

A fast GPU Monte Carlo radiative heat transfer implementation for coupling with direct numerical simulation

S. Silvestri, R. Pecnik*

Process and Energy Department, Delft University of Technology, Delft, the Netherlands

ARTICLE INFO

Article history:

Received 24 September 2018

Received in revised form 30 April 2019

Accepted 3 May 2019

Available online 15 May 2019

Keywords:

Radiative heat transfer

Monte Carlo simulation

Graphical processing units

Coupled radiation and convection

ABSTRACT

We implemented a fast Reciprocal Monte Carlo algorithm to accurately solve radiative heat transfer in turbulent flows of non-grey participating media that can be coupled to fully resolved turbulent flows, namely to Direct Numerical Simulation (DNS). The spectrally varying absorption coefficient is treated in a narrow-band fashion with a correlated- k distribution. The implementation is verified with analytical solutions and validated with results from literature and line-by-line Monte Carlo computations. The method is implemented on GPU with a thorough attention to memory transfer and computational efficiency. The bottlenecks that dominate the computational expenses are addressed, and several techniques are proposed to optimize the GPU execution. By implementing the proposed algorithmic accelerations, while maintaining the same accuracy, a speed-up of up to 3 orders of magnitude can be achieved.

© 2019 Published by Elsevier Inc. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

1. Introduction

Modeling radiative heat transfer is a challenging task due to the numerical complexity and the associated computational costs [1]. For example, radiative heat transfer is a six dimensional problem, which depends on spatial location, propagation direction and frequency of the electromagnetic wave. In addition, the calculation of radiative heat transfer poses a daunting challenge when it is coupled with convective and conductive heat transfer modes in turbulent flows. The computational cost of solving the radiative transfer equation makes it difficult to obtain an accurate description on how radiative transfer couples to a participating turbulent fluid flow. As a consequence, a complete view of the interplay between turbulence and radiation is missing.

Recently, several studies have addressed the problem with the aid of simplifying assumptions to ease the computational burden. In particular, Sakurai et al. [2] used the Optically Thin Approximation (OTA) to study the influence of radiative effects in a horizontally buoyant turbulent channel flow. They noticed that large scale buoyant structures are destroyed by the presence of non-local radiative heat transfer between walls and fluid. The OTA neglects self-absorption from the participating medium, leading to a constant incident radiation throughout the domain. This assumption greatly simplifies the description of radiative heat transfer. However, it does not allow the evaluation of incident radiation fluctuations and is, therefore, restricted to low values of absorption coefficient, κ , as demonstrated in Ref. [3].

A common approximation employed in solving the radiative transfer equation (RTE) in DNS coupled simulations consists in neglecting the spectral dependency of κ by assuming a grey gas. This assumption enables the efficient use of finite

* Corresponding author.

E-mail addresses: s.silvestri@tudelft.nl (S. Silvestri), r.pecnik@tudelft.nl (R. Pecnik).

difference schemes, such as the discrete ordinates method (DOM), which require a low computational effort and provide a high level of accuracy. Examples are given in Refs. [4–8], who have studied the influence of radiative heat transfer in turbulent flows using the grey gas approximation, coupled to either DNS or large eddy simulations (LES). These studies allow to investigate the effect of radiation on the turbulent temperature field and vice versa, to highlight and quantify the dissipative effect of radiative field fluctuations and the impact of Turbulence Radiation Interactions (TRI). However, these cases are highly idealized since radiation in real fluids is intrinsically non-grey.

If a spectral description of radiative heat transfer is to be included, the state of art involves the use of a Monte Carlo (MC) method. Compared to the above mentioned RTE solution methods, the Monte Carlo method can be considered the most accurate and flexible. Its solution time increases mildly with problem complexity, allowing a detailed spectral description or the simulations of complex geometries, which are challenging with other methods such as DOM. To the authors knowledge, the first instance of a Monte Carlo method coupled with DNS is reported by Wu et al. [9]. They developed a high resolution MC method, subsequently used by Deshmukh et al. [4] to study TRI in a statistically one dimensional premixed combustion system. In their study, they noticed that absorption TRI intensifies with an increase of optical thickness, while emission TRI is always relevant in reactive flows. The calculations were performed on a 64^3 mesh in a grey gas, but the use of a MC method coupled to DNS on finer grids and spectral medium was not investigated. On the other hand, more recently, Zhang et al. [10] and Vicquelin et al. [11] performed DNS of radiative channel flow using a MC code coupled to a narrow-band correlated- k spectral description. They investigated the effect of different radiative budgets to verify the modification of first and second order temperature and velocity statistics. Nonetheless, the computational expenses of the MC method prohibited the solution of the radiative heat transfer on a full DNS grid, which required an intermediate interpolation step between flow and radiation solution. Overall, the computational expenses of the Monte Carlo solver limit the possibilities of accurately solving coupled radiative heat transfer and turbulence together. However, since the Monte Carlo method is “embarrassingly parallelizable” (i.e., can be divided into a number of completely independent computations), it greatly benefits from the use of parallel architectures and in particular from the use of general purpose graphical processing units (GPGPU).

The use of GPUs for computational sciences has become increasingly investigated, especially for large parallelizable problems that are more efficiently mapped on many GPU parallel multiprocessors [12]. The development of NVIDIA CUDA, a versatile GPU programming language, has further popularized GPUs as accelerators alongside CPU computational clusters [13]. Several examples of GPU codes are available up to date, ranging from machine learning [14] to imaging [15] and computational biology [16]. Likewise, in the fluid mechanics field, Khajeh et al. [17] and Salvatore et al. [18] have been porting a Navier-Stokes solver on GPUs obtaining speedups up to $22 \times$. Additionally, many Monte Carlo codes have been developed on graphical processing units for many diverse fields and applications such as finance [19] and molecular dynamics [20]. On the other hand, to the authors knowledge, the only MC methods applied to the solution of thermal radiation implemented on GPU were developed independently by Humphrey et al. [21] for grey gas applications and by Heymann and Siebenmorgen [22] for dust radiation around active galactic nuclei. The code of the former showed excellent scaling capabilities up to 16834 GPUs, proving the feasibility of the GPU MC concept for thermal radiation. The latter focused on the optimization of the MC procedure, reaching a speed up of around $100 \times$.

Nonetheless, porting an application to GPU requires the exposure of the parallel portion of the application and algorithmic optimizations to improve the efficiency on a GPU architecture. Therefore, the objective of this work is to develop an optimized GPU Monte Carlo implementation, which can enable a fast and accurate solution of radiative heat transfer in largely fluctuating temperature fields typical of turbulent flows. We will include the spectral description of the absorption coefficient to have a complete and flexible solver. All the challenges involved in implementing an efficient GPU application are addressed in order to reduce the computational time and to improve the scaling with problem size.

2. The Monte Carlo method

In this section the details of the Monte Carlo methods are outlined for the sake of completeness. Within a domain containing a non grey absorbing and emitting medium, the radiative power emitted by cell i and absorbed within cell j is expressed, as in Tesse et al. [23], by

$$Q_{i \rightarrow j}^R = \int_0^\infty \kappa_\nu(T_i) I_{b\nu}(T_i) \int_{V_i} \int_{4\pi} \sum_{m=1}^{N_c} \tau_\nu(i \rightarrow j, m) \left[\int_0^{l_{j,m}} \kappa_\nu(T_j) e^{-\kappa_\nu(T_j)s_{j,m}} ds_{j,m} \right] d\Omega dV_i d\nu, \quad (1)$$

where ν is the wavenumber, κ_ν is the spectral absorption coefficient, τ_ν is the spectral transmissivity from cell i to the boundary of cell j following the path m , N_c is the number of paths that, from cell i , cross cell j , and $l_{j,m}$ is the distance travelled in cell j along the propagation direction. The volume integral V_j , as given in ref. [24] has been replaced by the integration over the solid angle Ω and the path length $s_{j,m}$ as done in ref. [23]. The integral in the square brackets represents the absorption within cell j , following path m . The analytical solution, considering cell j isothermal and homogeneous, is

$$\alpha_{\nu j,m} = 1 - e^{-\kappa_\nu(T_j)l_{j,m}}. \quad (2)$$

The spectral transmissivity $\tau_\nu(i \rightarrow j, m)$ is the result of the absorption by the finite volumes and surfaces crossed by path m , and can be calculated as

$$\tau_v(i \rightarrow j, m) = \prod_{k=i}^{j-1} (1 - \alpha_{v,k,m}) \times \prod_{c=1}^{N_r} (1 - \varepsilon_w), \quad (3)$$

where ε_w is the wall emissivity and N_r is the number of wall reflections that occurred for path m .

The Monte Carlo method consists in a statistical estimation of the integrals in equation (1) using a large number of samples that represent different paths and wavelengths. In particular, it is possible to develop probability distribution functions defined as

$$f_v = \frac{1}{V_i}, \quad f_\theta = \frac{\sin \theta}{2}, \quad f_\phi = \frac{1}{2\pi}, \quad f_v = \frac{\pi \kappa_v(T_i) I_{b_v}(T_i)}{\kappa_p(T_i) \sigma T_i^4}, \quad (4)$$

where $\kappa_p(T_i)$ is the Planck mean absorption coefficient of cell i , while θ and ϕ are the polar and azimuthal angles, respectively, with $d\Omega = \sin \theta d\theta d\phi$. Substituting the probability distribution functions in equation (1) leads to

$$Q_{i \rightarrow j}^R = Q^{R,e}(T_i) \int_0^\infty f_v \int_{V_i} f_v \int_0^{2\pi} f_\phi \int_0^\pi f_\theta A_{v,m,i \rightarrow j} d\theta d\phi dV_i dv, \quad (5)$$

where $Q^{R,e}(T_i)$ and $A_{v,m,i \rightarrow j}$ are the total radiative power emitted by cell i and the spectral energy fraction emitted by cell i and absorbed in cell j through path m , respectively. These are calculated using

$$Q^{R,e}(T_i) = 4V_i \kappa_p(T_i) \sigma T_i^4, \quad (6)$$

$$A_{v,m,i \rightarrow j} = \sum_{m=1}^{N_c} \tau_v(i \rightarrow j, m) \alpha_{v,j,m}. \quad (7)$$

A statistical estimation of the integrals in equation (5) involves launching several samples, referred hereafter as “rays” with properties sampled from the probability density functions given in (4).

The resulting discretized equation has then the form

$$\widetilde{Q}_{i \rightarrow j}^R = \frac{Q^{R,e}(T_i)}{N_r} \sum_{r=1}^{N_r} A_{r,i \rightarrow j}. \quad (8)$$

The tilde \sim denotes a statistical estimator and the subscript r indicates a ray, characterized by its wavenumber ν , and direction angles θ and ϕ (defining the path variable m), which are calculated inverting the following relations

$$\begin{aligned} R_\nu &= \int_0^\nu f_{\nu'}(T) d\nu' = \frac{\pi \int_0^\nu \kappa_{\nu'}(T) I_{b_{\nu'}}(T) d\nu'}{\kappa_p(T) \sigma T^4}, \\ R_\theta &= \int_0^\theta f_{\theta'} d\theta' = \frac{1 - \cos \theta}{2}, \\ R_\phi &= \int_0^\phi f_{\phi'} d\phi' = \frac{\phi}{2\pi}. \end{aligned} \quad (9)$$

R_ν , R_θ and R_ϕ are random numbers sampled from a uniform probability distribution function between 0 and 1.

