



**HAL**  
open science

**Efficient uncertainty propagation for photonics:  
combining Implicit Semi-analog Monte Carlo (ISMC)  
and Monte Carlo generalised Polynomial Chaos  
(MC-gPC)**  
Gaël Poëtte

► **To cite this version:**

Gaël Poëtte. Efficient uncertainty propagation for photonics: combining Implicit Semi-analog Monte Carlo (ISMC) and Monte Carlo generalised Polynomial Chaos (MC-gPC). 2020. hal-03017635

**HAL Id: hal-03017635**

**<https://hal.science/hal-03017635>**

Preprint submitted on 21 Nov 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Efficient uncertainty propagation for photonics: combining Implicit Semi-analog Monte Carlo (ISMC) and Monte Carlo generalised Polynomial Chaos (MC-gPC)

Gaël Poëtte<sup>a</sup>

<sup>a</sup>CEA DAM CESTA, F-33114 Le Barp, France

---

## Abstract

In this paper, we build wellposed intrusive generalised Polynomial Chaos (gPC) based reduced models for uncertain photonics. We solve the reduced models with a Monte-Carlo (MC) scheme. Care is taken to highlight under which condition a reduced model (gPC based or not) is wellposed. The analysis is carried out thanks to an analogy between the construction of reduced models for uncertainty quantification and the construction of reduced models for kinetic equations. In order to enforce the aforementioned wellposedness conditions, several strategies, inspired from the hyperbolicity-preserving ones [1, 2, 3, 4, 5, 6, 7, 8] are reviewed, adapted and analysed. The resolution of the reduced models is performed thanks to an astute combination of the Implicit Semi-analog MC (ISMC, see [9]) scheme for photonics and of MC-gPC (see [10]) for uncertainty propagation. This work demonstrates that MC-gPC can be efficiently applied to a stiff nonlinear set of partial derivative equations if the MC resolution allows a fast convergence with respect to both the time and spatial discretisation (the latter properties being allowed by ISMC). Several benchmarks are investigated in the last section, they allow putting forward important aspects of the new ISMC-gPC solver for uncertain photonics.

*Keywords:* Transport, Monte-Carlo, Numerical scheme, Photonics, ISMC, Uncertainties, Sensitivity Analysis, MC-gPC

---

## 1. Introduction

In this article, we are interested in propagating uncertainties through the time-dependent, non-linear, radiative transfer equations. The model has general form (see [11]):

$$\left\{ \begin{array}{l} \frac{1}{c} \partial_t I(x, t, \omega, \nu) + \omega \cdot \nabla I(x, t, \omega, \nu) + \sigma_t(T(E(x, t)), \nu) I(x, t, \omega, \nu) \\ \qquad \qquad \qquad = \sigma_a(T(E(x, t)), \nu) B(T(E(x, t)), \nu) + \sigma_s(T(E(x, t)), \nu) \int I(x, t, \omega', \nu) d\omega', \\ \partial_t E(T(x, t)) = \iint c \sigma_a(T(E(x, t)), \nu') (I(x, t, \omega', \nu') - B(T(E(x, t)), \nu')) d\omega' d\nu'. \end{array} \right. \quad (1)$$

---

*Email address:* [gael.poette@cea.fr](mailto:gael.poette@cea.fr) (Gaël Poëtte)

System (1) is relevant to describe the behaviour of photons incoming into cold media and is intensively used for, amongst others, astrophysical applications [11, 12, 13, 14], atmospheric physics [15] and optical imaging [16]. In the above equations,  $I$  and  $E$  are the unknowns of the system and stand respectively for the intensity of radiation energy and the material energy. Variables  $t \geq 0$ ,  $x \in \Omega \subset \mathbb{R}^3$ ,  $\omega \in \mathbb{S}^2$  and  $\nu \in ]0, \infty[$  are respectively the time, space, angle<sup>1</sup> and frequency variables. The opacities  $\sigma_t$ ,  $\sigma_a$  and  $\sigma_s$  are given functions of  $(x, t)$  or depend on  $E$ , often more commonly expressed *via*  $T$ , the material temperature. They stand for the total, absorption and scattering opacities. In particular, we have  $\sigma_t = \sigma_a + \sigma_s$ . Quantity  $c$  denotes the speed of light. Quantity  $B(T, \nu) = \frac{2h\nu^3}{c^2} (e^{\frac{h\nu}{kT}} - 1)^{-1}$  is the Planck function where  $h$  is the Planck constant and  $k$  is the Boltzmann constant. The density of internal energy  $E(T)$  depends on the material temperature  $T(x, t)$  *via* an equations of state (eos)  $dE = \rho C_v(T) dT$  with  $\rho$  the mass density and  $C_v > 0$  the heat capacity<sup>2</sup> (constant for a perfect gas). Initial and boundary conditions must be supplemented to system (1):

$$I(x, 0, \omega, \nu) = I^0(x, \omega, \nu), \quad E(x, 0) = E^0(x), \quad x \in \Omega, \quad \omega \in \mathbb{S}^2, \quad \nu \in ]0, \infty[, \quad (2)$$

$$I(x, t, \omega, \nu) = I_b(t, \omega, \nu), \quad t \geq 0, \quad x \in \partial\Omega, \quad \omega \cdot n_s(x) < 0, \quad \nu \in ]0, \infty[, \quad (3)$$

where  $n_s$  is the outward normal to  $\Omega$  at  $x$ . System (1) together with initial and boundary conditions (2) + (3) define the well-posed [17] mathematical problem we want to solve. In fact, in this paper, we are going to consider a slightly simpler model commonly called the *grey approximation*. It is obtained assuming the opacities are independent of  $\nu$  and focuses on unknowns  $I(x, t, \omega) = \int I(x, t, \omega, \nu) d\nu$  and  $E(x, t)$ :

$$\begin{cases} \frac{1}{c} \partial_t I(x, t, \omega) + \omega \cdot \nabla I(x, t, \omega) + \sigma_t(E(x, t)) I(x, t, \omega) \\ \qquad \qquad \qquad = \sigma_a(E(x, t)) B(E(x, t)) + \sigma_s(E(x, t)) \int I(x, t, \omega') d\omega', \\ \partial_t E(T(x, t)) = c \sigma_a(E(x, t)) \int (I(x, t, \omega') - B(E(x, t))) d\omega', \\ I(x, 0, \omega) = I^0(x, \omega), \quad E(x, 0) = E^0(x), \quad x \in \Omega, \quad \omega \in \mathbb{S}^2, \\ I(x, t, \omega) = I_b(t, \omega), \quad t \geq 0, \quad x \in \partial\Omega, \quad \omega \cdot n_s < 0. \end{cases} \quad (4)$$

For the grey approximation, the Planck function resumes to  $B(E, \nu) = B(T(E), \nu) = B(T(E)) = B(E) = aT^4(E)$  with  $a$  the radiative constant. In this paper, we are going to focus on (4) instead of (1): considering (4) eases the next analysis and calculations without loss of generalities.

Now, in this paper, we are even interested in *being able to accurately take uncertainties, in a broad sense<sup>3</sup>, into account* in system (4). When dealing with an uncertainty quantification problem, it is common to explicit the dependence of the solution with respect to the uncertain parameters denoted here by  $X$ . Without loss of generality in the following sections, we consider that  $X$  is a vector  $X = (X_1, \dots, X_Q)^t$  of  $Q$  independent random variables of probability measure  $d\mathcal{P}_X = \prod_{i=1}^Q d\mathcal{P}_{X_i}$ .

<sup>1</sup>Here, note that we use an abusive but concise notation in order to denote the integration on the unit sphere  $\mathbb{S}^2$  with  $d\omega = \frac{1}{|\mathbb{S}^2|} \mathbf{1}_{\mathbb{S}^2}(w) dw$ . It is such that  $\int d\omega = 1$ .

<sup>2</sup>The heat capacity being strictly positive,  $T \rightarrow E(T)$  defines a bijection and can be inverted. Function  $E \rightarrow T(E)$  denotes this inverse.

<sup>3</sup>geometrical, in the opacities, in the heat capacity, in the eos, in the boundary conditions etc.

It is always possible to come back to such framework<sup>4</sup>. As a result, solving the uncertain counterpart of (4) resumes to solving the Stochastic PDE (SPDE)

$$\left\{ \begin{array}{l} \frac{1}{c} \partial_t I(x, t, \omega, X) + \omega \cdot \nabla I(x, t, \omega, X) + \sigma_t(E(x, t, X), X) I(x, t, \omega, X) \\ \qquad \qquad \qquad = \sigma_a(E(x, t, X), X) B(E(x, t, X)) + \sigma_s(E(x, t, X), X) \int I(x, t, \omega', X) d\omega', \\ \partial_t E(x, t, X) = c \sigma_a(E(x, t, X), X) \int (I(x, t, \omega', X) - B(E(x, t, X))) d\omega', \\ I(x, 0, \omega, X) = I^0(x, \omega, X), \quad E(x, 0, X) = E^0(x, X), \quad x \in \Omega, \quad \omega \in \mathbb{S}^2, \\ I(x, t, \omega, X) = I_b(t, \omega, X), \quad t \geq 0, \quad x \in \partial\Omega, \quad \omega \cdot n_s < 0, \\ X \sim d\mathcal{P}_X. \end{array} \right. \quad (5)$$

System (5) is assumed wellposed<sup>5</sup>  $\forall X \sim d\mathcal{P}_X$ : in fact, we suppose that  $d\mathcal{P}_X$  is such that no realisation  $X$  can lead to a singular problem (5). Under such condition, system (5) then even respects a maximum principle, pointwise with respect to  $X \sim d\mathcal{P}_X$ , see [17]. In the above problem, we are mainly interested in the statistics<sup>6</sup> of  $X \rightarrow I(x, t, X) = \int I(x, t, \omega, X) d\omega = aT_r^4(x, t, X)$  where  $T_r$  is the radiative temperature,  $X \rightarrow E(x, t, X)$  and  $X \rightarrow T(x, t, X)$  at specified locations  $x \in \Omega$  and times  $t \in [0, t^*]$ .

Of course, different values of  $X$  correspond to different fully decoupled deterministic equations: in principle, there is no difficulty in solving such uncertain problems. The main issue comes from the fact that exact propagation of uncertainties is very expensive from the computational point of view: system (4) is often solved thanks to a Monte-Carlo (MC) scheme [24, 25, 26, 27, 28, 29]. MC schemes imply tracking ( $N_{MC}$ ) particles on a geometry of interest and tallying their contributions in order to solve (4). This resolution method is known to be efficient for high ( $3(x) + 1(t) + 2(\omega) + 1(\nu) = 7$ ) dimensional problems but costly. Running several deterministic MC computations for several ( $N$ ) values of  $X$  to propagate (non-intrusively) the uncertainties can consequently be prohibitive.

In [10, 30], a  $P$ -truncated generalised Polynomial Chaos (gPC) based reduced model of the *instationary uncertain linear Boltzmann equation* has been introduced<sup>7</sup>. It is solved thanks to an astute converging Monte-Carlo (MC) scheme [10]: the idea is to make the MC particles not only solve the physical fields  $(x, t, \omega)$  but also the uncertain one  $X$ . The uncertain counterpart is solved *on-the-fly* during the MC resolution: instead of tracking  $N \times N_{MC}$  particles discretising the physical space of variables  $(x, t, \omega)$  in a non-intrusive way<sup>8</sup>, the methodology ensures about the same order of accuracy with only  $N_{MC}$  particles by discretising the whole space of variables  $(x, t, \omega, X)$  with an MC method to approximate the gPC coefficients. Similar approaches have been developed for the Fokker-Planck equation [31] or for the quadratic Boltzmann equation [32, 33] or in an eigenvalue/eigenvector computation ( $k_{\text{eff}}$ ) context [34] and give promising results on other (usually MC solved) physical models. For the linear transport equation, the spectral (i.e. fast) convergence

---

<sup>4</sup>At the cost of more or less tedious pretreatments leading to a controled approximation [18, 19, 20] and decorrelation [21, 22].

<sup>5</sup>Existence and uniqueness of the solution is ensured.

<sup>6</sup>i.e. mean, variance, histogram, sensitivity indices [23] etc.

<sup>7</sup>Note that the instationary uncertain linear Boltzmann equation is a particular case of (5) when  $\partial_t E \sim 0$ , i.e.  $C_v \ll 1$  for example.

<sup>8</sup>Hence tensorising the set of points for  $X$  and for  $(x, t, \omega)$ .

of the built hierarchical models has been numerically [10] and theoretically [30] demonstrated. The approach in [10, 31, 33, 30], denoted by MC-gPC in the following, enters the class of intrusive<sup>9</sup> gPC methods. It consequently demands some code modifications. To sum-up, with MC-gPC (see [10, 30]), *important gains with respect to non-intrusive methods* have been observed in *low to moderate* stochastic dimensions<sup>10</sup>  $Q \sim 1 - 10$  with *simple* code modifications of an existing MC code and *without changing* the HPC strategy<sup>11</sup> of the code. The questions now are: *is it possible to generalise the methodology used in these papers to be able to solve a stiff nonlinear uncertain transport equation such as (5)? Is the intrusive MC-gPC [10, 31, 33] solver still more efficient than a non-intrusive solver for problem (5)?* This paper documents and presents our efforts in order to answer the latter questions. In this paper, we design an ISMC-gPC solver. It relies on the combined material of [10] (for MC-gPC) and [9] (for ISMC<sup>12</sup>) plus some additional key ingredients proper to the photonic problem we aim at addressing in this paper (mainly positiveness, boundedness and moment preserving strategies).

The paper is organized as follows: in section 2, we study the conditions we must enforce for the construction of a wellposed reduced model of (5) capturing the uncertainties. Wellposedness preserving strategies, i.e. numerical strategies allowing to ensure the latter conditions remains fulfilled during every computations, are also discussed in this section. Those are mainly inspired from the hyperbolicity-preserving strategies of the literature [1, 2, 3, 4, 5, 6, 7, 36, 8] as similar difficulties are encountered for the reduced models we build in this paper. In section 3, we recall the ISMC solver from [9] and explain in which way we think it is a good candidate to be combined to MC-gPC. The general sketch of an ISMC solver is also recalled. In section 4, the gPC based reduced model of (5) is built and its MC resolution is presented. Care is taken to put forward where, with respect to the algorithm described in section 3, modifications of a deterministic ISMC solver must be made to take into account uncertainties *on-the-fly* during the MC resolution. Following the descriptions of the next sections should allow the interested reader to perform the relevant modifications to his own ISMC implementation. Section 5 is devoted to benchmarks and numerical test-cases, section 6, to concluding remarks.

## 2. A well-posed reduced model for (5)

In [30], in a similar context but for the linear transport equation, efficient and competitive hierarchical  $P$ -truncated gPC based reduced models are built. To do so, the stochastic Galerkin gPC (sG-gPC) method [37, 38, 39, 40, 41, 42, 43, 44, 45, 46] is used. In the linear case, it is proven enough to ensure wellposedness of the  $P$ -truncated models independently of the truncation order  $P$ . On another hand, sG-gPC is known to fail for some nonlinear models, see [47, 8, 36, 5, 6, 7, 1, 2, 3, 4]. Numerical/robustness difficulties [8] are encountered and may question a code architecture if not *a priori* taken into account and clearly understood. The efficient strategies to ensure by construction the wellposedness of the reduced model are problem dependent, see [8, 5, 6, 1, 2, 3, 4]: they depend on the structure of the nonlinear problem (elliptic, parabolic, hyperbolic etc.) and on its properties

---

<sup>9</sup>It does not propagate the uncertainties by relying on several runs of a black-box code.

<sup>10</sup>MC-gPC is based on gPC which is sensitive to the curse of dimension, see [10, 30].

<sup>11</sup>The HPC strategy we have in mind is commonly called *replication domain*, see [35]. It consists in replicating the geometry on several processors and tracking several MC particles populations with different initial seeds in every replicated domains. At the end of the time steps, the contribution of every processors are averaged. This parallel strategy is particularly well suited to MC codes, taking advantage of the independence of the MC particles.

<sup>12</sup>ISMC is for Implicit Semi-analog Monte-Carlo.

(maximum principle, regularity of the solution etc.). The strategies to enforce some of the desired properties for the reduced model range from gPC-wavelet combination [47], Multi-Element gPC [48], basis adaptation [49], nonlinear transformations [8, 5, 1, 4], filtering [3, 50], weighting [51, 7]: in other words, many levers are now at hand, one mainly has to make sure choosing the relevant one for its problem<sup>13</sup>. In this section, we want to anticipate and avoid beforehand<sup>14</sup> potential numerical difficulties such as the ones encountered in the aforementioned papers. This is why we here aim at being able to build uncertainty capturing *wellposed* reduced model for (5). The questions now are: is it possible to do so for nonlinear system (5)? Are there compatibility conditions? If yes, which are they and how can we enforce them during the numerical resolution?

In order to answer the above questions, we are going to intensively use the material of paper [17]. Let us sum it up in the next lines. In [17], the author is interested in the wellposedness (existence, unicity and uniform convergence) of system (1). In order to prove the property, the author builds a Cauchy problem of unknown  $u = (I, E)$

$$\begin{cases} \partial_t u + \mathcal{A}u + \mathcal{B}u = 0, \\ u(t=0) = u_0, \end{cases} \quad (6)$$

where  $\mathcal{A}$  is an accretive operator (see [52, 17]) and  $\mathcal{B}$  is (at least) Lipschitz or (at best) Lipschitz and accretive. With those properties for  $\mathcal{A}, \mathcal{B}$ , Crandall-Liggett theorem [52] states that the solution of the time-discretised counterpart of (6) converges uniformly toward the *unique existing continuous solution* of (6) for any given final time  $t^* > 0$ . Now, in [17],  $\mathcal{A}$  and  $\mathcal{B}$  are not such that (6) exactly coincides with (1): this is due to the introduction of a *truncation*  $\psi(I)$  of  $I$  within operator  $\mathcal{B}$ . The author needs this truncation  $\psi$  for operator  $\mathcal{B}$  to have the desired properties: in a nutshell (see [17]), operator  $\mathcal{B}$  is Lipschitz if  $\psi$  is bounded and positive. This leads the author in [17] to choosing a truncation of the form  $\psi(I) = \max(0, \min(I, B(\nu, M)))$  where  $M > 0$  is a constant. As a consequence, if  $0 \leq I \leq B(\nu, M)$ , then  $I = \psi(I)$  and (6) coincides with (1). The author then only needs to be able to choose a relevant  $M$  in order to impose the previous inequality: thanks to a maximum principle for (6), the author identifies *a priori*, thanks to the initial and boundary conditions, the magnitude of the truncation constant  $M$  such that  $I = \psi(I)$ . Once a proper  $M$  chosen *a priori*, solving (6) is equivalent to solving (1) if fine enough computations are carried out, see [17]. In the next paragraph, we are going to use the same kind of reasoning in order to build a wellposed reduced model for (5).

Let us define what we call a reduced model: to reduce unknown  $(I(x, t, \omega, X), E(x, t, X))$ , we would like to look for unknowns of the form<sup>15</sup>

$$\begin{aligned} (I^P(x, t, \omega, X), E^P(x, t, X)) &= (I^P(\lambda^I(x, t, \omega), X), E^P(\lambda^E(x, t), X)), \\ \text{with } \lambda^I(x, t, \omega) &= (\lambda_0^I(x, t, \omega), \dots, \lambda_P^I(x, t, \omega)) \text{ and } \lambda^E(x, t) = (\lambda_0^E(x, t), \dots, \lambda_P^E(x, t)). \end{aligned} \quad (7)$$

In the previous expression,  $\lambda^I, X \rightarrow I^P(\lambda^I, X)$  and  $\lambda^E, X \rightarrow E^P(\lambda^E, X)$  are arbitrary at this stage of the discussion. They will be chosen at the appropriate moment in this paper to enforce wellposedness. Of course, amongst the possibilities for  $(I^P, E^P)$  we can count

<sup>13</sup>Which is far from being simple nor obvious.

<sup>14</sup>Note that this anticipation gains relevance in section 5: robustness difficulties are encountered in practice with sG-gPC if the material of this section is not applied, see section 5.3.

<sup>15</sup>In the next expression, we introduce equivalent notations which will be used all along the paper.

- polynomials, i.e. gPC based approximations (i.e. sG-gPC). After all, in [10, 30], we had no reason to resort to more elaborate reduced models. For linear problems, sG-gPC already builds wellposed ones and is theoretically (spectral convergence [30]) and numerically (thanks to MC-gPC [10]) efficient.
- But also more elaborate and more or less complicated ansatz such as nonlinear functions of polynomials [8, 1, 4, 6], relaxation schemes [50], filtered polynomials [1, 3], piecewise polynomials [48], weighted piecewise polynomials [51, 7], compositions of neurons (deep or not neural networks) [53] etc.

When looking for an unknown of the form (7), we implicitly aim at *trading dimensionality for size* in the sense that solving our reduced model now only supposes finding  $(\lambda^I(x, t, \omega), \lambda^E(x, t))$ : the vector of unknowns depends on  $(Q)$  less dimensions (i.e. they depend only on  $x, t, \omega$  instead of  $x, t, \omega, X$ ), but is now of size<sup>16</sup>  $(P + 1) \times 2 = (p_{1D} + 1)^Q \times 2$  instead of 2. Of course, in high stochastic dimensions,  $P$  may be important (curse of dimensionality). In the following,  $(I^P, E^P)$  is called a *truncation*<sup>17</sup> of  $(I, E)$  to insist on the analogy made with the material of [17].

Now, we would like truncation  $(I^P, E^P)$  to be close to  $(I, E)$  solution of (5). We want to identify conditions for  $(I^P, E^P)$  to be an existing and unique solution of reduced model (we drop the dependences with respect to  $x, t, \omega$  for conciseness here)

$$\left\{ \begin{array}{l} \left\{ \begin{array}{l} \frac{1}{c} \partial_t I^P(\lambda^I, X) + \omega \cdot \nabla I^P(\lambda^I, X) + \sigma_t(\lambda^E, X) I^P(\lambda^I, X) \\ \qquad \qquad \qquad = \sigma_a(\lambda^E, X) B(E^P(\lambda^E, X)) + \sigma_s(\lambda^E, X) \int I^P(\lambda^I, X) d\omega', \end{array} \right. \quad (8a) \\ \left\{ \begin{array}{l} \partial_t E^P(\lambda^E, X) = \int c \sigma_a(\lambda^E, X) (I^P(\lambda^I, X) - B(E^P(\lambda^E, X))) d\omega', \\ X \sim d\mathcal{P}_X. \end{array} \right. \quad (8b) \end{array} \right.$$

Of course, (8) must be supplemented by initial and boundary conditions but we abusively omit them in the later paragraphs for the sake of conciseness. In the above expression, we have only 2 equations and  $2 \times (P + 1)$  unknowns. We have to build the equations that unknowns  $\lambda_0^I, \dots, \lambda_P^I, \lambda_0^E, \dots, \lambda_P^E$  must satisfy for  $(I^P, E^P)$  to be a relevant reduced model for (5). Let us introduce the jacobian vectors  $\nabla_\lambda I^P$  and  $\nabla_\lambda E^P$  of general terms

$$[\nabla_\lambda I^P]_i = \partial_{\lambda_i^I} I^P(\lambda_0^I, \dots, \lambda_P^I, X) \quad \text{and} \quad [\nabla_\lambda E^P]_i = \partial_{\lambda_i^E} E^P(\lambda_0^E, \dots, \lambda_P^E, X). \quad (9)$$

Developping (8) and using the previous vector notation leads to

$$\left\{ \begin{array}{l} \left\{ \begin{array}{l} [\nabla_\lambda I^P]^t \frac{1}{c} \partial_t \begin{pmatrix} \lambda_0^I \\ \dots \\ \lambda_P^I \end{pmatrix} + [\nabla_\lambda I^P]^t \omega \cdot \nabla \begin{pmatrix} \lambda_0^I \\ \dots \\ \lambda_P^I \end{pmatrix} + \sigma_t I^P = \sigma_a B(E^P) + \sigma_s \int I^P d\omega', \\ [\nabla_\lambda E^P]^t \partial_t \begin{pmatrix} \lambda_0^E \\ \dots \\ \lambda_P^E \end{pmatrix} = \int c \sigma_a (I^P - B(E^P)) d\omega', \end{array} \right. \quad (10a) \\ X \sim d\mathcal{P}_X. \quad (10b) \end{array} \right.$$

<sup>16</sup>where  $p_{1D}$  denotes the order of the monodimensional in each stochastic directions.

<sup>17</sup>In an Uncertainty Quantification (UQ) context, it is more commonly called a metamodel or a surrogate model.

