

An Efficient SIMD Implementation of Pseudo-Verlet Lists for Neighbour Interactions in Particle-Based Codes

James S. Willis, Matthieu Schaller, Pedro Gonnet, Richard G. Bower &
Peter W. Draper
Durham University, ICC
September 14th 2017

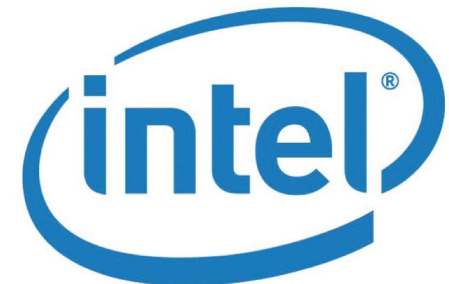
Team

This work is a collaboration between two departments at Durham University (UK):

- The Institute for Computational Cosmology,
- The School of Engineering and Computing Sciences,

with contributions from the astronomy group at the University of St. Andrews, University of Dublin, ETH Lausanne, the DiRAC software team and the Hartree Centre.

- This research is fully funded by an IPCC.

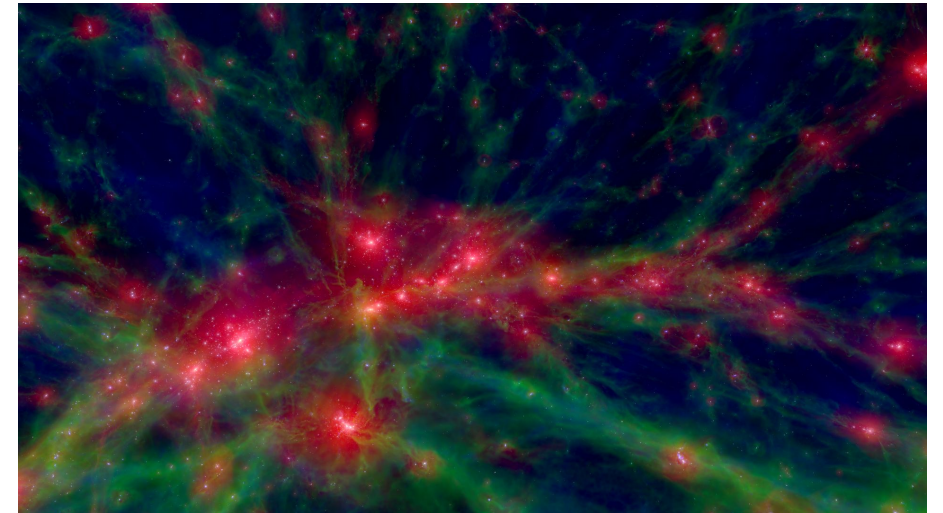


Overview

- Problem to solve
- Solution
 - Pseudo Verlet list
 - Particle sorting
- SIMD vectorisation strategy for particle based codes (MD, SPH, etc.)
 - Particle caches AoS to SoA
- Strategy applied to SWIFT
- Performance results
- Conclusions

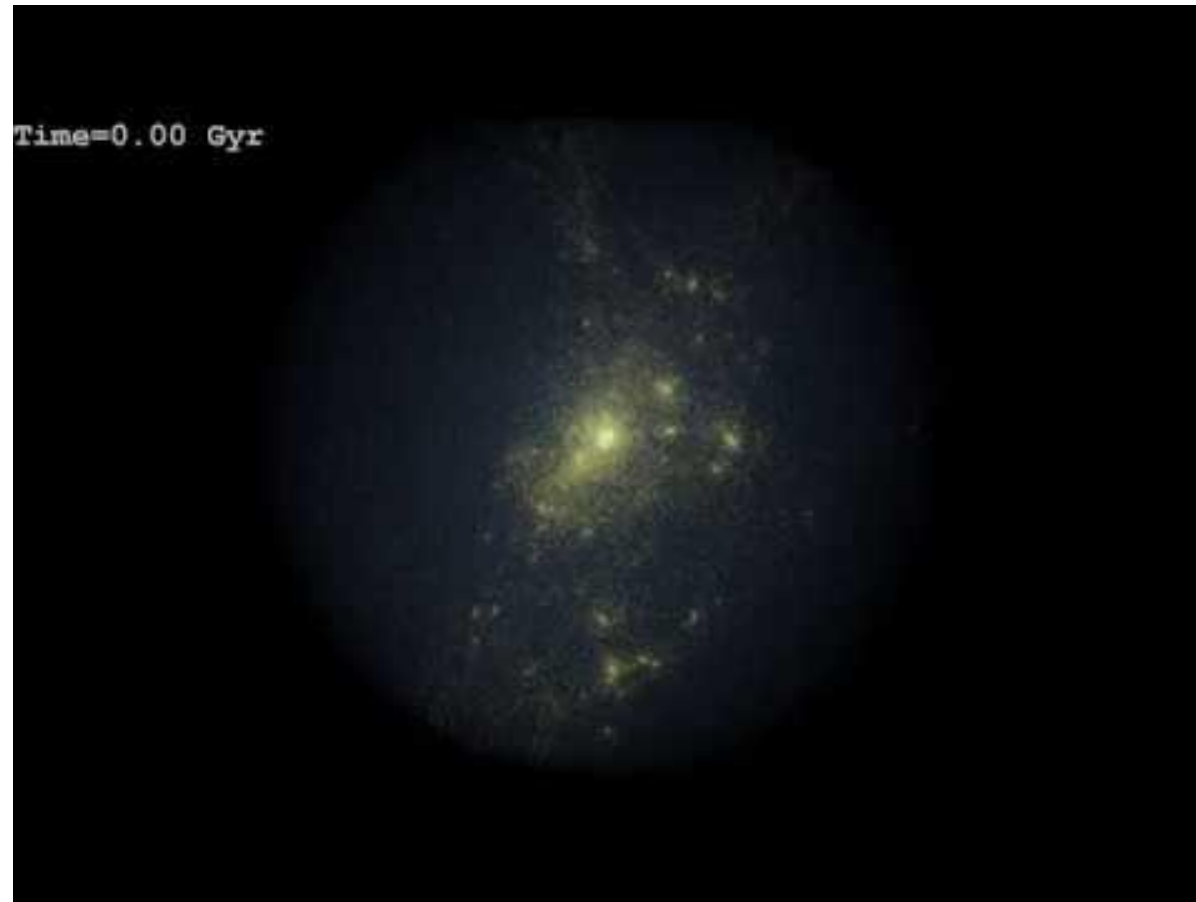
Motivation

- Create simulations of the formation and evolution of the Universe
- Update 10^9 particles using hydrodynamical and gravitational forces
- Simulate physical processes:
 - Cooling and heating of the gas due to the presence of stars and other emission
 - Formation of stars in cold and dense regions
 - Explosion of supernovae with injection of their energy in the surrounding gas
 - Formation of supermassive black holes



