

Simultaneous analysis of large *INTEGRAL/SPI*[☆] datasets: optimizing the computation of the solution and its variance using sparse matrix algorithms

L. Bouchet^{a,b,*}, P. Amestoy^c, A. Buttari^d, F.-H. Rouet^{c,e}, M. Chauvin^{a,b}

^aUniversité de Toulouse, UPS-OMP, IRAP, Toulouse, France

^bCNRS, IRAP, 9 Av. colonel Roche, BP 44346, F-31028 Toulouse cedex 4, France

^cUniversité de Toulouse, INPT-ENSEEIH-IRIT, France

^dCNRS-IRIT, France

^eLawrence Berkeley National Laboratory, Berkeley CA94720, USA

Abstract

Nowadays, analyzing and reducing the ever larger astronomical datasets is becoming a crucial challenge, especially for long cumulated observation times. The *INTEGRAL/SPI* X/γ-ray spectrometer is an instrument for which it is essential to process many exposures at the same time in order to increase the low signal-to-noise ratio of the weakest sources. In this context, the conventional methods for data reduction are inefficient and sometimes not feasible at all. Processing several years of data simultaneously requires computing not only the solution of a large system of equations, but also the associated uncertainties. We aim at reducing the computation time and the memory usage. Since the SPI transfer function is sparse, we have used some popular methods for the solution of large sparse linear systems; we briefly review these methods. We use the Multifrontal Massively Parallel Solver (MUMPS) to compute the solution of the system of equations. We also need to compute the variance of the solution, which amounts to computing selected entries of the inverse of the sparse matrix corresponding to our linear system. This can be achieved through one of the latest features of the MUMPS software that has been partly motivated by this work. In this paper we provide a brief presentation of this feature and evaluate its effectiveness on astrophysical problems requiring the processing of large datasets simultaneously, such as the study of the entire emission of the Galaxy. We used these algorithms to solve the large sparse systems arising from SPI data processing and to obtain both their solutions and the associated variances. In conclusion, thanks to these newly developed tools, processing large datasets arising from SPI is now feasible with both a reasonable execution time and a low memory usage.

Keywords: methods: data analysis, methods: numerical, techniques: imaging spectroscopy, techniques: miscellaneous, gamma-rays: general

1. Introduction

Astronomy is increasingly becoming a computationally intensive field due to the ever larger datasets delivered by observational efforts to map ever larger volumes and provide ever finer details of the Universe. In consequence, conventional methods are often inadequate, requiring the development of new data reduction techniques. The SPI X/γ-ray spectrometer, aboard the *INTEGRAL* observatory, perfectly illustrates this trend. The telescope is dedicated to the analysis of both point-sources and diffuse emissions, with a high energy resolution (Vedrenne et al., 2003). Its imaging capabilities rely on a coded-mask aperture and a specific observation strategy based on a dithering procedure (Jensen et al., 2003). After several years of operation, it also becomes important to be able to handle simultaneously all the data, in order, for example, to

get a global view of the Galaxy emission and to determine the contribution of the various emission components.

The sky imaging with SPI is not direct. The standard data analysis consists in adjusting a model of the sky and instrumental background to the data through a chi-square function minimization or a likelihood function maximization. The related system of equations is then solved for the intensities of both sources and background. The corresponding sky images are very incomplete and contain only the intensities of some selected sky sources but not the intensities in all the pixels of the image. Hence, images obtained by processing small subsets of data simultaneously cannot always be combined together (co-added) to produce a single image. Instead, in order to retrieve the low signal-to-noise ratio sources or to map the low surface brightness “diffuse” emissions (Bouchet et al., 2011), one has to process simultaneously several years of data and consequently to solve a system of equations of large dimension. Grouping all the data containing a signal related to a given source of the sky allows to maximize the amount of information on this specific source and to enhance the contrast between the sky and the background.

Ideally, the system of equations that connects the data to the

[☆]Based on observations with INTEGRAL, an ESA project with instruments and science data center funded by ESA member states (especially the PI countries: Denmark, France, Germany, Italy, Spain, and Switzerland), Czech Republic and Poland with participation of Russia and the USA.

*lbouchet@irap.omp.eu

sky model (where the unknown parameters are the pixels intensities) should be solved for both source intensities and variability timescales. This problem, along with the description and treatment of sources variability, is the subject of another paper (Bouchet et al., 2013).

It is mandatory, for example when studying large-scale and weak structures in the sky, to be able to process large amounts of data simultaneously. The spatial (position) and temporal (variability) description of sources leads to the determination of several tens of thousands of parameters, if ~ 6 years of SPI data are processed at the same time. Consequently, without any optimization, the systems to be solved can exceed rapidly the capacities of most conventional machines. In this paper we describe a technique for handling such large datasets.

2. Material and methods

2.1. The SPI spectrometer

SPI is a spectrometer provided with an imaging system sensitive both to point-sources and extended source/diffuse emission. The instrument characteristics and performance can be found in Vedrenne et al. (2003) and Roques et al. (2003). Data are collected thanks to 19 high purity Ge detectors illuminated by the sky through a coded-mask. The resulting Field-of-View (FoV) is $\sim 30^\circ$ and the energy ranges from 20 keV to 8 MeV. The instrument can locate intense sources with an accuracy of a few arc minutes (Dubath et al., 2005).

2.2. Functioning of the “spectro-imager” SPI

The coded mask consists of elements which are opaque (made of tungsten) or transparent to the radiation. Photons coming from a certain direction cast a shadow of the mask onto the detectors plane. The shadowgram depends on the direction of the source (Figure 1). The recorded counts rate in each detector of the camera is the sum of the contribution from all the sources in the FoV. The deconvolution consists of solving a system of equation which relates a sky model to the data through a transfer function. In the case of SPI, the imaging properties rely on the coded aperture, but also on a specific observing strategy: the dithering.

2.2.1. Dithering and sources variability

The reconstruction of all the pixels of the sky image enclosed in the FoV is not possible from a single exposure. Indeed, dividing the sky into $\sim 2^\circ$ pixels (the angular resolution), we obtain, for a 30° FoV, $\sim (30^\circ/2^\circ)^2 = 225$ unknowns. However, a single exposure contains only 19 data values which are the number of observed counts in the 19 Ge detectors and does not permit us to obtain the parameters necessary to determine the model of the sky and background. The related system of equations is thus undetermined. The dithering observation technique aims to overcome this difficulty.

By introducing multiple exposures for a given field that are shifted by an offset that is small compared to the size of the FoV, it is possible to increase the number of equations, by grouping

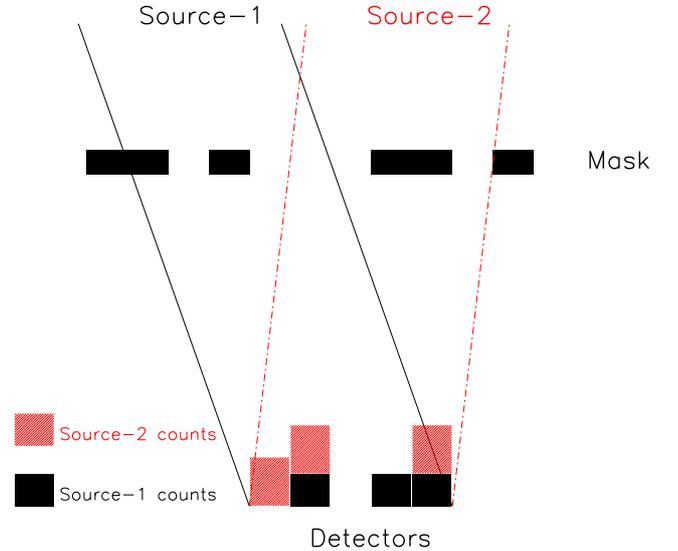


Figure 1: SPI imaging principle. The mask consists of elements transparent or opaque to the radiation. The opaque elements (made of tungsten) are shown in black. The shadowgram of the mask casts onto the detector plane (camera) depends on the source direction. Here the counts in the different detectors of source-1 and source-2 are shown in black and red. The counts recorded by the detectors are the sum of all the contributions from all the sources in the FoV.

exposures, until the system becomes determined and thus solvable. An appropriate dithering strategy (Jensen et al., 2003) has been used where the spacecraft continuously follows a dithering pattern throughout an observation. In general, the pointing direction varies around a target by steps of 2° within a five-by-five square or a seven-point hexagonal pattern. A pointing (exposure) lasts between 30 and 60 minutes. Thus, the dithering allows to construct a solvable system of equations.

