

## International Challenge to Predict the Impact of Radionuclide Releases from Medical Isotope Production on a Comprehensive Nuclear Test Ban Treaty Sampling Station

Paul W. Eslinger<sup>a</sup>, Ted W. Bowyer<sup>a</sup>, Pascal Achim<sup>b</sup>, Tianfeng Chai<sup>c</sup>, Benoit Deconninck<sup>d</sup>, Katie Freeman<sup>e</sup>, Sylvia Generoso<sup>b</sup>, Philip Hayes<sup>f</sup>, Verena Heidmann<sup>g</sup>, Ian Hoffman<sup>h</sup>, Yuichi Kijima<sup>i</sup>, Monika Krysta<sup>j</sup>, Alain Malo<sup>k</sup>, Christian Maurer<sup>l</sup>, Fantine Ngan<sup>c</sup>, Peter Robins<sup>e</sup>, J. Ole Ross<sup>m</sup>, Olivier Saunier<sup>n</sup>, Clemens Schlosser<sup>g</sup>, Michael Schöppner<sup>o</sup>, Brian T. Schrom<sup>a</sup>, Petra Seibert<sup>p</sup>, Ariel F. Stein<sup>c</sup>, Kurt Ungar<sup>h</sup>, Jing Yi<sup>h</sup>

### Author Affiliations:

<sup>a</sup> Pacific Northwest National Laboratory, 902 Battelle Boulevard, Richland, WA, 99352 USA

<sup>b</sup> Commissariat à l'Énergie Atomique, CEA, DAM, DIF, 91297 Arpajon, France

<sup>c</sup> NOAA/Air Resources Laboratory, College Park, Maryland

<sup>d</sup> Institut des Radioéléments, Fleurus, Belgium

<sup>e</sup> AWE, Aldermaston, Reading, RG7 4PR, United Kingdom

<sup>f</sup> Air Force Technical Applications Center, Patrick Air Force Base, Florida, USA

<sup>g</sup> Federal Office for Radiation Protection (*Bundesamt für Strahlenschutz, Bfs*), Freiburg, Germany

<sup>h</sup> Health Canada, Radiation Protection Bureau, Ottawa, Canada

<sup>i</sup> Japan Atomic Energy Agency, Tokai, Ibaraki, Japan

<sup>j</sup> Comprehensive Test Ban Treaty Organization (CTBTO), International Data Center, Vienna, Austria

<sup>k</sup> Environment Canada, Canadian Meteorological Centre, Dorval, Canada

<sup>l</sup> Zentralanstalt für Meteorologie und Geodynamik, Vienna, Austria

<sup>m</sup> Federal Institute for Geosciences and Natural Resources (BGR), Hannover, Germany

<sup>n</sup> French Institute for Radiation Protection and Nuclear Safety, Fontenay-aux-Roses, France

<sup>o</sup> Program on Science and Global Security, Princeton University, Princeton New Jersey USA

<sup>p</sup> University of Natural Resources and Life Sciences, Institute of Meteorology and University of Vienna, Department of Meteorology and Geophysics, Vienna, Austria

### Corresponding author:

Paul W. Eslinger

Tel.: +1 509 372 4392

E-mail address: paul.w.eslinger@pnnl.gov

Pacific Northwest National Laboratory, MSIN K7-76, 902 Battelle Boulevard, P.O. Box 999, Richland, Washington, USA

### Author Email Addresses:

Paul W. Eslinger paul.w.eslinger@pnnl.gov

Ted. W. Bowyer ted.bowyer@pnnl.gov

Brian T. Schrom brian.schrom@pnnl.gov

Pascal Achim pascal.achim@cea.fr

Tianfeng Chai tianfeng.chai@noaa.gov

Benoit Deconninck Benoit.Deconninck@ire-elit.eu

Sylvia Generoso sylvia.generoso@cea.fr

Philip Hayes philip.hayes.2@us.af.mil

Ian Hoffman ian.hoffman@hc-sc.gc.ca

Yuichi Kijima kijima.yuichi@jaea.go.jp

Monika Krysta monika.krysta@ctbto.org  
Christian Maurer christian.maurer@zamg.ac.at  
Fantine Ngan fantine.ngan@noaa.gov  
Peter Robins peter.robins@awe.co.uk  
Ole Ross ole.ross@bgr.de  
Olivier Saunier olivier.saunier@irsn.fr  
Clemens Schlosser cschlosser@bfs.de  
Michael Schoeppner schoeppner@princeton.edu  
Petra Seibert petra.seibert@univie.ac.at  
Ariel F. Stein ariel.stein@noaa.gov  
Kurt Ungar Kurt.Ungar@hc-sc.gc.ca