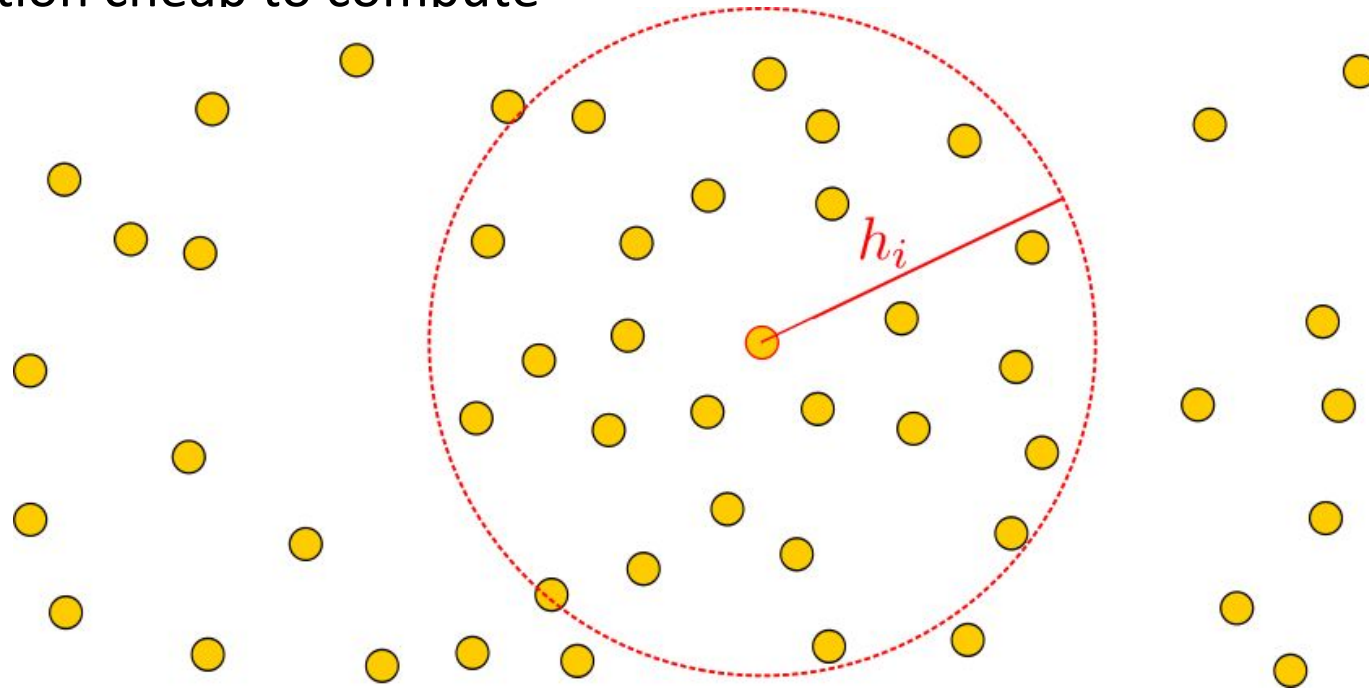
Motivation

- Dwarf galaxy simulation using SWIFT



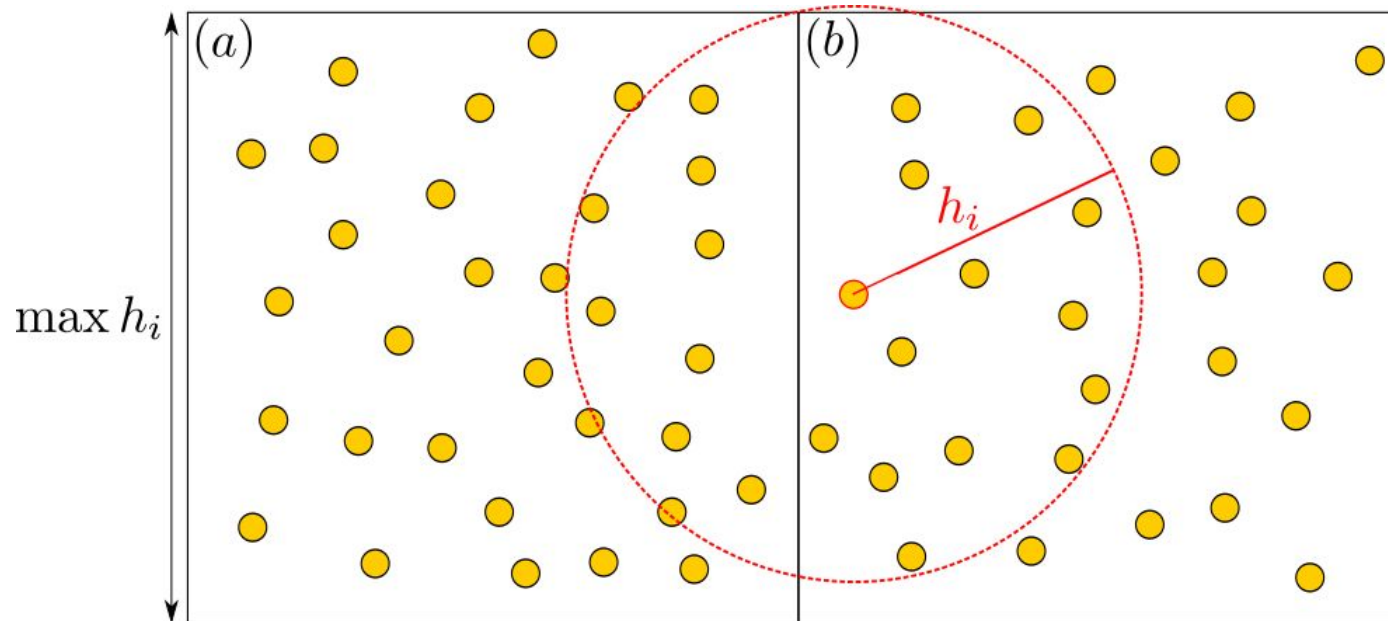
Problem

- We update each particle using SPH (Smoothed-Particle Hydrodynamics)
- Each particle interacts with its neighbours that are within a cut-off radius, h
- h varies depending on the particle density of the region
- Interaction cheap to compute



Problem

- The particles are divided up into cells of edge h_{\max} , where h_{\max} is the maximum particle cut-off radius in the simulation
- Computing the interactions of particles in two neighbouring cells would require a lot of unnecessary distance calculations
- The majority of particles will not be within range of each other



Naive Solution

Brute Force

- Perform a double for loop over all particles
- Interact particles that are within range of each other, $r < h$
- Trivial to vectorise

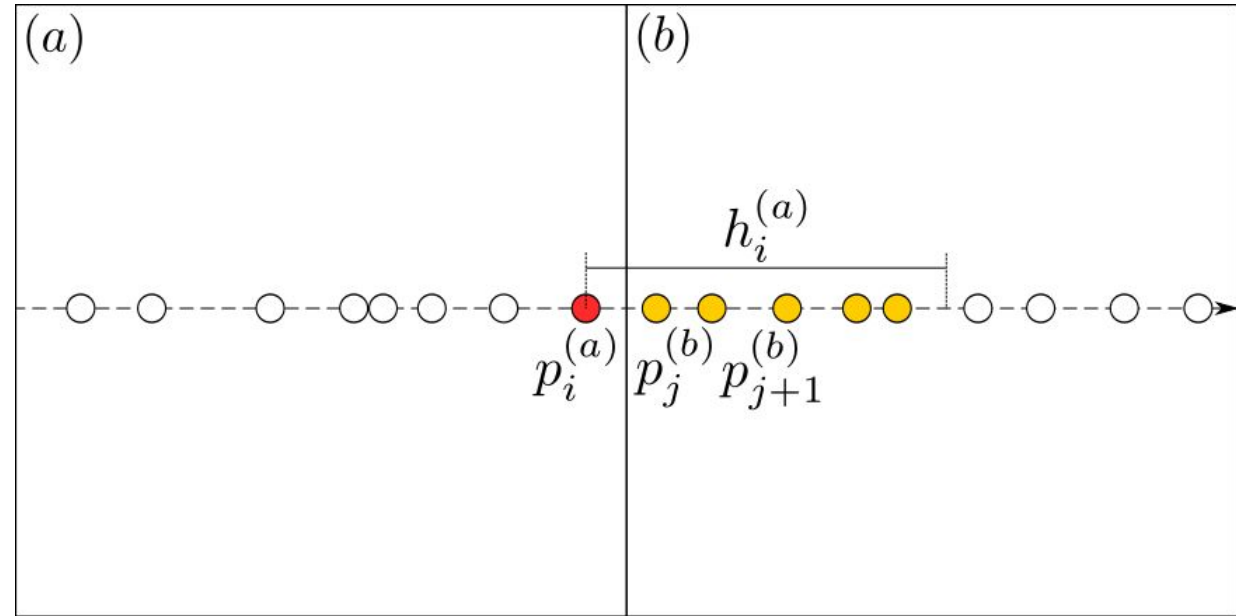
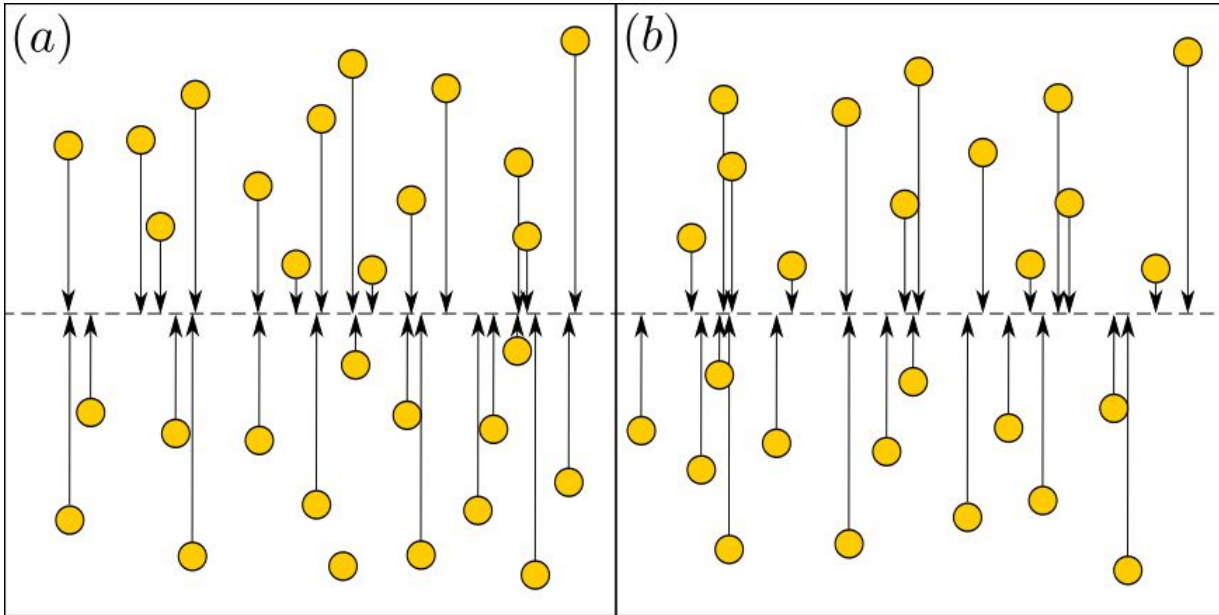
```
for (int i = 0; i < count_a; i++) {  
  
    const float h = parts_a[i].h;  
  
    for (int j = 0; j < count_b; j++) {  
  
        r = particle_dist(parts_a[i], parts_b[j]);  
  
        if (r < h) {  
            // Compute interaction.  
        }  
    }  
}
```


Smart Solution

Particle Sorting

- Place particles into a pseudo-Verlet list:
 - Project particles onto the axis joining the center of the two cells
 - Sort the particles on the axis based upon their position on the axis
 - Sorting performed using a merge sort
 - Only occurs when the particles have moved by a certain distance
- Reduces the number of candidates
- These particles are still tested so that they are within the 3D distance

