| AMD+M2090 | 748 | 37.8 (62.6)* | 19.8 (11.9)* |
| Intel®+K20M | 571 | 30.4 (81.9)* | 18.7 (7.0)* |

- Performance by each kernel (% of the total application time)

 **춘 Fermilab**
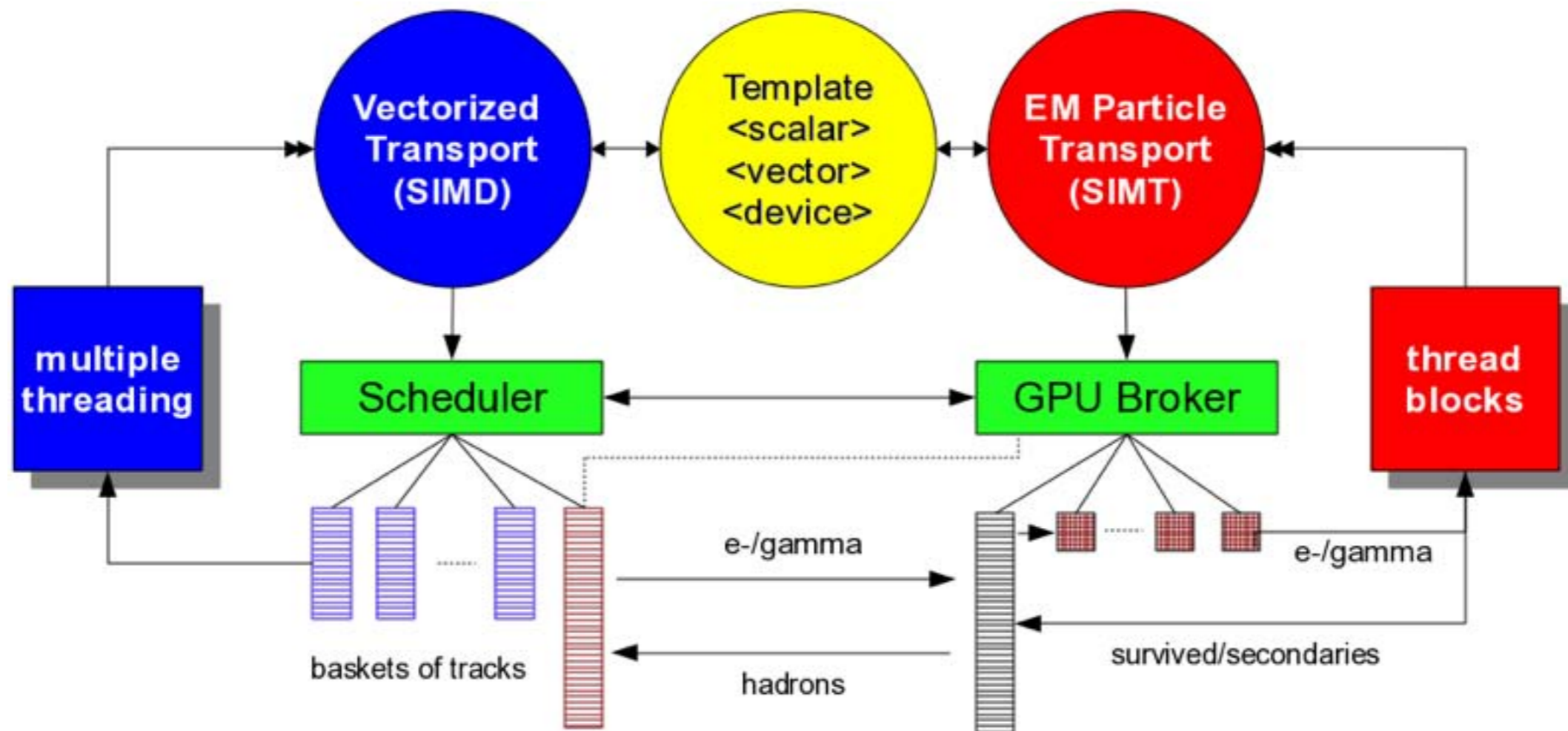
# Performance Issues and Considerations

- Observed issues (by the Nvidia profiler)


