

Author's Accepted Manuscript

Massively parallel forward modeling of scalar and tensor gravimetry data

M. Moorkamp, M. Jegen, A. Roberts, R. Hobbs

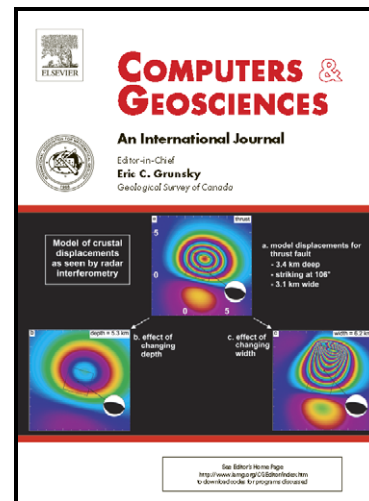
PII: S0098-3004(10)00057-9
DOI: doi:10.1016/j.cageo.2009.09.018
Reference: CAGEO 2305

To appear in: *Computers & Geosciences*

Received date: 16 June 2009
Revised date: 2 September 2009
Accepted date: 4 September 2009

Cite this article as: M. Moorkamp, M. Jegen, A. Roberts and R. Hobbs, Massively parallel forward modeling of scalar and tensor gravimetry data, *Computers & Geosciences*, doi:10.1016/j.cageo.2009.09.018

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting galley proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



www.elsevier.com/locate/cageo

Massively parallel forward modeling of scalar and tensor gravimetry data[☆]

M. Moorkamp^{*,a}, M. Jegen^a, A. Roberts^b, R. Hobbs^b

^a*IFM-GEOMAR, Wischhofstrasse 1-3, 24148 Kiel, Germany*

^b*University of Durham, Department of Earth Sciences, South Road, Durham DH1 3LE, UK*

Abstract

We present an approach to calculate scalar and tensor gravity utilizing the massively parallel architecture of consumer graphics cards. Our parametrization is based on rectilinear blocks with constant density within each blocks. This type of parametrization is well suited for inversion of gravity data or joint inversion with other datasets, but requires the calculation of a large number of model blocks for complex geometries. For models exceeding 10,000 cells we achieve an acceleration of a factor of 40 for scalar data and 30 for tensor data compared to a single thread on the CPU. This significant acceleration allows fast computation of large models exceeding 10^6 model parameters and thousands of measurement sites.

Key words: gravity modeling, cuda, parallel computing

[☆]Code available from server at <http://www.iamg.org/CGEditor/index.htm>

*Corresponding author

Email address: mmoorkamp@ifm-geomar.de (M. Moorkamp)

1 1. Introduction

2 Driven by the computer games industry graphics cards (GPUs) have
3 evolved into powerful computing devices that are geared towards a large
4 number of simultaneous calculations and high memory bandwidth (e.g. Ryoo
5 et al., 2008). In an attempt to broaden the scope of their products, the two
6 main consumer graphics cards manufacturers, Nvidia and AMD, have re-
7 leased programming interfaces for general purpose calculations to their cards.
8 So far massively parallel architectures were limited to specialized and costly
9 hardware. With these developments such an architecture becomes available
10 at low prices and makes the development of massively parallel algorithms
11 attractive.

12 The success of solving a numerical problem on a massively parallel archi-
13 tecture depends heavily on the anatomy of the algorithm. If the problem can
14 be split into independent parts that can be solved without having to transfer
15 information, parallelization is easy and we can expect good performance. If
16 conversely results have to be distributed globally during the calculation, par-
17 allelization becomes difficult and special care has to be taken to reduce the
18 amount of synchronization between the parallel threads of the program. The
19 challenge for GPU based computations is that the number of threads has to
20 be on the order of 10,000 or more to utilize the full computing power of the
21 architecture (Nickolls et al., 2008; Ryoo et al., 2008; Jeong and Whitaker,
22 2008; Komatitsch et al., 2009).

23 Modeling gravitational acceleration and its spatial derivatives is a com-
24 mon tool in geophysics to test models of the density distribution within the
25 subsurface. Often tectonic information or seismic models are used to de-

26 fine broad geological structures with a common density and these are then
27 parametrized as polygonal bodies within the numerical modeling scheme (e.g.
28 Götze and Lahmeyer, 1988). This type of approach has the advantage that
29 the number of bodies is kept low even for complex models which makes it easy
30 for the user to construct such a model and reduces the number of function
31 evaluations.