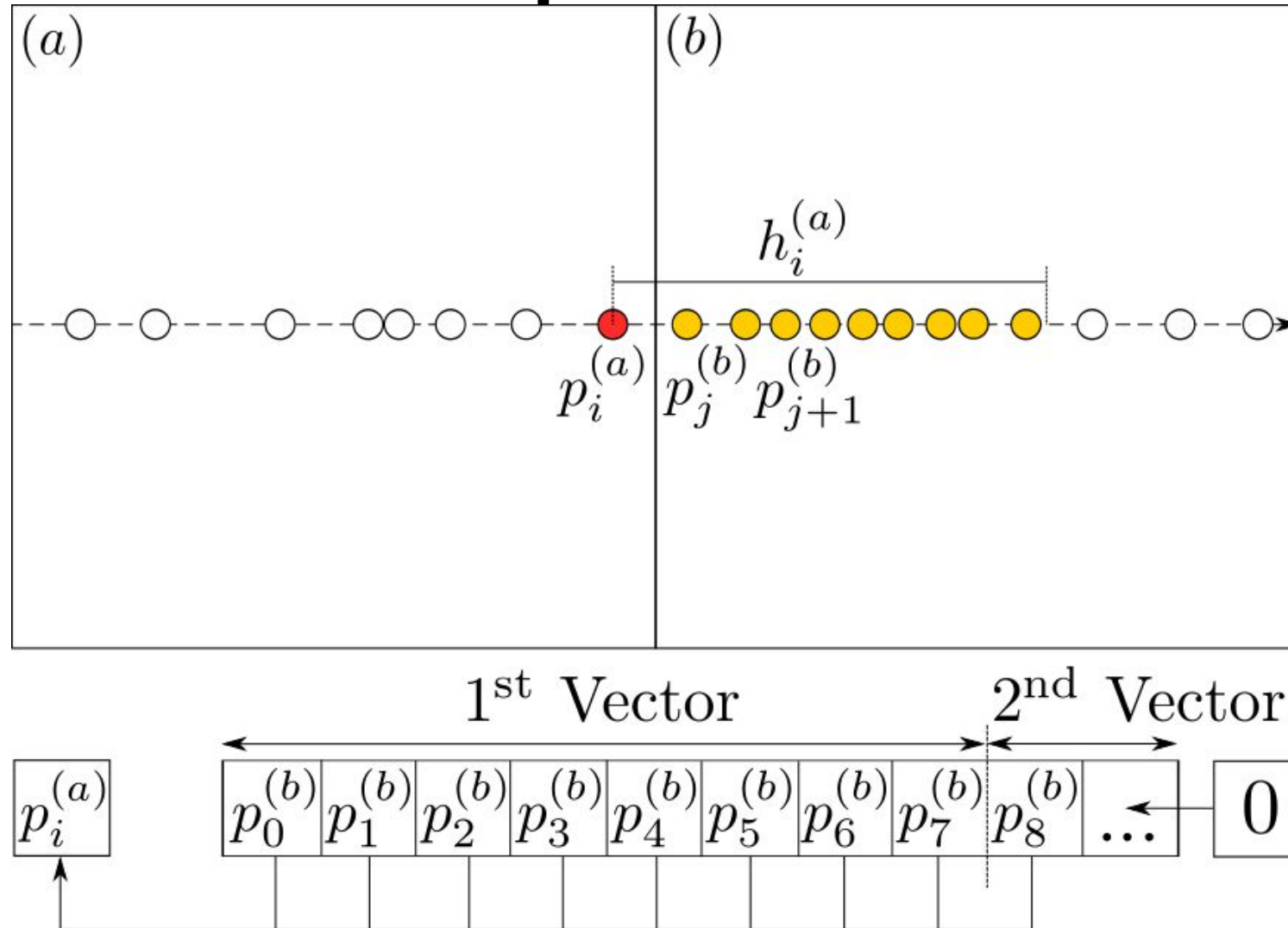
Smart Solution



Smart Solution

```
for (int i = 0; i < count_a; i++) {  
  
    const float h = parts_a[i].h;  
  
    for (int j = 0; j < count_b && dist_b[j] < dist_a[i] + h; j++) {  
  
        r = particle_dist(parts_a[index_a[i]], parts_b[index_b[j]]);  
  
        if (r < h) {  
            // Compute interaction.  
        }  
    }  
}
```

SIMD Optimisations

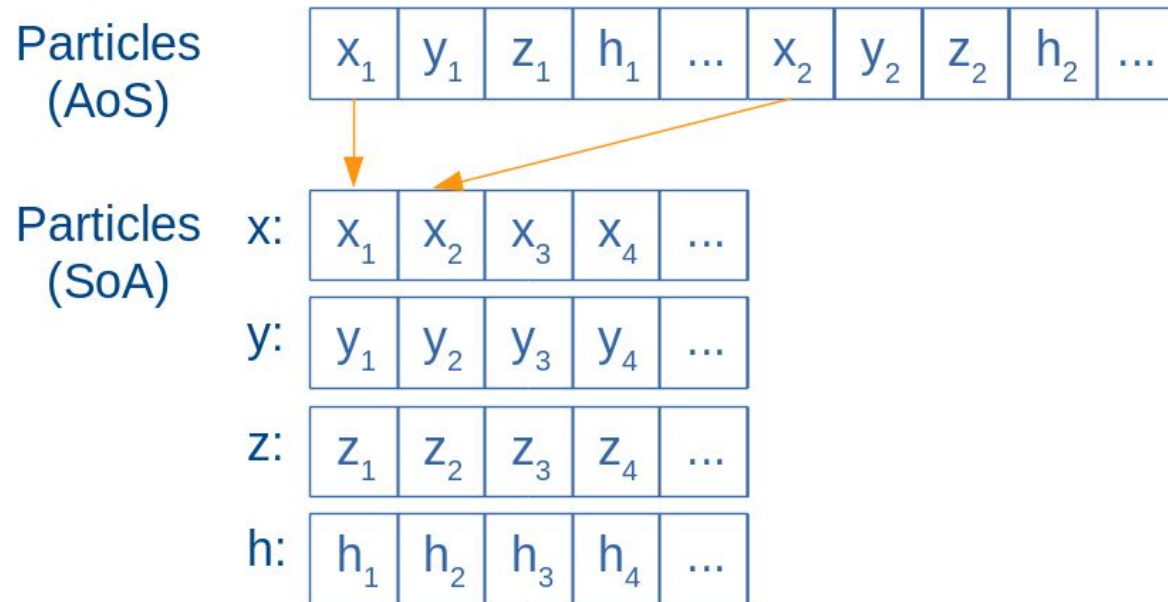


SIMD Optimisations

- Use local particle cache (AoS -> SoA)
- Only read particles that interact into cache.
- Calculate all interactions on a particle and store results in a set of intermediate vectors
- Perform horizontal add on intermediate vectors and update the particles with the result
- Pad caches to prevent remainders and mask out the result

Local Particle Cache

- The particles are stored in a global array of structs (AoS)
- Causes strided memory access when vectors are loaded
- We can improve performance by placing the required particle properties into a structure of arrays (SoA)