1 International Challenge to Predict the Impact of Radioxenon Releases from Medical Isotope Production  
2 on a Comprehensive Nuclear Test Ban Treaty Sampling Station

3 **Abstract**

4 The International Monitoring System (IMS) is part of the verification regime for the Comprehensive  
5 Nuclear-Test-Ban-Treaty Organization (CTBTO). At entry-into-force, half of the 80 radionuclide stations  
6 will be able to measure concentrations of several radioactive xenon isotopes produced in nuclear  
7 explosions, and then the full network may be populated with xenon monitoring afterward. An  
8 understanding of natural and man-made radionuclide backgrounds can be used in accordance with the  
9 provisions of the treaty (such as event screening criteria in Annex 2 to the Protocol of Treaty) for the  
10 effective implementation of the verification regime.

11 Fission-based production of  $^{99}\text{Mo}$  for medical purposes also generates nuisance radioxenon isotopes that  
12 are usually vented to the atmosphere. One of the ways to account for the effect emissions from medical  
13 isotope production has on radionuclide samples from the IMS is to use stack monitoring data, if they are  
14 available, and atmospheric transport modeling. Recently, individuals from seven nations participated in a  
15 challenge exercise that used atmospheric transport modeling to predict the time-history of  $^{133}\text{Xe}$   
16 concentration measurements at the IMS radionuclide station in Germany using stack monitoring data  
17 from a medical isotope production facility in Belgium. Participants received only stack monitoring data  
18 and used the atmospheric transport model and meteorological data of their choice.

19 Some of the models predicted the highest measured concentrations quite well. A model comparison rank  
20 and ensemble analysis suggests that combining multiple models may provide more accurate predicted  
21 concentrations than any single model. None of the submissions based only on the stack monitoring data  
22 predicted the small measured concentrations very well. Modeling of sources by other nuclear facilities

23 with smaller releases than medical isotope production facilities may be important in understanding how to  
24 discriminate those releases from releases from a nuclear explosion.

## 25 **Keywords**

26 Medical isotope production;  $^{133}\text{Xe}$ ; source-term estimation; atmospheric modeling; CTBTO

## 27 **1. Introduction**

28 The International Monitoring System (IMS) is part of the verification regime for the Comprehensive  
29 Nuclear-Test-Ban-Treaty Organization (CTBTO, 2014). The verification regime is designed to detect  
30 nuclear explosions no matter where they occur on the earth. When complete, 80 of the IMS stations will  
31 have aerosol measurement systems sensitive enough to detect releases from nuclear explosions at great  
32 distances. At entry-into-force, half of the 80 stations will also have equipment that measures  
33 concentrations of four radioactive xenon isotopes ( $^{131\text{m}}\text{Xe}$ ,  $^{133}\text{Xe}$ ,  $^{133\text{m}}\text{Xe}$ , and  $^{135}\text{Xe}$ ) produced in a nuclear  
34 explosion, and following entry-into-force, a plan to add xenon monitoring capabilities to the other 40  
35 stations will be reviewed (Comprehensive Nuclear-Test-Ban Treaty, 1996). An understanding of natural  
36 and man-made radionuclide backgrounds can also be used in accordance with the provisions of the treaty  
37 (such as event screening criteria in Annex 2 to the Protocol of Treaty) for the effective implementation of  
38 the verification regime.