32 Our forward modeling approach is geared towards usage within a joint
33 inversion algorithm that combines gravity, seismic and magnetotelluric data
34 (Heincke et al., 2006) and therefore we parametrize our model in terms of
35 rectilinear blocks (Hobbs and Trinks, 2005). This type of setup is also of-
36 ten used for inversion of gravity data alone (e.g. Li and Oldenburg, 1998;
37 Portniaguine and Zhdanov, 1999; Nagihara and Hall, 2001; Chasseriau and
38 Chouteau, 2003) and has the advantage that the equations for scalar and
39 tensor gravimetry are particularly simple, but requires the calculation of the
40 effect of a large number of blocks, as complex geometries have to be con-
41 structed from many small blocks. On a platform with no or only a low
42 degree of parallelism this leads to increased computational times compared
43 to the polygonal parametrization. However, the calculation of the effect of
44 many rectilinear block can be performed effectively on a massively parallel
45 architecture to compensate for the higher computational cost. This cost be-
46 comes particularly relevant when we have to calculate several large models
47 for which we cannot store the sensitivities in main memory or even on disk,
48 for example within a non-linear inversion.

49 Although gravity forward modeling is generally fast compared to other
50 methods and we restrict ourselves here to Nvidia's CUDA interface the con-

51 clusions and strategies for this relatively simple problem can be applied to
 52 other problems and other massively parallel architectures. Before we describe
 53 the details of our implementation we will discuss the basic equations of the
 54 gravimetry problem for rectilinear blocks. We will then show the performance
 55 of our approach for a number of scenarios and discuss the implications for
 56 forward modeling and inversion of gravimetric data.

57 2. Basic equations

The two quantities that are mainly used in gravimetry surveys, are the vertical gravitational acceleration U_z , i.e. the vertical derivative of the gravitational potential U and the gravitational tensor Γ , i.e. the tensor of second spatial derivatives,

$$\Gamma = \begin{pmatrix} U_{xx} & U_{xy} & U_{xz} \\ U_{yx} & U_{yy} & U_{yz} \\ U_{zx} & U_{zy} & U_{zz} \end{pmatrix}. \quad (1)$$

With the nomenclature shown in Figure 1 the equation for the effect of a single prism of density ρ on the vertical gravitational acceleration U_z is (Li and Chouteau, 1998)

$$U_z = -\gamma\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 \mu_{ijk} \left(x_i \ln(y_j + r_{ijk}) + y_j \ln(x_i + r_{ijk}) + z_k \arctan \frac{z_k r_{ijk}}{x_i y_j} \right), \quad (2)$$

and for two elements of the gravimetry tensor it is

$$U_{xx} = \gamma\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 \mu_{ijk} \arctan \frac{y_j z_k}{x_i r_{ijk}}, \quad (3)$$

$$U_{xy} = -\gamma\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 \mu_{ijk} \ln(z_k + r_{ijk}), \quad (4)$$

where

$$\begin{aligned} x_i &= x - \xi_i & y_j &= y - \eta_j & z_k &= z - \zeta_k \\ r_{ijk} &= \sqrt{x_i^2 + y_j^2 + z_k^2} \\ \mu_{ijk} &= (-1)^i (-1)^j (-1)^k. \end{aligned}$$

We can calculate all other elements of the gravimetry tensor by permutation of the coordinate axes (e.g. Li and Chouteau, 1998; Nagy et al., 2000), in addition the tensor is symmetric so that we only have to calculate 6 instead of all 9 tensor elements. Theoretically, we even only have to calculate 5 elements, as the diagonal terms of the tensor are related by Poisson's equation

$$\frac{\partial^2 U}{\partial x^2} + \frac{\partial^2 U}{\partial y^2} + \frac{\partial^2 U}{\partial z^2} = -4\pi\gamma\rho. \quad (5)$$

58 However, we calculate all three diagonal elements independently as this gives
59 us an indication of the numerical precision of the results.