The above system of equation is still of size 2. Let us now introduce additional hypothesis on the truncation  $(I^P, E^P)$ . Let us assume that we have

$$I^P(\lambda_0^I, \dots, \lambda_P^I, X) = f_I(\sum_{l=0}^P \lambda_l^I \phi_l(X)) \quad \text{and} \quad E^P(\lambda_0^E, \dots, \lambda_P^E, X) = f_E(\sum_{l=0}^P \lambda_l^E \phi_l(X)), \quad (11)$$

where<sup>18</sup>  $i \in \mathbb{R} \rightarrow f_I(i) \in \mathbb{R}$  and  $e \in \mathbb{R} \rightarrow f_E(e) \in \mathbb{R}$  are measurable scalar functions and  $(\phi_l(X))_{l \in \{0, \dots, P\}}$  is orthonormal<sup>19,20</sup> with respect to the scalar product defined by  $d\mathcal{P}_X$ . With this truncation hypothesis (10) becomes

$$\left\{ \begin{array}{l} \left( f'_I(\sum_l \lambda_l^I \phi_l)(\phi_0, \dots, \phi_P) \left( \frac{1}{c} \partial_t \begin{pmatrix} \lambda_0^I \\ \dots \\ \lambda_P^I \end{pmatrix} + \omega \cdot \nabla \begin{pmatrix} \lambda_0^I \\ \dots \\ \lambda_P^I \end{pmatrix} \right) + \sigma_t I^P = \sigma_a B(E^P) + \sigma_s \int I^P d\omega', \right. \\ \left. f'_E(\sum_l \lambda_l^E \phi_l)(\phi_0, \dots, \phi_P) \partial_t \begin{pmatrix} \lambda_0^E \\ \dots \\ \lambda_P^E \end{pmatrix} = \int c \sigma_a (I^P - B(E^P)) d\omega', \right. \\ \left. X \sim d\mathcal{P}_X. \right. \end{array} \right. \quad (12a)$$

Assume now that both  $f'_I, f'_E$  are non-zero, then we can rewrite (12) as

$$\left\{ \begin{array}{l} \left( \phi_0, \dots, \phi_P \right) \left( \frac{1}{c} \partial_t \begin{pmatrix} \lambda_0^I \\ \dots \\ \lambda_P^I \end{pmatrix} + \omega \cdot \nabla \begin{pmatrix} \lambda_0^I \\ \dots \\ \lambda_P^I \end{pmatrix} \right) + \Sigma_t I^P = \Sigma_a^I B(E^P) + \Sigma_s \int I^P d\omega', \\ \left( \phi_0, \dots, \phi_P \right) \partial_t \begin{pmatrix} \lambda_0^E \\ \dots \\ \lambda_P^E \end{pmatrix} = \int c \Sigma_a^E (I^P - B(E^P)) d\omega', \\ X \sim d\mathcal{P}_X, \end{array} \right. \quad (13a)$$

where  $\Sigma_a^I, \Sigma_a^E, \Sigma_t, \Sigma_s$  are related to  $\sigma_a, \sigma_t, \sigma_s$ . Due to the introduction of  $\Sigma_a^I$  and  $\Sigma_a^E$ , system (13) is not necessarily (this depends on the choice of  $f_I, f_E$ ) written in term of conservative variable anymore. Independently of the above choice for  $f_I, f_E$ , we can now easily perform a Galerkin projection by multiplying (13) by vector  $(\phi_0, \dots, \phi_l)^t$  and integrating the whole set of equations with respect to  $d\mathcal{P}_X$ . Vector  $(\lambda^I, \lambda^E)$  must satisfy

$$\left\{ \begin{array}{l} \frac{1}{c} \partial_t \begin{pmatrix} \lambda_0^I \\ \dots \\ \lambda_P^I \end{pmatrix} + \omega \cdot \nabla \begin{pmatrix} \lambda_0^I \\ \dots \\ \lambda_P^I \end{pmatrix} + \int \Sigma_t I^P \begin{pmatrix} \phi_0 \\ \dots \\ \phi_P \end{pmatrix} d\mathcal{P}_X = \\ \int \Sigma_a^I B(E^P) \begin{pmatrix} \phi_0 \\ \dots \\ \phi_P \end{pmatrix} d\mathcal{P}_X + \int \Sigma_s \int I^P d\omega' \begin{pmatrix} \phi_0 \\ \dots \\ \phi_P \end{pmatrix} d\mathcal{P}_X, \\ \partial_t \begin{pmatrix} \lambda_0^E \\ \dots \\ \lambda_P^E \end{pmatrix} = \iint c \Sigma_a^E (I^P - B(E^P)) \begin{pmatrix} \phi_0 \\ \dots \\ \phi_P \end{pmatrix} d\mathcal{P}_X d\omega'. \end{array} \right. \quad (14a)$$

<sup>18</sup>At this stage,  $f_I$  and  $f_E$  are from  $\mathbb{R}$  into  $\mathbb{R}$  but we will see that wellposedness imposes some restrictions on the output spaces.

<sup>19</sup>i.e. we have  $\int \phi_l(X) \phi_k(X) d\mathcal{P}_X = \delta_{k,l}, \forall (k, l) \in \{0, \dots, P\}^2$ .

<sup>20</sup>Note that  $(\phi_l)_{l \in \{0, \dots, P\}}$  may be a gPC basis or not. It only needs the orthonormality condition at this stage of the discussion.



System (14) has the structure of a *nonlinear multigroup transport system*, see [54, 26]. The system is closed once  $i \rightarrow f_I(i)$ ,  $e \rightarrow f_E(e)$  chosen and defined. Let us now sum up our logical chain of thinking: we want to be able to choose truncation<sup>21</sup>  $(I^P, E^P)$  such that reduced model (10), in which  $(\lambda^I, \lambda^E)$  satisfies (14), is wellposed. But before looking for the properties on  $(I^P, E^P)$  for wellposedness to hold, let us discuss on the interest of relying on a quite general methodology:

- If  $i \rightarrow I^P(i) = i$ ,  $e \rightarrow E^P(e) = e$  then the above methodology builds the same reduced model as in a classical sG-gPC framework [37, 38, 39, 40, 41, 42, 43, 44, 45, 46].
- If  $i \rightarrow I^P(i) = [\nabla_{\lambda^I} s]^{-1}(i, e)$ ,  $e \rightarrow [\nabla_{\lambda^E} s]^{-1}(i, e)$  where  $i, e \rightarrow s(i, e)$  is a strictly convex function, then the methodology recovers the Intrusive Polynomial Moment (IPM) method [8, 1, 4].
- If  $I^P, E^P$  are piecewise continuous functions then we recover ME-gPC [48, 49, 7].
- Other choices may lead to strategies closely related to the filtered gPC solver of [3, 50], the maximum-principle satisfying moment one of [1, 4], the weighted ones of [51, 7] or the hyperbolicity-preserving one of [5, 6]. Some of them will be discussed later on.

At this stage, the reader may wonder why all these remarks and analogies? Because if several choices are available, some code architecture may be easier to modify. And, at this stage of the discussion, we still do not know how to choose  $I^P, E^P$  for the resolution of the reduced model (8) to be *mathematically wellposed*<sup>22</sup>, *physically relevant*<sup>23</sup> and *numerically efficient*<sup>24</sup>.

We first would like to express conditions on truncation  $(I^P, E^P)$  for (8) (complemented with (7) and (14)) to be wellposed. For this, we intensively rely on an analogy between (8) and the frequency dependent model (1) whose wellposedness is discussed in [17]. Between (1) and (8), there are

- some similarities: the frequency  $\nu \in ]0, \infty[$  in the first equation of (1) is replaced by  $X \sim d\mathcal{P}_X$  in the first equation of (8).
- And some differences: the second equation of (1) is integrated with respect to  $\nu$  whereas the second equation of (8) is not integrated with respect to  $X$ .

This analogy between the uncertain reduced model (8) and the kinetic model (1) is going to help us identify relevant mathematical tools and theorems to identify the wellposedness conditions we aim at characterising. Note that the analogy between uncertainty propagation *via* reduced models and kinetic theory has already been put forward in several papers, see [1, 2, 33, 8, 36]. In a sense, this section is one additional argument encouraging such analogy in order to obtain theoretical results and have a better understanding of the uncertainty capturing models we build.

---

<sup>21</sup>or equivalently  $(f_I, f_E)$ .

<sup>22</sup>this is the purpose of the current section 2.

<sup>23</sup>We will only verify experimentally that we build physically relevant reduced models in the numerical section 5.

<sup>24</sup>The preservation of the desired ISMC and MC-gPC numerical properties with respect to efficiency will be discussed in sections 3–4 and will be numerically verified in section 5.

As explained above, in order to characterise the wellposedness conditions for our reduced model, we rely on an analogy between (1) and (8) and on the theory of accretive operators [17, 52]. The first step is to define operators  $\mathcal{A}, \mathcal{B}$  such that (8) can be put under the form

$$\begin{cases} \partial_t u^P + \mathcal{A}u^P + \mathcal{B}u^P = 0, \\ u^P(0) = u_0^P, \end{cases} \quad (15)$$

where  $u^P \equiv u^P(x, t, \omega, X) = (I^P(x, t, \omega, X), E^P(x, t, X))$ . Let us first define the norm we are going to work with:

$$\|u^P\|(t) = \frac{1}{|\Omega|} \iiint_{\Omega} |I^P(x, t, \omega, X)| \, d\omega \, dx \, d\mathcal{P}_X + \frac{1}{|\Omega|} \iint_{\Omega} |E^P(x, t, X)| \, dx \, d\mathcal{P}_X. \quad (16)$$

With the above choice of norm, having  $u^P$  bounded is equivalent to having<sup>25</sup>  $I^P \in L^1(\Omega, \mathbb{S}^2, d\mathcal{P}_X)$  and<sup>26</sup>  $E^P \in L^1(\Omega, d\mathcal{P}_X)$ . But more importantly, if  $I^P, E^P$  are both positive, the norm of  $u^P$  stands for the total energy at time  $t$  in the whole geometry  $\Omega$  of the system 'photon+matter' for reduced model (8).

According to [17, 52], it is enough characterising conditions for  $\mathcal{A}$  to be accretive and for  $\mathcal{B}$  to be at least Lipschitz, at best Lipschitz and accretive, for system (15) to be wellposed. Crandall-Liggett's theorem [52] then ensures that for any given  $t^* > 0$ , problem (15) has a unique solution  $u^P \in \mathcal{C}^0([0, t^*] : L^1(\Omega, \mathbb{S}^2, d\mathcal{P}_X) \times L^1(\Omega, d\mathcal{P}_X))$ . In other words, if the above properties hold for operators  $\mathcal{A}, \mathcal{B}$ , the solution of the built reduced model will be continuous: we would be in desirable conditions for gPC based reduced models to exhibit fast convergence rates [55, 30, 8].

Let us now decompose (15) into several operators and study them successively:

- Let us first rewrite  $\mathcal{A}$  as  $\mathcal{A} = \mathcal{L} - \mathcal{D}$  with

$$\mathcal{L}u^P = \begin{pmatrix} c\omega \cdot \nabla I^P \\ 0 \end{pmatrix} \text{ and } \mathcal{D}u^P = \begin{pmatrix} c\sigma_s \int I^P \, d\omega' - c\sigma_s I^P \\ 0 \end{pmatrix}.$$

Defining  $\mathcal{A}$  as above is convenient in practice because it corresponds exactly with operator  $\mathcal{A}$  of [17]: in [17], the author shows that  $\mathcal{L}$  is accretive and  $\mathcal{D}$  is Lipschitz, implying that  $\mathcal{A}$  is accretive.

**Remark 2.1.** Note that if  $\mathcal{B} = 0$  (i.e.  $\partial_t E = 0$ ), reduced model (15) coincides with the balanced ( $\sigma_a \equiv 0$ ) reduced model of the linear Boltzmann equation, see [10]. If  $\sigma_a \neq 0$ , it is enough working with truncation  $\tilde{I}^P = I^P e^{-c\sigma_a t}$  in  $\mathcal{A}$  to be in balanced conditions. In other words, if  $\mathcal{B} = 0$ , we are in the same conditions as in [10, 30]. The accretiveness of  $\mathcal{A}$  here introduces a new proof of the wellposedness of the sG-gPC reduced model studied in [10, 30]: no particular conditions are needed on the truncation  $I^P$  for wellposedness. It explains why we can choose  $I^P$  to be a gPC expansion  $I^P = \sum_{l=0}^P I_k \phi_l$  just as in [10, 30] without triggering particular numerical difficulties.

---

<sup>25</sup>We have  $I^P \in L^1(\Omega, \mathbb{S}^2, d\mathcal{P}_X)$  if  $\frac{1}{|\Omega|} \iiint_{\Omega} |I^P(x, t, \omega, X)| \, d\omega \, dx \, d\mathcal{P}_X < \infty$

<sup>26</sup>We have  $E^P \in L^1(\Omega, d\mathcal{P}_X)$  if  $\frac{1}{|\Omega|} \iint_{\Omega} |E^P(x, t, X)| \, dx \, d\mathcal{P}_X < \infty$ .

- Let us now define operator  $\mathcal{B}$ . Note that  $\mathcal{B}$  is different from the one of [17]. Some work consequently remains to be done. If  $\mathcal{B}$  is (at least) Lipschitz, then system (15) is accretive. Let us choose  $\mathcal{B}$  as

$$\mathcal{B}u^P = \mathcal{B} \begin{pmatrix} I^P \\ E^P \end{pmatrix} = \begin{pmatrix} \sigma_a(E^P)((I^P - B(E^P))) \\ - \int \sigma_a(E^P)((I^P - B(E^P)) d\omega \end{pmatrix}.$$

By definition, an operator  $\mathcal{R}$  is Lipschitz if there exists a constant  $0 < L < \infty$  such that  $\forall (u_1, u_2) \in L^1(\Omega, \mathbb{S}^2, d\mathcal{P}_X) \times L^1(\Omega, d\mathcal{P}_X)$ ,

$$\|\mathcal{R}u_1 - \mathcal{R}u_2\| \leq L\|u_1 - u_2\|.$$

Let  $(u_1, u_2) \in L^1(\Omega, \mathbb{S}^2, d\mathcal{P}_X) \times L^1(\Omega, d\mathcal{P}_X)$ , then

$$\begin{aligned} \|\mathcal{B}u_1 - \mathcal{B}u_2\|(t) &= \frac{1}{|\Omega|} \iiint_{\Omega} |\sigma_a(E_2^P)((I_2^P - B(E_2^P))) - \sigma_a(E_1^P)((I_1^P - B(E_1^P)))| d\omega dx d\mathcal{P}_X \\ &\quad + \frac{1}{|\Omega|} \iiint_{\Omega} \left| \int \sigma_a(E_2^P)((I_2^P - B(E_2^P))) d\omega - \int \sigma_a(E_1^P)((I_1^P - B(E_1^P))) d\omega \right| dx d\mathcal{P}_X, \quad (17) \\ &\leq \frac{2}{|\Omega|} \iiint_{\Omega} \underbrace{|\sigma_a(E_2^P)((I_2^P - B(E_2^P))) - \sigma_a(E_1^P)((I_1^P - B(E_1^P)))|}_{|q(u_1) - q(u_2)|} d\omega dx d\mathcal{P}_X. \end{aligned}$$

Now, let us assume that

**Hypothesis 1.**  $\exists C > 0$  and  $L_a > 0$  such that  $|\sigma_a(E^P)| \leq C$  and  $|\sigma_a(E_2^P) - \sigma_a(E_1^P)| \leq L_a |E_2^P - E_1^P|$ ,

then

$$\begin{aligned} |q(u_1) - q(u_2)| &= \left| \begin{aligned} &\sigma_a(E_2^P)((I_2^P - I_1^P) + (\sigma_a(E_2^P) - \sigma_a(E_1^P))I_1^P \\ &+ \sigma_a(E_2^P)(B(E_1^P) - B(E_2^P)) + (\sigma_a(E_1^P) - \sigma_a(E_2^P))B(E_1^P) \end{aligned} \right|, \\ &\leq \left| \sigma_a(E_2^P)(I_2^P - I_1^P) \right| + \left| (\sigma_a(E_2^P) - \sigma_a(E_1^P))I_1^P \right| \\ &\quad + \left| \sigma_a(E_2^P)(B(E_1^P) - B(E_2^P)) \right| + \left| (\sigma_a(E_1^P) - \sigma_a(E_2^P))B(E_1^P) \right|, \quad (18) \\ &\leq \left| \sigma_a(E_2^P) \right| |I_2^P - I_1^P| + \left| \sigma_a(E_2^P) - \sigma_a(E_1^P) \right| |I_1^P| \\ &\quad + \left| \sigma_a(E_2^P) \right| |B(E_1^P) - B(E_2^P)| + \left| \sigma_a(E_1^P) - \sigma_a(E_2^P) \right| |B(E_1^P)|, \\ &\leq C |I_2^P - I_1^P| + L_a [|I_1^P| + |B(E_1^P)|] \times |E_2^P - E_1^P| + C |B(E_1^P) - B(E_2^P)|. \end{aligned}$$

Let us first study the  $B(E_1^P) = aT^4(E_1^P)$  term. As  $\frac{dT}{dE} = \frac{1}{\rho C_v(T)} > 0$ ,  $E \rightarrow T(E)$  is strictly increasing and  $T \rightarrow B(T) = aT^4$  or  $E \rightarrow aT^4(E)$  are strictly increasing: nothing prevents them from going to infinity. Let us assume that

**Hypothesis 2.**  $\exists \beta_0, \beta_1$  such that  $0 < \beta_0 \leq \rho C_v \leq \beta_1 < \infty$ .

Then  $dE = \rho C_v dT$  ensures that  $E - E_0 = \int_{T_0}^T \rho C_v(\alpha) d\alpha$  so that

$$\begin{aligned} \beta_0 &\leq \rho C_v(\alpha) \leq \beta_1, \\ \beta_0(T - T_0) &\leq \int_{T_0}^T \rho C_v(\alpha) d\alpha \leq \beta_1(T - T_0), \\ \beta_0(T - T_0) &\leq E - E_0 \leq \beta_1(T - T_0). \end{aligned}$$

As a consequence, we have  $|T - T_0| \leq \frac{1}{\min(\beta_0, \beta_1)} |E - E_0|$ . Let us introduce  $\beta = \min(\beta_0, \beta_1)$ , we then have

$$\begin{aligned} |B(E_1^P) - B(E_2^P)| &= a |T^4(E_1^P) - T^4(E_2^P)|, \\ &= a |T(E_1^P) - T(E_2^P)| \times |T(E_1^P) + T(E_2^P)| (T^2(E_1^P) + T^2(E_2^P)), \\ &\leq \frac{a}{\beta} |E_1^P - E_2^P| \times |T(E_1^P) + T(E_2^P)| (T^2(E_1^P) + T^2(E_2^P)). \end{aligned}$$

Using the above inequality in (18) yields

$$\begin{aligned} |q(u_1) - q(u_2)| &\leq C |I_2^P - I_1^P| \\ &\quad + L_a \underbrace{\left[ |I_1^P| + |B(E_1^P)| + \frac{Ca}{\beta} |T(E_1^P) + T(E_2^P)| (T^2(E_1^P) + T^2(E_2^P)) \right]}_{K(E_1^P, E_2^P)} \times |E_2^P - E_1^P|. \end{aligned} \quad (19)$$

In order to have the desired Lipschitz property, we need  $K(E_1^P, E_2^P)$  to be finite and non-zero. For this, it is enough having

$$\min(|I_1^P|, |B(E_1^P)|, T(E_1^P), T(E_2^P)) > 0 \text{ and } \max(|I_1^P|, |B(E_1^P)|, T(E_1^P), T(E_2^P)) < \infty.$$

Using the fact that  $B(E) = B(T(E)) = aT^4(E)$ , we show that having

$$\boxed{\min(T(E^P)) > 0 \text{ and } \max(|I^P|, T(E^P)) < \infty,} \quad (20)$$

is proven enough.

**Remark 2.2.** In inequalities (20), the upperbounds may be easily earned: indeed, assume  $X$  is a uniform random variable on  $[0, 1]^{27}$ , then  $I^P$  and  $E^P$  are (nonlinear functions of) polynomials on a bounded interval and are consequently bounded. The upperbound in (20) mainly depends on  $E^P \rightarrow T(E^P) = T^P$ , just as the lower bound in (20). So, in a nutshell, the main effort consists in designing the truncation  $T^P = T(E^P)$  of the material temperature.

In order to impose relevant conditions on the truncation  $(I^P, E^P, T^P)$  for (8) to be a relevant reduced model for (5), we are going to intensively use the properties of system (5) and in particular the *maximum principle*. The maximum principle for (5) states that if  $\forall X \sim d\mathcal{P}_X$ , there exists  $E_m(X), E_M(X)$  such that

$$\begin{aligned} E_m(X) &\leq E^0(x, X) \leq E_M(X), \forall x \in \Omega, X \sim d\mathcal{P}_X, \\ B(E_m(X)) &\leq I^0(x, \omega, X) \leq B(E_M(X)), \forall x \in \Omega, \omega \in \mathbb{S}^2, X \sim d\mathcal{P}_X, \\ B(E_m(X)) &\leq I_b(t, \omega, X) \leq B(E_M(X)), \forall t \in [0, t^*], \omega \in \mathbb{S}^2, X \sim d\mathcal{P}_X, \end{aligned} \quad (21)$$

then

$$\begin{aligned} E_m(X) &\leq E(x, t, X) \leq E_M(X), \forall x \in \Omega, t \in [0, t^*], X \sim d\mathcal{P}_X, \\ B(E_m(X)) &\leq I(x, t, \omega, X) \leq B(E_M(X)), \forall x \in \Omega, t \in [0, t^*], \omega \in \mathbb{S}^2, X \sim d\mathcal{P}_X. \end{aligned} \quad (22)$$

We suggest using this *a priori* information to choose a relevant truncation  $I^P, E^P, T^P$  for (8).

---

<sup>27</sup>In practice we can always come back to such condition.

- From the maximum principle stated above, for any configuration of interest (i.e. for any set of initial and boundary conditions) we can define *a priori* bounds with respect to  $X \sim d\mathcal{P}_X$ :

$$E_m = \min_{X \sim d\mathcal{P}_X} (E_m(X), E_M(X)) \text{ and } E_M = \max_{X \sim d\mathcal{P}_X} (E_M(X), E_m(X)), \quad (23)$$

and use them in inequalities (20). Our candidate truncations  $(I^P, E^P, T^P)$  are going to be chosen according to both (20) and (23).

- Finally, the conditions for a wellposed reduced model are not very constraining: we only have to choose  $I^P, E^P, T^P$  such that

$$\begin{aligned} \lambda_0^I, \dots, \lambda_P^I, X &\rightarrow I^P(\lambda_0^I, \dots, \lambda_P^I, X), \text{ (no particular additional hypothesis, see remark 2.2),} \\ \lambda_0^E, \dots, \lambda_P^E, X &\rightarrow E^P(\lambda_0^E, \dots, \lambda_P^E, X), \text{ (no particular additional hypothesis, see remark 2.2),} \\ \lambda_0^E, \dots, \lambda_P^E, X &\rightarrow T^P(\lambda_0^E, \dots, \lambda_P^E, X), \text{ satisfies (20).} \end{aligned} \quad (24)$$

Final condition (24) on truncation  $(I^P, E^P, T^P)$  is general enough: many relevant, with respect to wellposedness, truncation candidates are at hand. We can not go through every of them. In this paper, we are mainly going to consider two different ones:

- the most classical and obvious one is certainly to use some gPC expansions for  $I^P$  and  $E^P$

$$\begin{aligned} I^P &= \sum_{l=0}^P \lambda_l^I \phi_l, \\ E^P &= \sum_{l=0}^P \lambda_l^E \phi_l, \\ T^P &= T(\max(E_m, E^P)), \end{aligned} \quad (25)$$

together with a truncation of  $T^P$  inspired by the one in [17].

**Remark 2.3.** *In truncation (25), we did not exactly take advantage of all the a priori knowledge the maximum principle offers. We can also limit  $I$  and  $E$  from above, for example with*

$$\begin{aligned} I^P &= \min(\max(I_m, \sum_{l=0}^P \lambda_l^I \phi_l), I_M), \\ E^P &= \min(\max(E_m, \sum_{l=0}^P \lambda_l^E \phi_l), E_M). \end{aligned} \quad (26)$$

*These strategies will be considered and investigated later on. With truncation (25), we wanted to insist on the less stringent conditions for our reduced model to be wellposed.*

With truncations (25) and even (26), if the gPC expansion of  $E$  is accurate enough, then  $E^P$  remains superior to  $E_m$  and  $T^P = T(E^P) > T_m > 0$ : in every configurations in which  $T^P = T(E^P)$  holds, and using truncations (25) or (26) is equivalent to applying sG-gPC (hence, in this case  $\lambda_k^I = I_k = \int I \phi_k d\mathcal{P}_X$ ,  $\lambda_k^E = E_k = \int E \phi_k d\mathcal{P}_X$ ,  $\forall k \in \{0, \dots, P\}$ ). Of course, as soon as the gPC expansions do not satisfy bounds (23), the truncation is activated and the methodology becomes singular, different from sG-gPC. This strategy can be compared to the one presented in [6]: in [6], an hyperbolicity-preserving strategy is suggested. The idea is to limit the gPC expansion by limiting the fluctuations around the mean when leading to excursions of the hyperbolicity set. The  $\theta$ -limitation used in [6] leads to a truncation of the form

$$\begin{aligned} I^P &= I_0 + \theta^I \sum_{l=1}^P \lambda_l^I \phi_l \text{ with } \theta^I \text{ such that } I_m \leq I^P \leq I_M, \\ E^P &= E_0 + \theta^E \sum_{l=1}^P \lambda_l^E \phi_l \text{ with } \theta^E \text{ such that } E_m \leq E^P \leq E_M. \end{aligned} \quad (27)$$

With truncations (25), (26) and (27), we, in a sense, designed several *accretive-preserving* strategies.