39 A number of studies of the release and transport of radioxenon from nuclear explosions, nuclear power  
40 plants, and medical isotope production facilities have been published (Becker et al., 2010; Eslinger et al.,  
41 2014; Hoffman et al., 2009; Kalinowski et al., 2008; Saey et al., 2010b; Wotawa et al., 2010; Wotawa et  
42 al., 2003; Zähringer et al., 2009). These studies confirm that fission-based production of  $^{99}\text{Mo}$  for medical  
43 purposes is the largest routine contributor of radioxenon to worldwide background levels. The  $^{99}\text{Mo}$  (half-  
44 life of 66 hours) decays into  $^{99\text{m}}\text{Tc}$  (half-life of 6 hours) and the resulting  $^{99\text{m}}\text{Tc}$  is used in approximately

45 30-40 million medical procedures per year (Peykov and Cameron, 2014) and the demand is expected to  
46 increase in the future.

47 A reduction in radioxenon releases to relatively low levels (Bowyer et al., 2013) has the potential to  
48 reduce background radioxenon to levels that don't significantly impact treaty verification activities.  
49 However, medical isotope production facilities meet regulatory release requirements and their releases  
50 don't pose public health risks, thus the operators have no financial incentive to reduce releases. Another  
51 way of mitigating the impact on treaty verification activities is to use stack monitoring data, if they are  
52 available, and atmospheric transport modeling. In the modeling context, one could attempt to model  
53 background sources accurately enough to subtract a background contribution from any sampled value.  
54 Given the uncertainties (source terms, modeling), simulated peaks may not accurately represent reality.  
55 Thus, alternately, when a xenon peak is observed, one could check whether the simulated background  
56 increases during the same period (synchronization in time). If that is the case, the observed peak could be  
57 linked to the rise of the radioxenon background.

58 Unfortunately, the details of the stack monitoring data needed, such as the time resolution, the accuracy,  
59 and whether or not local weather data are needed is not well known. There have been questions about  
60 whether stack data would be useful in a practical way at all, depending on the type of data made available  
61 and when it could be made available from a producer. To date, only one published study (Schöppner et  
62 al., 2013) has addressed the impacts the time resolution of stack monitoring data have on predicted  
63 concentrations at an IMS station location. The minimum source term resolution considered in that study  
64 was one day. Atmospheric modeling studies using inert tracers have been conducted since the early  
65 1980s (Ferber et al., 1986; Gudiksen et al., 1984). This study addresses the difficult nuance of whether  
66 atmospheric models currently in wide use can yield information on the accuracy and timing of the source  
67 term data needed to faithfully reproduce sampling data.

68 This paper describes a challenge exercise formulated to start to answer some of these questions. Namely,  
69 to ascertain the level of agreement that can be achieved between atmospheric transport models using stack  
70 monitoring data and xenon isotopic concentration measurements at IMS stations. An evaluation criterion  
71 is used to measure the level of agreement. However, the real value of the exercise is in discussions  
72 resulting from the challenge without over-analyzing the evaluation criterion. The challenge is expected to  
73 spark discussions on what techniques are best, what gaps exist in our knowledge, and what type of data  
74 fidelity is needed from stack monitors. In general, this challenge will help inform the international treaty  
75 verification community of the status of the current capability.

76 The general approach of the exercise was to challenge atmospheric transport modeling groups to  
77 reproduce the time-history of  $^{133}\text{Xe}$  measurements at an IMS station using stack monitoring data from a  
78 medical isotope production facility. Participants received stack monitoring data that included the location,  
79 UTC date and time of releases, the measured activity concentrations of  $^{133}\text{Xe}$  in  $\text{Bq m}^{-3}$ , an average stack  
80 flow rate ( $80,000 \text{ m}^3 \text{ hr}^{-1}$ ), and the height (m above ground level) of the release. All other data were  
81 gathered by the participants. Each participant used the atmospheric transport model and the associated  
82 meteorological data of their choice. The individuals participating in the challenge are identified in Table  
83 1. Participants were asked not to use the IMS sampling data, if they had access to them, until after  
84 completing the modeling exercise.