In a reciprocal Monte Carlo formulation, both emitted and absorbed power are statistically estimated as

$$Q_{i, RM}^R = \underbrace{\sum_{j=1}^{N_v+N_s} \widetilde{Q}_{i \rightarrow j}^R}_{Q_i^{R,e}} - \underbrace{\sum_{j=1}^{N_v+N_s} \widetilde{Q}_{j \rightarrow i}^R}_{Q_i^{R,a}}, \quad (10)$$

where $N_v + N_s$ are the numbers of volume and surfaces that interact with cell i . The reciprocal formulation employs the following principle

$$\frac{Q_{i \rightarrow j, v}^R}{I_{b_v}(T_i)} = \frac{Q_{j \rightarrow i, v}^R}{I_{b_v}(T_j)}, \quad (11)$$

to automatically satisfy the reciprocity condition. As a consequence, the above formulation avoids problems of large variance in case of low temperature gradients (i.e. non reactive flows) or high optical thickness that are typical of a forward Monte Carlo method. Depending on the estimated quantity, it is possible to distinguish between two reciprocity Monte Carlo formulations [23]. These are, the Absorption-based Reciprocity Monte Carlo (ARMC) which connects the emission to the statistical estimation of the absorbed power, and, vice-versa, the Emission-based Reciprocity Monte Carlo (ERMC) which estimates absorption based on the calculation of emitted power. While ARMC results in a lower variance in low temperatures zones, characterized by relevant absorption, ERMC is more accurate in the high temperature regions that are dominated by emission. The advantage of ERMC is that Q^R in i is calculated by the emission of the cell, requiring only the computation of the rays leaving the cell itself. The corresponding relation of an ERMC formulation is given as

$$Q_{i,ERMC}^R = \sum_{j=1}^{N_v+N_s} \widetilde{Q_{i \rightarrow j}^R} \cdot \left(1 - \frac{I_{bv}(T_j)}{I_{bv}(T_i)}\right). \quad (12)$$

In addition, in the ERMC formulation the volumetric heat source [kW/m³] is independent of the emitting cell volume. As a consequence, it is possible to emit all rays directly from the center of cell i , avoiding the sampling on f_v .

Recently, Zhang et al. [25] developed an optimized ERMC to reduce the variance in the low temperature regions. In the cold regions, Q^R is dominated by the absorption of radiation which originates from hot zones. Nevertheless, an ERMC entails the estimation of absorption based on the emission of the cell itself. Consequently, the wavelength of emission in colder regions will be higher than the actual wavelength of the absorbed radiation that follows Wien's displacement law. This leads to a large variance in cold spots, which is characteristic for an ERMC based method. Therefore, ref. [25] proposed to sample the wavenumber from the maximum temperature, which corresponds to a larger emission in the domain, using

$$f_v = \frac{\pi \kappa_v(T_{max}) I_{bv}(T_{max})}{\kappa_p(T_{max}) \sigma T_{max}^4}. \quad (13)$$

As a result, equation (8) has to be corrected with a prefactor R_I , resulting in

$$\widetilde{Q_{i \rightarrow j}^R} = \frac{Q^{R,e}(T_{max})}{N_r} \sum_{r=1}^{N_r} \underbrace{\left(\frac{\kappa_v(T_i) I_{bv}(T_i)}{\kappa_v(T_{max}) I_{bv}(T_{max})} \right)}_{R_I} A_{r,i \rightarrow j}. \quad (14)$$

2.1. Spectral discretization

In general, gas absorption spectra are characterized by discrete absorption lines, leading to a strong dependency on wavelength. In order to store the absorption coefficients and the probabilities associated with a line-by-line spectrum, comprised of more than a million spectral points, an excessive amount of memory is required. In addition, the high variability of the spectra translates in a lower convergence rate of the Monte-Carlo method. For this reason, we chose a narrow-band correlated- k model to couple with the Monte Carlo solver [26]. The narrow-band method constitutes of an accurate spectral representation, comparable to a line-by-line description if enough pseudo-spectral points are considered, with significantly lower memory requirements. In addition it is naturally adaptable to a simple implementation of species transport and wavenumber-dependent scattering, in case multiphase flows are considered. The line-by-line spectrum of common gasses, for a wide range of temperatures and pressures, can be found in accurate online spectroscopy databases. For this study, the data from HITRAN 2012 [27] and HITEMP 2010 [28] is used to develop the narrow-band pseudo-spectral coefficients.

Since the narrow-band correlated- k model divides the spectrum into narrow bands with assigned quadrature points, the wavenumber probability function in equation (4) is discretized using two "discrete" probability functions, one for the narrow-band and the other one for the quadrature point. The two variables associated with the wavenumber of the photon bundle are thus a narrow band index n and a quadrature point index g ,

$$\int_0^v f_{v'} dv' \approx \sum_{n'=1}^{n-1} f_{n'} + f_n \cdot \sum_{g'=1}^{g-1} f_{g'}(n), \quad (15)$$

where

$$f_n = \frac{\pi \Delta v_n I_{bn} \sum_{g'=1}^{N_q} \omega_{g'} k_{n,g'}}{\kappa_p \sigma T^4}, \quad f_g(n) = \frac{\omega_g k_{n,g}}{\sum_{g'=1}^{N_q} \omega_{g'} k_{n,g'}}, \quad (16)$$

and ω_g and N_q are the Gaussian weights associated with point g and the total number of quadrature points in a narrow-band, respectively. Since the quadrature points in a narrow-band all represent ideally the same wavenumber, the drawing of two independent random numbers is necessary in order to sample n and g ,

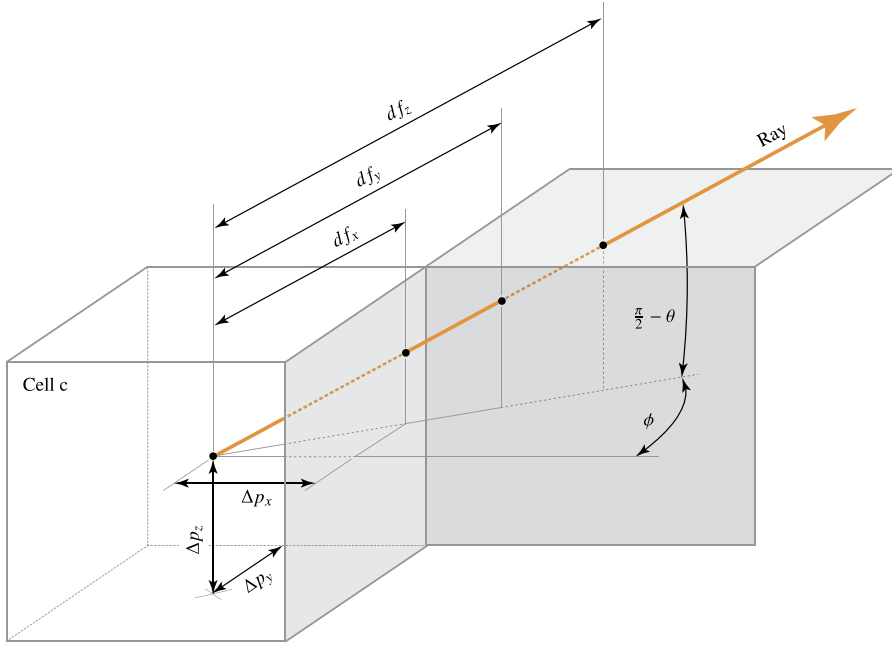


Fig. 1. Schematic displaying the marching ray procedure.

Algorithm 1 ERM CPU implementation.

```

1: for cell in Cells do
2:    $QE \leftarrow 4\kappa_P(T_{max})\sigma T_{max}^4 / \text{numberOfRays}$ 
3:   for ray in Rays do
4:     procedure INITIALIZE
5:        $R_\theta, R_\phi, R_n, R_g \leftarrow \text{Rand}(\text{uniform distribution})$ 
6:       ray.ind  $\leftarrow$  cell.ind
7:       ray.pos  $\leftarrow$  cell.center
8:       ray.dir  $\leftarrow \text{direction}(R_\theta, R_\phi)$ 
9:       ray.transmissivity  $\leftarrow 1.0$ 
10:      indDir  $\leftarrow \text{sign}(\text{ray.dir})$ 
11:       $n, g \leftarrow \text{findWavelength}(R_n, R_g)$ 
12:       $lb1 \leftarrow \text{interpBlackbody}(n, \text{temperature}(\text{ray.ind}))$ 
13:       $R_l \leftarrow lb1 \times \text{interpAbsorpCoeff}(n, g, \text{temperature}(\text{ray.ind}))$ 
14:       $R_l \leftarrow R_l / \text{interpBlackbody}(n, T_{max}) / \text{interpAbsorpCoeff}(n, g, T_{max})$ 
15:    end procedure
16:    procedure MARCH
17:      while ray.transmissivity > tol do
18:         $df \leftarrow \Delta p / \text{ray.dir}$ 
19:         $ds \leftarrow \min(df_x, df_y, df_z)$ 
20:         $\kappa \leftarrow \text{interpAbsorpCoeff}(n, g, \text{temperature}(\text{ray.ind}))$ 
21:         $\alpha \leftarrow 1 - \exp(-\kappa \times ds)$ 
22:         $lb2 \leftarrow \text{interpBlackbody}(n, \text{temperature}(\text{ray.ind}))$ 
23:        Absorption  $\leftarrow QE \times \text{ray.transmissivity} \times \alpha \times (lb2/lb1 - 1) \times R_l$ 
24:         $QR(\text{cell.ind}) \leftarrow QR(\text{cell.ind}) - \text{Absorption}$ 
25:        ray.pos  $\leftarrow \text{ray.pos} + ds \times \text{ray.dir}$ 
26:        ray.ind  $\leftarrow \text{ray.ind} + \text{indDir} \times (ds == [df_x, df_y, df_z])$ 
27:        ray.transmissivity  $\leftarrow \text{ray.transmissivity} \times (1 - \alpha)$ 
28:      end while
29:      Absorption  $\leftarrow QE \times \text{ray.transmissivity} \times (lb2/lb1 - 1) \times R_l$ 
30:       $QR(\text{cell.ind}) \leftarrow QR(\text{cell.ind}) - \text{Absorption}$ 
31:    end procedure
32:  end for
33: end for

```

▷ Loop over all finite volumes
 ▷ Cell emission $Q^{R,e}$ in equation (8)
 ▷ Loop over rays
 ▷ Draw random numbers for angles and indices n and g
 ▷ Initialize the ray with the cell index i, j , and k
 ▷ Initialize ray starting coordinates with cell center coordinates x, y , and z
 ▷ Find ray direction based on equation (9)
 ▷ Ray direction in terms of index i, j , and k
 ▷ Binary search on CDF with R_n and $CDF(n)$ with R_g
 ▷ blackbody intensity of initial cell c
 ▷ R_l in equation (14)
 ▷ Determine which face is crossed first (see Fig. 1)
 ▷ Shortest distance is where ray crosses face
 ▷ equation (2)
 ▷ equation (14)
 ▷ radiative heat source of initial cell c
 ▷ Update ray position
 ▷ Update cell index depending on which face has been intersected
 ▷ equation (3)
 ▷ Dump the residual energy into the initial cell

Table 1
Description of validation cases.

