

# A TIGR like atmospheric profile database for accurate radiative flux computation

F. Chevallier, A. Chédin,  
F. Chéruy and J J. Morcrette

Research Department

March 1999

This paper has not been published and should be regarded as an Internal Report from ECMWF.  
Permission to quote from it should be obtained from the ECMWF.



## Abstract

The availability of representative (adequately sampled) collections of atmospheric thermodynamic properties is essential, as pertinent a priori information, to the modelling of relevant processes. In practice, the sampling of associated variables (temperature, moisture, ...) is made difficult by the high dimension of the model space. This paper analyses different topological methods for sampling vertical profiles of heterogeneous variables, like atmospheric temperature and water vapour concentration, in connection with the approaches developed for the successive *Thermodynamic Initial Guess Retrieval* (TIGR) databases at LMD. The most recent is chosen and applied to select a 6,000 atmospheric profile dataset by the sampling of a much larger (1,350,000 profiles) dataset from the ECMWF short-range forecasts. The sampled dataset is then used for training a neural network-based radiative flux profile computation model (NeuroFlux).

## 1 Introduction

The importance of the role played by radiation in climate research results in a growing interest for accurate modelling of forward or inverse radiative transfer problems. A key for increased accuracy is an improved use of pertinent and representative (adequately sampled) a priori information on the systems considered. A major attempt to sample this kind of a priori information on global scales has been the constitution of successive versions of the *Thermodynamic Initial Guess Retrieval* database (TIGR: Chédin *et al.*, 1985 ; Achard, 1991 ; Escobar-Nunoz, 1993 ; Chevallier *et al.*, 1998a) from Laboratoire de Météorologie Dynamique (LMD). Each version groups together hundreds of soundings sampled from larger databanks of observation of the atmosphere: radiosonde reports and, for the latest version (TIGR-3) only, satellite-retrieved atmospheric profiles. Providing initial guess solutions and a priori information (covariance matrices) for the weakly non-linear problem of the retrieval of atmospheric temperature profiles from observations of *TIROS-N Operational Vertical Sounder* (TOVS) on-board the *National Oceanic and Atmospheric Administration* (NOAA) polar-orbiting satellites, has been the main purpose of the TIGR databases. This approach has been developed in the framework of the Improved Initialization Inversion (3I: Chédin *et al.*, 1985; Scott *et al.*, 1999). For stronger non-linear problems, like the retrieval of water vapor from TOVS or the design of fast forward radiative transfer models, TIGR has been used for training artificial neural networks (e.g., Escobar-Munoz *et al.*, 1993; Chérury *et al.*, 1996; Rieu *et al.*, 1996; Chaboureau *et al.*, 1998).

In practice, the high dimensionality of the atmospheric variable space makes the sampling problem delicate. From one TIGR version to the next, important choices have been made and improved, so as to fit with the applications of the database. New applications of the neural network-based techniques, like the computation of longwave (LW) flux profiles in General Circulation Models (GCMs) have led to further improvements of the sampling strategy.

This paper summarizes the characteristics of the TIGR sampling methods, and discusses these prospects. Section 2 describes the TIGR database and reviews its three versions. Section 3 analyses the characteristics of the TIGR sampling methods on the basis of a simple but representative problem. An extension of the most recent TIGR sampling method is presented. Its application for the sampling of atmospheric situations from the ECMWF short-range forecasts is shown in section 4. The qualities of this new sampled database have been estimated in the framework of a neural network-based LW radiative flux profile computation model (NeuroFlux:

Ch eruy *et al.*, 1996 ; Chevallier *et al.*, 1998a). The results are discussed in section 5, followed by conclusions in section 6.

## 2 The TIGR database

### 2.1 General description

The TIGR database consists of several subgroups, to serve as a database upon which several direct and inverse radiative transfer models rely. Each subgroup describes a particular aspect of the archived atmospheric soundings.

The subgroup of the geophysical parameters is the kernel of the database. To each sounding are associated a vertical temperature profile, a vertical water vapour concentration profile and a vertical ozone concentration profile. The vertical discretization refers to 40 levels between 0.05 and 1013 *hPa*, as reported in table 1. It has been chosen in the context of TOVS retrievals: for instance, the relatively coarse resolution in the lower troposphere is coherent with the instrument specifications. All soundings come from observations: the geographic location and the date of these are also archived. Unlike the temperature and water vapour profiles, the ozone profiles come from climatologies.

Various radiative characteristics of these soundings have been computed. Initially, only TOVS-related quantities were archived: the theoretical radiances for different satellite viewing angles and surface pressures, together with the corresponding atmospheric transmission profiles. With the new generation of vertical sounders, similar quantities have also been computed with reference to the Special Sensor Microwave/ Temperature (SSM/T) (Rieu *et al.*, 1996), to the Advanced Microwave Sounding Unit (AMSU) (Cabrera-Mercader and Staelin, 1995), and to the *Infrared Atmospheric Sounder Interferometer* (IASI) (Aires *et al.*, 1998). LW radiative flux profiles have also been added to the database for its use in a fast model for the computation of atmospheric LW flux profiles (NeuroFlux: Chevallier *et al.*, 1998a).

For practical use, the database is also currently classified into five statistically homogeneous air-mass classes: tropical, mid-latitude 1, mid-latitude 2, polar 1 and polar 2 (Achard, 1991).

### 2.2 Three successive versions

Up to now, there has been three different versions of the TIGR database. The first one, TIGR-1, groups together 1207 soundings (Moulinier, 1983). The initial database from which they were sampled was limited to 6600 radiosoundings, which was shown not to be enough for the use of TIGR in the framework of the 3I scheme (Flobert *et al.*, 1991). It has been followed by TIGR-2: the updated database contains 1761 situations sampled from a much larger set of 80,000 radiosoundings (Achard, 1991; Escobar-Munoz, 1993). The sampling methods for TIGR-1 and TIGR-2 are similar and rely on a temperature criterion. They are described in section 3.2. The necessity of an improved representativity of water vapour concentration led to the improvement of the method, and to the second update of the database (TIGR-3: Chevallier *et al.*, 1998a). The method is described in section 3.2. TIGR-3 gathers 2311 soundings and has been successfully used within the frame of the NOAA/ *National Aeronautics and Space Agency* (NASA) Pathfinder programme, for the re-analysis of the 20 years of TOVS observations (Scott

*et al.*, 1999).

## 2.3 Application to the computation of longwave flux profiles

TIGR-3 has been used as a training database for the neural network-based LW radiative transfer model NeuroFlux. NeuroFlux is a highly parameterized scheme that has been designed for computing LW flux profiles significantly more rapidly than the presently existing wide band models. It relies on a series of artificial neural networks, as defined by Rumelhart *et al.* (1986) (Multi-Layer Perceptrons). Directly from the geophysical parameters (TOVS retrievals or GCM soundings), they compute the contribution of clear sky radiation and of every cloudy layer to the fluxes. This design is at the basis of such a fast code: a gain of about one order of magnitude compared to the current ECMWF operational scheme (Morcrette, 1991 ; Zhong and Haigh, 1995) has been observed. One training database per neural network is needed: each one is issued from the same database, namely TIGR-3, and differs from the other ones by the pressure level of a black cloud layer. Due to the multilayer grey body algorithm (e.g., Washington and Williamson, 1977; Räisänen, 1998) used in NeuroFlux in conjunction with the artificial neural networks, the fluxes computed by the model take into account clouds as semi-transparent grey bodies, and not as black bodies. For a complete description of NeuroFlux, the reader is referred to Chevallier *et al.* (1998a).

