



HAL
open science

Wavelet analysis of X-ray spectroscopic data Part I. The method

P. Bury, N. Ennode, Jean-Marc C. Petit, Ph. Bendjoya, J.-P. Martinez, H. Pinna, J. Jaud, J.-L. Ballador

► To cite this version:

P. Bury, N. Ennode, Jean-Marc C. Petit, Ph. Bendjoya, J.-P. Martinez, et al.. Wavelet analysis of X-ray spectroscopic data Part I. The method. Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment, 1996, 383 (2-3), pp.572-588. <10.1016/S0168-9002(96)00721-8>. <hal-05471811>

HAL Id: hal-05471811

<https://hal.science/hal-05471811v1>

Submitted on 11 Feb 2026

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

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



Distributed under a Creative Commons CC BY-SA 4.0 - Attribution - ShareAlike - International License

Wavelet analysis of X-ray spectroscopic data. I: the method

P. Bury^{a,d}, N. Ennode^b, J.-M. Petit^c, Ph. Bendjoya^{e,d},
J.-P. Martinez^b, H. Pinna^b, J. Jaud^b and J.-L. Balladore^b

^a*Laboratoire de la Physique de la Matière Condensée, Université de Nice
Sophia-Antipolis, Avenue de Valrose, 06000 Nice Cedex, France*

^b*CNRS, CEMES/LOE, U.R.A. A8011, 29, rue Jeanne Marvig, B.P.4347, 31055
Toulouse Cedex*

^c*CNRS, U.R.A. 1362, Observatoire de la Côte d'Azur, B.P. 229,
06304 Nice Cedex 4, France*

^d*CNRS, U.R.A. 1360, Observatoire de la Côte d'Azur, B.P. 229,
06304 Nice Cedex 4, France*

^e*Département d'Astrophysique, Université de Nice Sophia-Antipolis, Avenue de
Valrose, 06000 Nice Cedex, France*

Abstract

We present a new filtering method based on wavelet analysis. This method works without any assumption about the noise, and we show that it gives reliable results on X-ray spectroscopic data. The originality of this method consist in the wavelet analysis: the filtering is done both in frequency space and in the spatial space, thus avoiding artifacts due to fourier filtering. We also show the reliability of the method using test signals.

All the software developped following this method is freely available on the network.

Key words: Wavelet analysis, signal processing, spectroscopy, X-ray diffraction

PROOF. Send proofs to P. Bury, B.P. 229, 06304 Nice cedex 4, France

1 Introduction

The knowledge of the matter chemical composition and its structure organisation in every scale and in every state, is a continuous quest. The X-ray diffraction allows us to own a precise matter structure image. To date, the technical ways to obtain these informations are photographs, position-sensitive proportional counters, and semiconductor detectors. The different methods present different advantages: good spatial resolution for the photographs, conviviality for the computer associated detector, fast acquisition for the diodes, ... They also present different disadvantages: high operational cost for the position-sensitive proportional counters, long exposure time for photographs, rather low signal to noise ratio for the semiconductor detectors, ... We tested four different devices (a photographic film, a proportional counter, a photodiode array and a silicon CCD imager) to define a processing method that avoid as much as possible the different disadvantage, while keeping as much as possible the advantages. During test the need of a powerfull processing method appeared. This led us to develop a new signal processing method, where classical methods (fourier filtering, adaptative filtering) have shown their inability to reduce noise and to keep the relevant part of the signal (i.e. position and amplitude of diffraction peaks).

The purpose of this paper is to present this new signal processing method based on the wavelet transform. This method allows us to detect, count peaks with a drastic decrease of noise. It also allows to give a level of significance to the detected structures against noise fluctuations, and also to estimate the error on the detected features. Since this method is essentially local, the main limitation of this method occurs when one is trying to detect weak structures close to strong ones. An iterative process of analysis is then needed. The method is not based on any modelling of the noise (Gaussian, Poissonian ...), and gives, as it will be shown, satisfying results even for very low signal-to-noise ratio signals.

In section 2 we briefly describe the physical experiment used, and also the detectors in use. We will also compare these experiments to a classical one in the same field, and explain the need of such accurate and efficient post-processing. In section 3 we briefly recall the principles of the wavelet transform. Then we describe the structure detection method and the reconstruction algorithm. In section 4 we apply this method on simulated signals. We show that human intervention is needed to split the signal in regions with large amplitude structures and regions with small amplitude structures. From the reconstructed signal, we extract some physically significant parameters such as the amplitude and location of peaks. In section 5, we apply this new technic to two different real signals obtained using photodiode arrays and CCD's and compare them with more complex experiments.

2 Experiment description

Classical experimental conditions usually make use of a film which is exposed for more than one or two days to an X-ray beam in the range 5 keV to 20 keV. The precise energy of the X-ray beam is chosen depending on the physical properties of the crystal under study. In order to reduce drastically (factor of a few hundreds) the exposure time while keeping the same quality of the physical results, we have tested two new acquisition methods.

The first experimental device is composed of a Guinier-type powder camera emitting a 8 keV quasi-monochromatic beam (Copper anticathode), a reflexion dispersive crystal, a set of slots, a sample-support that moves perpendicularly to the beam, and a detector-support that allows to move it on the focus cylinder. This assembly allows to obtain a Bragg diffraction diagram.

The second experimental device is composed of an EXAFS (Extended X-ray Absorption Fine Structure) emitting a quasi-monochromatic X-ray beam in the range 5 keV to 20 keV (silver anticathode), a transmission dispersive crystal, a set of slots, a fixed sample-support and a detector-support that allows us to adjust the sample-detector distance. This source will be used at 9 keV, which is close to the Guinier-type powder camera beam energy. This source will not be used at lower energy, where it generates its own diffraction peaks. This device is useful to obtain absorption diagrams.

2.1 Detectors (*principle and environment*)

The detectors used in our experiments are off-the-shelves ones. Each of these two detectors, which have their own electronic acquisition system, is controlled by a computer. They both have, in front of their sensitive zone, a Beryllium window that does not absorb these X-ray radiations. The sensor used with the EXAFS spectrometer, is a linear position gas ionization detector, produced by Braun. Its sensitive area is 50 mm \times 1 mm, and its spatial resolution 0.1 mm. The manufacturer gives an efficiency of 50 % at 9 keV.

Photodiodes, manufactured by Hamamatsu, are solid sensors using NMOS technology. The array, composed by 1024 pixels, has a 2.5 mm \times 25.6 mm sensitive area, i.e., for each diode, a 2.5 mm \times 25 μ m \times 300 μ m active zone (300 μ m is roughly the depletion region depth). The distance between diodes, hence the resolution, is 25 μ m.

The charge-coupled device (CCD), manufactured by Thomson, is an other kind of MOS capacity, and has a 9.73 mm \times 9.73 mm sensitive imaging area. The array is composed by 512 \times 512 pixels, about 19 μ m \times 19 μ m each. As

we use it in spectrometry, we only need to operate in vertical binning mode, with 512 columns each of 9.73 mm height.

These two detectors are used in the same experimental conditions. They are cooled, at -39°C , in order to reduce thermic dark current. Moreover, the sensor case is vacuum-kept in order to avoid any loss by convection and freeze deposits on the active part of the detector. Finally their electronic wiring was produced and optimised in our laboratory (CEMES) in order to operate swiftly and with little noise.

2.2 *Need of post-treatment*

Nevertheless, experiments operating with semiconductors are put through statistical fluctuations which introduce difficulties to interpret the results. Recorded spectra are pretreated, in order to remove the offset (dark current removing) and the non-uniform pixel response (flat fielding). Dark current removing just removes the mean but not the fluctuations and this noise is further amplified by flat fielding. Other sources of noise include: the photon random emission from a source, an acquisition noise, due to capacities, resistors and transistors, and a quantification noise, due to the information coding by means of an analogical-numerical converter. The detector cooling and the acquisition circuitry optimisation allow to decrease all these noises, without however making their contributions insignificant. Therefore, we need a signal processing method that can extract a signal from the remaining noise.

3 **Processing Method**

This section is divided in four parts. The first part gives a brief summary of the principles of the wavelet transform. More details can be found in Daubechies [?], Mallat [?] and Meyer [?]. The second part describes the *à Trous* algorithm [?,?], a fast numerical method, that will be used in this paper. The third part deals with the detection of significant information in the raw data [?]. Finally, we explain how to reconstruct a signal from which as much noise as possible has been removed.

3.1 *The wavelet transform*

We want to study a non-periodic signal $S(x)$ and obtain information on both the location and the size of the characteristic features (later on called struc-

tures). The usual way of dealing with a signal is to decompose it on an orthogonal basis. This gives a unique set of coefficients, and if the basis is correctly chosen, one can easily distinguish between “relevant” signal and noise. Extracting the information is therefore just a matter of summing the contribution of the coefficients corresponding to “relevant” signal.

The wavelet transform of $S(x)$ is the decomposition on a basis of functions $\psi_{a,b}(x)$, all derived, by translation and scaling, from a unique function $\psi(x)$, called the “mother wavelet”:

$$\psi_{a,b}(x) = a^{-1/2}\psi\left(\frac{x-b}{a}\right) \quad (1)$$

The “ a ” parameter is related to the scale, “ b ” to the translation. The coefficient $a^{-1/2}$ is a normalisation and is not mandatory. One can choose any other exponent for a , depending on the kind of normalization one is interested in. The exponent $1/2$ is useful for normalizing the energy (the square of the modulus of the coefficients).

In order to be a wavelet, the function $\psi(x)$ must satisfy a “admissibility condition”, i.e. it must belongs to $\mathcal{L}^2(R)$:

$$\int_0^\infty \frac{|\hat{\psi}(\omega)|^2}{|\omega|} d\omega < \infty, \quad (2)$$

where $\hat{\psi}(\omega)$ is the Fourier transform of $\psi(x)$. This means, in particular, that the integral of the wavelet should vanish. It must also be well localized both in space and scale. In order for the transform to be useful, a diadic approach is well suited and the collection

$$2^{j/2}\psi(2^j x - k), \quad j \in Z, \quad k \in Z \quad (3)$$

must be an orthonormal basis of $\mathcal{L}^2(R)$. Daubechies [?] and Mallat [?] give a precise mathematical description of the properties needed for $\psi(x)$ to be a wavelet. Daubechies [?] showed that there exist compactly supported functions $\psi(x)$ satisfying those conditions. Mallat [?] discovered a systematic approach to find $\psi(x)$.

