

# High-performance direct gravitational $N$ -body simulations on graphics processing units

Simon F. Portegies Zwart<sup>a,b,\*</sup>, Robert G. Belleman<sup>a</sup>, Peter M. Geldof<sup>a</sup>

<sup>a</sup> Section Computational Science, University of Amsterdam, Amsterdam, The Netherlands

<sup>b</sup> Astronomical Institute “Anton Pannekoek”, University of Amsterdam, Amsterdam, The Netherlands

Received 2 February 2007; received in revised form 24 April 2007; accepted 29 May 2007

Available online 5 June 2007

Communicated by M. van der Klis

---

## Abstract

We present the results of gravitational direct  $N$ -body simulations using the commercial graphics processing units (GPU) NVIDIA Quadro FX1400 and GeForce 8800GTX, and compare the results with GRAPE-6Af special purpose hardware. The force evaluation of the  $N$ -body problem was implemented in Cg using the GPU directly to speed-up the calculations. The integration of the equations of motions were, running on the host computer, implemented in C using the 4th order predictor–corrector Hermite integrator with block time steps. We find that for a large number of particles ( $N \gtrsim 10^4$ ) modern graphics processing units offer an attractive low cost alternative to GRAPE special purpose hardware. A modern GPU continues to give a relatively flat scaling with the number of particles, comparable to that of the GRAPE. The GRAPE is designed to reach double precision, whereas the GPU is intrinsically single-precision. For relatively large time steps, the total energy of the  $N$ -body system was conserved better than to one in  $10^6$  on the GPU, which is impressive given the single-precision nature of the GPU. For the same time steps, the GRAPE gave somewhat more accurate results, by about an order of magnitude. However, smaller time steps allowed more energy accuracy on the grape, around  $10^{-11}$ , whereas for the GPU machine precision saturates around  $10^{-6}$ . For  $N \gtrsim 10^6$  the GeForce 8800GTX was about 20 times faster than the host computer. Though still about a factor of a few slower than GRAPE, modern GPUs outperform GRAPE in their low cost, long mean time between failure and the much larger onboard memory; the GRAPE-6Af holds at most 256k particles whereas the GeForce 8800GTX can hold 9 million particles in memory.

© 2007 Elsevier B.V. All rights reserved.

*PACS:* 82.20.Wt; 83.10.Rs; 87.15.Aa; 87.64.Aa; 95.75.Pq; 95.30.Ky

*Keywords:* Gravitation; Stellar dynamics; Methods:  $N$ -body simulation; Methods: Numerical

---

## 1. Introduction

Since the first large scale simulations of self gravitating systems the direct  $N$ -body method has gained a solid footing in the research community. At the moment  $N$ -body techniques are used in astronomical studies of planetary systems, debris discs, stellar clusters, galaxies all the way to simulations of the entire universe (Hut, 2007). Outside

astronomy the main areas of research which utilize the same techniques are molecular dynamics, elementary particle scattering simulations, plate tectonics, traffic simulations and chemical reaction network studies. In the latter non-astronomical applications, the main force evaluating routine is not as severe as in the gravitational  $N$ -body simulations, but the backbone simulation environments are not very different.

The main difficulty in simulating self gravitating systems is the lack of antigravity, which results in the requirement of global communication; each object feels the gravitational attraction of any other object.

---

\* Corresponding author. Tel.: +31 20 5257510; fax: +31 20 5257484.  
E-mail address: [spz@science.uva.nl](mailto:spz@science.uva.nl) (S.F. Portegies Zwart).

The first astronomical simulation of a self gravitating  $N$ -body system was carried out by Holmberg (1941) with the use of 37 light bulbs and photoelectric cells to evaluate the forces on the individual objects. Holmberg spent weeks in order to perform this quite moderate 37-particle simulation. Over the last 60 or so years many different techniques have been introduced to speed-up the kernel calculation. Today, such a calculation requires about 50,000 integration steps for one dynamical time unit. At a speed of  $\sim 10$  GFLOP/s the calculation would be performed in a few seconds.

The gravitational  $N$ -body problem has made enormous advances in the last decade due to algorithmic design. The introduction of digital computers in the arena (von Hoerner, 1963; Aarseth and Hoyle, 1964; van Albada, 1968) led to a relatively quick evaluation of mutual particle forces. Advanced integration techniques, introduced to turn the particle forces in a predicted space–time trajectory, opened the way to predictable theoretical results (Aarseth and Lecar, 1975; Aarseth, 1999). One of the major developments in the speed-up and improved accuracy of the direct  $N$ -body problem was the introduction of the block-time step algorithm (Makino, 1991; McMillan and Aarseth, 1993).

In the late 1980s it became quite clear that the advances of modern computer technology via Moore’s law (Moore, 1965) were insufficient to simulate large star clusters by the new decade (Makino and Hut, 1988; Makino and Hut, 1990). This realization brought forward the initiatives employed around the development of special hardware for evaluating the forces between the particles (Applegate et al., 1986; Taiji et al., 1996; Makino and Taiji, 1998; Makino, 2001; Makino et al., 2003), and of the efficient use of assembler code on general purpose hardware (Nitadori et al., 2006; Nitadori et al., 2007).

One method to improve performance is by parallelising force evaluation (Eq. (1)) for use on a Beowulf or cluster computer (with or without dedicated hardware) (Harfst et al., 2007), a large parallel supercomputer (Makino, 2002; Dorband et al., 2003) or for grid operations (Guaularis et al., 2007). In particular, for distributed hardware it is crucial to implement an algorithm that limits communication as much as possible, otherwise the bottleneck simply shifts from the force evaluation to interprocessor communication.

A breakthrough in direct-summation  $N$ -body simulations came in the late 1990s with the development of the GRAPE series of special-purpose computers (Makino and Taiji, 1998), which achieve spectacular speedups by implementing the entire force calculation in hardware and placing many force pipelines on a single chip. The latest special purpose computer for gravitational  $N$ -body simulations, GRAPE-6, performs at a peak speed of about 64 TFLOP/s (Makino, 2001).

In our standard setup, one GRAPE-6Af processor board is attached to a host workstation, in much the same way that a floating-point or graphics accelerator card is used. We use a smaller version: the GRAPE-6Af which has

four chips connected to a personal workstation via the PCI bus delivering a theoretical peak performance of  $\sim 131$  GFLOP/s for systems of up to 128k particles at a cost of  $\sim \$6\text{K}$  (Fukushige et al., 2005). Advancement of particle positions [ $\mathcal{O}(N)$ ] is carried out on the host computer, while interparticle forces [ $\mathcal{O}(N^2)$ ] are computed on the GRAPE.

The latest developments in this endeavour is the design and construction of the GRAPE-DR, the special purpose computer which will break the PFLOP/s barrier by the summer of 2008 (Makino, 2007).<sup>1</sup> One of the main arguments to develop such a high powered and relatively diverse computer is to perform simulations of entire galaxies (Makino, 2005a; Hoekstra et al., 2007).