For the application of NeuroFlux in GCMs (Chevallier *et al.*, 1998b), difficulties arose mainly due to the insufficient description of the boundary layer in the vertical pressure grid on which the soundings are archived in TIGR-3 (see table 1). This induced systematic errors by NeuroFlux in the lower troposphere and therefore non negligible uncertainties in the GCM simulations using NeuroFlux to compute the LW radiative budget. As a consequence, it appeared that for that kind of application, the training database of NeuroFlux had to be redefined. The next section focusses at the TIGR sampling technique, in order to define an updated methodology for sampling the Earth's atmospheric profiles.

## 3 Choice of a sampling technique

### 3.1 The TIGR approaches

The successive TIGR databases were set up with similar two-step methods. The first step consists in filtering the infinity of possible profiles in the atmosphere, by gathering a high but representative of them. For example, 80,000 radiosonde reports have been used for TIGR-2. Let us call  $S$  this initial database. The sampling of  $S$  with a topological approach is the second step of the method. It relies on an index  $I$ , that measures the dissimilarity between two atmospheric situations. The process is iterative. At step one, a first atmospheric situation from  $S$  is randomly drawn and archived in a new set  $E$ . At step  $n$ , a  $n^{\text{th}}$  atmospheric situation is randomly drawn and archived in  $E$  if it is different enough from the already selected situations, relatively to criterion  $I$ . With that approach, the distribution of  $E$  over the space of the various atmospheric variables is smoother than that of  $S$ . In practice, restriction to some variables had to be made: for TIGR-1 and TIGR-2, index  $I$  only takes into account the information about vertical temperature, whereas the improvements in TIGR-3 mostly came from the combined

use of vertical temperature and water vapour concentration profiles.

### 3.2 Tests on a two dimensional problem

Assessing the quality of a sampling technique is rather difficult when dealing with a high dimensional space. In order to visualize the qualities of the TIGR sampling methods, a simplified problem has been studied. 450,000 tropical soundings have been used as initial set  $S$ : they had been originally gathered for the setting up of TIGR-3. Instead of processing all the geophysical variables, only the mean temperature and the water vapour content of the layer 850 - 1013  $hPa$  have been kept. Let us call  $x$  the first variable and  $y$  the second one.  $x$  goes from 265 to 310  $K$ , and  $y$  from weak values up to 4.5  $mm$ .

Four sampling experiments are carried out on this reduced set, using four different topological sampling methods. The results are evaluated with regards to the histograms of the sampled datasets. Indeed, with the aim of gathering a databank for regression parameter estimation, the histograms should be as regular as possible. Other applications of the sampling may lead to a different criterion, but this is not discussed here.

The first experiment reproduces the approach used for TIGR-1 and TIGR-2: only the information about temperature is taken into account. It will be referred to as A1. A1 uses the simple distance:

$$D_1(s_i, s_j) = (x_i - x_j)^2 \quad (1)$$

where  $s_i$  and  $s_j$  represent two soundings  $i$  and  $j$  of the database  $S$ .

The A1 dissimilarity index is the following: to be archived in  $E$ , a situation (different from the initial one) has to verify:

$$\text{Min}_{s_j \in E} D_1(s_i, s_j) > d \quad (2)$$

where  $d$  is an arbitrary parameter. Choosing  $d$  is equivalent to choosing  $n$ , the number of situations in the final database. In the following,  $n = 100$ .

Figure 1 shows the 100-class histograms of  $x$  and  $y$  in  $E$  after the sampling. An ideal method would have led to regular histograms: every class among the 100 would contain one and only one situation. This is nearly the case for  $x$ : 12 classes only are empty.  $y$  is far more irregularly distributed with 55 empty classes, and includes no situation with layered water vapour content in excess of 3  $mm$ .

The second experiment, A2, is inspired from the unsuccessful attempt from Escobar-Munoz (1993) to introduce the information about the water vapour concentration in the selection. The algorithm is the same than in A1, but  $D_1$  is replaced by  $D_2$ :

$$D_2(s_i, s_j) = \left(\frac{x_i - x_j}{\sigma_x}\right)^2 + \left(\frac{y_i - y_j}{\sigma_y}\right)^2 \quad (3)$$

where  $\sigma_x$  (resp.  $\sigma_y$ ) is the standard deviation of the variable  $x$  (resp.  $y$ ), computed in the TIGR-2 database. With this normalization, the quantities of equation (3) are of the same order of magnitude.

Figure 2 shows the resulting histograms for  $x$  and  $y$  corresponding to the 100 selected situations. Compared to A1, A2 improves to repartition of  $y$ , with only 37 empty classes and

a coverage of the highest values, but clearly degrades that of  $x$  with also 37 empty classes. Although the two quantities  $|\frac{x_i - x_j}{\sigma_x}|$  and  $|\frac{y_i - y_j}{\sigma_y}|$  are normalized, their respective variability is too different to allow a satisfactory dissimilarity index to be obtained from their sum (i.e.  $D_2$ ).

Thus, in a third approach, A3, the proximity recognition in the  $x$  space is separated from the one in the  $y$  space. This is the method used for TIGR-3 (Chevallier, 1998). For each situation  $s_i$  in selection phase, the two non Euclidian distances are defined:

$$D_x^m(s_i, E) = \text{Min}_{s_j \in E} \sqrt{\left(\frac{x_i - x_j}{\sigma_x}\right)^2} \quad (4)$$

$$D_y^m(s_i, E) = \text{Min}_{s_j \in E} \sqrt{\left(\frac{y_i - y_j}{\sigma_y}\right)^2} \quad (5)$$

In each space, this algorithm looks for the nearest neighbour, by computing the minimum distance of the considered situation to the ones already archived in  $E$ ,  $D_x^m$  and  $D_y^m$ . Then the two minimum distances, separately computed, are added and the criterion of the minimum distance is applied to the sum:

$$D_x^m(s_i, E) + D_y^m(s_i, E) > d \quad (6)$$

This criterion is more selective than the previous one, A2, because the sum of the minima is always smaller than the minimum sum of the distances. It enables to select a situation  $s_i$  in three cases: - if  $x_i$  is different enough from the  $x_j$ 's already archived, - or if  $y_i$  is different enough from the  $y_j$ 's already archived, - or if  $x_i$  and  $y_i$  are different enough from those already archived, even though none of the two differences is outstanding.

100 situations have been selected with A3. The histograms are shown in figure 3. 22 classes are empty for  $x$  and 27 for  $y$ . This shows improvements compared to A2, with a more regular spread of the two variables.

Finally, a last criterion has been tried: A4. It uses two separate proximity recognitions, as in A3, and a combined use of  $x$  and  $y$ , as in A2. A4 uses the two distances:

$$D^1 = D_x^m \cdot (D_x^m + D_y^m) \quad (7)$$

$$D^2 = D_y^m \cdot (D_y^m + D_x^m) \quad (8)$$

For clarity, the dependence of  $D^1$ ,  $D^2$ ,  $D_x^m$ , and  $D_y^m$  as a function of  $s_i$  and  $E$  has not been written (see equations (4) and (5)).  $D^1$  (respectively  $D^2$ ) is highly conditioned by  $D_x^m$  (respectively  $D_y^m$ ), but also takes  $D_y^m$  (respectively  $D_x^m$ ) into account.