Given $\psi(x)$, the wavelet coefficients of $S(x)$ are simply:

$$C_{j,k} = \int_{-\infty}^{\infty} S(x)\psi_{j,k}(x)dx, \quad (4)$$

and the unique way of recovering $S(x)$ from the coefficients is:

$$S(x) = \sum_j \sum_k C_{j,k} \psi_{j,k}(x). \quad (5)$$

The previous formula gives $S(x)$ as a sum of details from the finest scale to the largest one. One can also choose another orthonormal basis consisting of the collection $\psi_{j,k}(x)$, $j \in N$, $k \in Z$ (set of relative integers), together with the collection $\phi(x - k)$, $k \in Z$, where $\phi(x)$ is a new smooth function with a rapid decay at infinity. Obviously, the two functions $\phi(x)$ and $\psi(x)$ cannot be chosen independently. We choose ψ in such a way that:

$$\psi\left(\frac{x}{2}\right) = \phi(x) - \frac{1}{2}\phi\left(\frac{x}{2}\right) \quad (6)$$

In this basis, $S(x)$ can uniquely be written as:

$$S(x) = \sum_{k \in Z} f_k \phi(x - k) + \sum_{j \in N} \sum_{k \in Z} C_{j,k} \psi_{j,k}(x) = E_0(S) + \sum_{j \in N} D_j(S) \quad (7)$$

where $f_k = \int_{-\infty}^{\infty} S(x) \phi_k(x) dx = \int_{-\infty}^{\infty} S(x) \phi(x - k) dx$. This means that $S(x)$ is the sum of a smooth part $E_0(S)$ and a sequence of finer and finer fluctuations $D_j(S)$. The function $\phi(x)$ is a smoothing function that allows us to cut the doubly infinite sum of eq. (5). The $\phi(x - k)$'s account for the large scale fluctuations, the $\psi_{j,k}$'s for the small ones. Obviously, the resolution corresponding to the smooth part can be chosen arbitrarily.

3.2 The à Trouis Algorithm

Consider now a real sampled signal. The ranges of variation of j and k are now finite. The scale range is defined by the sampling step for the smallest one (finest resolution) and the length of the signal for the largest (coarsest resolution). The range of variation of k is also defined by the length of the signal. Among the many possibilities for the two functions $\phi(x)$ and $\psi(x)$, we choose:

$$\begin{aligned} \phi(x) &= 1 - |x| \quad \text{for } |x| \leq 1, \\ &= 0 \quad \text{otherwise,} \end{aligned} \quad (8)$$

and:

$$\psi(x) = \frac{|x|}{4} - \frac{1}{2} \quad \text{for } 1 < |x| \leq 2,$$

$$\begin{aligned}
&= \frac{1}{2} - \frac{3|x|}{4} \quad \text{for } |x| \leq 1, \\
&= 0 \quad \text{otherwise,}
\end{aligned} \tag{9}$$

Figure 1 shows the shape of $\phi(x)$ (long dashes), $\frac{1}{2}\phi(x/2)$ (short dashes) and $\psi(x)$ (solid line).

The process of computing the wavelet coefficients can be done iteratively. The first trivial step is to take the function $\phi(x)$ with the same resolution as the sampled signal. We only have to compute the f_k^0 's (the superscript refers to the number of the step in the iterative process, the subscript refers to the position). From the formula for $\phi(x)$, one can see that the f_k^0 's are simply the sampled values of $S(x)$ (the integrals are changed into summations since we are dealing with a discrete set of points). Then we consider a resolution twice as large. We choose the normalizing coefficients such that we replace $\phi(x)$ by $\frac{1}{2}\phi(x/2)$. There is only one level of fine fluctuations to add to the smoothed part f_k^1 to recover the signal, and it is defined by the set of C_k^1 . We can now repeat the same procedure to the smoothed part of level 1. This will give us two new sets of coefficients f_k^2 and C_k^2 . With our choice of $\phi(x)$ and $\psi(x)$, the formulae to obtain the f_k^j 's and the C_k^j 's are very simple:

$$f_k^j = \sum_{n=-1}^1 h_1(n) f_{k+n2^{j-1}}^{j-1} \tag{10}$$

$$= \frac{1}{2} f_k^{j-1} + \frac{1}{4} (f_{k-2^{j-1}}^{j-1} + f_{k+2^{j-1}}^{j-1}), \quad j \geq 1, \tag{11}$$

and:

$$C_k^j = f_k^{j-1} - f_k^j. \tag{12}$$

One can see that the f_k^j 's are obtained by applying a low-pass filter $h_1(n)$ to the f_k^{j-1} 's. The C_k^j 's being obtained by a difference between two levels of f_k^j 's are in fact the result of a band-pass filter applied to the signal. Figure 2 illustrates the whole procedure. Due to the shape of the function $\phi(x)$, the associated wavelet $\psi(x)$ is called the B-spline of order 1.

In operator notation, the smoothing part corresponds to the use of the low-pass filter H_j and the difference to the band-pass filter G_j defined by:

$$f^j = H_j(f^{j-1}), \tag{13}$$

$$G_j(f^{j-1}) = f^{j-1} - H_j(f^{j-1}), \tag{14}$$

$$\tag{15}$$

where f^j represents the set of f_k^j .

For the reconstruction formula, things are quite simple. We are interested in reconstructing the sampled values, not the signal $S(x)$. When adding the f_k^j 's and the C_k^j 's, we get the f_k^{j-1} . Hence, the reconstruction formula is:

$$S(x_k) = f_k^{j_{max}} + \sum_{j=1}^{j_{max}} C_k^j, \quad (16)$$

where j_{max} is the maximum scale compatible with the length of the signal.

3.3 Structure Detection

Our aim is to use the wavelet transform as a signal processing tool to extract the behavior of the signal while removing the noise. We also want to have the possibility to compute the statistical significance of detections. The main philosophy of the procedure consists in isolating, in the space–scale plane where the wavelet transform of the signal is unfolded, the patterns made by coefficients which are significant compared to noise fluctuations. These coefficients point out the structures in the signal. Then the signal is reconstructed from the selected coefficients.

The patterns of significant wavelet coefficients are exhibited by means of a double thresholding. A first one is very restrictive, leading to the detection of very significant structures within the signal, with a drastic decrease of the noise, and a second one which allows the addition of more and more local details as it is chosen more and more permissive. These thresholdings are made scale by scale by comparing the wavelet coefficients derived from the studied signal to the distribution of coefficients computed from a pseudo–random signal. The pseudo–random signal is generated from the signal. For each point x_k of the signal, we draw at random a value $S_{ps}(x_k)$ of the pseudo–random signal, according to the probability distribution computed from the real signal. This strategy has already been applied in asteroid family detection (Bendjoya *et al.* [?]) and allows us to take into account the intrinsic distribution of the noise instead of having to model it.

The wavelet transform of a white noise signal gives fluctuating coefficients, with no noticeable pattern in the space–scale plane (see fig. 3). Also of importance is the fact that almost no coefficients deviate greatly from the mean. Figure 3b displays a representation by a gray scale coding of the wavelet transform coefficients of the signal performed by means of the linear wavelet. For a noiseless signal, there is a definite pattern in the space–scale plane (see fig. 4b). Particularly, some coefficients deviate by a large amount from the mean, and they tend to the mean on each side of the maximum for each scale (this is due to the absence of structures on each side), with a width proportional to

the scale. For a noisy signal, we can still see definite patterns in the space-scale plane (see fig. 5b). In the following, in order to determine the mean and quantify a “large deviation”, we evaluate the mean and standard deviation, for each scale, of the coefficients of the wavelet transform of a pseudo-signal $S_{ps}(x)$. As mentioned above, this pseudo-signal is a random data set having the same cumulative distribution as the signal under study $S(x)$.

We use a pattern-recognition algorithm in order to isolate the significant patterns in the space-scale plane. This is done by a double thresholding. First, we search for points in the space-scale plane with a very high level of significance (the seeds). Those seeds are the points where the coefficients deviate by a large amount, i.e. by more than a given number times the standard deviation, from the mean. This defines a first threshold, the seed-threshold $T_{se}(a_j)$. Usually we take $T_{se}(a_j) = 4\sigma_{S_{ps}}(a_j)$, where $\sigma_{S_{ps}}(a_j)$ is the standard deviation of the coefficients of scale a_j of the pseudo-random signal. The thresholding is done by retaining the points where the coefficients minus the mean, defined above, are larger (resp. smaller) than T_{se} (resp. $-T_{se}$) for a bump (resp. hole) detection. The seeds give a rough estimate of the size, position and shape of significant structures in $S(x)$. In order to add details to these detected structures, we need to take into account some coefficients in the vicinity of the seeds.

We define a second threshold, the skeleton-threshold T_{sk} which determines the points that might be associated with the seeds. Usually $1.0\sigma_{S_{ps}}(a_j) \leq T_{sk}(a_j) \leq 2.5\sigma_{S_{ps}}(a_j)$. Among the points that are selected by this second weaker threshold, those that are connected by nearest neighbor relation to a seed actually belong to a pattern. In other words, a wavelet coefficient is included if the absolute value of its difference from the mean value is above the threshold, and if it is adjacent to a “seed” either in the location index or the scale index direction. These new points are new seeds. The procedure of extension is iterated until no new point is added. As seen before (Fig. 4b) a structure in $S(x)$ gives significant coefficients on a range of scales. The most significant ones, corresponding to the seeds, belong to a restricted range of scales (see Fig. 6a). Therefore we have to extent the seed along the scale axis. Furthermore, at a given scale, only the coefficients in the middle of the structure are selected by the seed-threshold. However the structure contributes to coefficients that are not selected by T_{se} but are significant according to T_{sk} (see Fig. 6b). Then we must consider them. At this stage of the procedure we have a set of coefficients, with a level of significance defined by T_{sk} but the cores of the patterns thus defined are significant at the level defined by T_{se} . In all the situations presented here, we use $T_{sk}(a_j) = 2.5\sigma_{S_{ps}}(a_j)$. For a description of an automated scheme to select the second threshold, see Bendjoya *et al.* [?].

In the following, we denote A the operator acting on a signal S and giving the wavelet transform coefficients in the skeleton as defined above, and zero

outside the skeleton. We also note $\nu = A(S)$.

3.4 Reconstruction

The argument underlying the thresholding procedure is that the so-called significant coefficients correspond to the “real” signal. The coefficients outside the skeleton represent noise. So we are looking for a signal whose wavelet transform coefficients equal those of the experimental data set on the skeleton. However, due to the over sampling, if we consider the set of coefficients of the skeleton surrounded by zeros, we do not have the wavelet transform of a signal. In other words, if we reconstruct the signal from the skeleton using equation (16), we get a signal whose wavelet transform is not the set of coefficients we used. In particular, coefficients outside the skeleton will not vanish.

Hence, we want to find the signal whose wavelet transform is as close as possible to the original one on the skeleton, and as small as possible outside the skeleton. This is a typical inverse problem. Several authors [?,?,?] have shown that iterative methods are well suited to solve this problem. After many trials, we choose to use the reconstruction method proposed by Rué *et al.* [?]: the conjugate gradient.