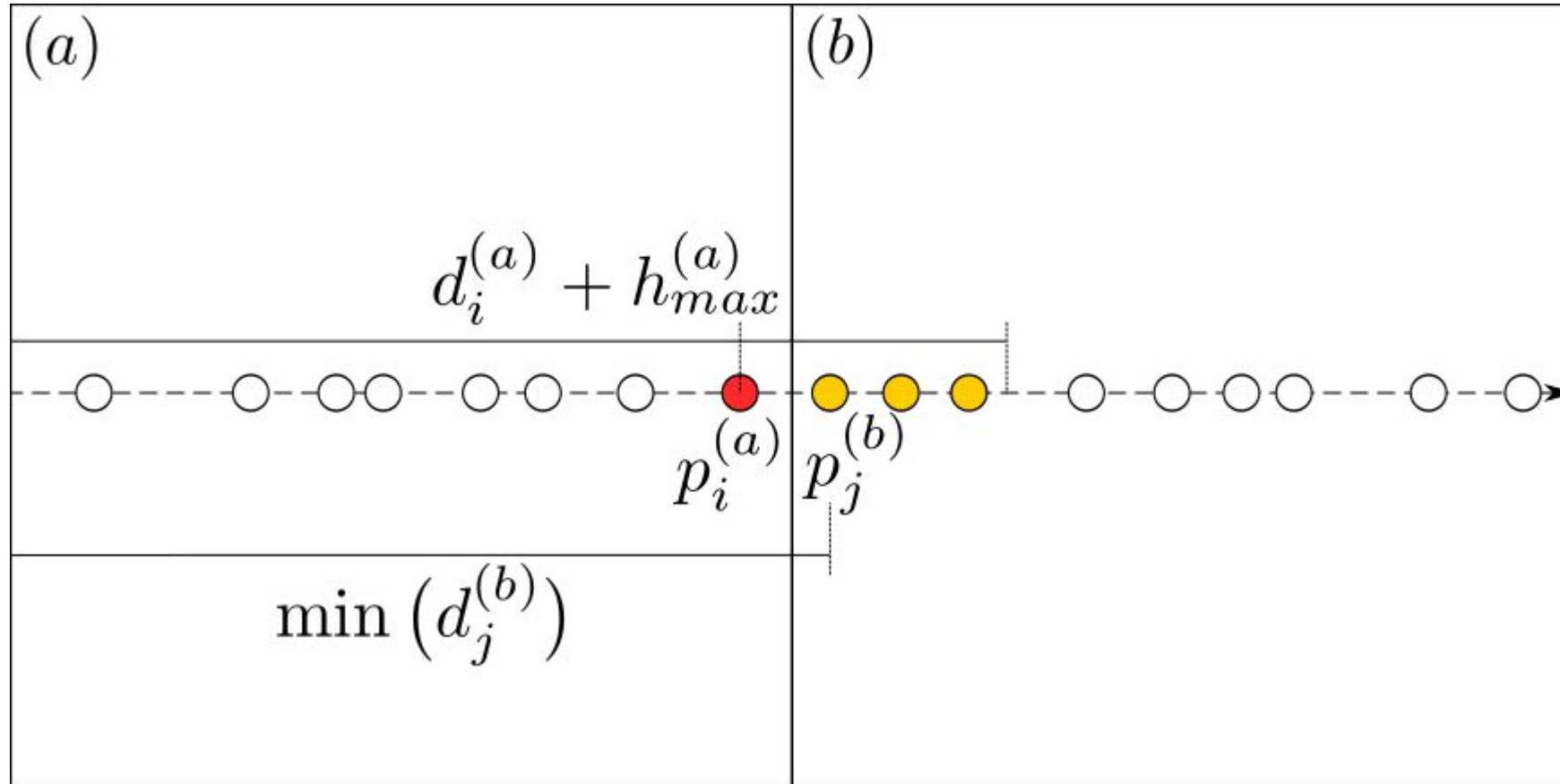


Populate Local Cache

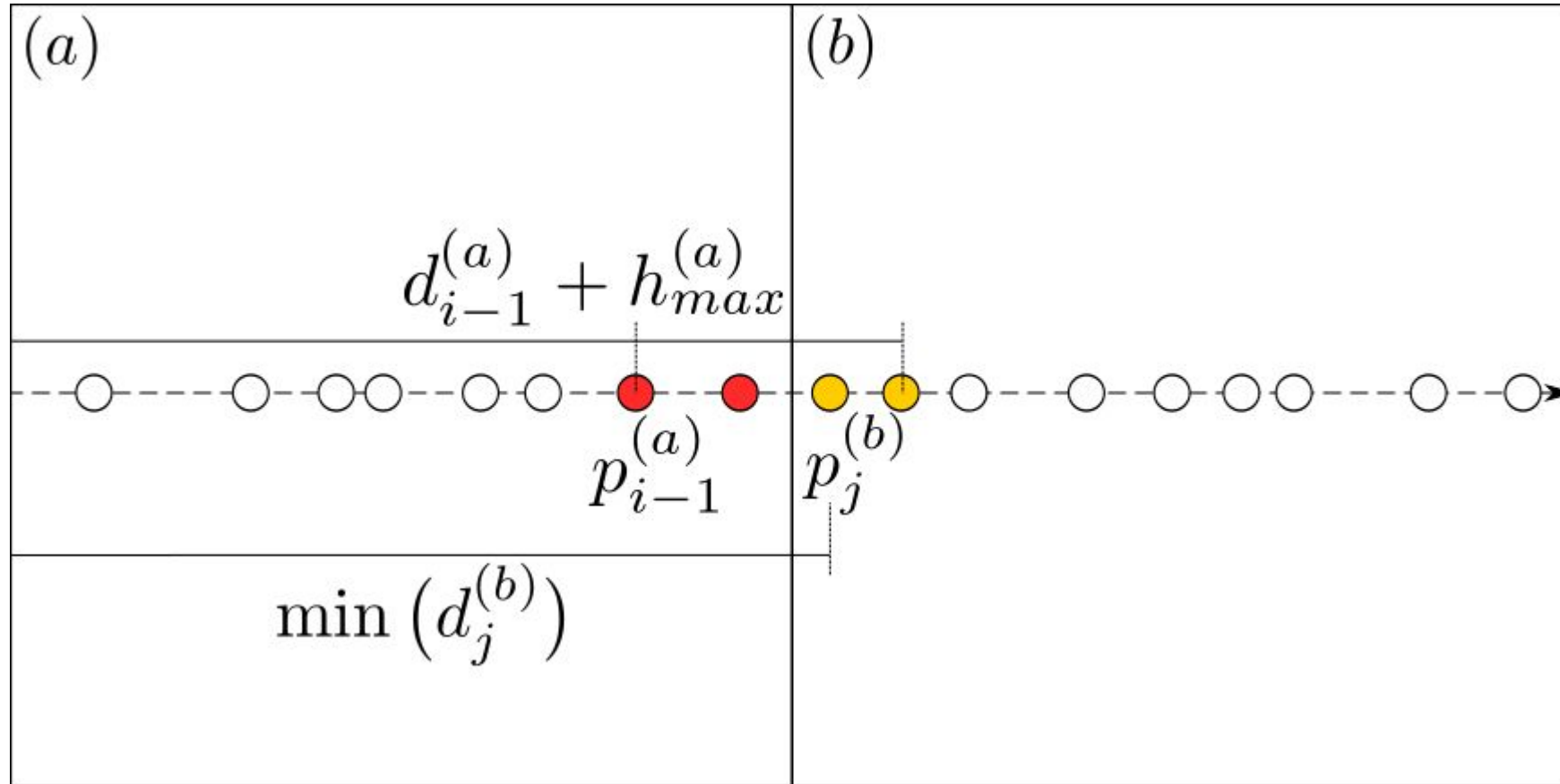
- In a uniform distribution of particles each cell pair orientation has a different number of interactions
- There are three cell pair orientations: *corner*, *edge* and *face*
- Number of interactions: *corner* < *edge* < *face*

- We want to reduce the cache overhead by only reading particles that are within range of each other
- Allows *edge* interactions to speedup instead of slowing down