The A4 selection criterion is:

$$D^1(s_i, E) > d \quad (9)$$

$$D^2(s_i, E) > d \quad (10)$$

If both are satisfied,  $s_i$  is selected, otherwise it is rejected.

This algorithm has been used to set up a 100 situation database, from the 450,000. The histograms (figure 4) are comparable to those from A3. 23 classes are empty both for  $x$  and  $y$ .

It can be noted that A4 selects less extreme values. For example, with A3, 11 situations are characterized by  $x < 275 K$ , only 7 with A4.

The results from A3 and A4 illustrate the non-uniqueness of the choice of a satisfactory sampling method. A good approach results from a compromise between the regular spread of  $x$  and that of  $y$ . The highest peak of the  $x$  histogram appears in the highest values of  $x$ , corresponding to the highest variability of  $y$ . Similarly, both A3 and A4 tend to select the situations in the weak values of the layered water vapour  $y$ , corresponding to the highest variability of the temperature  $x$ . In the following, the A3 process, that had already been chosen for TIGR-3, was preferred to A4 because of the poorer selection of extreme values noted for A4.

### 3.3 Update of the TIGR-3 sampling technique

In the previous example, only tropical-type situations have been considered. A wider range of soundings would not significantly change the conclusions, though it is obvious that the peaks of the histograms would be more pronounced: for instance, adding polar-type situations would increase the temperature variability in the weak values of the water vapour. In order to smooth out the non-symetries, the initial heterogeneous database  $S$  can be divided into subgroups that are more homogeneous, each subgroup is then sampled separately from the others. This approach is called "stratified sampling" (e.g., Cochran, 1977). The values of the  $\sigma_x$ 's and of the  $\sigma_y$ 's in equations (4) and (5) have to be adapted for each subset. For TIGR-3, the initial database was divided into tropical and non-tropical situations. The classification relied on a statistical analysis performed on TIGR-2 (Achard, 1991). Instead, as the natural variability of water vapour is much more heterogeneous than that of the temperature, and increases with the total water vapour content, the initial database can be divided into subgroups, with respect to water vapour contents only. For instance, one subset may contain the lowest values of the layered water vapour, a second one the highest values, and a third one the intermediate values. Then the final sampled database  $E$  results from the merging of the three sampled subgroups. Such a division enables to clearly separate between different ranges of variability of the water vapour.

Extrapolating results from the former two-dimensional problem to a  $2 \times N$ -dimensional problem, if  $N$  is the number of vertical layers in the soundings, is not trivial. Extending the A3 approach to vertical profiles, one can replace the distances given by equations (4) and (5) by:

$$D_{\theta}^m(s_i, E) = \text{Min}_{s_j \in E} \sqrt{\sum_{k=1}^N \left( \frac{\theta_i(k) - \theta_j(k)}{\sigma_{\theta}(k)} \right)^2} \quad (11)$$

$$D_w^m(s_i, E) = \text{Min}_{s_j \in E} \sqrt{\sum_{k=1}^N \left( \frac{w_i(k) - w_j(k)}{\sigma_w(k)} \right)^2} \quad (12)$$

where  $\theta(k)$  and  $w(k)$  respectively are the mean temperature and mean water vapour concentration in layer  $k$ . They correspond to variables  $x$  and  $y$  of equations (4) and (5).

The criterion of the minimum distance of equation (6) can be rewritten:

$$D_{\theta}^m(s_i, E) + \mu D_w^m(s_i, E) > d \quad (13)$$

where  $\mu$  is a weight.

The standard deviations  $\sigma_{\theta}(k)$  and  $\sigma_w(k)$  in equations (11) and (12) deal with the different ranges of the variables, whereas  $\mu$  is introduced as a weight that takes into account the difference in the vertical variability of temperature on the one hand, and of water vapour on the other. In the following  $\mu = 1/9$ . Tests have shown that the characteristics of the sampled database do not depend critically on the precise value of  $\mu$ : for example,  $\mu = 1/8$  gives results comparable to  $\mu = 1/9$ .

This technique is similar to the TIGR-3 approach, except that the TIGR-3 sampling took layered water vapour contents into account, rather than water vapour mixing ratios. This choice was influenced by the previous definition of a water vapour distance by Flobert *et al.* (1986) for pattern recognition, but induced an arbitrary screening of the information about water vapour.

## 4 A new database from the ECMWF model outputs

Given a large database of samples  $S$ , covering a wide range of atmospheric water vapour and temperature profiles, the approach described above enables the selection of a smaller sample  $E$ , the size of which is determined by the factor  $d$ . It is used here for the sampling of profiles generated from the ECMWF atmospheric model. In this application,  $S$  results from the aggregation of six days of profiles from the ECMWF short-range forecasts. The six days, namely the first day of the months of January, March, May, July, September and November 1997, include a complete description of the atmosphere on a 31-layer vertical grid from the top of the atmosphere to the surface and an horizontal  $1.125^{\circ} \times 1.125^{\circ}$  grid representation every six hours. The vertical grid is illustrated on table 2.

$S$  is divided into seven subgroups differing by the total precipitable water vapour content of the profiles: the first group ranges from 0 to 0.5 *cm*, the second from 0.5 to 1.5 *cm*, the third from 1.5 to 2.5 *cm*, and so on, until the seventh one that goes from 5.5 *cm* up to the highest values. Preliminary experiments revealed that the representativity of the first group was insufficient. Therefore, data from the first day of the months of February, April, June, August, October and December, with total water vapour contents below 0.5 *cm* (i.e. 150,000 data), are also added to the initial database  $S$ , that consists of 1,350,000 profiles. In each group, the standard deviations, i.e. the  $\sigma_{\theta}$ 's and the  $\sigma_w$ 's of equations (11) and (12), are directly computed. The sampling approach described in section 3.3 is used for the extraction of about 750 samples from each class, except for the first one, where 1500 profiles are extracted, in consideration of the higher temperature variability: this class includes all types of situations from polar to tropical. The whole sampled database includes about 6000 profiles.

Figure 5 shows the histograms of the sampled database  $E$  in layer 6, characterized by a mean pressure of about 800 *hPa* when the surface pressure equals 1000 *hPa*. Two symmetric peaks appear as on figure 4. As explained in section 3.2, the wing in the temperature (respectively water vapour) histogram, between 220 and 270 *K* (resp. 0.002 and 0.012 *g/g*), illustrates the weak variability of water vapour (respectively temperature) in this temperature (resp. water



vapour) range in the initial set  $S$ . Since the stratified sampling forced the representation of high water vapour contents (see section 3.3), the wing in the water vapour histogram is more regular than that of the temperature histogram.

The following section describes the application of this new database as a training database for NeuroFlux in replacement of TIGR-3 (see section 2.2).

## 5 Application to LW radiative flux modelling

### 5.1 The training datasets

All the input files of NeuroFlux training databases are derived from  $E$ . In addition to the temperature and water vapour profiles in database  $E$ , other variables are required for LW radiative flux computations: cloud and surface characteristics, ozone profiles and the mean  $CO_2$  concentration. Ozone profiles are obtained from the climatology of Fortuin and Langematz (1994).  $CO_2$  concentration, surface temperature and LW emissivity are obtained by random sampling (within a given range of variation), ensuring regular distributions that are essential for the neural network training. The training database of the clear sky neural network includes the 6000 profiles without any cloud, whereas the other training databases include the same 6000 profiles, associated with the presence of a black body in a particular layer.