We want to find the signal F that minimizes $\|\nu - A(F)\|$ for F . For a set of wavelet coefficients W , the distance $\|W\|$ is defined by:

$$\|W\| = \sqrt{\sum_j \sum_k (C_k^j)^2} \quad (17)$$

The distance is minimum if F is solution of:

$$\tilde{A}(\nu) = (\tilde{A} \circ A)(F) \quad (18)$$

The operator \tilde{A} is the join operator associated with A (A is the thresholding operator defined in section 3.3). A is the union of the vector operators A_j giving the wavelet transform coefficients of scale j . We have:

$$A_j = G_j \cdot H_{j-1} \dots H_1 \quad (19)$$

Hence the join operator \tilde{A}_j is simply:

$$\tilde{A}_j = \tilde{H}_1 \dots \tilde{H}_{j-1} \cdot \tilde{G}_j \quad (20)$$

And \tilde{A} is the sum of the \tilde{A}_j 's. Since

$$G_j = Id - H_j \quad (21)$$

the join operator is

$$\tilde{G}_j = Id - \tilde{H}_j \quad (22)$$

H_j is the correlation of the signal with the filter coefficients $h(n)$. Therefore, \tilde{H}_j is simply the convolution of the signal with the same coefficients. And $h(n)$ being symmetric, it turns out that $\tilde{H}_j = H_j$.

The operator \tilde{A} is applied to the thresholded wavelet coefficients C and gives a signal \tilde{F} :

$$\tilde{A}_j = (H_1 \dots H_{j-1}) \cdot G_j \quad (23)$$

$$\tilde{F} = \tilde{A}(C) = \sum_j \tilde{A}_j(C_k^j) \quad (24)$$

Since we have a symmetric operator, $\tilde{A} \circ A$, the gradient method can be applied [?].

We introduce at this stage two new quantities: the residual wavelet $C_r = C - A(F)$, which corresponds to the difference between the original wavelet coefficients and the wavelet coefficients of the reconstructed signal on the skeleton; the residual signal F_r which introduced in the conjugate gradient scheme and is related to $\tilde{A}(C_r)$.

We now detail the steps of the conjugate gradient method. $F^{(n)}$ indicates the estimation of F at step n .

- (1) Initialization step. The estimated signal $F^{(0)}$, the residual wavelet $C_r^{(0)}$ and residual signal $F_r^{(0)}$ are initialized :

$$\begin{aligned} F^{(0)} &= \tilde{A}(C) + f^{j_{max}} \\ C_r^{(0)} &= C - A(F^{(0)}) \\ F_r^{(0)} &= \tilde{A}(C_r^{(0)}) \end{aligned} \quad (25)$$

Since \tilde{A} operator acts only on the wavelet coefficients, $\tilde{A}(C)$ does not contain the information of the last smoothed signal $f^{j_{max}}$. In order to recover this information, we add it at initialization.

- (2) Computation of the convergence parameter $\alpha^{(n)}$:

$$\alpha^{(n)} = \frac{\|\tilde{A}(C_r^{(n)})\|^2}{\|A(F_r^{(n)})\|^2} \quad (26)$$

- (3) An iterative correction is applied to $F^{(n)}$:

$$F^{(n+1)} = F^{(n)} + \alpha^{(n)} F_r^{(n)} \quad (27)$$

- (4) Computation of the residual wavelet, inside the skeleton support:

$$C_r^{(n+1)} = C - A(F^{(n+1)}) \quad (28)$$

- (5) End of iteration criterium. We compute the distance between the reconstructed signalwavelet coefficients and the skeleton, inside the skeleton support. Outside of it, we compute the amplitude of the reconstructed signal wavelet coefficients. These two distances are weighted with the number of points in each domain. When this function reaches its minimum, we stop iterations. In other terms, at this minimum, the wavelet coefficients are as close as we can expect to be to the skeleton, without too much noise outside of it.

- (6) Computation of the convergence parameter $\beta^{(n+1)}$:

$$\beta^{(n+1)} = \frac{\|\tilde{A}(C_r^{(n+1)})\|^2}{\|\tilde{A}(C_r^{(n)})\|^2} \quad (29)$$

- (7) Residual signal computation:

$$F_r^{(n+1)} = \tilde{A}(C_r^{(n+1)}) + \beta^{(n+1)} F_r^{(n)} \quad (30)$$

- (8) Go to step (ii)

3.5 Complete algorithm

We now summarize the different steps of the complete procedure:

- (1) We generate a pseudo-signal whose cumulative distribution function matches those of the real data set;
- (2) we compute the wavelet transform of both signals;
- (3) we apply a first restrictive thresholding T_{se} to localize the strongly significant structures within the signal, and eliminate nearly all noise;
- (4) the seeds formed by the coefficients selected by T_{se} are then enriched along both the scale and the space axis by related coefficients coming from a second more permissive thresholding T_{sk} ;

- (5) using the skeleton (which is actually a mask), and the smoothed signal, we get a first estimate of the reconstructed signal from equation (16);
- (6) Using the previous initial estimate and the coefficients of the skeleton, we reconstruct the signal with the conjugate gradient method described above.

4 Tests on a simulated signal

In order to test the method, we use a synthetic signal, computer generated from the knowledge of the location, height and width of the peaks in a real signal (fig. 4a). The signal was generated assuming gaussian peaks. Positions and amplitudes of the peaks are those corresponding to the 16 most intense peaks of Norethisterone detected using an automatic diffractometer Siefert $\theta-2\theta$ experiment). We add a gaussian noise (Fig. 3a) whose amplitude (RMS fluctuation of 5.3 ADU) is comparable to the amplitude of noise obtained in a real signal sampled from a photodiode array (fig. 12).

4.1 Reconstruction

Figure 7 shows the reconstruction of the noisy synthetic signal (fig. 5a) as a whole. One can see that we miss the fine structures in the weak peaks. On the other hand, the intense peaks are correctly reconstructed with a good restitution of the fine structures. This effect is due to the generation of the pseudo-signal. The pseudo-signal has the same cumulative distribution function as the real signal. Hence there are many points with a large value (same amplitude as the intense peaks). Therefore, the standard deviation of the wavelet coefficients of the pseudo-signal will be large. It will eventually exceed the value achievable by the wavelet coefficients of the real signal in the region of weak peaks.

In order to avoid this “contamination” of the pseudo-signal in the regions of weak amplitude structures by the intense ones, we split the signal into pieces where the structures are of comparable amplitudes. Figure 9 shows the cumulative probability distribution function of values for both the complete signal (solid line) and the right most part ($x \in [23.5; 35]$, dashed line). There is no large value in the second one. Hence the wavelet coefficients of the corresponding pseudo-signal will be smaller as well as their standard deviation. Using a splitting into 5 pieces we reconstruct both the intense peaks and the weak ones.

When using directly this method on the different parts of the signal, we have

problems at the limits of the sub-signals. The different samples do not match exactly together, due to edge effects. To suppress these artifacts, we modify slightly the analysing method. At first, we compute the skeleton for the whole signal. Then we define three sub-signals ($x \in [10; 14]$, $x \in [17.5; 19]$, $x \in [23.5; 35]$) corresponding to the regions with weak peaks. We compute the thresholding for all the sub-signals and the corresponding skeletons. Then we add all the sub-sample skeletons to the complete skeleton, and this new skeleton is used to reconstruct the real signal (Fig. 8).

4.2 *Detection of peaks*

We now want to test the quality of this reconstruction. The reconstructed signal is good if we can perform the measurement we want on it: actually, recover the information needed to generate the noiseless signal. Here we want to measure the location and the amplitude of the peaks.

On a given reconstructed signal a peak is a local maximum. But there can be small fluctuations at the top of the peaks, yielding multiple detections. A point source gives a signal that has a non vanishing width due to imperfections of the instrument. So we do not expect to find peaks (or actually resolve peaks) that are closer than the width of the point spread function (PSF). So we select only the peaks that are isolated by at least the width of the PSF.

Next we drop all peaks with an amplitude too small compared to the amplitude of noise. Experimentally, we have verified that our method is reliably able to detect and reconstruct structures with amplitude equal to the standard deviation of the noise. For smaller amplitudes, the detection is more speculative. So we estimate the standard deviation of noise from the standard deviation of first scale of the wavelet coefficients of the real signal and drop all peaks with amplitude less than this standard deviation, corresponding roughly to a signal to noise ratio of 0.7.

The exact shape of the reconstructed signal vary with the pseudo-signal used. Hence the detected peaks may also vary. But the “real” peaks should all be detected in every reconstruction and they should be within a distance equal to the width of the PSF from each other. So we keep only these peaks that are grouped in a range less than the width of the PSF and appear at least in half the reconstructed signals. We use this rather soft criterium (half of the reconstructions showing the peaks) in order to detect peaks with a signal to noise ratio close to 1. Obviously, this increases the chance of false detections. Fig. 10 displays the peaks detected using 9 realizations of the pseudo-signal with error bars at 1σ of the distribution of amplitudes of the detected peaks. The dashed line represents the synthetic signal. One can note that the second

peak from the right is not detected, and that the five right most peaks are slightly, but consistently too weak. This is due to background removal. The background is a significant signal as well as the peaks. the only difference is the scale on which it lies. In the region $x \in [23.5; 30]$ the signal can be viewed as peaks superimposed on a large scale structure. This large scale structure is the “background” in that region and we remove it when automatically determining the peaks; using a manual detection of peaks, this effect can easily be suppressed.

In summary there are three criteria to select peaks: a PSF criterion (narrow peaks are rejected), an amplitude criterion (weak peaks are rejected) and a statistical criterion (peaks that are spurious against noise fluctuation are rejected).

5 Real data

Applying this method on real data allows us, after their classical pretreatment (dark current and diode non-uniform response removal), to take the remaining random noise off, which might hide the weakest peaks. Here are showed the results obtained with each experiment, a Guinier-type camera and an EXAFS spectrometer respectively.

The first experiment used is the Guinier-type camera, associated with the photodiode array. A pharmaceutical compound, the Norethisterone, presenting a lot of different intensity peaks, is used. The spatial resolution deduced from these measurements should only depend on our detector resolution.

In the case of an absorption study of a Copper sample, we test this method on a signal where peaks are superimposed on a large step background. The important informations lying near the step, we put in a prominent position the fact that the wavelet treatment can distinguish peaks with weak intensity close to peaks with strong one. This test is made using the EXAFS experiment, and the CCD camera.

5.1 Norethisterone

Figure 12 shows a 600 s acquisition using a Norethisterone sample. Actually, this signal is not really a single one, but the sum of two 300 s frames, in order to limit noise amplitude. These data are those corresponding to the pre-processed data, after dark current removal and flat-field corrections. The mean dark signal level is around 5000 ADU (Analog Digital Unit), while the

relative amplitude of the peaks varies from 10 to 160. So we are looking for pretty low intensity peaks.

One can see nine strong peaks, and two weak ones, lying at intensity levels close to the noise level. Our estimate of the noise standard deviation is 5.3 ADU.

Figure 13 shows the reconstructed signal, after removal of the background (dashed line). The open diamonds represent the location and intensity of the detected peaks, with error bars corresponding to $\pm\sigma$ of the distribution of intensities for the different pseudo-signals.