However, in addition to the variable instrumental background, sources are also variable on various timescales ranging from hours (roughly the duration of an exposure) to years. This is not a major problem at high energy ($E \gtrsim 100$ keV), since there are only few emitting sources, whose intensities are rather stable in time with respect to the statistics. At lower energies ($E \lesssim 100$ keV) and in most cases, the source intensities vary during the time spanned by the all the exposures. The chi-square, of the associated least-square problem, for this group can be relatively high, if sources intensity variations are not taken into account. In spite of this, it is possible to include a model of the source intensity variations in the formulation of the problem and to re-optimize the system of equations accordingly (Bouchet et al., 2013). Nevertheless, including sources variability in the system of equations increases the number of unknowns to determine (2.2.3) since intensities, in each “time-bin” (a segment of time where the intensity of a given source does not change statistically), are to be determined simultaneously along with the parameters which model the instrumental background.

2.2.2. Cases where it is better to process large amount of data simultaneously

It is impossible from a single exposure (19 data values) to obtain the sky image in the 30° FoV; only a coarse image containing at most 19 sources can be obtained. This coarse image is under-sampled and contains information on only 19 pixels (there is no information on the other pixels). Hence, images cannot always be combined together (co-added) to produce a single image. Furthermore, crowded regions like the Galactic Center contain hundreds of sources and thus a single exposure cannot provide the amount of information needed, even to build only a coarse image. The grouping of the exposures, by selecting all those that contain signal on a specific sky target, can provide the necessary information. The FoV spanned by these exposures is large (30° to 60°) and contains numerous sources.

2.2.3. Problem formulation

The signal (counts and energies) recorded by the SPI camera on the 19 Ge detectors is composed of contributions from each source in the FoV plus the background. For n_s sources located in the field of view, the data D_{dp}^{raw} obtained from detector d during an exposure (pointing) p , for a given energy band, can be expressed by the relation:

$$D_{dp}^{raw} = \sum_{j=1}^{n_s} R_{dp,j} I_{p,j}^s + B_{dp}^{bg} + \epsilon_{dp} \quad (1)$$

where $R_{dp,j}$ is the response of the instrument for source j (function of the source direction relative to the telescope pointing axis), $I_{p,j}^s$ is the flux of source j during pointing p and B_{dp}^{bg} the background both recorded during the pointing p for detector d . ϵ_{dp} are the measurement errors on the data D_{dp}^{raw} , they are assumed to have zero mean, to be independent and normally distributed with a known variance σ_{dp} ($\epsilon_{dp} \sim N(0, [\sigma_{dp}^2])$) and $\epsilon_{dp} = \sqrt{D_{dp}^{raw}}$.

For a given pointing p , D_{dp}^{raw} , $R_{dp,j}$, and B_{dp}^{bg} are vectors of n_d (say $n_d = 19$ detectors¹) elements. For a given set of n_p exposures, we have a system of $n_p \times n_d$ equations (Eq. 1). To reduce the number of free parameters related to background, we take advantage of the stability of relative counts rate between detectors and rewrite the background term as:

$$B_{dp}^{bg} = I_p^{bg} \times U^d \times t_{dp} \quad (2)$$

where I_p^{bg} is a normalization coefficient per pointing related to the background intensity, U^d is a background count rate pattern (uniformity map) on the SPI camera for detector d , and t_{dp} the effective observation time for pointing p and detector d . The number of parameters necessary to model the background reduces to n_p if U is assumed to be known². However, in some cases it can be determined while processing the data

¹The number of functioning detectors could be $n_d = 16, 17, 18$ or 19 for single events and up to 141 , when all the multiple events are used in addition to the single events (Roques et al., 2003).

²Derived from ‘‘empty-field’’ observations (Bouchet et al., 2010).

(Appendix A.4).

The two extreme cases, in terms of number of parameters to be determined, are

- First, when the sources and background intensities are assumed to be constant throughout all the observation (time spanned by the exposures), the relation between the data and the sky model can be written, omitting the detector indices, as

$$D_p^{raw} = \sum_{j=1}^{n_s} R_{pj} I_j^s + P_p I^{bg} + \epsilon_p \quad (3)$$

$$\text{with } P_p = t_{dp} \times U^d$$

The aim is to compute the intensities $I^s(j)$ of the n_s sources and the background relative intensity I^{bg} . Therefore, the relation can be written in matrix form, as

$$\begin{pmatrix} P_1 & R_{1,1} & \dots & R_{1,n_s} \\ P_2 & \ddots & \ddots & R_{2,n_s} \\ \vdots & \ddots & \ddots & \vdots \\ P_{n_p} & \dots & \dots & R_{n_p,n_s} \end{pmatrix} \begin{pmatrix} I^{bg} \\ I_1^s \\ I_2^s \\ \vdots \\ I_{n_s}^s \end{pmatrix} = \begin{pmatrix} D_1^{raw} \\ D_2^{raw} \\ \vdots \\ D_{n_p-1}^{raw} \\ D_{n_p}^{raw} \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_{n_p-1} \\ \epsilon_{n_p} \end{pmatrix}$$

We can rewrite the system in a more compact form as

$$y = H_0 x + \epsilon \text{ or } y_i = \sum_{j=1}^{n_s+1} h_{ij} x_j + \epsilon_i \text{ for } i = 1, \dots, M \quad (4)$$

where H_0 (elements h_{ij}) is an $M \times (n_s + 1)$ matrix and $M = n_d \times n_p$. The parameters to be determined, $x = (I^{bg}, I_1^s, \dots, I_{n_s}^s)$ is a vectors of length $n_s + 1$. The data $y = (D_1^{raw}, D_2^{raw}, \dots, D_{n_p}^{raw})$ and the associated statistical errors $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_{n_p})$ are vectors of length M .

- Second, if the background or the sources are variable on the exposure timescale, the number of unknowns (free parameters) of the set of $n_p \times n_d$ equations is then $(n_s + 1) \times n_p$ (for the n_s sources and the background intensities, namely I^s and I^{bg}). This leads, unless the number of sources is small, to an underdetermined system of equations.³

Fortunately, in real applications, many sources vary on timescales larger than the exposure. This leads to a further reduction of the number of parameters compared to the case where all sources vary on the exposure timescale. In addition, many point sources are weak enough to be considered as having constant flux within the statistical errors, especially for higher energies ($E \gtrsim 100$ keV). Then the $n_p \times n_s$ parameters related to sources will reduce into N_s^{eff} parameters and, similarly, N_b for the background. As these parameters have also a temporal connotation, they will hereafter be referred to as ‘‘time-bins’’.

³With the Compressed Sensing approach (Bobin et al. (2008); Wiaux et al. (2009) and references therein), it is possible to find a sparse solution even if the system is underdetermined and for systems in which the matrix is sparse.

If the source named or numbered J is variable, then the total duration covered by the n_p exposures is divided into K_J sub-intervals where the source intensity can be considered as stable/constant regarding the data statistics. The solution x_J is expanded in K_J segments, it takes the value “time-bins” s_k^J in segment k , and can be written in compact notation

$$x_j = \sum_{k=1}^{K_J} s_k^J I_k^J \text{ with } \begin{cases} I_k^J = 1 & \text{if } t \in [t_{k-1}^J, t_k^J[\\ I_k^J = 0 & \text{otherwise} \end{cases}$$

Actually the instants t_k^J correspond to the exposure acquisition time (exposure number), with $t_0=1$ and $t_k^J = n_p + 1$. There is at least one and at most n_p time segments for each source J

($x_J = [s_1^J, \dots, s_{K_J}^J]$ becoming a vector of length K_J). The matrix H_0 (eq. 4) is to be modified accordingly.

