# WEIGHTED PARTICLES IN THE FINITE POINTSET METHOD*

Michael Schreiner

Department of Mathematics
University of Kaiserslautern
W–6750 Kaiserslautern, Germany

## ABSTRACT

Using particle methods to solve the Boltzmann equation for rarefied gases numerically, in realistic streaming problems, huge differences in the total number of particles per cell arise. In order to overcome the resulting numerical difficulties the application of a weighted particle concept is well–suited. The underlying idea is to use different particle masses in different cells depending on the macroscopic density of the gas. Discrepance estimates and numerical results are given.

## 1  INTRODUCTION

The Boltzmann equation for rarefied gases [5] accuring in reentry problems in the earth's atmosphere, is usually numerically solved by particle methods, for instance the finite pointset method (FPM) [6], [7] or the Bird–algorithm [4]. Calculating relevant problems these algorithms are very time consuming and need a great amount of memory. Therefore improvements in the numerical realizations are desirable. One improvement hereby, the use of weighted particles, is presented here.

In order to explain the main ideas we first give a brief description of the FPM (cf. [6]): The initial conditions are approximated by a finite set of particles called an ensemble. Having chosen a time step $\Delta t$, the free flow reads as follows: For each particle of position $x_i$ and velocity $v_i$ we set $x_i \leftarrow x_i + \Delta t\, v_i$. Hereby one has to pay attention to the boundary conditions. The position space is divided into cells. In each cell there are determined pairs of collision partners. If and how they collide depends on a collision parameter which has to be calculated. After the collision process the free flow follows again, and so on.

---

Concerning a typical problem as the streaming of gas around an obstacle, it is clear that there are regions with a high macroscopic density and regions with a low one. On account of this there are cells with very many particles causing a great numerical effort, and cells of only a few ones which often results in a bad numerical accuracy.

The main idea is using particles of different masses in different regions, depending on the macroscopic density of the gas. Particles of smaller mass are used in low density regions and vice versa. Hence the number of particles is better suited both for numerical efficiency and accuracy.

To perform this ideas we introduce two procedures which can quite easily be fitted into existing programs: The procedure **MASSHA** (mass handling) calculates, dependent on the macroscopic density, a good choice for the desired particle mass in each cell. During the free flow the particles will change their cells. Therefore particles of different masses might be in the same cell. Since the collision should be performed with particles of the same mass, the procedure **SPLIPA** (split and paste) transforms a set of particles with different masses within a cell into an ensemble consisting of particles (nearly) all having the desired mass.

From a mathematical point of view the procedure **MASSHA** is not very challenging. Therefore we give just a few comments: In the algorithm we allow only particle masses of integer–values (giving some advantages in the implementation). The procedure **MASSHA** prescribes only such masses $m^*$ which have a representation $m^* = 2^j$ for a $j \in I\!N$. Therefore it is not always possible for **SPLIPA** to transform a given ensemble in such a way that in fact all particles have the prescribed mass. If particles are left they will not take part in the collision process.

The mathematical interesting part is the procedure **SPLIPA** investigated in the following.

The paper is organized as follows: After the introduction the second section will give some definitions. The splitting and pasting of particles will be considered in the third section. The topic of the fourth section is the **SPLIPA**–procedure and an estimation of discrepancy in it. Finally we will present some numerical results.

For more detailed proofs and more algorithmic aspects the reader is referred to [9].

## 2  DEFINITIONS AND NOTATIONS

As a *particle* we understand a tupel $\tau = (m, x, v)$ of the mass $m$, the position $x$ and the velocity $v$. Herein we assume that $m > 0$. If the position of the particle is not of interest, we just write $\tau = (m, v)$. An *ensemble of particles* is a finite family $((\tau_i)_{i=1,...,n})$ of particles, where we do not consider any ordering. Speaking formally we say that two families $((\tau_i)_{i=1,...,n})$ and $((\tilde{\tau}_i)_{i=1,...,n})$ are *equivalent*, if there exists a permutation $\pi \in S_n$ such that $\tau_{\pi(i)} = \tilde{\tau}_i$ for all $i = 1, \ldots, n$. Therefore an ensemble is just an equivalence class under this relation. As usual we identify the class with a member of it.

Given $\mathcal{E} = (\tau_1, \ldots, \tau_n)$ and $A \subset \{1, \ldots, n\}$ we call $((\tau_i)_{i \in A})$ a *sub–ensemble* of $\mathcal{E}$. With $n$ ensembles $\mathcal{E}_i = (\tau_i^1, \ldots, \tau_i^{l_i})$, $i = 1, \ldots, n$, we define $\mathcal{E} := \cup_{i=1}^n \mathcal{E}_i$ to be $\mathcal{E} = (\tau_1^1, \ldots, \tau_1,^{l_1}, \ldots, \ldots, \tau_n^1, \ldots, \tau_n^{l_n})$ and call $\mathcal{E} := \cup_{i=1}^n \mathcal{E}_i$ a *partition* of $\mathcal{E}$.

An ensemble $\mathcal{E} = (\tau_1, \ldots, \tau_n)$ of particles $\tau_i = (m_i, v_i)$ defines a discrete measure $\Phi_\mathcal{E} = \sum_{i=1}^n m_i \delta_{v_i}$. Two ensembles $\mathcal{E}$ and $\tilde{\mathcal{E}}$ shall be *equivalent* (written $\mathcal{E} \sim \tilde{\mathcal{E}}$) if

$\Phi_{\mathcal{E}} = \Phi_{\tilde{\mathcal{E}}}$. The following statements hold: Given two ensembles $\mathcal{E}$ and $\mathcal{F}$ we have

$$\Phi_{\mathcal{E} \cup \mathcal{F}} = \Phi_{\mathcal{E}} + \Phi_{\mathcal{F}} \qquad (1)$$

Given ensembles $\mathcal{E}, \tilde{\mathcal{E}}, \mathcal{F}, \tilde{\mathcal{F}}$ with $\mathcal{E} \sim \tilde{\mathcal{E}}$ and $\mathcal{F} \sim \tilde{F}$ we have

$$\mathcal{E} \cup \mathcal{F} \sim \tilde{\mathcal{E}} \cup \tilde{\mathcal{F}} \qquad (2)$$

For a later use let us define some physical quantities of ensembles: Given an ensemble $\mathcal{E} = ((m_i, v_i)_{i=1,\ldots,n})$, with $v_i = (v_i^1, \ldots, v_i^k) \in \mathbf{R}^k$ we set $M(\mathcal{E}) = \sum_{i=1}^n m_i$, $I^j(\mathcal{E}) = \sum_{i=1}^n m_i v_i^j$ and $E^j(\mathcal{E}) = \sum_{i=1}^n m_i (v_i^j)^2$. $M(\mathcal{E})$ is the *mass*, $I^j(\mathcal{E})$ the *momentum* of the $j$-th component and the *energy* is given by $(E^1(\mathcal{E}) + \ldots + E^k(\mathcal{E}))/2$. For one dimensional velocities we sometimes omit the indices in $I^j$ and $E^j$.

A great advantage (at least from a mathematical point of view) of the FPM is that (under certain assumptions) the convergence to solutions of the Boltzmann equation can be proven (cf. [2] and [3]). One tool herein is the concept of discrepancy, first introduced by H. Weyl [10]. A detailed description of this concept together with applications to several fields in mathematics can be found in [8]. We briefly give some definitions:

Let $x = (x^1, \ldots, x^k)$ and $y = (y^1, \ldots, y^k)$ denote elements of $\mathbf{R}^k$. We introduce the usual semi order on $\mathbf{R}^k$ by $x \leq y$ if and only if $x^i \leq y^i$ for all $i = 1, \ldots, k$. For $x \in \mathbf{R}^k$ we define $R(x) := \{y \in \mathbf{R}^k | y \leq x\}$.

**Definition 1** Let $\mu$ and $\omega$ be two measures on $\mathbf{R}^k$ with $\mu(\mathbf{R}^k) < \infty$ and $\omega(\mathbf{R}^k) < \infty$. For $x \in \mathbf{R}^k$ we define the *local discrepancy* of $\mu$ and $\omega$ in $x$ by

$$d_{\mu,\omega}(x) := \mu(R(x)) - \omega(R(x))$$

The *(extreme) discrepancy* of $\mu$ and $\omega$ is given by

$$D(\mu, \omega) := \sup_{x \in R^k} |d_{\mu,\omega}(x)|$$

For a more detailed introduction and the relations with the weak convergence of measures the reader is referred to [3].