Table 1 shows position and intensity of the peaks detected using our method, and the results of a $\theta - 2\theta$ experiment (used in section 4 to generate the test signal). The intensities are relative to that of the largest peak in each case. Lines denoted with † correspond to undetected peaks. Line denoted with †† corresponds to a false detection in the reconstruction of the real signal. Lines denoted with ‡ correspond to peaks that lie beyond the end of the phodiode array used for this frame. The relative amplitudes of the 2nd and 4rd (from the left) are not preserved, but this is expected from the signal we used. We observe a similar behaviour between the weak peaks on the right side of the signal. But here, the peaks have roughly the same amplitude (error bars overlap). The relative amplitude have in general the correct order of magnitude.

5.2 Copper

Figure 14 shows a 2400 s acquisition of a copper sample's absorption spectrum (EXAFS). Two problems are present here: first, the low angular resolution (only 430 points for the whole signal; the width of a peak is typically 10 pixels); second, a very large step occurring right in the middle of the signal.

Figure 15a shows the reconstructed signal using the whole data set, defining two sub-regions, one for the lower-left part, the other for the top-right part. When we reconstruct the large step, there appear artefacts just at the bottom and the top of the step, with amplitude as large as the amplitude of the peaks located there. Hence the structures we see at the bottom and the top of the step are not actual peaks (this can be seen by comparing the reconstructed signal to the real one). Figure 15b shows the reconstruction using two separate data sets. Here, we do not have the step. We find a significant peak just at the top of the step, but nothing significant at the bottom. Clearly the CCD signal suffers from a lack of spatial resolution (too little points to represent a peak). Moreover the weakness of some peaks with respect to the amplitude of the feature of the large scale background results in difficulties to get physical parameters.

6 Conclusions

In this paper we have described a new filtering method, for a mono-dimensional signal, based on wavelet analysis.

Using wavelet analysis, a well chosen thresholding scheme and an efficient reconstruction method, we can detect and analyse relevant features at very low (less than one) signal-to-noise ratio. This method also preserve shapes of detected structures.

We applied successfully this method to X-ray diffraction data, and results are in good agreement with other known results (obtained with much more expensive and accurate experiment). Even if this method have been developed for this peculiar use, it can be applied to any mono dimensionnal signal of that kind.

In a forthcoming paper the method is applied to a large set of data.

All the software developped and used for this analysis is freely available by anonymous ftp at location holst.obs-nice.fr:/pub/wavelet.tar.gz.

Acknowledgements

We wish to thank F. Rué for all the very usefull discussions about the choice of an appropriate reconstruction method.

Part of this work (N. Ennode) has been supported by a CIFRE grant number 614/92, supported by EDF.

We also wish to thank l'Observatoire de la Côte d'Azur for providing facilities to P. Bury.

Real signal		Theoretical signal	
Location	Intensity	Location	Intensity
13.406	6.21	13.406	10.35
14.488	37.31	14.495	40.26
14.906	100.00	14.936	100.00
15.257	48.49	15.243	36.62
16.573	64.73	16.557	47.44
18.271	8.61	18.298	10.14
†		18.533	7.49
19.854	13.26	19.809	16.54
20.850	42.78	20.843	38.56
21.486	2.64	††	
22.419	27.18	22.453	23.32
†		24.542	6.22
25.746	9.10	25.731	9.32
26.267	10.82	26.277	8.23
‡		26.925	8.83
‡		28.008	4.49
‡		29.251	5.13

Table 1

Position and relative amplitude of peaks. †: undetected peaks. ‡: theoretical peaks outside the range of the photodiode array. ††: false detection: the peak detected in the real signal correspond to no peak in the theoretical signal.

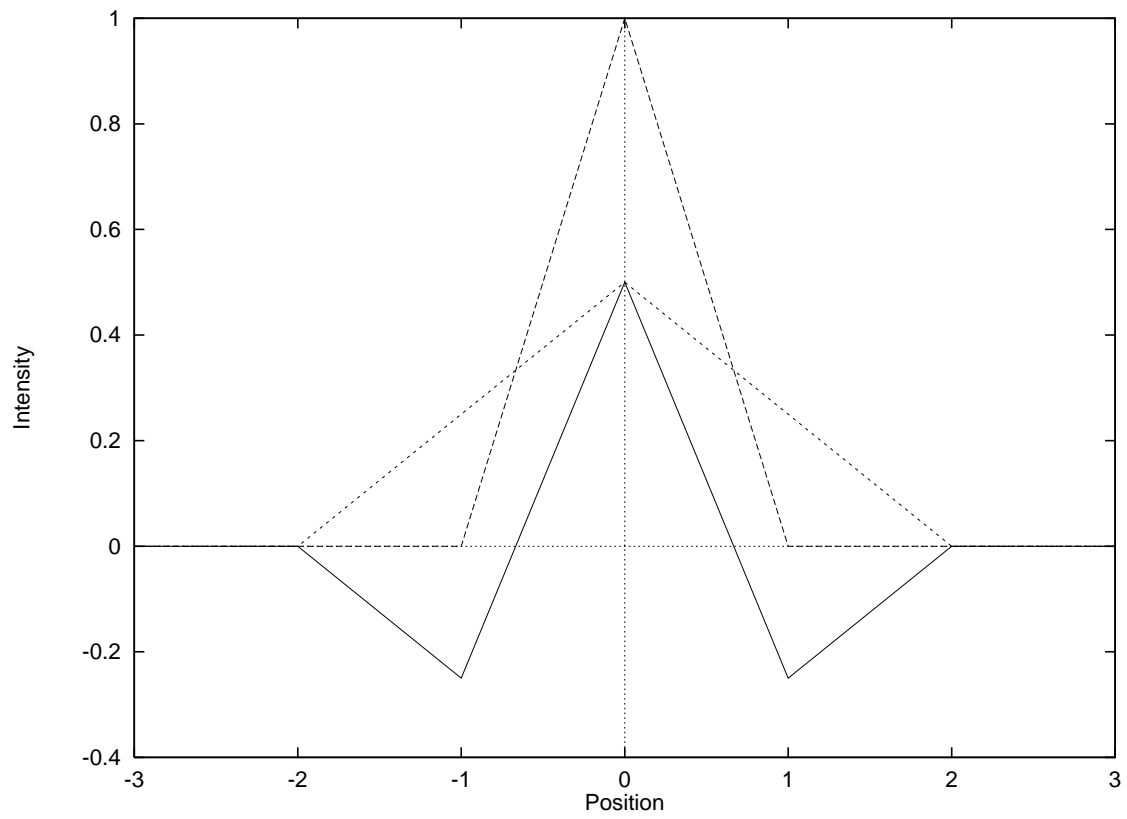


Fig. 1. Shape of the smoothing functions $\phi(x)$ (long dashes) and $\frac{1}{2}\phi(x/2)$ (short dashes), and the wavelet function $\psi(x)$ (solid line).

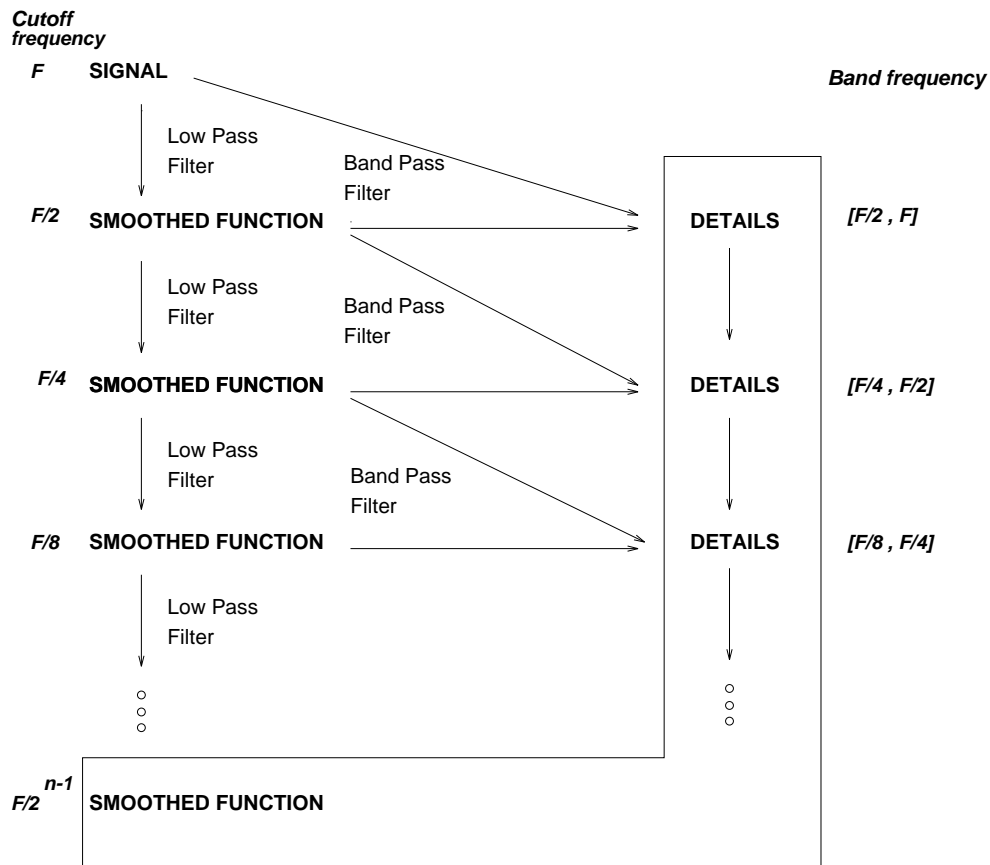


Fig. 2. Principle of the wavelet transform. The signal is smoothed using a low-pass filter. Then, in order to get the first wavelet plane, we compute the difference between the smoothed signal and the original one. We iterate this algorithm on the smoothed data, using a larger and larger filter (by a factor of 2 at each step) as the scale increases. To reconstruct the signal, one can see that we just have to sum all the wavelet plane (details data) and the last smoothed function.