The main disadvantages of these special purpose computers, however, are the relatively short mean time between failure, the limited availability, the limited applicability, the limited on-board memory to store particles, the simple fact that they are basically build by a single research team led by prof. J. Makino and the lack of competing architectures.

The gaming industry, though not deliberately supportive of scientific research, has been developing high power parallel vector processors for performing specific rendering applications, which are in particular suitable for boosting the frame-rate of games. Over the last 7 years graphics processing units (GPUs) have evolved from fixed function hardware for the support of primitive graphical operations to programmable processors that outperform conventional CPUs, in particular for vectorizable parallel operations. Regretfully, the precision of these processors is still 32-bit IEEE which is below the average general purpose processor, but for many applications it turns out that the higher (double) precision is not crucial or can be emulated at some cost. It is because of these developments, that more and more people use the GPU for wider purposes than just for graphics (Fernando, 2004; Pharr and Fernando, 2005; Buck et al., 2004). This type of programming is also called general purpose computing on graphics processing units (GPGPU)<sup>2</sup>. Earlier attempts to use a GPU for gravitational  $N$ -body simulations were carried out using approximate force evaluation methods with shared time steps (Nyland et al., 2004), but provide little improvement in performance. A 25-fold speed increase compared to an Intel Pentium IV processor was reported by Elsen et al. (2006), but details of their implementation of the force evaluation algorithm are yet unclear. Recently, Hamada and Iitaka (2007) proposed the ‘Chamomile’ scheme for running  $N$ -body simulations with a shared time-step algorithm on GPUs. Though, their method, using the CUDA programming environment, outperforms our implementations, the shared time step renders their code unpractical for simulating dense star clusters.

Using GPUs as a general purpose vector processor works as follows. Colours in computer graphics are represented by one or more numbers. The luminance can be represented by just a single number, whereas a coloured pixel

<sup>1</sup> See <http://grape.astron.s.u-tokyo.ac.jp/grape/computer/grape-dr.html>

<sup>2</sup> See <http://www.gpgpu.org>

may contain separate values indicating the amount of red, green and blue. A fifth value alpha may be included to indicate the amount of transparency. Using this information, a pixel can be drawn. For general purpose computing, the colour information of a pixel is used to represent attributes of the computation. There are many pixels in a frame, and ideally, these should be updated all at the same time and at a rate exceeding the response time of the human eye. This requires fast computations for updating the pixels when, for example, a camera moves or a new object comes into view. Such operations usually have an impact on many or even all pixels and therefore fast computations are required. As the majority of pixels do not require information from other pixels, processing can be done efficiently in parallel. All information required to build a pixel should go through a series of similar operations, a technique which is better known as single instruction, multiple data (SIMD). There are many different kinds of operations this information needs to go through. The stream programming model has been designed to make the information go through these operations efficiently, while exposing as much parallelism as possible. The stream programming model views all informations as “streams” of ordered data of the same data type. The streams pass through “kernels” that operate on the streams and produce one or more streams as output.

In this paper we report on our endeavour to convert a high precision production quality  $N$ -body code to operate with graphics processing units. In Section 2, we explain the adopted  $N$ -body integration algorithm, in Section 3, we address the programming environment we used to program the GPU. In Sections 4 and 5, we present the results on two GPUs and compare them with GRAPE-6Af and we discuss a model to explain the GPU’s performance. In Section 6, we summarize our findings, and in [Supplementary material](#) we present a snippet of the source code in Cg.

## 2. Calculating the force and integrating the particles

The gravitational evolution of a system consisting of  $N$  stars with masses  $m_j$  and at position  $\mathbf{r}_j$  is computed by the direct summation of the Newtonian force between each of the  $N$  stars. The force  $\mathbf{F}_i$  acting on particle  $i$  is then obtained by summation of all other  $N - 1$  particles

$$\mathbf{F}_i \equiv m_i \mathbf{a}_i = m_i G \sum_{j=0, j \neq i}^{N-1} m_j \frac{\mathbf{r}_i - \mathbf{r}_j}{|\mathbf{r}_i - \mathbf{r}_j|^3}. \quad (1)$$

Here  $G$  is the Newton constant. For further readability we omit the particle index  $i$  and present vectors in boldface.

A cluster consisting of  $N$  stars evolves dynamically due to the mutual gravity of the individual stars. For an accurate force calculation on each star a total of  $\frac{1}{2}N(N - 1)$  partial forces have to be computed. This  $\mathcal{O}(N^2)$  operation is the bottleneck for the gravitational  $N$ -body problem.

The GPU scheme described in this paper is implemented in the  $N$ -body integrator. Here particle motion is calculated using a fourth-order, individual-time step “Hermite”

predictor–corrector scheme (Makino and Aarseth, 1992). This scheme works as follows. During a time step the positions ( $\mathbf{x}$ ) and velocities ( $\mathbf{v} \equiv \dot{\mathbf{x}}$ ) are first predicted to fourth order using the acceleration ( $\mathbf{a} \equiv \ddot{\mathbf{x}}$ ) and the “jerk” ( $\mathbf{k} \equiv \dot{\mathbf{a}}$ , the time derivative of the acceleration) which are known from the previous step.

At the start of each simulation an initial time step is calculated

$$dt = v \frac{\mathbf{a}}{\mathbf{k}}. \quad (2)$$

Here we introduce  $v$  as an accuracy control parameter ( $v = 0.01$  for most of our simulations, see also Eq. (9)).

The predicted position ( $\mathbf{x}_p$ ) and velocity ( $\mathbf{v}_p$ ) are calculated for all particles

$$\mathbf{x}_p = \mathbf{x} + \mathbf{v} dt + \frac{1}{2} \mathbf{a} dt^2 + \frac{1}{6} \mathbf{k} dt^3, \quad (3)$$

$$\mathbf{v}_p = \mathbf{v} + \mathbf{a} dt + \frac{1}{2} \mathbf{k} dt^2. \quad (4)$$

The acceleration ( $\mathbf{a}_p$ ) and jerk ( $\mathbf{k}_p$ ) are then recalculated at the predicted time from  $x_p$  and  $v_p$  using direct summation. Finally, a correction is based on the estimated higher-order derivatives:

$$\ddot{\mathbf{a}} = -6\Delta\mathbf{a}/dt^2 - (4\mathbf{k} + 2\mathbf{k}_p)/dt, \quad (5)$$

$$\ddot{\mathbf{k}} = 12\Delta\mathbf{a}/dt^3 + 6(\mathbf{k} + \mathbf{k}_p)/dt^2. \quad (6)$$

Here  $\Delta\mathbf{a} = \mathbf{a} - \mathbf{a}_p$ . Which then leads to the new position and velocity at time  $t + dt$ .

$$\mathbf{x} = \mathbf{x}_p + \frac{\ddot{\mathbf{a}}}{24} dt^4 + \frac{\ddot{\mathbf{k}}}{120} dt^5, \quad (7)$$

$$\mathbf{v} = \mathbf{v}_p + \frac{\ddot{\mathbf{a}}}{6} dt^3 + \frac{\ddot{\mathbf{k}}}{24} dt^4. \quad (8)$$

The new timestep is calculated using a new predicted second derivative of the acceleration  $\ddot{\mathbf{a}}_p = \ddot{\mathbf{a}} + \ddot{\mathbf{k}} dt$  for each particle  $i$  individually with (Aarseth, 1985)

$$dt = \left( v \frac{|\mathbf{a}_p| |\ddot{\mathbf{a}}_p| + \mathbf{k}^2}{|\mathbf{k}| |\ddot{\mathbf{k}}| + \ddot{\mathbf{a}}_p^2} \right)^{1/2}. \quad (9)$$

Here we use for accuracy parameter  $v = 0.01$ .