- Another possibility is to apply IPM, see [8, 36]. We can choose what we call an entropy<sup>28</sup>  $s$  defined as

$$s(\alpha) = (\alpha - \alpha_m) \ln(\alpha - \alpha_m) - \alpha + \alpha_m + (\alpha_M - \alpha) \ln(\alpha_M - \alpha) - \alpha_M + \alpha, \forall \alpha \in \{I, E\}. \quad (28)$$

The above function is commonly called the Bounded-Barrier entropy and is thoroughly studied in [1]. The above entropy  $s$  is strictly convex in  $[E_m, E_M]$ . Minimising entropy (28) under constraints  $E_0, \dots, E_P$  (defined as the gPC coefficients of the expansion) leads to looking for  $I^P, E^P, T^P$  as

$$\begin{aligned}
I^P(\lambda_0^I, \dots, \lambda_P^I, X) &= \frac{I_m + I_M \exp \left[ \sum_{l=0}^P \lambda_l^I \phi_l(X) \right]}{1 + \exp \left[ \sum_{l=0}^P \lambda_l^I \phi_l(X) \right]} \in [I_m, I_M], \\
\text{with } \lambda_0^I, \dots, \lambda_P^I \text{ such that } &\begin{cases} I_0 = \int I^P(\lambda_0^I, \dots, \lambda_P^I, X) \phi_0(X) d\mathcal{P}_X, \\ \dots, \\ I_k = \int I^P(\lambda_0^I, \dots, \lambda_P^I, X) \phi_k(X) d\mathcal{P}_X, \\ \dots, \\ I_P = \int I^P(\lambda_0^I, \dots, \lambda_P^I, X) \phi_P(X) d\mathcal{P}_X, \end{cases} \\
E^P(\lambda_0^E, \dots, \lambda_P^E, X) &= \frac{E_m + E_M \exp \left[ \sum_{l=0}^P \lambda_l^E \phi_l(X) \right]}{1 + \exp \left[ \sum_{l=0}^P \lambda_l^E \phi_l(X) \right]} \in [E_m, E_M], \\
\text{with } \lambda_0^E, \dots, \lambda_P^E \text{ such that } &\begin{cases} E_0 = \int E^P(\lambda_0^E, \dots, \lambda_P^E, X) \phi_0(X) d\mathcal{P}_X, \\ \dots, \\ E_k = \int E^P(\lambda_0^E, \dots, \lambda_P^E, X) \phi_k(X) d\mathcal{P}_X, \\ \dots, \\ E_P = \int E^P(\lambda_0^E, \dots, \lambda_P^E, X) \phi_P(X) d\mathcal{P}_X, \end{cases} \\
T^P(\lambda_0^E, \dots, \lambda_P^E, X) &= T(E^P(\lambda_0^E, \dots, \lambda_P^E, X)).
\end{aligned} \quad (29)$$

Truncation strategy (29) is closely related to the maximum-principle preserving strategy suggested in [1] for scalar nonlinear hyperbolic equations. It does define an accretive-preserving strategy as, by construction, it leads to having  $0 < T_m \leq T^P$ . The main difference with the accretive-preserving strategies of the previous bullet is that with truncation (29), efforts are made for the truncations  $I^P, E^P, T^P$  to preserve the moments  $(I_k)_{k \in \{0, \dots, P\}}$  of  $I$  and the ones of  $(E_k)_{k \in \{0, \dots, P\}}$  of  $E$ . Care will be taken to highlight this difference in the numerical examples of section 5.

---

<sup>28</sup>Note that it is not a *mathematical* entropy as in [8].

Many other candidates could be imagined for the truncation. We suggest focusing on the three ((26), (27), (29)) above as they can be put quite easily in the same code framework and implemented within the same MC architecture (at least for the photonic system we consider in this paper).

We would like to end this section by two concluding remarks:

- first, we insist on the fact that the reduced model we build here is *reduced* in the sense one cannot expect to satisfy the maximum principle as stated for (5): indeed, the maximum principle for (5) states, for example for  $E$ , that

$$E_m(X) \leq E(x, t, X) \leq E_M(X), \forall x \in \Omega, t \in [0, t^*], X \sim d\mathcal{P}_X,$$

whereas the solution of reduced model (8) with condition (20) only ensures

$$E_m \leq E^P(x, t, X) \leq E_M, \forall x \in \Omega, t \in [0, t^*], X \sim d\mathcal{P}_X.$$

In other words, if for some particular  $X_0 \sim d\mathcal{P}_X$ ,  $E_m \ll E_m(X_0)$ , then nothing prevents  $E^P$  to predict  $E_m \leq E^P(x, t, X_0) \leq E_m(X_0) \leq E(x, t, X_0)$ .

**Remark 2.4.** *We can not expect the pointwise maximum property with the truncations/reduced models considered in this paper.*

- Finally, this section focused on conditions in order to build wellposed reduced models. But nothing ensures these reduced models will be relevant for the physical applications of interest.

**Remark 2.5.** *To convince oneself, consider a bad choice of  $E_m, E_M$  (for example chosen independently of the maximum principle). Then  $\mathcal{B}$  is still Lipschitz. The reduced model is wellposed: the solution exists and is unique. But it is not necessarily converging toward the solution of (5) as  $P$  grows. This will be emphasized in the numerical section 5.*

The theoretical proof of the converging behaviour of reduced model (8) with truncations respecting (20) is beyond the scope of this paper. We rely on the experiments of section 5 to numerically tackle this point.

Until now, we intensively focused on conditions the truncation  $I^P, E^P, T^P$  must satisfy for reduced model (8) to be wellposed. Now that they are characterised, we are going to focus, in sections 3–4, on how the coefficients  $\lambda^I, \lambda^E$  can be computed in practice.

### 3. The Implicit Semi-analog Monte-Carlo (ISMC) solver

As explained in the previous section, we, until now, focused on the conditions (20) the truncations  $\lambda^I, X \rightarrow I^P(\lambda^I, X), \lambda^E, X \rightarrow E^P(\lambda^E, X)$  must satisfy in order to make sure wellposed reduced models are built for any chosen gPC order  $P$ . We are going to focus on more practical considerations here and tackle how the  $\lambda^I, \lambda^E$  are calculated. In particular, we would like to be able to apply the MC-gPC solver and preserve some of its interesting properties, namely:

- MC-gPC in [10] can be implemented with simple modifications of an existing MC solver. We would like to preserve this property. Now, if for the linear Boltzmann equation, every MC schemes are almost equivalent<sup>29</sup>, this is not necessarily the case for MC codes solving the

---

<sup>29</sup>The semi-analog (intensively used for neutronic applications) and the non-analog (intensively used in photonic applications) schemes are almost equivalent in term of accuracy, parallel efficiency, and code architecture for the linear Boltzmann equation [36].

photonic system (1). For such nonlinear problem, a relevant linearisation must be selected before relying on an MC scheme. The choice of the linearisation considerably affects the structure of the code and the choice of the MC scheme. Several linearisations can be chosen (explicit MC [56], IMC [57, 58], tilted-IMC [59, 29, 27], SMC [60], ISMC [9]). We would like, the more possible, to preserve the relevant properties of the chosen linearisation once combined with MC-gPC.

- In [10], the author put forward the fact that the cost of the *tracking* of the uncertain MC particles is relatively independent of the stochastic dimension  $Q$ . What is strongly sensitive to the number of uncertain parameters  $Q$  is the *tallying* phase: each uncertain MC particle contribution must be tallied in an array of size  $P = (p_{1D} + 1)^Q$  in each cell. The highest the stochastic dimension  $Q$  or the polynomial order per direction  $p_{1D}$ , the more costly the tallying phase and the parallel reduction. So for MC-gPC to be efficient, we need to be able to avoid the more possible the number of tallying phases and of parallel reductions. This can be done in practice if we can take large enough time steps. For this reason, explicit MC solvers [56, 60] are discarded (note that they usually are, even in a deterministic context).
- Furthermore, as explained in the previous point, the number of coefficients per cell  $P$  can be large. This can have a strong impact on the memory consumption, especially if the linearisation/solver demands a large number of cells in order to reconstitute accurate results. For this reason, the seminal IMC [57] linearisation is discarded.

As a result, amongst the list of possible linearisations on which MC-gPC is going to be based, only remains the recent nssIMC solver [58], tilted IMC ones [59, 28, 25, 59, 27] and the ISMC one [9]. In the following lines, we explain why we think ISMC is the best candidate and which interesting properties of ISMC we would like to preserve for our ISMC+MC-gPC=ISMC-gPC solver.

The ISMC solver is based on a particular linearisation of (4) which integrates the source term  $\sigma_a B(T)$  into the scattering part. It is based on rewriting (4) with respect to  $E$  as (we drop the dependences for the sake of conciseness)

$$\begin{cases} \frac{1}{c} \partial_t I + \omega \cdot \nabla I + \sigma_t I = \sigma_a \eta(T(E)) E + \sigma_s \int I d\omega', \\ \partial_t E = \int c \sigma_a (I - \eta(T(E)) E) d\omega'. \end{cases} \quad (30)$$

In (30), we have  $B(T(E)) = \frac{B(T(E))}{E} E = \eta(T(E)) E$ . If we now introduce variable  $e(x, t, \omega)$  defined by  $\int e(x, t, \omega) d\omega = E(x, t)$ , system (30) can then be rewritten in term of unknowns  $(I, e)$  as:

$$\begin{cases} \frac{1}{c} \partial_t I + \omega \cdot \nabla I = +\sigma_a \eta(T(E)) \int e d\omega' - \sigma_t I + \sigma_s \int I d\omega', \\ \partial_t e = -c \sigma_a \eta(T(E)) e + \int c \sigma_a I d\omega'. \end{cases} \quad (31)$$

This system is still nonlinear. The ISMC solver linearises (31) by choosing and fixing astutely  $\eta$  during time step  $t \in [t^n, t^n + \Delta t]$ . The equation satisfied by  $\eta$  is

$$\partial_t \eta(T(E)) = \zeta(T(E)) c \sigma_a \left( \int I - \eta(T(E)) E \right), \text{ with } \zeta(T(E)) = E \eta'(T(E)). \quad (32)$$



In [9], an **explicit-implicit** time discretization is chosen for  $\zeta^n, \eta^{n+1}$  leading to

$$\partial_t \eta = \zeta^n c \sigma_a^n \left( \frac{1}{E} \int I d\omega - \eta^{n+1} \right). \quad (33)$$

Once the above expression integrated and inversed, we obtain an estimation of  $\eta^{n+1}$  given by

$$\eta^{n+1} = \eta^n \chi^n + (1 - \chi^n) \frac{1}{E} \int I, \quad (34)$$

with  $\chi^n = \frac{1}{1 + c \sigma_a^n \zeta^n \Delta t}$  called the *modified Fleck factor* [57, 9]. With this approximation, one can rewrite the transport equation using the new estimation of  $\eta$  on the time step:

$$\begin{cases} \frac{1}{c} \partial_t I + \omega \cdot \nabla I + \sigma_t^n I & = \chi^n \sigma_a^n \eta^n \int e + ((1 - \chi^n) \sigma_a^n + \sigma_s^n) \int I d\omega', \\ \partial_t e & + \chi^n c \sigma_a^n e = \chi^n c \sigma_a^n \int I d\omega. \end{cases} \quad (35)$$

System (35) is closed, linear, explicit (i.e. only quantities at time  $t^n$  appear) and exactly conservative in total energy in time step  $[t^n, t^n + \Delta t]$ .

**Remark 3.1 (First point in favor of ISMC to be combined with MC-gPC).** *As briefly explained in the above lines, ISMC relies on an implicit hypothesis. In practice, this astute implicit strategy allows taking bigger time step with respect to the (explicit) SMC scheme of [60] (which is unaffordable) and bigger time steps than for nssIMC [58]. In [61, 35], some HPC studies have been carried out with respect to the replication domain<sup>30</sup> parallel strategy (intensively used in [10] in an MC-gPC context). One of the main conclusion is that in order to make sure replication domain is efficient, one has to be able to take big time steps. Otherwise, the rythm of the parallel reduction (i.e. communication) may become prohibitive, especially if the vector to be reduced is of important size ( $P = (p_{1D} + 1)^Q$  here).*

Let us now build a new unknown  $\psi(t, x, \omega, v)$  for system (35), depending on one more dimension and on unknowns  $(I, e)$  solutions of (35). Variable  $v$  is chosen such that  $\psi(x, t, \omega, v) = I(x, t, \omega) \delta_c(v) + e(x, t) \delta_0(v)$ <sup>31</sup>. In fact,  $v$  is nothing more than a velocity which can be  $c$  for photons or 0 for matter. The equation satisfied by  $\psi$  is given by

$$\partial_t \psi(x, t, \omega, v) + v \omega \nabla \psi(x, t, \omega, v) + c \Sigma_t^n(x, v) \psi(x, t, \omega, v) = \int_{\mathcal{V}} \int c \Sigma_s^n(x, v, v') \psi(x, t, \omega', v') dv' d\omega', \quad (36)$$

<sup>30</sup>Replication domain consists in replicating the geometry on several processors and tracking several MC particles populations with different initial seeds in every replicated domains. At the end of the time steps, the contribution of every processors are averaged. This parallel strategy is particularly well suited to MC codes, taking advantage of the independence of the MC particles.

<sup>31</sup>In the latter expression,  $\delta_0, \delta_c$  are such that

$$\int_{\{V\}} \delta_c(v) dv = \delta_{V,c} \text{ and } \int_{\{V\}} \delta_0(v) dv = \delta_{V,0},$$

where  $\delta_{V,k}$  is the Kronecker symbol, i.e. is such that  $\delta_{V,k} = 0$  if  $V \neq k$  and  $\delta_{V,k} = 1$  if  $V = k$  and  $\{V\}$  denotes the singleton  $V$ .

where  $\mathcal{V} = \{0, c\}$  and

$$\begin{aligned}
\Sigma_t^n(x, v) &= \sigma_t \delta_c(v) + \sigma_a^n(x) \chi^n(x) \eta^n(x) \delta_0(v), \\
\Sigma_s^n(x, v, v') &= \Sigma_s^n(x, v) P_s^n(x, v, v'), \\
\Sigma_s^n(x, v) &= \chi^n(x) \sigma_a^n(x) \eta^n(x) \delta_c(v) + ((1 - \chi^n(x)) \sigma_a^n(x) + \sigma_s^n(x)) \delta_c(v) \\
P_s^n(x, v, v') &= \delta_0(v) \delta_c(v') + \delta_c(v) \frac{[\sigma_s^n(x) + (1 - \chi^n(x)) \sigma_a^n(x)] \delta_c(v') + \sigma_a^n(x) \chi^n(x) \delta_0(v')}{\sigma_t^n(x)}.
\end{aligned}$$

Note that in the above expression, the dependence with respect to  $x$  is piecewise constant per cell: for example, introduce the grid/set of  $N_x$  non-overlapping cells  $\mathcal{D} = \cup_{i=1}^{N_x} \mathcal{D}_i$ , then  $\sigma_a^n(x) = \sum_{i=1}^{N_x} \sigma_{a,i}^n \mathbf{1}_{\mathcal{D}_i}(x)$  where  $\mathbf{1}_{\mathcal{D}_i}(x)$  denotes the indicatrix of cell  $\mathcal{D}_i$ . The same applies for  $\eta^n(x)$ ,  $\chi^n(x)$ ,  $\sigma_s^n(x)$ ,  $\sigma_t^n(x)$ .

The identification of  $\Sigma_t^n, \Sigma_s^n, P_s^n$  as above allows being in the conditions of theorem 3.2.1 of [24]: it ensures we can build a converging<sup>32</sup> MC scheme toward system (35) on time step  $[t^n, t^n + \Delta t]$ . Equation (36) has the structure of a linear (multigroup) transport equation just as the very first equation in [10, 30] on which MC-gPC is derived. More importantly in this paper, it allows being in the conditions of theorem 1 of [30] for (36) for which spectral convergence is ensured.

**Remark 3.2 (Second point in favor of ISMC to be combined with MC-gPC).** *The second point in favor of ISMC concerns the fact that even once implicated, the system we have to solve during time step  $[t^n, t^n + \Delta t]$  can be put under the form of a linear Boltzmann equation (35): during any time step  $[t^n, t^n + \Delta t]$ , we are in the conditions of [10, 30] for which MC-gPC proved to be numerically and theoretically efficient. In other words, the material of [10, 30] is almost straightforward to apply. Of course, this does not mean spectral convergence necessarily holds for the full nonlinear uncertain problem nor that numerical efficiency will be easily earned. These properties will (only) numerically be investigated in section 5.*

Finally, system (4) is often studied together with a particular regime, commonly called the *equilibrium diffusion limit*. Introduce  $\delta \sim 0$ , a small parameter together with a characteristic length  $\mathcal{X}$ , a characteristic time  $\mathcal{T}$  and a characteristic collision rate  $\lambda$ . Then the equilibrium diffusion regime is characterised by

$$\begin{cases} c \frac{\mathcal{T}}{\mathcal{D}} = \mathcal{O}(\frac{1}{\delta}), \\ c \sigma \frac{\mathcal{T}}{\lambda} = \mathcal{O}(\frac{1}{\delta^2}). \end{cases} \quad (37)$$

Under condition (37), system (4) behaves, at leading order with  $\delta$ , like the nonlinear diffusion equation on  $\Phi_r(T_r) = aT_r^4 = \int I \, d\omega$  given by

$$\begin{cases} \partial_t(\Phi_r(T_r) + E(T_r)) - \nabla \cdot \left( \frac{c}{3\sigma_t} \nabla(\Phi(T_r)) \right) = \mathcal{O}(\delta), \\ \Phi_r(T_r) = \int I \, d\omega = B(T) + \mathcal{O}(\delta). \end{cases} \quad (38)$$

With  $\int B(T) \, d\omega = aT^4$  and  $\Phi_r(T_r) = aT_r^4$ , the second equation is equivalent to  $T = T_r$ : the radiative and matter temperatures are at equilibrium. For more details on the stakes of being able to capture efficiently this regime we refer to [11, 13, 62, 27, 25].

---

<sup>32</sup>with respect to the number  $N_{MC}$  of MC particles which, at this stage of the discussion, remain to be defined.

**Remark 3.3 (Third point in favor of ISMC to be combined with MC-gPC).** *The last point in favor of ISMC to be combined with MC-gPC concerns its fast (spatial) converging property in the equilibrium diffusion regime<sup>33</sup> [9]. We hope the ISMC-gPC solver to inherit this good property. This will be verified numerically in section 5.*

We finally end this section with few words on the general sketch of an ISMC solver. This will ease the identification of the modifications needed in order to modify the original ISMC code in order to implement the MC resolution of the gPC based reduced model of section 2. Algorithms 1–2–3 present the tracking phase. They assume a population of MC particles discretizing the initial solution has been built. We insist on the fact that both the radiation intensity and the matter energy are described by MC particles (of respectively, velocities  $c$  and  $0$ ).

---

<sup>33</sup>In particular, it is not sensitive to the *teleporation error*, see [25, 26, 27, 28, 29].

```

#BEGINNING OF TIME STEP [t^n, t^n + Δt]
for i ∈ {1, ..., N_x} do
  # Keep the material energy per cell of the previous time step in memory
  E_0^i = E^i
  #Set to zero the (mesh) arrays in which will be tallied the MC particle contributions
  U^i = 0, E^i = 0
end
for p ∈ {1, ..., N_MC} do
  set s_p = 0 #this will be the current time of particle p
  #i_p is such that 1_{Ω_{i_p}}(x_p(s_p)) = 1 (current cell for particle p)
  while s_p < Δt do
    if x_p ∉ D then
      | apply_boundary_conditions(x_p, s_p, v_p)
    end
    #sample the collision time from an uniform sampling U
    τ = - ln(U) / (cΣ(x_p, s_p, v_p, E_0^{i_p}))
    if τ > Δt then
      | #move the particle p, update s_p to end the treatment of the current particle
      | x_p = x_p + v_p ω_p × (t - τ), s_p ← t
      | #tally the contribution of particle p in the cell array in which it ends:
      | if v_p == c then
      | | U^{i_p} += w_p
      | end
      | else
      | | E^{i_p} += w_p
      | end
    end
    else
      | #move the particle p, update the life time of particle p
      | x_p ← x_p - v_p ω_p τ, s_p ← s_p + τ < t
      | #Sample the angle W' and 'group' of particle p after the collision
      | W', V' = sample_angle_and_group(x_p, s_p, ω_p, v_p, E_0^{i_p})
      | ω_p = W', v_p = V'
    end
  end
end

```

**Algorithm 1:** Pseudo-code for an ISMC implementation.

Algorithm 1 presents the tracking of the MC particles: note that the structure of the algorithm is close to the one described in [10] for the linear Boltzmann equation, see remark 3.2. In algorithm 1, the material energy in each cell is stored at the beginning of the time step and the radiation intensity and the material energy per cells arrays are initialised to 0. The particle contributions are going to be tallied in these arrays. Then comes the treatment of the MC particles: the life time of the particles is set to zero: it is going to be gradually incremented until reaching the value  $\Delta t$ . While  $s_p < \Delta t$ , a collision/emission time is sampled depending on the particle field (i.e. whether it is a photon or a material particle). The opacity used in order to sample the collision/emission time

depends on the nature of the particle: it is detailed in algorithm 2 and will be discussed later on.

**Data:**  $x_p, s_p, v_p, E_0^{i_p}$

**Result:**  $\Sigma(x_p, s_p, v_p, E_0^{i_p})$ , the opacity  $\Sigma$  seen by particle  $p$

**begin**

$\Sigma = 0$

    # $E_0^{i_p}$  is the material energy in cell  $i_p$  at the beginning of the time step.

$T^{i_p} = T(E_0^{i_p})$

**if**  $v_p == c$  **then**

        |  $\Sigma = \sigma_t(x_p, s_p, T^{i_p})$

**end**

**else**

        #Here,  $\eta$  is explicited in the case of a perfect gas

$$\eta = a \frac{(T^{i_p})^3}{\rho_{i_p} C_v^{i_p}}$$

        # $\chi^n$  is the modified Fleck factor [9]

$$\chi = \frac{1}{1 + 3c\sigma_a(x_p, s_p, T^{i_p})\eta\Delta t}$$

$$\Sigma = \eta\sigma_a(x_p, s_p, T^{i_p})\chi$$

**end**

**end**

**Algorithm 2:** Pseudo-code for the computation of the opacity seen by particle  $p$ .

If the particle collides/is emitted, its life time is updated, it moves of a distance<sup>34</sup>  $v_p\tau$  in the direction  $\omega_p$ , it changes of nature (from photon to matter or matter to photon) and its angle is sampled as detailed in algorithm 3. Finally, when  $s_p$  reaches  $\Delta t$ , the particle moves at its last location for the time step. The cell in which it ends is  $i_p$ , its contribution is tallied in the arrays  $E^{i_p}$  or  $I^{i_p}$  depending on the value of its velocity field  $v_p$ .

---

<sup>34</sup>For material particles,  $v_p = 0$ , for photons  $v_p = c$ .

**Data:**  $x_p, s_p, \omega_p, v_p$   
**Result:** angle  $W'$  and velocity  $V'$  of particle  $p$  after a collision

```

begin
  # sample_angle_and_group
  #  $E_0^{i_p}$  is the material energy in cell  $i_p$  at the beginning of the time step.
   $T^{i_p} = T(E_0^{i_p})$ 
  # Here,  $\eta$  is explicited in the case of a perfect gas
   $\eta = a \frac{(T^{i_p})^3}{\rho_{i_p} C_v^{i_p}}$ 
  #  $\chi^n$  is the modified Fleck factor [9]
   $\chi = \frac{1}{1 + 3c\sigma_a(x_p, s_p, T^{i_p})\eta\Delta t}$ 
  if  $v_p = c$  then
    #  $U'$  is an uniform random variable in  $[0, 1]$ 
    if  $U' \times \sigma_t(x_p, s_p, T^{i_p}) < \sigma_s(x_p, s_p, T^{i_p}) + (1 - \chi)\sigma_a(x_p, s_p, T^{i_p})$  then
      # Scattering:  $v_p$  does not change, it remains equal to  $c$ 
       $W' = \text{sample\_scattering\_angle}(x_p, s_p, \omega_p, T^{i_p})$ 
    end
  else
    # Absorption:  $v_p = c$  becomes  $v_p = 0$ 
     $V' = 0$ 
  end
  else
    # Emission:  $v_p = 0$  becomes  $v_p = c$ 
     $W' = \text{sample\_emission\_angle}(x_p, s_p, \omega_p, T^{i_p})$ 
     $V' = c$ 
  end
end
end

```

**Algorithm 3:** Pseudo-code for the computation of the angle and group of particle  $p$  after a collision.

Algorithm 2 focuses on the computation of the opacity  $\Sigma$  used to sample the collision/emission time for any MC particle. First, a material temperature  $T^{i_p}$  in cell  $i_p$  is built from  $E_0^{i_p}$ . Then, depending on the nature (photon or matter) of the particle, the total opacity (photon case) or an artificial scattering (matter case) opacity is reconstructed from the modified Fleck factor  $\chi$  and the local temperature  $T^{i_p}$ . The reconstructed opacity is then simply used to sample a classical exponential collision/emission time  $\tau$  during the tracking phase (see algorithm 1). Algorithm 3 presents how an MC particle is scattered/absorbed (photon) or emitted (matter) when  $\tau < \Delta t$ . The computation once again needs the reconstruction of a local temperature in the cell  $i_p$  in which MC particle  $p$  encounters a collision/emission/absorption. Of course, this reconstruction (together with the one of the modified Fleck factor) is redundant and can be mutualised in practice with what is done in algorithm 2. A photon particle can either encounter an artificial scattering or an absorption (i.e. become matter) whereas matter can only be emitted. If a photon is (artificially) scattered, its angle is resampled (see function `sample_scattering_angle`, the sampling can be isotropic or not without additional difficulties) and it remains in a photon state. If it is absorbed, it stops at the collision location and becomes matter. A matter particle can be emitted: its angle is resampled (see function `sample_emission_angle`, the emission can be isotropic or not without additional difficulties)

and it becomes a photon particle.

This ends the brief description of the general sketch of an ISMC resolution. We now have to find a way to combine efficiently ISMC and MC-gPC while keeping the aforementioned properties (fast gPC convergence, big stable time steps, equilibrium diffusion limit, simplest modifications possible).

#### 4. Combining ISMC [9] to MC-gPC [10, 30, 31, 33]: the ISMC-gPC solver

In this section, we highlight the modification one has to perform to an ISMC implementation in order to take uncertainties into account *on-the-fly* during the MC resolution. The material is a combination of both the MC-gPC scheme described in [10] for the transport equation and of the truncation of the previous section 2. In order to explain how we combine the construction of a wellposed reduced model (14) (see section 2) and the ISMC scheme (see section 3), we need to insist on two points:

- on how an ISMC implementation is modified to become an ISMC-gPC with minimal development effort, just as in [10].
- On how the truncation preserving wellposedness is introduced.