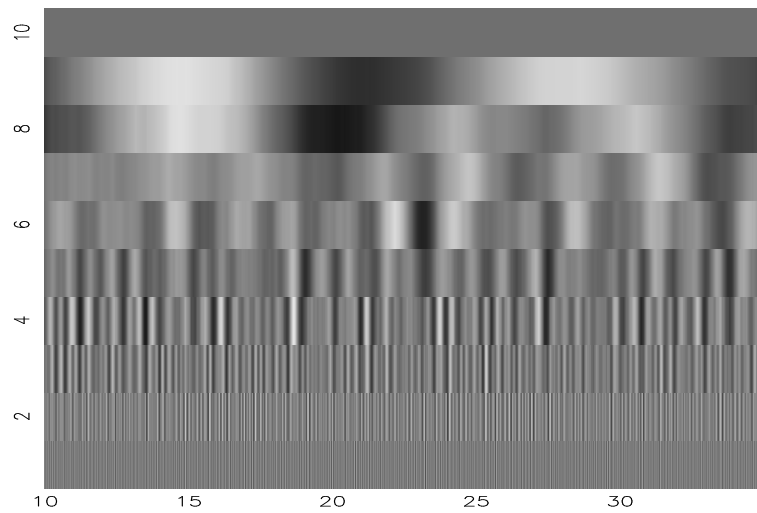
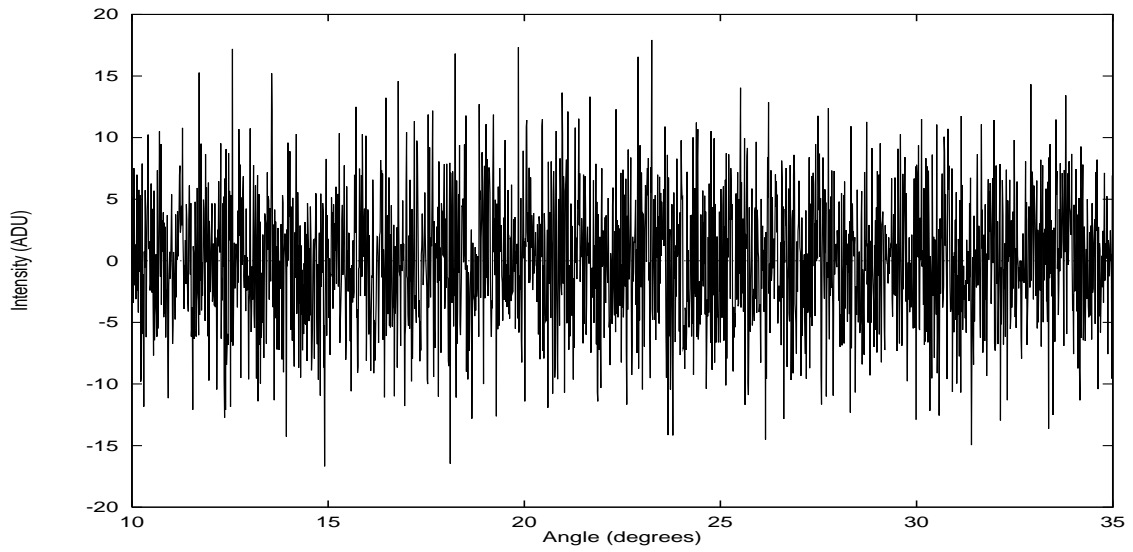


Fig. 3. (a) The noise signal, issued from 2000 drawn, with a gaussian distribution. The x-axis is the position, and the y-axis is the signal intensity $S(x)$. (b) Gray scale representation of the wavelet transform coefficients of the noise signal (a) using a linear B-spline wavelet. The x-axis is the position, the y-axis is the scale number. Large scales are on the top of the figure, and small ones at the bottom. For a given scale, the value of the coefficient is coded by a gray level which varies linearly from black to white. The black color corresponds to the lower negative value, while the white corresponds to the higher positive value. This correspondence is made scale by scale.

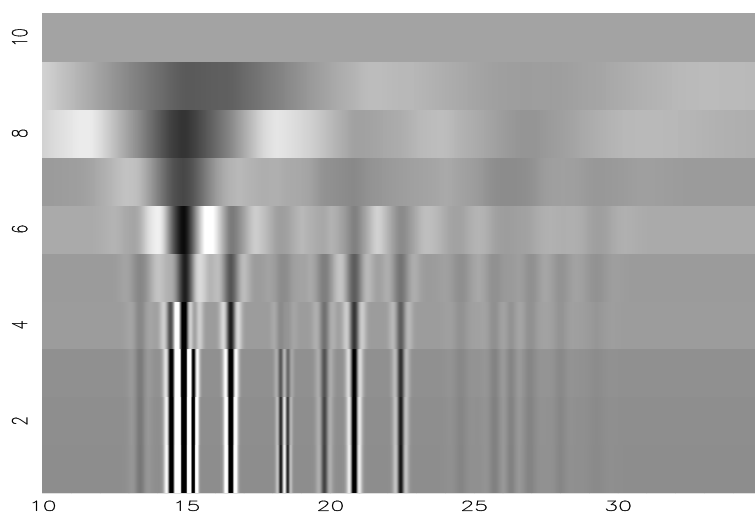
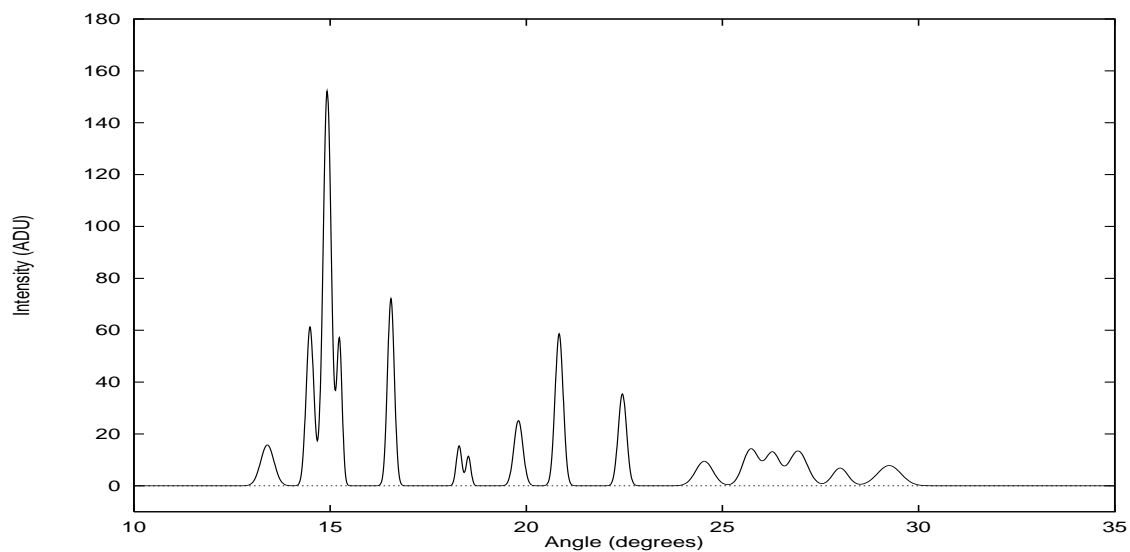


Fig. 4. (a) A simulated signal artificially generated with 16 peaks. Positions and intensity are issued from a $\theta-2\theta$ analysis of the Norethisterone sample. The peaks are gaussian shaped with σ estimated from the same analysis. (b) Same representation as 3b for the simulated signal (a).

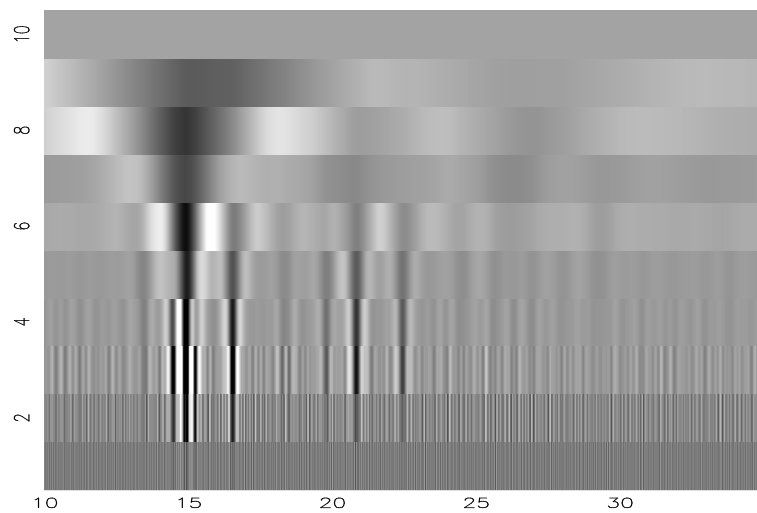
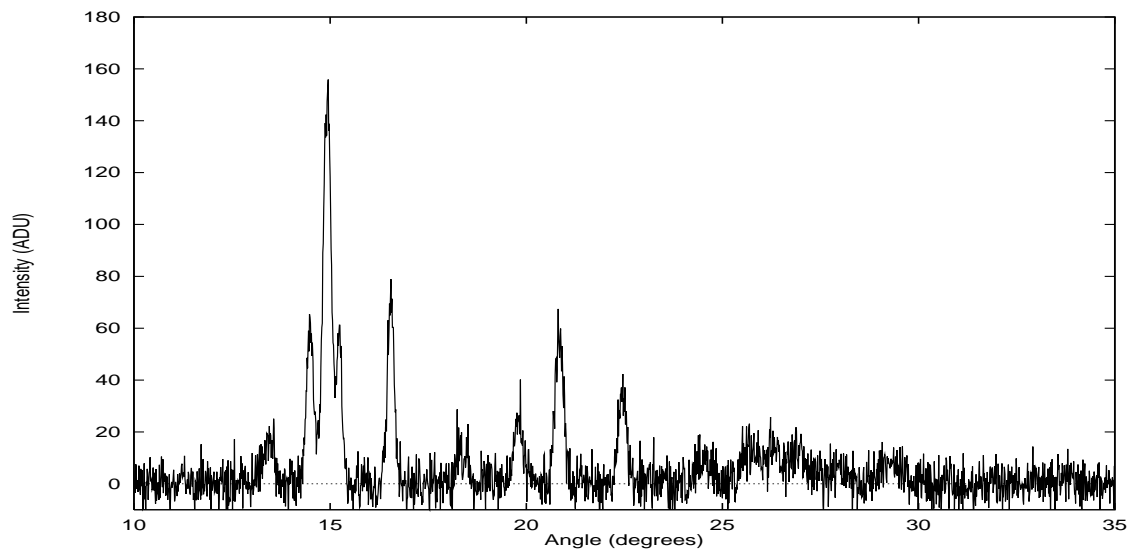


Fig. 5. (a) Addition of the simulated signal (Fig. 4a) and the noise (Fig. 3a). The amplitude of the noise is chosen to be similar to the amplitude of noise in real photodiode array spectra ($\sigma = 5.3$ in this case). (b) Same as Fig. 3b for the signal of (a).

