# Tutorial 02 – C++ Threaded Programming

Dennis-Florian Herr, M.Sc.

Chair for Computer Architecture and Parallel Systems   (CAPS)

Technical University Munich

May 07, 2024



TUM Uhrenturm

Assignment 1: VV-AES

# General Ideas to optimize sequential code

- Improve algorithmic performance (think about asymptotic complexity / Big O)
- Think about the right data structure for the problem
- Optimize cache usage (cache access is orders of magnitude faster than main memory access)
- Reuse previously computed values (cf. dynamic programming)
- Consider precomputing an often recomputed small amount of data
- Do not reinvent the wheel (check *at least* the standard library for existing solutions)

# Improve substitute_bytes

Slow approach by repeatedly searching through the originalCharacter array

Before

```
1   void substitute_bytes() {
2     // For each byte in the message
3     for (int column = 0; column < BLOCK_SIZE; column++) {
4       for (int row = 0; row < BLOCK_SIZE; row++) {
5         int index = -1;
6         for (int i = 0; i < UNIQUE_CHARACTERS; i++) {
7           if (originalCharacter[i] == message[row][column]) {
8             index = i;
9           }
10        }
11        message[row][column] = substitutedCharacter[index];
12      }
13    }
14  }
```

# Improve substitute_bytes

Create a lookup map reducing the search from up to 256 iterations to 1

Optimized

```
1  uint8_t substituteMap[UNIQUE_CHARACTERS];
2  void create_substitute_map(){
3      for(int i = 0; i < UNIQUE_CHARACTERS; i++){
4          substituteMap[originalCharacter[i]] = substitutedCharacter[i];
5      }
6  }
7  void substitute_bytes() {
8      // For each byte in the message
9      for (int column = 0; column < BLOCK_SIZE; column++) {
10         for (int row = 0; row < BLOCK_SIZE; row++) {
11             message[row][column] = substituteMap[message[row][column] ];
12         }
13     }
14 }
```

# Improve substitute_bytes

And we're done.... (oops)

| Runtime | Speedup | Status |
|---------|---------|--------|
| 8.11421 | 8.25712 | ☑ Passed |

# Improve shift_rows

Don't reinvent the wheel. Use the std library rotate function

optimized

```cpp
#include <algorithm>
void shift_rows() {
    // Shift each row, where the row index corresponds to how many columns the data is shifted.
    for (int row = 0; row < BLOCK_SIZE; ++row) {
        std::rotate(std::begin(message[row]), std::begin(message[row]) + row, std::end(message[row]));
    }
}
```

# Improve shift_rows

## Alone it's dwarfed by substitute_bytes

| | | |
|---|---|---|
| 59.94560 | 1.11768 | ⏱ Too Slow |

## But together it's a significant boost

| Runtime | Speedup | Status |
|---|---|---|
| 3.77289 | 17.75829 | ☑ Passed |

# Improve mix_columns

Precompute all possible powers of 0 to 255

optimized

```
1   int powers[256][BLOCK_SIZE + 1];
2   // precomputed powers of all possible message values (256)
3   for (int i = 0; i < 256; i++) {
4       for (int j = 1; j <= BLOCK_SIZE; j++) {
5           powers[i][j] = power(i, j);
6       }
7   }
8   void mix_columns() {
9       for (int column = 0; column < BLOCK_SIZE; ++column) {
10          for (int row = 0; row < BLOCK_SIZE; ++row) {
11              int result = 0;
12              for (int degree = 0; degree < BLOCK_SIZE; degree++) {
13                  result += polynomialCoefficients[row][degree] * powers[message[degree][column]][degree + 1];
14              }
15              message[row][column] = result;
16          }
17      }
18  }
```

# Improve substitute_bytes

| Runtime | Speedup | Status |
|---------|---------|--------|
| 3.04771 | 21.98369 | ☑ Passed |