Given ensembles $\mathcal{E}, \mathcal{F}$ and $\mathcal{G}$ as before with particles $(m_i, v_i)$, $v_i \in \mathbf{R}^k$, some simple calculations yield the following statements:

$$D(\Phi_{\mathcal{E} \cup \mathcal{G}}, \Phi_{\mathcal{F} \cup \mathcal{G}}) \;=\; D(\Phi_{\mathcal{E}}, \Phi_{\mathcal{F}}) \qquad (3)$$

$$D(\Phi_{\mathcal{E} \cup \mathcal{G}}, \Phi_{\mathcal{F}}) \;\leq\; D(\Phi_{\mathcal{E}}, \Phi_{\mathcal{F}}) + M(\mathcal{G}) \qquad (4)$$

To conclude this section let us explain the use of the concept of discrepancy in the SPLIPA–procedure. The SPLIPA–procedure replaces one ensemble $\mathcal{E}$ by a new one $\mathcal{F}$ with certain properties. As we have convergence of the FPM we can assume $D(\Phi_{\mathcal{E}}, f dv)$ to be small, whereas $f$ is the solution of the Boltzmann equation. By use of the triangle inequality we have

$$D(\Phi_{\mathcal{F}}, f dv) \leq D(\Phi_{\mathcal{E}}, \Phi_{\mathcal{F}}) + D(\Phi_{\mathcal{E}}, f dv)$$

If we can achieve that $D(\Phi_{\mathcal{E}}, \Phi_{\mathcal{F}})$ is small enough we can ensure the convergence of the new method with weighted particles. Hence, the aim of SPLIPA is clear: replace $\mathcal{E}$ by $\mathcal{F}$ such that the discrepancy $D(\Phi_{\mathcal{E}}, \Phi_{\mathcal{F}})$ is small. If we have $D(\Phi_{\mathcal{E}}, \Phi_{\mathcal{F}}) \to 0$ as the number of particles increases, the convergence of the new algorithm is obvious.

# 3 SPLITTING AND PASTING OF PARTICLES

As a preparation for the SPLIPA-procedure we show in this section how to transform an arbitrary ensemble in another one consisting of two particles. Hereby we consider first the one dimensional case and then generalize the results to higher dimensions. Since the SPLIPA-procedure will work within a cell where spatial homogeinity is always assumed, we do not care about the positions of the particles.

## 3.1 The One Dimensional Case

Given an ensemble $\mathcal{E} = ((m_i, v_i)_{i=1,\ldots,n})$ with $m_i > 0$ and $v_i \in \boldsymbol{R}$ the task is to transform $\mathcal{E}$ into an ensemble $\mathcal{F} = ((\mu_1, \nu_1), (\mu_2, \nu_2))$. Therefore we want to preserve the collision invariants (cf. [5]), namely the mass, the momentum and the energy. If we set $M := M(\mathcal{E})$, $I := I(\mathcal{E})$ and $E := E(\mathcal{E})$ we have

$$M(\mathcal{F}) = M \tag{5}$$

$$I(\mathcal{F}) = I \tag{6}$$

$$E(\mathcal{F}) = E \tag{7}$$

Because of (5) we can write $\mu_1 = \lambda M$ and $\mu_2 = (1 - \lambda)M$ for a $\lambda \in (0, 1)$. Assuming $\lambda$ is given, we have

$$\lambda \nu_1 + (1 - \lambda)\nu_2 = \frac{I}{M}$$

$$\lambda(\nu_1)^2 + (1 - \lambda)(\nu_2)^2 = \frac{E}{M}$$

Simple calculations show the general existence of solutions of this quadratic equation, since

$$\frac{E}{M} - \frac{I^2}{M^2} = \frac{1}{M} \sum m_i \left( v_i - \frac{I}{M} \right)^2 \geq 0$$

In detail we have:

**Theorem 2** *For an ensemble* $\mathcal{E} = ((m_i, v_i)_{i=1,\ldots,n})$ *with* $m_i > 0$ *and* $v_i \in \boldsymbol{R}$ *set* $M := M(\mathcal{E})$, $I := I(\mathcal{E})$ *and* $E := E(\mathcal{E})$. *Given* $\lambda \in (0, 1)$ *let* $\mu_1 = \lambda M$ *and* $\mu_2 = (1 - \lambda)M$. *Define*

$$\nu_{11} = \frac{I}{M} + \sqrt{\frac{1 - \lambda}{\lambda} \left( \frac{E}{M} - \frac{I^2}{M^2} \right)}$$

$$\nu_{21} = \frac{I}{M} - \sqrt{\frac{\lambda}{1 - \lambda} \left( \frac{E}{M} - \frac{I^2}{M^2} \right)}$$

$$\nu_{12} = \frac{I}{M} - \sqrt{\frac{1 - \lambda}{\lambda} \left( \frac{E}{M} - \frac{I^2}{M^2} \right)}$$

$$\nu_{22} = \frac{I}{M} + \sqrt{\frac{\lambda}{1 - \lambda} \left( \frac{E}{M} - \frac{I^2}{M^2} \right)}$$

*Then for* $\mathcal{F}_1 = ((\mu_1, \nu_{11}), (\mu_2, \nu_{21}))$ *and* $\mathcal{F}_2 = ((\mu_1, \nu_{12}), (\mu_2, \nu_{22}))$ *the following holds:*

4

*(i) With the given masses the ensembles $\mathcal{F}_1$ and $\mathcal{F}_2$ are the only ones which fulfill (5)–(7).*

*(ii) $\mathcal{F}_1 = \mathcal{F}_2$ if and only if $v_1 = \ldots = v_n$ or $\lambda = 1/2$.*

For the important case $\lambda = 1/2$ we introduce the following notations: $\mu := M/2$, $\delta := \sqrt{E/M - I^2/M}$, $\nu_1 := I/M - \delta$ and $\nu_2 := I/M + \delta$. Hereby the only solution of (5)–(7) is $\mathcal{F} = ((\mu, \nu_1), (\mu, \nu_2))$.

The following definition takes care of this situation:

**Definition 3**    For a given $\mathcal{E} = ((\tau_i)_{i=1,\ldots,n})$ as in Theorem 2 (with $\lambda = 1/2$) we have $\mathcal{F} = ((\mu, \nu_1), (\mu, \nu_2))$ as described above. With the use of $a := \min\{v_1, \ldots, v_n, \nu_1\}$ and $b := \max\{v_1, \ldots, v_n, \nu_2\}$ we define

$$U(\mathcal{E}) = U(\tau_1, \ldots, \tau_n) := [a, b]$$

As a preparation for the estimation of discrepancy, we prove

**Theorem 4**    *Given an ensemble $\mathcal{E} = (\tau_i)_{i=1,\ldots,n}$ let $M$, $I$ and $E$ be as before. The ensemble $\mathcal{F} = ((\mu, \nu_1), (\mu, \nu_2))$ and $U(\mathcal{E}) = [a, b]$ are defined as described earlier ($\lambda = 1/2$). Then for all $w \in [a, b]$:*

$$\left| \sum_{v_i \leq w} m_i - \#\{v_i \,|\, v_i \leq w\} \,\mu \right| \leq \frac{1}{2} M = \mu \tag{8}$$

**Proof:**    We give an outline of the proof. First one can show that it is sufficient to prove the inequalities

$$\sum_{v_i < \nu_1} m_i \leq \mu \tag{9}$$

$$\sum_{v_i \leq \nu_2} m_i \geq \mu \tag{10}$$

For symmetry reasons it is even enough to show (9). Furthermore it is obvious that one can assume $I = 0$.