When expanding matrix H_0 , column J is expanded in K_J new columns, hence the number of intensities (unknowns) increases. Schematically H_0 ($M \times (n_s + 1)$) is mapped into a matrix H ($M \times N$), N being the sum of all sources intervals ($N = \sum_{J=0}^{n_s} K_J$), that is the number of “time-bins” (the index $J=0$ correspond to the background). Matrix $H(1 : M, 1 : K_0)$ is related to the background while $H(1 : M, K_0 + 1 : N)$ is related to the sources response. Parameters $x(1 : K_0)$ and $x(K_0 + 1 : N)$ are related to background and source intensity variations with the time (number of exposures). Box I illustrates schematically how the matrix H is derived from the matrix H_0 .

$$H_0 = \begin{pmatrix} h_{11} & h_{12} & h_{13} & \dots & h_{1N} \\ h_{21} & h_{22} & h_{23} & \ddots & h_{2N} \\ h_{31} & h_{32} & h_{33} & \ddots & h_{3N} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ h_{M1} & h_{M2} & h_{M3} & \dots & h_{MN} \end{pmatrix}$$

$$\mapsto H = \begin{pmatrix} h_{11} & 0 & 0 & 0 & h_{12} & 0 & 0 & h_{13} & 0 & h_{1N} & \dots & 0 \\ 0 & h_{21} & 0 & 0 & h_{22} & 0 & 0 & 0 & h_{23} & h_{2N} & \dots & 0 \\ 0 & 0 & h_{31} & 0 & 0 & h_{32} & 0 & 0 & h_{33} & 0 & \dots & h_{3N} \\ \vdots & \ddots & \vdots \\ 0 & 0 & 0 & h_{M1} & 0 & 0 & h_{M2} & 0 & h_{M3} & 0 & \dots & h_{MN} \end{pmatrix}$$

Finally, the relation between the data and the sky model, similarly as in eq. 4, is

$$Hx = y + \epsilon \quad (5)$$

Physically, H corresponds to the transfer function or matrix, y to the data and x to the unknown intensities (sources plus background) to be determined (a vector of length N).

Taking into account the variability of sources and instrumental background increases the size of the system of equation and the number of unknowns, but also increases the sparsity of the matrix related to the system of equations, which means that the underlying matrices have very few non-zero entries. In our application, the matrix H_0 is sparse, thus matrix H is even sparser. Objective methods to construct the matrix H from H_0 are described in Bouchet et al. (2013).

To give an idea, for the dataset which corresponds to (~ 6 years of data, the number of free parameters $N = N_s^{eff} + N_b$ to be determined are between $N \sim 5\,000$ and $N \sim 90\,000$ depending on the energy band considered and hypotheses made on sources and background variability timescale (2.3).

2.3. Material

The material is related to the analysis of data accumulated between 2003 and 2009 with the spectrometer SPI. The astrophysical application is the study of diffuse emission of the Galaxy. The details can be found in Bouchet et al. (2011). The goal is to disentangle the “diffuse” emission (modeled with 3 spatially extended components) from the point-sources emission and instrumental background. We need to sample these large-scale structures efficiently over the entire sky and consequently use the maximum amount of data simultaneously, since a single exposure covers only one-hundredth of the total sky area. The datasets consist of 38 699 exposures that yield $M = 649\,992$ data points. In most cases considered here, the background intensity is considered to be quite stable on a ~ 6 hours timescale, which corresponds to $N_b \simeq 5870$ unknowns.

- The highest energy bands ($E \gtrsim 100$ keV) are less problematic in terms of number of parameters to determine, as illustrated by the 200–600 keV band. The sky model contains only 29 sources which are essentially not variable in time (given the instrument sensitivity). The number of unknowns is $N \simeq 5900$.
- The lowest energy bands ($E \lesssim 100$ keV) are more problematic. We use the 25–50 keV band. The sky model contains 257 sources variable on different timescales. When the background intensity is assumed to vary on ~ 6 hours timescale, $N \simeq 22\,500$ “time-bins” intensity are to be determined. In some configurations, essentially used to assess the results, background intensity and/or strongest variable sources vary on the exposure timescale, and the number of unknowns could be as high as $N \simeq 55\,000$ to $N \simeq 90\,000$. Nevertheless, the matrices associated with these problems remain relatively structured.
- To avoid excessively structured matrices, we generate also matrices H , with a variable number of columns, the number of segments K_J for a particular source being a random number between 1 and n_p . This results in a different number of parameters N .

Another astrophysical application is the study of a particular source or sky region, here the crowded central region of the Galaxy. In this case, it is possible to use a smaller number of exposures. We use 7147 exposures which cover a sky region of radius 30° around the Galactic center. We measure the intensity variations of a set of 132 sources. The number of parameters to determine $N = 3\,578$ is relatively small. Details can be found in Bouchet et al. (2013). A second matrix, used for verification purposes, has $N = 9\,437$. It corresponds to the case where some sources are forced to vary on shorter timescales.

The material consists of rectangular matrices H and symmetric square matrices A ($A = H^T H$) related to the above physical

problems (2.2.3). The characteristics of some of these matrices are described in Table 1.

The system we use in the experiments consists of an Intel i7-3517U processor with 8 GB main memory. We ran the experiments on a single core, although our algorithms are amenable to parallelism.

Table 1: Sparsity of matrices H and $H^T H$.

N	$\rho(H)(\%)$	$\rho(A)(\%)$	
3578	2.67	2.96	Central Galaxy (27-36 keV)
9437	1.01	1.05	
5900	0.12	0.13	Diffuse emission 200-600 keV
22503	0.18	0.28	Diffuse emission 25-50 keV
55333	0.07	0.09	
149526	0.03	0.04	Simulation (25-50 keV)

$\rho(\text{Matrix})$ is the so-called *density* of the matrix: the ratio between the number of non-zero elements in the matrix and the total number of elements in the matrix ($M \times N$ for H and N^2 for $A = H^T H$, where M is the number of rows of H). The matrix H arising from the diffuse emission study have $M = 672\,495$ rows. The number of non-zero elements is constant $n_z = 27\,054\,399$ for the matrices with $N \geq 22\,503$ corresponding to the 25-50 keV band and $n_z = 4\,677\,821$ for the 200-600 keV band. The matrix with $N = 3\,578$ has $M = 124\,921$ rows and $n_z = 11\,948\,840$ non-zero elements.

2.4. Methods

The mathematical problem described in Section 2.2.3 and developed in 2.4.1 requires the solution of several algebraic equations. First, if the chi-square statistic is used, a linear least-squares problem has to be solved to estimate the parameters of the model. Second, elements (entries) of the inverse of a matrix have to be computed in order to determine the error bars (variances of these parameters). Third, in some cases, the parameters are adjusted to the data through a multi-component algorithm based on likelihood tests (Poisson statistics); this problem leads to a non-linear system of equations (Appendix A.1).

These three problems can be reduced to solving a linear system with a square matrix: a linear least-squares problem $\min_x \|Hx - y\|$ can be transformed into a square system $Ax = b$ by use of the normal equations⁴ ($A = H^T H$ and $b = H^T y$). Similarly, computing entries of the inverse of a matrix amounts to solving many linear systems, as described in detail in Section 3.3.1. For the above mentioned non-linear problem, we chose a Newton-type method; this involves solving several linear systems as well. Our problems are large, but sparse (cf. Table 1), which justifies the use of sparse linear algebra techniques. In Section 3.1, we describe how we selected a method suitable for our application.

2.4.1. The least-square solution (LSQ)

The system is, in most cases, overdetermined (there are more equations - or measures here - than unknowns), therefore there

⁴For clarity, we omit to weight the matrix H and the data by the inverse of the data standard deviation, see Section 2.4.1

is (generally) no exact solution, but a “best” solution, motivated by statistical reason, obtained by minimizing the following merit function, which is the chi-square⁵:

$$\chi^2 = \sum_{i=1}^M \left[\frac{y_i - \sum_{j=1}^N h_{ij} x_j}{\sigma_i} \right]^2 \quad (6)$$

$y = (y_1, \dots, y_M)$ is vector of length M representing the data, $[\Sigma]$ a diagonal matrix of order M whose diagonal is $(\sigma_1, \dots, \sigma_M)$, where σ_i is the measurement error (standard deviation) corresponding to the data point y_i . These quantities are assumed to be known (formally $\sigma_i = \sqrt{y_i}$). $H = h_{ij}$ is a matrix of size $M \times N$. The least-square solution $x = (x_1, \dots, x_N)$ is obtained by solving the following normal equation:

$$(H^T [\Sigma^{-2}] H)x = H^T [\Sigma^{-2}] y \text{ or as } Ax = b \quad (7)$$

Once the solution has been computed, the uncertainties on the estimated solution x are needed as well. The corresponding variance can be obtained by computing the diagonal of A^{-1} :

$$\text{Var}(x_i) \propto a_{i,i}^{-1} \text{ where } a_{i,i}^{-1} \text{ refers to } (A^{-1})_{i,i} \quad (8)$$

3. Theory

3.1. Processing large datasets: efficient solution of large sparse systems of equations

Sparse matrices appear in numerous industrial applications (mechanics, fluid dynamics, ...), and the solution of sparse linear systems has been an active field of research since the 1960s. Many challenges still arise nowadays, because industrial applications involve larger and larger number of unknowns (up to a few billions nowadays), and because hardware architectures are increasingly complex (multi-core, multi-GPU, etc.).