# Pushing further

## Use CISC PSHUFB: Packed Shuffle Bytes

```
1  uint8_t shift_row_mask[64] = {
2      0, 1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 8,
3      2, 3, 4, 5, 6, 7, 0, 1, 11, 12, 13, 14, 15, 8, 9, 10,
4      4, 5, 6, 7, 0, 1, 2, 3, 13, 14, 15, 8, 9, 10, 11, 12,
5      6, 7, 0, 1, 2, 3, 4, 5, 15, 8, 9, 10, 11, 12, 13, 14
6  };
7
8  inline void shift_rows() {
9      __asm__ (
10     "leaq    shift_row_mask(%rip), %rax\n"
11     "vmovdqa message(%rip), %xmm0\n"
12     "vpshufb (%rax), %xmm0, %xmm0\n"
13     "vmovaps %xmm0, message(%rip)\n"
14     "vmovdqa 16+message(%rip), %xmm1\n"
15     "vpshufb 16(%rax), %xmm1, %xmm1\n"
16     "vmovaps %xmm1, 16+message(%rip)\n"
17     "vmovdqa 32+message(%rip), %xmm2\n"
18     "vpshufb 32(%rax), %xmm2, %xmm2\n"
19     "vmovaps %xmm2, 32+message(%rip)\n"
20     "vmovdqa 48+message(%rip), %xmm3\n"
21     "vpshufb 48(%rax), %xmm3, %xmm3\n"
22     "vmovaps %xmm3, 48+message(%rip)\n"
23     );
24  }
```

# Leaderboard

A word on the leaderboard

- Friendly competition. Completely optional!
- Code/algorithm can be changed freely (as long as the result is accepted)
- ...but please don't try to hack the system!
- ...and don't write scripts to submit hundreds of times (I will delete and disqualify these)

How to compete?

- Find optimal parallelization
- Optimize sequential performance
- Skip unnecessary work
- Approximate solutions might be good enough!

# Need for Speed

*Need for Speed*

Share your cool solutions (after the deadline)

https://zulip.in.tum.de/#narrow/stream/2330-ParProg-24/topic/Need.20for.20Speed

Theory

# C++ Threads : Passing Data

## Creating Threads

```
1   #include <thread>
2   std::thread thread_object( callable )
3   std::thread thread_object( callable , arguments)
```

# C++ Threads : Passing Data

## Creating Threads

```
1                    #include <thread>
2                    std::thread thread_object( callable )
3                    std::thread thread_object( callable , arguments)
```

## Pass by Pointer

Main

```
1  ...
2  int argument = 1;
3  std::thread thread( kernel , &argument )
4  ...
```

Kernel

```
1  void* kernel( int* args) {
2      ++(*args);
3      ...
4  }
```

# C++ Threads : Passing Data

## Creating Threads

```
1    #include <thread>
2    std::thread thread_object( callable )
3    std::thread thread_object( callable , arguments)
```

## Pass by Pointer

Main

```
1  ...
2  int argument = 1;
3  std::thread thread( kernel, &argument )
4  ...
```

Kernel

```
1  void* kernel( int* args) {
2      ++(*args);
3      ...
4  }
```

## Pass by Reference

Main

```
1      ...
2  int argument = 1;
3  std::thread thread( kernel, std::ref(argument) )
4  ...
```

Kernel

```
1  void* kernel( int &args) {
2      ++args;
3      ...
4  }
```

# C++ Threads: Join

Joining Threads

```
1    #include <thread>
2    std::thread::join()
```

# C++ Threads: Join

Parallel Programming 2024
Tutorial 02

## Joining Threads

```
1    #include <thread>
2    std::thread::join()
```

## Example

<div align="center">Main</div>

```
1  ...
2  // Fork a thread
3  std::thread t(my_kernel);
4
5  // Join the thread
6  t.join();
7  ...
```

<div align="center">Kernel</div>

```
1  void my_kernel () {
2      sleep(1);
3      ...
4  }
```

Homework Solution        **Theory**        Quiz        In-Class Exercise        Homework        Questions        16

# C++ Threads: Pitfalls

## Example 1

Main

```
1  ...
2  // Start numThreads threads
3  std::thread threads[numThreads];
4  for (int i = 0 ; i< numThreads ; ++i){
5      threads[i] = std::thread(oops_kernel, &i);
6  }
7  for (int i = 0 ; i< numThreads ; ++i){
8      threads[i].join();
9  }
10 ...
```

Kernel

```
1  void oops_kernel (int *args) {
2      int argument = *args;
3      ...
4  }
```

# C++ Threads: Pitfalls

Example 1

<div style="text-align:center">Main</div>

```
1  ...
2  // Start numThreads threads
3  std::thread threads[numThreads];
4  for (int i = 0 ; i< numThreads ; ++i){
5      threads[i] = std::thread(oops_kernel, &i);
6  }
7  for (int i = 0 ; i< numThreads ; ++i){
8      threads[i].join();
9  }
10 ...
```