Scalar and tensor gravity calculation are well known linear problems and therefore in both cases a term that is purely determined by the geometry of the cell is multiplied by the density of the cell (e.g. Nagy et al., 2000). Also, the effect of several prisms is simply the sum of the contributions of a single cell. We can therefore write the forward calculation as a vector-matrix multiplication between the model vector of density values \mathbf{m} and the geometric sensitivities G

$$\mathbf{d} = G\mathbf{m}. \quad (6)$$

60 Here each row of G corresponds to one observed quantity, i.e. a measurement
61 of the vertical acceleration or an element of the gravimetric tensor. The
62 resulting data vector \mathbf{d} contains the data resulting from the model. We

63 therefore have two parts in the calculation of the forward problem, 1) the
64 calculation of the elements of G and 2) the evaluation of the matrix vector
65 product.

66 3. Implementation

67 Before we describe the details of our implementation we have to clarify the
68 standard nomenclature for the CUDA interface and briefly explain the archi-
69 tecture. A function that can be executed on the GPU is called a kernel and is
70 described by the extended C-syntax `kernelname<<<dimGrid,dimBlock>>>(Parameters)`.
71 Here `dimGrid` and `dimBlock` are variables that describe the number of inde-
72 pendent thread blocks in the computing grid and the number of threads in
73 each block, respectively (see Figure 2). The number of threads in a single
74 block is determined by the specifications of the GPU and is typically between
75 64 and 512 to optimize memory access by the hardware (nvidia, 2009). In
76 principle different threads within a block can share information, but we will
77 not use this feature in our implementation. The size of the grid depends on
78 the size of the problem, in our case the number of model parameters M , and
79 each block can be computed independently and in any order. During the par-
80 allel execution of the kernel the implementation determines the sub-problem
81 to work on from the two variables `blockIdx` and `blockDim`. The values of
82 these variables is set by the GPU depending on the current block index and
83 thread index for the calculation. In principle this index can have several
84 dimensions, we only use the first dimension `blockIdx.x` and `blockDim.x`,
85 respectively.

86 As we can calculate each element of the sensitivity matrix independently

87 and with relatively few input parameters, this part can be performed very
88 efficiently. We parallelize over the number of grid cells M , i.e. a single
89 row of the matrix G . In principle, it would be possible to also parallelize
90 over the number of measurements N to obtain $N * M$ independent threads.
91 However, for large models, for which the parallelization makes most sense, M
92 already exceeds one million or more and therefore we can utilize the threading
93 capabilities of all currently available GPUs. By only parallelizing over the
94 grid cells, we avoid additional administrative overhead and also avoid having
95 to store the full sensitivity matrix if we do not need it, instead we only have
96 to store a single row at a time. The following listing shows the core algorithm
97 using NVidia's CUDA API.

```
98 __global__ void CalcScalarMeas(const double x_meas, const double y_meas,  
99     const double z_meas, const double *XCoord, const double *YCoord,  
100     const double *ZCoord, const double *XSizes, const double *YSizes,  
101     const double *ZSizes, const int nx, const int ny, const int nz,  
102     double *Grow)  
103  
104 //calculate memory offset from execution parameters  
105 const unsigned int offset = blockIdx.x * blockDim.x + threadIdx.x;  
106 int xindex, yindex, zindex;  
107 //if the offset is within the model size  
108 if (offset < nx * ny * nz)  
109 {  
110     //calculate the coordinate indices for all three directions  
111     OffsetToIndex(offset, ny, nz, xindex, yindex, zindex);
```



```
112 //calculate and assign the geometric term to the
113 //row of the sensitivity matrix
114 Grow[offset] = CalcGravBoxTerm(x_meas, y_meas, z_meas,
115     XCoord[xindex], YCoord[yindex], ZCoord[zindex], XSizes[xindex],
116     YSizes[yindex], ZSizes[zindex]);
117
118 }
```

119 We generate the storage offset for the results within the current row of
120 the sensitivity matrix from the built-in variables `blockId.x`, `blockDim.x` and
121 `threadId.x`. As mentioned above, the values of these variables are set by the
122 hardware for each executed thread. Therefore each offset is unique within one
123 calculation of the sensitivities. The optimum number of blocks `blockDim.x`
124 depends on the register use and the ability to load data from global memory
125 to local memory in a coalesced fashion. The CUDA programming guide
126 (nvidia, 2009) recommends a minimum number of 64 blocks or a multiple
127 of this number. We will investigate the impact of the block size in the
128 performance section. Depending on the block size and the model size, we
129 might have some extra threads in the last block for which we do not need to
130 perform any calculations. We therefore have to check whether the offset is
131 smaller than the dimension of the model `nx*ny*nz`.