Exploiting sparsity can significantly reduce the number of operations and the amount of memory needed to solve a linear system. Let us take the example of a partial differential equation to be solved on a 2D physical domain; the domain can be discretized on a $k \times k$ 2D grid and using, say, finite differences, the equation can be transformed into a sparse linear system with $N = k \times k$ unknowns. Without exploiting sparsity, this system would be solved in $O(N^3)$ operations (using an exact method), with a memory usage in $O(N^2)$. It has been shown that, for this particular case, the number of arithmetic operations can be reduced to $O(N^{3/2})$, and space complexity to $O(N \log N)$ by exploiting the sparsity of the matrix (Hoffman et al., 1973).

Many methods exist for solving sparse linear systems (Duff et al., 1989; Saad, 1996). Two main classes can be distinguished: *direct methods*, that rely on a matrix factorization (e.g., $A = LU$), and *iterative methods*, that build a sequence of iterates that hopefully converges to the solution.

⁵The number of counts per detector is high enough to use the Gaussian statistics.

Direct methods are known to be numerically robust but often have large memory and computational requirements, while iterative methods can be less memory-demanding and often faster but are less robust in general. Iterative methods often need to be *preconditioned*, i.e., to be applied to a modified system $M^{-1}Ax = M^{-1}b$ for which the method will converge more easily; a trade-off has to be found between the cost of computing and using the preconditioner M and how the preconditioner improves the convergence. The choice of a method is often complicated and strongly depends on the application. In our case, we choose to use a direct method for the following reasons:

- Memory usage is often a bottleneck that prevents the use of direct methods, but with the matrices arising from our application, direct and iterative methods have roughly the same memory footprint. This is explained in the next section.
- The matrices from our application are numerically challenging; we found that unpreconditioned iterative methods (we tried GMRES) have difficulties converging and that a direct method that does not implement robust numerical features is also likely to fail (we illustrate this in Section 5).
- We need to compute error bars, which amounts to solving a large number ($O(N)$) of linear systems with different right-hand sides but the same matrix. This is particularly suitable for direct methods; indeed, once the matrix of the system is factored (e.g., $A = LU$), the factors can be reused to solve for different right-hand sides at a relatively inexpensive cost. We describe this in Section 3.3.1.

In this work, we use the *MUMPS* (Multifrontal Massively Parallel Solver) package. *MUMPS* (Amestoy et al., 2001, 2006) aims at solving large problems on parallel architectures. It is known to be very robust numerically, by offering a large variety of numerical processing operations, and provides a large panel of features. In the following section, we briefly describe how sparse direct methods work. We introduce the basic material needed to understand the algorithm used for the computation of error bars (described in Section 3.3.1).

3.2. Sparse direct methods

Direct methods are commonly based on Gaussian elimination, with the aim to factorize the sparse matrix, say A , of the linear system $Ax = b$ into a product of “simpler” matrices called *factors*. Typically, A can be factored into $A = LU$ where L and U are lower and upper triangular matrices respectively, or $A = LDL^T$, where D is a diagonal matrix if A is symmetric (which is the case in our application).

Sparse direct methods depend on the non-zero pattern of the matrix and are optimized in that sense; specialized mathematical libraries for tridiagonal, banded, cyclic matrices are common. If the pattern is more complex, then the method usually consists of three phases: *analysis*, *factorization* and *solution*.

3.2.1. Analysis

The *analysis* phase applies numerical and structural preprocessing to the matrix, in order to optimize the subsequent phases. One of the main preprocessing operations, called *re-ordering*, aims at reducing the *fill-in*, i.e., the number of non-zero elements which appear in the factors but do not exist in the initial matrix; this step consists in permuting the rows and columns of the initial matrix in such a way that less fill-in will occur in the permuted matrix. Table 2 shows the amount of fill-in for different problems coming from our astrophysical application when the matrices are permuted using the nested-dissection method. For each matrix, the number of non-zero elements in the original matrix A and in the L factor of the LDL^T factorization of A are reported. Note that in our application, the fill-in is not very large: the number of non-zero elements in the factors is of the same order of magnitude as in the original matrix. As a result, the use of sparse, direct methods is likely to provide a good scalability with respect to the size of the matrix produced by the application. Moreover, this implies that, for our application, direct and iterative methods will have roughly the same memory requirements; indeed, in an unpreconditioned iterative method, the memory footprint is mainly due to the storage of the matrix A , while the major part of memory requirements of direct methods comes from the factors. Note that, while our application exhibit low amount of fill-in, this not the case in other applications; in many problems, especially those involving PDEs on 3D physical domains, the number of non-zero coefficients in the factors can be as big as one hundred times more than in the original matrix. In this case, using an iterative method can be interesting memory-wise.

Matrix size	3578	9437	22503	55333	149526
$nz(A)$	378475	932143	1436937	2705492	9379127
$nz(L)$	519542	1380444	2885821	9189447	14432264

Table 2: Number of non-zeros in the original matrix A and in the L factor of the $A = LDL^T$ factorization for different problems of our experimental set.

An important step of the analysis phase is the *symbolic factorization*: this operation computes the non-zero pattern of the factors, on which the numerical factorization and the solution will rely. The symbolic factorization computes the structure of the factors by manipulating graphs, and also a structure called the *elimination tree*, a tree-shaped graph with N vertices. This tree represents tasks dependencies for the factorization and the solution phases. We describe in more details the elimination tree since it is a key structure to explain and understand (see Section 3.3.1) how to accelerate the solution phase since computing entries in the inverse of the matrix corresponds to incomplete traversals of the elimination tree. Figure 2(b) shows an elimination tree and we use it to illustrate some definitions: one of the nodes is designated to be the *root*; in the figure, this is node 6. For our purpose, the root is the node corresponding to the variable of the linear system that is eliminated last. An *ancestor* of a vertex v is a vertex on the path from v to the root. The *parent* (or *father*) of v is its first ancestor; all the nodes but the root have a parent. For example, on the figure, nodes 6 and 5 are ancestors of 4; 5 is the parent of 4. A *child* of a vertex v is

a vertex of which v is the parent. For example, 4 and 3 are the children of 5. A vertex without children is called a *leaf*; 1 and 2 are leaves. *Descendants* of a vertex v are all the nodes in the subtree rooted at v ; for example, 1, 2, 3 and 4 are descendants of 5.

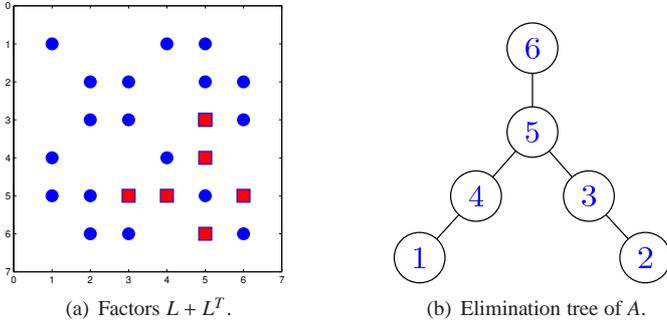


Figure 2: The factors and the elimination tree of a symmetric matrix A . (a) pattern of the $L + L^T$ factors of A with filled-in entries shown with squares, (b) the elimination tree of A where the children of a node are drawn below the node itself.

In the following subsections (*factorization* and *solution phase*), we describe briefly how a sparse direct solver uses elimination trees; we will also rely on this notion in Section 3.3.1 for the computation of error bars. Further details about the construction and the role of elimination trees in sparse solvers are given in Liu (1990).

3.2.2. Factorization

After the preprocessing performed during the analysis phase, the numerical factorization takes place and the matrix A is transformed into a product of factors (e.g., LU). The factorization consists in traversing the elimination tree following a *postorder*, that is a *topological ordering* (i.e. each parent is visited after its children) where the nodes in each subtree are visited consecutively. In Figure 2(b), 1-4-2-3-5-6 is, for example, a postorder. At each node, a partial factorization of a dense matrix is performed. Note that nodes that belong to different branches can be processed independently, which is especially useful in a parallel setting.

The factorization phase tries to follow as much as possible the preparation from the analysis phase, but sometimes, because of numerical issues (typically, division by a “bad pivot”, i.e. a very small diagonal entry that could imply round-off errors), it has to adapt dynamically: the structure of the factors and the scheduling of the tasks can be modified on the fly.