Case	κ	Temp.	Dimensions	Domain	ε_w	τ / t	Comparison
Case 1	1 [m ⁻¹]	lin1	1D	1 [m]	1 (all walls)	1 / 0.368	analytical solution
Case 2	1 [m ⁻¹]	parab1	1D	1 [m]	1 (all walls)	1 / 0.368	analytical solution
Case 3	0.5 [m ⁻¹]	sin	3D	1 [m ³]	1 (all walls)	0.5 / 0.607	analytical solution [29]
Case 4	5 [m ⁻¹]	sin	3D	1 [m ³]	1 (all walls)	5 / 0.007	analytical solution [29]
Case 5	H ₂ O	1000 [K]	1D	0.1 [m]	1 (all walls)	0.51 / 0.81	Kim et al. [30]
Case 6	H ₂ O	1000 [K]	1D	1 [m]	1 (all walls)	5.1 / 0.54	Kim et al. [30]
Case 7	CO ₂	lin2	1D	1 [m]	1 (all walls)	25.5 / 0.82	Cherkaoui et al. [24]
Case 8	CO ₂	lin2	1D	1 [m]	0, 1	25.5 / 0.82	Cherkaoui et al. [24]
Case 9	CO ₂	lin2	1D	1 [m]	0.1, 0.1	25.5 / 0.82	Cherkaoui et al. [24]
Case 10	H ₂ O	parab1	1D	1 [m]	1 (all walls)	4.98 / 0.49	Line-by-Line MC
Case 11	CO ₂	parab1	1D	1 [m]	1 (all walls)	22.11 / 0.821	Line-by-Line MC
Case 12	mix	parab2	1D	0.2 [m]	1 (all walls)	0.30 / 0.76	Tesse et al. [23]
Case 13	mix	parab2	1D	4 [m]	1 (all walls)	5.93 / 0.036	Tesse et al. [23]
Case 14	H ₂ O	3dimens	3D	1 [m ³]	1 (all walls)	4.98 / 0.49	Line-by-Line MC

$$R_n = \sum_{n'=1}^{n-1} f_{n'}, \quad R_g = \sum_{g'=1}^{g-1} f_{g'}(n). \quad (17)$$

2.2. Algorithm

To ease the understanding of the GPU ERMC implementation, we first describe a standard CPU implementation in Algorithm 1. The first loop (line 1) is performed over all finite volumes in the computational domain. Each finite volume is described by its index (i,j,k), and coordinates of the center and the surrounding faces. For each finite volume, a predefined number of rays (numberOfRays) are launched. The variable *ray* is a data structure that contains the current position (*pos*) of the ray and the index of the corresponding cell (*ind*), as well as the direction vector (*dir*) and the current transmissivity (*transmissivity*). The MC method mainly consists of two routines, the initialization (line 4) and the marching of the ray (line 16). In the first routine, the necessary random numbers are drawn and the properties of the ray are initialized accordingly. To accommodate a narrow-band correlated-*k* description, two independent random numbers are drawn R_n and R_g , which lead to two different indices n and g that specify the narrow-band and the quadrature point within. Marching the ray consists in finding the distances $\Delta p_x, \Delta p_y, \Delta p_z$, between the current location of the ray and the cell faces in direction *ray.dir*, specified by the angles ϕ and θ . The minimum distance, ds , determines which plane is crossed by the ray first. A schematic is displayed in Fig. 1, for which the ray intersects the *x*-normal plane first, such that the minimum distance ds will be equal to df_x . The radiative power of the initial cell (QR) is then calculated in a reciprocal fashion. Furthermore, the new ray position and cell index are updated accordingly. If the transmissivity drops under a certain tolerance *tol* (line 17), the ray is terminated and the remaining energy is dumped into the initial cell (line 29). The on-the-fly calculation of the blackbody intensity from Planck's law is prohibitive due to the excessive computations involved. To overcome this issue, the blackbody intensity is precomputed for the narrow band wavelengths and discrete points in the required temperature range and then stored in a suitable 2D table. The functions *interpBlackbody* and *interpAbsorptionCoeff* (lines 12-14, 20 and 22) perform linear interpolations of the spectral blackbody intensity and the absorption coefficient from the corresponding tables, respectively.

2.3. Verification and validation

To ensure a correct implementation, the algorithm is first verified and validated for a CPU implementation using a combination of grey and non-grey gases in 1D and 3D. In total, 14 cases are used which are summarized in Table 1. All the cases, 1D and 3D alike, are calculated on a 3D grid. To mimic the one-dimensionality, two periodic directions are employed in the 1D cases, over which the results are then averaged. Beside the case names in column 1, the second column shows values of the absorption coefficient κ for the grey gas cases (cases 1 to 4), and the names of the non-grey gases (H₂O, CO₂ or mix) for cases 5 to 14. The entry "mix", found in cases 12 and 13, refers to a mixture of 11.6% CO₂ and 15.5% H₂O at 1 [atm]. All other cases are enclosures filled with a pure substance (molar fraction of 1) at 1 [atm]. The other columns indicate the prescribed temperature distribution (linear, parabolic, etc.), the spatial inhomogeneous dimensions (1D or 3D), the wall emissivities ε_w , the global optical thickness and global transmissivity (τ / t), and the source used for the verification or the validation. Further details are given in the subsequent discussions of the individual cases.

The grey gas cases 1, 2, 3 and 4 are used to verify the correctness of the ray marching procedure and are compared to existing analytical solutions. Although $\kappa \neq f(\nu)$, the spectral (narrow-band) description shown in Algorithm 1 is retained with precomputed probability functions based on a grey gas absorption coefficient. Two different geometries are examined, namely a 1 [m] parallel slab (1D) and a 1 [m³] cube (3D). The walls are considered black with $\varepsilon_w = 1$. For the 1D cases, two different temperature profiles (lin1 and parab) are considered, given as

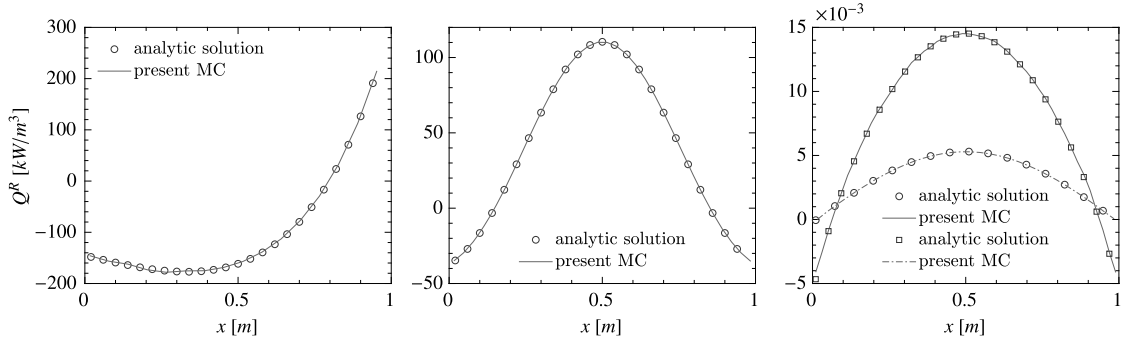


Fig. 2. Verification of the present MC code (lines) for a grey gas in comparison with analytic solution (symbols). Left: case 1, center: case 2, right: circles and dashed dotted line case 3, squares and solid line case 4 (both at $y = z = 0.5$ [m]).

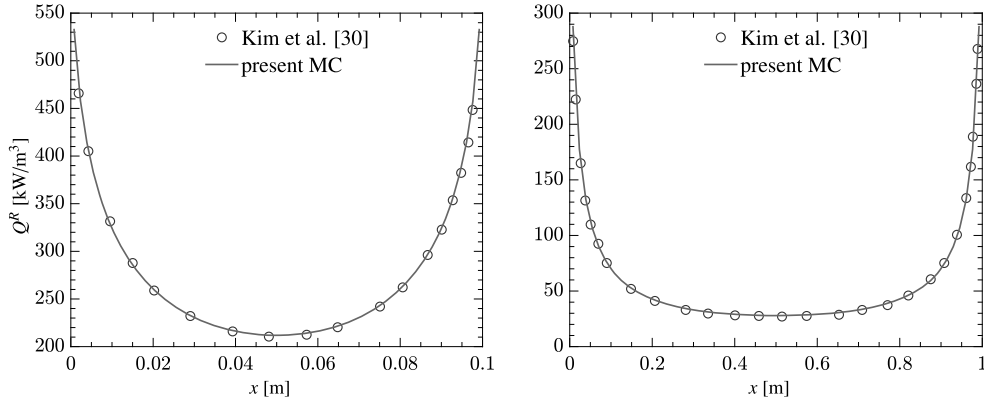


Fig. 3. Validation of the present MC code (lines) for H_2O in the isothermal case in comparison with values from [30] (circles). Left: case 5, right: case 6.

$$\text{lin1:} \quad T_m = 500 + 1000x \text{ [K]}, \quad T_{w1} = 500 \text{ [K]}, \quad T_{w2} = 1500 \text{ [K]}, \quad (18)$$

$$\text{parab1:} \quad T_m = 500 - 2000x^2 + 2000x \text{ [K]}, \quad T_{w1} = T_{w2} = 500 \text{ [K]}, \quad (19)$$

where T_m is the temperature of the medium and T_{w1} and T_{w2} are the temperatures at the left and right wall, respectively. For the 3D cases, the walls are cold (0 [K]), and the temperature profile is given as

$$\sin: T_m = (\sin\pi x \cdot \sin\pi y \cdot \sin\pi z \cdot \pi/\sigma)^{0.25} \text{ [K]}, \quad (20)$$

in order to compare the results with the quasi-analytic solution derived by Sakurai et al. [29]. The absorption coefficient for the 1D slab has a value of 1 [m^{-1}], while for the 3D domain the two cases have different absorption coefficients of $\kappa = 0.5$ and 5 [m^{-1}] (case 3 and 4, respectively). For these four cases, the results are obtained on a 32^3 grid with 2000 rays per cell. Note that for the 1D cases, an averaging was performed along the periodic directions. As can be seen in Fig. 2, the MC implementation is accurately able to reproduce the analytic solutions with adequate precision.

To validate the spectral discretization, a combination of isothermal and non-isothermal cases with H_2O and CO_2 have been used. 119 and 139 narrow bands were selected for H_2O and CO_2 , respectively, with each band containing 16 quadrature points. The radiative power of a 1D slab filled with water vapour at 1 [atm] and 1000 [K], bounded by two cold black walls at a distance of 0.1 (case 5) and 1 [m] (case 6), has been compared with data presented in Kim et al. [30] as shown in Fig. 3. The results for the 1D slab filled with CO_2 at 1 [atm] and three different wall emissivities (cases 7, 8 and 9) are shown in Fig. 4. The temperature profiles for the CO_2 cases are linear with the left wall at 295 [K] and the right wall at 305 [K] (lin2). The radiative power is compared to data presented in Cherkaoui et al. [24]. In all cases (cases 5–9) the comparison clearly demonstrates the high accuracy of the spectral discretization.

Two additional cases (case 10 and 11) are proposed to validate the spectral discretization and the MC implementation with a line-by-line version of the present MC code. The radiative power is calculated for H_2O and CO_2 at 1 [atm] with parabolic temperature profiles (equation (19)). The results obtained with the narrow-band correlated- k MC, shown in Fig. 5, are in close agreement with the line-by-line benchmark to again prove the correct implementation. Fig. 5 also shows the results for the H_2O - CO_2 mixture (cases 12 and 13) in comparison with the results from Tesse et al. [23]. For these cases the temperature profile was set to

$$\text{parab2:} \quad T_m = 500 - 8000(x/L)^2 + 8000(x/L) \text{ [K]}, \quad T_{w1} = T_{w2} = 500 \text{ [K]}, \quad (21)$$

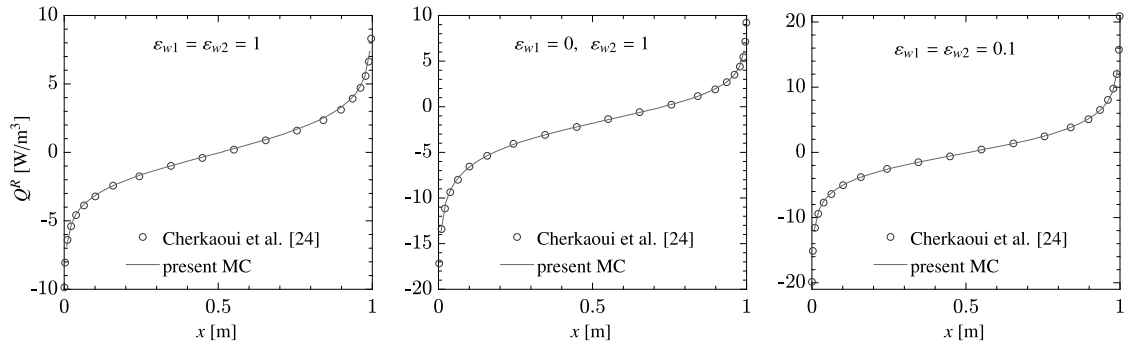


Fig. 4. Validation of the present MC code for CO₂ at 1 [atm] in comparison with results from [24]. Linear temperature profile $T = 295 + 10x$ [K]. $T_{w1} = 295$ [K] (lin2). $T_{w2} = 305$ [K]. Cases 7, 8 and 9 at the left, right and center, respectively.