132 If the current thread is active, we calculate the indices of the current
133 cell in x-direction, y-direction and z-direction, respectively, from the offset
134 and the total size of the model in y-direction and z-direction. The function
135 `CalcGravBoxTerm` is a straightforward implementation of the geometric term
136 in Equation 2 and takes the three components of the measurement position,

137 the three coordinates of the upper left front corner of the current cell and
138 the sizes of the cell in the three coordinate directions as arguments. After
139 the API has executed the above code we have obtained a single row of the
140 sensitivity matrix.

141 The further computational strategy depends on the context in which the
142 calculation is performed. For pure forward modeling the most efficient ap-
143 proach is to perform a scalar multiplication between the current row of the
144 sensitivity matrix and the vector of densities on the GPU to obtain the cur-
145 rent datum and then discard the sensitivity information. In this case we min-
146 imize both the storage requirements and the number of transfers between the
147 memory of the GPU and the main memory. In an inversion context however
148 it is beneficial to store the sensitivity matrix, if possible, for two reasons.
149 First, as long as the geometry does not change we can calculate the data
150 for models with varying density distributions by a matrix-vector product as
151 shown in Equation 6. We will show the acceleration we can achieve with this
152 below. Second, we can use the sensitivity matrix to perform Gauss-Newton
153 type inversion. We therefore always transfer the current row of the sensitiv-
154 ity matrix from the GPU to main memory, then perform the scalar product
155 on the CPU and let the main application decide whether this row should be
156 stored for later use or discarded. In the performance section we will assess
157 the cost of the additional transfers.

158 The implementation for the gravimetric tensor is similar to the scalar
159 implementation. We only have to replace the calculation of the geometric
160 term with the appropriate mathematical expressions and preserve the 6 inde-
161 pendent rows of the sensitivity matrix when copying from the GPU to main

162 memory.

163 4. Performance

164 In this section we demonstrate the performance gain we can achieve with
165 our GPU based implementation. All tests were run on a Intel Q6600 with
166 2.4GHz, 4GB of main memory and a NVidia GTX260 graphics card which
167 has 192 processor cores and 896 MB of onboard memory with a bandwidth
168 of 111.9 GB/s. This is the cheapest graphics card that can handle double
169 precision computations that we use throughout the comparison and is readily
170 available in standard consumer PCs.

171 We compiled the main code with the GNU compiler collection version
172 4.3.3 under Ubuntu 09/04 using the “-O3” optimization flag and the GNU
173 openmp implementation. For the CUDA code we used NVidia’s `nvcc` in
174 Version 2.1 with the driver version 180.44. In all cases we average over
175 5 independent runs to obtain the calculation time. In each run we use a
176 different density model where each cell of the model is randomly assigned a
177 density between 0.1 – 3.0 g/cm³ and the cell sizes randomly vary between 1
178 and 11 km.

179 First, we examine the impact of the execution block size on the perfor-
180 mance. For three different model sizes we vary the number of threads per
181 execution block between 64 and 256. In Figure 3 we plot the time relative to
182 the fastest run for each model size in order to make the results for the three
183 model sizes comparable. For the chosen model and block sizes we observe
184 that the performance varies by only 17% between the fastest and the slowest
185 configuration. Depending on the model size 64, 128 or 256 threads per block

186 result in the highest performance. Between these three configurations the
187 maximum difference in performance is only 6%. We therefore choose a block
188 size of 128 for all subsequent experiments and do not attempt to optimize
189 this value.

190 Figure 4 shows the computation time for varying model sizes between 8
191 and 1 million model cells and 30 stations for computation of scalar gravity
192 data on one CPU core, 4 CPU cores and the graphics card, respectively. To
193 illustrate the benefits of storing the sensitivity matrix for later computations
194 we also show the time it takes to evaluate the matrix vector product using
195 the ATLAS linear algebra library (Whaley et al., 2001).

196 As expected, for a single core of the CPU the time increases linearly with
197 model size. There is very little overhead to the computation and profiling
198 shows that most time is spent evaluating the trigonometric and natural loga-
199 rithm functions in Equation 2. When using all 4 cores of the CPU we observe
200 that for models with less than 1,000 model cells there is some administrative
201 overhead associated with the parallelization. For larger models, however, we
202 achieve the same linear increase with model size. For these large models
203 the acceleration compared to a single core is close to the theoretical maxi-
204 mum of a factor of 4. This demonstrates that the problem can be efficiently
205 parallelized for multi-core architectures.