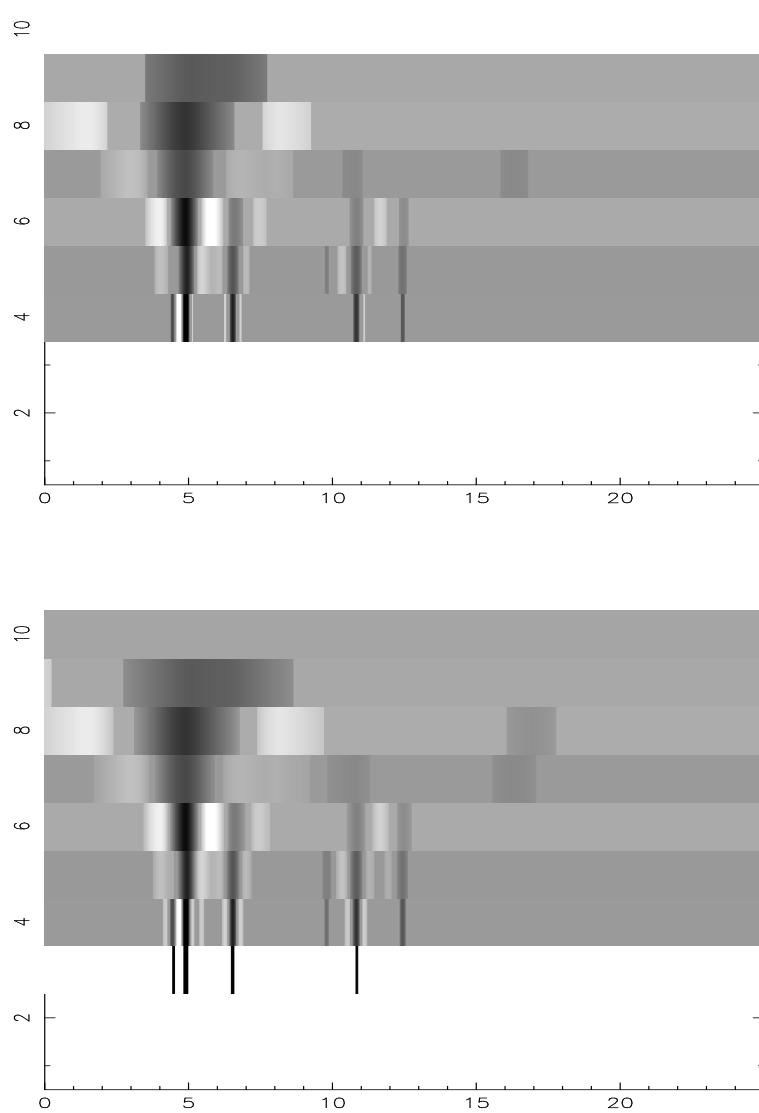


Fig. 6. (a) The seeds obtained after a thresholding of the noise signal at a 4σ significance level. The threshold is estimated separately for each scale. (b) The skeleton, we select all the point in all wavelet planes that are higher in intensity than 2.5σ and that are connected (in terms of spatial position and scale) to the 4σ seeds.

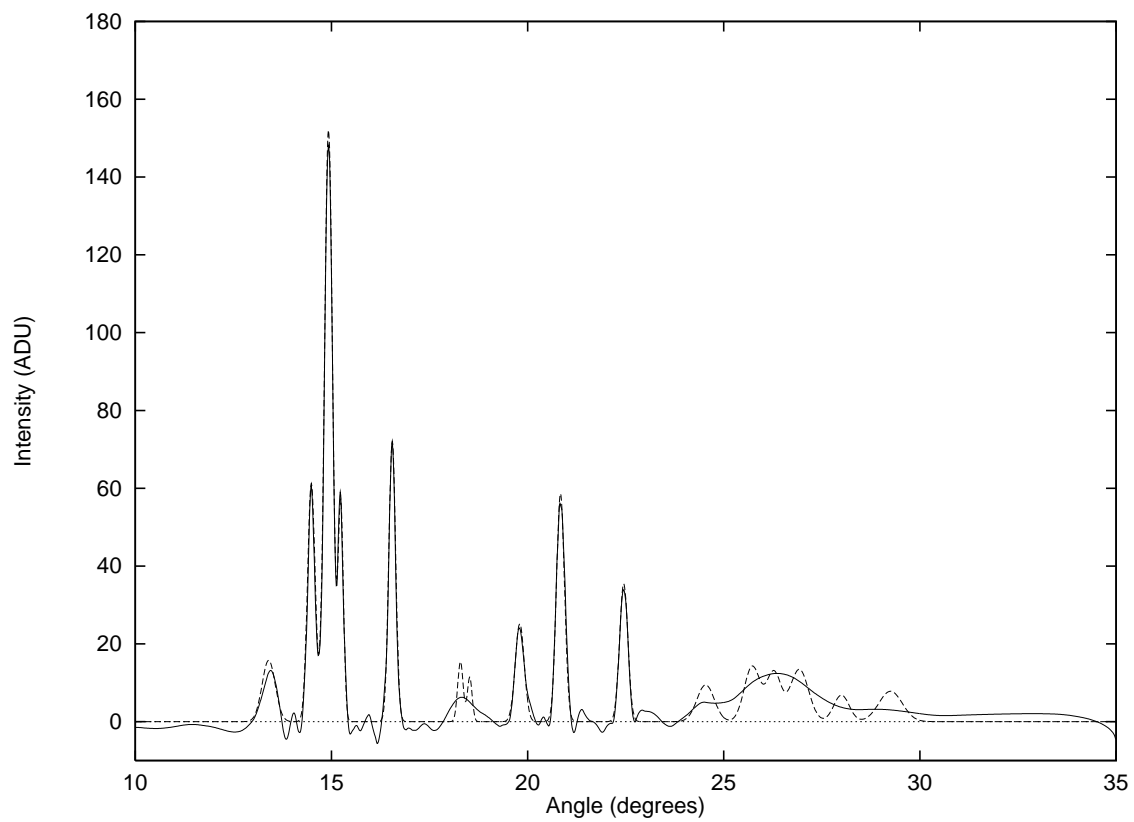


Fig. 7. Solid line: the signal reconstructed using a seed of 4σ and a skeleton at 2.5σ of the signal showed Fig. 5. Dashed line: the initial simulated signal. One can see that we miss the weaker peaks, which are detected as grouped into a large one. We also miss double peaks that are very close and of small amplitude.

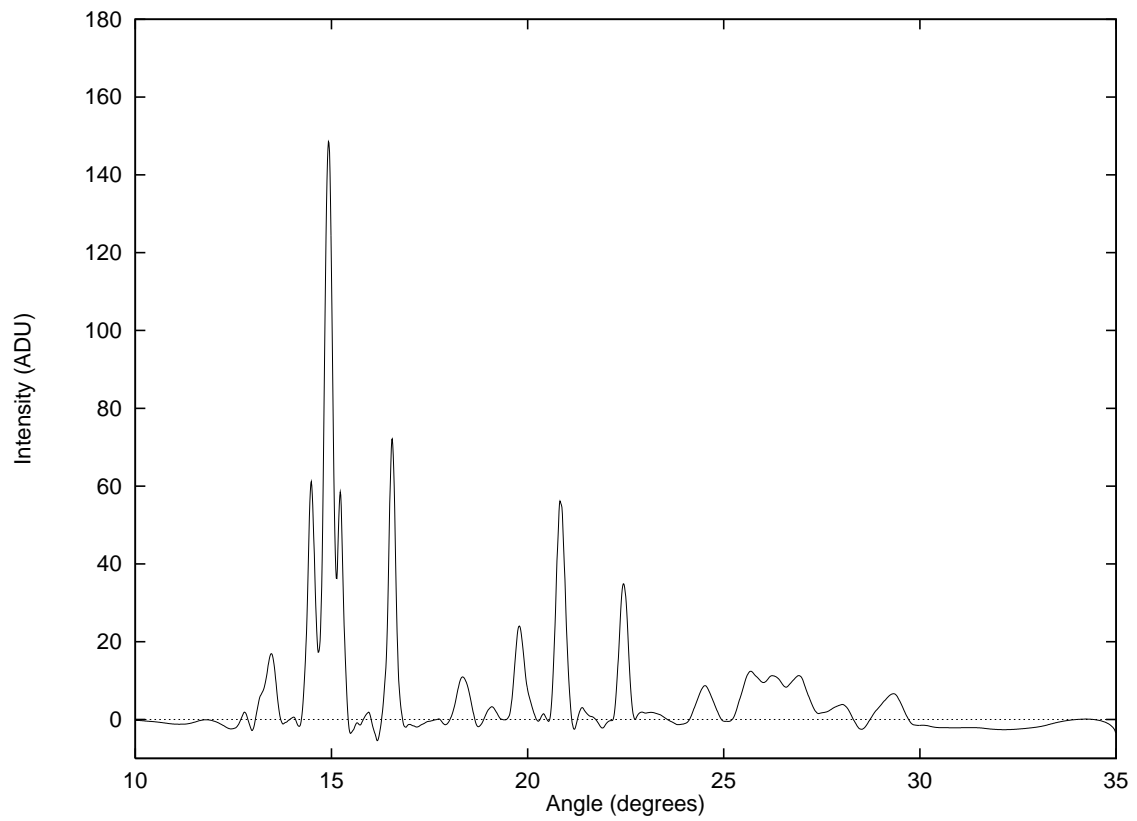


Fig. 8. Reconstruction of the signal by parts. The skeleton has been computed separately for each of the three regions with low amplitude peaks (10 to 14, 17.5 to 19 and 23.5 to 35). We then applied the reconstruction algorithm to this composite skeleton.

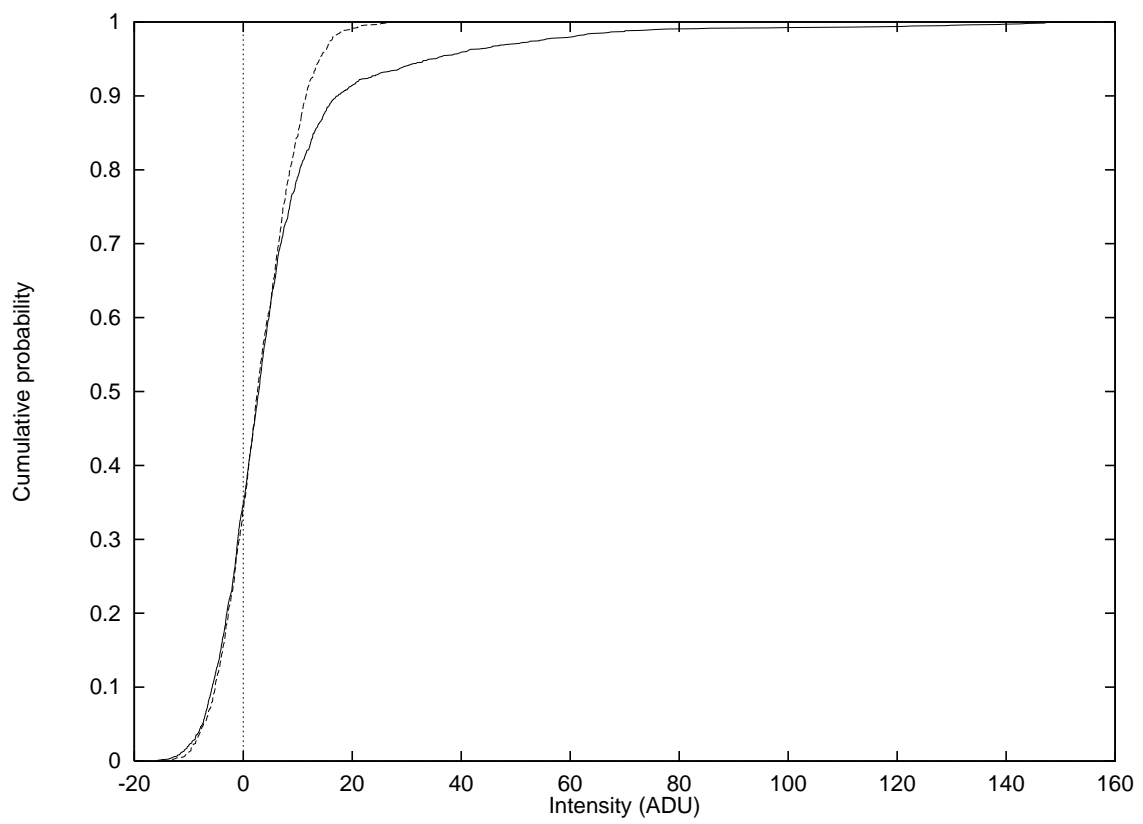


Fig. 9. Solid line: cumulative probability distribution function (PDF) of the whole signal. Dashed line: PDF of the right most part (23.5 to 35) of the signal.