3.2.3. Solution phase

Once the matrix has been factored, the linear system is solved. For example, in the case of the LU factorization, the system $Ax = b$ becomes $LUx = b$ and is solved in two steps (two solutions of triangular systems):

$$\begin{cases} z = L^{-1}b & \text{“Forward substitution”} \\ x = U^{-1}z & \text{“Backward substitution”} \end{cases}$$

The forward substitution follows a bottom-up traversal of the elimination tree as in the factorization, while the backward substitution follows a top-down traversal of the tree. At each node, one component of the solution is computed, and some updates are performed on the dependent variables (ancestor nodes for the forward phase, descendant nodes for the backward phase).

3.3. Computing error bars: partial computation of the inverse of a sparse matrix

In our astrophysical application, once the solution, either for the linear or non-linear problem, has been found, it is necessary to compute the variances of the parameters of the fitted function. In the case of multiple regressions such as least squares problems, the standard deviation of the solution can be obtained by inverting the Hessian or covariance matrix. However, since the inverse of a sparse matrix is structurally full, it is impractical to compute or store it (Duff et al., 1988). In our case, the whole inverse of the covariance matrix is not required: since we only want the variances of the parameters (not their covariances), we only need to compute the diagonal elements of the inverse.

Some work has been done since the 1970s in order to compute a subset of elements of the inverse of a sparse matrix. One of the first works is Takahashi et al. (1973) which has been extended in Campbell & Davis (1995); this approach relies on a direct method (i.e. on a factorization). An iterative method has been proposed in Tang & Saad (2009) for matrices with a decay property. Some methods have also been developed for matrices arising from specific applications; a more detailed survey is given in Amestoy et al. (2010). Many of these methods provide sophisticated ideas and interesting performance on specific problems, but no software package is publicly available, with the exception of the approach implemented within MUMPS solver, that we describe in the next section.

3.3.1. MUMPS A^{-1} feature

The A^{-1} feature in MUMPS has been described in (Slavova, 2009) and was motivated by the *INTEGRAL/SPI* application, among other applications that require the computation of inverse entries, or, more generally, applications that involve sparse right-hand sides (as explained in this section). This feature is able to compute any set of entries of A^{-1} , relying on a traditional solution phase, i.e. by computing every required entry a_{ij}^{-1} as:

$$a_{ij}^{-1} = (A^{-1}e_j^\dagger)_i$$

Using the LU factors of A , this amounts to solving two triangular systems:

$$\begin{cases} Lz = e_j^\dagger \\ a_{ij}^{-1} = (U^{-1}z)_i \end{cases}$$

The first triangular system in the equation above is particular because its right-hand side e_j^\dagger is very sparse (only one non-zero entry). Furthermore, we do not need the whole solution of the second triangular system, but only one component. This

information can be exploited to reduce the traversal of the elimination tree; while a regular solution phase would consist in visiting the whole elimination tree twice (a bottom-up traversal followed by a top-down traversal), computing a_{ij}^{-1} consists in two partial traversals of the tree: the first triangular system is solved by following the path from node j to the root node, and the second triangular system is solved by following the path from the root node to node i ; this is referred to as *pruning* the elimination tree. Since each node of the tree corresponds to operations to be performed (arithmetic operations, or expensive accesses to the factors in the out-of-core case), this leads to significant improvements in computation time. Moreover, since we do not have to manipulate dense solution vectors, this also leads to significant savings in memory usage.

We illustrate this technique in Figure 3: entry a_{23}^{-1} is required, thus the only nodes of tree that have to be visited lie on the path from node 3 to the root node (6) and on the path from the root node to node 2. Therefore, one does not have to perform operations at nodes 4 and 1.

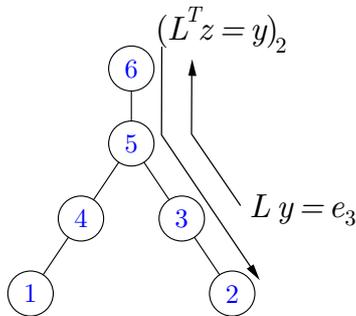


Figure 3: Computation of a_{23}^{-1} . The traversal of the tree is reduced to the path from 3 to 6 and the path from 6 to 2; no computation is performed at nodes 1 and 4.

When many entries of the inverse are requested, they cannot generally be computed all at once (mainly because of memory usage), but they can be computed by blocks, which allows to take advantage of efficient dense linear algebra kernels. Work has been performed in order to find optimal ways to form the blocks in different contexts (Amestoy et al., 2010) and to improve the parallel efficiency.

4. Calculation

A substantial time is spent in computing $A = H^T H$ with a basic algorithm. The use of an appropriate algorithm to perform the operation $A = H^T H$ helps to reduce the computation time (see Section 4.1). The MUMPS solver is used to solve the system of equations as described in Section 4.2. Finally, the error bars on the solution are computed, which means the calculation of the diagonal elements of inverse matrix. The new A^{-1} feature of MUMPS is compared with several algorithms, in terms of computation time in Section 4.3.

4.1. Improvements of the computation of $A = H^T H$

The computation of the normal equation $A = H^T H$ is of paramount importance in many problems, yet is a very chal-

lenging operation due to the considerable amount of symbolic operations needed to compute the sparsity structure of A . For this reason efficient algorithms have been developed in the past. To perform this operation we decided to use part of a larger code developed by Puglisi (1993) for computing the QR factorization of sparse matrices. The used part was originally developed to compute only the structure of the A matrix and, thus, we had to extend it in order to compute the coefficients values. This was possible thanks to the help of the original code developer.

One important feature of this code offers the possibility to update the elements of A that are changed after modification of some numerical values of the columns of H without recomputing the whole matrix (the technique used to compute simultaneously the solution and the background pattern in the algorithm is described in 2.2.3 and Appendix A.4).

Table 3: Time for the computation of $A = H^T H$

Matrix	Improved algorithm ^a	Simple algorithm ^b
$A = H^T H$		
$N = 22\,503$		
Full matrix	28.2	5 779
5000 H columns modified	0.43	258
$N = 149\,526$		
Full matrix	41.2	18 585
5000 H columns modified	0.13	31.7

Times are in seconds. H is an N by M matrix; here $M = 672\,495$ and H has 27 054 399 non-zero elements. ^a Based on an original package from Puglisi (1993) and improved as suggested by the author. ^b N matrix vector product are used following the Compressed Column Storage scheme, but for each operation a dense vector of length M (with many zero element), that represents a column of H is built in place.

Table 3 shows the time reduction achieved for both the computation of $A = H^T H$ and its update after the modifications of some columns of H . The results in the first column are obtained with the code extracted from the software package by Puglisi (1993) and improved as suggested by the author. The results on the second column, instead, are obtained by computing N matrix vector products where, for each product, a dense vector of length M (with many zero elements) corresponding to a column of H is built in place.

The gain over a simple basic algorithm is significant (a factor ~ 300) and demonstrates the interest of using specialized libraries dedicated to sparse matrix computations.

4.2. Solving a sparse linear system

Here we briefly illustrate the interest of exploiting sparsity of the matrix when solving a linear system. In Table 4, we compare the time for solving linear systems arising from our application using a dense solver (LAPACK (Andersen et al., 1990)) and a sparse solver (MUMPS). Times are in seconds and include the LDL^T factorization of a symmetric matrix A of order N and the computation of the solution x of the system $Ax = b$ (x and b are vectors of length N). In the results related to MUMPS, the time for the analysis phase is included. In the second row

of the table, instead, the matrix is treated as dense, hence its full storage is used and no analysis phase is performed. For the largest two problems, the dense algorithm cannot be used as the memory requirements are roughly 23 GB and 167 GB respectively. We can extrapolate that on this system, the run time would be around 22 hours for the largest problem (instead of 6.7 seconds using a sparse algorithm).

Table 4: Times (in seconds) for the computation of the solution.

Matrix size	3578	9437	22503	55333	149526
Sparse	0.2	0.7	1.6	8.0	6.7
Dense	1.2	20.1	169.9	/	/

The results in Table 4 confirm that sparse, direct solvers achieved a good scalability on the problems of our target application whereas dense linear algebra kernels quickly exceed the limit of modern computing platforms.