The output files of NeuroFlux training databases are the LW fluxes from the top of the atmosphere to the surface, that are associated to those profiles. In the present study, the current ECMWF operational LW code (Morcrette, 1991 ; Zhong and Haigh, 1995) is used to compute them. In the following, this code will be referred to as EC-OPE.

It should be noted that the variability of databases like NeuroFlux training sets, is far larger than the expected changes of the atmospheric variables during the next century (e.g., Houghton *et al.*, 1990).

### 5.2 Validation with code-by-code comparisons

The accuracy of NeuroFlux has been tested on various data from either the ECMWF analyses or the ECMWF forecasts, for different periods of time. Examples are given here of comparisons based on the re-analysis archives (Gibson *et al.*, 1997) for the first of December 1987, and on the short-range forecasts archives for the first of June 1998. Corresponding to these two dates, radiative fluxes obtained either by using NeuroFlux or EC-OPE have been compared for the whole globe. The radiative computations were for no cloud. Global data for the four synoptic times (00, 06, 12 and 18 UTC) at an horizontal resolution of  $1.125^\circ \times 1.125^\circ$  for the re-analysis archives or  $0.5625^\circ \times 0.5625^\circ$  for the 1998 forecasts were taken into account in the statistics: i.e. 200,000 atmospheric situations for the re-analysis archives and 800,000 for the forecasts. Results are presented for three latitude classes. The tropical class covers the  $30^\circ N$  -  $30^\circ S$  region. The mid-latitude class covers the  $30$ - $60^\circ N$  and  $30$ - $60^\circ S$  regions. The polar class covers North of  $60^\circ N$  and South of  $60^\circ S$ . For the three latitude classes, biases and standard deviations of the differences between the radiative calculations of NeuroFlux and those of EC-OPE were computed, as well as the maximum absolute differences. Since this version of NeuroFlux simulates EC-OPE, the differences are expected to be as small as possible.

Results for LW cooling rates are shown on figures 6 and 7. The two sets of comparisons indicate a similar behaviour for NeuroFlux, even though the nature of the profiles (analyses or forecasts), the version of the ECMWF GCM used for obtaining them (the 1995 13r4 cycle for the analyses and the 1998 18r5 cycle for the forecasts) and the horizontal resolution differ. The standard deviations and the biases in absolute value are less than  $0.3 K.d^{-1}$  and  $0.4 K.d^{-1}$  respectively, except in the lowest layer where the bias reaches  $0.6 K.d^{-1}$  in the polar class, and the standard deviation is around  $1.6 K.d^{-1}$  in the three classes. The maximum absolute error reaches  $10 K.d^{-1}$  both in the middle troposphere in the tropical class and in the lowest layer in the three classes. The higher standard deviation for the cooling rates in the lowest layer also exists when the test is performed on the training database (result not shown). Sensitivity tests have shown that it originates from the correlation between surface temperature and surface-air temperature. Such a problem could be reduced with more complex neural networks, but the code would then be computationally less efficient.

Similar statistics are presented for the outgoing longwave radiation (OLR) and the surface net flux in tables 3 and 4. The surface net flux is defined as the upward flux minus the downward flux, both positively defined. Standard deviations and absolute biases are less than  $2.0 W.m^{-2}$  and  $2.8 W.m^{-2}$  respectively. Maximum values reach  $8 W.m^{-2}$  for the OLR and  $10 W.m^{-2}$  for the surface net flux. Since the cloudy sky neural networks work on restricted parts of the atmospheric column, they have fewer inputs, fewer outputs and therefore are given fewer neurons. This leads to faster convergence and slightly better results than those presented here for clear sky (Chevallier, 1998): in particular the maximum errors are reduced (results not shown).

The most erroneous cases for NeuroFlux obviously correspond to situation types that are less represented in the initial 1,350,000 sounding database: in particular the high elevations in the Himalayas. Further reduction of the error maxima will lead to the extension of the initial database.

To the authors knowledge, no radiative transfer scheme has been validated with reference computations on such a high number of profiles. In particular, the maximum errors performed by EC-OPE, with reference to the real values, are not known. The previous version of NeuroFlux, using TIGR-3 in the training phase and a 19-layer vertical grid has been validated on a 1032 radiosonde database and on a 15000 sounding database from the LMD GCM: the biases and standard deviations were comparable to those shown here (Chevallier, 1998). But here the maximum errors are significantly reduced in the stratosphere and in the lower layers.

With the new sampled database, NeuroFlux has been adapted to a higher vertical resolution than previously, while simultaneously increasing its accuracy. Compared to EC-OPE, this version of NeuroFlux is 7 times faster. The performances of this new version in the framework of GCM simulations is discussed in a companion paper (Chevallier *et al.*, 1999). In particular, the latter examines the effect of the higher error of NeuroFlux in the lowest atmospheric layer.

## 6 Conclusions and prospects

Collecting pertinent - i.e. properly sampled - a priori information about the atmosphere is an essential need for atmospheric modelling. An important step has been the constitution of the TIGR database. In this paper, the characteristics of successive versions developed at LMD



since 1983 were summarized. The recent application to the computation of longwave radiative flux profiles in GCMs has brought evidence that there was a need of further improvements of the sampling strategy. Therefore a flexible sampling method was defined, based on the approach used in the latest TIGR database, TIGR-3. Given a large database of atmospheric situations, covering a wide range of atmospheric temperature and water vapour profiles, the technique described enables selecting a smaller sample, defined as a regular mesh of the initial database. An application to the sampling of atmospheric situations from the ECMWF forecast model outputs has been presented. The sampled database has been used for inferring the parameters of artificial neural networks, in a longwave flux profile computation model (NeuroFlux: Ch ery *et al.*, 1996, Chevallier *et al.*, 1998a). With this new database, NeuroFlux has been extended to a higher vertical resolution than previously, while simultaneously increasing its accuracy in terms of code-by-code comparisons. Further results in the ECMWF GCM are presented in a companion paper (Chevallier *et al.*, 1999).

The sampling method described here may be applied to the development of TIGR-type databases adapted to applications involving a higher vertical variability than the present TIGR database. An example is the adaptation of NeuroFlux to the forthcoming 60-level ECMWF model, with increased vertical resolution in both the stratosphere and the boundary layer. In the same way, other databases are needed for application to the retrieval of thermodynamic variables from the *Infrared Atmospheric Sounder Interferometer* (IASI) and the *Advanced Infrared Radiometric Sounder* (AIRS) instruments. They could simply be derived by interpolating the 60-level ECMWF database, once it is available, to a coarser grid.

## Acknowledgments

Authors are grateful to N. A. Scott and R. Armante for substantial help in the design of the TIGR-3 database at LMD. They also wish to thank M. Miller for careful review of the manuscript.