The best way to do so, in our opinion, remains to highlight where we need to modify the algorithms of section 3 to describe both the reduced model and the MC scheme.

```

#BEGINNING OF TIME STEP [t^n, t^n + Δt]
for i ∈ {1, ..., N_x} do
  for k ∈ {0, ..., P} do
    # Keep the previous material energy gPC coefficients per cell in memory
    E_{k,0}^i = E_k^i
    #Set to zero the (mesh) arrays in which will be tallied the MC particle contributions
    U_k^i = 0, E_k^i = 0
  end
end
end
for p ∈ {1, ..., N_{MC}} do
  set s_p = 0 #this will be the current time of particle p
  #i_p is such that 1_{Ω_{i_p}}(x_p(s_p)) = 1 (current cell for particle p)
  while s_p < Δt do
    if x_p ∉ D then
      | apply_boundary_conditions(x_p, s_p, v_p, X_p)
    end
    #sample the collision time from an uniform sampling U
    τ = - ln(U) / cΣ(x_p, s_p, v_p, E_{0,0}^{i_p}, ..., E_{P,0}^{i_p}, X_p)
    if τ > Δt then
      #move the particle p, update s_p to end the treatment of the current particle
      x_p = x_p + v_p ω_p × (t - τ), s_p ← t
      #tally the contribution of particle p in the in which it ends: for k ∈ {0, ..., P} do
        if v_p == c then
          | U_k^{i_p} += w_p φ_k(X_p)
        end
        else
          | E_k^{i_p} += w_p φ_k(X_p)
        end
      end
    end
    end
    else
      #move the particle p, update the life time of particle p
      x_p ← x_p - v_p ω_p τ, s_p ← s_p + τ < t
      #Sample the angle W' and 'group' of particle p after the collision
      W', V' = sample_angle_and_group(x_p, s_p, ω_p, v_p, E_{0,0}^{i_p}, ..., E_{P,0}^{i_p}, X_p)
      ω_p = W'
      v_p = V'
    end
  end
end
end

```

**Algorithm 4:** Pseudo-code for an ISMC-gPC implementation (only the tracking of MC particles is detailed).

Algorithm 4 presents the modifications to algorithm 1 that one has to do in order to transform its ISMC implementation into an ISMC-gPC one. Those modifications are highlighted in blue. First,



just as in [10], an uncertain MC particle has an additional field  $X_p \sim d\mathcal{P}_X$  besides the more classical ones which are  $x_p, s_p, v_p, \omega_p$ . The main modification of this tracking step consists in dealing with arrays instead of scalars, i.e. dealing with  $(I_k^i, E_k^i)_{k \in \{0, \dots, P\}, i \in \{1, \dots, N_x\}}$  instead of  $(I^i, E^i)_{i \in \{1, \dots, N_x\}}$ . Note that in  $(I_k^i, E_k^i)_{k \in \{0, \dots, P\}, i \in \{1, \dots, N_x\}}$  are tallied the projections on the components of the gPC basis  $(\phi_k(X))_{k \in \{0, \dots, P\}}$  of the contribution of each uncertain MC particle at the end of the time step. The calls to `apply_boundary_conditions`,  `$\Sigma$`  and `sample_angle_and_group` have one more argument  $X_p$  and the scalar dependency with respect to  $E^{i_p}$  becomes vectorial, with respect to  $E_0^{i_p}, \dots, E_P^{i_p}$ .

**Data:**  $x_p, s_p, v_p, E_{0,0}^{i_p}, \dots, E_{P,0}^{i_p}, X_p$

**Result:**  $\Sigma(x_p, s_p, v_p, E_{0,0}^{i_p}, \dots, E_{P,0}^{i_p}, X_p)$ , the opacity  $\Sigma$  seen by particle  $p$

**begin**

$\Sigma = 0$

    # $E_{0,0}^{i_p}, \dots, E_{P,0}^{i_p}$ : gPC coefficients of the material energy in cell  $i^p$  at the beginning of the time step.

$T^{i_p} = T^P(E_{0,0}^{i_p}, \dots, E_{P,0}^{i_p}, X_p)$

**if**  $v_p == c$  **then**

        |  $\Sigma = \sigma_t(x_p, s_p, T^{i_p}, X_p)$

**end**

**else**

        #Here,  $\eta$  is explicited in the case of a perfect gas

$$\eta = a \frac{(T^{i_p})^3}{\rho_{i_p}(X_p) C_v^{i_p}(X_p)}$$

        # $\chi$  is the modified uncertain Fleck factor [9]

$$\chi = \frac{1}{1 + 3c\sigma_a(x_p, s_p, T^{i_p}, X_p)\eta\Delta t}$$

$$\Sigma = \eta\sigma_a(x_p, s_p, T^{i_p}, X_p)\chi$$

**end**

**end**

**Algorithm 5:** Pseudo-code for the computation of the uncertain opacity seen by particle  $p$ .

Let us now explicit the dependences of  $\Sigma$  with respect to  $(E_{k,0}^i)_{k \in \{0, \dots, P\}, i \in \{1, \dots, N_x\}}$ . These are detailed in algorithm 6 which focuses on the modifications one has to perform to algorithm 3: the truncation of the material temperature  $T^P$  is used mainly in this function. From the gPC coefficients of the material energy  $(E_{k,0}^i)_{k \in \{0, \dots, P\}, i \in \{1, \dots, N_x\}}$  at the beginning of the time step (i.e. the scheme remains explicit) and the uncertain parameter  $X_p$  of particle  $p$ , an uncertain material temperature  $T^{i_p}$  is reconstructed. Once this step done, it remains to add one additional dependence with respect to  $X_p$  on  $C_v, (\sigma_\alpha)_{\alpha \in \{s,t,a\}}, \rho$  depending on whether they are considered uncertain or not.

**Data:**  $x_p, s_p, \omega_p, v_p, E_{0,0}^{i_p}, \dots, E_{P,0}^{i_p}, X_p$

**Result:** angle  $W'$  and velocity  $V'$  of particle  $p$  after a collision

**begin**

```

# sample_angle_and_group
#  $E_{0,0}^{i_p}, \dots, E_{P,0}^{i_p}$  gPC coefficients of the material energy in cell  $i^p$  at the beginning of the
  time step.
 $T^{i_p} = T^P(E_{0,0}^{i_p}, \dots, E_{P,0}^{i_p}, X_p)$ 
# Here,  $\eta$  is explicited in the case of a perfect gas
 $\eta = a \frac{(T^{i_p})^3}{\rho_{i_p}(X_p) C_v^{i_p}(X_p)}$ 
#  $\chi$  is the uncertain modified Fleck factor [9]
 $\chi = \frac{1}{1 + 3c\sigma_a(x_p, s_p, T^{i_p}, X_p)\eta\Delta t}$ 
if  $v_p = c$  then
  #  $U'$  is an uniform random variable in  $[0, 1]$ 
  if  $U' \times \sigma_t(x_p, s_p, T^{i_p}, X_p) <$ 
     $\sigma_s(x_p, s_p, T^{i_p}, X_p) + (1 - \chi(x_p, s_p, T^{i_p}, X_p))\sigma_a(x_p, s_p, T^{i_p}, X_p)$  then
      # Scattering:  $v_p$  does not change, it remains equal to  $c$ 
       $W' = \text{sample\_scattering\_angle}(x_p, s_p, \omega_p, T^{i_p}, X_p)$ 
    end
  else
    # Absorption:  $v_p = c$  becomes  $v_p = 0$ 
     $V' = 0$ 
  end
  else
    # Emission:  $v_p = 0$  becomes  $v_p = c$ 
     $W' = \text{sample\_emission\_angle}(x_p, s_p, \omega_p, T^{i_p}, X_p)$ 
     $V' = c$ 
  end
end
end

```

**Algorithm 6:** Pseudo-code for the computation of the angle and group of uncertain particle  $p$  after a collision.

The dependences with respect to  $(E_{k,0}^i)_{k \in \{0, \dots, P\}, i \in \{1, \dots, N_x\}}$  of `sample_angle_and_group` are detailed in algorithm 6. The modifications needed for this function are of the same nature as for algorithm 5. Note that the lines computing  $\eta$  and  $\chi$  may be redundant as they can easily be mutualised during the computation of  $\Sigma$  and avoid few additional operations per treatment of an uncertain MC particle. Still, they are recalled here in order to ease the understanding of the algorithm.

As expected (see remark 3.2), the modifications needed to go from an ISMC implementation to an ISMC-gPC one are fairly the same as the one in [10] for the linear Boltzmann equation. In a sense, our first objective to suggest as simple modifications as in [10] to an ISMC solver is reached. The truncation is only needed locally during the tracking phase. In the next section, we suggest presenting few results obtained with the ISMC-gPC solver we just described.

## 5. Numerical results

In this last section, we present numerical results obtained with the ISMC-gPC solver we described in the previous sections. The results with ISMC-gPC are compared with the results obtained with non-intrusive ISMC (ni-ISMC) simulations. ni-ISMC supposes choosing a quadrature rule  $(X_i, w_i)_{i \in \{1, \dots, N\}}$  consistently discretising  $(X, d\mathcal{P}_X)$  and running  $N$  times an ISMC code solving (4) at the prescribed points  $(X_i)_{i \in \{1, \dots, N\}}$ . Once the runs performed, it only remains to post-process the results  $((I(x, t, \omega, X_i), E(x, t, X_i), T(x, t, X_i), T_r(x, t, X_i)), w_i)_{i \in \{1, \dots, N\}}$  in order to estimate the statistics of the observables of interest  $X \rightarrow (I(x, t, \omega, X), E(x, t, X), T(x, t, X), T_r(x, t, X))$  by numerical integration. Now, MC codes are known to be efficient but costly and relying on many ( $N \gg 1$ ) runs is not always possible. As a consequence, the benchmarks of this section remain quite simple in order to make sure we can have converged statistics<sup>35</sup> with relatively small<sup>36</sup>  $N$ .

The next benchmarks are, to our opinion, progressive in complexity. The first test-cases correspond to uncertain relaxation problems. For these, it is easy producing accurate reference solutions and performing convergence studies: indeed, in these cases, problem (5) degenerates toward solving several times a set of Ordinary Differential Equations (ODEs). The second one is built from the commonly known *Heaviside* problem [14, 9, 58]. It is made uncertain by considering a fluctuating opacity. This test-case allows testing the capabilities of the ISMC-gPC solver with respect to the *equilibrium diffusion limit*. The last test-case is an uncertain Marshak wave on which are going to be tested and analysed several of the admissible (with respect to wellposedness) truncations of section 2. For each benchmark, care is taken to carry out fair performance studies. In all the test-cases, the opacities and equations of state satisfy hypothesis 1 and 2.

### 5.1. Uncertain infinite medium problems: uncertain relaxations

In this section, we consider uncertain infinite medium problems without scattering ( $\sigma_s \equiv 0$ ). In such conditions, (5) degenerates toward the Stochastic system of ODEs

$$\begin{cases} \partial_t I(t, X) = c\sigma_a(E(t, X), X)(B(E(t, X)) - I(t, X)), \\ \partial_t E(t, X) = c\sigma_a(E(t, X), X)(I(t, X) - B(E(t, X))), \end{cases} \quad (39)$$

with  $c = 1$ ,  $I(t = 0, X) = I_0(X) = I_0 = 1$ ,  $E(t = 0, X) = E_0(X) = E_0 = 10^{-3}$ . Furthermore, we assume that  $C_v(T, X) = C_v(X)$  (uncertain perfect gas) so that  $E(t, X) = C_v(X)T(t, X)$ . Finally,  $B(E(T)) = aT^4$  with  $a = 1$ .

System (39) can easily be solved with an explicit Euler scheme with a fine time discretisation for several values of  $X \sim d\mathcal{P}_X$ . In practice we take  $\Delta t = 10^{-6}$  for each  $(X_i, w_i)_{i \in \{1, \dots, N\}}$  with  $N = 15$  Gauss-Legendre quadrature points. This set-up is proven enough to produce accurate results on the means and variances of  $I, E, T, T_r$ . In the next paragraphs and figures, we compare the reference results obtained with non-intrusive ISMC (ni-ISMC $_N$ ) with the  $N$  Gauss-Legendre points to ISMC-gPC $_P$ , for several time discretisations  $\Delta t$ , several number of  $N_{MC}$  of MC particles and several gPC orders  $P$ .

In section 5.1.1, we consider an uncertain absorption opacity  $\sigma_a(E(t, X), X) = \sigma_a(X) = \bar{\sigma}_a + \hat{\sigma}_a X$  with  $\bar{\sigma}_a = 1$ ,  $\hat{\sigma}_a = \frac{1}{2}$  and<sup>37</sup>  $X \sim \mathcal{U}([-1, 1])$  together with a deterministic heat capacity

<sup>35</sup>i.e. reliable references.

<sup>36</sup>Note that this is also a point in favor of an easy reproducibility of the results of this paper.

<sup>37</sup> $X \sim \mathcal{U}([-1, 1])$  must be read:  $X$  follows an uniform distribution on  $[-1, 1]$ .

$C_v(X) = C_v = \frac{3}{2}$ . In section 5.1.2, the opacity is deterministic, given by  $\sigma_a(E(t, X), X) = \sigma_a = \bar{\sigma}_a$  with an uncertain heat capacity  $C_v(X) = \bar{C}_v + \hat{C}_v X$  with  $\bar{C}_v = \frac{3}{2}$ ,  $\hat{C}_v = \frac{1}{2}$  and  $X \sim \mathcal{U}([-1, 1])$ . In section 5.1.3, a 2D ( $Q = 2$ ) stochastic test problem is considered in which the two previous uncertain parameters are combined.

### 5.1.1. Uncertain infinite medium problems: uncertain absorption opacity $\sigma_a$

In this section, we consider an uncertain absorption opacity  $\sigma_a(E(t, X), X) = \sigma_a(X) = \bar{\sigma}_a + \hat{\sigma}_a X$  with  $\bar{\sigma}_a = 1$ ,  $\hat{\sigma}_a = \frac{1}{2}$  and  $X \sim \mathcal{U}([-1, 1])$ . Figure 1 compares the results obtained from the reference

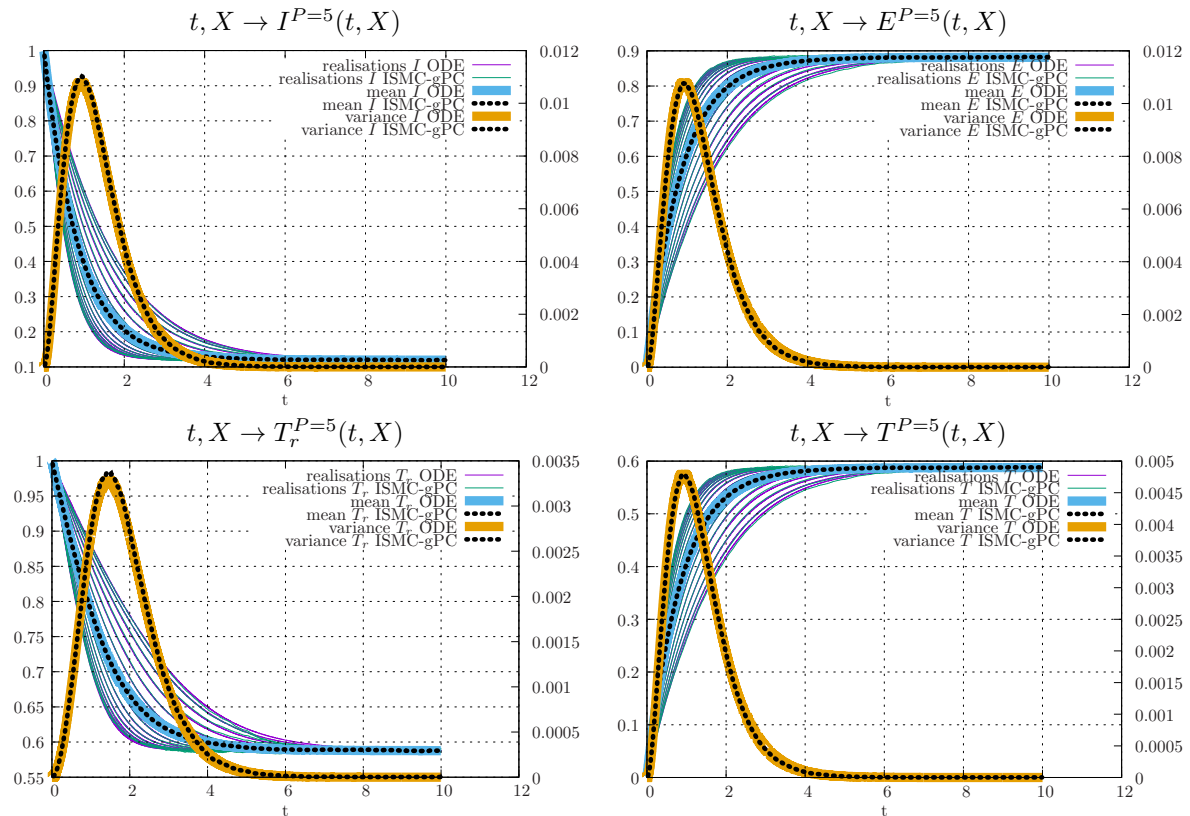


Figure 1: Mean, realisations at the  $N = 15$  Gauss-Legendre points and variance with respect to time of the ni-ISM approximation and the ISMC-gPC $_{P=5}$  ones. In particular, we present the results in term of intensity of radiation  $I$ , matter energy  $E$  and material and radiative temperatures  $T, T_r$ . The numerical parameters for ISMC-gPC are  $\Delta t = 10^{-2}$ ,  $N_{MC} = 10^6$ ,  $P = 5$ .

ni-ISM $_{N=15}$  and ISMC-gPC $_{P=5}$  for  $\Delta t = 10^{-2}$  and  $N_{MC} = 10^6$ . Several statistical quantities are displayed:

- the mean with respect to time of the radiation intensity, the matter energy and the material

and radiative temperatures:

$$\begin{aligned} t \rightarrow \mathbb{E}[I](t) &= \int I(t, X) d\mathcal{P}_X, \\ t \rightarrow \mathbb{E}[E](t) &= \int E(t, X) d\mathcal{P}_X, \\ t \rightarrow \mathbb{E}[T_r](t) &= \int T_r(t, X) d\mathcal{P}_X, \\ t \rightarrow \mathbb{E}[T](t) &= \int T(t, X) d\mathcal{P}_X. \end{aligned}$$

- The realisations  $t \rightarrow \alpha(t, X_i)$  for  $\alpha \in \{I, E, T, T_r\}$  at the  $N = 15$  Gauss Legendre points for ni-ISM and ISMC-gPC $_{P=5}$ .
- The variances of  $I, E, T, T_r$  with respect to time

$$\begin{aligned} t \rightarrow \mathbb{V}[I](t) &= \int (I(t, X) - \mathbb{E}[I](t))^2 d\mathcal{P}_X, \\ t \rightarrow \mathbb{V}[E](t) &= \int (E(t, X) - \mathbb{E}[E](t))^2 d\mathcal{P}_X, \\ t \rightarrow \mathbb{V}[T_r](t) &= \int (T_r(t, X) - \mathbb{E}[T_r](t))^2 d\mathcal{P}_X, \\ t \rightarrow \mathbb{V}[T](t) &= \int (T(t, X) - \mathbb{E}[T](t))^2 d\mathcal{P}_X. \end{aligned}$$

Note that the vertical axis on the right hand side must be used for the variances.

Let us first describe briefly the test-case: the uncertainty on the opacity makes the transient regime uncertain while leaving the stationary one deterministic: the variance drops to zero after  $t \approx 6$  here. The temperatures  $T_r, T$  are at equilibrium after this time, independently of the realisation  $X$ . For early times, there is a quite important variability of the different quantities. Now, the two different solvers (the non-intrusive and the intrusive one we suggest in this paper) are in excellent agreement: for the means, the variances, but also for the pointwise approximations at the Gauss-Legendre points.

Figure 2 presents a qualitative convergence study with respect to  $P \in \{0, \dots, 7\}$  on the radiation intensity  $I$ . For  $P = 0$ , figure 2 (top left) the ISMC-gPC $_{P=0}$  predicts a zero variance and only the mean  $t \rightarrow \mathbb{E}[I^{P=0}](t)$  is (nonetheless accurately) captured by the ISMC-gPC $_{P=0}$  approximation. For  $P = 1$ , figure 2 (top right), the mean  $t \rightarrow I(t, X)$  is still well captured and the solver predicts a non-zero variance. But it is not accurate enough to reconstitute equivalent results as ni-ISM. As  $P$  increases, the results are better and better in term of variance and pointwise realisations. As soon as  $P = 3$ , the results from the two solvers are not anymore discernable on the figures attesting for a qualitatively fast (spectral?) convergence of the ISMC-gPC solver with respect to  $P$ .

Let us now consider a more quantitative convergence study: to this purpose, in figure 3, we display the curves<sup>38</sup>  $P \rightarrow \ln(\|\alpha - \alpha^P\|)(t)$  for different  $\alpha \in \{I, E, T, T_r\}$ , different times  $t \in \{1, 2, 4\}$ , different values of the numerical parameters  $\Delta t \in \{10^{-2}, 10^{-3}\}$  and different MC discretisations  $N_{MC} \in \{10^4, 10^5, 10^6\}$ . Figure 3 (top left) presents a convergence study with respect to  $P$  on  $I, E, T, T_r$  at time  $t = 1$  for  $\Delta t = 10^{-3}$  and  $N_{MC} = 10^6$ . Independently of the observable of interest, i.e.  $\forall \alpha \in \{I, E, T, T_r\}$ , the curves  $P \rightarrow \ln(\|\alpha - \alpha^P\|)$  present two regimes: a converging regime for  $P \in \{0, \dots, 3\}$ , and a stagnating one for  $P > 3$ . The errors reach a plateau which is in agreement with the MC error, i.e. parameter  $N_{MC}$ . This is emphasized by figure 3 (top right) where<sup>39</sup>  $P \rightarrow \ln(\|T_r - T_r^P\|)$  is plotted for several MC discretisations  $N_{MC} \in \{10^4, 10^5, 10^6\}$ : the finer the MC resolution, the lower the stagnation plateau. Note that the MC error is known to be

<sup>38</sup>The norm is as defined in section 2, see (16).

<sup>39</sup>We insist that the behaviour is the same for the other observables  $I, E, T$ . Plotting them would be redundant.

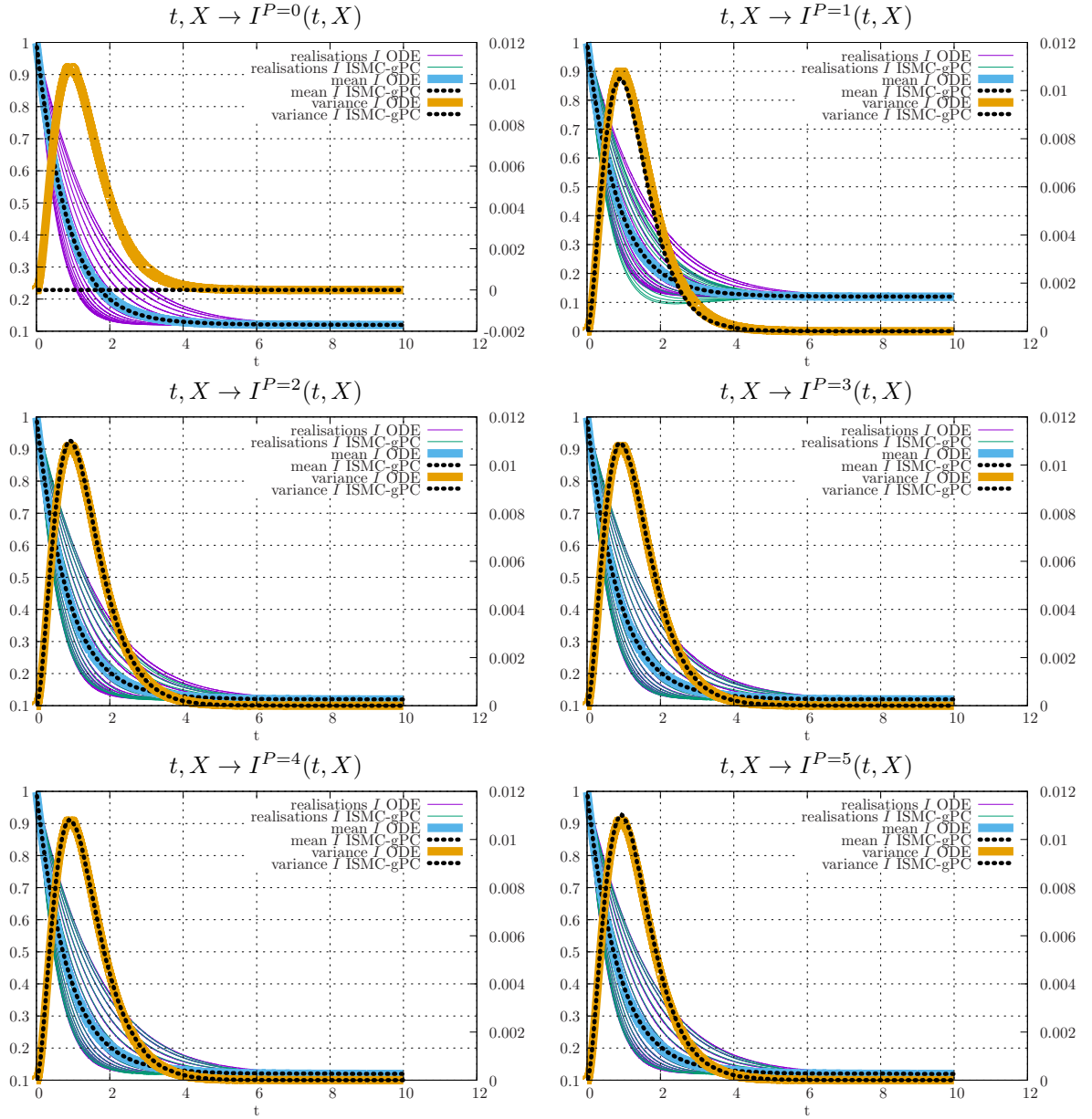


Figure 2: A qualitative convergence study with respect to  $P$  on the mean, realisations at the  $N = 15$  Gauss-Legendre points and variance with respect to time of the ISMC-gPC $_P$  approximations of the radiation intensity  $I$ . The remaining numerical parameters for ISMC-gPC are  $\Delta t = 10^{-2}$ ,  $N_{MC} = 10^6$ .