A single integration step in the integrator proceeds as follows:

- Determine which stars are to be updated. Each star has an individual time ( $t_i$ ) associated with it at which it was last advanced, and an individual time step ( $dt_i$ ). The list of stars to be integrated consists of those with the smallest  $t_i + dt_i$ . Time steps are constrained to be powers of 2, allowing “blocks” of many stars to be advanced simultaneously (McMillan and Aarseth, 1993).
- Before the step is taken, check for system reinitialization, diagnostic output, termination of the run, storing data.
- Perform low-order prediction of all particles to the new time  $t_i + dt_i$ . This operation may be performed on the GPU or GRAPE, whatever is available.

- Recompute the acceleration and jerk on all stars in the current block (using the GPU or GRAPE, if available), and correct their positions and velocities to fourth-order.

Note that this scheme is rather simple as it does not include treatment for close encounters, binaries or higher order (hierarchical or democratic) stable multiple systems.

### 3. The programming environment

The part of the algorithm that executes on the GPU (the force evaluation) is implemented in the Cg programming language (C for graphics, [Fernando and Kilgard \(2003\)](#), see [Supplementary material](#)), which has a syntax quite similar to C. The Cg programming environment includes a compiler and run-time libraries for use with the open graphics library (OpenGL)<sup>3</sup> and DirectX<sup>4</sup> graphics application programming interfaces. Though originally developed for the creation of real-time special effects without the need to program directly to the graphics hardware assembly language, researchers soon recognized the potential of Cg and started to apply it not only to high-performance graphics but also to a wide variety of “general-purpose computing” problems ([Fernando, 2004](#); [Pharr and Fernando, 2005](#)).

#### 3.1. Mapping the $N$ -body problem to a GPU

The challenge in the implementation of an efficient  $N$ -body code on a GPU lies in the mapping of the algorithm and the data to graphical entities supported by the Cg language. Particle data arrays are represented as “textures”. Normally, textures are used to represent pixel colour attributes with one single component (luminance, red, green, blue or alpha), three components (red, green and blue) or four components (red, green, blue and alpha). In our implementation we use multiple textures to represent the input and output data of  $N$  particles, as follows:

- Input: position ( $3N$ ), velocity ( $3N$ ), mass ( $N$ ), acceleration ( $3N$ ), jerk ( $3N$ ) and potential ( $N$ ).
- Output: acceleration ( $3N$ ), jerk ( $3N$ ) and potential ( $N$ ).

All values are represented as single-precision (32-bit) floating point numbers for a total of 21 floats or 84 bytes per particle. In [Supplementary material](#), we present a snippet of the source code in Cg, showing the implemented force evaluation routine. With the 768 Mbyte on-board memory of the GeForce 8800GTX it can store about 9 million particles, whereas the GRAPE-6Af can store only 128k (see [Table 1](#)).

Transferring data from CPU to GPU is accomplished through the definition of textures, which can either be read-only or write-only, but not both at the same time.

Table 1  
Detailed information on the hardware used in our experiments

Data	GRAPE-6Af	8800GTX	FX1400	Xeon	Unit
$n_{\text{pipe}}$	24	128	12	1	
$f_{\text{cycle}}$	90	575	350	3400	MHz
$1/t_{\text{bus}}$	0.133	2.0	2.0	NA	GB/s
Memory	128k	9362k	1562k	–	Particles

The first column gives the parameter followed by the four different hardware setups (GRAPE-6Af, GeForce 8800GTX, Quadro FX1400) and information about the host computer. The information for the GRAPE is taken from [Makino et al. \(2003\)](#), the GPU information is from <http://www.nvidia.com>. The hardware details are the number of processor pipelines ( $n_{\text{pipe}}$ ), the processor’s clock frequency ( $f_{\text{cycle}} \equiv 1/t_{\text{cycle}}$ ), the memory bandwidth for communication between host and attached processor ( $1/t_{\text{bus}}$ ), the amount of memory (in number of particles, one particle requires 84 bytes, here we adopt  $1\text{k} \equiv 1024$ ). For measured hardware parameters, see [Table 3](#). Note that the measured internal communication speed between the device memory and the GPU for the 8800GTX and the FX1400 are 86.4 Gbyte/s and 19.2 Gbyte/s, respectively.

The data structures in the CPU are then copied onto appropriately defined textures in the graphics card’s memory. Obtaining the results from GPU to CPU is done by reading back the pixels from the appropriate rendering targets into data structures on the host CPU. Therefore the output textures (acceleration, jerk and potential) are represented by a double-buffered scheme, where after each GPU computation the textures are swapped between reading and writing. There is some additional overhead (of order  $\mathcal{O}(N)$ ) for this operation which has to be performed every block time step.

Conventionally, graphics cards render into a “frame buffer”, a special memory area that represents the image seen on a display. However, a frame buffer is unsuitable for our purposes as the data elements in this buffer are “clamped” to value ranges that map the capabilities of the display. Invariably this means that 32-bit real vectors are reduced in resolution and therefore in accuracy too. This is perfectly fine for visual displays where the number of colours after clamping are still  $2^{24}$  ( $\approx 16$  million), sufficient to make two neighbouring colours indiscernible to the human eye. However, this is unacceptable for scientific production calculations. The workaround is to create an off-screen frame buffer object and instruct GPU programs to render into these rather than to the screen. Off-screen frame buffers support 32-bit floating point values and are not clamped and therefore preserve their precision.

The GPU has two main kernel operations available in programmable graphics pipelines, these are a “vertex shader” and a “fragment shader”. Our implementation only makes use of the fragment shader pipeline as it is better suited for the kind of calculations in the  $N$ -body problem and because the fragment pipeline in general provides more processing power.<sup>5</sup> The host CPU is responsible for allocating

<sup>5</sup> Before the 8800GTX family of GPUs, vertex programs and fragment programs had to execute on distinct processing units. The 8800GTX is the first generation of GPUs where this distinction no longer exists and the two are unified.

<sup>3</sup> See <http://www.opengl.org>

<sup>4</sup> See <http://www.microsoft.com/directx>

the input textures and frame buffer objects, copy the data between CPU and GPU, and binding textures that are to be processed by kernels. The lower order prediction and correction of the particle positions is also done on the host CPU. In Table 1 we summarize the hardware properties of the two adopted GPUs and the GRAPE.