## 85 **2. Atmospheric transport models and meteorological data**

86 The participants used several transport codes and several different sources for meteorological data.  
87 Several participants submitted results for more than one model. Some of the submissions were averages  
88 of other models or low and high resolution runs for the same model. Model metadata are provided in  
89 Table 2. Although the analysis considers all twenty six submissions, a subset of the submissions was  
90 selected to discuss common model characteristics. The reduced set of submissions is identified in the last  
91 column of Table 2. Some submissions are not specifically identified in Table 2. The submission Hof 3

92 was an average of the submissions Hof1 and Hof2. Submissions Sei4, Sei5 and Sei6 were slight  
93 variations, including different release height assumptions, on submissions Sei 1, Sei2, and Sei3. Ros2  
94 was a low resolution (smaller number of particles) version of submission Ros1 and Mau1 was a low  
95 resolution version of Mau3.

96 The participants used five different atmospheric transport models. The models, in order of the number of  
97 uses by participants are the following: FLEXPART (Stohl et al., 2005; Stohl et al., 1998), a Lagrangian  
98 particle dispersion model; HYSPLIT (Draxler and Hess, 1998, 2010) a hybrid single particle Lagrangian  
99 integrated trajectory model; Eulerian IdX (Tombette et al., 2014) which is part of IRSN's (French  
100 Institute for Radiation protection and Nuclear Safety) C3X operational platform; the Weather Research  
101 and Forecasting (WRF) model (Done et al., 2004; Michalakes et al., 2001) and MLDP0 (D'Amours et al.,  
102 2015; D'Amours et al., 2010) a Lagrangian particle dispersion model designed for long-range problems  
103 associated with events of regional, continental and global consequences.

104 The participants used six different meteorological data sets, some of which are available in different  
105 spatial and time resolutions. Meteorological analysis data are created by assimilation of a forecast model  
106 to observational data. Reanalysis data (i.e. GDAS) are produced later to have a consistent standardized  
107 gridded product of past weather patterns.

108 Thirteen of the submissions used global analysis data from the European Centre for Medium-Range  
109 Weather Forecasts (ECMWF) (Simmons et al., 1989). The U.S. National Oceanic and Atmospheric  
110 Administration's (NOAA) National Weather Service's National Centers for Environmental Prediction  
111 (NCEP) (Environmental Modeling Center, 2003) produces operational forecasts and a series of computer  
112 analyses. NCEP's Global Forecast System (GFS) produces pressure level data that can be used in  
113 FLEXPART (NCEP tag in Table 2). It also produces the GDAS (Global Data Assimilation System)  
114 reanalysis data which can be used in HYSPLIT (Kanamitsu et al., 1991). Five submissions used NCEP  
115 data and three submissions used GDAS data. Two submissions used the Weather Research and

116 Forecasting (WRF) model (Done et al., 2004; Michalakes et al., 2001; Skamarock et al., 2008). One  
117 participant used the global model ARPEGE (Action de Recherche Petite Echelle Grande Echelle) from  
118 the French meteorological office (Météo-France) (Déqué et al., 1994; Déqué and Piedelievre, 1995). One  
119 participant used the global meteorological analyses provided by the Canadian Meteorological Centre  
120 (CMC). CMC runs operationally a complete integrated suite of numerical weather prediction (NWP)  
121 models under an infrastructure called the Global Environmental Multiscale (GEM) system (Côté et al.,  
122 1998). The GEM system executed in a global configuration is called the GDPS: Global Deterministic  
123 Prediction System (Buehner et al., 2015; Buehner et al., 2013; Charron et al., 2012). The GDPS includes  
124 a 4D vibrational data assimilation system and is run twice a day (00 and 12 UTC) with a horizontal grid  
125 mesh defined at ~25 km (0.23° horizontal resolution). This global meteorological analyses database is  
126 used to drive MLDP0.

127 The spatial resolution of the meteorological grids in Table 2 is typically expressed in units of degrees. A  
128 1° grid for meteorological data in this region of the world has a north-south spacing of approximately 111  
129 km and an east-west spacing of 78 km. Similarly, a 0.5° grid has a spacing of 55 and 39 km, and a 0.2°  
130 grid has a spacing of about 22 and 16 km.