Now, assume that (9) is wrong which means ($I = 0$)

$$\sum_{v_i \leq -\delta} m_i > \mu$$

After sorting the $v_i$'s, we have the existence of $l \in \mathbb{N}$ so that

$$
\begin{aligned}
v_i &\leq -\delta &&\text{for} \quad i = 1, \ldots, l \\
v_i &> -\delta &&\text{for} \quad i = l+1, \ldots, n
\end{aligned}
$$

With $M_l := \sum_{i=1}^{l} m_i$ we can prove

$$\left( \sum_{i=l+1}^{n} \frac{m_i}{M - M_l} v_i \right)^2 > \delta^2$$

5

On the other hand we have

$$\sum_{i=l+1}^{n} \frac{m_i}{M - M_l} {v_i}^2 \leq \delta^2$$

and using the convexity of the mapping $x \mapsto x^2$ we get a contradiction. Hence the proof is complete. ∎

## 3.2 The Higher Dimensional Case

Now we assume that the particles of a given ensemble $\mathcal{E} = ((m_i, v_i)_{i=1,\ldots,n})$ have $k$–dimensional velocities. Essential for the calculations in the one dimensional case were the conservation equalities (5)–(7). Now, we generalize these equations as follows: For $\mathcal{E}$ we set $M := M(\mathcal{E})$, $I^j := I^j(\mathcal{E})$ and $E^j := E^j(\mathcal{E})$ for $j = 1, \ldots, k$. We search for an ensemble $\mathcal{F}$ with

$$M(\mathcal{F}) = M \tag{11}$$

$$I^j(\mathcal{F}) = I^j \quad \text{for} \quad j = 1, \ldots, k \tag{12}$$

$$E^j(\mathcal{F}) = E^j \quad \text{for} \quad j = 1, \ldots, k \tag{13}$$

Since we do not only have the conservation of the energy, but the energy per component, the results of the last section can be used easily in this situation. We have:

**Theorem 5** *Given an ensemble $\mathcal{E} = ((m_i, v_i)_{i=1,\ldots,n})$ with masses $m_i > 0$ and velocities $v_i = (v_i^1, \ldots, v_i^k) \in \mathbf{R}^k$ we set $M := M(\mathcal{E})$, $I^j := I^j(\mathcal{E})$ and $E^j := E^j(\mathcal{E})$ for $j = 1, \ldots, k$. By $\lambda \in (0,1)$ we have the masses $\mu_1 = \lambda M$ and $\mu_2 = (1 - \lambda)M$. For a given $\pi = (\pi^1, \ldots, \pi^k) \in \{-1, +1\}^k$ we define the velocities $\nu_1 = (\nu_1^1, \ldots, \nu_1^k)$ and $\nu_2 = (\nu_2^1, \ldots, \nu_2^k)$ by*

$$\nu_1^j = \frac{I^j}{M} + \pi^j \sqrt{\frac{1 - \lambda}{\lambda} \left( \frac{E^j}{M} - \frac{(I^j)^2}{M^2} \right)}$$

$$\nu_2^j = \frac{I^j}{M} - \pi^j \sqrt{\frac{\lambda}{1 - \lambda} \left( \frac{E^j}{M} - \frac{(I^j)^2}{M^2} \right)}$$

*for $j = 1, \ldots, k$ and the ensemble $\mathcal{F}_\pi = ((\mu_1, \nu_1), (\mu_2, \nu_2))$. Then:*

*(i) For all $\pi \in \{-1, +1\}^k$ the ensemble $\mathcal{F}_\pi$ solves the equations (11)–(13).*

*(ii) For $\pi_1, \pi_2 \in \{-1, +1\}^k$ with $\pi_1 \neq \pi_2$ it is true that $\mathcal{F}_{\pi_1} = \mathcal{F}_{\pi_2}$ if and only if $\lambda = 1/2$ and $\pi_1 = -\pi_2$, or if $v_1^i = v_2^i = \ldots = v_n^i$ for all $i \in \{1, \ldots, k\}$ with $\pi_1^i = \pi_2^i$.*

*(iii) The ensembles $\mathcal{F}_\pi$ given by $\pi \in \{-1, +1\}^k$ are the only ensembles with the given masses which fulfill (11)–(13).*

Later on it can be seen, that the choice of $\pi \in \{-1, +1\}^k$ does not influence the estimation of discrepancy. But from a practical point of view it is very important not to

6

use the same $\pi$ everytime. Numerical experiments show that it is a good choice selecting the $\pi$ randomly.

For the case $\lambda = 1/2$ we introduce the following definition:

**Definition 6**    Let $\mathcal{E} = ((\tau_i)_{i=1,\ldots,n})$ be an ensemble of particles $\tau_i = (m_i, v_i)$ whereby $m_i > 0$ and $v_i = (v_i^1, \ldots, v_i^k) \in \mathbf{R}^k$. For $\pi \in \{-1, +1\}^k$ and $\lambda = 1/2$ we have $\mathcal{F} = ((\mu, \nu_1), (\mu, \nu_2))$ as in Theorem 5. For $j = 1, \ldots, k$ we set

$$a^j := \min\left\{v_1^j, \ldots, v_n^j, \nu_1^j, \nu_2^j\right\}$$
$$b^j := \max\left\{v_1^j, \ldots, v_n^j, \nu_1^j, \nu_2^j\right\}$$

and define

$$U(\mathcal{E}) = U(\tau_1, \ldots, \tau_n) := [a^1, b^1] \times [a^2, b^2] \times \ldots \times [a^k, b^k]$$

If there is no possibility of misunderstanding we write sometimes $U(v_1, \ldots, v_n)$ instead of $U(\tau_1, \ldots, \tau_n)$.

Of course we have a natural generalization of Theorem 4:

**Theorem 7**    *Let $\mathcal{E} = ((\tau_i)_{i=1,\ldots,n})$ be an ensemble consisting of particles $\tau_i = (m_i, v_i)$ with masses $m_i > 0$ and velocities $v_i \in \mathbf{R}^k$, and $M = \sum_{i=1}^n m_i$. For a given $\lambda = 1/2$ and an arbitrary $\pi \in \{-1, +1\}^k$ let the ensemble $\mathcal{F}_\pi = ((\mu, \nu_1), (\mu, \nu_2))$ be calculated as in Theorem 5. Let $U(\mathcal{E}) = [a^1, b^1] \times \ldots \times [a^k, b^k]$. Then we have for all $j \in \{1, \ldots, k\}$ and $w \in [a^j, b^j]$:*

$$\left| \sum_{v_i^j \leq w} m_i - \#\{v_i \,|\, v_i \leq w\}\,\mu \right| \leq \frac{1}{2} M = \mu$$

**Proof:**    The proof follows immediately from Theorem 4 applied to each component.

## 3.3    On the Positions of the New Velocities

An important aspect in the later estimation of the discrepancy in the **SPLIPA**–procedure is that the velocities of the new particles are close to those of the old ones. This problem is the main topic of this section.

**Definition 8**    Consider an ensemble $\mathcal{E} = (\tau_1, \ldots, \tau_n)$ with particles $\tau_i = (m_i, v_i)$.

(i) If $v_i \in \mathbf{R}$, set $a := \min\{v_1, \ldots, v_n\}$ and $b := \max\{v_1, \ldots, v_n\}$ and define

$$U_0(\mathcal{E}) = U_0(v_1, \ldots, v_n) := [a, b]$$

(ii) If $v_i \in \mathbf{R}^k$, set $a^j := \min\{v_1^j, \ldots, v_n^j\}$ and $b^j := \max\{v_1^j, \ldots, v_n^j\}$ for $j \in \{1, \ldots, k\}$ and define

$$U_0(\mathcal{E}) = U_0(v_1, \ldots, v_n) := [a^1, b^1] \times \ldots \times [a^k, b^k]$$

Therefore the above stated problem turns into the question "how large" is $U_0(\mathcal{E})$ compared with $U(\mathcal{E})$.