<div style="text-align:center">Kernel</div>

```
1  void oops_kernel (int *args) {
2      int argument = *args;
3      ...
4  }
```

✘ Use of freed memory / Race Condition

# C++ Threads: Pitfalls
## Example 1

Main

```
1  ...
2  // Start numThreads threads
3  std::thread threads[numThreads];
4  for (int i = 0 ; i< numThreads ; ++i){
5      threads[i] = std::thread(oops_kernel, &i);
6  }
7  for (int i = 0 ; i< numThreads ; ++i){
8      threads[i].join();
9  }
10 ...
```

Kernel

```
1  void oops_kernel (int *args) {
2      int argument = *args;
3      ...
4  }
```

✘ Use of freed memory / Race Condition

Kernel

```
1  void wicked_kernel (int *args) {
2      (*args) = 0;
3      ...
4  }
```

## Example 2

Main

```
1  ...
2  // Start numThreads threads
3  int ids[numThreads];
4  std::thread threads[numThreads];
5  for (int i = 0 ; i< numThreads ; ++i){
6      ids[i] = i;
7      threads[i] = std::thread(a_kernel, &ids[i]);
8  }
9  for (int i = 0 ; i< numThreads ; ++i){
10     threads[i].join();
11 }
12 ...
```

Kernel

```
1  void a_kernel (int* args) {
2      int argument = *args;
3      ...
4  }
```

# C++ Threads: Pitfalls

## Example 2

Main

```
1   ...
2   // Start numThreads threads
3   int ids[numThreads];
4   std::thread threads[numThreads];
5   for (int i = 0 ; i< numThreads ; ++i){
6       ids[i] = i;
7       threads[i] = std::thread(a_kernel, &ids[i]);
8   }
9   for (int i = 0 ; i< numThreads ; ++i){
10      threads[i].join();
11  }
12  ...
```

Kernel

```
1   void a_kernel (int* args) {
2       int argument = *args;
3       ...
4   }
```

✔ No Problems

# C++ Threads: Pitfalls

Example 3

Main

```
1  std::thread threads[numThreads];
2  void lets_spawn_threads(){
3
4    int ids[numThreads];
5    for (int i = 0 ; i< numThreads ; ++i){
6      ids[i] = i;
7      threads[i] = std::thread(noooo_kernel, &ids[i]);
8    }
9  }
10
11 int main(){
12   lets_spawn_threads();
13   for (int i = 0 ; i< numThreads ; ++i){
14     threads[i].join();
15   }
16
17 }
```
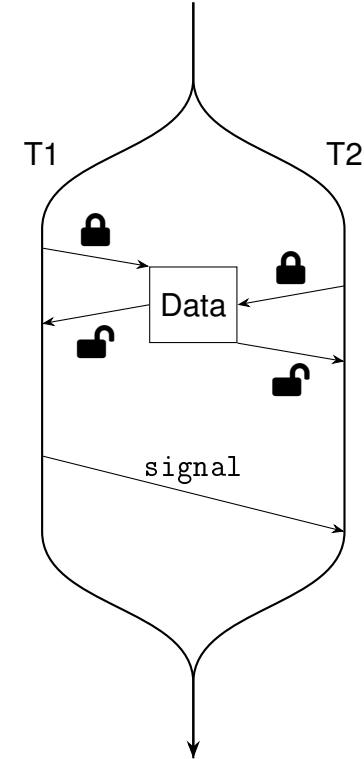
Kernel

```
1  void noooo_kernel (int* args) {
2    int argument = *args;
3    ...
4  }
```

# C++ Threads: Pitfalls

Example 3

Parallel Programming 2024
Tutorial 02

Main

```cpp
std::thread threads[numThreads];
void lets_spawn_threads(){

  int ids[numThreads];
  for (int i = 0 ; i< numThreads ; ++i){
    ids[i] = i;
    threads[i] = std::thread(noooo_kernel, &ids[i]);
  }
}

int main(){
  lets_spawn_threads();
  for (int i = 0 ; i< numThreads ; ++i){
    threads[i].join();
  }

}
```

Kernel

```cpp
void noooo_kernel (int* args) {
    int argument = *args;
    ...
}
```

✖ Use of freed memory

Homework Solution    **Theory**    Quiz    In-Class Exercise    Homework    Questions    19