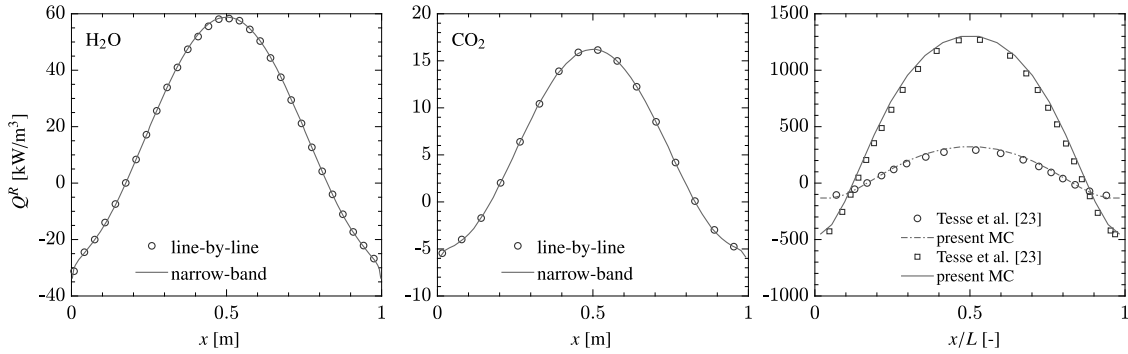


Fig. 5. Left and right show the comparison of MC implementation with the line-by-line solution for H₂O (left) and CO₂ (center) at 1 [atm] with a parabolic temperature profile, equation (19). The right figure shows the comparison of the H₂O-CO₂ mixture with results from Tesse et al. [23]. Circles and dashed-dotted line: case 13. Squares and solid line: case 12.

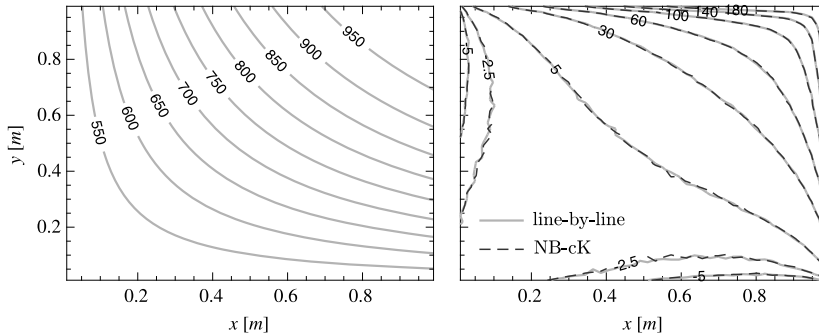


Fig. 6. 3D non grey case. Left: temperature profile at $z = 0.5$ [m] in [K]. Right: radiative power at $z = 0.5$ [m] in [kW/m²].

with L the length of the domain. Due to the high temperatures involved ($T_c = 2500$ [K]), the spectrum for the mixture was retrieved using the HITEMP 2010 database [28]. The spectra between 0 and 25000 [cm⁻¹] were then divided into 999 narrow-bands and 16 quadrature points per band for a total of 15984 pseudo-spectral points. As demonstrated in Fig. 5, also in case of a mixture, the present MC shows satisfactory agreement with the published literature.

The last validation case (14) consists of 1 [m³] cube with black walls filled with H₂O. The temperature profile is given by

$$3\text{dims: } T_m = 500 - 2000 \cdot (x \cdot y \cdot z)^2 + 2000 \cdot x \cdot y \cdot z \text{ [K]}, \quad T_w = 500 \text{ [K]}. \quad (22)$$

Fig. 6 shows the results of the 3D non-grey case compared to a line-by-line version of the current MC. The left contour shows the temperature at $z = 0.5$ [m], while the right plot shows the comparison of the results obtained with the narrow-band correlated- k method and the line-by-line benchmark at the same location (shown in [kW/m²]). The solution is again in excellent agreement with the line-by-line benchmark case.

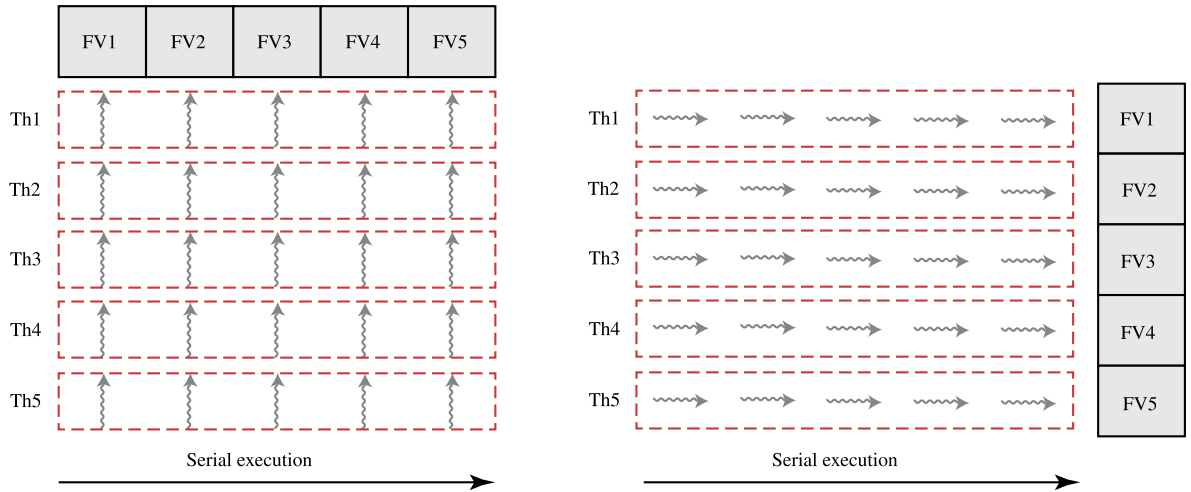


Fig. 7. Schematic showing the concept of ray parallelization (left) and domain parallelization (right). This simple example is composed of five control volumes (FV1 – FV5), five rays per cell (grey lines) and five threads working in parallel (Th1 – Th5).

In the following sections, the 1D H₂O parallel slab case with parabolic temperature profile (case 10) will be used as a test case to compare the computational performances of the different implementations. As stated before, although the case is 1D in nature, it is calculated on a 3D grid with two periodic directions to mimic the computations for a DNS of a fully developed turbulent channel flow.

3. GPU implementation

Graphical processing units have an architecture that, differently from CPUs, promote compute bound, highly parallelizable algorithms. The smallest parallel GPU units, called threads, run concurrently and are organized in thread blocks. All blocks can read and write into a global memory. The global memory is the “main” memory of the GPU, comparable to the heap in a C program, and has the slowest I/O access. Threads are grouped into groups of 32, termed “warps”, which are executed by a single scheduling unit and thus follow a Single Instruction Multiple Thread (SIMT) execution model. Hence, all threads belonging to a particular warp execute the same instruction simultaneously. Due to these features, the objective of porting an application from CPU to GPU, is to increase parallelization to favour the SIMT execution. A further level of parallelization is obtained by using “streams”. With this GPU feature, a device function, called “kernel” can be subdivided into parallel streams that run concurrently and independently, i.e. in a Multiple Instruction Multiple Data (MIMD) fashion, similar to multicore CPU computation (MPI parallelization). In compute bound problems, the use of streams is always recommended, since parallel MIMD execution is preferred to SIMT execution due to the absence of branch divergence (see section 4.2).

There are two main approaches to parallelize a radiative Monte Carlo algorithm on a GPU. Consider an example with a computational domain of five finite volumes (FV), each one sending five rays to march through the overall domain and five threads (Th) that can execute the marching of the rays. A schematic of this configuration is outlined in Fig. 7. The algorithm can then be parallelized by either ray parallelization or domain parallelization, which are outlined in more detail below.

In the first approach, each thread calculates one ray per finite volume. In this case, within each thread the finite volumes are executed in serial, while the rays per cell are parallelized. The solution will, therefore, be obtained by adding the partial results of each thread. The drawback of this approach is the continuous use of expensive atomic reductions (different threads have to read/write in the same memory location). On the other hand, if an ERMIC formulation is employed, it is possible to use the second approach, which consists in having a single thread calculate all the rays belonging to an individual finite volume. This is possible due to the fact that, in a reciprocal formulation, the only information required to calculate the radiative source in a point are the rays leaving the latter. In the schematics of Fig. 7, the ray decomposition and the domain decomposition approaches are displayed on the left and on the right, respectively. It is important to note that it is also possible to combine the two methods by exploiting the block/thread arrangement. Namely, divide the domain through different blocks and implement a ray parallelization for the finite volumes contained by the block. This configuration would speed up the necessary atomic reductions by enabling the use of shared memory that can be accessed by the whole block. However, shared memory is limited in size and for this reason this approach cannot scale efficiently to larger grids. Given these reasons, we conclude that the domain parallelization approach is more suitable for coupling a GPU Monte Carlo code with DNS.

Algorithm 2 displays the GPU implementation of the ERMIC based on domain parallelization. The implementation closely resembles the one displayed in Algorithm 1, with the difference that the routine now consists of two different GPU functions (or kernels) highlighted in light blue. The first one is in charge of initiating the calculation on the GPU, which immediately returns the control to the CPU, while the second routine retrieves the results. This approach enables a completely asyn-

Algorithm 2 ERMCM GPU implementation.

```

1: __device__ solution[stream_max][Ncell/stream_max]                                     ▷ global device variable
2: cudaMemcpyAsync(Temperature T, absCoeff κ, Grid, CopyFromCPUtoGPU)                 ▷ memory copy to device (GPU)
3: for s=0; s<stream_max do                                                         ▷ loop over streams
4:   procedure KICKOFF(thread t, block b, stream s)                                ▷ First kernel for stream number s
5:     Shared state = cuRandInit                                                    ▷ cuRand variable in shared memory
6:     tid ← threadIdx.x + blockIdx.x × blockDim.x
7:     for idx=tid; idx < Ncells; idx = idx + blockDim.x × gridDim.x do             ▷ Grid-stride loop over the GPU grid structure
8:       cell.ind.i ← idx/(kmax × jmax) + 1 + s × imax/stream_max                  ▷ Mapping thread index to mesh
9:       cell.ind.j ← idx/kmax + 1 - (cell.ind.i - 1 - s × imax/stream_max) × jmax
10:      cell.ind.k ← idx - kmax × (cell.ind.j - 1 + (cell.ind.i - 1 - s × imax/stream_max) × jmax + 1)
11:      QE ← 4κP(Tmax)σ Tmax4/numberOfRays
12:      for ray in Rays do
13:        procedure INITIALIZE
14:          Rθ, Rφ, Rn, Rg ← cuRand(Uniform distribution, state)                ▷ As in the CPU algorithm, but with cuRand instead
15:          Lines 6 – 14 in Algorithm 1
16:        end procedure
17:        procedure MARCH
18:          Lines 17 – 23 in Algorithm 1
19:          solution[s][idx] ← solution[s][idx] - Absorption                      ▷ device global variable that allows asynchronous computations
20:          Lines 25 – 30 in Algorithm 1
21:        end procedure
22:      end for
23:    end for
24:  end procedure
25: end for
26: Perform other tasks
27: for s=0; s<stream_max do                                                         ▷ loop over streams
28:   procedure RETURN(thread t, block b, stream s)                                ▷ Second kernel for stream number s
29:     tid ← threadIdx.x + blockIdx.x × blockDim.x
30:     for idx=tid; idx < Ncells; idx = idx + blockDim.x × gridDim.x do
31:       QR[idx] ← solution[s][idx]
32:     end for
33:   end procedure
34:   cudaMemcpyAsync(Solution QR, CopyFromGPUtoCPU)                                ▷ memory copy to host (CPU)
35:   cudaDeviceReset()                                                            ▷ clear device memory allocations
36: end for