#### 4. Results

To test the various implementations of the force evaluator we perform several tests on different hardware. For clarity we perform each test with the same realization of the initial conditions. For this we vary the number of stars from  $N = 256$  particles with steps of two to half a million stars (see Table 2). Not each set of initial condition is run on every processor, as the Intel Xeons, for example, would take a long time and the scaling with  $N$  is unlikely to change as we clearly have reached the CPU-limited calculation regime (see Section 5).

The initial conditions for each set of simulations were generated by randomly selecting the stellar positions and velocities according to the Plummer (1911) distribution using the method described by Aarseth et al. (1974). Each of the stars were given the same mass. The initial particle representations were scaled to virial equilibrium before starting the calculation.

Each set of initial conditions is run from  $t = 0$  to  $t = 0.50$  time units (Heggie and Mathieu, 1986),<sup>6</sup> but the performance is measured only over the last quarter of an  $N$ -body time unit to reduce the overhead for reading the snapshot and the initialization of the integrator. The maximum time step for the particles was 0.125, to guarantee that each particle was evaluated at least twice during the course of the simulation. For the minimum timestep we adopted  $1/2^{26}$ , and all time steps were calculated using  $v = 0.01$  in Eqs. (2) and (9). The force calculations were performed by adopting a softening of  $1/256$  for all simulations.

##### 4.1. Performance measurements

For our performance measurements we have used four nodes of a Hewlett-Packard xw8200 workstation cluster, each with a dual Intel Xeon CPU running at 3.4 GHz and either the GRAPE, a Quadro FX1400 or GeForce 8800GTX graphics card in the PCI Express (16x) bus. The cluster nodes were running a Linux SMP kernel version 2.6.16, Cg version 1.4, graphics card driver version 1.0-9746 and the OpenGL 2.0 bindings.

In Fig. 1, we show the timing results of the  $N$ -body simulations. The FX1400 is slower than the general purpose computer over the entire range of  $N$  in our experiments. The bad performance of the FX1400 is mainly attributed to the additional overhead in communication and memory allocation. For  $N \lesssim 10^4$  the GeForce 8800GTX GPU is

Table 2

Results of the performance measurements for a Plummer sphere with  $N$  equal mass particles initially in virial equilibrium for  $0.25N$ -body time units (from  $t = 0.25$  to  $t = 0.5$ ) using a softening of  $1/256$

$N$	GRAPE-6Af	8800GTX	FX1400	Xeon
256	0.07098	2.708	3.423	0.1325
512	0.1410	8.777	10.59	0.5941
1024	0.3327	17.46	20.20	2.584
2048	0.7652	45.27	54.16	10.59
4096	1.991	128.3	157.8	50.40
8192	5.552	342.7	617.3	224.7
16,384	16.32	924.4	3398	994.0
32,768	51.68	1907	13180	4328
65,536	178.2	3973	40560	19,290
131,072	–	8844	–	–
262,144	–	22,330	–	–
524,288	–	63,960	–	–

In the first column we list the number of particles, followed by the timing results of the GRAPE in seconds. In the last column we give the timing results for the calculation without an attached processor. The GRAPE (second column) was measured up to 64k particles, because the on-board memory did not allow for larger simulations. Simulations on the FX1400 and the host computer were limited to the same number for practical reasons.

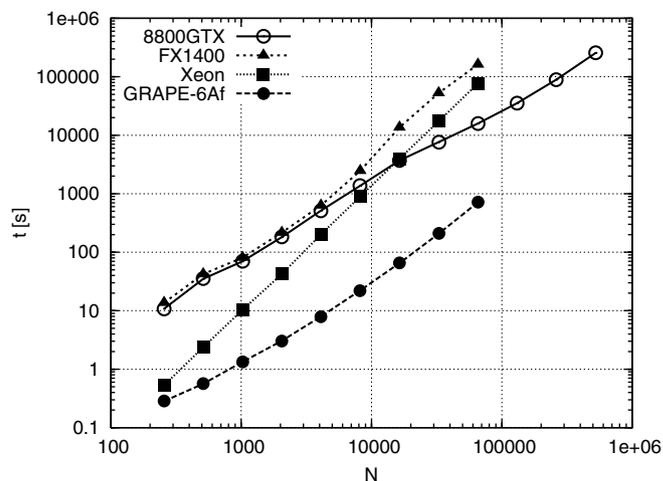


Fig. 1. Timing of several implementations of the gravitational  $N$ -body simulations for  $N = 256$  particles to  $N = 512k$  particles (only for the 8800GTX, the others up to 64k) over one  $N$ -body time unit. The 8800GTX are represented with open circles connected with a solid curve, the GRAPE is given by bullets with dashed line. The thin dashed (triangles) line and thin dotted (squares) lines give the results of the calculations with the FX1400 and with only the host computer. Note that the timings in Table 2 were multiplied by a factor of four to estimate the compute time for one dynamical time unit, rather than the  $1/4^{\text{th}}$  over which the timing calculations were performed.

slower than the host computer but continues to have a relatively flat scaling, comparable to the GRAPE-6, whereas the host has a much worse ( $\propto N^2$ ) scaling. The scaling of the compute time of the GPU is proportional to that of the GRAPE ( $\propto N^{3/2}$ ), but the latter has a smaller offset by about an order of magnitude. This is mainly caused by the efficient use of the GRAPE pipeline, which requires fewer clock cycles per force evaluation compared to the GPU (see Section 6).

<sup>6</sup> See also [http://en.wikipedia.org/wiki/Natural\\_units#N-body\\_units](http://en.wikipedia.org/wiki/Natural_units#N-body_units)

#### 4.2. Accuracy measurements

The GPU has a lower accuracy compared to the GRAPE or the host workstation. The GRAPE-6 uses 24-bit mantissa for calculating the differential position, and 64-bit fixed point format for accumulation. The pipeline for the time derivative is designed with 20-bit mantissa and 32-bit fixed-point notation for the final accumulation (Makino et al., 2003).

With the literature value of the mantissa the NVIDIA architecture should at most be able to reach an accuracy in relative energy ( $|\Delta E|/E$ ) of about  $1/10^6$ , whereas the GRAPE will be able to reach much higher accuracy. We tested this by running the initial conditions for 1k particles and 16k particles on the 8800GTX as well as on the GRAPE for a range of accuracy parameters ( $\nu$  in Eqs. (2) and (9)) ranging from  $\nu = 0.1$  (very low accuracy, but fast calculation) down to  $\nu = 0.1/2^8$  (very accurate but slow calculation) (Aarseth, 1985).