# Synchronization

- Needed for accesses to shared resources
- Drawback: Overhead
- Frequent reasons for synchronizing:
  - Prevention of concurrent access
  - Signal passing
- C++ standard library provides following mechanisms (not exhaustive):
  - Mutexes
  - Condition Variables
  - Barriers (not covered)
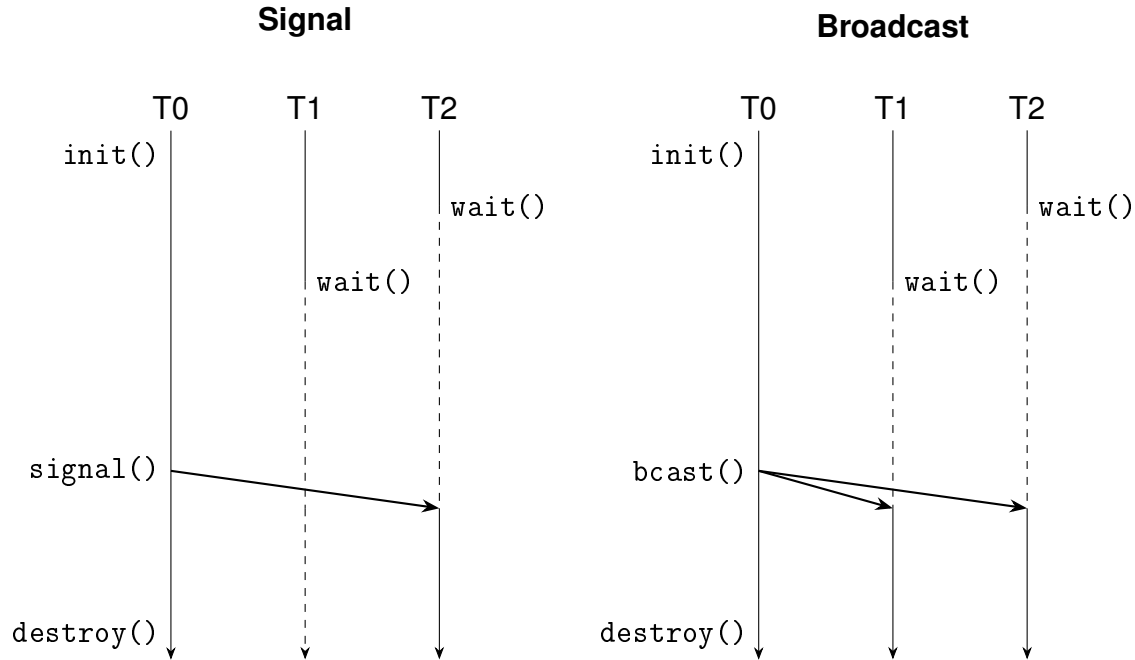  - Semaphores (not covered)

# Mutex (Mutual Exclusion Lock)

- The simplest and most primitive synchronization method
- Uses atomic (hardware) operations
- Ensures absolute owner of (critical) code section
- Threads can lock and unlock mutexes

# Condition Variables – Signaling



**Signal**

T0    T1    T2

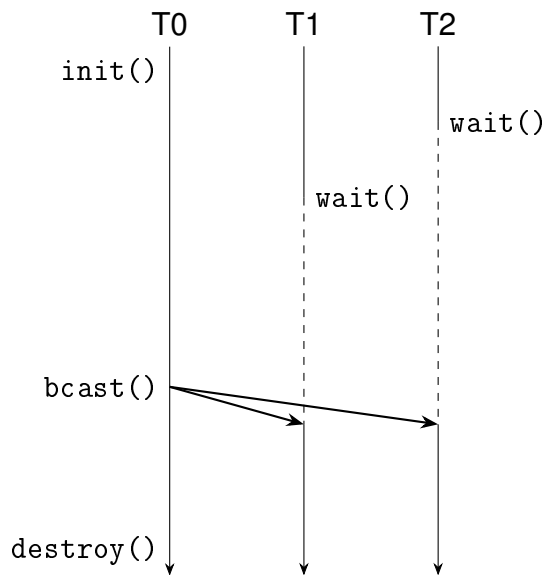init()

wait()

wait()

signal()

destroy()