## References

- Achard, V., 1991: Trois problèmes clés de l'analyse 3D de la structure thermodynamique de l'atmosphère par satellite : mesure du contenu en ozone ; classification de masses d'air ; modélisation hyper rapide du transfert radiatif. PhD thesis, University Paris VI, 168 pp. [Available from LMD, Ecole Polytechnique, 91128 Palaiseau cedex, France].
- Aires, F., R. Armante, A. Chédin, and N. A. Scott, 1998: A regularised neural net approach to the surface and the atmospheric temperature retrieval with the high resolution interferometer IASI. In *Proceedings of the Am. Meteor. Soc. Conference*, Paris, France.
- Cabrera-Mercader, C. R., and D. H. Staelin, 1995: Passive microwave relative humidity retrievals using feedforward neural networks. *IEEE Trans. Scien. Rem. sens.*, **33:6**, 1324-1326.
- Chaboureaud, J.-P., Chédin, A., and N. A. Scott, 1998 : Remote Sensing of The Vertical Distribution of Atmospheric Water Vapor From the TOVS Observations. Method and Validation. *J. Geophys. Res.*, **103**, 8743-8752.
- Chédin, A., N. A. Scott, C. Wahiche and P. Moulinier, 1985: The Improved Initialization Inversion method : a high resolution physical method for temperature retrievals from satellites of the TIROS-N series. *J. Climate Appl. Meteor.*, **24**, 128-143.
- Chéruy, F., F. Chevallier, J.-J. Morcrette, N. A. Scott, A. Chédin, 1996 : Une méthode utilisant les techniques neuronales pour le calcul rapide de la distribution verticale du bilan radiatif thermique terrestre. *C. R. Acad. Sci. Paris*, **322:IIb**, 665-672, in French.
- Chevallier, F., 1998: La modélisation du transfert radiatif à des fins climatiques : une nouvelle approche fondée sur les réseaux de neurones artificiels. PhD thesis, University Paris VII, 230 pp. [Available from LMD, Ecole Polytechnique, 91128 Palaiseau cedex, France].
- Chevallier, F., F. Chéruy, N. A. Scott, and A. Chédin, 1998a: A neural network approach for a fast and accurate computation of longwave radiative budget. *J. Appl. Meteor.*, **37**, 1385-1397.
- Chevallier, F., F. Chéruy, Z. X. Li, and Scott, N. A., 1998b: A fast and accurate neural network-based computation of longwave radiative budget application in a GCM. *Proceedings of the Am. Meteor. Soc. Conference*, Paris, France, in press.
- Chevallier, F., J.-J. Morcrette, F. Chéruy, and N. A. Scott, 1999: Use of a neural network-based longwave radiative transfer scheme in the ECMWF atmospheric model. Submitted to *Quart. J. Roy. Meteor. Soc.*
- Cochran, W.G., 1977: Sampling techniques, third edition, John Wiley & Sons, 428 pp.
- Escobar-Munoz, J., 1993 : Base de données pour la restitution de variables atmosphériques à l'échelle globale. Étude sur l'inversion par réseaux de neurones des données des sondeurs verticaux atmosphériques satellitaires présents et à venir. PhD thesis, Univ. Paris VII, 190 pp. [Available from LMD, Ecole Polytechnique, 91128 Palaiseau cedex, France].
- Escobar-Munoz, J., A. Chédin, F. Chéruy, and N. A. Scott, 1993: Réseaux de neurones multi-couches pour la restitution de variables thermodynamiques atmosphériques à l'aide de sondeurs verticaux satellitaires. *C. R. Acad. Sci. Paris*, **317:IIb**, 911-918, in French.



- Flobert, J.-F., E. Andersson, A. Chédin, A. Hollingsworth, G. Kelly, J. Pailleux, and N. A. Scott, 1991: Global data assimilation and forecast experiments using the Improved Initialization Inversion method for satellite soundings. *Mon. Wea. Rev.*, 119:8, 1881-1914.
- Flobert, J.-F., N. A. Scott, and A. Chédin, 1986: A fast model for TOVS radiances computation. In *Proceedings of the 6th conference on atmospheric radiation*. Williamsburg, USA, p. 186-189.
- Fortuin, J. P. F. and Langematz, U., 1994: An update on the global ozone climatology and on concurrent ozone and temperature trends. *Proceedings SPIE*, 2311, 207-216.
- Gibson, J. K., P. Källberg, S. Uppala, A. Hernandez, A. Nomura, and E. Serrano, 1997: ECMWF Re-analysis. 1. ERA description. ECMWF Project Report Series.
- Houghton, J.T., G.J. Jenkins, and J.J. Ephraums, Eds, 1990: Climate change. The IPCC scientific assessment. Cambridge University Press.
- Morcrette, J.J., 1991 : Radiation and Cloud Radiative Properties in the European Center for Medium Range Weather Forecasts forecasting system. *J. Geophys. Res.*, **96:D5**, pp 9121-9132.
- Moulinier, P., 1983: Analyse statistique d'un vaste échantillonnage de situations atmosphériques sur l'ensemble du globe. LMD Internal note 123, 30 pp., in French [Available from LMD, Ecole Polytechnique, 91128 Palaiseau cedex, France].
- Räsänen, P., 1998: Effective longwave cloud fraction and maximum-random overlap clouds - a problem and a solution. accepted for publication in the *Mon. Wea. Rev.*
- Rumelhart, D. E., G. E. Hinton, and R. J., Williams, 1986: Learning internal representations by error propagation. *Parallel distributed processing: Explorations in the macrostructure of cognition 1*, D. E. Rumelhart and McClelland, Eds., MIT Press, 318-362.
- Scott, N. A., A. Chédin, R. Armante, J. Francis, C. Stubenrauch, J.-P. Chaboureau, F. Chevallier, C. Claud, and F. Chérury, 1998: Characteristics of the TOVS Pathfinder Path-B database. *Submitted to the Bull. Amer. Meteo. Soc.*
- Washington, W. M. and D. L. Williamson, 1977 : A description of the NCAR GCM's in General circulation models of the atmosphere. *Methods in Computational Physics*, J. Chang. Ed., 17, Academic Press, 111-172.
- Zhong, W., and J. D. Haigh, 1995: Improved broadband emissivity parameterization for water vapor cooling rate calculations. *J. Atmos. Sci.*, **52:1**, 124-138.



level	pressure (hPa)	level	pressure (hPa)	level	pressure (hPa)	level	pressure (hPa)
1	0.05	11	7.43	21	131.20	31	471.86
2	0.09	12	11.11	22	161.99	32	525.00
3	0.17	13	16.60	23	200.00	33	584.80
4	0.30	14	24.73	24	222.65	34	651.04
5	0.55	15	37.04	25	247.90	35	724.78
6	1.00	16	45.73	26	275.95	36	800.00
7	1.50	17	56.46	27	307.20	37	848.69
8	2.23	18	69.71	28	341.99	38	900.33
9	3.33	19	86.07	29	380.73	39	955.12
10	4.98	20	106.27	30	423.86	40	1013.00

Table 1: The vertical grid on which the TIGR databases are archived.



level	pressure (hPa)	level	pressure (hPa)	level	pressure (hPa)	level	pressure (hPa)
1	0.00	11	222.93	21	610.47	31	1005.18
2	20.00	12	253.69	22	656.28	32	1013.00
3	40.00	13	286.57	23	702.57		
4	60.00	14	321.46	24	748.95		
5	80.00	15	358.23	25	794.90		
6	100.16	16	396.75	26	839.75		
7	121.16	17	436.87	27	882.58		
8	143.63	18	478.45	28	922.23		
9	167.95	19	521.35	29	957.21		
10	194.35	20	565.42	30	985.63		

Table 2: Boundary pressures of the 31 layer-vertical grid of the ECMWF model, when the surface pressure equals 1013 *hPa*. The general formulation depends on surface pressure.