$\mathcal{O}(\frac{1}{\sqrt{N_{MC}}})$  and the errors on figure 3 are all stagnating around  $10^{-3}$  which is in agreement with the MC error being  $\mathcal{O}(\frac{1}{\sqrt{N_{MC}}}) \approx \frac{1}{\sqrt{10^6}} = 10^{-3}$ . Figure 3 (bottom left) presents the convergence curves<sup>40</sup>

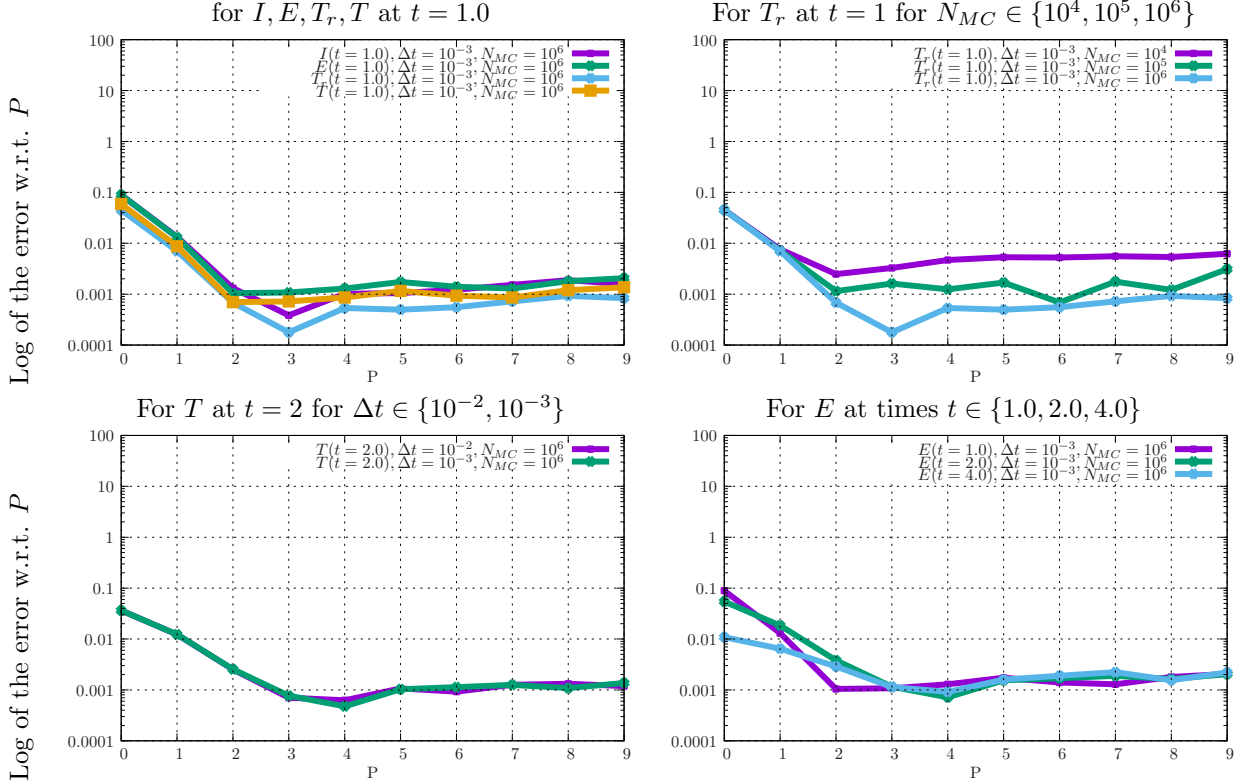


Figure 3: A quantitative convergence study with respect to  $P$  of the ISMC-gPC $_P$  approximations of  $I, E, T_r, T$  for several times  $t \in \{1.0, 2.0, 4.0\}$ , several time discretisations  $\Delta t \in \{10^{-2}, 10^{-3}\}$  and several number of MC particles  $N_{MC} \in \{10^4, 10^5, 10^6\}$ .

for  $P \rightarrow \ln(\|T - T^P\|)$  for several time discretisations  $\Delta t \in \{10^{-2}, 10^{-3}\}$  at time  $t = 2.0$ : the results of the convergence study are not sensitive to a change of time step: the time step is fine enough and is not the constraining numerical parameter here ( $N_{MC}$  is). Note that the time step is way coarser than for the explicit Euler resolution used as a reference: this is the purpose of ISMC<sup>41</sup> [9] to be able to provide affordable stable and accurate time steps for the resolution of (1). *We recover this desirable property for ISMC-gPC here.* Finally, figure 3 (bottom right) presents the convergence curves for  $E$  at several times. The spectral convergence, i.e. linear curve  $P \rightarrow \ln(\|E - E^P\|)$ , is discernable for early polynomial orders  $P \in \{0, \dots, 4\}$  before reaching the plateau related to the MC resolution  $\mathcal{O}(\frac{1}{\sqrt{N_{MC}}})$ . It is interesting noticing that this plateau is the same for every times

<sup>40</sup>We insist that the behaviour is the same for the other observables  $I, E, T_r$  at the other times  $t \in [0, t^* = 10]$ . Plotting them would be redundant.

<sup>41</sup>The 'I' of ISMC stands for 'Implicit'.

$t \in \{1.0, 2.0, 4.0\}$ : it attests that, at least for this problem, the ISMC-gPC solver is not sensitive to the long-term behaviour of gPC<sup>42</sup> (encountered numerically in several publications [63, 49, 30] and theoretically recovered in [30] for the uncertain linear Boltzmann equation in multiplicative media).

With this set of convergence studies, we notice that the MC accuracy is the constraining one: *MC-gPC is interesting in an MC context because the fast gPC convergence with respect to  $P$  ensures reaching the MC constraining numerical accuracy with low polynomial orders.* Of course, the previous statement tacitly assumes that the  $P$ -truncated reduced model we solve here converges as  $P$  grows: spectral convergence with respect to  $P$  has been proved for the uncertain linear equation, see [30], but the proof for the nonlinear photonic system remains out of the scope of this paper.

For this test-case, truncation (26) has been used. During the run of ISMC-gPC, we monitored whether or not the truncation is activated (i.e. whether or not sG-gPC would have been enough). With the numerical parameters used in this section, the *truncation has not been activated*.

Let us finish by simple performance considerations. For this test case, the average cost of one ISMC run is approximately  $\approx 1845.4s$ . whereas one run of ISMC-gPC <sub>$P=3$</sub>  costs  $\approx 2023.9s$ . As a consequence, in order to produce the same results as in figure 2 with ni-ISMC <sub>$N=15$</sub>  and ISMC-gPC <sub>$P=3$</sub>  with equivalent accuracies, the gain is of about  $\times \frac{N \times 1845.4}{2023.9} = \frac{15 \times 1845.4}{2023.9} = 13.67$  in favor of the intrusive ISMC-gPC solver. We can see that an ISMC-gPC run costs more than an (average) ISMC one: this is mainly due to the additional cost induced by the evaluation of the truncation of  $T$  in order to evaluate the opacities during the tracking together with an overall cost of the tallying phase (just as in [10]). Still, a factor  $\times 13.67$  is gained as only one run of ISMC-gPC is needed instead of  $N = 15$  in this case.

In the next section, we consider the same test-case but with an uncertain heat capacity  $C_v$  instead of an uncertain absorption opacity.

### 5.1.2. Uncertain infinite medium problems: uncertain heat capacity $C_v$

In this section, we slightly change the previous test-case: the opacity is now deterministic, given by  $\sigma_a(E(t, X), X) = \sigma_a = \bar{\sigma}_a$  but the heat capacity is considered uncertain, given by  $C_v(X) = \bar{C}_v + \check{C}_v X$  with  $\bar{C}_v = \frac{3}{2}$ ,  $\check{C}_v = \frac{1}{2}$  and  $X \sim \mathcal{U}([-1, 1])$ .

Figure 4 presents the results obtained with ni-ISMC <sub>$N=15$</sub>  and ISMC-gPC <sub>$P=3$</sub>  on this test-problem. Here, the steady state is not anymore deterministic and the problem ends with non-zero variances for the different quantities of interest  $I, E, T, T_r$ . Figure 4 (top left) presents the time evolution of the mean, the variance and some realisations of the radiation intensity  $t \rightarrow I(t, X)$ . In average, the radiation intensity decreases until it reaches a steady state. On another hand, as time passes, the variance of the radiation intensity grows until reaching a plateau. The figure also displays the realisations  $t \rightarrow I(t, X_i)$  at the  $N = 15$  Gauss-Legendre points to give an idea of the dispersion around the mean. Note that both ni-ISMC and ISMC-gPC are in good agreement. This is also the case for the other outputs of interest,  $E, T, T_r$ , displayed in the other pictures of figure 4. The middle-right picture of figure 4 presents the results on the material temperature  $T$ . It is interesting noticing the singular behaviour of this observable with respect to uncertainty: a peak of uncertainty is observed during the transient regime. Finally, the last line of figure 4 presents some convergence studies with respect to  $P$ . The behaviour is overall the same as in the previous test-case (fast convergence for  $P \in \{0, \dots, 3\}$  then a plateau of level  $\approx \mathcal{O}(\frac{1}{\sqrt{N_{MC}}})$  is reached)

---

<sup>42</sup>Characterised by a degradation of the gPC accuracy as time increases and the need to resort to higher polynomial order  $P$  for later times  $t$  in order to reconstitute the same level of accuracy, see [63, 49].



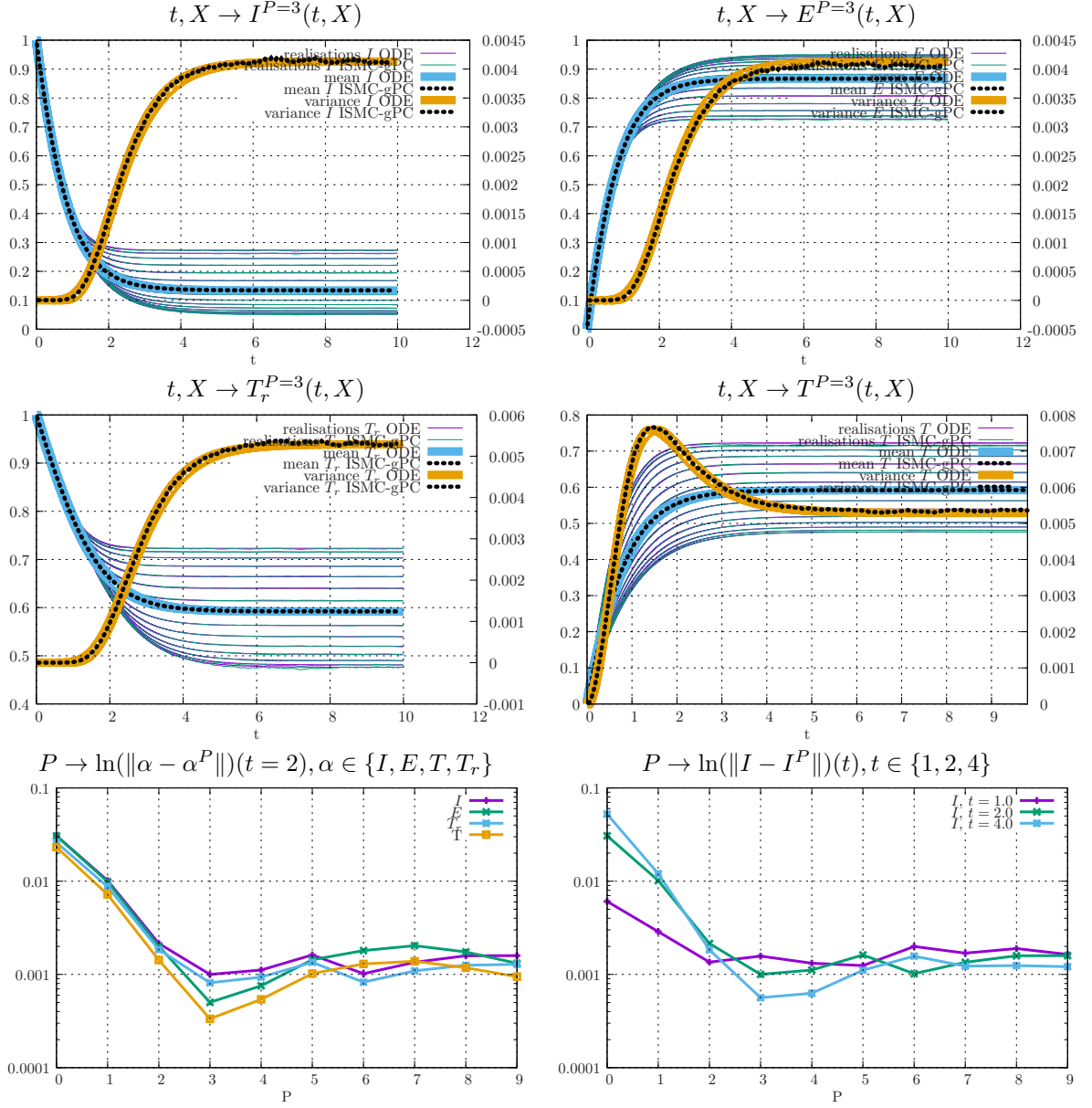


Figure 4: The four top pictures: mean, realisations at the  $N = 15$  Gauss-Legendre points and variance with respect to time of the ni-ISMC approximations and the ISMC-gPC $_{P=3}$  ones. In particular, we present the results in term of intensity of radiation  $I$ , matter energy  $E$  and material and radiative temperatures  $T, T_r$ . The numerical parameters for ISMC-gPC are  $\Delta t = 10^{-2}$ ,  $N_{MC} = 10^6$ ,  $P = 3$ . The two bottom pictures: quantitative convergence study with respect to  $P$  of the ISMC-gPC $_P$  approximations of  $I, E, T_r, T$  (bottom-left) and for  $I$  for several times  $t \in \{1.0, 2.0, 4.0\}$  (bottom-right).

except for the behaviour of the convergence curves with respect to times, see figure 4 (bottom-right). For this benchmark, the error seems to increase for later times: for example for  $P = 1$ ,  $\|I - I^{P=1}\|(t = 1) \leq \|I - I^{P=1}\|(t = 2) \leq \|I - I^{P=1}\|(t = 4)$ . It seems that the long-term behaviour of gPC depends strongly on the configuration of interest (this is in agreement with the observations of [49, 30]). Still, fast convergence rates are observed and the MC constraining error (the plateau) is reached as soon as  $P \geq 3$ .

### 5.1.3. Uncertain infinite medium problems: uncertain absorption opacity $\sigma_a$ and heat capacity $C_v$

The test-case of this section is a combination of the two previous ones: both the opacity and the heat capacity are uncertain. They are given by  $\sigma_a(E(t, X), X) = \sigma_a = \bar{\sigma}_a + \hat{\sigma}_a X_1$  and  $C_v(X) = \bar{C}_v + \hat{C}_v X_2$  with  $\bar{C}_v = \frac{3}{2}$ ,  $\hat{C}_v = \frac{1}{2}$  with  $X_1, X_2 \sim \mathcal{U}([-1, 1])$ . In other words, this problem is 2D ( $Q = 2$ ) with respect to the stochastic dimension.

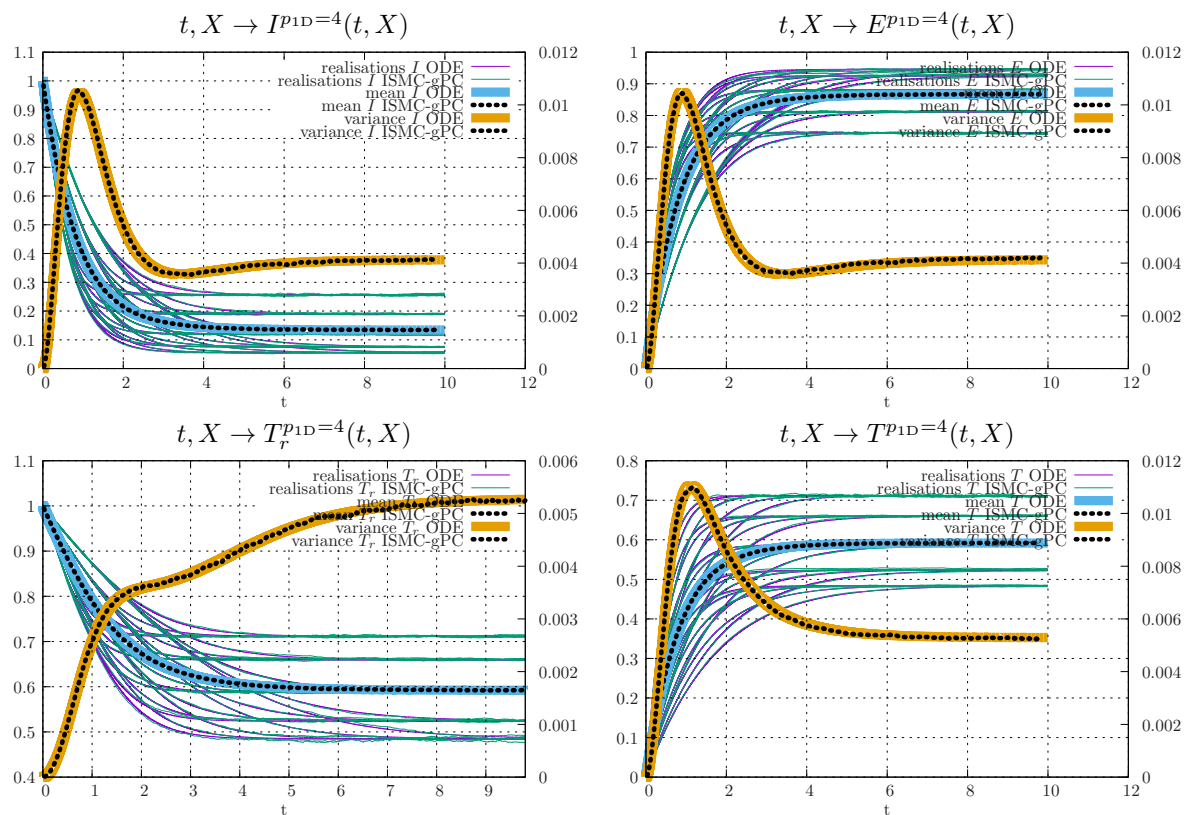


Figure 5: Time evolutions of the mean  $t \rightarrow \mathbb{E}[\alpha^{p_{1D}=4}](t)$ , the variance  $t \rightarrow \mathbb{V}[\alpha^{p_{1D}=4}](t)$  and the 25 realisations  $t \rightarrow \alpha^{p_{1D}=4}(t, X_i)$  (at the Gauss-Legendre points  $i \in \{1, \dots, 25\}$ ) of the radiation intensity, the material energy, the radiation temperature and of the material temperature (i.e. for  $\alpha \in \{I, E, T_r, T\}$ ).

Figures 5–6 present the results obtained with ni-ISM $C_{N_{1D}=5}$  and ISMC-gPC $_{p_{1D}=4}$  on this test-problem. Note that  $N_{1D} = 5$  for ni-ISM $C$  means that  $N = N_{1D}^Q = 5^2 = 25$  tensorised Gauss-Legendre points are used in practice. For ISMC-gPC,  $p_{1D} = 4$  means that  $P = (p_{1D} + 1)^Q = (4 + 1)^2 = 25$  gPC coefficients are evaluated during the intrusive resolution.

Figure 5 presents the time evolutions of the mean  $t \rightarrow \mathbb{E}[\alpha^{p_{1D}=4}](t)$ , the variance  $t \rightarrow \mathbb{V}[\alpha^{p_{1D}=4}](t)$  and the 25 realisations  $t \rightarrow \alpha^{p_{1D}=4}(t, X_i)$  (at the Gauss-Legendre points  $i \in \{1, \dots, 25\}$ ) of the radiation intensity, the material energy, the radiation temperature and of the material temperature (i.e. for  $\alpha \in \{I, E, T_r, T\}$ ) obtained by both ni-ISMC and ISMC-gPC. For this test-case, neither the transient state nor the steady state are anymore deterministic. The variance is non-zero as soon as  $t > 0$  for the different quantities of interest  $I, E, T, T_r$ . If one compares the results (for example the variances) of the two previous sections, one can observe that the variances in the transient regime are close to the ones of figure 1 whereas the variances in the steady state phase are more like the ones of figure 4. Let us analyse more precisely this behaviour thanks to figure 6.

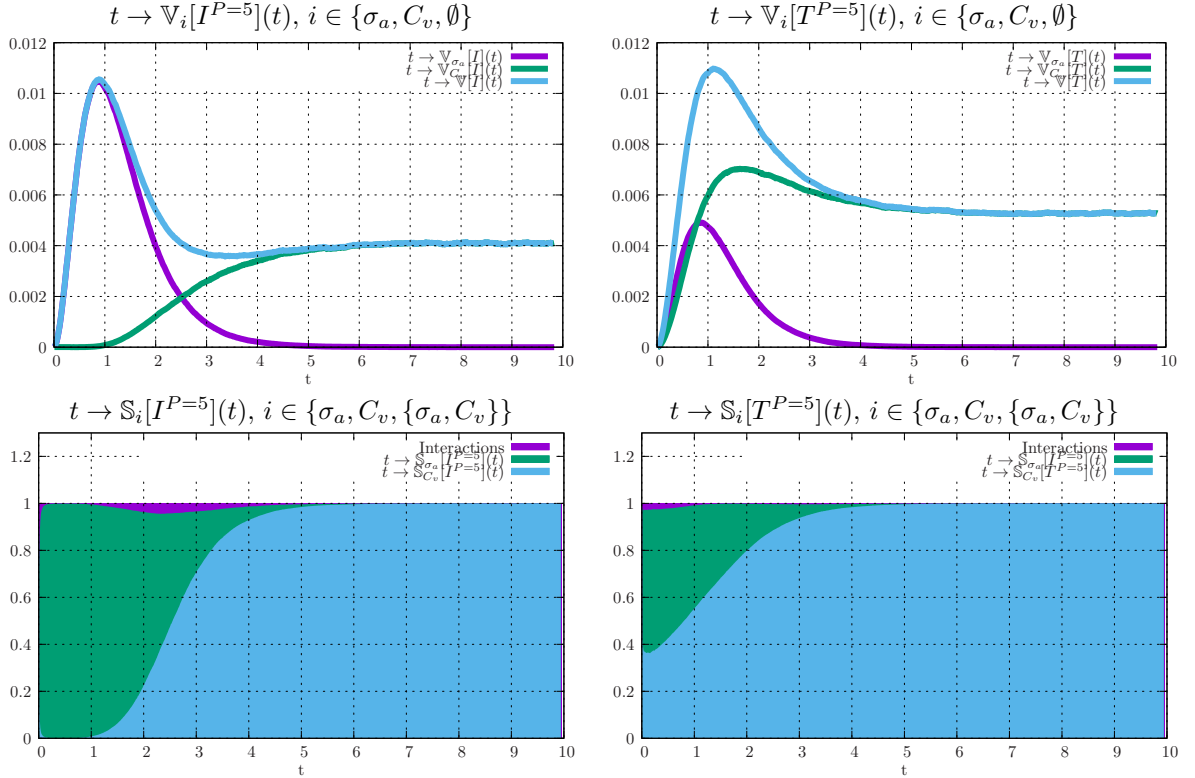


Figure 6: Time evolutions of the partial variances  $t \rightarrow \mathbb{V}_{\sigma_a}[\alpha^{p_{1D}=4}](t)$ ,  $t \rightarrow \mathbb{V}_{C_v}[\alpha^{p_{1D}=4}](t)$  and the total variance  $t \rightarrow \mathbb{V}[\alpha^{p_{1D}=4}](t)$  and the Sobol indices  $t \rightarrow \mathbb{S}_{\sigma_a}[\alpha^{p_{1D}=4}](t)$  and  $t \rightarrow \mathbb{S}_{C_v}[\alpha^{p_{1D}=4}](t)$  of the radiation intensity and of the material temperature (i.e. for  $\alpha \in \{I, T\}$ ).

Figure 6 (top) presents the time evolutions of the partial variances, cf. [23, 64, 10],  $t \rightarrow \mathbb{V}_{\sigma_a}[\alpha^{p_{1D}=4}](t)$  and  $t \rightarrow \mathbb{V}_{C_v}[\alpha^{p_{1D}=4}](t)$  together with the one of the total variance  $t \rightarrow \mathbb{V}[\alpha^{p_{1D}=4}](t)$  for  $\alpha \in \{I, T\}$ . The time evolutions of the partial variance with respect to  $\sigma_a$  is the same as in section 5.1.1 whereas the time evolutions of the partial variance with respect to  $C_v$  is the same as in section 5.1.2. They sum-up to the total variance.

Figure 6 (top) displays the Sobol indices  $t \rightarrow \mathbb{S}_{\sigma_a}[\alpha^{p_{1D}=4}](t)$  and  $t \rightarrow \mathbb{S}_{C_v}[\alpha^{p_{1D}=4}](t)$  of the radiation intensity and of the material temperature (i.e. for  $\alpha \in \{I, T\}$ ). The Sobol indices, cf. [23, 64, 10], express the percentage of the total variance explained by the two uncertain parameters

$\sigma_a, C_v$ . For this test-case, the two uncertain parameters only slightly interact during a very narrow window of time  $t \in [1, 4]$ . Otherwise, before  $t \in [0, 1]$ , the variance is only explained by the uncertainty on  $\sigma_a$  whereas it is only explained by  $C_v$  for  $t \in [4, 10]$ . Such statistical tools are very powerful: thanks to them, for example, we can help design experiments to calibrate  $\sigma_a$  and  $C_v$ . Thanks to them, we know  $\sigma_a$  (early times) and  $C_v$  (later times) can be identified with the same configuration and without too much interferences.