206 The curve for the GPU based computations shows some interesting be-
207 havior. For models with less than 100 cells the computation time is higher
208 than for both CPU based calculations. This demonstrates the overhead asso-
209 ciated with initializing the GPU and transferring data between main memory
210 and the memory of the graphics card. Furthermore, for less than 3,000 simul-

211 taneous threads the calculation time is independent of model size illustrating
212 the massively parallel architecture. For fewer than a few thousand model pa-
213 rameters we do not utilize all available computing units on the card. For
214 more than 10,000 parameters we again achieve a linear dependency of com-
215 putation time on the model size. Within the linear domain the acceleration
216 compared to a single core of the CPU is approximately a factor of 40. This is
217 a significant increase in performance that allows to calculate the response of
218 large models within a few seconds. In our case the number of measurement
219 sites is relatively low and therefore even the calculation time of 70 s at 10^6
220 model parameters for the single CPU core is not problematic, for large sur-
221 veys with hundreds of sites however the acceleration provided by the GPU
222 marks an important step.

223 Our performance comparison also shows the time for calculations with
224 pre-computed sensitivities as it could be done within a non-linear inversion,
225 e.g. when combining gravity with other data (Heincke et al., 2006). Given
226 enough RAM we only have to perform the full computation in the first it-
227 eration and can then benefit from the accelerated evaluation with the atlas
228 library. In this case the acceleration factor is 1,000 for large models. This
229 makes the calculation of the model response essentially instantaneous, but
230 requires large amounts of memory. The storage of the sensitivity matrix in
231 double precision requires $8 \times N \times M$ bytes which corresponds to about 240 MB
232 for our largest test case, but exceeds the memory of current computers for
233 larger models or more measurement sites.

234 The graph for the full tensor calculation in Figure 5 shows the same
235 general behavior as for the scalar data. Although we now calculate 6 elements

236 of the tensor, the calculation time only increases by a factor of two compared
237 to the scalar data. The reason for this is the simpler structure of the equations
238 for the elements of the tensor. This result also shows that the calculations on
239 the CPUs are essentially dominated by the evaluation of the mathematical
240 functions and not by memory transfers. As before we observe a nearly linear
241 increase of the calculation time with model size for the calculations with
242 one processor and, apart from some overhead for small models, also for four
243 processors.

244 The transition from constant calculation time to linear increase for the
245 GPU calculation again occurs at a model size of 3,000 parameters. This is
246 because we calculate the 6 elements of the tensor in strictly serial order. The
247 structure of the calculation in terms of parallelization is therefore the same
248 as for the scalar case. The acceleration through the GPU for the tensor case
249 is a factor of 30 compared to 1 processor of the CPU. Due to the simpler
250 structure of the equations and the larger amount of data we have to transfer,
251 the acceleration is not quite as high in this case as for the scalar case, but
252 still significant.

253 As the FTG calculations require the most transfers of sensitivity infor-
254 mation between the GPU and general memory, we use these calculations to
255 assess the cost of the memory transfers. For each independent element of the
256 gravimetric tensor, we transfer a row of the sensitivity matrix from the GPU
257 to the CPU. Profiling shows that for models with 10^6 parameters the code
258 only spends 1% of its time for these memory transfers and this behavior is
259 therefore not critical for the performance.

260 Finally, we examine the numerical precision of the results. Figure 6 shows

261 a histogram of the relative difference between the results from the CPU and
262 the GPU for FTG calculations with 5 random models with 30 sites each. The
263 histogram shows a clear peak around zero with most values concentrated be-
264 tween $-5 \cdot 10^{-13}$ and $5 \cdot 10^{-13}$, the minimum and maximum relative difference
265 are $-5 \cdot 10^{-10}$ and $1 \cdot 10^{-10}$, respectively. This shows that for practical pur-
266 poses the results are identical. Also the trace of the tensor agrees with the
267 theoretical value within numerical precision.