# Condition Variables – Signaling

# Condition Variables – Signaling

**Signal**

T0     T1     T2

`init()`

`wait()`

`wait()`

`signal()`

`destroy()`

**Broadcast**

T0     T1     T2

`init()`

`wait()`

`wait()`

`bcast()`

`destroy()`

**Timed Wait**

T0     T1

`init()`

`timedwait()`

`timedwait()`

`signal()`

`destroy()`

# Work Distribution Pitfalls

### Undercounting

```
 1  int blockSize = SIZE / NUM_THREADS;
 2
 3
 4
 5
 6  for(int i = 0; i < blockSize; i++) {
 7      int dataIndex = threadId * blockSize + i;
 8      ...
 9  }
10  ...
```

Data



### Overcounting

```
 1  int blockSize = 8;
 2  int myBlock = 0;
 3
 4  while(myBlock * blockSize < SIZE) {
 5      myBlock = getNextBlock();
 6      for(int i = 0; i < blockSize; i++) {
 7          int dataIndex = myBlock * blockSize + i;
 8          ...
 9      }
10  }
```

Data

# Moodle Quiz

`https://www.moodle.tum.de/mod/quiz/view.php?id=2976059`

# In-Class Exercise

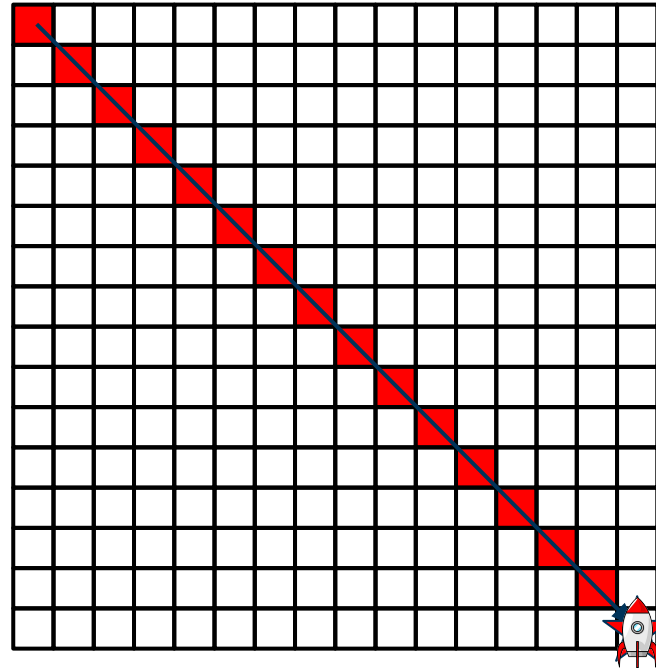# Help to calculate the starship's damage

The classic asteroids game
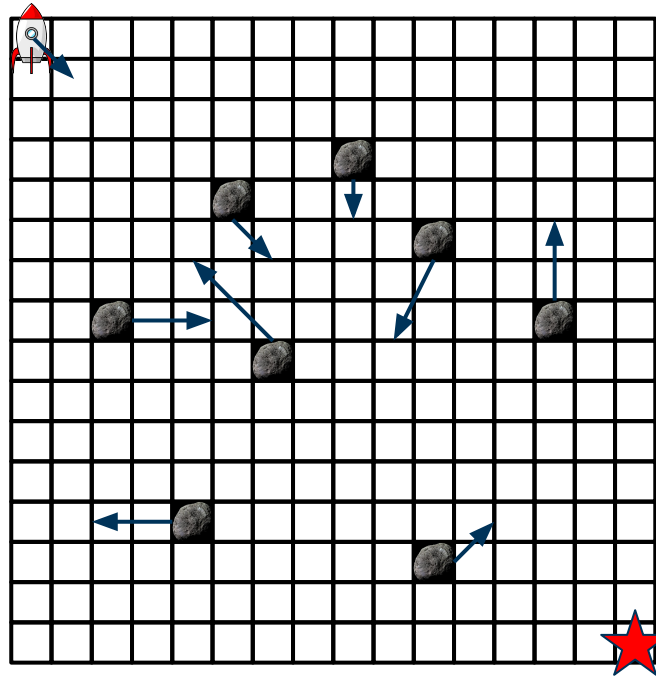
# Help to calculate the starship's damage