As a preparation for the proof of Theorem 10 we give the following lemma:

7

**Lemma 9**   *For $n \geq 3$ let $\alpha_1 := 0$, $\alpha_n := 1$ and define the function*

$$f_n : \underbrace{[0,1] \times \ldots \times [0,1]}_{n-2} \longrightarrow I\!\!R$$

*by*

$$f_n(\alpha_2, \ldots, \alpha_{n-1}) := \frac{1}{n} \sum_{i=1}^{n} \alpha_i + \sqrt{\frac{1}{n} \sum_{i=1}^{n} \left( \alpha_i - \frac{1}{n} \sum_{j=1}^{n} \alpha_j \right)^2}$$

*Then we have for all $\alpha_2, \ldots, \alpha_{n-1} \in [0,1]$ the estimation*

$$f_n(\alpha_2, \ldots \alpha_{n-1}) \leq \frac{1}{7}\left(6 + \sqrt{6}\right)$$

**Proof:**   Proving this theorem we first search for the maximum of $f_n$ for a given $n$. By the use of the partial derivatives $\frac{\partial f_n}{\partial \alpha_i}$, it can be seen easily that $f_n$ is monotonically increasing in each component. Hence, the maximum is

$$f_n^{max} = f_n(1, \ldots, 1) = \frac{n-1}{n} + \frac{1}{n}\sqrt{n-1}$$

To find an upper bound for $f_n^{max}$ with respect to $n$ we consider

$$f^{max}(x) := \frac{x-1}{x} + \frac{1}{x}\sqrt{x-1} \quad \text{for } x \in [1, \infty)$$

As this function attains its maximum at $x_{max} = 4 + 2\sqrt{2} \approx 6.8$, we have

$$f_n^{max} \leq \max\{f_6^{max}, f_7^{max}\} = \frac{1}{7}(6 + \sqrt{6})$$

which completes the proof.   ∎

With the use of this theorem we are able to give an estimation of the size of $U(\mathcal{E})$ with respect to $U_0(\mathcal{E})$. Here we are considering only the ensemble $\mathcal{E}$ consisting of particles with the same mass $m$. This is no restriction, since in the algorithm all particles have masses of integer values. It can easily be seen, that such an ensemble can be replaced by one which all particles have the same mass in (just take the greatest common divisor of the masses).

**Theorem 10**   Let $\mathcal{E} = ((m, v_i)_{i=1,\ldots,n})$ $(n \geq 3)$ be an ensemble and set $\kappa = \frac{1}{7}(6 + \sqrt{6}) - 1 \approx 0.207$.

(i) *If $v_i \in I\!\!R$ let $U(\mathcal{E}) = [a, b]$ and $U_0(\mathcal{E}) = [c, d]$. Then*

$$b \leq d + \kappa(d - c)$$
$$a \geq c - \kappa(d - c)$$

(ii) *For $v_i = (v_i^1, \ldots, v_i^k) \in I\!\!R^k$, $k \geq 2$, one has with $U(\mathcal{E}) = [a^1, b^1] \times \ldots \times [a^k, b^k]$ and $U_0(\mathcal{E}) = [c^1, d^1] \times \ldots \times [c^k, d^k]$ that for all $j \in \{1, \ldots, k\}$*

$$b^j \leq d^j + \kappa(d^j - c^j)$$
$$a^j \geq c^j - \kappa(d^j - c^j)$$

8

**Proof:** It is sufficient to prove (i), since (ii) follows by application of (i) to each component. Because $U_0(\mathcal{E})$ is the smallest closed interval containing all $v_i$ $(i = 1, \ldots, n)$, there exist $i, j \in \{1, \ldots, n\}$ so that $v_i = c$ und $v_j = d$. Without loss of generality let $v_1 = c$ and $v_n = d$. Then

$$v_i = c + \alpha_i(d - c) \quad \text{for } \alpha_i \in [0, 1]$$

where $\alpha_1 = 0$ and $\alpha_n = 1$. So the proof follows by straightforeward calculations from Lemma 9. ∎

# 4   THE ESTIMATION OF DISCREPANCY

In this section we restrict ourselves to particles having three dimensional velocities (the physical relevant case). Furthermore we assume that the masses of the particles are of integer values (for the reasons described in the introduction). The prescribed masses $m^*$ which are determined by the procedure **MASSHA** have the representation $m^* = 2^j$ for a $j \in \mathbb{N}$.

Now consider the following situation: Given $\mathcal{E} = ((m_i, v_i)_{i=1,\ldots,n})$ consisting of particles within a cell. Since spatial homogeneity is assumed, we do not care about the positions of the particles. The procedure **MASSHA** has determined a prescribed mass $m^*$ for the particles. Therefore we are looking for an ensemble $\mathcal{F}$ consisting of particles with the mass $m^*$ so that $M(\mathcal{E}) = M(\mathcal{F})$, $I^j(\mathcal{E}) = I^j(\mathcal{F})$ and $E^j(\mathcal{E}) = E^j(\mathcal{F})$ for $j = 1, 2, 3$ and the discrepancy $D(\Phi_\mathcal{E}, \Phi_\mathcal{F})$ is "small".

## 4.1   Further Definitions

Let $x = (x^1, x^2, x^3) \in \mathbb{R}^3$. For the interval $R(x)$ we give the following

**Definition 11**   Let $x = (x^1, x^2, x^3) \in \mathbb{R}^3$. For $j \in \{1, 2, 3\}$ set

$$R_j(x) := \left\{ y = (y^1, y^2, y^3) \in \mathbb{R}^3 \, \middle| \, y \leq x \text{ und } y^j = x^j \right\}$$

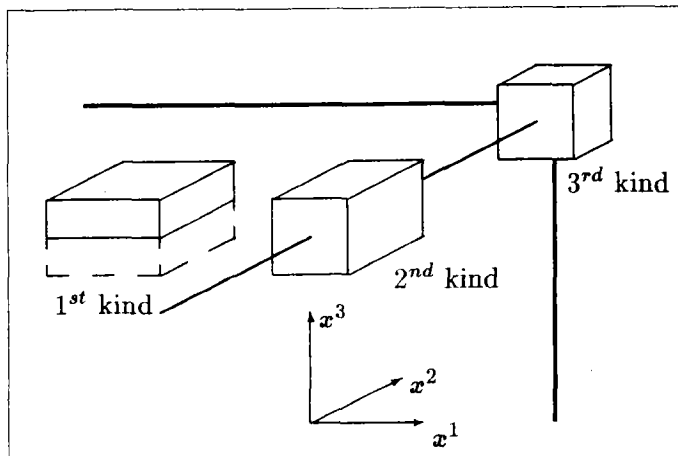**Definition 12**   Let $U = [a^1, b^1] \times [a^2, b^2] \times [a^3, b^3]$ be an interval, $x \in \mathbb{R}^3$.

(i) If $U \subset R(x)$, we call $U$ an *interval of 0. kind with respect to* $x$.

(ii) If $U \not\subset R(x)$, but $U \cap R(x) \neq \emptyset$, let

$$j := \# \left\{ k \in \{1, 2, 3\} \, | \, U \cap R_k(x) \neq \emptyset \right\}$$

and we call $U$ an *interval of j-th kind with respect to* $x$.

Figure 1 illustrates this definition.

**FIGURE 1**
Intervals of Different Kind for $R(x)$



**Definition 13**  Let $U_1, \ldots, U_n \subset \mathbb{R}^3$ be closed intervals and set $\mathcal{U} = \{U_1, \ldots, U_n\}$. For $x \in \mathbb{R}^3$ and $j \in \{0, 1, 2, 3\}$ define

$$\mathcal{U}_j(x) := \{U_i \in \mathcal{U} \,|\, U_i \text{ is interval of } j\text{-th kind w. r. t. } x\}$$

and

$$\mathcal{U}_C(x) := \mathcal{U} \setminus (\mathcal{U}_0(x) \cup \mathcal{U}_1(x) \cup \mathcal{U}_2(x) \cup \mathcal{U}_3(x))$$

Of course it holds for all $x \in \mathbb{R}^3$

$$\mathcal{U} = \mathcal{U}_C(x) \cup \mathcal{U}_0(x) \cup \mathcal{U}_1(x) \cup \mathcal{U}_2(x) \cup \mathcal{U}_3(x)$$

and this union is disjoint.

## 4.2  The Use of Splitting and Pasting

Given an ensemble $\mathcal{E} = ((m_i, v_i)_{i=1,\ldots,n})$ and a prescribed mass $m^*$ we search for an ensemble $\mathcal{F}$ which (nearly) all particles have the mass $m^*$ in. Furthermore we have the conservation equalities

$$
\begin{aligned}
M(\mathcal{E}) &= M(\mathcal{F}) & &(14)\\
I^j(\mathcal{E}) &= I^j(\mathcal{F}) & \text{for} \quad j = 1, \ldots, 3 & &(15)\\
E^j(\mathcal{E}) &= E^j(\mathcal{F}) & \text{for} \quad j = 1, \ldots, 3 & &(16)
\end{aligned}
$$

and desire the discrepancy $D(\Phi_\mathcal{E}, \Phi_\mathcal{F})$ to be "small". For this purpose we make use of the techniques described in section 3.