```

Table 2

Comparison between standard CPU and GPU implementation.

grid size	16 ³	32 ³	48 ³	64 ³	96 ³	128 ³	160 ³
CPU	269.4 s	2921.1 s	13182.7 s	39313.3 s	271844.4 s	(920183.3) s	(2452230.3) s
GPU	11.8 s	84.8 s	394.0 s	1169.5 s	6143 s	19623 s	47539 s
Speedup	22.8×	34.4×	33.5×	33.6×	44.3×	(46.9)×	(51.6)×
rays per cell	6 · 10 ²	1.5 · 10 ³	6 · 10 ³	1.5 · 10 ⁴	3 · 10 ⁴	6 · 10 ⁴	1.5 · 10 ⁵
CPU	369.7 s	961.7 s	3928.6 s	10432.2 s	19661.1 s	39313.3 s	132641.3 s
GPU	14.1 s	31.8 s	119.2 s	294.4 s	585.8 s	1170 s	2916 s
Speedup	26.2×	30.2×	33.0×	35.4×	33.6×	33.6×	45.5×

chronous computation of the GPU and allows to perform other tasks on the CPU (line 26) that would otherwise remain idle. Each kernel is executed `stream_max` times and computes $(1/\text{stream_max})$ th of the domain. The stream loops (lines 3 and 27) contain only non-blocking statements that enable a parallel stream execution. The core of the domain parallelization consists in mapping the thread index to a specific finite volume (lines 8–10). The for-loop over the computational cells is then replaced by a GPU-grid-stride loop that runs over the thread index (line 7) and covers all cells in the domain. The random number generation is performed on-the-fly by employing the CUDA library `cuRand`. The solution is stored in a global device variable `solution`, which is then retrieved by the second kernel once the computations are complete.

The GPU implementation is tested for case 10 (see Table 1, plane parallel slab of 1 [atm] H₂O with parabolic temperature profile) on a Tesla K40M. The execution speed is benchmarked against the CPU implementation executed on an Intel Xeon E5-2680 @ 2.40GHz. Table 2 shows the computational time required as a function of mesh size and number of rays per cell. In all the test cases, the maximum allowed number of streams (16) is used, while the number of blocks and threads per block are calculated such that the GPU resources are fully utilized. The default values for the parameters that are not varied are a grid size of 64³ and 6 · 10⁴ rays per cell. The results in Table 2 show that the speedup obtained with a straightforward GPU implementation, using domain decomposition, is already relatively high. Nonetheless, with the increase of problem size, the speedup does not show a satisfying improvement, reaching values of around ~50×. This apparent limit is caused by the finite resources of the GPU. Being a compute bound algorithm, the scarce resource is the amount of registers per thread that sets the maximum number of threads running concurrently. If the number of registers is increased, the scheduling units serialize the execution of the exceeding warps. As a consequence, no further gain is observed when increasing the

Table 3

Execution time with classical versus textured memory approach.

grid size	16^3	32^3	48^3	64^3	96^3	128^3	160^3
classic	11.8 s	84.8 s	394.0 s	1169.5 s	6143 s	19623 s	47539 s
texture	8.5 s	48.6 s	260.0 s	716.0 s	4381 s	12343 s	34047 s
Speedup	1.38×	1.74×	1.52×	1.64×	1.40×	1.59×	1.40×
rays per cell	$6 \cdot 10^2$	$1.5 \cdot 10^3$	$6 \cdot 10^3$	$1.5 \cdot 10^4$	$3 \cdot 10^4$	$6 \cdot 10^4$	$1.5 \cdot 10^5$
classic	14.1 s	31.8 s	119.2 s	294.4 s	585.7 s	1170 s	2916 s
texture	9.6 s	20.6 s	74.3 s	180.7 s	358.6 s	714.5 s	1784 s
Speedup	1.47×	1.54×	1.60×	1.63×	1.63×	1.64×	1.64×

mesh size or the number of rays per cells. Note that the values in parenthesis for the CPU execution time in Table 2 are extrapolated from the scaling of the other results and, as such represent an estimation only.

4. Algorithm acceleration

A naive GPU implementation, as demonstrated in the section above, is useful to provide a certain level of speedup, but is certainly not enough to address the computational requirements of a DNS simulation. In particular, the main problems and bottlenecks of such an algorithm are the slow memory access and the large inactivity of the threads due to the SIMT execution model. For this reason, we will address these issues by implementing acceleration techniques that will significantly reduce the execution time and thus enable a full coupling between DNS and the GPU Monte Carlo code.

4.1. Texture memory

Due to the GPU architecture, memory input and output is heavily affected by the access pattern of the threads. In particular, the global memory of a GPU is optimized for coalesced access. A coalesced memory transaction is one in which all of the threads in a half-warp access global memory at the same time. That is to say, consecutive threads should access consecutive memory addresses in the global memory to obtain efficient memory loads/stores. To avoid penalties associated with uncoalesced transactions, it is possible to store variables in registers (the memory associated with the single thread) or shared memory, which is fast-access memory common to all threads in a block. Unfortunately, these two memory types are severely limited in size (on a tesla K40M shared memory consists of only 49 kB per multiprocessor for a total of ~ 735 kB). Therefore, after all the fast memory resources have been depleted, it is necessary to store the bulk of the variables in the global memory. Since most memory fetches depend on the drawing of random numbers, it is not straightforward to predict which address consecutive threads might access. As a consequence, coalesced memory transactions are impossible to achieve in a Monte Carlo simulation. An easy way to optimize memory input and output is hence to employ texture memory. Texture memory is a type of read-only memory, which has been developed for graphical applications. Instead of storing variables linearly, as global memory does, texture memory is designed to optimize the spatial locality of memory access. In other words, each point is associated to a coordinate, and the most efficient memory fetch occurs when consecutive threads access adjacent coordinates in the texture memory instead of consecutive addresses. This scenario is much more likely in a domain parallelized Monte Carlo simulation. The input values to access a texture memory location are float coordinates, while the value returned from the memory is a linear (or trilinear in case of a 3D texture) interpolation of the adjacent values. This feature is extremely useful as it provides fast linear interpolation, which is repeatedly required in a spectral MC code (lines 12, 13, 14, 20 and 22 in Algorithm 1).

Variables that were residing in the global memory (temperature, blackbody intensity and absorption coefficient), are therefore relocated to the texture memory. The results of the texture memory implementation are shown in Table 3 in comparison to the standard GPU implementation. The use of texture memory results in a computational gain for all the different settings. Nonetheless, the speedup tends to decrease with mesh size. This behaviour could be caused by the reduced spatial locality of memory access for contiguous threads on a finer grid (i.e. the ray travels further, distancing itself from the aligned source cells). On the other hand, the speedup increases if more memory transactions are performed (i.e., increasing the numbers of rays per cell). While this implementation has been performed for a structured mesh, the benefits will be equivalent if the same warp behaviour is ensured on an unstructured grid (adjacent source point for consecutive threads).

4.2. Narrow band sorting

The SIMT execution model can lead to a severe performance loss, known as “branch divergence”. A warp executes one common instruction at a time, so the threads must wait until the execution is terminated for the entire warp before proceeding. For this reason, full efficiency is realized when all 32 threads of a warp follow the same execution path (i.e., execute the same instruction at the same time). If threads of a warp diverge due to a data-dependent conditional branch (if or while statements), the warp executes all the paths entirely, disabling threads that are not on that path. For the

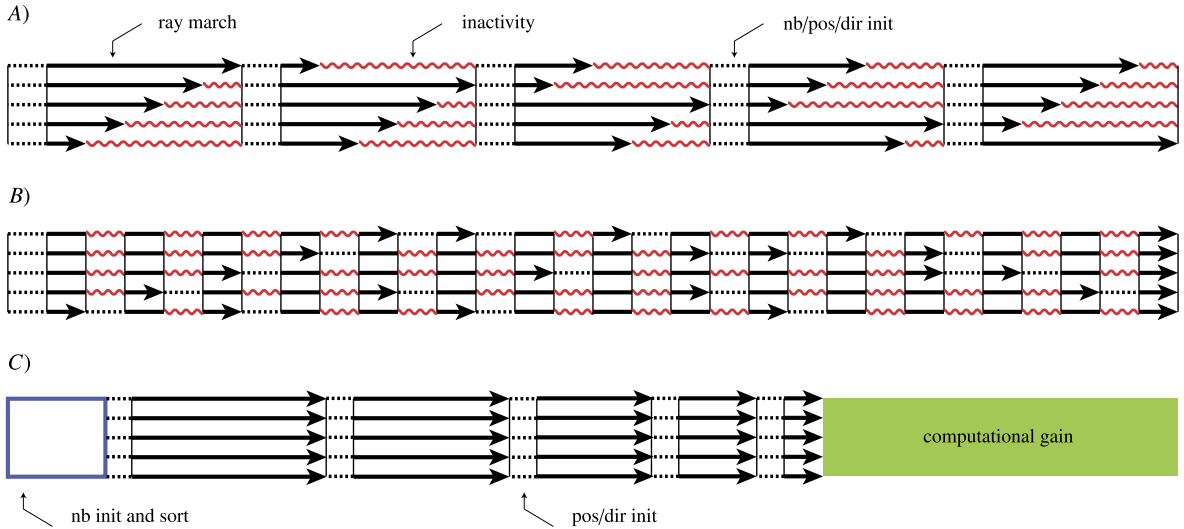


Fig. 8. Example of a marching procedure for different GPU MC schemes. (A): standard MC implementation; (B): reinitialization MC; (C): sorting MC. The different rows represent the sequential execution of different threads in a warp. We show here only 5 threads and 5 rays to simplify the scheme, but in reality there are 32 threads in a warp and tens of thousand rays per thread. Note that the length of the arrows and the dashed lines (representing marching and initialization) are always preserved among the three schemes. On the other hand, the position and direction initialization time (dashed lines) is shorter in the last scheme (C), since the wavelength has already been chosen in the preprocessing step (blue box).

Table 4

Speedup using the narrow band sorting. The values of the speedup are referred to the textured execution times of Table 3.

grid size	16^3	32^3	48^3	64^3	96^3	128^3	160^3
Time	6.4 s	43.2 s	180.8 s	573.0 s	2866 s	9597 s	22189 s
Speedup	1.32×	1.13×	1.44×	1.25×	1.52×	1.29×	1.53×
rays per cell	$6 \cdot 10^2$	$1.5 \cdot 10^3$	$6 \cdot 10^3$	$1.5 \cdot 10^4$	$3 \cdot 10^4$	$6 \cdot 10^4$	$1.5 \cdot 10^5$
Time	10.3 s	21.7 s	73.2 s	165.5 s	308.0 s	573.2 s	1313 s
Speedup	0.93×	0.95×	1.01×	1.1×	1.16×	1.25×	1.36×

purpose of correctness, the SIMT execution model can be essentially ignored. However, in terms of code efficiency, thread divergence is a serious issue and has to be addressed, particularly in the case of a MC simulation, where the abundant while loops and if conditions cause large thread inactivity.

A simple and straightforward approach to reduce inactivity, would be to re-initialize the ray whenever a marching is terminated within the warp. On the other hand, re-initializing the ray on a particular thread forces to temporarily disable the threads that have not yet completed the marching, serializing the initialization procedure. As a consequence, the execution time of multiple initializations might become longer than the benefit obtained by the lower inactivity during the marching procedure.