Our starship wants to go to the goal along the diagonal of the grid map
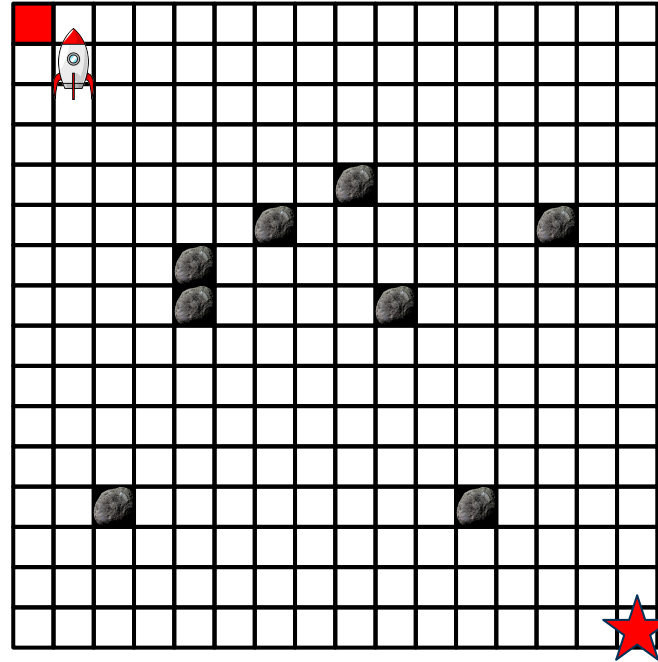
# Help to calculate the starship's damage

The red blocks are the intended trajectory of the starship

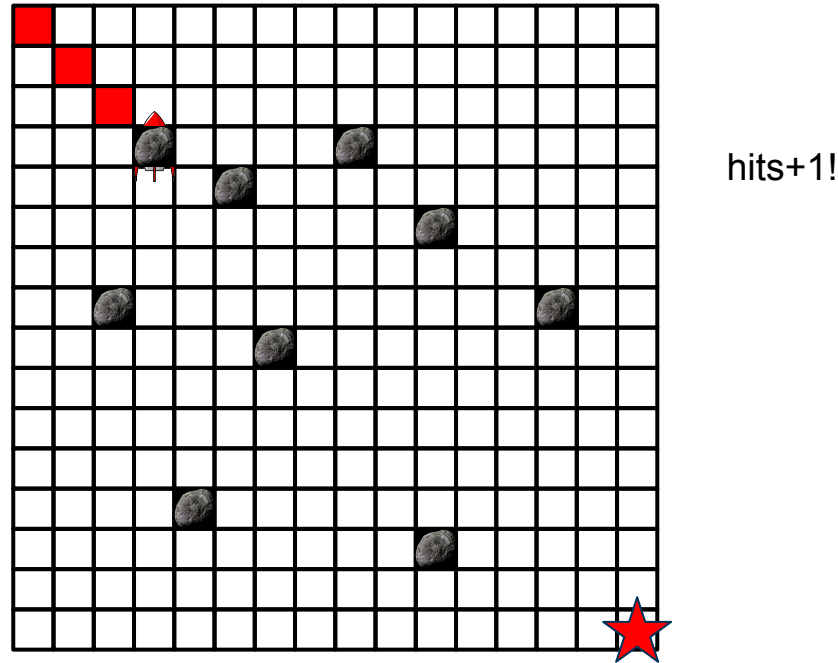# Help to calculate the starship's damage

But there are many rocks floating in the space, having very complicated velocity function
We can call *compute_vel( )* function to calculate the velocity of a rock

# Help to calculate the starship's damage

At every time step, all the rocks in the map and the starship will update their location

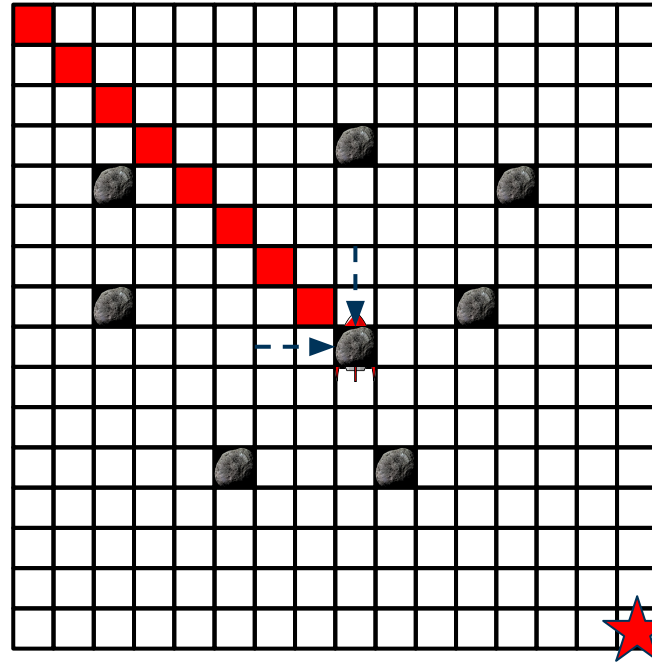We can call *update_rock( )* to update a rock's locations

# Help to calculate the starship's damage

hits+1!