In Fig. 2, we present the degree to which energy  $E$  is conserved (in units of  $|\Delta E|/E$ ) for the simulations running on GRAPE (dashed lines) and the 8800GTX (solid curves) as a function of the accuracy parameter  $\nu$ . For relatively large values of  $\nu \gtrsim 0.02$  the GRAPE and 8800GTX produce similar energy errors, indicating that we are in the regime where the accuracy is limited by the integrator. For smaller values of  $\nu$ , however, the GRAPE has far superior energy conservation compared to the GPU. The error for the GRAPE continues to decrease for smaller values of  $\nu$ , until about  $|\Delta E|/E \sim 10^{-11}$ , whereas for the GPU the energy error does not drop below  $|\Delta E|/E \sim 10^{-7}$ . Note that for all the performance calculations in Section 4.1 we adopted  $\nu = 0.01$ , which produces about maximum accuracy achievable for the GPU. For the GRAPE, however, we could have used

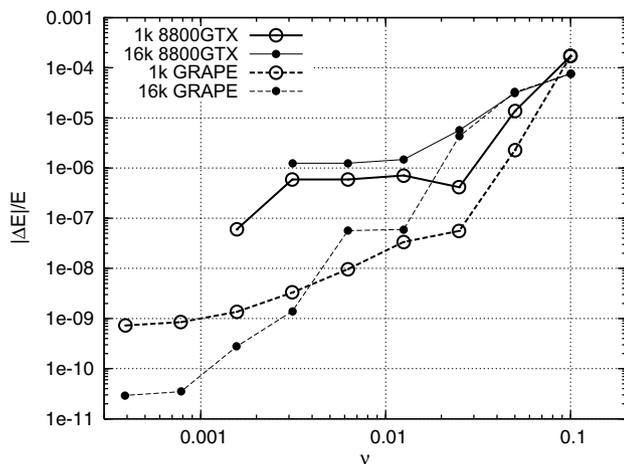


Fig. 2. Absolute value of the relative error in the energy as a function of the timestep control parameter  $\nu$ . For clarity we only present the result for simulations with 1k (thick curves with circles) and 16k (thin curves with bullets) particles with the 8800GTX (solid curves) and GRAPE (dashes). The simulations were performed over 0.25  $N$ -body time units using identical input realizations as adopted for the performance measurements.

much smaller integration time steps, resulting in considerably smaller energy errors. The price to pay, however, would be a much longer compute time.

#### 5. Performance modelling of the GPU

In modelling the performance of the GPU we adopt the model proposed by Makino (2002), Harfst et al. (2007) but tailored to the host plus GPU and to the GRAPE architecture.

The wall clock time required for advancing the  $n_{\text{block}}$  particles in a single block time step in the  $N$ -body systems is

$$t_{\text{step}} = t_{\text{host}} + t_{\text{force}} + t_{\text{comm}}. \quad (10)$$

Here  $t_{\text{host}} = t_{\text{pred}} + t_{\text{corr}}$  is the time spent on the host computer for predicting and correcting the particles in the block,  $t_{\text{force}}$  is the time spent on the attached processor and  $t_{\text{comm}}$  is the time spent communicating between the host and the attached processor. We now discuss the characteristics of each of the elements in the calculation for  $t_{\text{step}}$ .

##### 5.1. Host operation

The predictions and corrections of the particles are calculated on the host computer, and the time for this operation is directly related to the speed of the host processor  $t_{\text{cpu}}$ , the number of operations in the prediction step  $\eta_{\text{pred}}$  and in the correction step  $\eta_{\text{corr}}$ . The total time spent per block step then yields

$$t_{\text{pred}} \simeq \eta_{\text{pred}} t_{\text{cpu}} N, \quad (11)$$

for the prediction and

$$t_{\text{corr}} \simeq \eta_{\text{corr}} t_{\text{cpu}} n_{\text{block}}, \quad (12)$$

for the correction. The number of clock cycles per prediction step  $\eta_{\text{pred}} \approx 900$  and for the correction  $\eta_{\text{corr}} \approx 16,000$ . This operation could be performed on the GPU, though the GRAPE is not designed for the predictor and corrector calculation. For a fair comparison between the GRAPE and the GPU and in order to preserve high accuracy we performed these calculations on the host.

##### 5.2. Communication

The time spent communicating between the host and the attached processor is expressed by the sum of the time needed to send  $n_{\text{send}}$  particles to the acceleration hardware and the time needed to receive  $n_{\text{rec}}$  particles from the acceleration hardware:

$$t_{\text{comm}} = \eta_{\text{send}} t_{\text{send}} n_{\text{send}} + \eta_{\text{rec}} t_{\text{rec}} n_{\text{send}}. \quad (13)$$

Here  $\eta_{\text{send}} t_{\text{send}}$  and  $\eta_{\text{rec}} t_{\text{rec}}$  are the time needed to send and receive one particle, respectively. For the computation without the hardware acceleration  $t_{\text{comm}} = 0$ , since the forces between all particles are calculated locally. For the

GRAPE and the GPUs, however, a considerable amount of time is spent in communication. For the GRAPE sending data are equally fast as receiving data, i.e.,  $t_{\text{send}} = t_{\text{rec}} = t_{\text{bus}}$ . Sending data to the GPU are considerably slower than receiving data (see Table 3).

The two send and receive efficiency factors  $\eta_{\text{send}}$  and  $\eta_{\text{rec}}$  are the product of the overhead  $\eta_0$  and the number of bytes per particle that has to be sent or received. The overhead  $\eta_0 = 188$  (Fukushige et al., 2005) for each of the attached processors. Since for the GRAPE the send and receive operation are equally expensive we can just count the number of bytes that has to be transported per particle, which for the GRAPE hardware is 72 bytes (Fukushige et al., 2005; Harfst et al., 2007). For the GRAPE we then write  $\eta_{\text{send}} t_{\text{send}} + \eta_{\text{rec}} t_{\text{rec}} = 72 \times 188 t_{\text{bus}}$ .

For the GPU  $\eta_{\text{send}} > \eta_{\text{rec}}$  (see Table 3 for the measured values). The additional overhead  $\eta_0$  is the same as for the GRAPE, but per particle the number of bytes to send is different from the number to receive. As we discussed in Section 3.1 a total of 56 bytes has to be sent from the host to the GPU, whereas only 28 bytes are received. For the GPU we then write  $\eta_{\text{send}} = 56 \times 188$ , whereas  $\eta_{\text{rec}} = 28 \times 188$ .

In addition to the difference in the speeds for sending and receiving data, the GPU suffers from an additional penalty. The GRAPE sends the particles in the block ( $n_{\text{send}} = n_{\text{block}}$ ), whereas due to internal memory management the GPU has to send and receive all particles  $n_{\text{send}} = N$  (see Section 3.1). This efficiency loss is quite substantial, but is reduced when we use CUDA as programming environment (see Section 6).