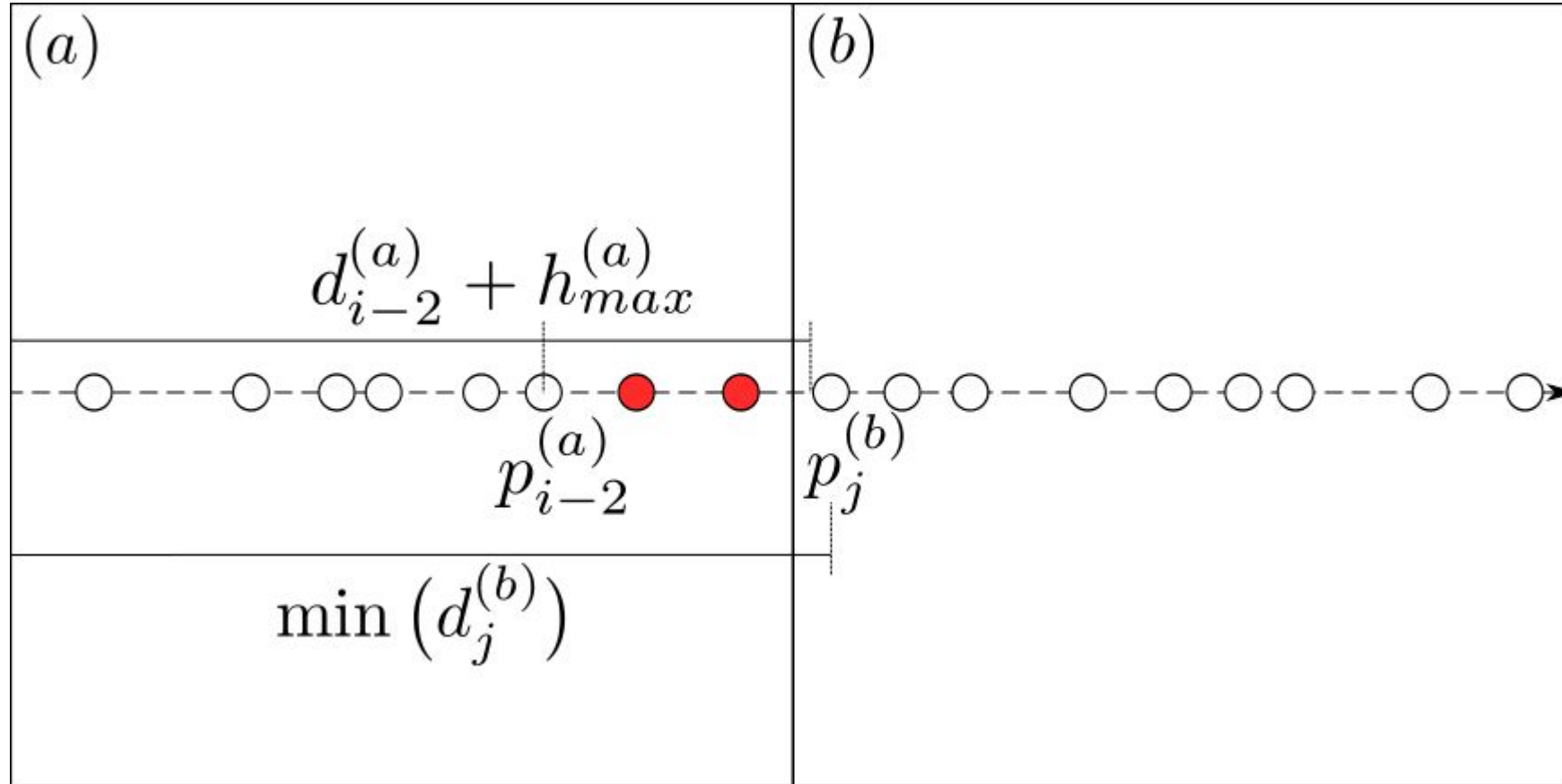
Populate Local Cache



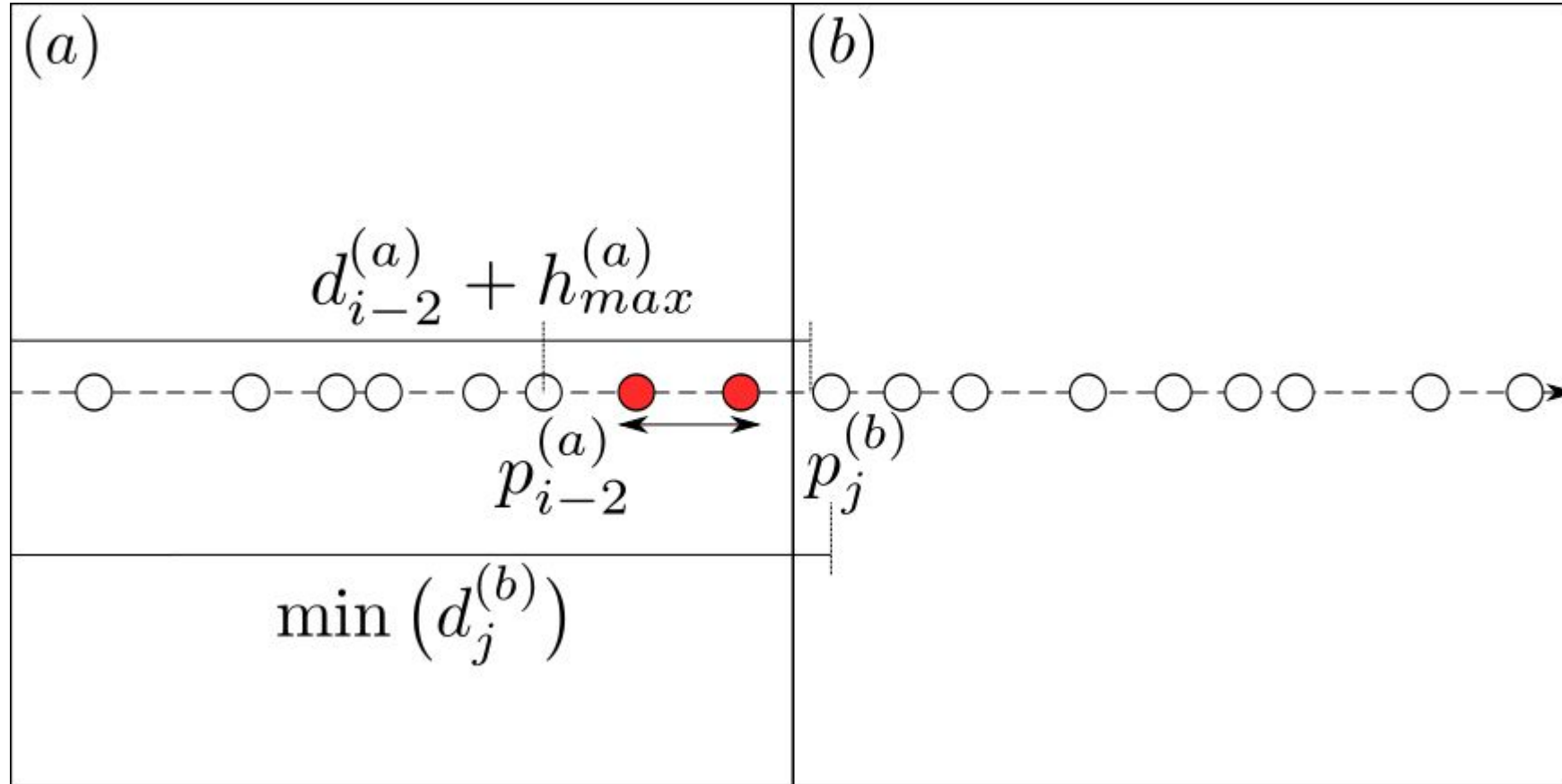
Populate Local Cache



Populate Local Cache



Populate Local Cache

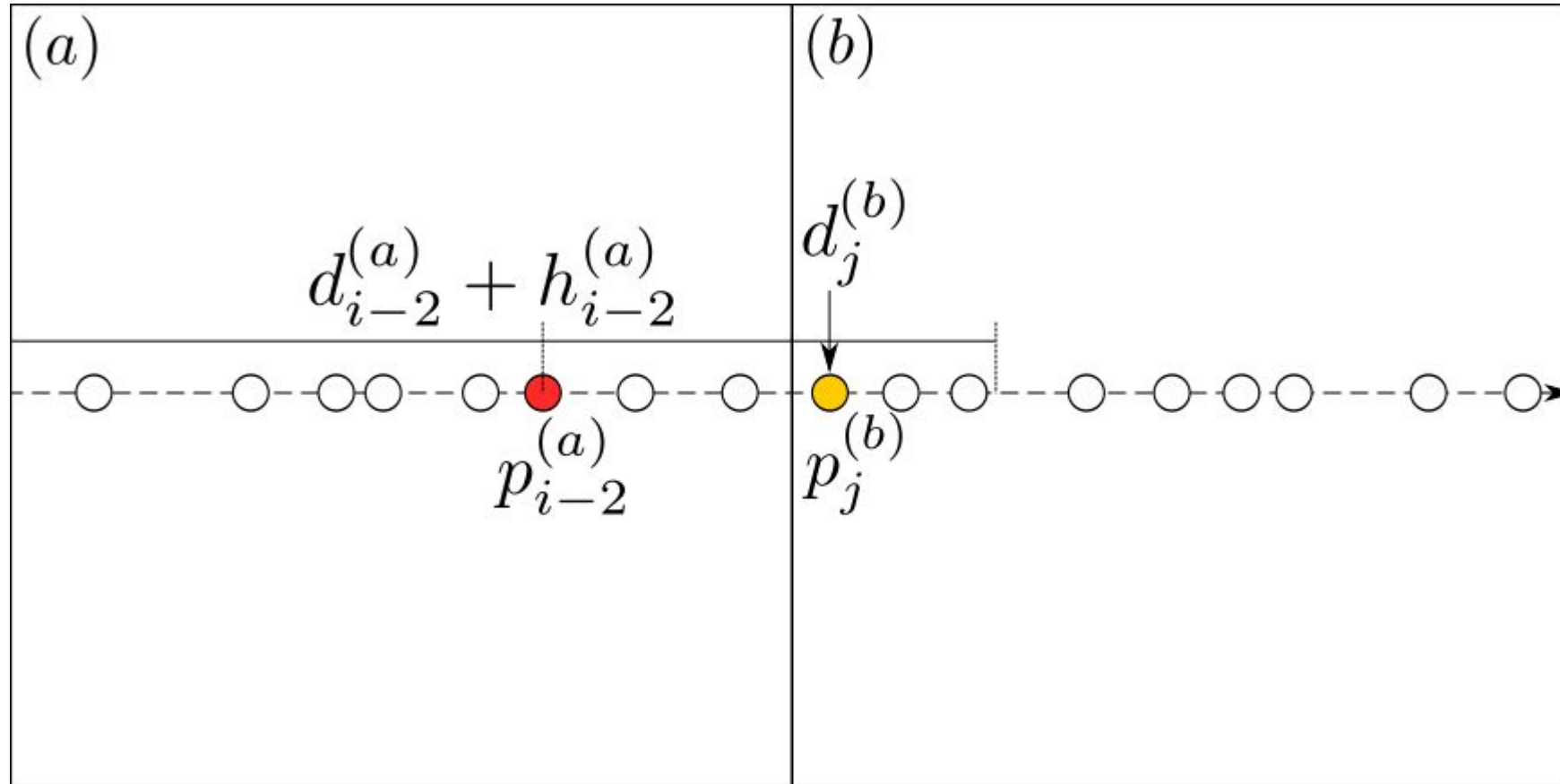


Limit Loop Bounds

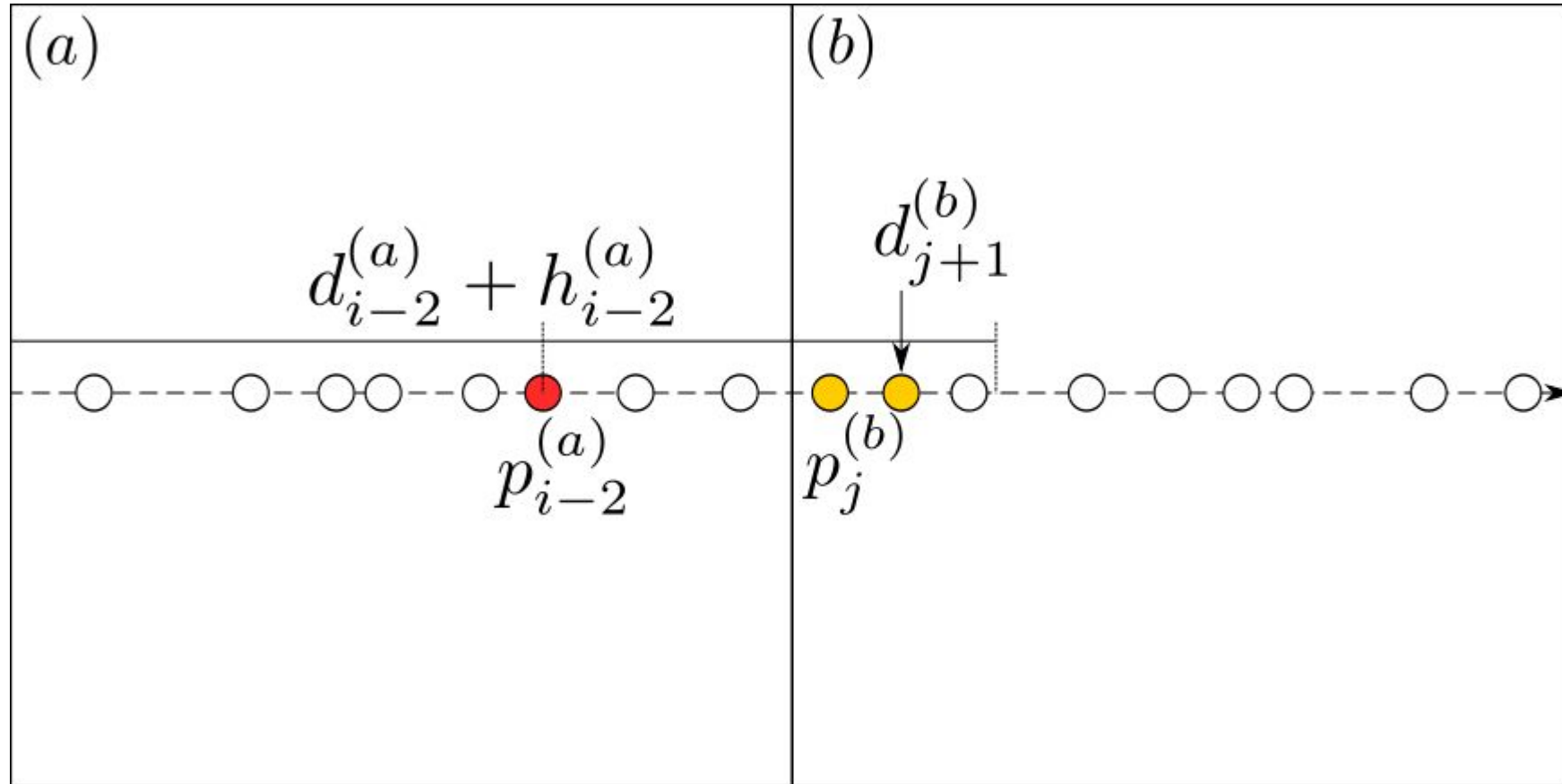
- For each particle we loop over every candidate in the neighbouring cell
- Even though most will be out of range as we move further from the interface between the two cells
- We want to reduce the number of distance calculations even further