(a)

	polar			mid-latitude			tropical		
	m	$\sigma$	M	m	$\sigma$	M	m	$\sigma$	M
NeuroFlux - EC-OPE	-0.92	1.56	7.48	0.34	-1.21	6.11	-0.75	1.20	8.05

(b)

	polar			mid-latitude			tropical		
	m	$\sigma$	M	m	$\sigma$	M	m	$\sigma$	M
NeuroFlux - EC-OPE	-0.74	2.73	8.78	-0.61	1.58	6.98	-0.25	1.08	6.84

Table 3: Mean (m), standard deviation ( $\sigma$ ) and maximum absolute error (M) of the comparisons between NeuroFlux and EC-OPE for the computation of the OLR (a), and the net flux at the surface (b). Clouds were not taken into account in the computations. ECMWF re-analyses. 1<sup>st</sup> December 1987, 00, 06, 12 and 18 UTC.  $1.125^\circ \times 1.125^\circ$  horizontal resolution. Fluxes in  $W.m^{-2}$ . Results are shown by latitude class.



(a)

	polar			mid-latitude			tropical		
	m	$\sigma$	M	m	$\sigma$	M	m	$\sigma$	M
NeuroFlux - EC-OPE	0.35	1.07	4.94	0.00	1.14	5.38	-0.44	1.16	7.17

(b)

	polar			mid-latitude			tropical		
	m	$\sigma$	M	m	$\sigma$	M	m	$\sigma$	M
NeuroFlux - EC-OPE	-0.40	1.86	8.07	-0.64	1.36	6.73	0.06	1.15	9.34

Table 4: Mean (m), standard deviation ( $\sigma$ ) and maximum absolute error (M) of the comparisons between NeuroFlux and EC-OPE for the computation of the OLR (a), and the net flux at the surface (b). Clouds were not taken into account in the computations. ECMWF 6-hour forecasts. 1<sup>st</sup> June 1998, 00, 06, 12 and 18 UTC.  $0.5625^\circ \times 0.5625^\circ$  horizontal resolution. Fluxes in  $W.m^{-2}$ . Results are shown by latitude class.

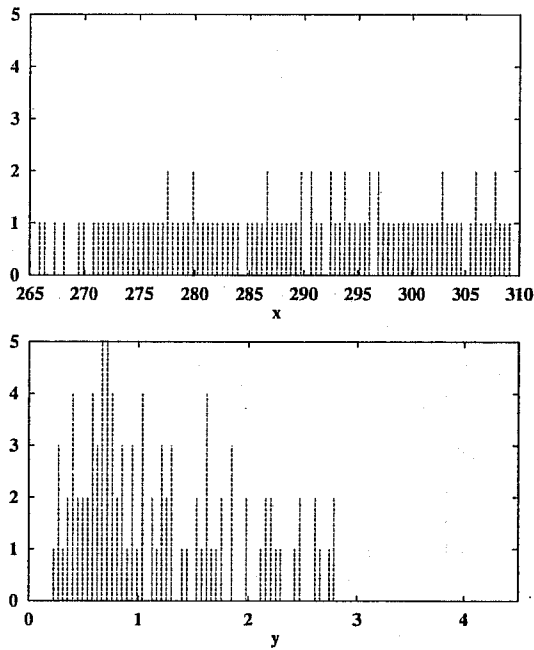


Figure 1: Histograms (100 classes) of  $x$  and  $y$  after the A1 sampling.  $x$  in  $K$ ,  $y$  in  $mm$ .

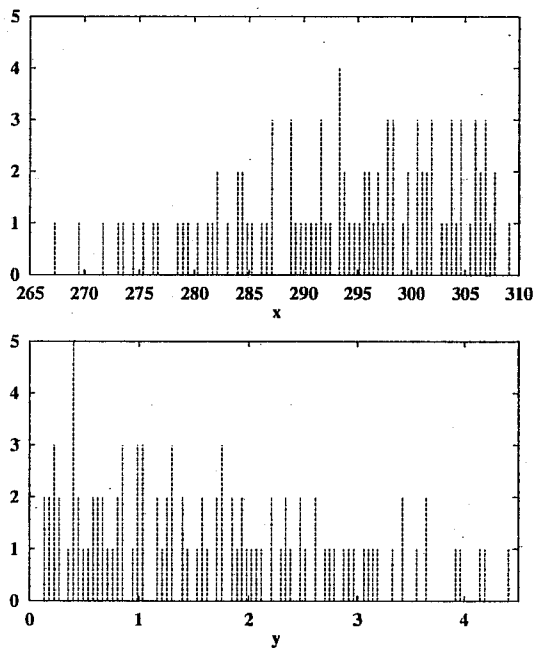


Figure 2: Histograms (100 classes) of  $x$  and  $y$  after the A2 sampling.  $x$  in  $K$ ,  $y$  in  $mm$ .

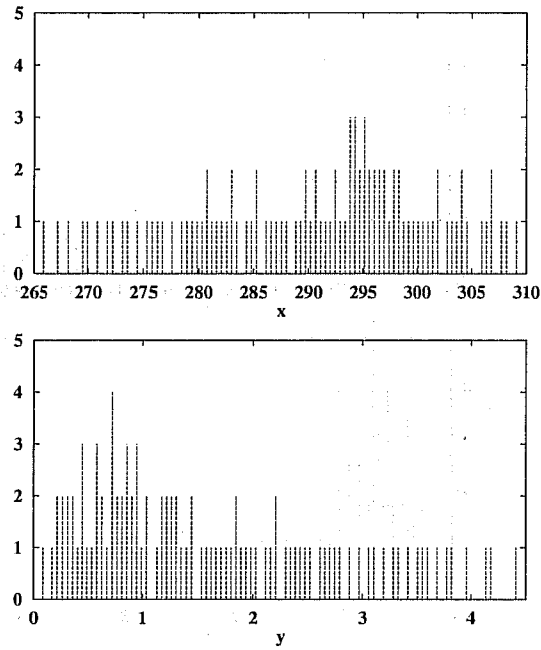


Figure 3: Histograms (100 classes) of  $x$  and  $y$  after the A3 sampling.  $x$  in  $K$ ,  $y$  in  $mm$ .

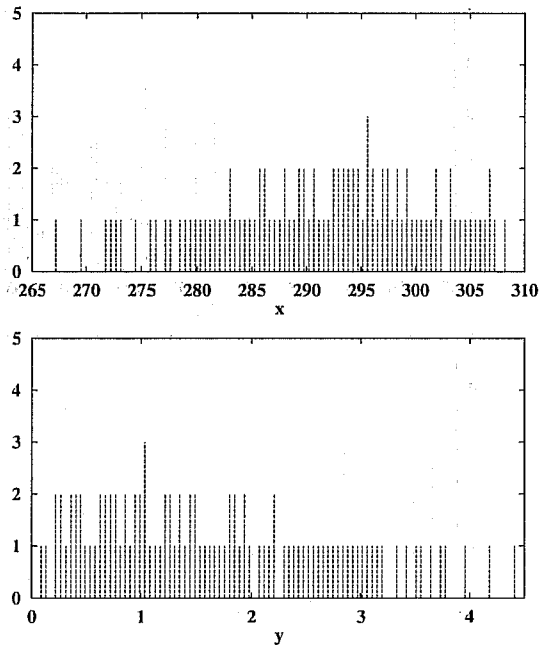


Figure 4: Histograms (100 classes) of  $x$  and  $y$  after the A4 sampling.  $x$  in  $K$ ,  $y$  in  $mm$ .