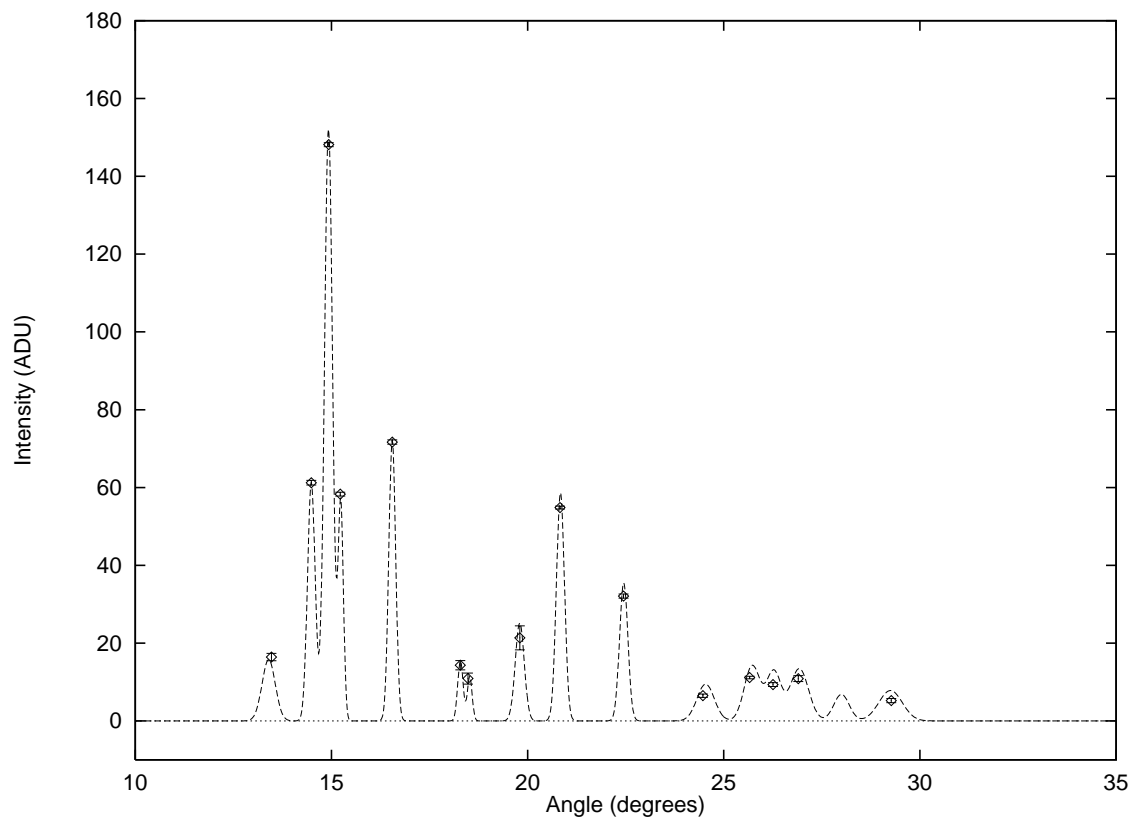


Fig. 10. Dashed line: the simulated signal. The symbols indicates the mean peak position for the noisy signal and 9 realisations of the pseudo signal. The errors bars (at σ) are issued from the data sets.

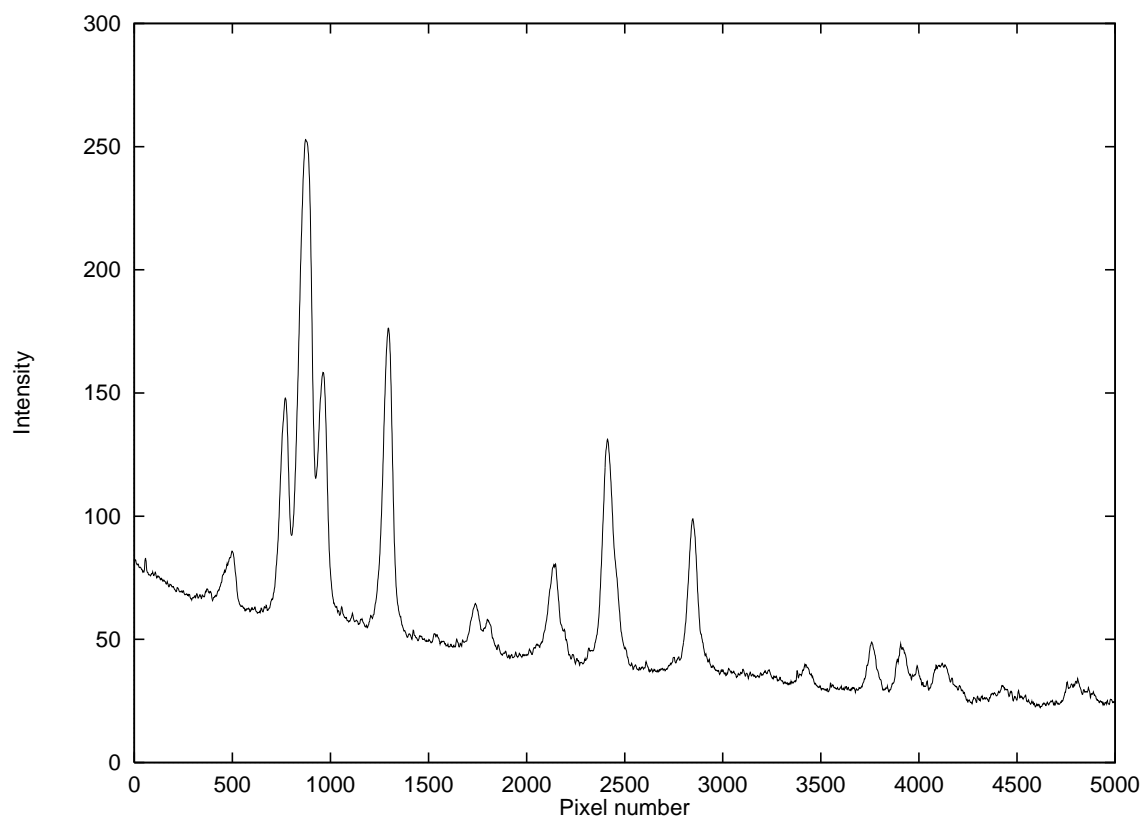


Fig. 11. Scan of a photographic image of norethisterone spectra obtained using a Guinier-type camera and Kodak HP5 film. The negative is scanned at a resolution of 2540 dpi ($10\ \mu\text{m}$ pixel size). One can see that the weak peaks on the right part of the spectrum are not precisely defined, even with this 64 hour exposure.

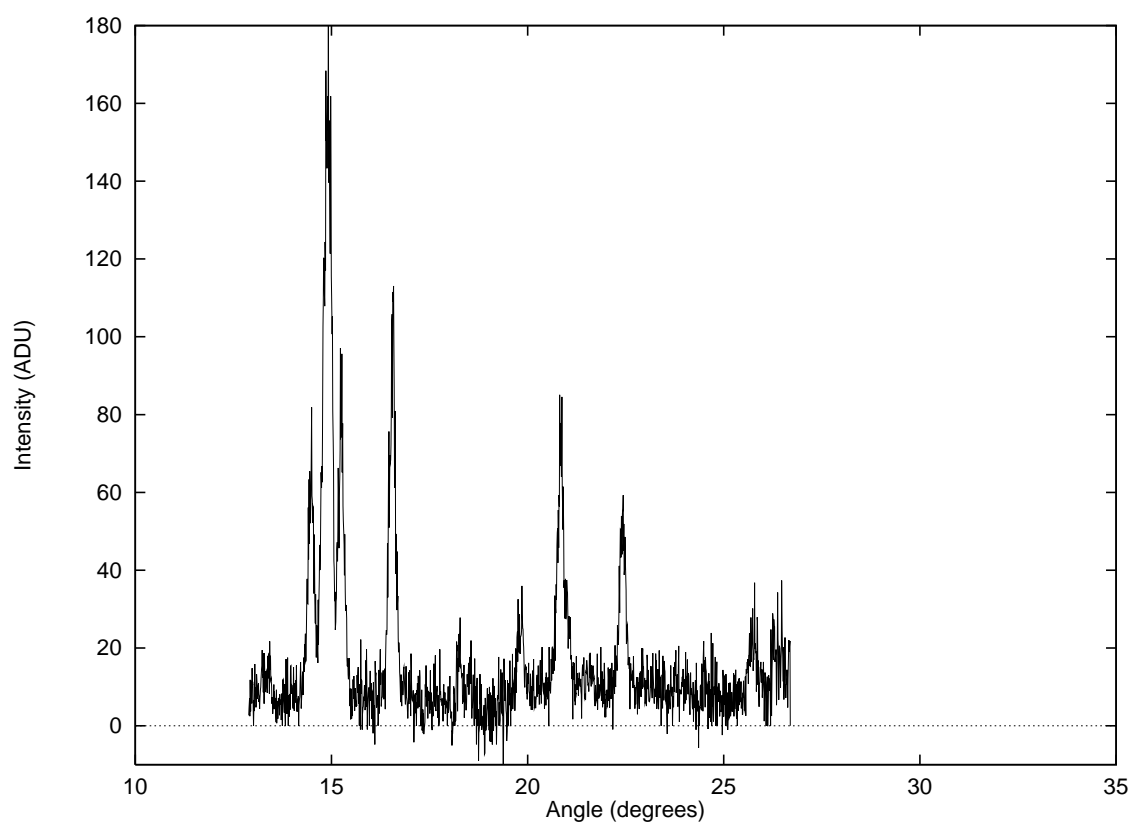


Fig. 12. Norethisterone 600 seconds exposure, pre-processed data. The X-axis corresponds to an angular position, and the Y-axis to intensity, expressed in ADU.