- Considerations

**Fermilab**

# Challenges and New Strategies

- HEP detector simulation (ex. Geant4) is a giant

- Coprocessor architectures

- Top-down approach

EM Particle Transportation on the GPU - S.Y. Jun @Annual Corcurrent Forum Meeting                    April 2, 2014

**Fermilab**

# Incorporating into the Vector prototype

- The vector prototype started at CERN (talk by F. Carminati)

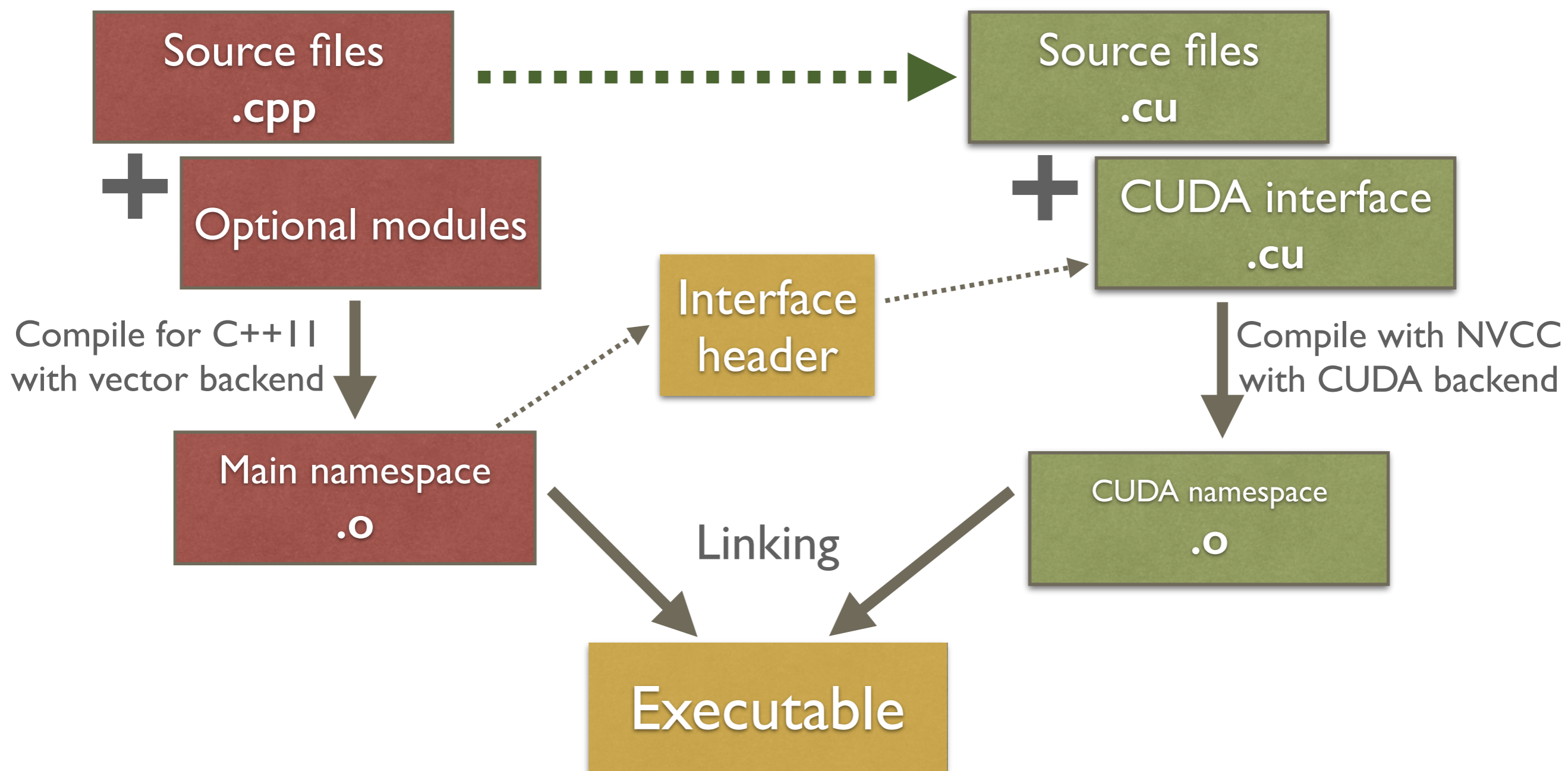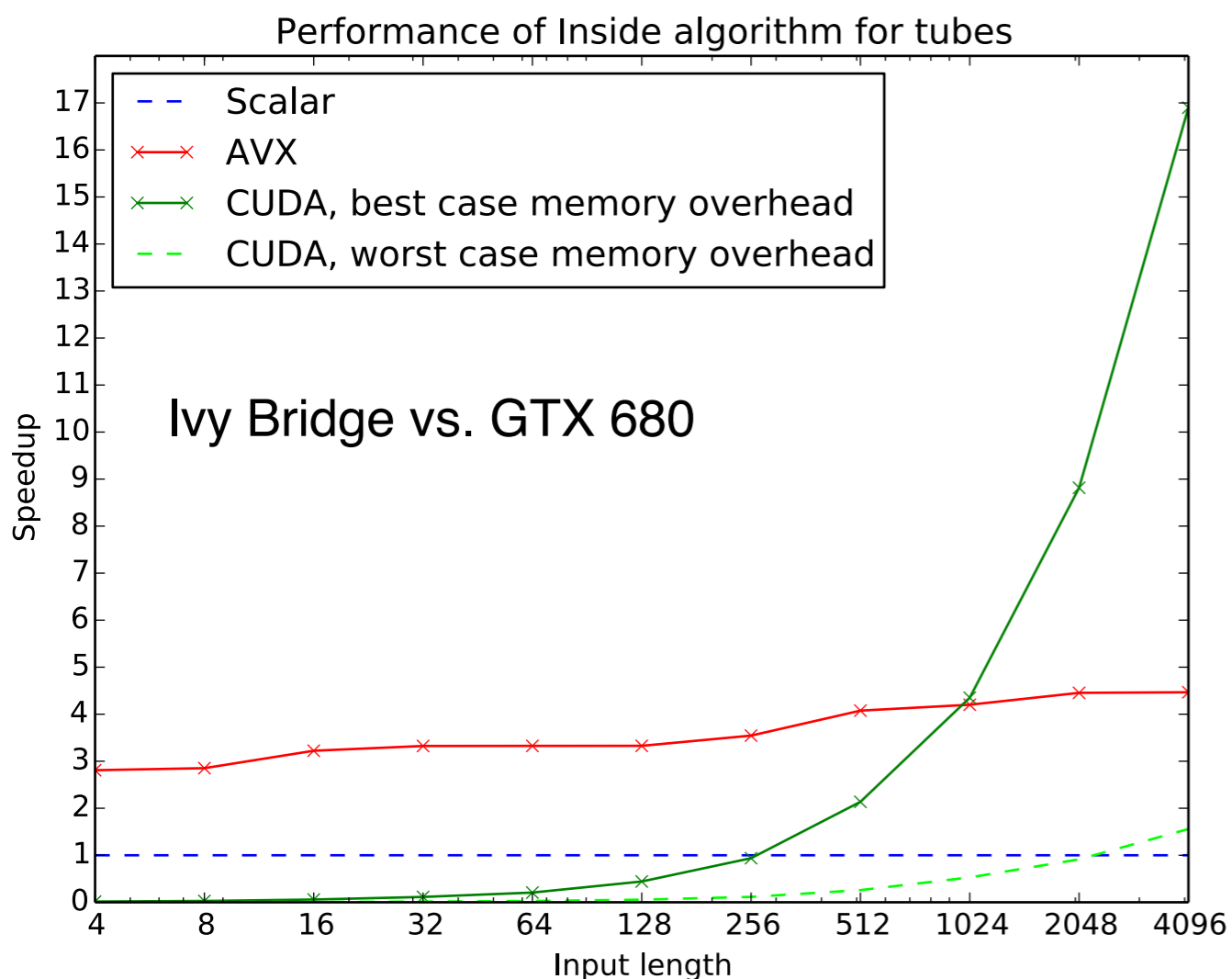- Integration to the vector prototype (GPU broker)