268 We also examine the possibility of performing the calculations in single
269 precision on the GPU. Until recently GPUs were only capable of single pre-
270 cision calculations and their performance is significantly higher for this type
271 of calculations. Compared to the double precision calculations we observe
272 an acceleration factor of roughly 4, more than 100 times faster than calcu-
273 lations on the CPU. However, the numerical precision is problematic. When
274 comparing the results to the double precision calculations in most cases the
275 relative difference stays below $1 \cdot 10^{-3}$, a satisfactory value for practical pur-
276 poses. However, more than 10% of the results show a relative difference of 0.1
277 or more, most likely due to accumulated rounding errors (Li and Chouteau,
278 1998). Such a difference impacts on the result of an inversion or the inter-
279 pretation of a forward model and thus is not acceptable for reliable forward
280 modeling.

281 5. Conclusions

282 The calculation of the scalar and tensorial forward response of large den-
283 sity models can be efficiently parallelized and accelerated by performing the
284 calculation on a standard consumer GPU. Our tests show that it is important

285 perform the calculations with double precision to obtain reliable results. In
286 this case we achieve accelerations of a factor of 40 for scalar data and a factor
287 of 30 for tensorial data with more than 3,000 model parameters, respectively.
288 For the tested cases the number of threads per execution block has only a
289 minor impact on the performance.

290 This is a significant improvement, particularly when considering the rel-
291 atively low cost of these graphics cards. Our approach allows to quickly
292 calculate the response for different density distributions as required, for ex-
293 ample, in a joint inversion without storing sensitivity information. Although
294 utilizing the sensitivity information accelerates the calculation further, even
295 modern computers cannot store the sensitivity matrix for large models. Fur-
296 thermore, even then we have to calculate the sensitivities once which can be
297 performed using the GPU based algorithm.

298 **6. Acknowledgments**

299 We thank Colin G. Farquharson and an anonymous reviewer for con-
300 structive comments on the manuscript. This work was funded by Chevron,
301 ExxonMobil, Nexen, RWE Dea, Shell, Statoil and Wintershall within the
302 JIBA consortium.

303 **References**

- 304 Chasseriau, P., Chouteau, M., 2003. 3D gravity inversion using a model of
305 parameter covariance. *Journal of Applied Geophysics* 52, 59–74.
- 306 Götze, H.-J., Lahmeyer, B., 1988. Application of three-dimensional interac-
307 tive modeling in gravity and magnetics. *Geophysics* 53, 1096–1108.

- 308 Heincke, B., Jegen, M., Hobbs, R., 2006. Joint inversion of mt, gravity and
309 seismic data applied to sub-basalt imaging. *SEG Technical Program Ex-*
310 *expanded Abstracts* 25 (1), 784–789.
- 311 Hobbs, R., Trinks, I., 2005. Gravity modelling based on small cells. In: *Pro-*
312 *ceedings of 67th EAGE (European Association of Geoscientists and Engi-*
313 *neers) Conference and Exhibition*. Madrid, Spain, p. 347.
- 314 Jeong, W.-K., Whitaker, R. T., 2008. A fast iterative method for eikonal
315 equations. *SIAM Journal on Scientific Computing* 30 (5), 2512–2534.
- 316 Komatitsch, D., Michéa, D., Erlebacher, G., 2009. Porting a high-order finite-
317 element earthquake modeling application to nvidia graphics cards using
318 cuda. *Journal of Parallel and Distributed Computing* 69 (5), 451 – 460.
- 319 Li, X., Chouteau, M., 1998. Three-dimensional gravity modeling in all space.
320 *Surveys in Geophysics* 19, 339–368.
- 321 Li, Y., Oldenburg, D. W., 1998. 3-D inversion of gravity data. *Geophysics*
322 63, 109–119.
- 323 Nagihara, S., Hall, S. A., 2001. Three-dimensional gravity inversion using
324 simulated annealing: Constraints on the diapiric roots of allochthonous
325 salt structures. *Geophysics* 66, 1438–1449.
- 326 Nagy, D., Papp, G., Benedek, J., 2000. The gravitational potential and its
327 derivatives for the prism. *Journal of Geodesy* 74, 552–560.
- 328 Nickolls, J., Buck, I., Garland, M., Skadron, K., 2008. Scalable parallel pro-
329 gramming with cuda. *Queue* 6 (2), 40–53.