For the adopted (Hermite predictor–corrector block time-step) integration scheme the number of particles in a single block  $n_{\text{block}}$  cannot be determined implicitly, though theoretical arguments suggest  $n_{\text{block}} \propto N^{2/3}$ . Instead of using this estimate, we fitted the average number of particles in a block time step. This fit was done with the equal mass Plummer sphere initial conditions running on GRAPE and run over one dynamical ( $N$ -body) time unit. The average number of particles in a single block is then

$$n_{\text{block}} \simeq 0.20N^{0.81}. \quad (14)$$

Table 3  
Time measurements (in seconds) for the various hardware operations used in this paper

Param	GRAPE-6Af	8800GTX	FX1400	Xeon
$t_{\text{host}}$	$3.82 \times 10^{-7}$	$3.82 \times 10^{-7}$	$3.82 \times 10^{-7}$	$3.82 \times 10^{-7}$
$\eta_{\text{fe}} t_{\text{cycle}} / n_{\text{pipe}}$	$4.63 \times 10^{-10}$	$8.15 \times 10^{-10}$	$1.43 \times 10^{-8}$	$5.29 \times 10^{-8}$
$\eta_{\text{send}} t_{\text{send}}$	$8.00 \times 10^{-7}$	$1.76 \times 10^{-5}$	$1.89 \times 10^{-5}$	NA
$\eta_{\text{rec}} t_{\text{rec}}$	$8.00 \times 10^{-7}$	$5.97 \times 10^{-6}$	$5.98 \times 10^{-6}$	NA

The first line gives the time spent by the host computer for predicting and correcting a single particle. The second row is for calculating the force between two particles. The last two rows give the time to send a single particle to, and to receive a single particle from the attached hardware. For the calculations with only the host computer this operation is not available. In particular the communication with the GPUs turns out to be relatively slow.

### 5.3. Calculation

The time spent by the hardware acceleration ( $t_{\text{force}}$ ) is directly related to the speed of the dedicated processor ( $t_{\text{cycle}}$ ), the number of pipelines per processor ( $n_{\text{pipe}}$ ) and the number of operations for one force evaluation ( $\eta_{\text{fe}} \simeq 60$ ).

$$t_{\text{force}} = \eta_{\text{fe}} N n_{\text{block}} t_{\text{cycle}} / n_{\text{pipe}}. \quad (15)$$

The details of the different hardware are presented in Table 1 and the measured values are in Table 3. The GRAPE has a vector pipeline for each processor which allows a more efficient force evaluation than the GPU, the number of operations per force evaluation for the GRAPE is therefore  $\eta_{\text{fe}} \simeq \mathcal{O}(1)$ .

In order to enable hardware acceleration on our  $N$ -body code we had to introduce a number of additional operations, like reallocating arrays, which give rise to an extra computation overhead. For the calculations with the host computer without hardware acceleration we adopt  $\eta_{\text{fe}} \simeq 180$ , a factor of three larger than for the GPUs.

### 5.4. Total performance

The total wall-clock time spent per dynamical ( $N$ -body) time unit is then

$$t = n_{\text{steps}} t_{\text{step}}. \quad (16)$$

Here we fitted the number of block steps per dynamical ( $N$ -body) time units. According to Makino and Hut (1988), Makino and Hut (1990)  $n_{\text{steps}} \propto n^{1/3}$ . We measured the number of block time steps using the equal mass Plummer distributions as initial conditions, using the GRAPE enabled code and fitted the result:

$$n_{\text{steps}} \simeq 247N^{0.35}. \quad (17)$$

In Fig. 3, we compare the results of the performance model with the measurements on the workstation without additional hardware (squares) and with three attached processors; a single GRAPE-6Af processor board (bullets), an FX1400 (triangles) and the newer GeForce 8800GTX (circles). Note that the measurements in Table 2 were multiplied by a factor four to compensate for the fact that we performed our timings only over a quarter  $N$ -body time unit. Though these curves are not fitted, they give a satisfactory comparison.

The largest discrepancy between the performance model and the measurements can be noticed for the FX1400 GPU, which, for  $N \gtrsim 10^4$  seems to perform considerably less efficient than expected according to the performance model. Part of this discrepancy, though not explicitly mentioned in Section 5.2, is the result of a hysteresis effect in the communication of both GPUs. For the 8800GTX, however, this effect is less evident, but still present. Both GPUs tend to have a maximum communication speed for blocks of 0.5 Mbyte (about 6000 particles). An additional effect which causes performance loss on the FX1400 is

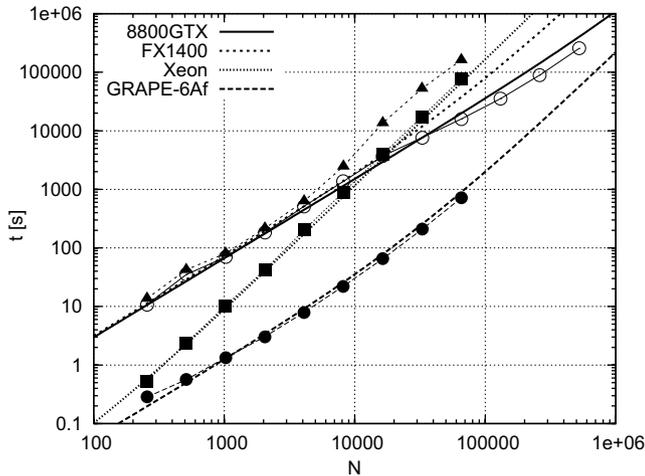


Fig. 3. The results of the above described performance model (thick lines) over-plotted with the results of the measurements for the three attached processors (symbols). The bullets represents the results from a single GRAPE-6Af processor board, the squares give the host workstation, the circles are for the GeForce 8800GTX and the triangles give the FX1400 graphics processor.

the increase in the number of block time steps. This number continued to increase beyond our measurements performed with GRAPE (see Eq. (14)).

The numbers listed in Table 3, and used in our performance model, are the optimum values. The communication speed drops by about a factor of two for much larger amounts of data transfer to and from the GPU. For the FX1400, this drop in communication is considerable, whereas for the 8800GTX it results in a smaller performance loss (mainly due to the larger number of processor pipelines). The discrepancy for the GRAPE calculation with low  $N$  is the result of neglecting the limited size of the processor pipeline in the performance model and due

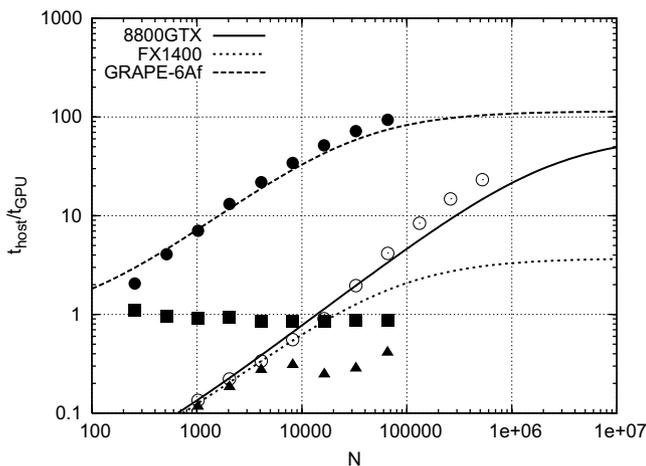


Fig. 4. The speed-up of the two GPUs and the GRAPE with respect to the host workstation as a function of the number of particles. The (lower) dotted curve is for the Quadro FX1400, the solid curve (middle) gives the timing for the GeForce 8800GTX and the top line (dashes) represents the GRAPE.

to the irregular behavior of the number of particles in each block time step.