4.3. Time to compute error bars

In this section we present experimental results related to the computation of error bars or, equivalently, of the diagonal entries of the inverse matrix A^{-1} . Our approach that relies on the pruned tree, presented in Section 3.3.1, is compared to the basic, left-looking approach described in (Stewart, 1998). In the case of a symmetric matrix, this approach computes the diagonal entries of the inverse matrix as

$$a_{ii}^{-1} = \sum_{k=1}^N \frac{1}{d_{kk}} (l_{ki}^{-1})^2$$

where we denoted with a_{ij}^{-1} and l_{ij}^{-1} the coefficients of A^{-1} and L^{-1} , respectively. This amounts to computing, one at a time, the columns of L^{-T} and then summing the corresponding contribution onto the a_{ii}^{-1} coefficients. In this algorithm, the sparsity of the right-hand side and of the factor matrix L is exploited but not completely, and the experimental results discussed below show that this results in a higher execution time. Furthermore, because of memory issues, this simple algorithm does not allow to simultaneously compute many diagonal entries of A^{-1} ; clearly this is also a limiting factor for performance. Our implementation of this method is based on the LDL package (Davis, 2005). As a second term of comparison we also provide experimental results for a brute force approach with no exploitation of sparsity of the right-side and solution vectors. For this purpose, we use directly the MUMPS package and solve several systems of equations in order to compute the inverse matrix. In addition, we analyze the influence of grouping the computation of the diagonal entries (1 right-hand-side (RHS) at a time or 128 at a time).

The experimental results for the three methods described above are reported in Table 5. For the sake of this comparison, all these methods are executed in sequential mode although the code of the brute force approach and of the MUMPS A^{-1} feature are parallel. The experiments were carried out on the above-mentioned system.

These results show that the brute force algorithm becomes competitive with respect to the simple algorithm when the entries are processed by blocks. The MUMPS A^{-1} feature described in 3.3.1 is significantly faster than all other approaches and the gain increases with the size of the problem. Pruning leads to clear gains over a strict traditional solution phase. The gain is even larger for the largest problem due to the good scalability of the A^{-1} algorithm with respect to the problem size. The simple, left-looking approach shows reasonable performance for small problems, but could not be tested on our largest matrix because numerical pivoting, not available in LDL package, is needed during factorization to obtain an accurate solution.

Table 5: Time to compute the diagonal elements of the inverse of a symmetric matrix.

Matrix size	3578	9437	22503	55333	149526
Left-looking	28.2	376.1	2567.9	489.1	/
MUMPS (1 RHS)	3.77	38.4	204.1	1324.9	8230.5
MUMPS (128 RHS)	1.32	7.34	45.5	245.6	2833.5
MUMPS A^{-1}	0.28	0.9	4.9	36.0	9.5

Execution times (in seconds) for the computation of all the diagonal entries of the A^{-1} matrix with the left-looking, brute force and MUMPS A^{-1} methods. For the brute force approach results are provided for blocks of size 1 and 128.

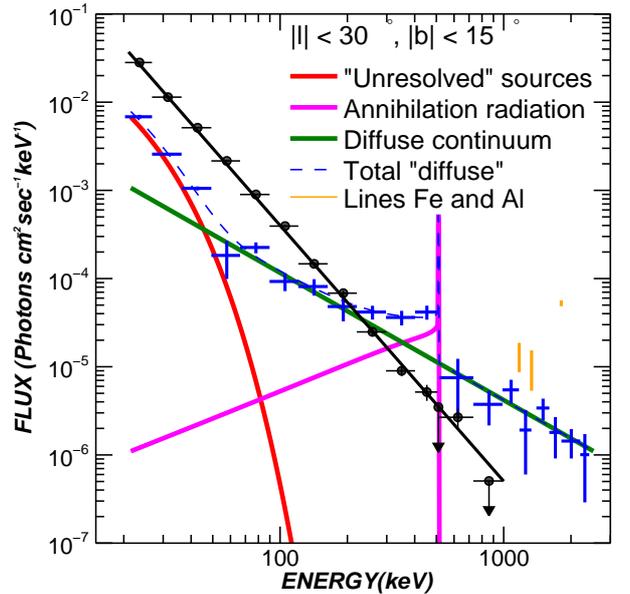


Figure 4: Different contributions to the total emission at hard X-ray and soft gamma-ray energies in the central radian of the galaxy. The data points shown in black (plus filled circle) correspond to the contribution due to 270 point sources. The average spectrum of these sources can be viewed at <http://sigma-2.cesr.fr/integral/>. The data points shown in blue correspond to the diffuse emission.

5. Results and discussion

The MUMPS solver and its A^{-1} functionality are the core tools to solve systems of equations related to the measurements

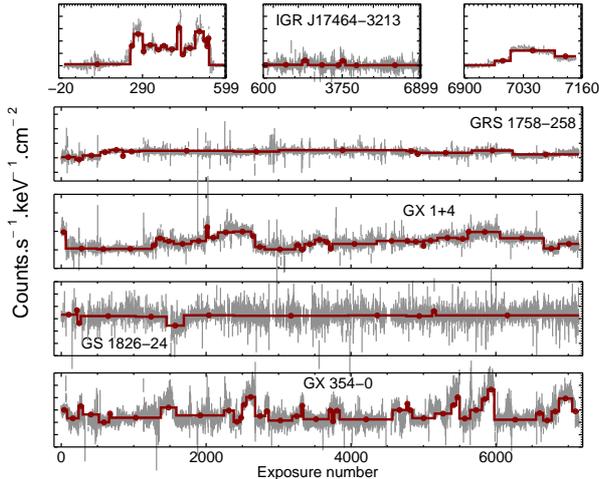


Figure 5: Intensity evolutions (Red) of IGR J17464-3213, GRS 1758-258, GX 1+4, GS 1826-24 and GX 354-0 in the 27-36 keV band. These intensity variations are compared to the time-series (“Quick-look” analysis) obtained with the IBIS instrument (Ubertini et al., 2003) aboard the *INTEGRAL* observatory. The time-series (30-40 keV) is shown in gray.

of the sources intensity. Figure 4 shows the application to the determination of the different components of the Galaxy spectrum. The related analysis is performed in 24 consecutive energy bands in order to extract counts spectra. The counts spectra are then converted into photon spectra. The details are given in Bouchet et al. (2011). Another application is the study of the intensity variations of a peculiar source or sky region. Figure 5 shows the intensity in function of the time (exposure) of some of the sources located in the central crowded region of the Galaxy. For this application, the end and start of the “time-bins” are determined by a segmentation algorithm, which is based on the efficient LDL^T factorization of symmetric matrix provided by MUMPS, details can be found in Bouchet et al. (2013).

We have demonstrated that even for the basic operation such as sparse matrix product, it is better to use dedicated algorithms or libraries (4.1).

The MUMPS solver is very effective on the sparse matrix structure arising from astrophysical problems encountered with SPI. This solver is robust and the matrix factorization is stable against rounding errors. It provides many numerical preprocessing options and implements robust pivoting strategies, which make it one of the most numerically stable solvers available. The matrices arising from the *INTEGRAL/SPI* application are symmetric and indefinite; they are not

extremely challenging numerically, but they do require two-by-two numerical pivoting for stability (in Table 5 the LDL package could not provide an accurate solution on our largest matrix). The proposed approach not only leads to substantial time reduction but also enables the resolution of large sparse system of equations which could not be solved using basic algorithms.

Other sparse linear systems solvers exists and have been used

to validate the performance reported in the experimental section; for an exhaustive list see for example Bai et al. (2000), but they all lack a function to compute also the error bars on the solution quickly, which is mandatory in our astrophysical application.

The A^{-1} feature in MUMPS (computation of selected inverse entries) did not exist before the beginning of this study, the *INTEGRAL/SPI* application was actually one of the motivating applications for developing techniques for the computation of inverse entries, and for releasing a publicly available code. This functionality allows to compute easily and rapidly the error bars on the solution. The gain in time over already optimized algorithms is important.

Among other methods to solve the problem completely, solution and error bars, one should mention alternative methods such as Monte Carlo Markov Chains (Metropolis et al., 1953; Hastings, 1970; Neal, 1993) or Simulated Annealing (Kirkpatrick et al., 1983). Such advanced statistical tools can provide the best fit and the variances of the solution of both linear and non-linear systems of equations. In particular MCMC methods could be useful when computing error bars, in case of complex constraints on the function. However, these methods may be very prohibitive in time, especially if high precision on the parameters is required; they have in general a weak or non-guaranteed convergence and are not the best suited for our needs, given the complexity of our problem.