# CUDA in GeantV geometry

- Enable dispatching to CUDA as a co-processor
- Run separate to scalar and vectorized code
- Use same codebase
- Achieved by abstracted, templated algorithms

# Separate compilation of backends

# Shape measurements



Performance of Inside algorithm for tubes

Ivy Bridge vs. GTX 680

- Scalar, vector and CUDA code templated from <u>same</u> abstracted algorithm
- Dispatch on the fly to optimal processor
- Typical GPU scaling; high minimum input threshold

# Where are we now?

- Scheduler
  - The new version, hopefully improved of the scheduler has been committed and we are testing it
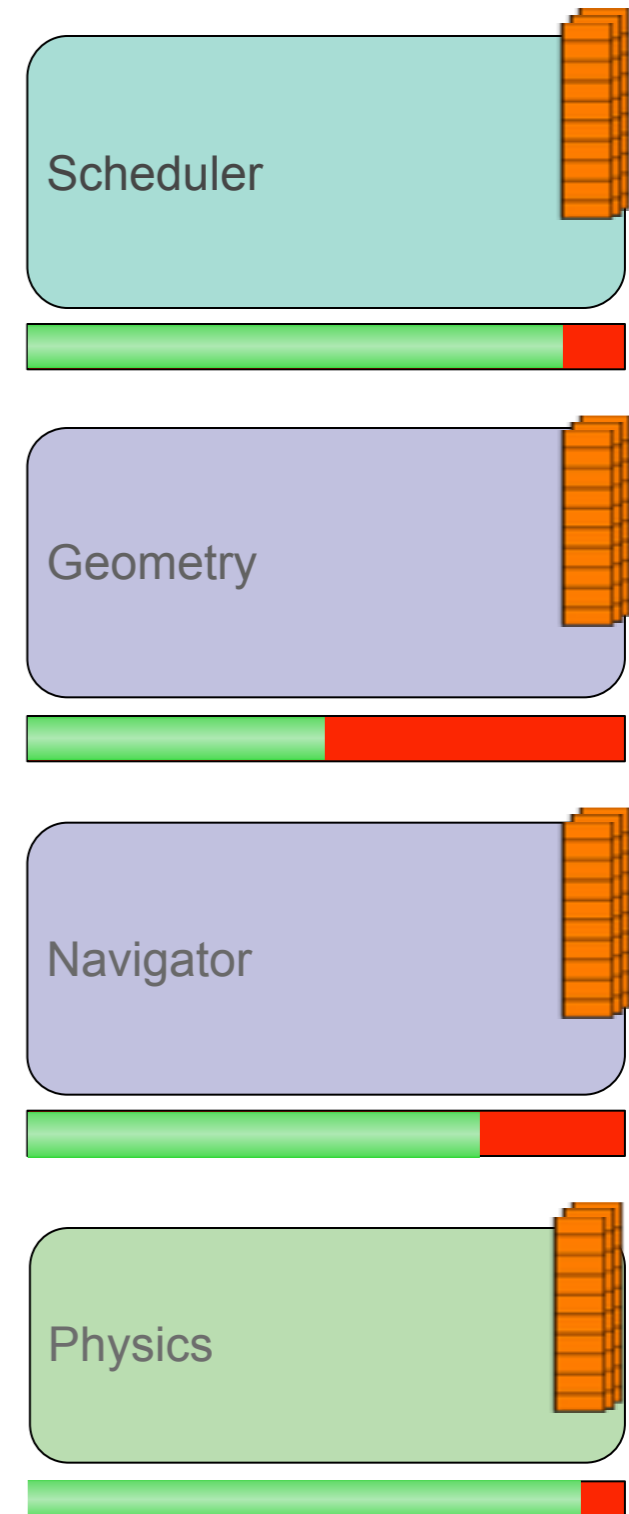- Solids
  - The proof or principle that we can achieve large speedups (3-5+) is there, however a lot of work lays ahead
- Navigator
  - "Percolating" vectors through the navigator is a difficult business. We have a simplified navigator that achieves that, but more work is needed here
- Physics
  - Can generate x-secs and final states and sample them, but there are still many points to be clarified with Geant4 experts