For $i = 1, \ldots n$ let $n_i^*, m_i' \in \mathbb{N}$ such that

$$m_i = n_i^* \cdot m^* + m_i' \quad \text{with} \quad 0 \le m_i' < m^* \tag{17}$$

and set
$$\mathcal{E}_i := (\underbrace{(m^*, v_i), \ldots, (m^*, v_i)}_{n_i^*}, (m_i', v_i)) \quad \text{if} \quad m_i' > 0$$

respectively
$$\mathcal{E}_i := (\underbrace{(m^*, v_i), \ldots, (m^*, v_i)}_{n_i^*}) \quad \text{if} \quad m_i' = 0$$

Hence,
$$\mathcal{E}_i \sim ((m_i, v_i))$$

Putting all particles of $\mathcal{E}_1, \ldots \mathcal{E}_n$ with mass $m^*$ together, one obtains
$$\mathcal{E}^* := (\underbrace{(m^*, v_1), \ldots, (m^*, v_1)}_{n_1^*}, \ldots, \underbrace{(m^*, v_n), \ldots, (m^*, v_n)}_{n_n^*})$$

where $n^* := \sum_{i=1}^n n_i^*$ is the number of particles in $\mathcal{E}^*$. The other particles of $\mathcal{E}_1, \ldots \mathcal{E}_n$, form the ensemble
$$\mathcal{E}' := ((m_i', v_i)_{i \in N'})$$

where $N' := \{i \mid m_i' > 0\}$. Hence,
$$\mathcal{E} \sim \mathcal{E}^* \cup \mathcal{E}'$$

Since the particles of $\mathcal{E}^*$ have already the desired mass $m^*$, we have to consider $\mathcal{E}'$ consisting of particles with $m_i' < m^*$. Since all masses are of integer values, it makes sense to define
$$m'' := \gcd\{m_i' \mid i \in N'\}$$

Then there exist $n_i'' \in I\!\!N$ so that $m_i' = n_i'' \cdot m''$. Therefore we are able to define
$$\mathcal{E}_i'' := (\underbrace{(m'', v_i'), \ldots, (m'', v_i')}_{n_i''})$$

which fulfills $((m_i', v_i)) \sim \mathcal{E}_i''$. If we set $\mathcal{E}'' := \bigcup_{i \in N'} \mathcal{E}_i''$ we have $\mathcal{E}'' \sim \mathcal{E}'$ and the number of particles in $\mathcal{E}''$ is $n'' := \sum_{i \in N'} n_i''$. Putting things together we have
$$\mathcal{E} \sim \mathcal{E}^* \cup \mathcal{E}''$$

where $\mathcal{E}^*$ consists of particles of mass $m^*$ and $\mathcal{E}''$ of particles of mass $m''$. Since $m''$ is a divisor of $m^*$ we can define $p := 2m^*/m''$ and therefore have
$$n'' = l \cdot p + r \quad \text{with} \quad l, r \in I\!\!N, r < p$$

Assume the ensemble $\mathcal{E}''$ is divided in an ensemble $\mathcal{E}'''$ with $l \cdot p$ particles and in $\mathcal{E}'''_{left}$ with $r$ particles. Then we have

$$\mathcal{E}'' = \mathcal{E}''' \cup \mathcal{E}'''_{left} \tag{18}$$
$$M(\mathcal{E}''') = l \cdot p \cdot m'' = 2m^* \cdot l \tag{19}$$
$$M(\mathcal{E}'''_{left}) = r \cdot m'' < p \cdot m'' = 2m^* \tag{20}$$

Now assume we have a partition

$$\mathcal{E}''' = \bigcup_{i=1}^{l} \mathcal{E}_i'''$$

each $\mathcal{E}_i'''$ consisting of $p$ particles. This partition is essential for the later estimation of the discrepancy and we will give an example how to choose it well in the next section.

By application of Theorem 5 for each sub–ensemble $\mathcal{E}_i'''$ with $\lambda = 1/2$, we get ensembles $\mathcal{F}_i'''$. Setting $\mathcal{F}''' := \bigcup_{i=1}^{l} \mathcal{F}_i'''$ the ensemble $\mathcal{F}''' \cup \mathcal{E}^*$ consists of particles of the desired mass $m^*$. With

$$\mathcal{F} := \mathcal{F}''' \cup \mathcal{E}^* \cup \mathcal{E}_{left}'''$$

the ensemble $\mathcal{F}$ fulfills the equations (14)–(16) and it yields

$$M\left(\mathcal{E}_{left}'''\right) < 2m^*$$

Hence, the total mass of those particles not having the mass $m^*$ is bounded by $2m^*$.

## 4.3 The Estimation of Discrepancy

We start with

**Definition 14**    Let the ensemble $\mathcal{E} = ((m, v_i)_{i=1,\ldots,l\cdot p})$ with $l \cdot p$ particles be divided in

$$\mathcal{E} = \bigcup_{i=1}^{l} \mathcal{E}_i$$

where each $\mathcal{E}_i$ consists of $p$ particles. For $i = 1,\ldots,l$ let $U_i = U(\mathcal{E}_i)$ and $\mathcal{U} = \{U_i \,|\, i \in \{1,\ldots,l\}\}$. Let $q_1$, $q_2$, $q_3 \in \mathbb{N}$ be given such that for all $x \in \mathbb{R}^3$ and all $j \in \{1,2,3\}$ it holds that

$$\#\left\{U_i \in \mathcal{U} \,|\, U_i \text{ is an interval of } j\text{–th kind w. r. t. } x\right\} \leq q_j$$

Then we say the ensemble $\mathcal{E}$ with partition $\mathcal{E} = \bigcup_{i=1}^{l} \mathcal{E}_i$ has the *property (V) with* $(q_1, q_2, q_3)$.

**Theorem 15**    *For a given ensemble $\mathcal{E} = ((m_i, v_i)_{i=1,\ldots,n})$ with $v_i \in \mathbb{R}^3$, $m_i > 0$ and a prescribed mass $m^*$ construct the equivalent ensemble*

$$\mathcal{E}^* \cup \mathcal{E}''' \cup \mathcal{E}_{left}$$

*as described above. Assume $\mathcal{E}'''$ has a partition*

$$\mathcal{E}''' = \bigcup_{i=1}^{l} \mathcal{E}_i'''$$

*with property (V) with $(q_1, q_2, q_3)$. Putting*

$$\mathcal{F} = \mathcal{E}^* \cup \bigcup_{i=1}^{l} \mathcal{F}_i''' \cup \mathcal{E}_{left}$$

*it holds that*

*(i)*

$$M(\mathcal{F}) = M(\mathcal{E})$$
$$I^j(\mathcal{F}) = I^j(\mathcal{E}) \quad for \quad j = 1, 2, 3$$
$$E^j(\mathcal{F}) = E^j(\mathcal{E}) \quad for \quad j = 1, 2, 3$$

*(ii)*
$$M(\mathcal{E}_{left}) < 2m^*$$

*(iii)*

$$D(\Phi_\mathcal{E}, \Phi_\mathcal{F}) \le (q_1 + 2q_2 + 2q_3)m^*$$

**Proof:** The statements (i) und (ii) are evident by construction. To prove (iii) we have by (1) and (2) with $\mathcal{F}''' := \bigcup_{i=1}^l \mathcal{F}_i'''$

$$D(\Phi_\mathcal{E}, \Phi_\mathcal{F}) = D(\Phi_{\mathcal{E}'''}, \Phi_{\mathcal{F}'''})$$

Assume $\mathcal{F}_i'''$ are given by $\mathcal{F}_i''' = ((m^*, \nu_{i1}), (m^*, \nu_{i2}))$ then the local discrepancy is

$$d_{\Phi_{\mathcal{F}'''}, \Phi_{\mathcal{E}'''}}(w) = m^* \cdot \# \{\nu_{ij} \,|\, \nu_{ij} \in R(w)\} - m'' \cdot \# \{v_{ij} \,|\, v_{ij} \in R(w)\}$$