Taking into account the properties of the ray, leads to a more effective solution. For example, when two different threads in the same warp are marching rays with different wavelength, they handle different absorption coefficients. The ray with a higher κ_v will complete the marching quicker than the one with lower κ_v , due to the shorter path length. Since a Monte Carlo routine requires random draws of the wavelength based on a probability distribution function, it is a common scenario that threads are handling absorption coefficients of different order of magnitude. Due to the SIMT execution model, the time required for the warp to complete the current ray tracing is dictated by the thread with the lowest absorption coefficient. It is therefore beneficial to have threads handling absorption coefficient of similar value at all times, such that the tracing might complete simultaneously. To achieve this, it is necessary to precompute all wavelengths for each ray in each finite volume and sort them based on the magnitude of κ_v . Consequently, threads will always march rays from the lowest to the highest κ_v . While these values might be slightly different for different threads, the order of magnitude of κ_v will be similar, thus significantly reducing the branch divergence of the warp.

The different configurations are outlined in Fig. 8. The first scheme is a standard MC that does not account for any branch divergence reduction technique. Scheme B shows a re-initialization scheme in which, wherever a thread in the warp completes the marching, the ray is immediately re-initialized. It is clear that this scheme is successful only if the cost of initializing a ray is smaller than the tracing of the shortest ray. This is not the case in a medium with a high absorption, where rays can be terminated within 5 steps. Scheme C shows the advantage of reordering the rays based on their absorption coefficient which aligns the ray marching executions.

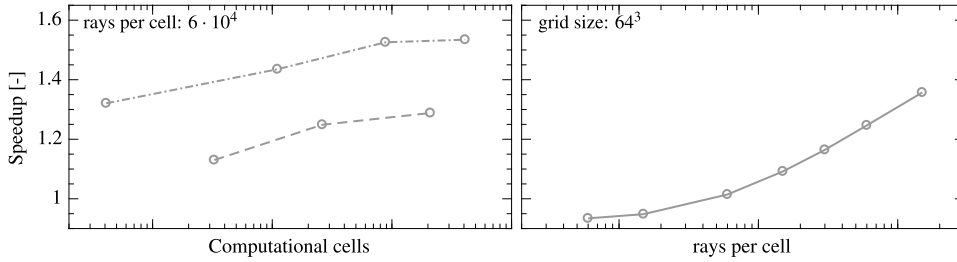


Fig. 9. Speedup obtained with the narrow-band sorting technique. In the figure on the left, the dashed line connects the points characterized by a mesh which is $2^n \cdot \text{streams}$ (16³, 32³ and 64³), while the dashed dotted line all the other points.

The results of the tests for a narrow-band sorting algorithm are shown in Table 4 and Fig. 9. The speedup obtained with sorting the narrow bands is larger when the grid is not $2^n \cdot \text{streams}$ (32, 64, 128). This is caused by an inefficient mapping of the grid onto the device resources, which in this case are powers of 2. Indeed, by sorting the narrow bands, it is possible to correct the penalties associated with an inadequate mapping. It is possible to notice that the speedup increases with increasing the number of mesh points, until it reaches a plateau for large mesh sizes. On the other hand, if the number of rays per cell are too small, the advantage of a lower warp inactivity is overshadowed by the cost of the sorting procedure. Contrarily, increasing the number of rays per cell leads to an linear growth of the speedup, since the warp inactivity is efficiently replaced by the ray marching computation.

It is interesting to notice the difference between the speedup of the narrow-band sorting scheme with respect to mesh size and the speedup using a texture memory approach only. While the first one increases, the latter decreases with grid size. This difference shows the interplay between memory transactions and computations as the mesh size increases, highlighting the larger relative importance of compute statements with increasing mesh size.

4.3. Multigrid

The radiative intensity is absorbed exponentially as function of the absorption coefficient and the travelled distance. Therefore, the intensity absorbed by traversing a cube of size Δx^3 will be roughly proportional to

$$I_{abs} \sim (1 - \exp(-\kappa_v C \Delta x)) , \quad (23)$$

with C a constant depending on orientation. Consequently, the intensity of the ray leaving the cell is

$$I_{out} = I_{in} - I_{abs} \sim \exp(-\kappa_v C \Delta x) , \quad (24)$$

which signifies that, for a low κ_v , the intensity gradient of the propagating ray will be mild and the required cell size Δx can be relatively large. Vice versa, if κ_v is large, a lower Δx is necessary to capture the steep intensity gradient. If an adequate Δx is chosen as a pre-processing step (as it could be done in a grey gas medium) the mesh will be over-resolved for the rays with low absorption, resulting in an inefficient ray tracing. Nevertheless, since a high κ_v ray will be terminated fairly quickly, it requires a high resolution only on a small zone around the source point. On the contrary, a ray with low absorption will propagate far into the domain. By combining these two features of rays with different absorption coefficient, it is possible to construct a mesh strategy that optimizes the ray tracing, while retaining a high accuracy. The

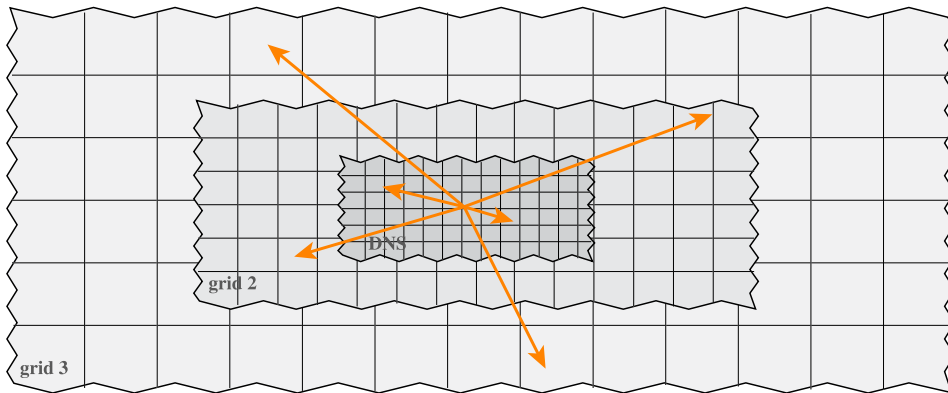


Fig. 10. Schematic showing the concept of the mesh coarsening scheme. The orange lines symbolize the marched rays. Several grids are overlayed one on top of each other. The ray falls onto the coarsened mesh when it reaches the maximum number of steps in the current grid. The concept is shown here in two dimensions for simplicity.

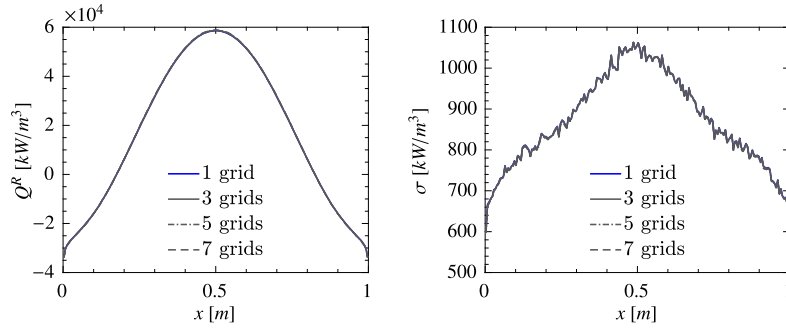


Fig. 11. Comparison of the results for the parabolic H₂O case for different numbers of overlayed coarsened mesh used. 5 steps are allowed in each mesh. a): Radiative heat source, b): standard deviation.

Table 5

Absorbed fraction (p) in seven successive grids with 5 steps per grid. In case less than seven grids are used, the last grid absorbs the remaining intensity. For example, if three grids are used for $\kappa_{\text{ref}} = \kappa_p$, the first grid absorbs $p_1 = 12.4\%$, the second $p_2 = 20.3\%$, while the third the remaining intensity $p_3 = (100 - 12.4 - 20.3)\%$.

	κ_{ref}	p_1	p_2	p_3	p_4	p_5	p_6	p_7
κ_p	(5.07 [m ⁻¹])	12.4%	20.3%	27.6%	25.9%	12.1%	1.6%	0.1%
$\bar{\kappa}_{\text{max}}$	(98.1 [m ⁻¹])	92.22%	7.73%	0.05%	0.0%	0.0%	0.0%	0.0%
$\bar{\kappa}_n$	(8.80 [m ⁻¹])	20.5%	29.2%	30.2%	16.9%	3.1%	0.09%	0.01%

objective is to have a grid that is fine close to the starting cell and gradually coarser as the ray travels further away from the initial point. To obtain this effect, it is possible to overlay several meshes characterized by different cell sizes. The temperature values will be interpolated on the coarser meshes from the DNS solution which represents the finest mesh level (radiative heat transfer does not introduce new spatial wavenumbers, so the smallest radiative length scales are as small as the Batchelor scales). For all finite volumes, the ray tracing commences on the DNS mesh and the ray is allowed to step onto the current mesh a fixed number of times. If the ray is not exhausted, it falls into a coarser mesh and so forth, until the last mesh is reached. The last (and coarsest) mesh will trace the ray until depletion. The only added overhead is the cost of the interpolation onto coarser meshes, which is completely irrelevant compared to the gain in computational speed obtained. A similar method, involving patches of interest, was previously implemented by Humphrey et al. in two different occasions. Namely, in a parallel CPU Monte Carlo implementation [31] and in a grey gas GPU implementation [21]. They used this technique to reduce computation and communication time. On the other hand, we highlight the additional benefits that such a method has in a non-grey GPU implementation, other than the reduction of computations. Here, due to the SIMT execution model, threads characterized by higher absorption rays will remain inactive, while waiting for the lower absorption rays to complete the computations. In this case, it is possible to tailor the method towards a pure reduction of thread inactivity by targeting the resolution of the higher impact, high absorption, rays. If this is done, the solution retains its accuracy and the parallel efficiency is greatly enhanced.

To optimize the number of steps in subsequent meshes, it is possible to rely on the transmissivity of the grids by specifying a fraction of intensity absorbed in the n th grid (p_n). As a consequence, the number of steps in the n th grid (s_n) will be

$$s_n = \text{int} \left(\frac{\ln \left(1 - \sum_{m=1}^{n-1} p_m \right)}{\kappa_{\text{ref}} \Delta x_n} - \frac{\ln \left(1 - \sum_{m=1}^n p_m \right)}{\kappa_{\text{ref}} \Delta x_n} \right). \quad (25)$$

Here, κ_{ref} is a reference absorption coefficient and can be specified in different ways. Using the Planck-mean absorption coefficient ($\kappa_{\text{ref}} = \kappa_p$) is the safest option to ensure unaltered accuracy. If execution speed has to be minimized, the largest average narrow-band absorption coefficient in the spectrum ($\kappa_{\text{ref}} = \bar{\kappa}_{\text{max}}$) can be employed. On the other hand, a more balanced choice is the average absorption coefficient of the most influential band ($\kappa_{\text{ref}} = \bar{\kappa}_n$) calculated as

$$\bar{\kappa}_n = \sum_{g=1}^{N_g} \omega_g k_{n,g}, \quad \text{where } n \text{ is the band containing } \bar{\nu} = \sum_{n'=1}^{N_b} \nu_{c,n'} \cdot f'_n, \quad (26)$$

Table 6

Speedup using multiple overlayed grids. 5 steps per grid.

grid number	1	2	3	4	5	6	7
Speedup	1×	1.4×	2.6×	4.2×	5.8×	6.6×	7.1×

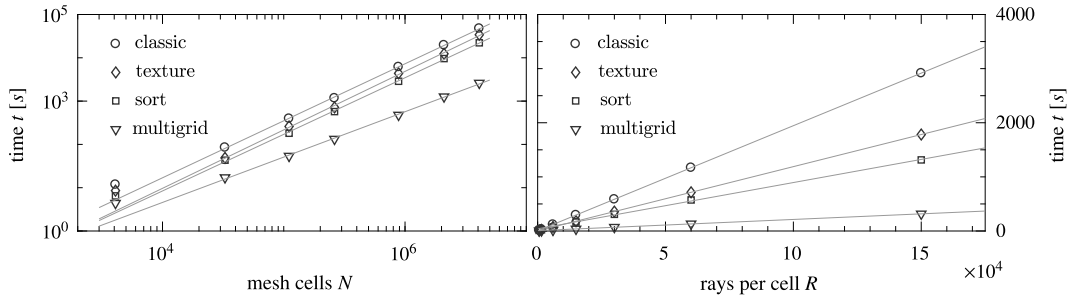


Fig. 12. Scaling of the code with grid cells and rays per cell.

and $\nu_{c,n'}$ is the central wavenumber of band n' . Since the wavenumber is sampled from the maximum temperature of the considered system, κ_p , $\bar{\kappa}_{max}$ and $\bar{\kappa}_n$ can be calculated from the maximum temperature spectrum only.