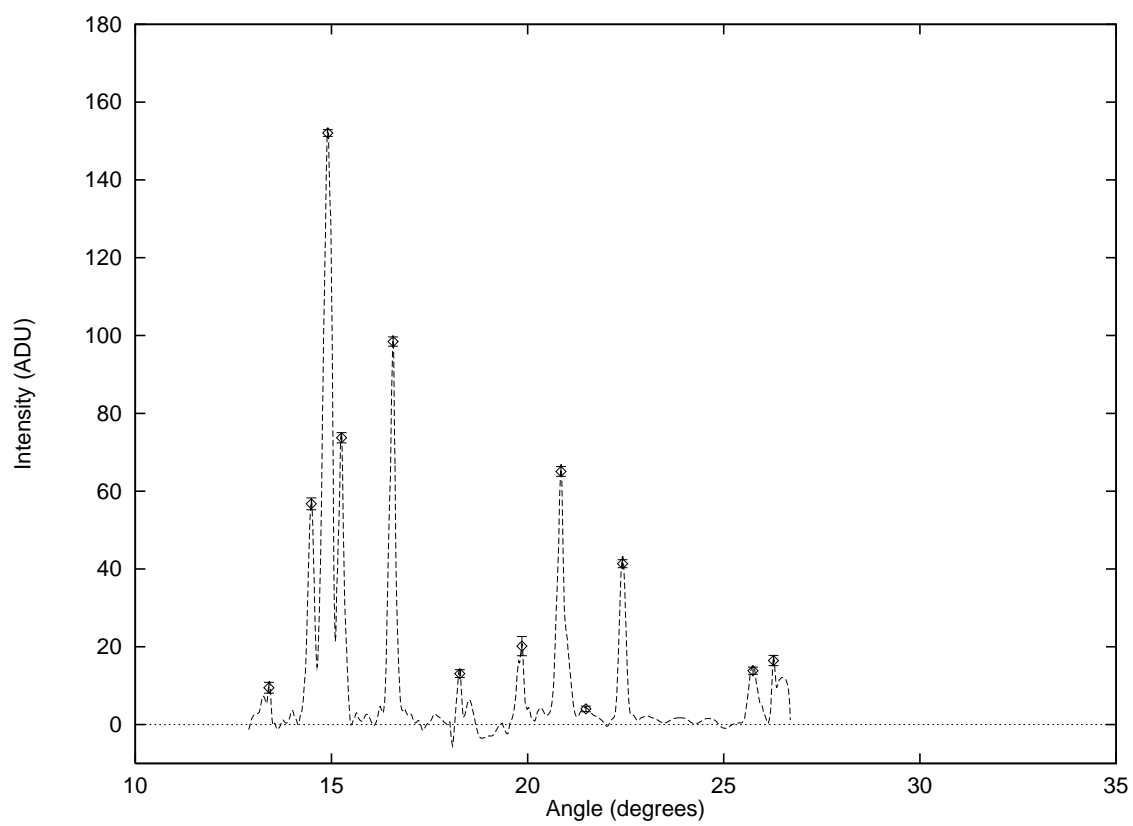


Fig. 13. Dashed line: Norethisterone 600 seconds, after treatment. The symbols indicates the mean peak positions. The errors bars (at 1σ) are issued from the data sets.

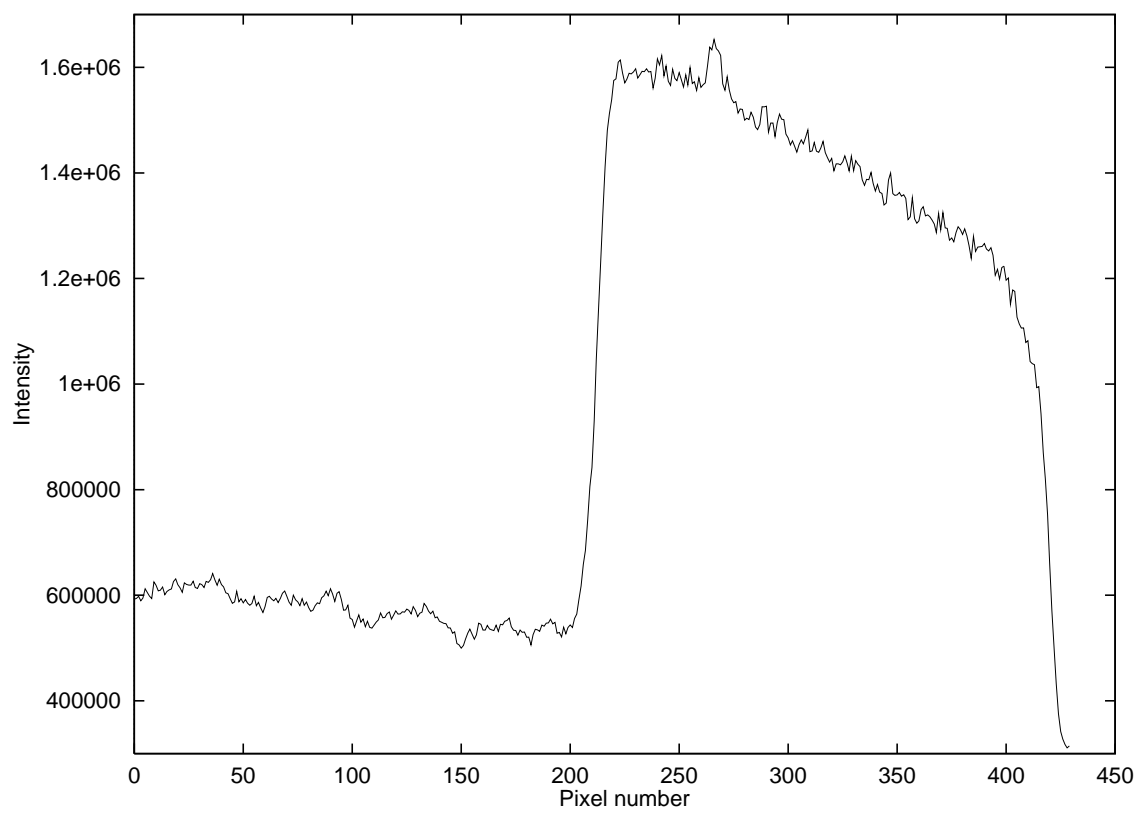


Fig. 14. Copper spectrum obtained with the EXAFS experiment, with six 300 s co-added frames. Note the importance of the step occurring at the middle of the spectra, and the poor resolution of the image.

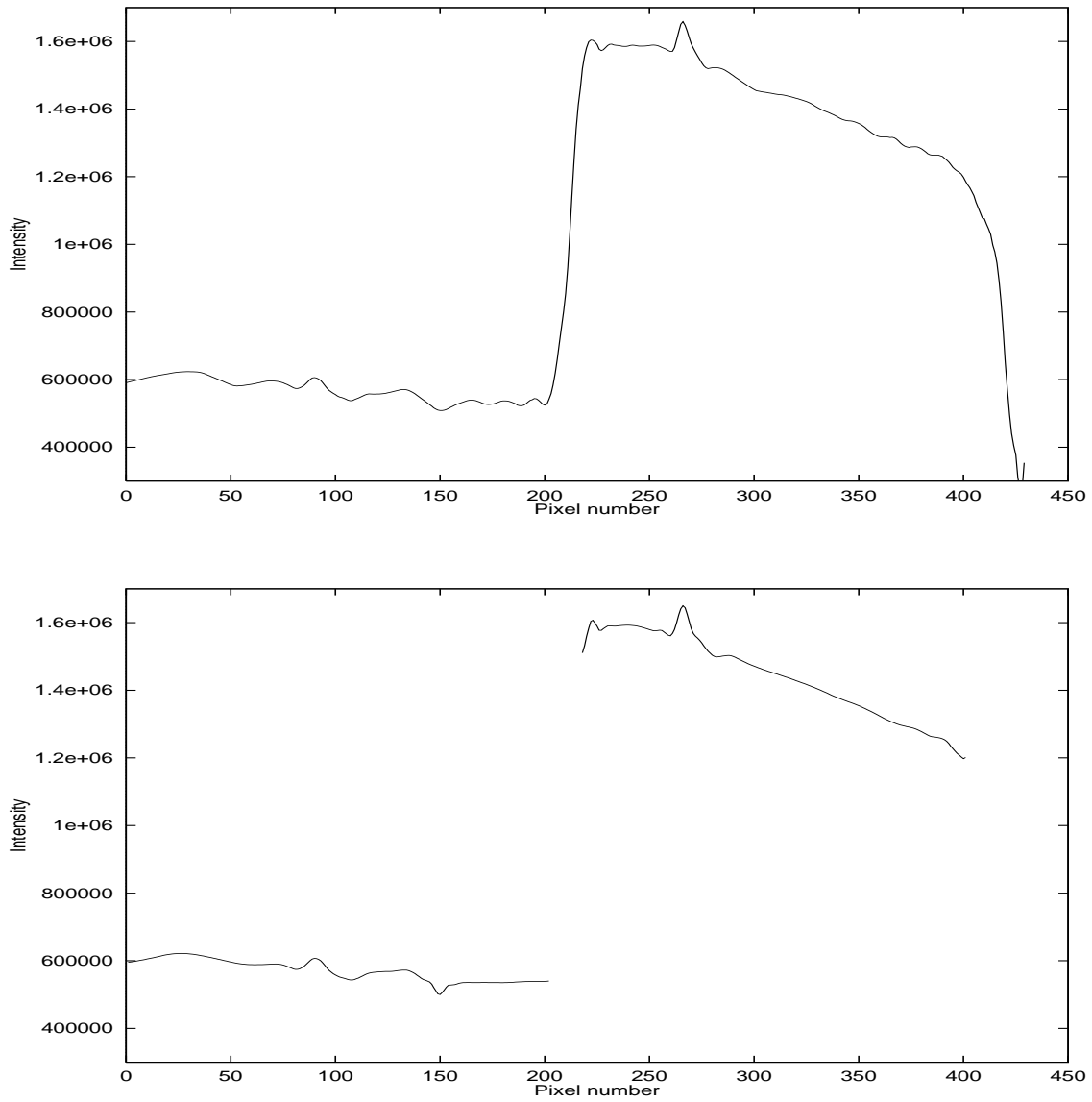


Fig. 15. (a) Spectrum after filtering, using the whole data set and two sub-regions ($x \in [1; 201]$ and $x \in [218; 400]$). Weak peaks near the step are due to reconstruction artifacts. In order to avoid this, we processed separately the left and right parts of the signal. (b) shows the results in that case, and present a much better agreement with the original (noisy) sample.