By the definition of $U_i = U(\mathcal{E}_i''')$ we have

$$d_{\Phi_{\mathcal{F}'''}, \Phi_{\mathcal{E}'''}}(w) = \sum_{i=1}^l (m^* \cdot \# \{\nu_{ij} \,|\, \nu_{ij} \in R(w) \cap U_i\} -$$
$$m'' \cdot \# \{v_{ij} \,|\, v_{ij} \in R(w) \cap U_i\})$$
$$= \sum_{U_i \in \mathcal{U}_0(w)} + \sum_{U_i \in \mathcal{U}_1(w)} + \sum_{U_i \in \mathcal{U}_2(w)} + \sum_{U_i \in \mathcal{U}_3(w)} + \sum_{U_i \in \mathcal{U}_C(w)}$$
$$(m^* \cdot \# \{\nu_{ij} \,|\, \nu_{ij} \in R(w) \cap U_i\} - m'' \cdot \# \{v_{ij} \,|\, v_{ij} \in R(w) \cap U_i\})$$

It is clear that the first sum is equal to zero. The $U_i$ in the second sum are of $1^{st}$ kind. Therefore we have (Theorem 7)

$$\left| m^* \cdot \# \{\nu_{ij} \,|\, \nu_{ij} \in R(w) \cap U_i\} - m'' \cdot \# \{v_{ij} \,|\, v_{ij} \in R(w) \cap U_i\} \right| \le \frac{1}{2}M = m^*$$

Looking at the $U_i$ in the third and fourth sum the trivial estimate holds

$$\left| m^* \cdot \# \{\nu_{ij} \,|\, \nu_{ij} \in R(w) \cap U_i\} - m'' \cdot \# \{v_{ij} \,|\, v_{ij} \in R(w) \cap U_i\} \right| \le 2m^*$$

Since the $U_i$ of the fifth sum fulfill $U_i \cap R(w) = \emptyset$ this sum is equal to zero. Putting things together we obtain

$$d_{\Phi_{\mathcal{F}'''}, \Phi_{\mathcal{E}'''}}(w) \le 0 + \sum_{U_i \in \mathcal{U}_1(w)} m^* + \sum_{U_i \in \mathcal{U}_2(w)} 2m^* + \sum_{U_i \in \mathcal{U}_3(w)} 2m^* + 0$$
$$\le (q_1 + 2q_2 + 2q_3)\, m^*$$

which completes the proof. ∎

The following example shows how the partition $\mathcal{E}'''$ can be chosen so that the numbers $q_i$ and therefore the discrepancy $D(\Phi_\mathcal{E}, \Phi_\mathcal{F})$ are small.

13

**Example 16** Given an ensemble $\mathcal{E} = ((m, v_i)_{i=1,\ldots,l\cdot p})$ we assume that for all $j \in \{1, 2, 3\}$ the components of the velocities $v_1^j, \ldots, v_{l\cdot p}^j$ are pairwise distinct and (for simplicity) that $k := \sqrt[3]{l} \in \mathbf{N}$.

Because of the distinction of the first components of the velocities there exist

$$-\infty = w_0 < w_1 < \ldots < w_{k-1} < w_k = \infty$$

so that

$$\# \left\{ v_i \, \middle| \, w_{r-1} < v_i^1 \le w_r \right\} = k^2 p \text{ for } r = 1, \ldots, k$$

With the same arguments there exist for every $r \in \{1, \ldots, k\}$ numbers

$$-\infty = w_0^r < w_1^r < \ldots < w_{k-1}^r < w_k^r = \infty$$

so that for all $s \in \{1, \ldots, k\}$

$$\# \left\{ v_i \, \middle| \, w_{r-1} < v_i^1 \le w_r \text{ and } w_{s-1}^r < v_i^2 \le w_s^r \right\} = kp$$

Finally, we have for all $r \in \{1, \ldots, k\}$ and $s \in \{1, \ldots, k\}$

$$-\infty = w_0^{rs} < w_1^{rs} < \ldots < w_{k-1}^{rs} = \infty$$

so that for every interval

$$Q_{rst} := (w_{r-1}, w_r] \times (w_{s-1}^r, w_s^r] \times (w_{t-1}^{rs}, w_t^{rs}]$$

it is

$$\# \left\{ v_i \, \middle| \, w_{r-1} < v_i^1 \le w_r \text{ and } w_{s-1}^r < v_i^2 \le w_s^r \text{ and } w_{t-1}^{rs} < v_i^3 \le w_t^{rs} \right\} = p$$

Hence there exists a partition of $\mathbb{R}^3$ in $k^3 = l$ intervals $Q_{rst}$, $r, s, t \in \{1, \ldots, k\}$, everyone of them containing $p$ particles. Now, we construct a partition of $\mathcal{E}$ by

$$\mathcal{E}_{rst} = ((m, v_1^{rst}), \ldots, (m, v_p^{rst}))$$

where $\mathcal{E}_{rst}$ contains all particles of $\mathcal{E}$ which have velocities in $Q_{rst}$. If $U(\mathcal{E}_{rst}) \subset Q_{rst}$ is assumed (which is no strong assumption because of Theorem 10) it can easily be shown that this partition has the property (V) with $(3(k-1)^2, 3(k-1), 1)$.

It can easily be seen that in this situation the order of convergence is like

$$D(\Phi_{\mathcal{E}}, \Phi_{\mathcal{F}}) = \mathcal{O}\left(\frac{1}{\sqrt[3]{n}}\right)$$

where $n$ is the number of particles of $\mathcal{F}$.

# 5 COMPUTATIONAL EXPERIMENTS

In this section we present some numerical results gained by the new algorithm using weighted particles. But first we will give the following remark concerning algorithmical aspects of the **SPLIPA**-procedure:

As it was mentioned before, the partition of the ensembles $\mathcal{E}'''$ in sub–ensembles $\mathcal{E}'''_i$ is essential for the estimation of discrepancy. Nevertheless the strategy described in Example 16 would need too much computing time. Therefore we follow another way: Let $\mathcal{E}''' = ((m'', v_i)_{i=1,\dots,n})$ be the former ensemble and $m^*$ the prescribed mass. Then $l := \left[\frac{M(\mathcal{E}''')}{2m^*}\right]$ is the number of pairs of particles that will be created. With

$$k := \left\{ \begin{array}{ll} \sqrt[3]{l} & \text{if } \sqrt[3]{l} \in \boldsymbol{N} \\ \left[\sqrt[3]{l} + 1\right] & \text{if } \sqrt[3]{l} \notin \boldsymbol{N} \end{array} \right.$$

and with

$$\begin{aligned} a^j &:= \min_{i=1,\dots,n} \{v_i^j\} \\ b^j &:= \max_{i=1,\dots,n} \{v_i^j\} \end{aligned}$$

for $j = 1, 2, 3$ we divide the intervals $[a^j, b^j]$ by

$$x_p^j := a^j + \frac{p}{k}\left(b^j - a^j\right) \quad \text{for} \quad p = 0, \dots, k$$

and therefore the interval $[a^1, b^1] \times [a^2, b^2] \times [a^3, b^3]$ in $k^3$ sub intervals. The order of the particles is choosen by going through this $k^3$ sub intervals, so that succeding sub intervals have a common side. An illustration can be found in Figure 2.

### FIGURE 2
### A possible Way through the Sub Intervals



For the numerical experiments we consider the following two dimensional streaming problem: A rectangle with edges $l^1 = 80\,\text{m}$ and $l^2 = 60.8\,\text{m}$ is divided into quadratic

cells with $l_{\text{cell}} = 0.8\,\text{m}$. Hence, there are $100 \times 76$ cells. At the center of this rectangle is an ellipse with semi–major axis $16\,\text{m}$ and semi–minor axis $3.2\,\text{m}$ with the angle of attack $\alpha = 90°$.

We consider the mono atomic gas argon, with

gas constant: $\quad R \;=\; 208\frac{\text{J}}{\text{kg}\cdot\text{K}}$

molecule mass: $\quad m \;=\; 6.6378 \cdot 10^{-26}\,\text{kg}$

molecule diameter: $\quad d \;=\; 3.66 \cdot 10^{-10}\,\text{m}$

Far away from the ellipse we assume a thermodynamical equilibrium with

velocity in $x^1$–direction: $\quad v^1_\infty \;=\; -4983.17\,\frac{\text{m}}{\text{s}}$

velocity in $x^2$–direction: $\quad v^2_\infty \;=\; 0$

temperature: $\quad T_\infty \;=\; 279.8\,\text{K}$

mean free path: $\quad \lambda_\infty \;=\; 0.8\,\text{m}$

On the surface of the ellipse there is a temperature $T_0 = 1442.3\,\text{K}$. The boundary conditions are simulated by diffuse reflexion (cf. [5]).