The previous statistical tools, i.e. the partial variances and the Sobol indices, give precious information concerning the relative behaviour of the uncertain parameters on the outputs of interest. Their interpretation is easy and straightforward. Sobol indices are amongst the most efficient and reliable sensitivity indices but also amongst the most costly, see [23]. Now, with ISMC-gPC, we can have access to accurate evaluations of these indicators with only one run of a code. For this study, one ISMC-gPC run costs  $\approx 1050s$ . The pick-and-freeze strategy needed to approximate the Sobol indices (see [65]) in a non-intrusive manner implies  $(Q + 2) \times N$  runs of a code with  $N$  ranging from 100 to 1000, see [23]. With respect to such integration scheme, the gain is important. But the comparison would not be fair: in our case, accurate results on the same statistical observables can be obtain by applying non-intrusive gPC (ni-gPC) with  $p_{1D} = 4$ : it consists in using the Gauss-Legendre experimental design to estimate gPC coefficients and deduce the partial variances and Sobol indices from them (see [66]). In this case, suppose one needs  $p_{1D}$  as a polynomial order in each direction, then ni-gPC needs  $N^Q = (p_{1D} + 1)^Q$  non-intrusive runs. Let us take  $p_{1D} = 4$ , just as for ISMC-gPC, then ni-gPC needs  $(4 + 1)^2 = 25$  independent runs of a black-box ISMC code. For this study, the cost of the ni-ISMIC runs ranges from  $\approx 360s$ . to  $\approx 401s$ . ISMC-gPC hence ensure a gain ranging from  $\times \frac{25 \times 360}{1050} \approx 8.57$  to  $\times \frac{25 \times 401}{1050} \approx 9.54$  for this study. Of course if one has access to 25 computational units than the restitution time for ni-ISMIC becomes  $401s$ , the maximum time. But with several processors, ISMC-gPC can also be accelerated, thanks, for example, to domain replication (this will be emphasized in the next sections).

In the next section, we consider an uncertain spatial benchmark based on the Heaviside problem [14, 9, 58]. This benchmark is especially used to stress the capabilities of the MC solver with respect to what is commonly called the *teleportation error* [28, 29, 25, 59, 27, 26] in the equilibrium diffusion limit.

## 5.2. Uncertain Heaviside

Let us here consider a new configuration. Every details (initial conditions, numerical parameter choices and test-case justifications) of the test-problem are presented in Appendix A. The initial condition is a Heaviside of (relaxed  $T = T_r$ ) temperatures in the center of the spatial domain, see figure A.17. The conditions are exactly the same as in [9, 58] except that the absorption opacity is uncertain with  $\sigma_a(X) = \bar{\sigma}_a + \hat{\sigma}_a X$  and  $X \sim \mathcal{U}([-1, 1])$ . In practice, we take  $\bar{\sigma}_a = 1800$  and  $\hat{\sigma}_a = \bar{\sigma}_a \times 75\%$ . For  $X = 0$ , the benchmark is the same as in [9, 58]. In other words, for  $X = 0$ , the equilibrium diffusion limit (38) is valid. Otherwise, for this benchmark, equilibrium must be fulfilled  $\forall X \sim d\mathcal{P}_X$  but not necessarily the diffusion limit. Note that in all the results of this section, *truncation* (26) is considered.

Figure 7 displays the curve  $x \rightarrow T(x, t^* = 10^{-8}, X = 0)$  obtained from

- a deterministic reference code solving directly equation (38) for the equilibrium diffusion limit ( $N_x = 2000$  cells),

- a non-intrusive application of ISMC (ni-ISM) for  $X = 0$  with  $\Delta t = 10^{-12}$ ,  $N_{MC} = 3.9 \times 10^7$  and  $N_x = 40$  cells,
- and the ISMC-gPC $_{P=5}$  solver of this paper taken at  $X = 0$  in the same numerical conditions as ni-ISM.

First, ISMC is able to recover the equilibrium diffusion limit on a coarse mesh (this was already put forward in [9]). Now, ISMC-gPC $_{P=5}$  at  $X = 0$  allows recovering the results obtained with ni-ISM

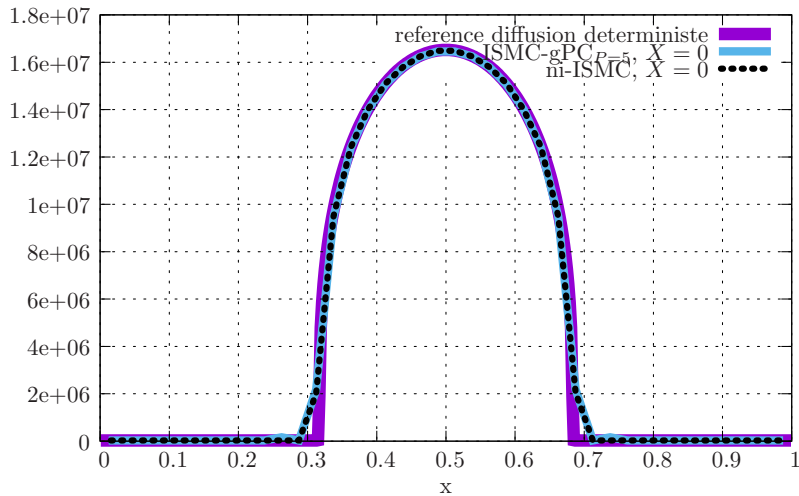


Figure 7: Comparison of the curve  $x \rightarrow T(x, t^* = 10^{-8}, X = 0)$  obtained from a deterministic reference code solving directly equation (38) for the equilibrium diffusion limit, a non-intrusive application of ISMC for  $X = 0$ , and the ISMC-gPC $_{P=5}$  solver of this paper taken at  $X = 0$ .

at  $X = 0$  with a very good agreement, in the same numerical conditions. With this figure, we put forward the fact that ISMC-gPC inherits the fast convergence rate of ISMC with respect to spatial discretisation (see remark 3.3) in the equilibrium diffusion regime: *ISM-gPC (as ISMC) is not sensitive to the teleportation error*. Note also that, this will be discussed more in details later on when commenting table 1, the computational time ( $\approx 8h40min$ ) for ni-ISM with  $X = 0$  is only slightly smaller than the one of ISMC-gPC ( $\approx 10h$ , see table 1): the restitution times and accuracies are comparable together with ISMC-gPC allowing much richer capabilities in term of uncertainty analysis.

Figure 8 compares the results obtained with ni-ISM $_{N=15}$  to the ones obtained with ISMC-gPC $_{P=5}$  in term of mean (left axis) and variance (right axis) spatial profiles  $x \rightarrow \mathbb{E}[\alpha](x, t^* = 10^{-8})$ ,  $x \rightarrow \mathbb{V}[\alpha](x, t^* = 10^{-8})$  for  $\alpha \in \{I, E, T_r, T\}$ . In term of mean spatial profiles, the solution at the final time exhibits steep propagation fronts. Less steep than in the deterministic case (see figure 7 for example) as the averaging process tends to smooth out the spatial profiles but the gradients remain sharp. The spatial profiles of the variance tend to show that the uncertainty is mainly located in the center of the domain  $x = 0.5$  and in the vicinities of the propagation fronts ( $x \approx 0.3$  and  $x \approx 0.7$  for this time of interest). The uncertainty does not affect all quantities in the same manner: if for  $E, T_r, T$  the variance is mainly important for the propagation front, this is not necessarily the case for  $I$ : the fluctuations on  $I$  due to  $X$  are more important in the vicinity

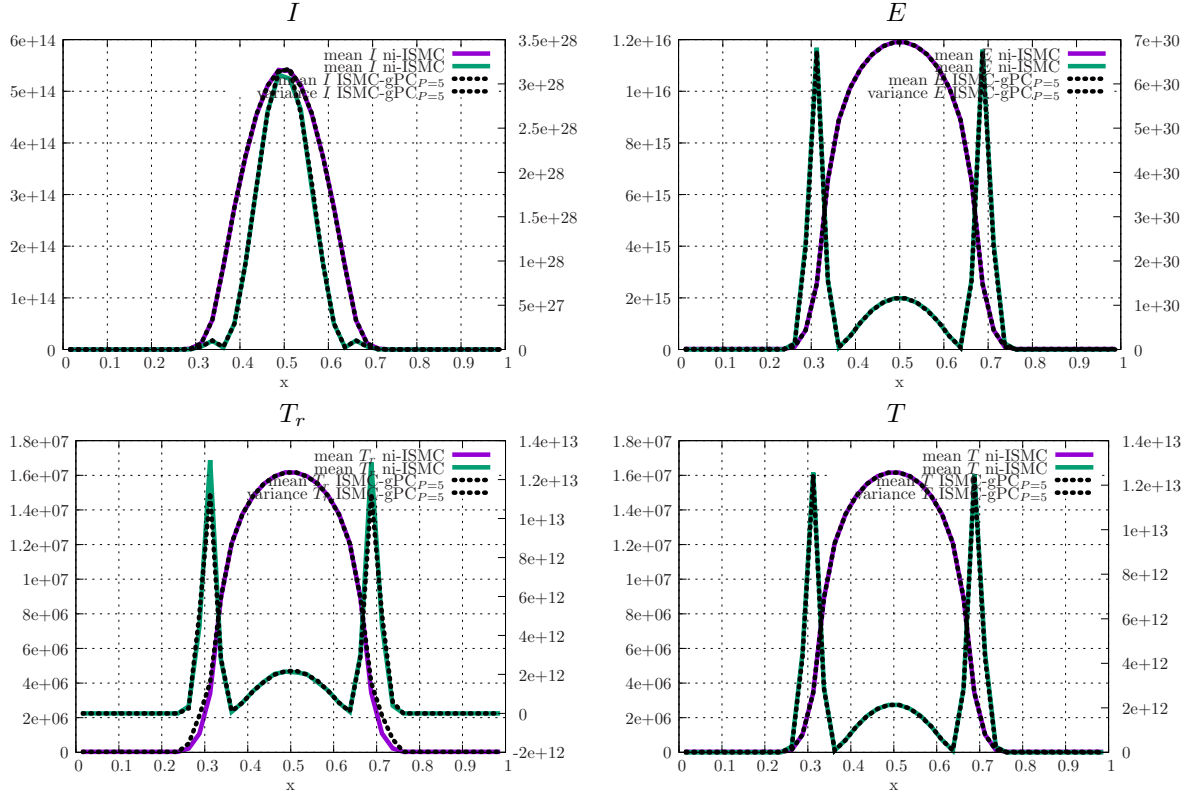


Figure 8: Comparison of the results obtained with ni-ISMC $_{N=15}$  to the ones obtained with ISMC-gPC $_{P=5}$  in term of mean (left axis) and variance (right axis) spatial profiles  $x \rightarrow \mathbb{E}[\alpha](x, t^* = 10^{-8}), x \rightarrow \mathbb{V}[\alpha](x, t^* = 10^{-8})$  for  $\alpha \in \{I, E, T_r, T\}$ .

$x = 0.5$ . Now, in term of mean and variance spatial profiles, ni-ISM $C_{N=15}$  and ISMC-gPC $_{P=5}$  present very good agreements for every observables of interest  $I, E, T_r, T$ .

Figure 9 compares the 15 realisations  $x \rightarrow \alpha(x, t^* = 10^{-8}, X_i), \alpha \in \{I, E, T_r, T\}$  obtained with ni-ISM $C$  at the  $N = 15$  Gauss-Legendre points  $(X_i)_{i \in \{1, \dots, 15\}}$  to the reconstructed, *via* truncation (26), ones from ISMC-gPC $_{P=5}$  at the same points. The averaged spatial profiles  $x \rightarrow \mathbb{E}[\alpha](x, t^*), \alpha \in \{I, E, T_r, T\}$  are also displayed. Figure 9 attests, for this test-case at least, that very accurate

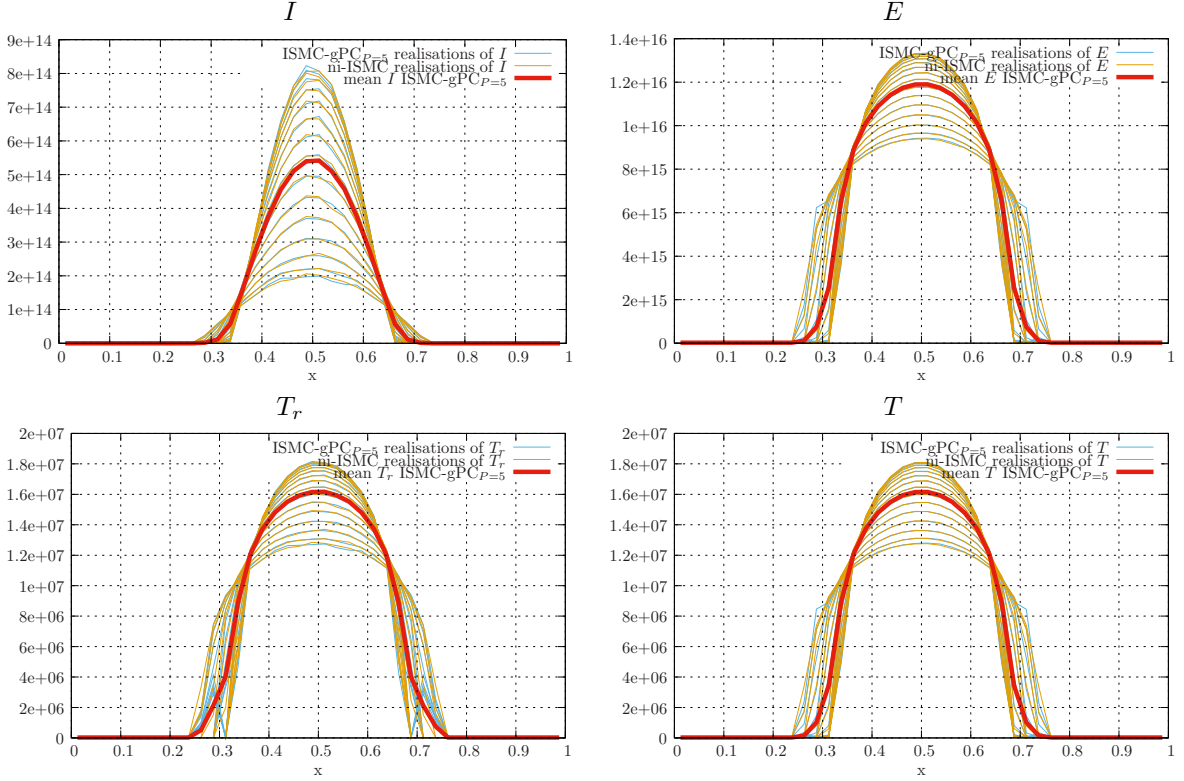


Figure 9: Comparison of the 15 realisations  $x \rightarrow \alpha(x, t^* = 10^{-8}, X_i), \alpha \in \{I, E, T_r, T\}$  obtained with ni-ISM $C$  at the  $N = 15$  Gauss-Legendre points  $(X_i)_{i \in \{1, \dots, 15\}}$  to the reconstructed (*via* truncation (26)) ones from ISMC-gPC $_{P=5}$  at the same points. The means  $x \rightarrow \mathbb{E}[\alpha](x, t^*), \alpha \in \{I, E, T_r, T\}$  are also displayed.

pointwise approximations can be reconstructed from the solution of our reduced model.

With figures 8–9, we can also see that similar accuracies ( $\approx 1\%$  difference on the profiles of the realisations of  $I$  and  $E$ ) are recovered with ni-ISM $C$  and ISMC-gPC on several statistical (mean, variance, realisations) and physical ( $I, E, T_r, T$ ) observables of interest. As a consequence, fair performance comparisons can be made on this benchmark: table 1 presents the sequential runtimes for the ni-ISM $C_{N=15}$  and ISMC-gPC $_{P=5}$ . First, of course, with ISMC-gPC, only one run of the code is necessary. Second, the more important the opacity, the more collisional/diffusive the medium and the longer the computations: the sequential runtimes range from about 3 hours to more than 15 hours. The reader familiar with MC methods will recognise the classical behaviour of MC codes in the diffusion regime (characterised by an important number of collisions per time steps). The

$X$	$\sigma_a(X)$	ni-ISM $C_{N=15}$
$X_1 = -0.98799252$	466.21	03h03min18s
$X_2 = -0.93727339$	534.68	03h15min14s
$X_3 = -0.84820658$	654.92	03h40min52s
$X_4 = -0.72441773$	822.03	04h21min30s
$X_5 = -0.57097217$	1029.1	05h10min53s
$X_6 = -0.39415135$	1267.8	06h10min25s
$X_7 = -0.20119409$	1528.3	07h19min35s
$X_8 = 0.00000000$	1800.0	08h40min29s
$X_9 = +0.20119409$	2071.6	10h07min27s
$X_{10} = +0.39415135$	2332.1	11h26min47s
$X_{11} = +0.57097217$	2570.8	12h39min56s
$X_{12} = +0.72441773$	2777.9	13h56min31s
$X_{13} = +0.84820658$	2945.0	15h32min23s
$X_{14} = +0.93727339$	3065.3	15h26min49s
$X_{15} = +0.98799252$	3133.7	15h42min09s
Total time for the ni-ISM $C$ study		5days 16h34min11s
Average time for the ni-ISM $C$ study		09h06min45s

ISM $C$ -gPC $P=5$
09h54min32s

Table 1: Sequential runtimes for ni-ISM $C_{N=15}$  and for ISM $C$ -gPC. The  $N = 15$  Gauss-Legendre points used for the ni-ISM $C$  and the absorption opacities at those points are displayed.

ISM $C$ -gPC run is only a little more costly in term of runtime than the average computational time of ni-ISM $C$ . In term of sequential restitution time, the gain between ni-ISM $C$  and ISM $C$ -gPC is of  $\times 13.68$  (from about 5 days and 10 hours to only 10 hours). Of course, if one has access to many more computational units (which is common ground), the  $N = 15$  runs may be run at the same time: in this case, the gain is only of  $\approx 1.58$  (maximum runtime over the ISM $C$ -gPC one). But in this case, the ISM $C$ -gPC restitution time can also easily be accelerated thanks to, for example, replication domain [35]: when using 15 replicated domains to perform the same computations (i.e. there are  $2.6 \times 10^6$  MC particles per replicated domains/processors) the ISM $C$ -gPC is accelerated of a factor  $\times 13$ . In other words, when comparing ni-ISM $C_{N=15}$  and ISM $C$ -gPC on 15 computational units, the gain is of about  $\times \frac{56529s.}{2744s.} = 20.6$ .

The benchmark of this section presents some promising results: the ISM $C$ -gPC solver seems to inherit important properties of both the ISM $C$  solver [9] (stable affordable time steps, capture of the equilibrium diffusion limit on coarse meshes) and of the gPC based reduced model [10] (fast convergence with respect to  $P$  in the uncertain space, computational gain with respect to ni-ISM $C$ ).

Finally, for this benchmark, we only used truncation (26): the truncation is activated in practice, mainly in the vicinities of the steep propagation fronts. The effect of the truncation activation is mainly visible on the radiation temperature spatial profiles of figure 9 (bottom-left). In the next sections, other choice of truncation will be made (see section 5.4.1). In particular, we present the behaviour of the reduced model if a bad minoration/majoration/truncation choice is made.



### 5.3. Uncertain test-cases with some bad truncation choices

In this section, we briefly revisit the test-cases of sections 5.1.1–5.2 but instead of relying on truncation (26) with maximum-principle-based lower and upper bounds, we are going to consider some bad choices of the extremal bounds. Of course, in practice, the bounds are going to be totally irrelevant. We here mainly want to highlight the fact that even if bad choices are made, the wellposedness of the built reduced model is not questioned. But its physical relevance, on another hand, may be.

#### 5.3.1. Uncertain relaxation case with a bad truncation choice

Figure 10 (top left) presents the results obtained for  $t, X \rightarrow T(t, X)$  with ni-ISM $C_{N=15}$  and ISMC-gPC $_{P=5}$  i.e. in the same conditions as in figure 1 but with a bad truncation: for this test-case, the initial conditions are within

$$I_m = 10^{-3} \leq I^0(x, \omega, X) \leq 1 = I_M, E_m = 10^{-3} \leq E^0(x, X) \leq 1 = E_M, \forall x \in \Omega, \omega \in \mathbb{S}^2, X \sim d\mathcal{P}_X.$$

Of course, due to the maximum principle holding for (5), the solutions must satisfy the above inequalities  $\forall t \in [0, t^*]$ . The above values were the one used in truncation (26) for the computations of section 5.1.1. Here, instead of choosing  $E_m = 10^{-3} = I_m$  and  $E_M = I_M = 1$ , we suggest, on purpose, wrongly choosing half the majorant for  $E_M^{bad} = \frac{1}{2}E_M, I_M^{bad} = \frac{1}{2}I_M$  in truncation (26). In such conditions, see figure 10, the ISMC-gPC $_{P=5}$  reduced model does not fail to reconstitute some results. The reduced model, using a bounded truncation, is wellposed. But it does fail to reconstitute relevant physical solutions: the mean solution is truncated as soon as  $T$  reaches  $T_M^{bad} = \frac{E_M^{bad}}{C_v} = \frac{1}{2}$  whereas it should reach a plateau at  $T \approx 0.57$  as testifies the reference (non-intrusive) solution. It induces some errors also on the variance with respect to time  $t \rightarrow \mathbb{V}[T](t)$  and on the reconstructed realisations at the Gauss-points  $t \rightarrow T(t, X_i)$  for  $i \in \{1, \dots, 15\}$  as can be seen in figure 10 (top-left).

#### 5.3.2. Heaviside case with a bad truncation choice

In this section, we reconsider the test-case of section 5.2 with two bad truncations: the first truncation we consider is the sG-gPC one. The second one is similar to the one used in the previous section.

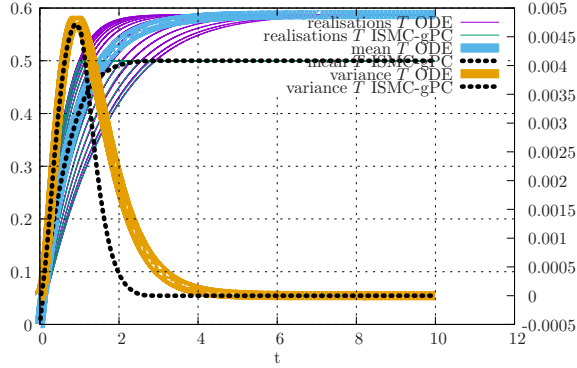
Recovering the sG-gPC reduced model from truncation (26) consists in choosing  $E_m = -\infty = I_m$  and  $E_M = \infty = I_M$ . In this case, nothing prevents the truncation  $T^P$  of the material temperature from going below zero. *If we do not enforce positiveness of  $T^P$ , i.e. when applying sG-gPC, the code crashes on this Heaviside case.* This testifies of the relevance of the discussion of section 2.

Now, let us, instead of choosing  $E_m = 9.2 \times 10^{11}, I_m = 2.798 \times 10^5, E_M = 9.2 \times 10^{14}, I_M = 2.798 \times 10^{16}$  as in section 5.2, choose to take  $E_M^{bad} = \frac{1}{2}E_M, I_M^{bad} = \frac{1}{2}I_M$  in truncation (26). The results are displayed on the top-right and bottom pictures of figure 10. Due to the poor choice of bounds in the truncation, the solution of the reduced model, even if relatively converged with respect to the number  $N_{MC}$  of MC particle (as the results are far from being noisy) and with respect to  $P$ , does not capture the physical solution.

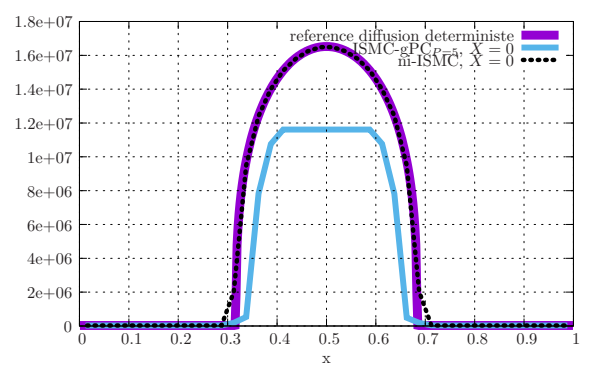
To sum-up, with the benchmarks of this section, we put forward the fact that

- sG-gPC is not enough to produce wellposed reduced models (leading to robustness issues, just as for hyperbolic systems of conservation laws, see [8]),
- choosing a bad truncation (bad bounds) can lead to wellposed reduced model (no crash of the code nor numerical instabilities) but with poor physical relevance.

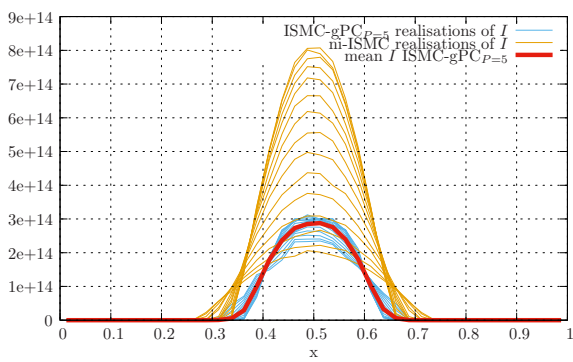
Relaxation (section 5.1.1) with a bad truncation



Heaviside case  $T$  (section 5.2) with a bad truncation



Heaviside case  $I$  (section 5.2) with a bad truncation



Heaviside case  $T_r$  (section 5.2) with a bad truncation

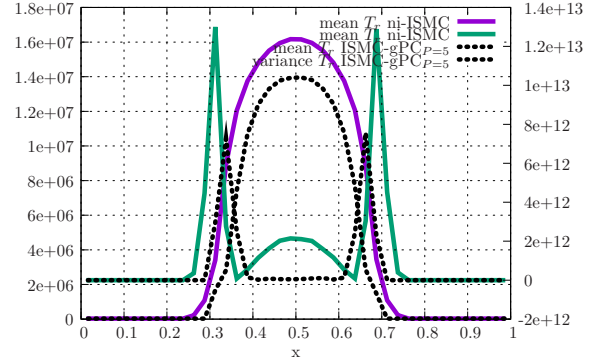


Figure 10: The consequences of using a bad truncation on the relaxation problem (see section 5.1.1) on the top left picture and on the Heaviside problem (see section 5.2) on the other picture: using a bad truncation may lead to a wellposed problem but its solution may be physically irrelevant.

From now, *a priori* relevant (maximum principle based) bounds will be chosen in the benchmarks. In the next section, we consider one last test-case which allows putting forward one additional important aspect of the choice of the truncation.