- Form array of maximum indices into neighbouring cell
- Use array to limit the number of particles looped over in the neighbouring cell

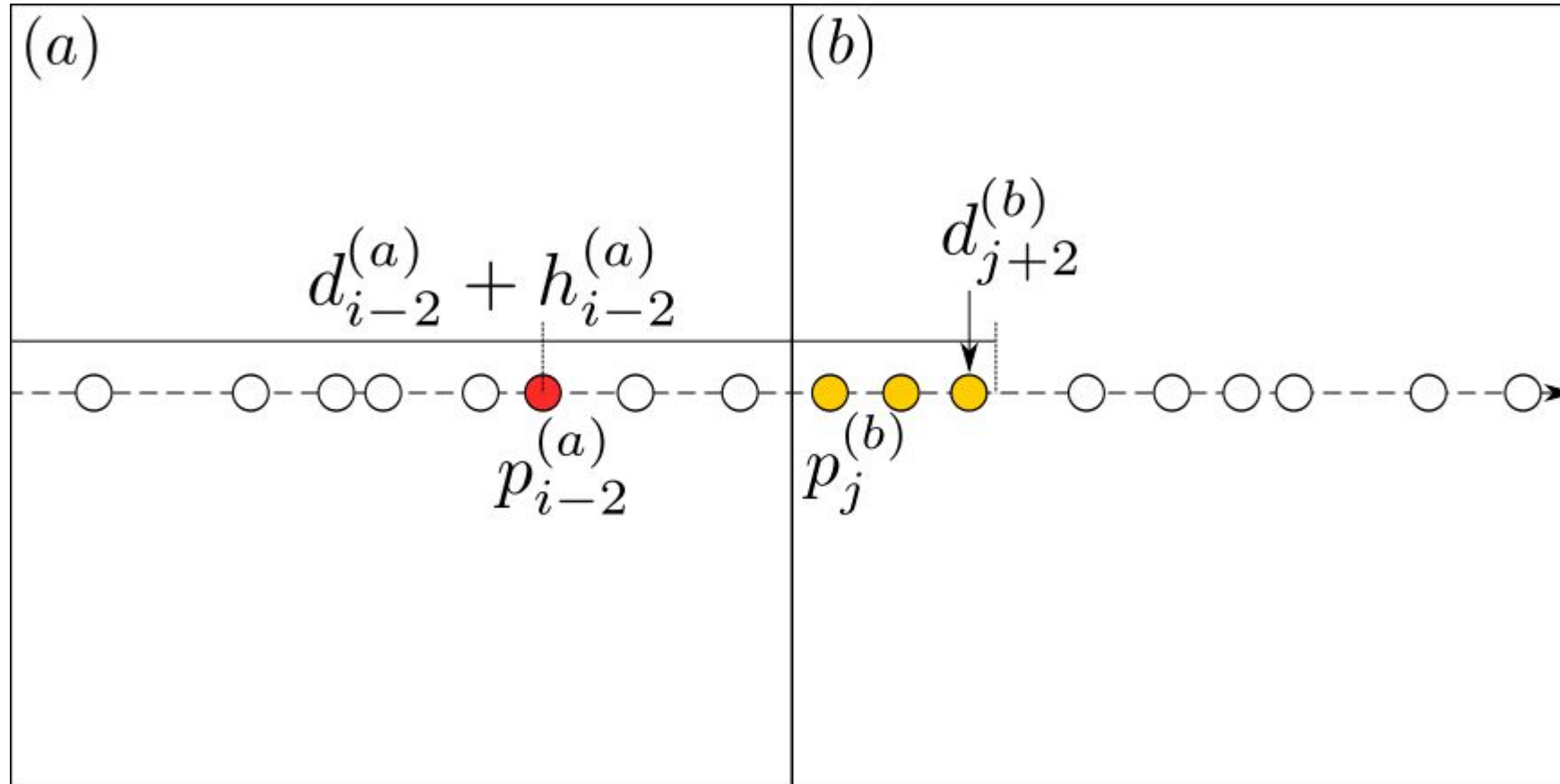
Limit Loop Bounds



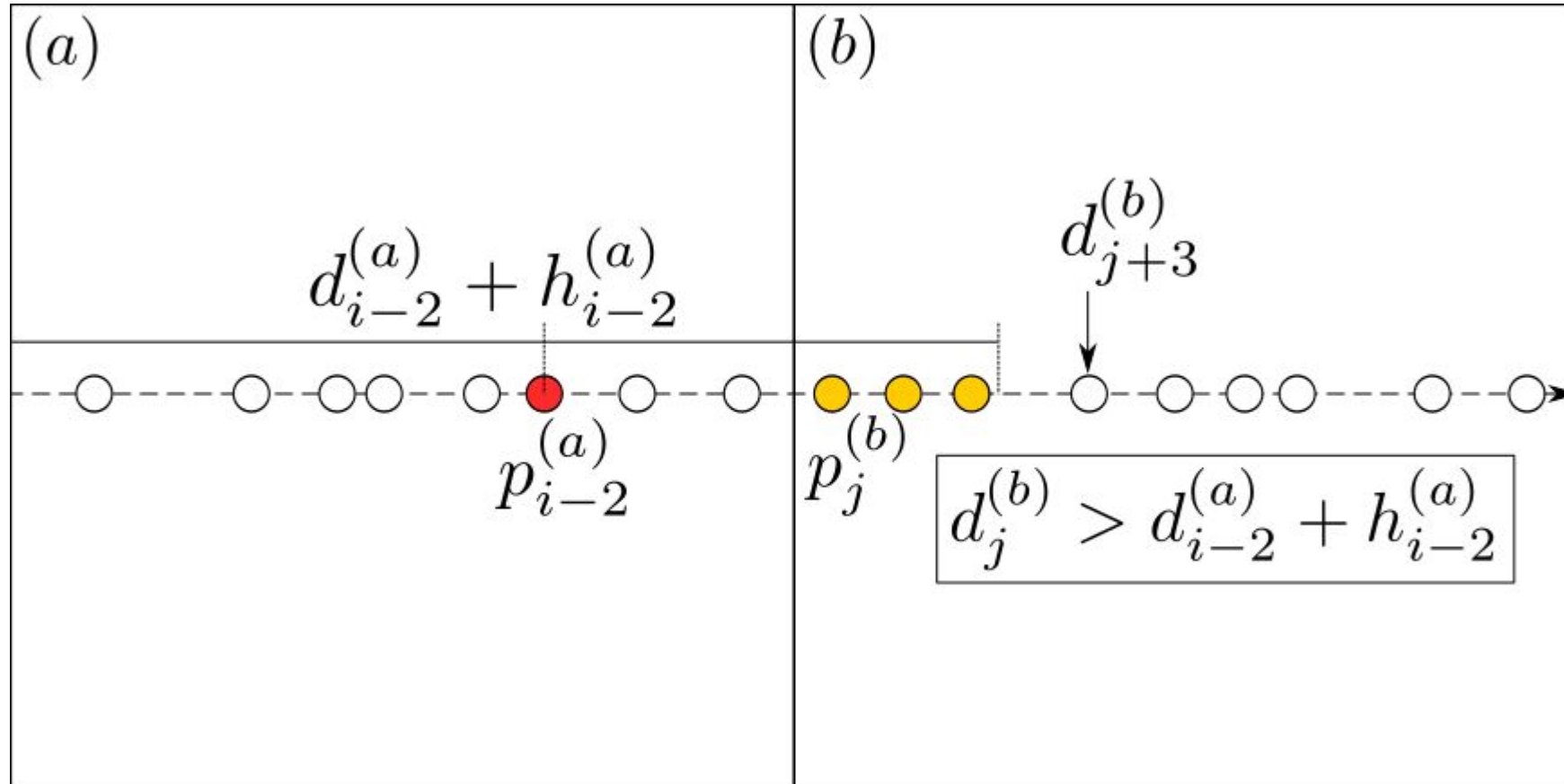
Limit Loop Bounds



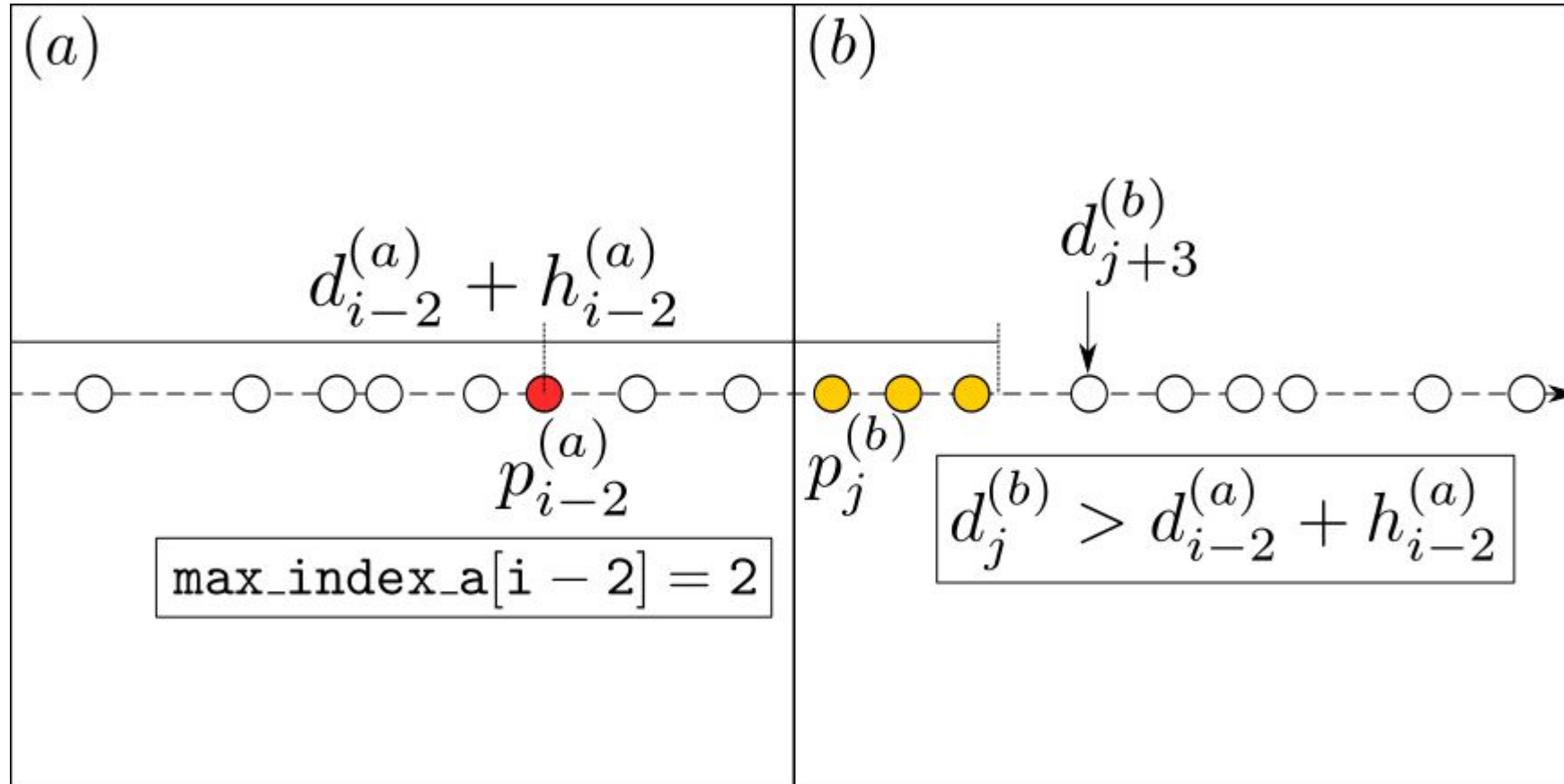
Limit Loop Bounds