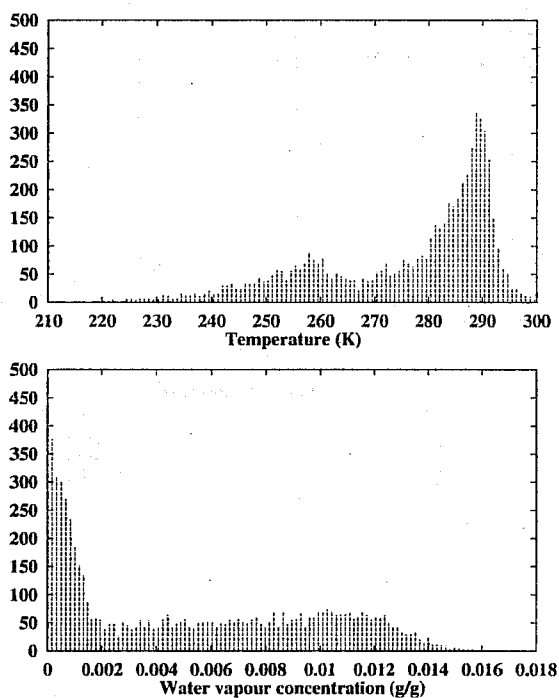


Figure 5: Histograms (100 classes) of temperature (in  $K$ ) and water vapour concentration (in  $g/g$ ) in the database sampled using criterion A3 (section 3.2) from ECMWF generated profiles. Layer 6, characterized by a mean pressure of about  $800 hPa$  when the surface pressure equals  $1000 hPa$ .

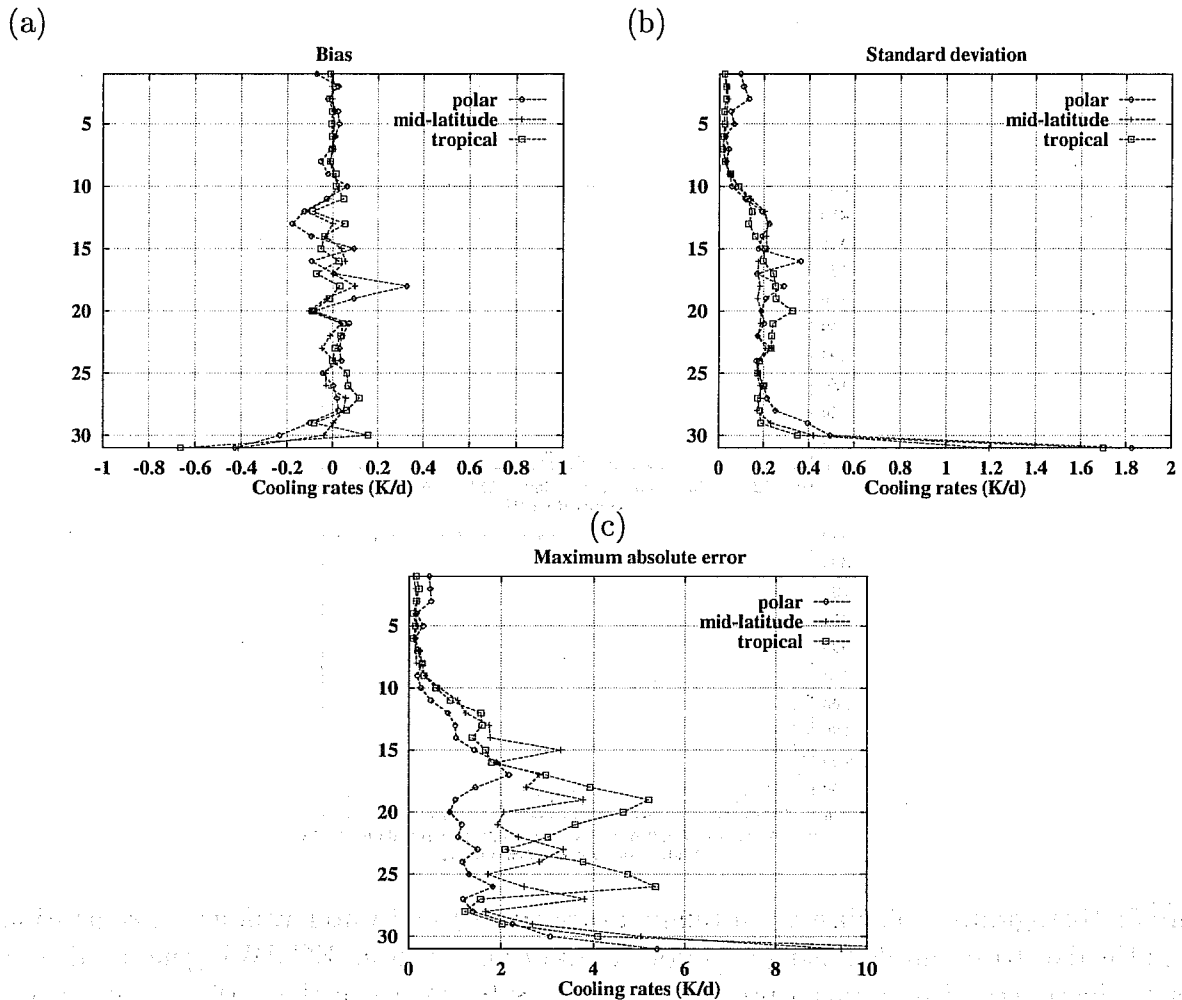


Figure 6: Comparison between the computations of NeuroFlux and those of EC-OPE: cooling rates from NeuroFlux minus cooling rates from EC-OPE, in  $K.d^{-1}$ . Clouds were not taken into account in the computations. ECMWF 6-hour forecasts. 1<sup>st</sup> December 1987, 00, 06, 12 and 18 UTC.  $1.125^\circ \times 1.125^\circ$  horizontal resolution. Fluxes in  $W.m^{-2}$ . Results are shown by latitude class.