In Fig. 4, we present the speed-up for the various hardware configurations, compared to running on the host workstation. Here it is quite clear that for low  $N$  the GPUs do not give a appreciable speedup, but for a large number of particles, the GeForce 8800GTX gives a speed-up of at least an order of magnitude, but not as much as the GRAPE. The latter, however, will not be able to perform simulations of more than 128k particles.<sup>7</sup>

## 6. Discussion

We have successfully implemented the direct gravitational force evaluation calculation using Cg on two graphics cards, the NVIDIA Quadro FX1400 and the NVIDIA GeForce 8800GTX, and compared their performance with the host workstation and the GRAPE-6Af special purpose computer.

For  $N \lesssim 10^4$  particles the workstation outperforms the GPUs. This is mainly due to additional overhead introduced by the communication to the GPU and memory allocation on the GPU. For a larger number of particles the more modern GPU (8800GTX) outperforms the workstation by up to about a factor of 50 (for 9 million particles). Such a large number of particles cannot be simulated on the GRAPE-6Af, due to memory limitations. For up to 256k, the maximum number of particles that can be stored on the GRAPE, the 8800GTX is slower than the GRAPE by a factor of a few. Still, at this particle number the GPU is faster than the workstation by an order of magnitude.

For the adopted accuracy  $v = 0.01$ , the average mean error in the energy measured over 0.25  $N$ -body time unit is  $|\Delta E|/E = (1.7 \pm 1.6) \times 10^{-6}$  for the 8800GTX and  $(5.1 \pm 0.56) \times 10^{-6}$  for the FX1400 (averaged over the simulations for  $N = 256$  to  $N = 64k$ ), whereas for the GRAPE we measured  $(1.9 \pm 1.2) \times 10^{-7}$ , which is comparable to the mean error on the host. For the adopted accuracy, both the host and GRAPE produce an energy error which is about an order of magnitude smaller than that of the GPUs. For smaller values of  $v$  the energy errors in the GRAPE continue to decrease whereas for the GPU this is not the case, as we have reached the precision of the hardware. For many applications an energy error of  $|\Delta E|/E = \mathcal{O}(10^{-7})$  may be satisfactory.

In the release notes of CUDA version 0.8, NVIDIA announced that GPUs supporting 64-bit double precision floating point arithmetic in hardware will become available in late 2007. In the meantime, we could improve the accuracy of the GPU by sorting the forces on size before adding them, summing the smallest forces first.

<sup>7</sup> Due to a defective chip on our GRAPE-6 the on-board memory was reduced from 128k particles to 64k particles. The latest GRAPE-6Af are equipped with 256k particles of memory.

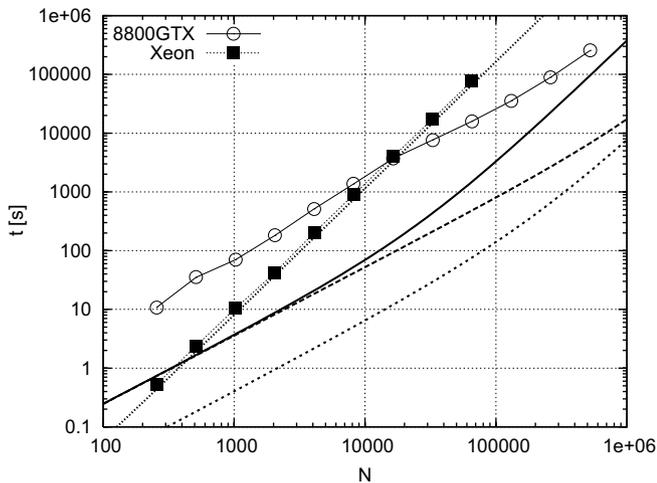


Fig. 5. Prospective of future CPU and GPU performance, based on the model from Section 5. The two thin curves with squares and circles give the measured performance of the CPU and 8800GTX GPU, respectively. The thick solid curve gives a prediction of the performance for the 8800GTX in which only blocks of particles are communicated with the GPU. The dashed curves gives in addition the effect of carrying the predictor/corrector calculation to the GPU. The dotted curve gives the performance of a hypothetical 8800GTX-like architecture for which in addition the processor pipeline would be used more efficiently ( $\eta_{ic} = 1$ ).

The GRAPE-6 is much more efficient in using its clock cycles, allowing effectively one operation per clock cycle, whereas the NVIDIA architecture requires more cycles. This turns out to be an important reason why the 8800GTX is slower than the GRAPE-6.

The main advantage of the GPU over that of the dedicated GRAPE hardware, is the much larger memory, the wider applicability and the much lower cost of the former. The large memory on the GPU allows simulations of up to about 9 million particles, though one has to wait for about two years for one dynamical time scale.

In theory the 8800GTX should be able to outperform the GRAPE-6Af, but due to relatively inefficient memory access and additional overhead cost, which is not present in the GRAPE hardware, many clock cycles seem to get lost. With a more efficient use of the hardware the GPU could, in principle, improve performance by about two orders of magnitude. For the next generation of GPUs we hope that this efficiency bottleneck will be lifted. In that case, the GPU would outperform GRAPE by almost an order of magnitude. Note, however, that the GRAPE-6 is based on 5 year old technology, and the next generation GRAPE is likely to outperform modern GPUs by a sizable margin.

These current bottlenecks in the GPU may be reduced using the Compute Unified Device Architecture (CUDA)<sup>8</sup> programming environment, which is supposed to provide an improved environment for general purpose programming on the GPU. In Fig. 5, we present the possible future

performance assuming that the additional communication overhead on the GPU is lifted, the clock cycles are used more efficiently without any assumptions of improved hardware speed. In the first step we simply reduce communication to blocks rather than having to transport all particles each block time step (solid curve). This relatively simple improvement has recently been carried out using CUDA (Hamada and Iitaka, 2007; Bédorf et al., in preparation). The second optimization (dashed curve in Fig. 5) is achieved when, in addition to reducing the communication we also carry the predictor and corrector steps to the GPU. This improvement, however, may be associated with a quite severe accuracy penalty. For both improvements we used the performance data for the current design 8800GTX. Further improvement can be achieved when, in addition to more efficient communication and force computing pipeline is further optimized. The result of this hypothetical case would improve performance by more than a factor 100 compared to the workstation over the entire range of  $N$ .

## Acknowledgments

We are grateful to Mark Harris and David Luebke of NVIDIA for supplying us with the two NVIDIA GeForce 8800GTX graphics cards on which part of the simulations were performed. We also thank Jeroen Bédorf, Derek Groen, Alessia Gualandris and Jun Makino for numerous discussions, and the referee Piet Hut for pointing us to the importance of discussing the accuracy of the GPU. This work was supported by NWO (via Grants #635.000.303 and #643.200.503) and the Netherlands Advanced School for Astrophysics (NOVA). The calculations for this work were done on the Hewlett-Packard xw8200 workstation cluster and the MoDeStA computer in Amsterdam, both are hosted by SARA computing and networking services, Amsterdam.