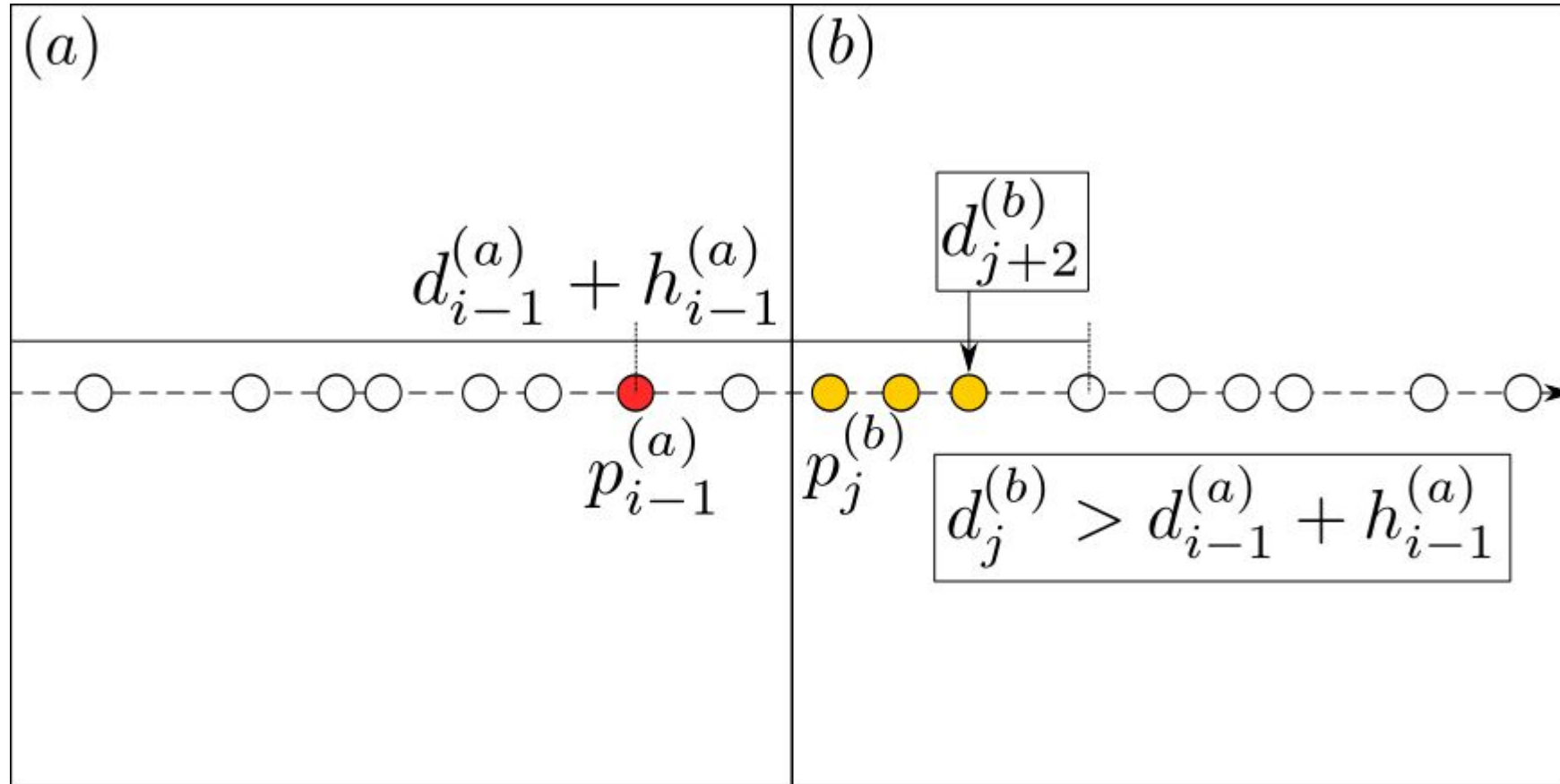
Limit Loop Bounds



Limit Loop Bounds



Limit Loop Bounds



Calculating Interactions

- Particle density interactions are calculated using:

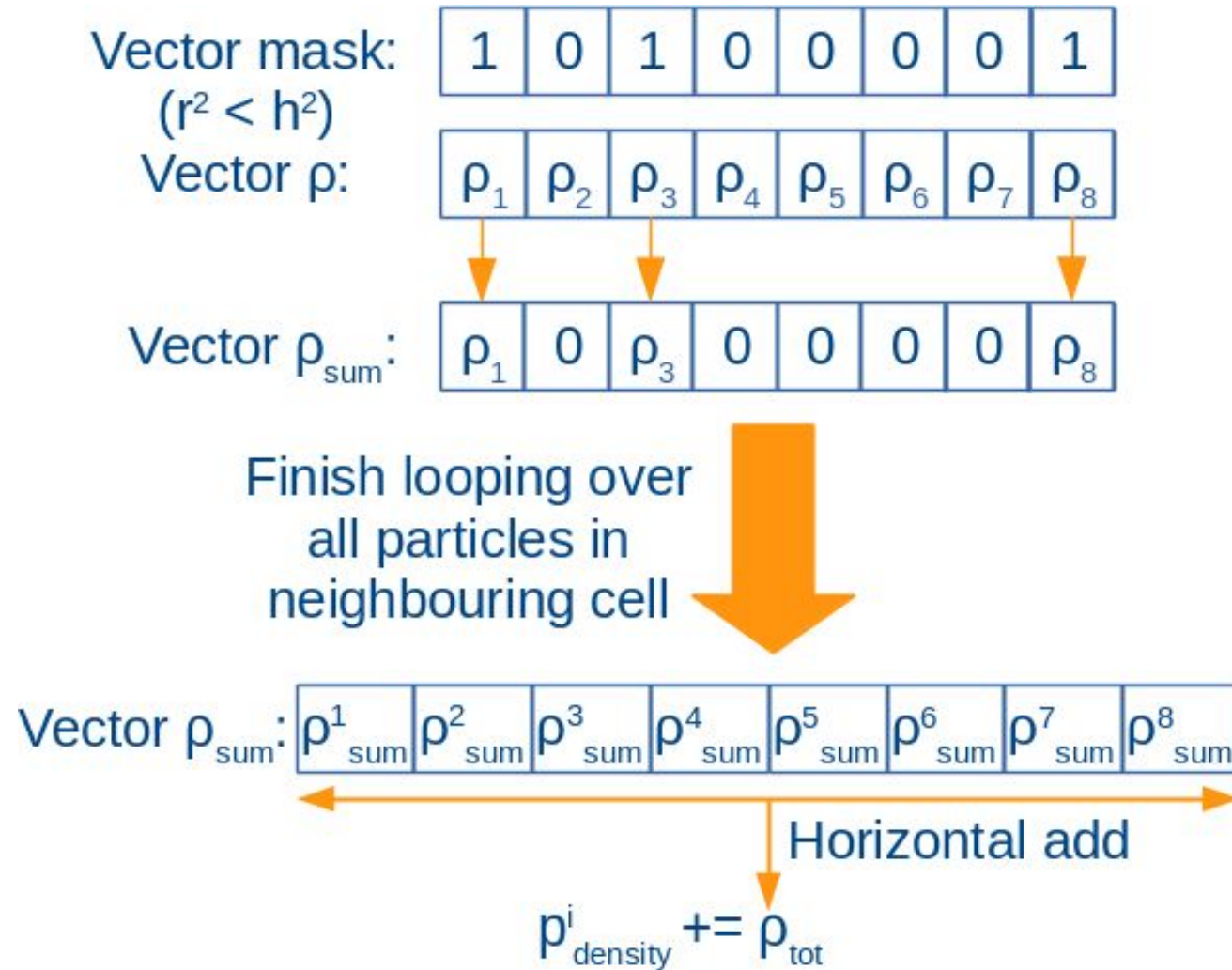
$$\rho(\mathbf{r}) = \sum_{b=1}^{N_{neigh}} m_b W(\mathbf{r} - \mathbf{r}_b, h)$$