- 330 nvidia, 2009. Nvidia CUDA programming guide, Version 2.1.
331 URL <http://www.nvidia.com/cuda>, [accessed 1 December 2009]
- 332 Portniaguine, O., Zhdanov, M. S., 1999. Focusing geophysical inversion im-
333 ages. *Geophysics* 64, 874–887.
- 334 Ryoo, S., Rodrigues, C. I., Bbagsorkhi, S. S., Stone, S. S., Kirk, D. B., Hwu,
335 W., 2008. Optimization principles and application performance evaluation
336 of a multithreaded gpu using cuda. In: *Proceedings of the 13th ACM SIG-*
337 *PLAN Symposium on Principles and Practice of Parallel Programming.*
338 Salt Lake City, UT, pp. 73–82.
- 339 Whaley, C. R., Petitet, A., Dongarra, J. J., 2001. Automated empirical opti-
340 mization of software and the ATLAS project. *Parallel Computing* 27 (1–2),
341 3–35.

Figure 1: Nomenclature and parameterization for gravity forward problem. Position of the measurement is described by the coordinate triple (x, y, z) . Model is divided into rectilinear blocks of constant density ρ , for clarity we only show a single block. Coordinates of corners of the block can be completely described by two coordinate triples (ξ_1, η_1, ζ_1) and (ξ_2, η_2, ζ_2) for opposing corners of the block.

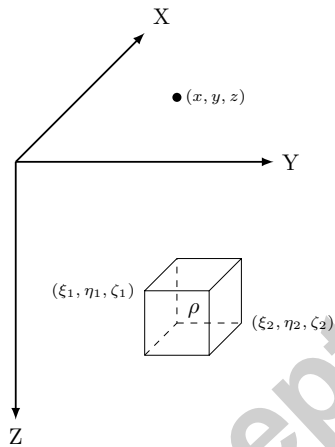
Figure 2: Overview of CUDA execution model and mapping of sensitivities. Execution grid consists of independent blocks that can be executed in any order. In turn each block consists of a number of threads. Each element of the sensitivity vector for the current measurement is mapped onto a different thread.

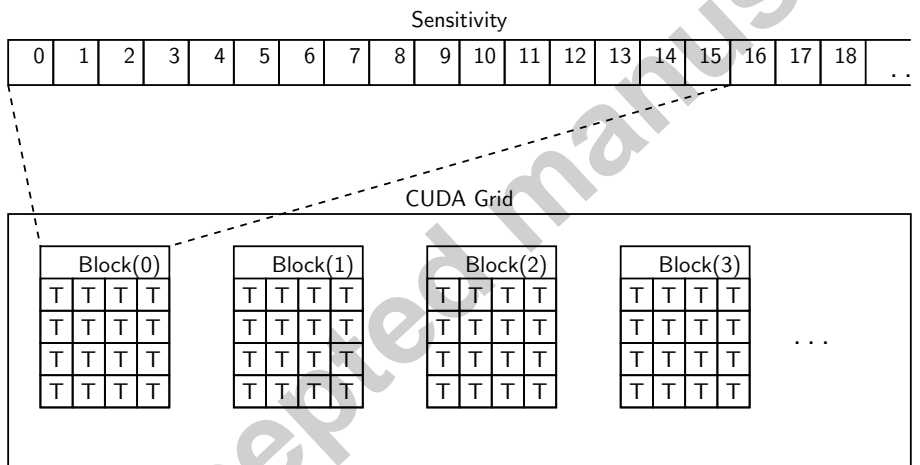
Figure 3: Dependency of execution time on number of threads per block. For each block size we measure execution time of models with $40 \times 40 \times 40$, $60 \times 60 \times 60$ and $80 \times 80 \times 80$ model cells, respectively. To make results comparable we divide by the time for the fastest execution for each model size. Execution time is relatively similar for all block sizes but shows minima at 64, 128, 192 and 256, respectively.

Figure 4: Calculation times for different size models for scalar gravity data for a single CPU thread (Q6600), 4 CPU threads and GPU (GTX260). For comparison we also show the time to evaluate the matrix vector product with the ATLAS library when the sensitivity matrix has been calculated.