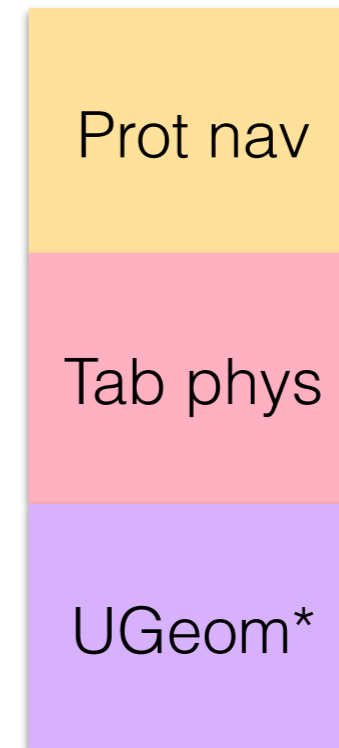Scheduler

Geometry

Navigator

Physics

# Testing - benchmarking

Geant4

Prototype

| G4 nav |
| :---: |
| G4 phys |
| G4 geom |

| Prot nav |
| :---: |
| Tab phys |
| UGeom* |

With small, large & venti geometries

*UGeom == USolid + navigation

# Testing - benchmarking

Geant4

Prototype

| G4 nav | G4 nav |
|--------|--------|
| G4 phys | Tab phys |
| G4 geom | G4 geom |

| Prot nav |
|----------|
| Tab phys |
| UGeom* |

With small, large & venti geometries

29 *UGeom == USolid + navigation

# Testing - benchmarking

## Geant4

| | | |
|---|---|---|
| G4 nav | G4 nav | G4 nav |
| G4 phys | Tab phys | G4 phys / V phys |
| G4 geom | G4 geom | G4 geom |

## Prototype

| |
|---|
| Prot nav |
| Tab phys |
| UGeom* |

With small, large & venti geometries

*UGeom == USolid + navigation

# Testing - benchmarking

Geant4

Prototype

| G4 nav | G4 nav | G4 nav | G4 nav | Prot nav |

| G4 phys | Tab phys | G4 phys / V phys | Tab phys / V phys | Tab phys |

| G4 geom | G4 geom | G4 geom | G4 geom | UGeom* |

With small, large & venti geometries

29 *UGeom == USolid + navigation

# Testing - benchmarking

## Geant4

| G4 nav | G4 nav | G4 nav | G4 nav |
|--------|--------|--------|--------|
| G4 phys | Tab phys | G4 phys / V phys | Tab phys / V phys |
| G4 geom | G4 geom | G4 geom | G4 geom |

## Prototype

| Prot nav | Prot nav |
|----------|----------|
| Tab phys | Tab phys / V phys |
| UGeom* | UGeom* |

## With small, large & venti geometries

29 *UGeom == USolid + navigation

# Conclusion

- Work on a full prototype is progressing
- The first performance figures are very encouraging
- The template specialisation technique seems to provide a reasonable model for portable HPC
- We hope to have a demonstrator by the end of the year