#### 5.4. Uncertain Marshak wave

In this section, we consider an uncertain Marshak wave built from the benchmark of [27]. The test-case corresponds to the study of a 1D Marshak wave [11] with dimensionless units. A black body heats the left boundary of the domain  $x \in \mathcal{D} = [0, 4]$  with temperature  $T(x = 0) = 1$ . The radiation constant is  $a = 1$  and so is the speed of light  $c = 1$ . There is no scattering (i.e.  $\sigma_s = 0$ ) and  $\sigma_t(T_m) = \sigma_a(T_m) = \frac{\sigma_a}{T_m^3} = \frac{10}{T_m^3}$ . Note that this benchmark will demonstrate that the ISMC-gPC solver can be used with temperature dependent opacities. Besides, the test-problem considers a perfect gas eos with  $\rho = 1$  and  $C_v = 7.14$ . The medium is initially cold as  $T(x, t = 0) = T_0(x) = 10^{-2} \forall x \in \mathcal{D} = [0, 4]$ . We are here interested in the (material and radiative) temperature profiles at  $t^* = 500$ . The previously described configuration is made uncertain by considering an uncertain absorption opacity  $\sigma_a(T_m, X) = \frac{\bar{\sigma}_a + \hat{\sigma}_a X}{T_m^3}$  with  $X \sim \mathcal{U}([-1, 1])$ . In practice, we take  $\hat{\sigma}_a = 60\% \bar{\sigma}_a = 60\% \times 10 = 6$ .

In the next three subsections, three different admissible truncations (using relevant bounds) are going to be considered on this same benchmark.

##### 5.4.1. Uncertain Marshak wave with uncertain $\sigma_a$ with truncation (26)

In this section, we consider the test-case described above and compare the results obtained by ni-ISMC $_{N=15}$  and with ISMC-gPC $_{P=5}$  with truncation (26).

Figure 11 compares the results obtained with ni-ISMC $_{N=15}$  to the ones obtained with ISMC-gPC $_{P=5}$  in term of mean (left axis) and variance (right axis) spatial profiles  $x \rightarrow \mathbb{E}[\alpha](x, t^* = 10^{-8})$ ,  $x \rightarrow \mathbb{V}[\alpha](x, t^* = 10^{-8})$  for  $\alpha \in \{I, E, T_r, T\}$  on the Marshak wave test-problem. For  $I, E, T$ , the means and variances of ni-ISMC and ISMC-gPC are in good agreement. But for  $T_r$ , a singular behaviour is observed. The mean of  $T_r$  is overestimated whereas its variance is underestimated: the differences can not be explained by numerical noise and increasing the polynomial order  $P$  improves only slightly the results. It is easier understanding what happens by analysing figure 12. Figure 12 compares the means and  $N = 15$  Gauss-Legendre realisations obtained by both ni-ISMC $_{N=15}$  and ISMC-gPC $_{P=5}$  for the spatial profiles of  $I, E, T, T_r$ . If we focus on the  $N = 15$  realisations, we can see that for  $I, E, T$ , even if some spurious oscillations seem to appear for ISMC-gPC (and not for ni-ISMC), the reconstructed realisations of the spatial profiles of these quantities are quite well captured by the  $P$ -truncated reduced model. Note that those spurious oscillations are closely related to the fact that no pointwise (in the  $X$ -space) maximum principle is ensured by the reduced model, see remark 2.4. Now, for the spatial profile of  $T_r$ , the reconstructed realisations have a much more oscillating behaviour in the vicinity  $x \in [1.5, 2.5]$ . In this vicinity, truncation (26) is often activated in order to obtain a positive quantity under the exponent in  $T_r(I) = (\frac{I}{a})^{\frac{1}{4}}$ : in fact, truncation (26) can not, natively, preserve the mean (or the higher order moments) of  $T_r$ . The main risk with using such truncation here is to think that, on average,  $T$  and  $T_r$  have not yet relaxed to the same value whereas, in practice, this is wrong (as testify the reference ni-ISMC results). This is all the more unfortunate that truncation (26) is computationally fast to apply and gives good results in other situations (see for example the benchmark of section 5.2). Let us explore, in the next subsections, other truncations for the same benchmark.

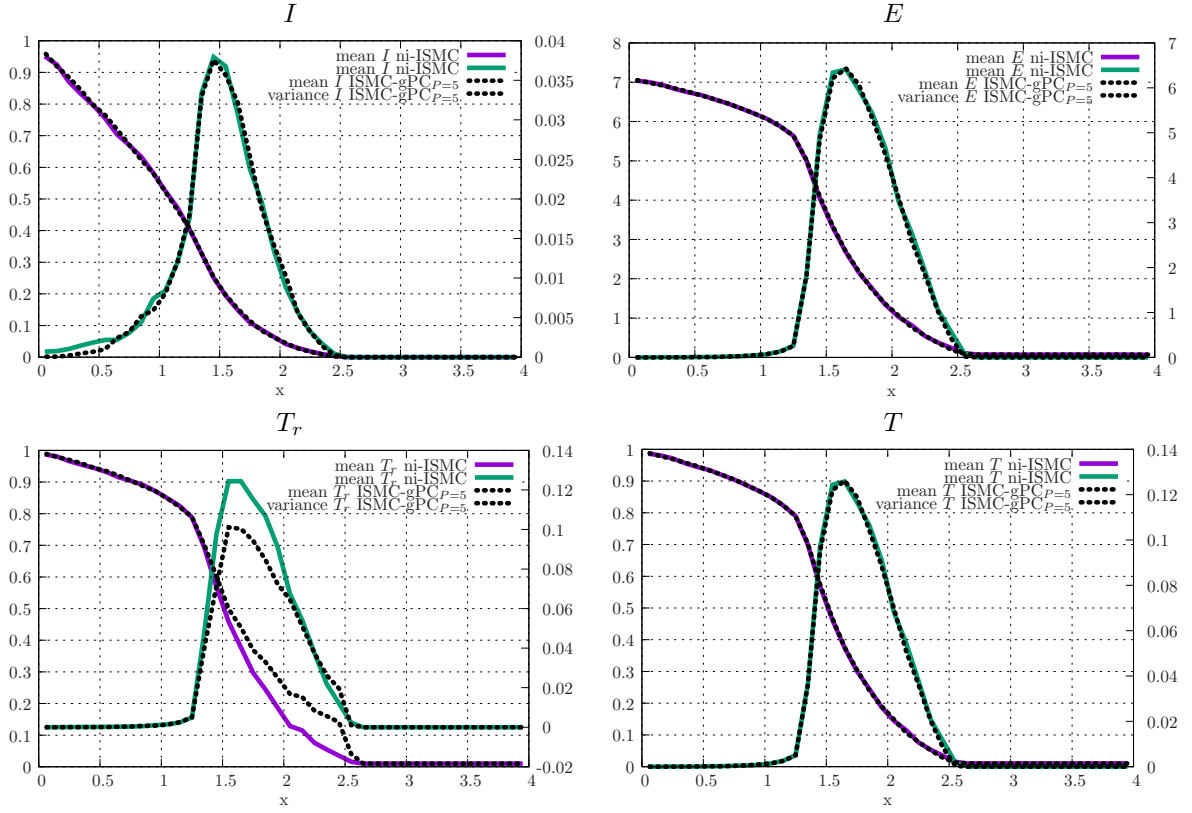


Figure 11: Comparison of the results obtained with ni-ISMCM $_{N=15}$  to the ones obtained with ISMCM-gPC $_{P=5}$  in term of mean (left axis) and variance (right axis) spatial profiles  $x \rightarrow \mathbb{E}[\alpha](x, t^* = 10^{-8}), x \rightarrow \mathbb{V}[\alpha](x, t^* = 10^{-8})$  for  $\alpha \in \{I, E, T_r, T\}$  on the Marshak wave test-problem.

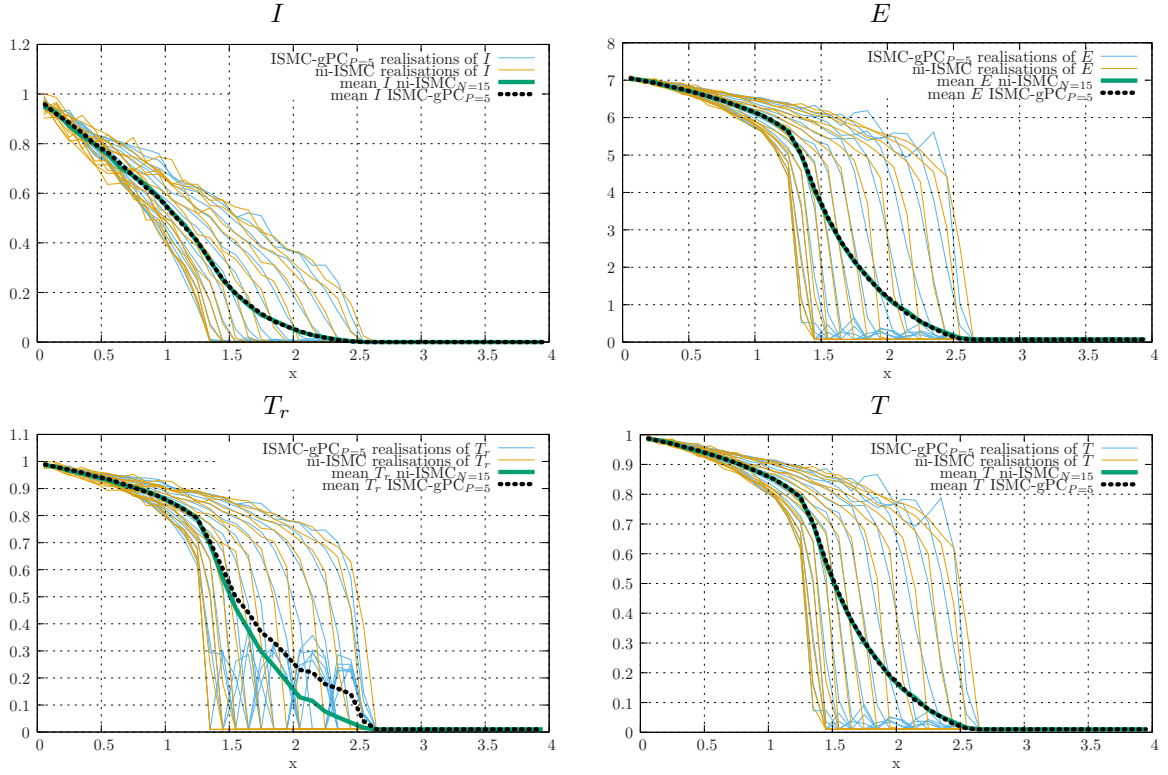


Figure 12: Comparison of the 15 realisations  $x \rightarrow \alpha(x, t^* = 10^{-8}, X_i), \alpha \in \{I, E, T_r, T\}$  obtained with ni-ISMC at the  $N = 15$  Gauss-Legendre points  $(X_i)_{i \in \{1, \dots, 15\}}$  to the reconstructed (via truncation (25)) ones from ISMC-gPC $_{P=5}$  at the same points. The means  $x \rightarrow \mathbb{E}[\alpha](x, t^*), \alpha \in \{I, E, T_r, T\}$  are also displayed.

#### 5.4.2. Uncertain Marshak wave with uncertain $\sigma_a$ with truncation (27)

In this section, we revisit the previous benchmark but with the  $\theta$ -truncation (27) instead of truncation (26). Truncation (27) is built to preserve the first moment (i.e. the mean) of  $I$  and

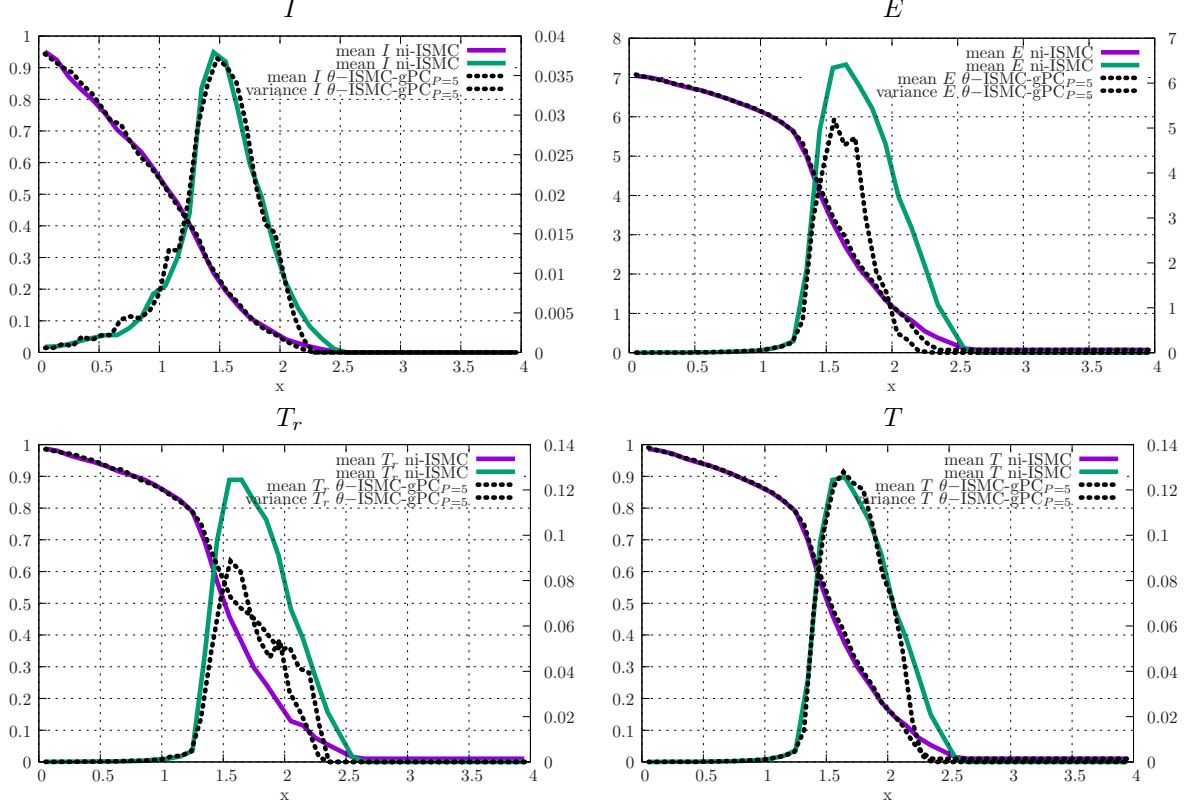


Figure 13: Comparison of the results obtained with ni-ISMCM $_{N=15}$  to the ones obtained with ISMCM-gPC $_{P=5}$  in term of mean (left axis) and variance (right axis) spatial profiles  $x \rightarrow \mathbb{E}[\alpha](x, t^* = 10^{-8}), x \rightarrow \mathbb{V}[\alpha](x, t^* = 10^{-8})$  for  $\alpha \in \{I, E, T_r, T\}$  on the Marshak wave test-problem. Truncation (27) is used

$E$ , as we have  $\int I^P d\mathcal{P}_X = \lambda_0^I$  and  $\int E^P d\mathcal{P}_X = \lambda_0^E$ . But the  $\theta$ -truncation does not necessarily preserve their higher order moments. For example, the variances are given by  $\int (I^P)^2 d\mathcal{P}_X = (\theta^I)^2 \sum_{k=1}^P (\lambda_k^I)^2$  and  $\int (E^P)^2 d\mathcal{P}_X = (\theta^E)^2 \sum_{k=1}^P (\lambda_k^E)^2$ : they are preserved only if  $\theta^\alpha = 1, \forall \alpha \in \{I, E\}$ . This will typically not be true if the truncation is activated.

In practice, the  $\theta$ -limitation on  $I$  and  $E$  is built at the beginning of every time step (i.e.  $\theta^\alpha, \alpha \in \{I, E\}$  is cell and time step dependent). The additional computational is negligible in comparison to the MC resolution (the media is very collisional). The results in term of mean and variance spatial profiles are displayed in figure 13 and in term of spatial realisations in figure 14. Astonishingly, the results are worse than with truncation (26): this is because the  $\theta$ -limitation, when activated in a cell during a time step, operates on every realisations  $X_p$  of each MC particle  $p$  crossing the cell during the time step. On another hand, truncation (26) of section 5.4.1 was only activated for MC

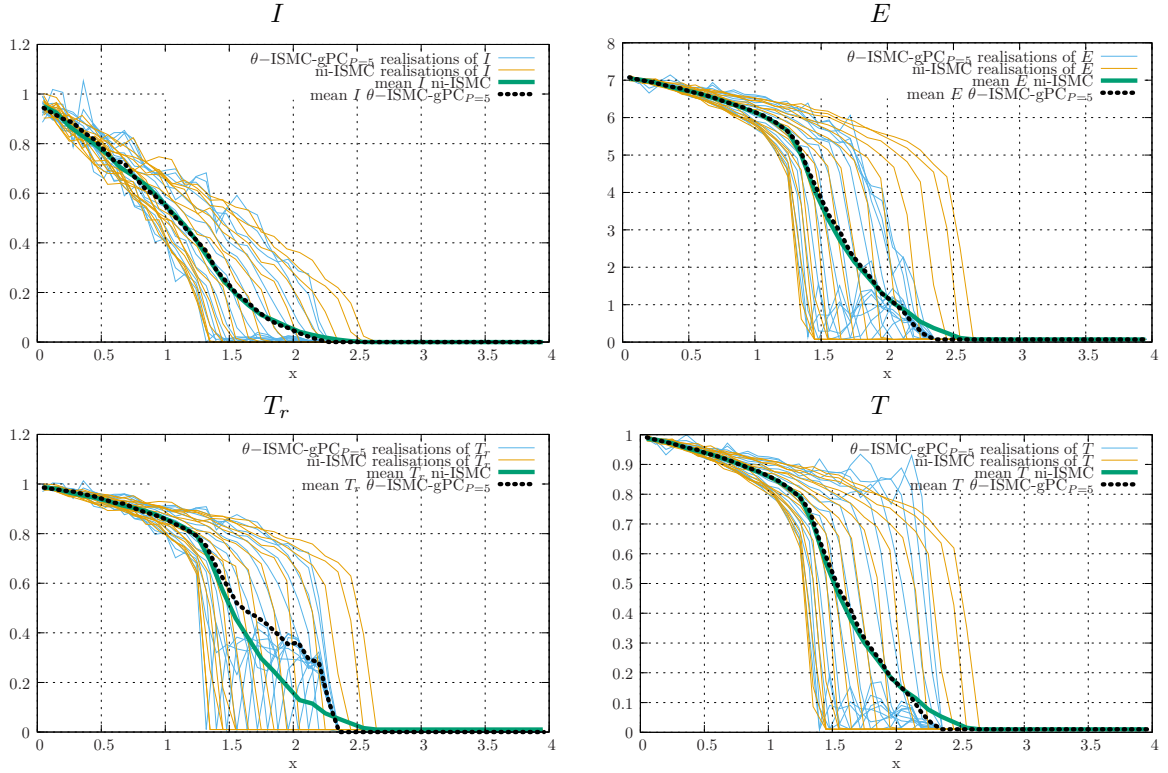


Figure 14: Comparison of the 15 realisations  $x \rightarrow \alpha(x, t^* = 10^{-8}, X_i), \alpha \in \{I, E, T_r, T\}$  obtained with ni-ISM at the  $N = 15$  Gauss-Legendre points  $(X_i)_{i \in \{1, \dots, 15\}}$  to the reconstructed (via truncation (27)) ones from ISMC-gPC $_{P=5}$  at the same points. The means  $x \rightarrow \mathbb{E}[\alpha](x, t^*), \alpha \in \{I, E, T_r, T\}$  are also displayed.

particles  $p$  of field  $X_p$  such that  $I^P, E^P, T^P$  did not respect<sup>43</sup> (20).

In the next paragraph, we study the effect of truncation (29) which both preserves accretiveness and the polynomial moments of the different quantities.

#### 5.4.3. Uncertain Marshak wave with uncertain $\sigma_a$ with truncation (29)

Finally, in this section, we revisit the same benchmark as in the two previous sections but using, in an original way which will be detailed later on, truncation (29). Let us first comment on the results of figures 15 and 16 before giving few implementation details. On figure 15, we can see that

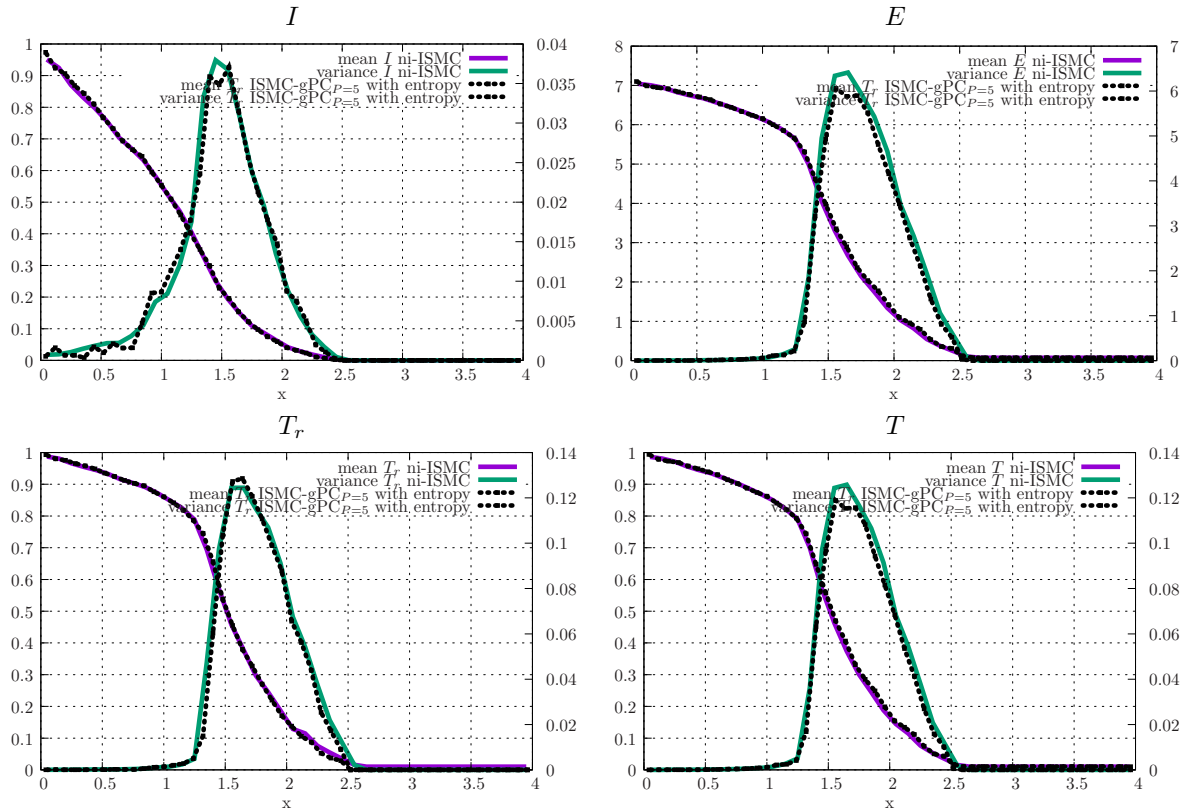


Figure 15: Comparison of the results obtained with ni-ISMCM $_{N=15}$  to the ones obtained with ISMCM-gPC $_{P=5}$  in term of mean (left axis) and variance (right axis) spatial profiles  $x \rightarrow \mathbb{E}[\alpha](x, t^* = 10^{-8}), x \rightarrow \mathbb{V}[\alpha](x, t^* = 10^{-8})$  for  $\alpha \in \{I, E, T_r, T\}$  on the Marshak wave test-problem. Truncation (29) is used

the truncation allows recovering, with a numerically acceptable agreement, the same mean and variance as the reference ni-ISMCM solver, even on the spatial profile of  $T_r$  (which was particularly problematic in the two previous sections 5.4.1–5.4.2). Figure 16 displays the spatial profiles of the realisations at the  $N = 15$  Gauss-Legendre points with ni-ISMCM and with the reconstructions

<sup>43</sup>Note that applying a particle dependent  $\theta$ -limitation is strictly equivalent to choosing truncation (26).



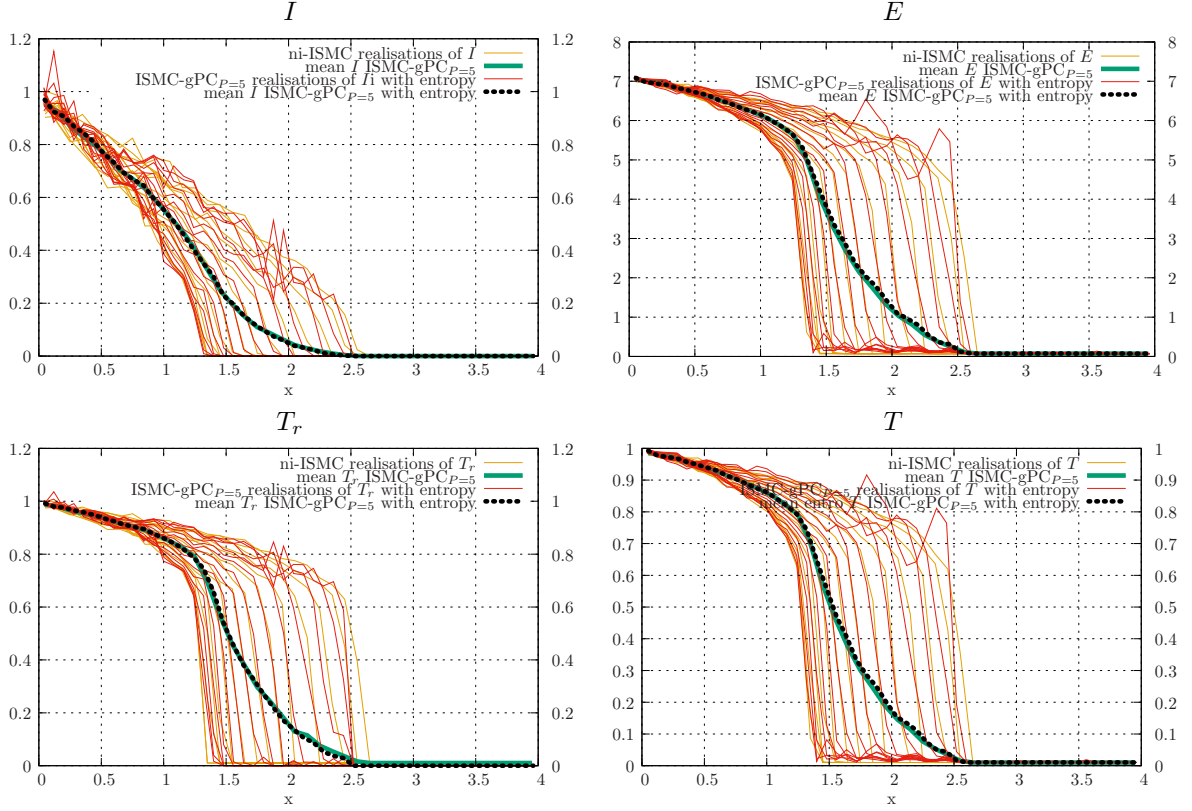


Figure 16: Comparison of the 15 realisations  $x \rightarrow \alpha(x, t^* = 10^{-8}, X_i), \alpha \in \{I, E, T_r, T\}$  obtained with ni-ISM at the  $N = 15$  Gauss-Legendre points  $(X_i)_{i \in \{1, \dots, 15\}}$  to the reconstructed (via truncation (29)) ones from ISMC-gPC $_{P=5}$  at the same points. The means  $x \rightarrow \mathbb{E}[\alpha](x, t^*), \alpha \in \{I, E, T_r, T\}$  are also displayed.

obtained from ISMC-gPC $_{P=5}$ . We can see that even if some spurious oscillations are still observable, especially on the spatial profiles of  $E$  and  $T$ , the oscillations of  $T_r$  are better controlled in the vicinity of  $x \in [1.5, 2.5]$  for low radiative temperatures.