- W is the weight function which is a low order polynomial

SIMD Implementation

- Use intermediate vectors to accumulate sum of particle updates in interaction function
- Perform horizontal add on these vectors and update the particles
- Decreases the amount of writes to memory

Calculating Interactions



Calculating Interactions

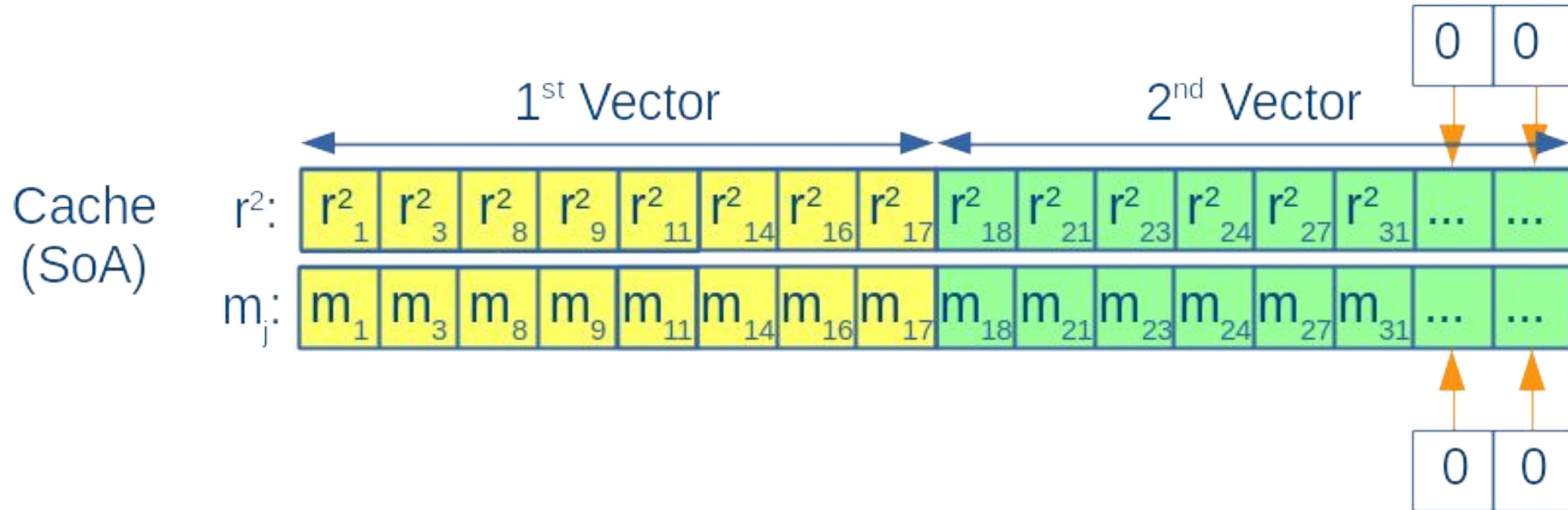
```
vector densitySum;
density = setzero();

for (int pjd = 0; pjd < icount; pjd+=VEC_SIZE) {
    INTERACT(&r2[pjd], &dx[pjd], &dy[pjd],
            &dz[pjd], &m[pjd], &v[pjd],
            &densitySum);
}

VEC_HADD(densitySum,pi); // _mm_hadd_ps
```

Padding Local Cache

- Pad vectors to remove the serial remainders and mask the result



Performance Results

- Vectorisation performance was measured using AVX, AVX2 and AVX-512 instruction sets on the following hardware:
 - Intel Xeon CPU E5-4640 @ 2.4GHz (Sandy Bridge)
 - Intel Xeon CPU E5-2650 @ 2.2GHz (Broadwell)
 - Intel Xeon Phi CPU 7210 @ 1.3GHz (Knights Landing)
 - Configured in *Flat-Quadrant* mode
- Intel Compiler 17.0.2

Performance Results

Cell-pair Orientation	Pseudo-Verlet List Scalar Time [ms]	Pseudo-Verlet List SIMD Time [ms]	Speed-up
Corner	0.00035	0.00070	0.49x
Edge	0.0052	0.0035	1.48x
Face	0.082	0.034	2.41x

Performance Results

Cell-pair Orientation	Pseudo-Verlet List Scalar Time [ms]	Pseudo-Verlet List SIMD Time [ms]	Speed-up
Corner	0.00035	0.00070	0.49x
Edge	0.0052	0.0035	1.48x
Face	0.082	0.034	2.41x

Performance Results

CFLAGS	Speed-up of raw particle interactions over serial version	Speed-up over serial pseudo-Verlet list
-O3 -xAVX -no-prec-sqrt -fp-model fast=2	5.66x	2.24x
-O3 -xCORE-AVX2 -no-prec-sqrt -fp-model fast=2	6.77x	2.43x
-O3 -xMIC-AVX512 -no-prec-sqrt -fp-model fast=2	21.30x	4.07x

Performance Results

CFLAGS	Speed-up of raw particle interactions over serial version	Speed-up over serial pseudo-Verlet list
-O3 -xAVX -no-prec-sqrt -fp-model fast=2	5.66x	2.24x
-O3 -xCORE-AVX2 -no-prec-sqrt -fp-model fast=2	6.77x	2.43x
-O3 -xMIC-AVX512 -no-prec-sqrt -fp-model fast=2	21.30x	4.07x

Performance Results

CFLAGS	Speed-up of raw particle interactions over serial version	Speed-up over serial pseudo-Verlet list
-O3 -xAVX -no-prec-sqrt -fp-model fast=2	5.66x	2.24x
-O3 -xCORE-AVX2 -no-prec-sqrt -fp-model fast=2	6.77x	2.43x
-O3 -xMIC-AVX512 -no-prec-sqrt -fp-model fast=2	21.30x	4.07x

Performance Results

- Compare our improvements against the naive implementation of the algorithm
- Assuming the naive algorithm could achieve maximum speed due to SIMD vectorisation

Performance Results

CFLAGS	Naive Solution Scalar Time [ms]	Naive Solution SIMD Time [ms] (Max SIMD Speedup)	Pseudo-Verlet List SIMD Time [ms]
-O3 -xAVX -no-prec-sqrt -fp-model fast=2	24.49	3.06	0.25
-O3 -xCORE-AVX2 -no-prec-sqrt -fp-model fast=2	24.88	3.11	0.20
-O3 -xMIC-AVX512 -no-prec-sqrt -fp-model fast=2	70.88	4.43	0.49

Conclusions and Insights

- Increased performance of algorithm using a pseudo Verlet list and particle sorting
- Implemented a local particle cache (SoA)
- Implemented a vectorisation strategy
 - Only read particles into cache that interact
 - Calculate all interactions on a particle and store results in a set of intermediate vectors
 - Perform horizontal add on intermediate vectors and update the particles with the result
 - Pad caches to prevent remainders and mask out the result
- Obtained speed-up on AVX, AVX2 and AVX512 instruction sets

Future Work

- Reduce the impact of overheads even further
- Improve vectorisation efficiency to obtain speedup closer to 8x and 16x for AVX and AVX-512 instruction sets

Questions

- Thank you for your attention
- Any questions?
- Website: www.icc.dur.ac.uk/swift/

Calculating Interactions

```
// AVX intrinsics
vector interactionMask
vector v_densitySum;
vector v_mj, v_wi, v_r2, v_hi2;

// Form mask
interactionMask = _mm256_cmp_ps(v_r2, v_hi2, _CMP_LT_OQ);

// Mask and add to density sum
v_densitySum = _mm256_add_ps(v_densitySum, _mm256_and_ps(interactionMask, _mm256_mul_ps(v_mj, v_wi)));

// AVX-512 intrinsics
__m512 interactionMask

// Form mask
interactionMask = _mm512_cmp_ps_mask(v_r2, v_hi2, _CMP_LT_OQ);

// Mask and add to density sum
v_densitySum = _mm512_mask_add_ps(v_densitySum, interactionMask, _mm512_mul_ps(v_mj, v_wi), v_densitySum);
```