6. Conclusions

We have developed algorithms to process years of data and to enhance the production of *INTEGRAL* hard X/soft γ -ray survey catalogs. These methods have been successfully applied to a set of ~ 6 years of data (Bouchet et al., 2011). We have shown that, for processing efficiently and accurately years of data, it is critical to use algorithms that take advantage of the sparse structure of the transfer function (matrix), such as those implemented in the MUMPS software⁶. It was also demonstrated that error bars can be obtained at a relatively inexpensive cost (the same order of magnitude as a simple problem solution) thanks to a recently developed algorithmic feature that efficiently computes selected entries of the inverse of a matrix. In addition, thanks to many efforts in optimization, gains are achieved both in memory usage and in computation time. Hence, it is possible to process large datasets in a reasonable time and to compute the diagonal of the covariance matrix, and thus error bars, in a rather short time. More generally, the ideas described here can be applied to a large variety of problems. Finally, we are today able to solve sparse linear systems with millions, sometimes billions, of unknowns. Although we have not used explicitly parallel computing but instead performed many sequential computations at the same time (for each energy band, etc.), the proposed approach can also be used directly in a parallel setting on massively parallel machines.

⁶Available at <http://mumps.enseeiht.fr/> or <http://graal.ens-lyon.fr/MUMPS/>

In the near future, instruments will commonly create datasets of a few tens to a few hundreds of Terabytes for a single observation project. Many of the current tools and techniques for managing large datasets will not scale easily to meet this challenge. Surveys of the sky already require parallel computing in order to be performed. New surveys will demand orders of magnitude increases in the available data and therefore in data processing capabilities. It is also a challenge for scientists who need to extract a maximum of science from the data. Exciting scientific breakthroughs remain to be achieved as astronomers manipulate and explore massive datasets, but they will require advanced computing capabilities, infrastructure and algorithms.

Acknowledgments

The *INTEGRAL/SPI* project has been completed under the responsibility and leadership of CNES. We are thankful to ASI, CEA, CNES, DLR, ESA, INTA, NASA and OSTC for support. We are very grateful to Chiara Puglisi, research engineer at INPT-IRIT, for her contribution to improve the performance of the $H^T H$ product in Section 4.1.

Appendix A. Schematic solution of the system of equations

Appendix A.1. Maximum Likelihood Estimator (MLE) of the solution

In the case of a low number of counts, it is recommended to use the MLE of the solution instead of the χ^2 solution. Following Cash (1979), we maximize the likelihood function,

$$L = -2 \times \left(\sum_{i=1}^M e_i - y_i \ln e_i \right) \quad (\text{A.1})$$

where e_i is the model of the data obtain through the relation $e = Hx$.

Appendix A.2. Optimization of the non-linear problem

To optimize this non-linear problem, potentially with bound constraints (such as positivity of the solution), there are at least four approaches:

- (a) Newton type methods (or modified Newton methods): they use the Hessian matrix to define a search direction and hence need the solution of a large linear system of equations at least at each few iterations. They are the most powerful methods available and can find the solution in a few iterations.
- (b) Quasi-Newton methods: they build an approximation of the Hessian at each iteration. They optimize a quadratic function in at most n iterations (n being the number of unknowns).
- (c) Conjugate-gradient methods: unlike the Newton-type and quasi-Newton methods, conjugate gradients methods do not require the storage of an n by n Hessian matrix and thus are suited to solve large problems. The gradient of the

function (a vector of length n) is used to define the search direction. They are not usually as reliable or efficient as Newton-type methods and often a relatively large number of iterations has to be performed before obtaining an acceptable solution.

- (d) Simplex (Nelder & Mead, 1965), simulated annealing (Kirkpatrick et al., 1983) or Monte Carlo Markov Chain (MCMC) (Neal, 1993) can also be considered, but they are often prohibitive in time.

Methods (a) and (b) are known as order-2 optimization methods (gradient and Hessian used), (c) as an order-1 optimization method (gradient used), while method (d) can use only the function value.

To use a Newton type method (order-2), we need to compute the gradient G and the Hessian H_{ess} of the function

$$G_j = \frac{\partial L}{\partial x_j} = 2 \times \sum_{i=1}^M H_{ij} \left(1 - \frac{d_i}{e_i} \right) \text{ for } j = 1, \dots, N \quad (\text{A.2})$$

and $H_{ess} = \frac{1}{2} \frac{\partial^2 L}{\partial^2 x} = H^T \left[\frac{d}{e^2} \right] H$

$\left[\frac{d}{e^2} \right]$ is a diagonal matrix of order M whose diagonal is $\left(\frac{d_1}{e_1^2}, \dots, \frac{d_M}{e_M^2} \right)$. As for the LSQ case, the variance of the solution is obtained thanks to the inversion of the Hessian matrix (note that in the limit $\lim_{e_i \rightarrow d} \frac{d_i}{e_i^2} = \frac{1}{d_i} = \frac{1}{\sigma_i^2}$, the likelihood (H_{ess}) and chi-square (A) Hessian matrices are the same). A guess solution to this non-linear optimization problem is the LSQ solution.

Appendix A.3. Codes for non-linear optimization

The fitting algorithm, based on the likelihood test statistic, is a non-linear optimization problem. To optimize a non-linear problem, potentially with bound constraints, a Newton type method, known for its efficiency and reliability can be used, as we already have a solver for large sparse systems at hand. A software package for large-scale non-linear optimization such as IPOPT⁷ (Interior Point OPTimizer) can be used. IPOPT uses a linear solver such as MUMPS or MA57 (Duff & Red, 2004) as a core algorithm. For more details on this large-scale non-linear algorithm, see Wächter & Biegler (2006). A few similar software packages for large-scale non-linear optimization exist, among them LANCELOT (Conn et al., 1996), MINOS (Murtagh et al., 1982) and SNOPT (Gill et al., 1997).

Appendix A.4. “Empty-field” auto-computation

Sometimes the “empty-field” or “uniformity map” U has to be computed with the solution. In order to preserve the linearity of the problem, we have adopted the algorithm described below. We consider that if the solution x is known,

$$y_i = \sum_{j=1}^{K_0} h_{ij} x_j + \sum_{j=K_0+1}^N h_{ij} x_j = y_i^B + y_i^S \quad (\text{A.3})$$

⁷IPOPT is available at <https://projects.coin-or.org/Ipopt>

Coming back to the detector and pointing number

$$D_{dp}^{raw} = D_{dp}^S + U^d I_p^{bg} t_{dp} \quad (\text{A.4})$$

In the above formula $y_i^S \equiv D_{dp}^S$ is the counts due to the sources, assumed to be known. $y_i^B \equiv U^d I_p^{bg} t_{dp}$ is the background contribution, I_p^{bg} is assumed to be known and U is to be estimated. At this stage, using the model of the sky described by A.4, a rough estimate of the pattern is $U^d \simeq \frac{\sum_{p=1}^{n_p} (D_{dp}^{raw} - D_{dp}^S)}{\sum_{p=1}^{n_p} t_{dp}}$.

Appendix A.4.1. Expression for the detector pattern

For the LSQ statistic, we wish to minimize the following quantities for each of the working detectors,

$$\chi^2(d) = \sum_{p=1}^{n_p} \left(\frac{D_{dp}^{raw} - D_{dp}^S - U^d I_p^{bg} t_{dp}}{\sigma_{dp}} \right)^2 \text{ for } d = 1, \dots, n_p \quad (\text{A.5})$$

The LSQ solution $U^{LSQ}(d)$ is

$$U^{LSQ}(d) = \frac{\sum_{p=1}^{n_p} (D_{dp}^{raw} - D_{dp}^S) \times I_p^{bg} t_{dp} / \sigma_{dp}^2}{\sum_{p=1}^{n_p} (I_p^{bg} t_{dp})^2 / \sigma_{dp}^2} \quad (\text{A.6})$$

For the MLE statistic, we do not have to preserve the linearity of the problem and hence the computation of the improved “empty-field” pattern can be done during the non-linear optimization process. On another side, the algorithm is simplified if we proceed similarly as in the LSQ case. Then, we wish to maximize the following quantities for each of the working detectors,

$$L(d) = -2 \left(\sum_{p=1}^{n_p} D_{dp}^S + U^d I_p^{bg} t_{dp} - D_{dp}^{raw} \ln [D_{dp}^S + U^d I_p^{bg} t_{dp}] \right) \quad (\text{A.7})$$

for $d = 1, \dots, n_p$

The MLE solution $U^{MLE}(d)$ is

$$U^{MLE}(d) = \frac{\sum_{p=1}^{n_p} (D_{dp}^{raw} - D_{dp}^S)}{\sum_{p=1}^{n_p} I_p^{bg} t_{dp}} \quad (\text{A.8})$$

One should mention that it is possible to compute, similarly, an “empty-field” pattern on some restricted time interval instead of the whole dataset; the best “empty field” for pointing intervals p_k to p_{k+1} is then,

$$\begin{cases} U^{MLE}(d, k) = \frac{\sum_{p=p_k}^{p_{k+1}} (D_{dp}^{raw} - D_{dp}^S)}{\sum_{p=p_k}^{p_{k+1}} I_p^{bg} t_{dp}} \\ U^{LSQ}(d, k) = \frac{\sum_{p=p_k}^{p_{k+1}} (D_{dp}^{raw} - D_{dp}^S) \times I_p^{bg} t_{dp} / \sigma_{dp}^2}{\sum_{p=p_k}^{p_{k+1}} (I_p^{bg} t_{dp})^2 / \sigma_{dp}^2} \end{cases} \quad (\text{A.9})$$

Appendix A.5. “Empty-field” schematic construction

A sub-optimal algorithm to obtain both the sources and the background fluxes, as well as the improved “empty-field” pattern is described in Algorithm 1. We start with an approximation U_0 and apply some iterative refinement. In practice, the algorithm converges in a few iterations.