In practice, the correct fractions p_n that leave the results unaltered are not known a-priori. Therefore, to ensure an optimal number of steps, a trial and error procedure has to be employed until the results (both in terms of Q^R and σ) are method-independent. It is possible then to assess the implementation of the grid coarsening by calculating the fractions p_n from equation (25) and use these for future implementations. In addition, more conservative settings (larger p_n in finer grids) have to be adopted in case of steeper temperature gradients and inhomogeneous cases.

Fig. 10 shows a 2D representation of the mesh coarsening concept, while Fig. 11 shows the solutions of the test case employing the multigrid technique. In particular, the results shown in Fig. 11 have been obtained with a maximum of 7 overlayed grids corresponding to $192^3 \rightarrow 96^3 \rightarrow 48^3 \rightarrow 24^3 \rightarrow 12^3 \rightarrow 6^3 \rightarrow 3^3$. The rays were allowed to travel a maximum of 5 steps in each grid, while proceeding until termination on the last one. The absorbed fraction for different choices of κ_{ref} , based on equation (25), is shown in Table 5. When using κ_p or $\bar{\kappa}_n$ as a reference, the largest absorption occurs in the third grid. On the other hand, the rays represented by $\bar{\kappa}_{max}$ are absorbed almost completely within the first grid.

As shown in Fig. 11, the results of the test case are unaffected by the grid coarsening technique, both in terms of radiative power source and its standard deviation. The speedup obtained, defined as t_1/t_n , where t_1 is the time required for completing the calculation with one grid while t_n with using n grids, is shown in Table 6. Despite the negligible influence in the results, by employing the multigrid technique, it is possible to reduce the computational cost by a factor that is roughly equal to the number of grids used. The reduction of thread inactivity using this method is clearly visible by comparing Table 5 and Table 6. By changing the implementation from 4 to 7 grids, the intensity fraction that is computed on the coarser grids is only 13.9%, 0.0% and 3.2% for κ_p , $\bar{\kappa}_{max}$ and $\bar{\kappa}_n$, respectively. Nonetheless, the speedup increases by almost a factor of 2.

5. Overall performance increase

An overview of the scaling performance using different acceleration techniques is given in Fig. 12 for varying problem sizes. Note that the implementations are additive (i.e., sorting employs texture memory allocations and multigrid performs also a narrow-band sorting). A coarsening ratio of 2 has been employed for successive grids in the multigrid implementation. The smallest allowed mesh had a size of 3^3 , resulting in 3 grids for 16^3 , 4 for 32^3 , 5 for 48^3 and 64^3 , and 6 grids for 96^3 , 128^3 and 160^3 . Again, only 5 steps were allowed in each level. The scaling of all implementations is well described by power functions of mesh cells N and linear functions of the number of rays R . The grey lines depicted in Fig. 12 take the following form

- classic $t \propto N^{1.32}$, $t \propto 0.98R$,
- texture $t \propto N^{1.35}$, $t \propto 0.96R$,
- sorting $t \propto N^{1.31}$, $t \propto 0.7R$,
- multigrid $t \propto N^{1.05}$, $t \propto 0.7R$,

Table 7

Comparison between standard CPU implementation and optimized GPU implementation.

grid size	16^3	32^3	48^3	64^3	96^3	128^3	160^3
CPU	269.4 s	2921.1 s	13182.7 s	39313.3 s	271844.4 s	(920183.3) s	(2452230.3) s
GPU	4.4 s	17.1 s	53.7 s	132.8 s	476.6 s	1262 s	2612 s
Speedup	61.2×	170.8×	245.5×	296.0×	570.4×	(729.1)×	(938.8)×
rays per cell	$6 \cdot 10^2$	$1.5 \cdot 10^3$	$6 \cdot 10^3$	$1.5 \cdot 10^4$	$3 \cdot 10^4$	$6 \cdot 10^4$	$1.5 \cdot 10^5$
CPU	369.7 s	961.7 s	3928.6 s	10432.2 s	19661.1 s	39313.3 s	132641.3 s
GPU	4.1 s	6.3 s	17.1 s	37.5 s	69.9 s	132.4 s	316.0 s
Speedup	90.2×	152.7×	229.7×	278.2×	281.3×	296.9×	419.8×

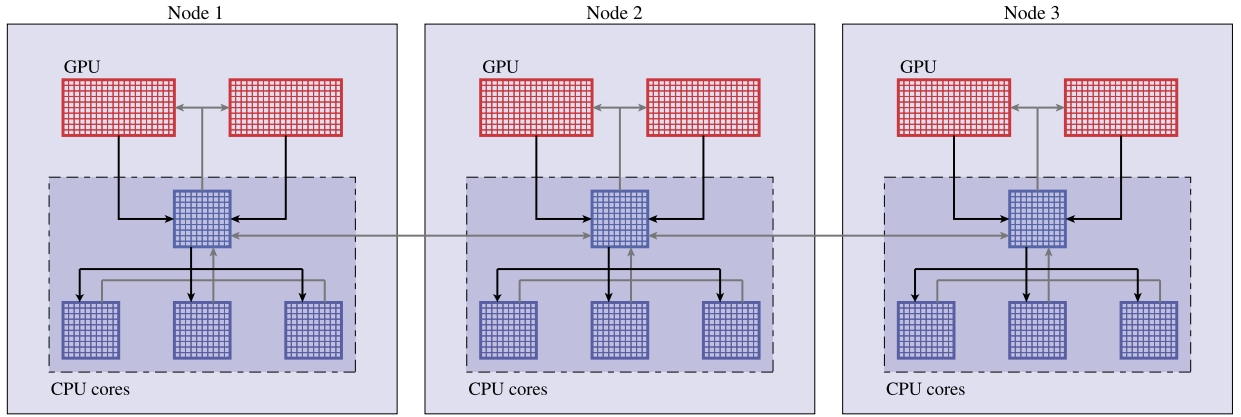


Fig. 13. Schematic representing the multi GPU implementation. The domain is decomposed on different nodes. In each node one CPU core communicates to the GPUs the entire temperature domain and returns the computed radiative heat source to the CPUs within the node.

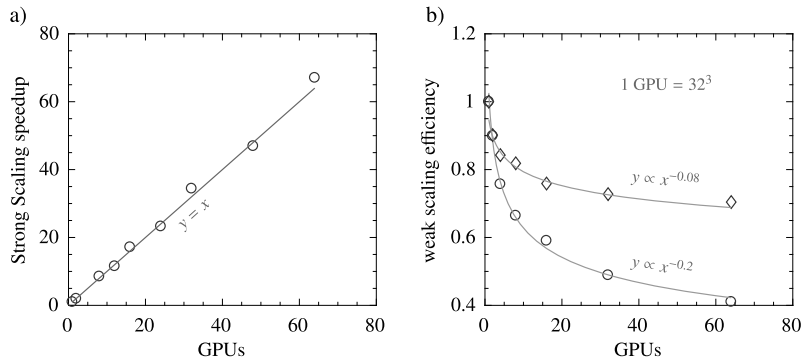


Fig. 14. Performance of a multi GPU implementation. a): strong scaling speedup (t_1/t_n). b): weak scaling efficiency (t_1/t_n). Circles are results obtained without using multigrid and diamonds employing the multigrid technique. t_1 is the execution time on 1 GPU and t_n is the execution time on n GPUs.

While a texture memory allocation has large benefits for the investigated cases, the computational gain is bound to decrease when the grid size increases (as seen in section 4.1) as given by the larger exponent when compared to the classic implementation ($1.35 > 1.32$). On the contrary, with a multigrid scheme it is possible to obtain a quasi-linear scaling Monte Carlo code with mesh size (exponent ≈ 1). Moreover, the narrow band sorting procedure allows a scaling greater than ideal with respect to the rays per cell ($0.7 \cdot R$). With more rays being launched, the drawn absorption coefficients fill the whole spectrum space efficiently, replacing the inactivity by aligning more effectively the thread marching.

It is demonstrated that, by employing these optimization techniques, it is possible not only to reduce the computational time, but also to significantly improve the scaling of the code with problem size. The performances of the optimized GPU Monte Carlo code, compared to a serial CPU Monte Carlo implementation executed on an Intel Xeon E5-2680 @ 2.40GHz, is shown in Table 7. It has to be reminded that, while texture memory allocation and narrow band sorting only improve computational speed on a GPU, multigrid, although less effective, can be also implemented for a code that runs on a CPU, leading to an increase of code efficiency. The maximum speedup achieved was $570.4\times$ for a grid size of 96^3 . For the largest problems, the CPU computational time was estimated from the scaling. Based on this estimation, we expect a impressive increase of speedup, differently from what is observed in Table 2 (potentially we could achieve $938.8\times$ for a 160^3 grid).

6. Multi GPU and DNS coupling

The optimized GPU Monte Carlo version can be used to efficiently couple radiative heat transfer with a DNS code in order to study the interactions between radiative heat transfer and turbulent mixing. The coupling is implemented with the use of MPI libraries that handle communications between CPU cores. Each node has a master core which communicates with the available GPUs on the node. Thanks to the reciprocal formulation, the GPUs calculate the radiative source term only on the domain handled by the associated node. On the other hand, to perform ray tracing and to avoid boundary communication, all GPUs require the complete temperature field. A schematic of the multi GPU implementation is shown in Fig. 13. The grey arrows show the communication of the temperature field, while the black arrows show the path of the computed Q^R . The memory transfer to and from the GPU is completely asynchronous, such that the CPUs proceed to calculate additional

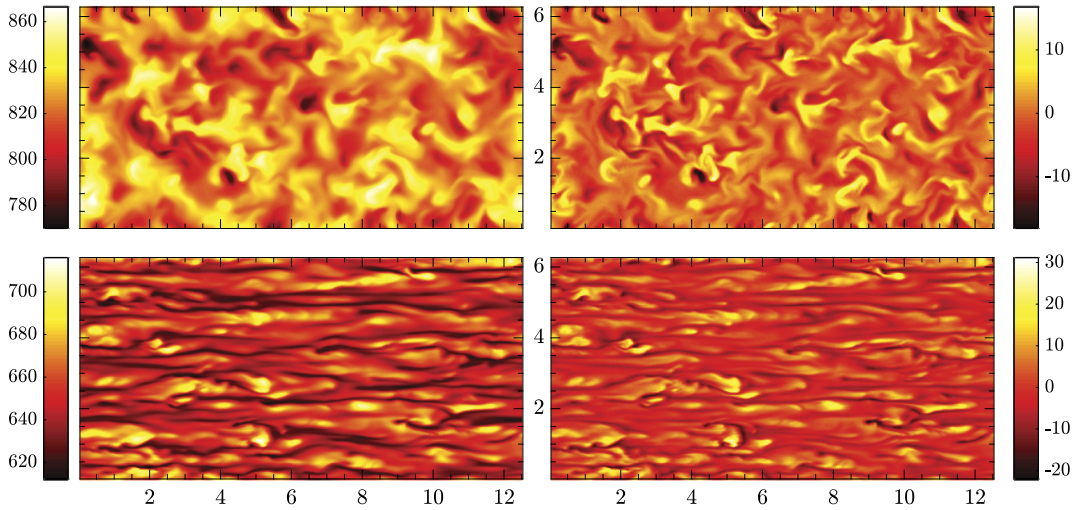


Fig. 15. Instantaneous snapshots on a wall-parallel plane ($x-z$) at $y = 1.1$ (top) and $y = 1.97$ (bottom). Left: temperature T [K]. Right: Radiative power Q^R [kW/m³].

fluid time steps, while the GPUs compute the radiative heat source. As a consequence the CPU computation is completely hidden by the radiative power calculation.