When our starship moves towards the goal, it could crash rocks
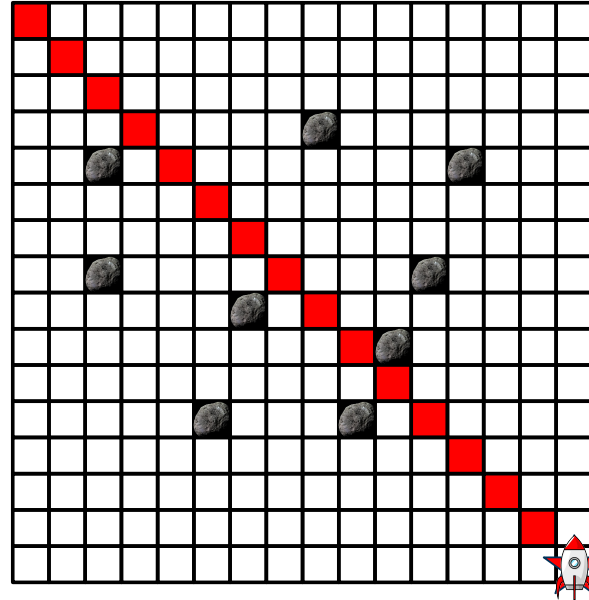We can call *calc_hits( )* to calculate how many rocks hit the starship at this time step

# Help to calculate the starship's damage



hits+2!

Several rocks can hits our starship at the same time

# Help to calculate the starship's damage



Help us to calculate the total hits number faster!
Although we implemented the sequential way to calculate this, but it is too slow!

Use C++ standard library threads to help us calculate this faster !

# Where to find the exercise ?

- Go to the following repository to get the exercise:

  `https://gitlab.lrz.de/lrr-tum/teaching/parprog/ss2024/published-assignments`

- Use git to clone the exercise to your local machine:

  cd your_folder

  git clone `https://gitlab.lrz.de/lrr-tum/teaching/parprog/ss2024/published-assignments.git`

- You can pull from this repository every time a new exercise is published

- Go to folder in-class-2 for this week's task, you can also find a README.md there

# Week 2 exercise introduction:

Use C++ standard library threads

- You will find `student_submission.cpp` the partially implemented parallel code
- Complete the //TODO section to
  - include necessary library
  - create threads for parallelization (pass called function, arguments etc.)
  - join threads to terminate parallelization
  - use mutex locks to avoid data racing

- achieve a speed up of 12 on the sever.

- Our server has 16 cores and 2 way hyperthreading (i.e. 32 threads)

# In-Class Exercise: Solution

# TODO#1: Include thread and mutex

```
1  // ###################### TODO: include library for enabling thread and mutex ######################
2  #include <thread> // For std::thread
3  #include <mutex> // For std::mutex
4  // ###################### TODO END ######################
```

```
1  // TODO: uncomment once you added the correct headers
2  std::thread threads[THREAD_NUM];
3  std::mutex mutex;
```
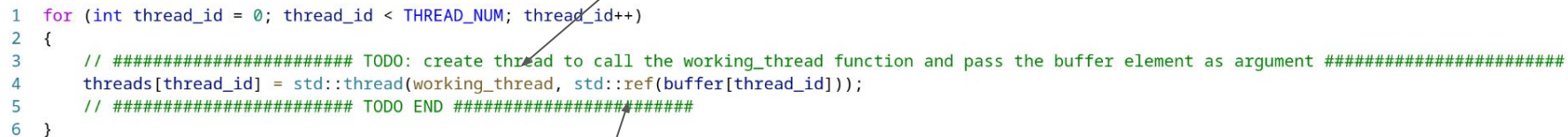
# TODO#2: Implement synchronized access to global variable

```
1  // ####################### TODO: Copy task id from the global variable and increment it #######################
2  // ####################### DO NOT FORGET TO LOCK AND UNLOCK! #######################
3  mutex.lock();
4  local_task_id = task_id++;
5  mutex.unlock();
6  // ####################### TODO END #######################
```