In each run there are calculated 500 time steps of size $\Delta t = 1.6 \cdot 10^{-4}\,\text{s}$. After 250 time steps we start averaging the moments.

To compare the accuracy and computing time we did four runs: the original program, where in the beginning were 25 respectively 64 particles in each cell. These runs are called A25 resp. A64. We present two runs with the new algorithm using weighted particles: In B25–3 there are three allowed particle masses, namely the initial mass $m_{\text{init}}$, $1/2\,m_{\text{init}}$ and $1/4\,m_{\text{init}}$. In B25–4 there are four allowed particle masses namely the initial mass $m_{\text{init}}$, $2m_{\text{init}}$, $1/2\,m_{\text{init}}$ and $1/4\,m_{\text{init}}$ (cf. Table 1).

**TABLE 1**

|  | particles per cell | used particle masses |
|---|---|---|
| **A64** | 64 | $m_{\text{init}}$ |
| **A25** | 25 | $m_{\text{init}}$ |
| **B25–3** | 25 | $1/4\,m_{\text{init}}$ $\quad 1/2\,m_{\text{init}}$ $\quad m_{\text{init}}$ |
| **B25–4** | 25 | $1/4\,m_{\text{init}}$ $\quad 1/2\,m_{\text{init}}$ $\quad m_{\text{init}}$ $\quad 2\,m_{\text{init}}$ |

The number of particles in the equilibrium state and the consumed computing time can be seen in Table 2.

**TABLE 2**

|  | CPU time: | number of particles in the equilibrium state |
|---|---|---|
| **A64** | 44'41" | 706000 |
| **A25** | 24'52" | 275000 |
| **B25–3** | 29'13" | 334000 |
| **B25–4** | 26'17" | 248000 |

The appendix presents results for the density, the Mach number and the temperature in the rows 10, 20 and 38, counted from bottom to top. Hereby the following can be seen:

- In the density and Mach number there is a good agreement of A25, B25-3 and B25-4 with A64.

- In the temperature there are great differences between A25 and A64 in the rows 20 and 38 (where is a small density of the gas), where the results of B25-3 and B25-4 are really better.

- As expected, the results in regions of high density are sometimes better in B25-3 than in B25-4.

Taking all results into consideration, it can be seen that the new algorithm creates better numerical results with a comparable amount of time and memory.