The code has been tested on the Cartesius' cluster located in Amsterdam, The Netherlands, on the accelerator island composed of 60 nodes containing 2 Tesla K40M each. The scaling of the code was examined up to 64 GPUs. The results are shown in Fig. 14. The strong scaling of the code is calculated by keeping the grid size constant (192^3 in this case) and increasing the number of GPUs. The quantity shown in Fig. 14(a) is the time required for one time step to complete on 1 GPU over the time required for N GPUs. As expected by the computational nature of the code, the scaling is almost ideal. Moreover, Fig. 14(b) shows the weak scaling efficiency, tested with and without the use of the multigrid scheme. In this case the grid size is increased proportionally to the number of GPUs used, with one GPU always computing on a 32^3 mesh. Since the problem size increases with the number of GPUs used, the code greatly benefits from the multigrid scheme, which improves the weak scaling efficiency from $\propto GPU^{-0.2}$ to $\propto GPU^{-0.08}$.

To prove the level of accuracy achievable in an acceptable time span, the radiative power is calculated for a turbulent temperature field obtained from a DNS. The DNS represents a fully developed turbulent channel flow with a bulk Reynolds number of $Re = 7500$ and isothermal walls at 955 and 573 [K] at the bottom and top, respectively. The flow is periodic in the streamwise and spanwise directions. The radiative properties of the medium are those of water vapour at 1 [atm]. The Planck mean absorption coefficient varies roughly from 5.5 [m⁻¹] near the hot wall to 15 [m⁻¹] near the cold wall and, therefore, can be considered optically thick. In such conditions, the turbulent radiative power spectrum is characterized by short length scales, comparable to the largest wavenumbers of the turbulent temperature spectrum. Therefore, the radiative heat source requires to be accurate on the full DNS mesh. The mesh is composed of 192^3 elements, while the box dimensions are 2 , 2π and 4π [m] in the wall normal (y), span wise (z) and stream wise (x) directions, respectively. $6 \cdot 10^4$ rays per cells were used to calculate the radiative power. Snapshots of the radiative field are shown in Fig. 15. The left contours show the temperature field (in [K]), while the contours on the right are the calculated radiative power in [kW/m³]. The top figures show the fields at a y location of 1.1 [m] (roughly at the center of the channel), while the bottom figures show a wall normal plane located near the cold wall ($y \approx 1.97$ [m]). As seen from the figures, the radiative field is solved with a high accuracy, matching quite closely the turbulent structures of the temperature field as expected for a highly participating medium. In addition, as predicted in [8], in the center of the channel, the turbulent radiative field filters the large turbulent wavenumbers, due to the action of incident radiation acting on the isotropic temperature structures.

7. Conclusions

A reciprocal Monte Carlo formulation for radiative heat transfer calculation has been ported to GPU using NVIDIA programming language CUDA. The naive GPU implementation already showed a speedup of almost two order of magnitude compared to a classical CPU implementation. The efforts were focussed on improving the GPU implementation by overcoming the bottlenecks typical of a Monte Carlo code. In particular, the memory access has been enhanced by employing a texture memory for the storage of all read-only variables. This approach allows random memory access and speeds up the computation of the constantly required linear interpolations. Furthermore, the warp inactivity has been significantly reduced using a combination of narrow-band sorting procedures and a multigrid approach. Using this technique the accuracy of the MC solver is retained while the computational expenses are significantly reduced. Therefore, by solving these issues a speedup of up to 3 orders of magnitude when compared to the initial CPU implementation, was achieved. In addition

the scaling of the code with problem size (grid cells and rays per cells) was thoroughly studied, demonstrating that the optimized implementation shows a superior scaling when compared to the classical implementation. Indeed the speedup plateau noticed with the standard GPU implementation, is far from being reached even for the largest problems considered.

Moreover, a multi-GPU implementation was performed, showing an efficient strong and weak scaling up to 64 GPUs. While the strong scaling is ideal due to the computational nature of a MC code, the weak scaling benefits largely from the multigrid approach. The coupling with DNS shows the capability of achieving accurate results also for challenging problems as optically thick turbulent flows.

Declaration of Competing Interest

There is no competing interest.

References

- [1] M. Modest, *Radiative Heat Transfer*, 3rd edition, 2013.
- [2] A. Sakurai, K. Matsubara, K. Takakuwa, R. Kanbayashi, Radiation effects on mixed turbulent natural and forced convection in a horizontal channel using direct numerical simulation, *Int. J. Heat Mass Transf.* 55 (2012) 2539–2548, <https://doi.org/10.1016/j.ijheatmasstransfer.2012.01.006>.
- [3] P. Coelho, O. Teerling, D. Roekaerts, Spectral radiative effects and turbulence/radiation interaction in a non-luminous turbulent jet diffusion flame, *Combust. Flame* 133 (2003) 75–91.
- [4] K. Deshmukh, M. Modest, D. Haworth, Direct numerical simulation of turbulence-radiation interactions in a statistically one-dimensional nonpremixed system, *J. Quant. Spectrosc. Radiat. Transf.* 109 (2008) 2391–2400.
- [5] S. Ghosh, R. Friedrich, C. Stemmer, Contrasting turbulence-radiation interaction in supersonic channel and pipe flow, *Int. J. Heat Fluid Flow* 48 (2014) 24–34, <https://doi.org/10.1016/j.ijheatfluidflow.2014.04.002>.
- [6] S. Ghosh, R. Friedrich, Effects of radiative heat transfer on the turbulence structure in inert and reacting mixing layers, *Phys. Fluids* 27 (2015) 055107, <https://doi.org/10.1063/1.4920990>.
- [7] A. Gupta, M. Modest, D. Haworth, Large-eddy simulation of turbulence-radiation interactions in a turbulent planar channel flow, *J. Heat Transf.* 131 (6) (2009) 1–8, <https://doi.org/10.1115/1.3085875>.
- [8] S. Silvestri, A. Patel, D. Roekaerts, R. Pecnik, Turbulence radiation interaction in channel flow with various optical depths, *J. Fluid Mech.* 834 (2018) 359–384.
- [9] Y. Wu, M. Modest, D. Haworth, A high-order photon monte carlo method for radiative transfer in direct numerical simulation, *J. Comput. Phys.* 223 (2007) 898–922.
- [10] Y. Zhang, R. Vicquelin, O. Gicquel, J. Taine, Physical study of radiation effects on the boundary layer structure in a turbulent channel flow, *Int. J. Heat Mass Transf.* 61 (2013) 642–666.
- [11] R. Vicquelin, Y. Zhang, O. Gicquel, J. Taine, Effects of radiation in turbulent channel flow: analysis of coupled direct numerical simulations, *J. Fluid Mech.* 753 (2014) 360–401, <https://doi.org/10.1017/jfm.2014.368>.
- [12] J. Owens, D. Luebke, N. Govindaraju, M. Harris, J. Krüger, A. Lefohn, T. Purcell, A survey of general-purpose computation on graphics hardware, in: *State of the Art Report, Eurographics* (2005) 21–51.
- [13] NVIDIA Corporation, *NVIDIA CUDA C Programming Guide*, 2011.
- [14] Gpu machine learning, <http://www.nvidia.com/object/machine-learning.html>.
- [15] Gpu medical imaging, http://www.nvidia.com/object/medical_imaging.html.
- [16] M. Nobile, P. Cazzaniga, A. Tangherloni, D. Besozzi, Graphics processing units in bioinformatics, computational biology and systems biology, *Brief. Bioinform.* 18 (5) (2016) 870–885.
- [17] A. Khajeh-Saeed, J. Perot, Direct numerical simulation of turbulence using gpu accelerated supercomputers, *J. Comput. Phys.* 235 (2013) 241–257, <https://doi.org/10.1016/j.jcp.2012.10.050>.
- [18] F. Salvadore, M. Bernardini, M. Botti, Gpu accelerated flow solver for direct numerical simulation of turbulent flows, *J. Comput. Phys.* 235 (2013) 129–142, <https://doi.org/10.1016/j.jcp.2012.10.012>.
- [19] V. Cvetanoska, T. Stojanovski, Using high performance computing and monte carlo simulation for pricing american options, in: *The 9th Conference for Informatics and Information Technology*.
- [20] Y. Liang, X. Xing, Y. Li, A gpu-based large-scale monte carlo simulation method for systems with long-range interactions, *J. Comput. Phys.* 338 (2017) 252–268.
- [21] A. Humphrey, D. Sunderland, T. Harman, M. Berzins, Radiative heat transfer calculation on 16384 gpus using a reverse monte carlo ray tracing approach with adaptive mesh refinement, in: *2016 IEEE International Parallel and Distributed Processing Symposium Workshops*, 2016, pp. 1222–1231.
- [22] F. Heymann, R. Siebenmorgen, GPU-based Monte Carlo dust radiative transfer scheme applied to active galactic nuclei, *Astrophys. J.* 751 (27) (2012).
- [23] L. Tesse, F. Dupoirieux, B. Zamuner, J. Taine, Radiative transfer in real gases using reciprocal and forward monte carlo methods and a correlated-k approach, *Int. J. Heat Mass Transf.* 45 (2002) 2797–2814.
- [24] M. Cherkaoui, J. Dufresne, R. Fournier, J. Grandpeix, A. Lahellec, Monte Carlo simulation of radiation in gases with narrow-band model and a net-exchange formulation, *J. Heat Transf.* 118 (1996) 401–407.
- [25] Y. Zhang, O. Gicquel, J. Taine, Optimized emission-based reciprocity monte carlo method to speed up computation in complex systems, *Int. J. Heat Mass Transf.* 55 (2012) 8172–8177.
- [26] A. Soufiani, J. Taine, High temperature gas radiative property parameters of statistical narrow-band model for H₂O, CO₂ and CO, and correlated-k model for H₂O and CO₂, *Int. J. Heat Mass Transf.* 40 (4) (1996) 987–991.
- [27] L. Rothman, I. Gordon, A. Barbe, D. Benner, P. Bernath, M. Birk, The HITRAN 2012 molecular spectroscopic database, *J. Quant. Spectrosc. Radiat. Transf.* 130 (2013) 4–50.
- [28] L. Rothman, I. Gordon, R. Barber, H. Dothe, R. Gamache, A. Goldman, V. Perevalov, S. Tashkun, J. Tennyson, HITEMP, the high-temperature molecular spectroscopic database, *J. Quant. Spectrosc. Radiat. Transf.* 111 (2010) 2139–2150.
- [29] A. Sakurai, T. Song, S. Maruyama, H. Kim, Comparison of radiation element method and discrete ordinates interpolation method applied to three-dimensional radiative heat transfer, *Numer. Heat Transf.* 48 (2) (2005) 259–264.
- [30] T. Kim, J. Menart, H. Lee, Nongrey radiative gas analyses using the S-N discrete ordinates method, *J. Heat Transf.* 113 (1991) 946–952.
- [31] A. Humphrey, T. Harman, M. Berzins, P. Smith, A scalable algorithm for radiative heat transfer using reverse Monte Carlo ray tracing, in: *International Conference on High Performance Computing*, 2015, pp. 212–230.