## Appendix A. Supplementary material

The  $N$ -body code presented in this paper consists of a part implemented in C (running on a CPU) and a part implemented in Cg (running on the GPU). In this appendix, we show the routine that evaluates the acceleration, jerk and potential in Cg (which was based on a tutorial available from Göddeke, 2005). The C code which handles communication between CPU and GPU and supporting data structures is not presented here. A copy of the entire working version of the code is available via <http://modesta.science.uva.nl>. The online version contains the full source for the sample  $N$ -body code with Hermite individual-timestep scheme. The supplementary data associated with this article can be found, in the online version, at Supplementary data associated with this article can be found, in the online version, at [doi:10.1016/j.newast.2007.05.004](https://doi.org/10.1016/j.newast.2007.05.004).

<sup>8</sup> See <http://developer.nvidia.com/cuda>

## References

- Aarseth, S.J., 1999. *PASP* 111, 1333.
- Aarseth, S.J., Henon, M., Wielen, R., 1974. *A&A* 37, 183.
- Aarseth, S.J., Hoyle, F., 1964. *ApNr* 9, 313.
- Aarseth, S.J., 1985. Direct methods for  $N$ -body simulations. In: Brackbill, J.C., Cohen, B.I. (Eds.), *Multiple Time Scales*, p. 377.
- Aarseth, S.J., Lecar, M., 1975. *ARA&A* 13, 1.
- Applegate, J.H., Douglas, M.R., Gürsel, Y., Hunter, P., Seitz, C.L., Sussman, G.J., 1986. In: Hut, P., McMillan, S.L.W. (Eds.), *LNP Vol. 267: The Use of Supercomputers in Stellar Dynamics*, p. 86.
- Bédorf, P., Belleman, R., Portegies Zwart, S. In: *The International Conference for High Performance Computing, Networking, Storage, and Analysis* (in preparation).
- Buck, I., Foley, T., Horn, D., Sugerma, J., Mike, K., Pat, H., 2004. *Brook for GPUs: stream computing on graphics hardware*.
- Dorband, E.N., Hemsendorf, M., Merritt, D., 2003. *J. Comput. Phys.* 185, 484.
- Elsen, E., Houston, M., Vishal, V., Darve, E., Hanrahan, P., Pande, V., 2006. In: *N-Body Simulation on GPUs Proceedings of the ACM/IEEE SC2006 Conference on High Performance Networking and Computing*, November 11–17, 2006. Tampa, FL, USA, ACM Press, New York, p. 188.
- Fernando, R., 2004. *GPU Gems (Programming Techniques, Tips, and Tricks for Real-Time Graphics)*. Addison-Wesley, ISBN 0-321-22832-4.
- Fernando, R., Kilgard, M.J., 2003. *The Cg Tutorial (The Definitive Guide to Programmable Real-Time Graphics)*. Addison-Wesley, ISBN 0-321-19496-9.
- Fukushige, T., Makino, J., Kawai, A., 2005. *PASJ* 57, 1009.
- Göddeke, D., 2005. *GPGPU – basic math tutorial*. *Ergebnisberichte des Instituts für Angewandte Mathematik*, Nummer 300.
- Gualandris, A., Portegies Zwart, S., Tirado-Ramos, A., 2007. *PARCO* 33, 159.
- Hamada, T., Iitaka, T., 2007. *NewA*. ArXiv Astrophysics e-prints (astro-ph/0703100) (submitted for publication).
- Harfst, S., Gualandris, A., Merritt, D., Spurzem, R., Portegies Zwart, S., Berczik, P., 2007. *NewA* 12, 357.
- Heggie, D.C., Mathieu, R.D., 1986. In: Hut, P., McMillan, S.L.W. (Eds.), *LNP Vol. 267: The Use of Supercomputers in Stellar Dynamics*, p. 233.
- Hoekstra, A., Portegies Zwart, S., Bubak, M., Sloot, P., 2007. CRC Press LLC. ArXiv Astrophysics e-prints (astro-ph/0703485) (submitted for publication).
- Holmberg, E., 1941. *ApJ* 94, 385.
- Hut, P., 2007. Presented at *A Life With Stars (Conference in Honor of Ed van den Heuvel)*, Amsterdam, August, 2007. ArXiv Astrophysics e-prints (astro-ph/0601232).
- Makino, J., 1991. *ApJ* 369, 200.
- Makino, J., 2001. In: Deiters, S., Fuchs, B., Just, A., Spurzem, R., Wielen, R. (Eds.), *ASP Conference Series 228: Dynamics of Star Clusters and the Milky Way*, p. 87.
- Makino, J., 2002. *NewA* 7, 373.
- Makino, J., 2005a. *JKoAS* 38, 165.
- Makino, J., 2007. ArXiv Astrophysics e-prints (astro-ph/0509278).
- Makino, J., Aarseth, S.J., 1992. *PASJ* 44, 141.
- Makino, J., Fukushige, T., Koga, M., Namura, K., 2003. *PASJ* 55, 1163.
- Makino, J., Hut, P., 1988. *ApJS* 68, 833.
- Makino, J., Hut, P., 1990. *ApJ* 365, 208.
- Makino, J., Taiji, M., 1998. Scientific simulations with special-purpose computers: The GRAPE systems. In: Junichiro Makino, Makoto Taiji (Eds.), *Scientific Simulations with Special-purpose Computers: The GRAPE Systems*. Wiley, Chichester, Toronto.
- McMillan, S.L.W., Aarseth, S.J., 1993. *ApJ* 414, 200.
- Moore, G.E., 1965. *Electronics* 38 (8).
- Nitadori, K., Makino, J., Hut, P., 2006. *NewA* 12, 169.
- Nitadori, K., Makino, J., Abe, G., 2007. In: *Conference Proceedings of Computational Science 2006*. ArXiv Astrophysics e-prints (astro-ph/0606105) (in press).
- Nyland, L., Harris, M., Prins, J., 2004. Poster presented at *The ACM Workshop on General Purpose Computing on Graphics Hardware*, August 7–8, Los Angeles, CA.
- Pharr, M., Fernando, R., 2005. *GPU Gems 2 (Programming Techniques for High-Performance Graphics and General-Purpose Computation)*. Addison-Wesley, ISBN 0-321-33559-7.
- Plummer, H.C., 1911. *MNRAS* 71, 460.
- Taiji, M., Makino, J., Fukushige, T., Ebisuzaki, T., Sugimoto, D., 1996. In: Hut, P., Makino, J. (Eds.), *IAU Symposium 174: Dynamical Evolution of Star Clusters: Confrontation of Theory and Observations*, p. 141.
- van Albada, T.S., 1968. *BAN* 19, 479.
- von Hoerner, S., 1963. *ZA* 57, 47.