## ACKNOWLEDGEMENT

## APPENDIX

**FIGURE 3**

Density in Row 10
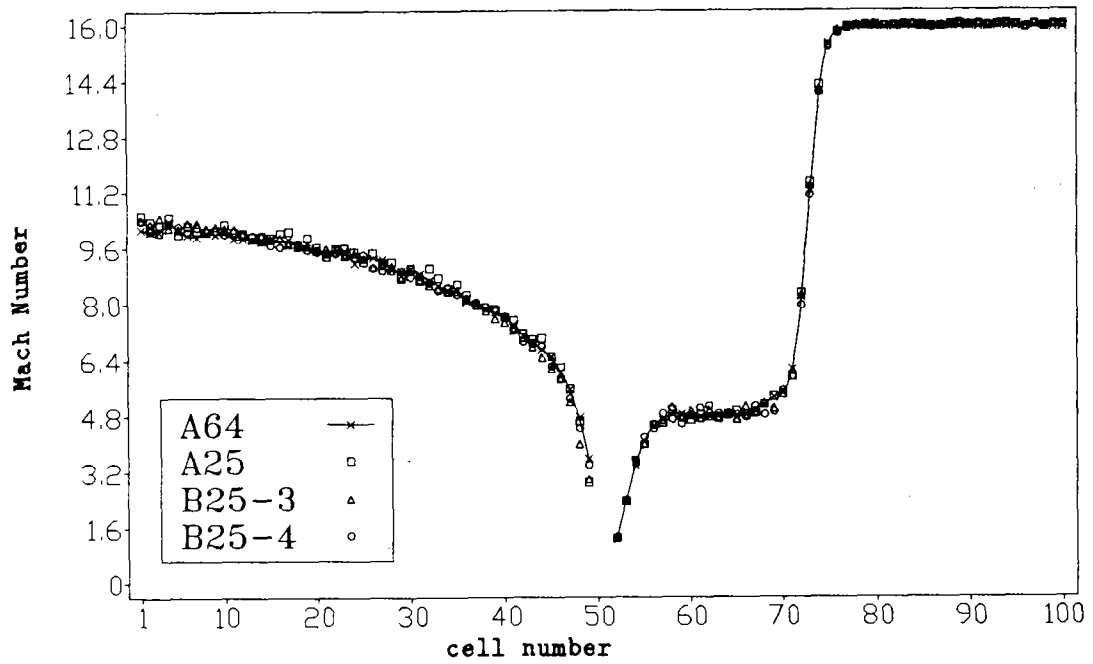
**FIGURE 4**

Density in Row 20
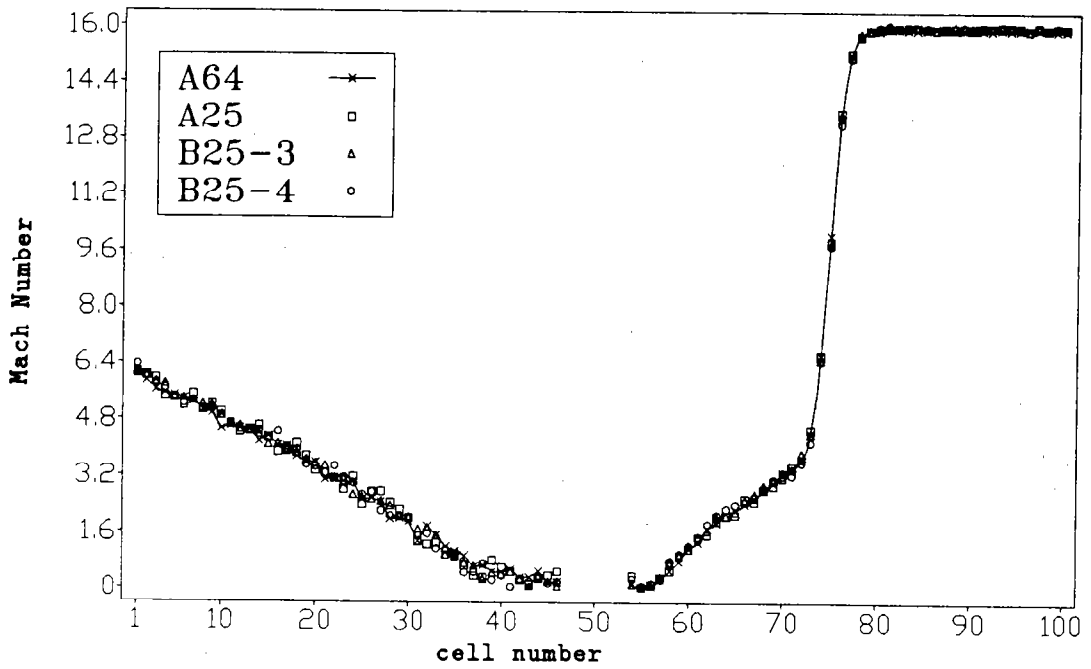


**FIGURE 5**

Density in Row 38

**FIGURE 6**

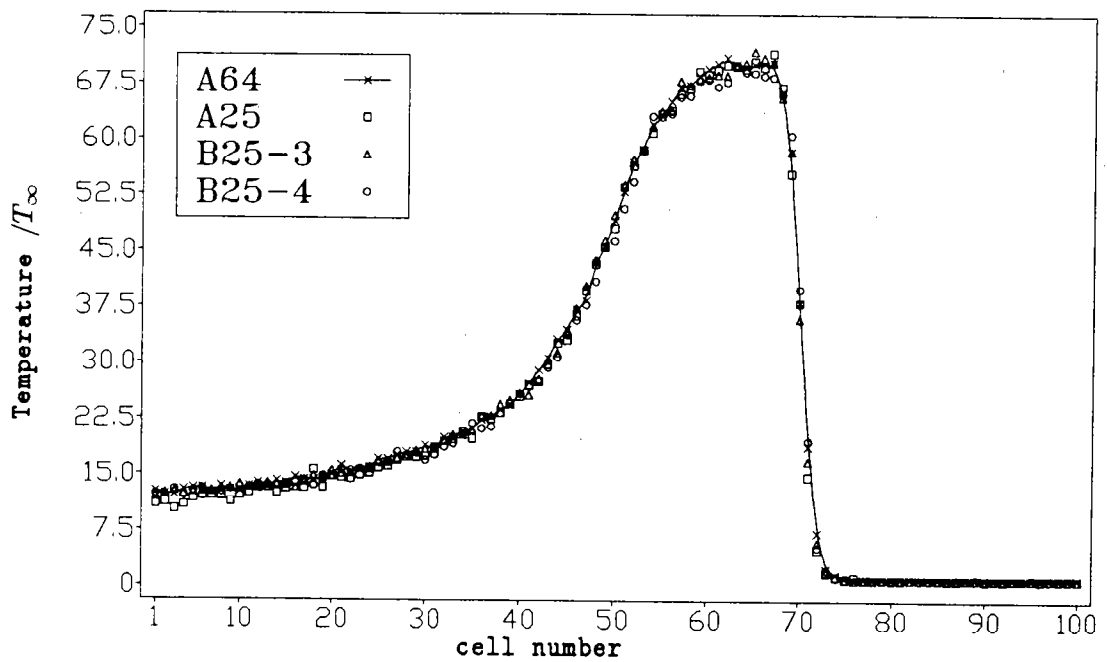Mach Number in Row 10



**FIGURE 7**

Mach Number in Row 20

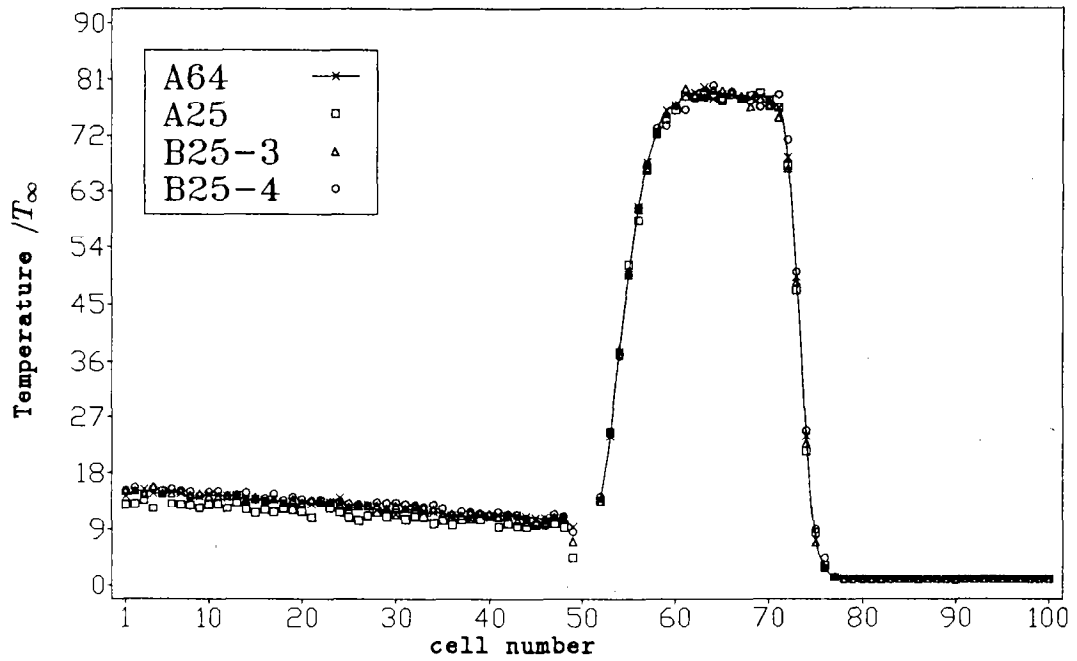**FIGURE 8**

Mach Number in Row 38
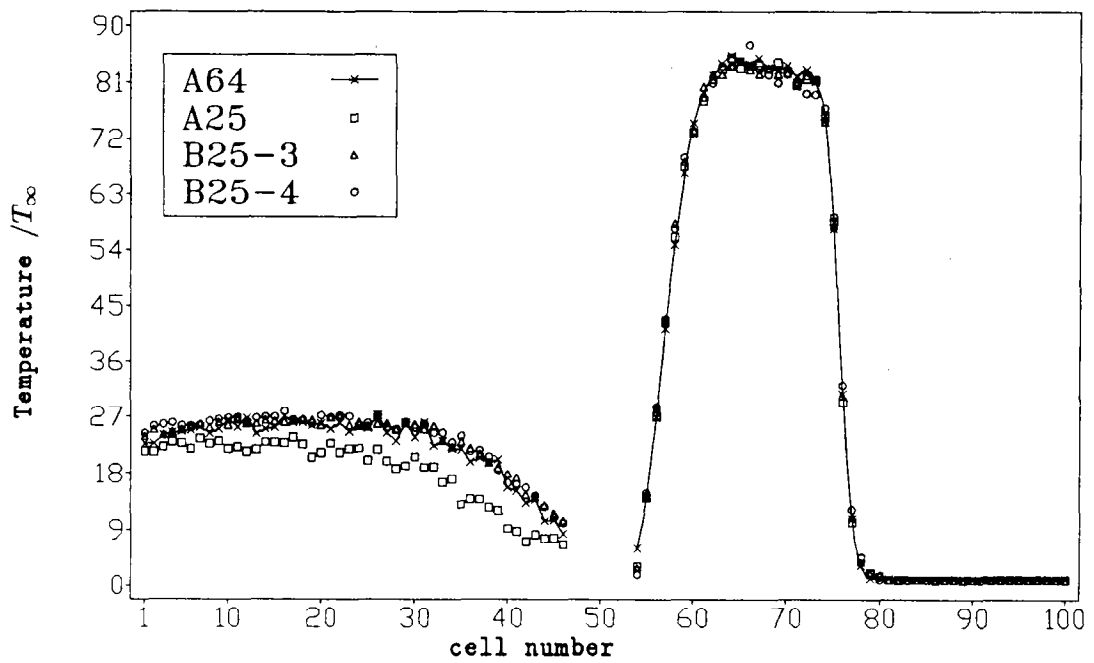


**FIGURE 9**

Temperature in Row 10

**FIGURE 10**

Temperature in Row 20



**FIGURE 11**

Temperature in Row 38

# References

[1] H. Babovsky, *On a Simulation Scheme for the Boltzmann Equation,* Math. Meth. Appl. Sci., 8 (1986), pp. 223 – 233

[2] H. Babovsky, *A Convergence Proof for Nanbu's Boltzmann Simulation Scheme,* Eur. J. Mech., B/Fluids, 8, No. 1 (1989), pp. 41 – 55

[3] H. Babovsky and R. Illner, *A convergence Proof for Nanbu's Simulation Method for the full Boltzmann Equation,* SIAM J. Numer. Anal., Vol. 26, No. 1 (1989), pp. 45 – 65

[4] G. A. Bird, *Molecular Gas Dynamics,* Clarendon Press, Oxford, 1976

[5] C. Cercignani, *Theory and Application of the Boltzmann Equation,* Scot. Acad. Press, Edinburgh, 1975

[6] H. Neunzert, F. Gropengiesser, J. Struckmeier, *Computational Methods for the Boltzmann Equation,* in *The State of the Art in Applied and Industrial Mathematics,* R. Spigler Ed, Kluwer Academic Publishers, Dordrecht, 1990

[7] R. Illner and H. Neunzert, *On Simulation Methods for the Boltzmann Equation,* Transport Theory Statist. Phys., 16 (1987), pp. 141–154

[8] L. Kuipers and H. Niederreiter, *Uniform Distribution of Sequences,* John Wiley & Sons, New York, 1974

[9] M. Schreiner, *Gewichtete Teilchen in der Methode der finiten Punktmenge,* Diplomarbeit, Universität Kaiserslautern, 1991

[10] H. Weyl, *Über die Gleichverteilung von Zahlen mod 1,* Math. Ann. 77 (1916)