From the results of figures 15–16, we can attest that the use of

- an accretiveness-preserving strategy is mandatory for mathematical wellposedness (otherwise the code crashes),
- but a moment-matching one may be important for physical relevance (accuracy).

The IPM (truncation (29)) can ensure both properties (wellposedness and moment-matching) but is known to be costly [4, 1, 5, 8, 36] as it induces an additional computational cost: truncation (29) needs the computation of  $(\lambda^I, \lambda^E)$  from the moments of  $(I, E)$  which needs a nonlinear inversion. For the computations of figures 15–16, we claim that *the overall cost of the simulations is not strongly affected by the use of the IPM truncation (29): the cost is the same as the one of the results of section 5.4.1 obtained with truncation (26)*. Let us explain how the results of figures 15–16 are produced in practice:

- first, the calculations are performed with truncation (26), just as in section 5.4.1. This truncation is cheap and ensures wellposedness.
- Second, truncation (29) is used only for the *output* profiles for which accurate results are needed (i.e. it is applied *offline/non-intrusively*). It is used in order to preserve the moments of the population of MC particles at the times and locations of interest.

In other words, the inversion needed by truncation (29) *is not* systematically performed within each cell at each time step as in [8] for example. The idea to resort to the cheapest wellposedness-preserving strategy (here truncation (26)) together with a moment-matching procedure (truncation (29)) only where and when needed. In a sense, this kind of optimised strategy is inspired by both [6] (with a cheap and efficient wellposedness-preserving strategy) and [4] (in which computations on reduced models obtained from hyperbolic systems are accelerated by only applying IPM when relevant).

**Remark 5.1 (Few implementation details on the *offline* inversion for truncation (29)).** *Even if performed offline, the inversion strategy may be difficult: we encountered few issues which deserve to be tackled for the sake of reproducibility of the numerical results of this section: due to the noisy numerical MC integration in certain cells and times of interest, the hessian matrix within the newton may be ill-conditioned for too important orders  $P$ . When this situation occurs (it is detected if too many newton iterations are performed), the inversion is carried out on less polynomial moments ( $P \leftarrow P - 1$ ), i.e. on a smaller hessian matrix, until convergence is ensured.*

It is interesting noticing that the Marshak wave problem of this section put forward difficulties which were not encountered in a hyperbolicity-preserving context [5, 6, 1, 4]: this is probably due to the strong stiffness of the nonlinear system we consider here. Still, an astute combination of the existing techniques (hyperbolicity-preserving strategy [5, 6] together with a moment-matching one [8, 1, 4]) allows obtaining mathematically wellposed reduced model which are physically relevant together with being numerically efficient.

## 6. Conclusion

In this paper, we build wellposed generalised Polynomial Chaos (gPC) based reduced models for photonics and solve them with a Monte-Carlo (MC) scheme. In the first part of the paper, care is taken to highlight under which conditions a reduced model (gPC based or not) is wellposed. The analysis is carried out thanks to an astute analogy between the construction of reduced models for uncertainty quantification and the construction of reduced models for kinetic equations. In particular, the analysis leads to quite simple conditions for wellposedness: bounds on the material temperature and on the radiative intensity must be satisfied. Several strategies are tested and analysed in order to control those bounds. They are mainly inspired from the literature on uncertainty quantification for hyperbolic systems of conservation laws. For the resolution of the truncated reduced models, an astute combination of the ISMC scheme and of MC-gPC is performed. The description of the scheme is made by highlighting where, in an ISMC implementation, the MC-gPC modifications must be made. These modifications are simple and efficient in practice: we verify, at least numerically, that the ISMC-gPC scheme preserve interesting properties from both ISMC (no teleportation error, good behaviour in the equilibrium diffusion limit on affordable meshes, affordable time steps) and MC-gPC (fast convergence rates with respect to the truncation order, gain with respect to efficient non-intrusive methods). In a sense, the work of this paper also demonstrates that MC-gPC can be efficiently applied to a stiff nonlinear set of partial derivative equations and that the same parallel strategies as in the linear case can be applied if the MC resolution allows a fast convergence with respect to both the time and spatial discretisation. Several benchmarks are investigated in the last section. They emphasize the importance of relying on

- an ansatz preserving the wellposedness to obtain mathematically relevant surrogates (and avoid robustness difficulties),
- and on the relative importance of relying on a moment-preserving ansatz to obtain physically relevant models (converging toward the physical solution),

in order to implement computationally efficient codes for uncertainty propagation for photonics in low to moderate uncertain dimensions.

## Acknowledgments

The authors would like to thank Xavier Valentin for valuable discussions on the ISMC solver.

## References

- [1] J. Kusch, G. W. Alldredge, M. Frank, Maximum-principle-satisfying second-order intrusive polynomial moment scheme, arXiv preprint arXiv:1712.06966 (2017).
- [2] J. Kusch, M. Frank, Intrusive methods in uncertainty quantification and their connection to kinetic theory, *International Journal of Advances in Engineering Sciences and Applied Mathematics* (2018) 1–16.
- [3] J. Kusch, R. G. McClarren, M. Frank, Filtered stochastic galerkin methods for hyperbolic equations, *J. Comput. Phys.* (2018).

- [4] J. Kusch, J. Wolters, M. Frank, Intrusive acceleration strategies for uncertainty quantification for hyperbolic systems of conservation laws, *Journal of Computational Physics* 419 (2020) 109698. doi:<https://doi.org/10.1016/j.jcp.2020.109698>.  
URL <http://www.sciencedirect.com/science/article/pii/S0021999120304721>
- [5] L. Schlachter, F. Schneider, A hyperbolicity-preserving stochastic galerkin approximation for uncertain hyperbolic systems of equations, *arXiv preprint arXiv:1710.03587* (2017).
- [6] J. Dürrwächter, T. Kuhn, F. Meyer, L. Schlachter, F. Schneider, A hyperbolicity-preserving discontinuous stochastic galerkin scheme for uncertain hyperbolic systems of equations, *Journal of Computational and Applied Mathematics* 370 (2020) 112602. doi:10.1016/j.cam.2019.112602.  
URL <http://dx.doi.org/10.1016/j.cam.2019.112602>
- [7] L. Schlachter, F. Schneider, O. Kolb, Weighted essentially non-oscillatory stochastic galerkin approximation for hyperbolic conservation laws, *Journal of Computational Physics* 419 (2020) 109663. doi:<https://doi.org/10.1016/j.jcp.2020.109663>.  
URL <http://www.sciencedirect.com/science/article/pii/S002199912030437X>
- [8] B. Després, G. Poëtte, D. Lucor, Robust Uncertainty Propagation in Systems of Conservation Laws with the Entropy Closure Method, Vol. 92 of *Lecture Notes in Computational Science and Engineering, Uncertainty Quantification in Computational Fluid Dynamics*, 2013.
- [9] G. Poëtte, X. Valentin, A new implicit monte-carlo scheme for photonics (without tele-  
portation error and without tilts), *Journal of Computational Physics* 412 (2020) 109405.  
doi:<https://doi.org/10.1016/j.jcp.2020.109405>.  
URL <http://www.sciencedirect.com/science/article/pii/S0021999120301790>
- [10] G. Poëtte, A gPC-intrusive Monte-Carlo scheme for the resolution of the uncertain  
linear Boltzmann equation, *Journal of Computational Physics* 385 (2019) 135 – 162.  
doi:<https://doi.org/10.1016/j.jcp.2019.01.052>.  
URL <http://www.sciencedirect.com/science/article/pii/S002199911930110X>
- [11] D. Mihalas, B. W. Mihalas, *Foundations of Radiation Hydrodynamics*, Dover Publications,  
1999.
- [12] J. Castor, *Radiation hydrodynamics*, Cambridge University Press, 2004.
- [13] G. C. Pomraning, *The equations of radiation hydrodynamics*, Dover Publications, 1973.
- [14] W. Li, C. Liu, Y. Zhu, J. Zhang, K. Xu, Unified gas-kinetic wave-particle methods  
iii: Multiscale photon transport, *Journal of Computational Physics* 408 (2020) 109280.  
doi:<https://doi.org/10.1016/j.jcp.2020.109280>.  
URL <http://www.sciencedirect.com/science/article/pii/S0021999120300541>
- [15] A. Marshak, A. Davis, *3D Radiative Transfer in Cloudy Atmospheres, Physics of Earth and  
Space Environments*, Springer Berlin Heidelberg, 2005.  
URL <https://books.google.fr/books?id=EdgBysLXndwC>
- [16] A. D. Klose, U. Netz, J. Beuthan, A. H. Hielscher, Optical tomography using the  
time-independent equation of radiative transfer — part 1: forward model, *Journal  
of Quantitative Spectroscopy and Radiative Transfer* 72 (5) (2002) 691 – 713.

doi:[https://doi.org/10.1016/S0022-4073\(01\)00150-9](https://doi.org/10.1016/S0022-4073(01)00150-9).

URL <http://www.sciencedirect.com/science/article/pii/S0022407301001509>

- [17] B. Mercier, Application of accretive operators theory to the radiative transfer equations, *SIAM Journal on Mathematical Analysis* 18 (2) (1987) 393–408. arXiv:<https://doi.org/10.1137/0518030>, doi:10.1137/0518030. URL <https://doi.org/10.1137/0518030>
- [18] R. A. Todor, C. Schwab, Karhunen-Loève approximation of random fields by generalized fast multipole methods, *J. Comp. Phys.* 217 (1) (2006) 100–122.
- [19] M. Meyer, H. Matthies, Efficient model reduction in non-linear dynamics using the Karhunen-Loève expansion and dual-weighted-residual methods, *Comp. Meth. Appl. Mech. Eng. Informatikbericht 2003-08*, TU Braunschweig, Germany (2004).
- [20] J. Mercer, Functions of Positive and Negative Type and their Connection with the Theory of Integral Equations, *Philos. Trans. Roy. Soc.* 209 (1909).
- [21] R. Lebrun, A. Dutfoy, A Generalization of the Nataf Transformation to Distributions with Elliptical Copula, *Prob. Eng. Mech.* 24,2 (2009) 172–178.
- [22] R. Lebrun, A. Dutfoy, An Innovating Analysis of the Nataf Transformation from the Copula viewpoint, *Prob. Eng. Mech.* 24,3 (2009) 312–320.
- [23] B. Iooss, P. Lemaître, A Review on Global Sensitivity Analysis Methods, Dellino, Gabriella and Meloni, Carlo, Springer US, Boston, MA, 2015, pp. 101–122. doi:10.1007/978-1-4899-7547-8\_5. URL [http://dx.doi.org/10.1007/978-1-4899-7547-8\\_5](http://dx.doi.org/10.1007/978-1-4899-7547-8_5)
- [24] B. Lapeyre, E. Pardoux, R. Sentis, Méthodes de Monte Carlo pour les équations de transport et de diffusion, no. 29 in *Mathématiques & Applications*, Springer-Verlag, 1998.
- [25] A. B. Wollaber, Four decades of implicit monte carlo, *Journal of Computational and Theoretical Transport* (2016).
- [26] M. A. Cleveland, N. Gentile, Mitigating teleportation error in frequency-dependent hybrid implicit monte carlo diffusion methods, *Journal of Computational and Theoretical Transport* (2014).
- [27] A. G. Irvine, I. D. Boyd, N. A. Gentile, Reducing the spatial discretization error of thermal emission in implicit monte carlo simulations, *Journal of Computational and Theoretical Transport* 45 (1-2) (2016) 99–122. arXiv:<https://doi.org/10.1080/23324309.2015.1060245>, doi:10.1080/23324309.2015.1060245. URL <https://doi.org/10.1080/23324309.2015.1060245>
- [28] J.-F. Clouet, G. Samba, Asymptotic diffusion limit of the symbolic monte-carlo method for the transport equation, *Journal of Computational Physics* 195 (1) (2004) 293 – 319. doi:<https://doi.org/10.1016/j.jcp.2003.10.008>. URL <http://www.sciencedirect.com/science/article/pii/S0021999103005333>
- [29] M. S. McKinley, E. D. B. III, A. Szoke, Comparison of implicit and symbolic implicit monte carlo line transport with frequency weight vector extension, *Journal of Computational Physics* (2003).

- [30] G. Poëtte, Spectral convergence of the generalized Polynomial Chaos reduced model obtained from the uncertain linear Boltzmann equation, Preprint submitted to Mathematics and Computers in Simulation (2019).
- [31] J. A. Carrillo, M. Zanella, Monte Carlo gPC methods for diffusive kinetic flocking models with uncertainties, Vietnam Journal of Mathematics 47:931-954 (2019).
- [32] L. Pareschi, M. Zanella, Monte carlo stochastic galerkin methods for the boltzmann equation with uncertainties: Space-homogeneous case, Journal of Computational Physics 423 (2020) 109822. doi:<https://doi.org/10.1016/j.jcp.2020.109822>.  
URL <http://www.sciencedirect.com/science/article/pii/S0021999120305969>
- [33] L. Pareschi, An introduction to uncertainty quantification for kinetic equations and related problems (2020). arXiv:2004.05072.
- [34] G. Poëtte, E. Brun, Efficient uncertain  $k_{\text{eff}}$  computations with the Monte Carlo resolution of generalised Polynomial Chaos Based reduced models, Preprint (2020).
- [35] D. Dureau, G. Poëtte, Hybrid Parallel Programming Models for AMR Neutron Monte-Carlo Transport, in: Joint International Conference on Supercomputing in Nuclear Applications + Monte-Carlo, no. 04202 in Parallelism and HPC, Monte-Carlo, 2013.
- [36] G. Poëtte, Contribution to the mathematical and numerical analysis of uncertain systems of conservation laws and of the linear and nonlinear boltzmann equation, Habilitation à diriger des recherches, Université de Bordeaux 1 (Sep. 2019).  
URL <https://hal.archives-ouvertes.fr/tel-02288678>
- [37] P. Spanos, R. G. Ghanem, Stochastic Finite Element Expansion for Random Media, ASCE J. Eng. Mech. 115 (5) (1989) 1035–1053.
- [38] R. G. Ghanem, W. Brzakala, Stochastic finite-element analysis of soil layers with random interface, ASCE J. Eng. Mech. 122 (1996) 361–369.
- [39] R. G. Ghanem, J. Red-Horse, Propagation of Uncertainty in Complex Physical Systems using a Stochastic Finite Elements Approach, Physica D 133 (1999) 137–144.
- [40] L. Mathelin, O. P. L. Maître, A Posteriori Error Analysis for Stochastic Finite Element Solutions of Fluid Flows with Parametric Uncertainties, ECCOMAS CFD (2006).
- [41] M. Jardak, C. Su, G. Karniadakis, Spectral polynomial chaos solutions of the stochastic advection equation, Journal of Scientific Computing 17 (1) (2002) 319–338.
- [42] H. N. N. O.P. Le Maître, O. M. Knio, R. G. Ghanem, A Stochastic Projection Method for Fluid Flow I: Basic Formulation, J. Comp. Phys. 173 (2001) 481–511.
- [43] H. N. N. O.P. Le Maître, O. M. Knio, R. G. Ghanem, A Stochastic Projection Method for Fluid Flow II: Random Process, J. Comp. Phys. 181 (2002) 9–44.
- [44] H. Matthies, A. Keese, Galerkin methods for linear and nonlinear elliptic stochastic PDEs, Comp. Meth. Appl. Mech. Eng. 31 (1-2) (2003) 179–191.

- [45] M. K. Deb, I. M. Babuska, J. T. Oden, Solution of Stochastic Partial Differential Equations using Galerkin Finite Element Techniques, *Comp. Meth. Appl. Mech. Engrg.* 190 (2001) 6359–6372.
- [46] D. Gottlieb, D. Xiu, Galerkin Method for Wave Equations with Uncertain Coefficients, *Commun. Comp. Phys.* 3 (2008) 505–518.
- [47] O. P. L. Maître, O. M. Knio, Uncertainty Propagation using Wiener-Haar Expansions, *J. Comp. Phys.* 197 (2004) 28–57.
- [48] X. Wan, G. Karniadakis, Multi-Element generalized Polynomial Chaos for Arbitrary Probability Measures, *SIAM J. Sci. Comp.* 28(3) (2006) 901–928.
- [49] M. I. Gerritsma, J.-B. van der Steen, P. Vos, G. E. Karniadakis, Time-dependent generalized polynomial chaos., *J. Comput. Physics* (2010) 8333–8363.
- [50] B. Després, B. Perthame, Uncertainty propagation;intrusive kinetic formulations of scalar conservation laws, *SIAM/ASA Journal on Uncertainty Quantification* 4 (1) (2016) 980–1013.  
URL <http://hal.upmc.fr/hal-01146188>
- [51] R. Abgrall, A Simple, Flexible and Generic Deterministic Approach to Uncertainty Quantifications in Non Linear Problems: Application to Fluid Flow Problems, *Rapport de Recherche INRIA* (2007).
- [52] M. G. Crandall, T. M. Liggett, Generation of semi-groups of nonlinear transformations on general Banach spaces., *Am. J. Math.* 93 (1971) 265–298.
- [53] M. Raissi, P. Perdikaris, G. Karniadakis, Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations, *Journal of Computational Physics* 378 (2019) 686 – 707.  
doi:<https://doi.org/10.1016/j.jcp.2018.10.045>.  
URL <http://www.sciencedirect.com/science/article/pii/S0021999118307125>
- [54] E. E. Lewis, W. F. M. Jr., *Computational Methods of Neutron Transport*, John Wiley and Son New York, 1984.
- [55] Ernst, Oliver G., Mugler, Antje, Starkloff, Hans-Jörg, Ullmann, Elisabeth, On the convergence of generalized polynomial chaos expansions, *ESAIM: M2AN* 46 (2) (2012) 317–339.  
doi:[10.1051/m2an/2011045](https://doi.org/10.1051/m2an/2011045).  
URL <https://doi.org/10.1051/m2an/2011045>
- [56] J. A. F. Jr., The calculation of nonlinear radiation transport by a monte carlo method, *Tech. rep.*, Lawrence Radiation Laboratory, University of California (1961).
- [57] J. A. Fleck, J. D. Cummings, An implicit monte-carlo scheme for calculating time and frequency dependent nonlinear radiation transport, *Journal of Computational Physics* (1971).
- [58] G. Poëtte, X. Valentin, A. Bernede, Canceling teleportation error in legacy imc code for photonics (without tilts, with simple minimal modifications), *Journal of Computational and Theoretical Transport* 49 (4) (2020) 162–194. arXiv:<https://doi.org/10.1080/23324309.2020.1785893>, doi:[10.1080/23324309.2020.1785893](https://doi.org/10.1080/23324309.2020.1785893).  
URL <https://doi.org/10.1080/23324309.2020.1785893>

- [59] R. P. Smedley-Stevenson, R. G. McClarren, Asymptotic diffusion limit of cell temperature discretisation schemes for thermal radiation transport, *Journal of Computational Physics* 286 (2015) 214 – 235. doi:<https://doi.org/10.1016/j.jcp.2013.10.038>.  
URL <http://www.sciencedirect.com/science/article/pii/S0021999113007146>
- [60] C. Ahrens, E. Larsen, A semi-analog monte carlo method for grey radiative transfer problems, in: *Proceedings of the ANS Topical Meeting: International Conference on Mathematical Methods to Nuclear Applications*, 2001.
- [61] A. Bernede and G. Poëtte, An Unsplit Monte-Carlo solver for the resolution of the linear Boltzmann equation coupled to (stiff) Bateman equations, *Journal of Computational Physics* 354 (2018) 211 – 241. doi:<https://doi.org/10.1016/j.jcp.2017.10.027>.  
URL <http://www.sciencedirect.com/science/article/pii/S0021999117307805>
- [62] J. D. Densmore, E. W. Larsen, Asymptotic equilibrium diffusion analysis of time-dependent monte carlo methods for grey radiative transfer, *Journal of Computational Physics* (2004).
- [63] P. Vos, Time dependent polynomial chaos, Master of science thesis, Delft University of technology, Faculty of Aerospace engineering (2006).
- [64] I. Sobol, Sensitivity estimates for nonlinear mathematical models, *Matematicheskoe Modelirovanie* 2 (1990) 112–118, in Russian, translated in English in Sobol, I. (1993). Sensitivity analysis for non-linear mathematical models. *Mathematical Modeling and Computational Experiment (Engl. Transl.)*, 1993, 1, 407–414.
- [65] A. Saltelli, Making best use of model evaluations to compute sensitivity indices, *Computer Physics Communications* 145 (2) (2002) 280 – 297. doi:[https://doi.org/10.1016/S0010-4655\(02\)00280-1](https://doi.org/10.1016/S0010-4655(02)00280-1).  
URL <http://www.sciencedirect.com/science/article/pii/S0010465502002801>
- [66] G. Blatman, B. Sudret, Efficient computation of global sensitivity indices using sparse polynomial chaos expansions, *Rel. Eng. Syst. Saf.* 95 (2010) 1216–1229.

## Appendix A. Details about the Heaviside problem of section 5.2 ( initial and boundary conditions, test-case justifications)

In this paper, we intensively make use of the configuration presented in this appendix. The initial and boundary conditions together with the problem justifications are provided here for both, the sake of conciseness of the paper and of reproducibility of the results.

The Heaviside problem considered in section 5.2 can be described as follows: let us consider a 1D spatial domain such that  $x \in \Omega = [0, 1]$ . The domain is filled with a diffusive media  $\sigma_t = 2000$ , with no (physical) scattering, i.e.  $\sigma_s = 0$  and  $\sigma_t = \sigma_a$ . Initially, a Heaviside of temperatures at equilibrium is set in the middle  $[0.4, 0.6]$  of domain  $\Omega = [0, 1]$ . In other words, we have at  $t = 0$ :

$$T_m(x, t = 0) = T_r(x, t = 0) = 2.3 \times 10^7 \mathbf{1}_{[0.4, 0.6]}(x) + 2.3 \times 10^4 \mathbf{1}_{[0, 1] \setminus [0.4, 0.6]}(x).$$

Note that  $\mathbf{1}_\Omega(x)$  denotes the indicatrix of domain  $\Omega$ . The initial condition is displayed in figure A.17 together with the solution of system (38) at final time  $T = 10^{-8}$ . This reference solution has been obtained solving (38) with a deterministic solver (with a fine mesh). Note that for time



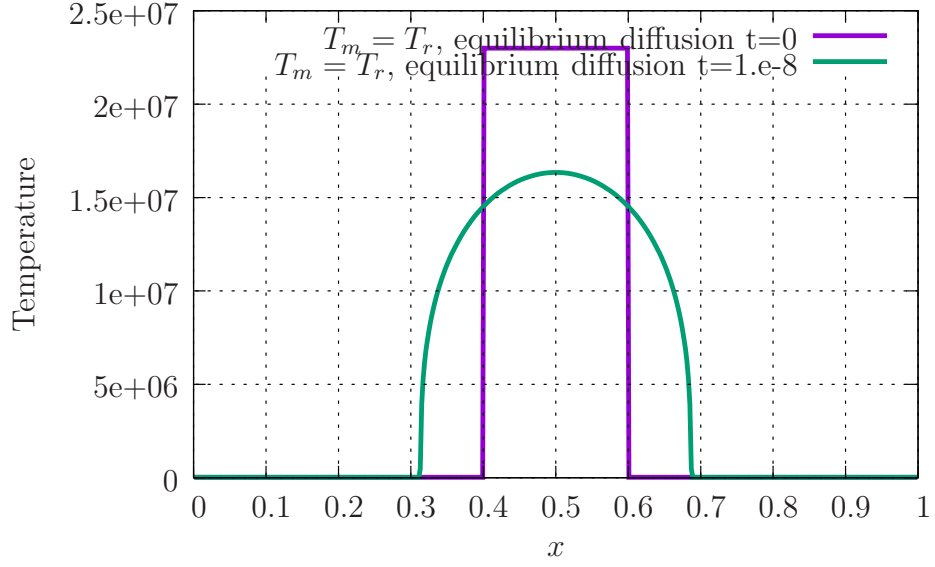


Figure A.17: Initial and final spatial profile of the temperatures ( $T_m = T_r$ ) in the equilibrium diffusion limit for the Heaviside test-problem of this paper.

$t \in [0, T]$ , the solution does not reach the boundaries. The radiative constant is set to  $a = 10^{-14}$ , the speed of light to  $c = 3 \times 10^{10}$ . A perfect gas is considered so that  $E(T_m) = \rho C_v T_m$  with  $\rho = 20$ ,  $C_v = 4 \times 10^7$ .