Algorithm 1 Computation of the “Empty field”, the solution and its variance

- 1: $U = U_0$, compute the structure of the Hessian (A or H_{ess})
 - 2: **for** $i=1$ **to** itermax **do** {Iterative computation of U and x }
 - 3: Compute LSQ or MLE solution
 - 4: Compute a new approximation of U by minimizing again LSQ or maximizing MLE statistics
 - 5: Update H (The first K_0 columns of H and update the new Hessian matrix (Sec. 4.1))
 - 6: If χ^2 stops decreasing or the likelihood function stops increasing then go to step 8
 - 7: **end for**
 - 8: Compute H at the solution (if not already done) and the diagonal of H^{-1} to obtain the uncertainties on the solution
-

References

- Amestoy, P. R., Duff, I. S., Koster, J. & L’Excellent, J.-Y., 2001, “A fully asynchronous multifrontal solver using distributed dynamic scheduling”, *SIAM Journal of Matrix Analysis and Applications* 23 (1): 15–41.
- Amestoy, P. R., Guermouche, A., L’Excellent, J.-Y. & Pralet, S., 2006, “Hybrid scheduling for the parallel solution of linear systems”, *Parallel Computing* 32 (2): 136–156.
- Amestoy, P. R., Duff, I. S., L’Excellent, J.-Y., Robert, Y., Rouet, F.-H. & Uçar, B., 2012, “On computing inverse entries of a sparse matrix in an out-of-core environment”, *SIAM journal on Scientific Computing* 34 (4): 1975–1999.
- Anderson, E., Bai, Z., Dongarra, J., Greenbaum, A., McKenney, A., Du Croz, J., Hammerling, S., Demmel, J., Bischof, C & and Sorensen, D., 1990, “LAPACK: A portable linear algebra library for high-performance computers”, *Proceedings of the 1990 ACM/IEEE conference on Supercomputing*, 2–11.
- Bai, Z., Demmel, J., Dongarra, J., Ruhe, A. & van der Vorst, H. 2000, “Templates for the Solution of Eigenvalue Problems: A Practical Guide”, *SIAM*, Philadelphia, 2000.
- Bobin J., Starck J.-L. & Ottensamer R., 2008, *IEEE Sel. Top. Signal Proc.*, 2, 718
- Bouchet, L., Roques, J. P. & Jourdain, E., 2010, *ApJ*, 720, 177.
- Bouchet, L., Strong, A., Porter, T.A, et al., 2011, *ApJ*, 739, 29.
- Bouchet, L., Amestoy, P. R., Buttari, A., Rouet, F.-H. & Chauvin, M., 2013, A&A, in press (<http://dx.doi.org/10.1051/0004-6361/201219605>)
- Campbell, Y. E. & Davis, T. A., 1995, “Computing the sparse inverse subset: an inverse multifrontal approach”, *Tech. Rep. TR-95-021*, CISE Dept., Univ. of Florida.
- Cash, W., 1979, *ApJ*, 228, 939.
- Conn, A. R., Gould, I. M., & Toint, P. L., 1996, “Numerical experiments with the LANCELOT package (Release A) for large-scale nonlinear optimization”, *Mathematical Programming*, 73(1), 73–110.
- Davis, T. A., 2005, “User guide for LDL, a concise sparse Cholesky package”, *Tech. Rep.*, CISE Dept., Univ. of Florida.
- Dubath, P., Knödseder, J., Skinner, G. K., et al., 2005, *MNRAS*, 357, 420.
- Duff, I. S., Erisman, A. M., Gear, C. W. & Reid, J. K., 1988, “Sparsity structure and Gaussian elimination”, *ACM SIGNUM Newsletter* 23 (2), 2–8.
- Duff, I. S., Erisman A. M. & Reid, J. K., 1989, “Direct methods for sparse matrices”. Oxford University Press.
- Duff, I. S. & Reid, J. K., 2004, “MA57 - A code for the solution of indefinite sparse symmetric linear systems”, *ACM Transactions on Mathematical Software*, 30, 118.

- Gill, P. E., Murray, W., & Saunders, M. A., 1997, "SNOPT: an SQP algorithm for large-scale constrained optimization". Technical Report SOL97-3, Department of EESOR, Stanford University, Stanford, California 94305, USA, 1997.
- Hastings, W. K., 1970, *Biometrika* 57 (1): 97–109.
- Hoffman, A. J., Martin, M. S. & Rose, D. J., 1973, "Complexity bounds for regular finite difference and finite element grids", *SIAM Journal on Numerical Analysis* 10 (2): 364–369.
- Jensen, P. L., Clausen, K., Cassi, C., et al., 2003, *A&A*, 411, L7.
- Kirkpatrick, S., Gelatt, C. D. & Vecchi, M. P., 1983, *Science, New Series*, 220 (4598), 671–680.
- Liu, J. W. H., 1990, "The role of elimination trees in sparse factorization", *SIAM Journal on Matrix Analysis and Applications*, 11 (1): 134–172.
- Metropolis, N., Rosenbluth, A.W., Rosenbluth, M.N., Teller, A.H. & Teller, E., 1953), *Journal of Chemical Physics* 21 (6): 1087–1092.
- Murtagh, B. A., & Saunders, M. A., 1982, "A projected Lagrangian algorithm and its implementation for sparse nonlinear constraints", *Mathematical Programming Study* 16 (Constrained Optimization), 84-117.
- Neal, R. M., 1993, "Probabilistic Inference Using Markov Chain Monte Carlo Methods", Technical Report CRG-TR-93-1, Department of Computer Science, University of Toronto.
- Nelder, J., & Mead, R., 1965, *Computer Journal*, vol. 7, no 4, 308–313.
- Puglisi, C., 1993, "QR factorization of large sparse matrix overdetermined and square matrices with the multifrontal method in a multiprocessing environment", PhD thesis, Institut National Polytechnique de Toulouse, Toulouse, France.
- Roques, J. P., Schanne, S., Von Kienlin, A., et al., 2003, *A&A*, 411, L91.
- Saad, Y., 1996, "Iterative methods for sparse linear systems": PWS Publications.
- Slavova, T., 2009, "Parallel triangular solution in an out-of-core multifrontal approach for solving large sparse linear systems", PhD thesis, Institut National Polytechnique de Toulouse, France.
- Stewart, G. W., 1998, "Matrix algorithms", SIAM Press, Pennsylvania, 1998.
- Takahashi, K., Fagan, J., and Chen, M.S., 1973, "Formation of a sparse bus impedance matrix and its application to short circuit study", in *Power Industry Computer Applications Conference*, 63–69.
- Tang, J., & Saad, A., 2009, "A probing method for computing the diagonal of the matrix inverse", Tech. Report umsi-2010-42, Minesota Supercomputer Institute, University of Minnesota, Minneapolis, MN, 2009.
- Ubertini, P., Lebrun, F., Di Cocco, G., et al., 2003, *A&A*, 411, L131.
- Vedrenne, G., Roques, J. P., Schonfelder, V., et al., 2003, *A&A*, 411, L63.
- Wächter, A. and Biegler, L. T., 2006, *Mathematical Programming*, 106(1), 25–57.
- Wiaux Y., Jacques L., Puy G.m Scaife A. M. M., Vanderghenst P, 2009, *A&A*, 395, 1733

Appendix B. Highlights

- *INTEGRAL/SPI X*/ γ -ray spectrometer data analysis
- Large astronomical data sets arising from the simultaneous analysis of years of data.
- Resolution of a large sparse system of equations; solution and its variance.
- The Multifrontal Massively Parallel Solver (MUMPS) to solve the equations.
- MUMPS A^{-1} feature to compute selected inverse entries (variance of the solution, . . .).