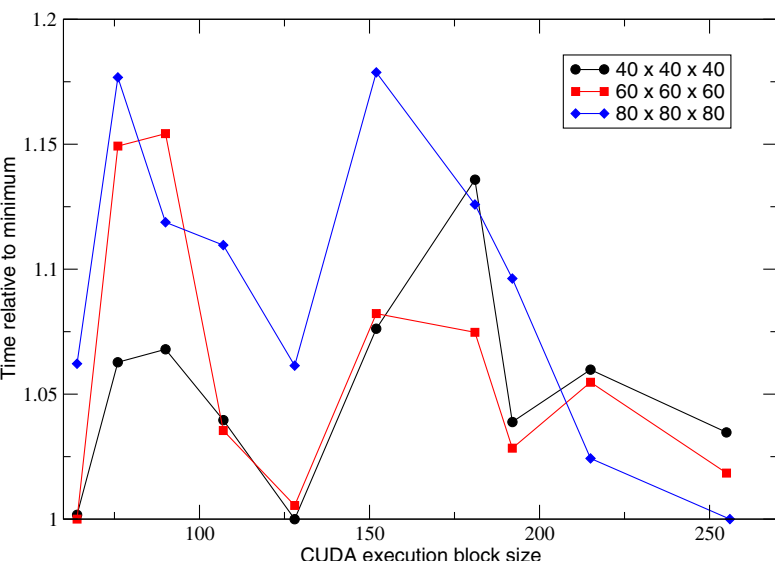
Figure 5: Calculation times for different size models for FTG data for a single CPU thread (Q6600), 4 CPU threads and GPU (GTX260).

Figure 6: Relative difference between FTG calculations performed on CPU and on GPU in double precision, respectively. Maximum relative deviation between results is $-5 \cdot 10^{-10}$.





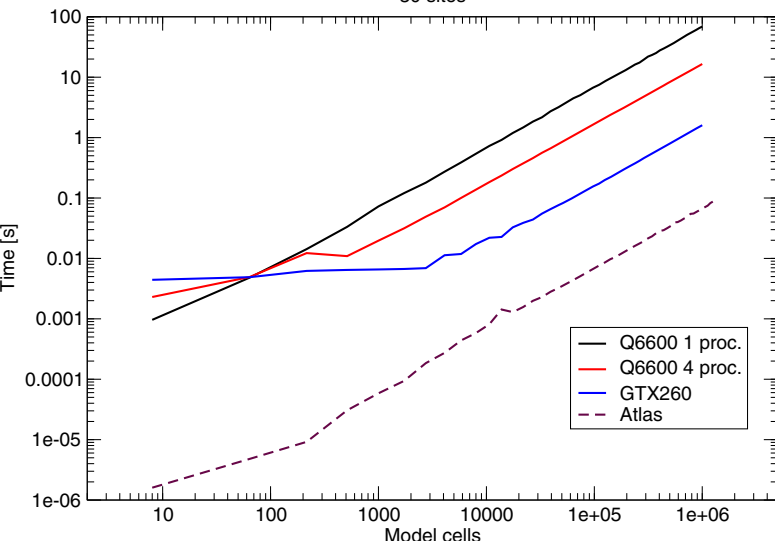
Accepted Manuscript



Accepted manuscript

Forward calculation for scalar data

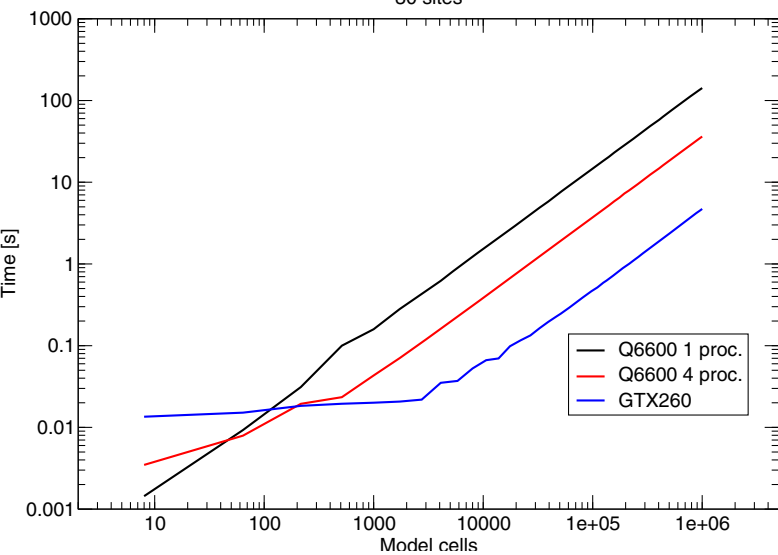
30 sites



Accepted manuscript

Forward calculation for FTG data

30 sites



Accepted manuscript