# TODO#3: Create the threads

Thread kernels need to have return type void!

```
1  for (int thread_id = 0; thread_id < THREAD_NUM; thread_id++)
2  {
3      // ##################### TODO: create thread to call the working_thread function and pass the buffer element as argument #####################
4      threads[thread_id] = std::thread(working_thread, std::ref(buffer[thread_id]));
5      // ##################### TODO END #####################
6  }
```

Ensures we pass the reference and avoids (rare) compile errors

# TODO#4: Wait for threads to complete and accumulate result

```
1  for (int thread_id = 0; thread_id < THREAD_NUM; thread_id++)
2  {
3      // ##################### TODO: join thread to terminate thread, get the returned value from the buffer and add it to the (total) crashed_count #####################
4      threads[thread_id].join();
5      crashed_count += buffer[thread_id];
6      // ##################### TODO END #####################
7  }
```
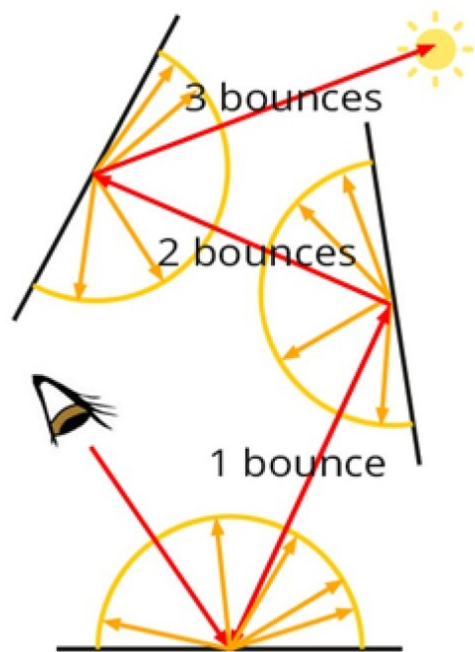
Homework

# Assignment: Raytracer

- Task: Render a 3D scene consisting of metallic spheres using raytracing
  - The program creates a random scene based on the seed read from `stdin.`
- The renderer sends rays for each pixel.
- When a ray hits a surface it gets reflected. A ray can only be reflected a certain number of times.
- Color of the pixel is determined by the materials of the surfaces a ray hits.
- Each pixel is sampled multiple times to reduce noise.



© www.scratchapixel.com

Figure 1: Rendering with a raytracer.

Figure 2: Raytraced scene.

- Parallelize the sequential implementation with C++ threads.
- Your speedup should be $\geq 10$.
- You don't have to read `maths.h` or `raytracer.h` to work on your solution.
- Evaluation command: `./student_submission -n ≪ <seed>`
- The output is a PPM image file.
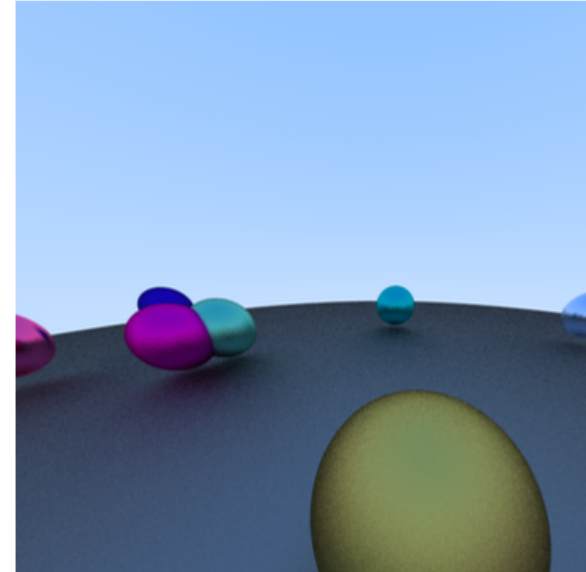  - Try xdg-open or use `http://paulcuth.me.uk/netpbm-viewer` to view PPM files.

# Recap & Questions

Covered today:

- Launching & Joining Threads
- Synchronization

# Questions

This slide is intentionally left blank.