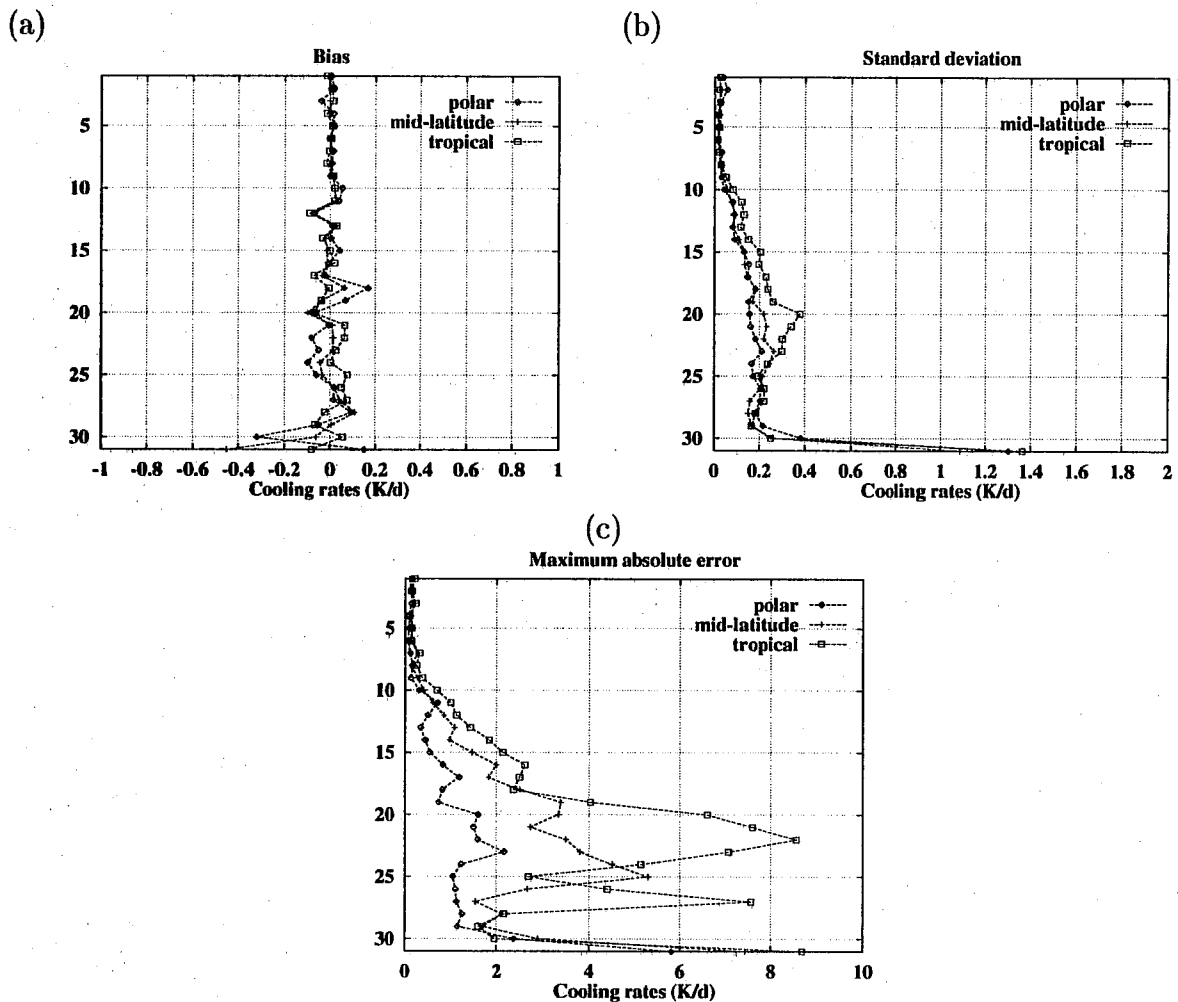
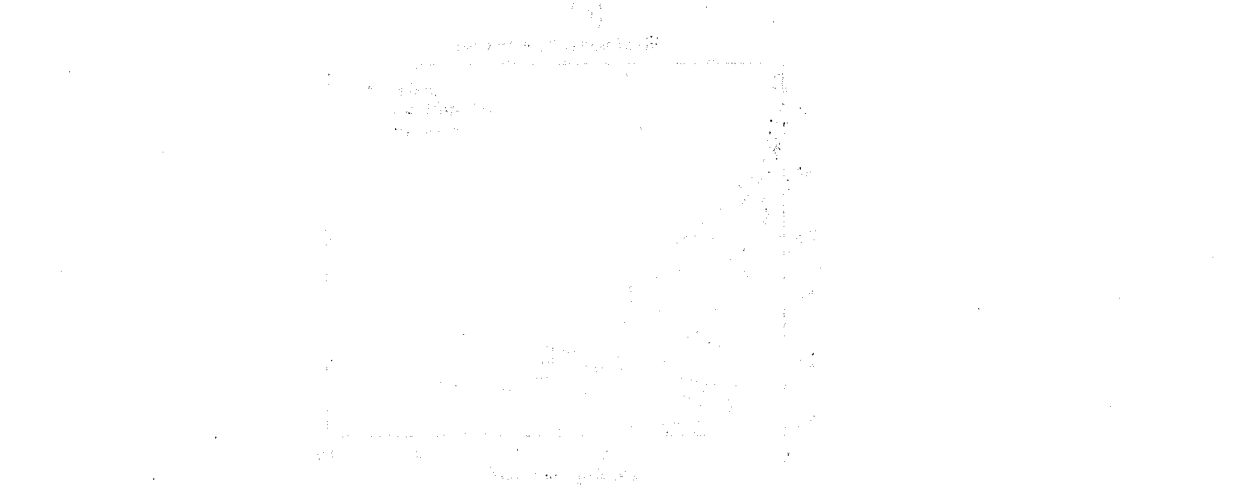
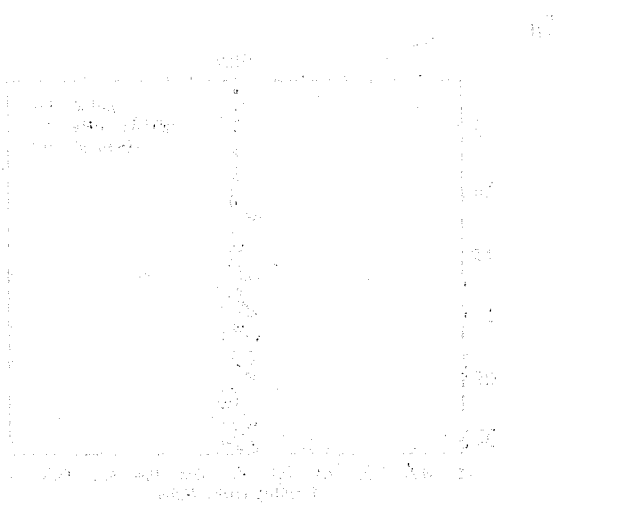
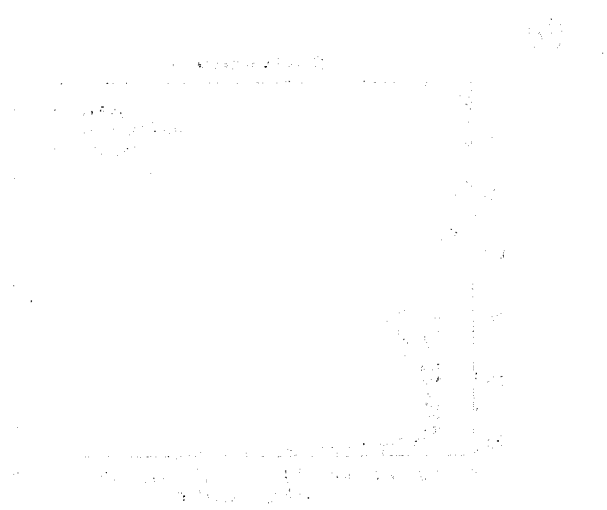


Figure 7: Comparison between the computations of NeuroFlux and those of EC-OPE: cooling rates from NeuroFlux minus cooling rates from EC-OPE, in  $K.d^{-1}$ . Clouds were not taken into account in the computations. ECMWF re-analyses. 1<sup>st</sup> June 1998, 00, 06, 12 and 18 UTC.  $0.5625^\circ \times 0.5625^\circ$  horizontal resolution. Fluxes in  $W.m^{-2}$ . Results are shown by latitude class.

Bar chart showing the number of students who took part in different sports. The Y-axis represents the number of students (0 to 100). The X-axis represents the sports (Football, Basketball, Tennis, Badminton, Table Tennis, and Table Tennis).



2. The following table shows the number of students who took part in a school sports day. The data is presented in a bar chart. Complete the table below.

Bar chart showing the number of students who took part in different sports. The Y-axis represents the number of students (0 to 100). The X-axis represents the sports (Football, Basketball, Tennis, Badminton, Table Tennis, and Table Tennis).