### 131 **3. Comparison measures**

132 The purpose of this challenge was to ascertain the level of agreement one can achieve between simulated  
133 concentrations and IMS measurements using only the stack data and an atmospheric transport model, as  
134 might be expected for situations in which there was a detection of radionuclides at an IMS station and very  
135 little other information. Concentration estimates from this modeling exercise are expected to be quite  
136 variable (Draxler et al., 2015), thus it is useful to explore the general characteristics of the models with  
137 the closest agreement with the sampled data. Researchers have proposed a number of different  
138 performance measures for comparing the outputs of atmospheric transport models. For purposes of this



139 analysis, five statistical measures described by other researchers (Chang and Hanna, 2004; Draxler, 2006)  
140 are used.

141 A brief introduction of each statistical measure is provided here. Additional information is given in the  
142 Appendix. The fractional bias (FB) is a measure of the bias between measured and predicted values. The  
143 correlation coefficient R is used to represent the linear relationship between measured and predicted  
144 values. The fraction of predicted values within a factor of five of the measured value (F5) is also used.  
145 The Kolmogorov–Smirnov (KS) statistic quantifies the differences between the distribution of unpaired  
146 measured and predicted values. The normalized mean square error (NMSE) is a measure of the difference  
147 between paired measured and predicted values.

148 The five statistical model comparison measures implicitly assume that all of the  $^{133}\text{Xe}$  measured at the  
149 IMS sampling station in originated from the IRE facility. Although IRE is the largest emitter of  $^{133}\text{Xe}$  in  
150 the region, it is not the only one. Nuclear power plants emit low levels of  $^{133}\text{Xe}$  (Kalinowski and Tuma,  
151 2009; Saey, 2009) and a number of nuclear power plants in Europe were in operation during this time  
152 period. Another medical isotope production facility in the Netherlands (Tyco Healthcare) releases about  
153 0.1% of the amount of  $^{133}\text{Xe}$  (Saey, 2009) as released from IRE on an annual basis. The medical isotope  
154 production facility in Chalk River, Canada, annually releases from three to four times as much  $^{133}\text{Xe}$   
155 (Saey, 2009) as IRE and under suitable meteorological conditions, may produce a measurable  
156 contribution to the  $^{133}\text{Xe}$  levels across Europe. In spite of these other sources, this is a realistic test case  
157 when data are only available from a single facility. In other words, for real world scenarios, we are testing  
158 the hypothesis that a single larger emitter may dominate the concentrations observed at an IMS facility.

159 Based on approaches suggested by other researchers (Chang and Hanna, 2004; Draxler, 2006), we  
160 combine four of the statistics into a single model ranking parameter as follows:

161 
$$Rank = R^2 + \left(1 - \frac{|FB|}{2}\right) + F5 + (1 - KS)$$

162 The model rank ranges from 0 (a model with no predictive ability) to 4 (a perfect model).

163 It is desirable to have contributors to an overall rank that measure different aspects of disparity. For  
164 example, a data set could have an  $R^2$  value of 1.0 but have a large magnitude of FB. There is some  
165 concern that FB and F5 measure similar aspects of disparity. However, for this data set, other than the  
166 four submissions with the lowest F5, the values for F5 and FB do not seem to be correlated.

#### 167 **4. Release and detection data**

168 Participants in the modeling challenge received  $^{133}\text{Xe}$  stack emission data from the Institut des  
169 Radioéléments (IRE) radiopharmaceutical plant in Fleurus, Belgium. Releases from IRE have a  
170 measurable influence on  $^{133}\text{Xe}$  concentrations collected at DEX33 (Saey et al., 2010a) which is located  
171 376 km from the IRE stack. The emission data covered the period 10 Nov 2013 through 8 Dec 2013. The  
172 measured concentration values for the stack data are based only on the 81 keV decay energy and have an  
173 uncertainty (one sigma) of approximately 10% of the measured values. The stack air flow rate was  $8 \times 10^4$   
174  $\text{m}^3 \text{h}^{-1}$ , without any uncertainty estimate. The concentrations of  $^{133}\text{Xe}$  in the exhaust stack air were  
175 provided for 2784 contiguous 15-min release periods. The amount released (concentration multiplied by  
176 the air flow rate) in each 15 minute period is shown in Fig. 1. Release quantities may vary by as much as  
177 two orders of magnitude for different 15-min duration periods in the same day.

178 The German national authority Bundesamt für Strahlenschutz (BfS) provided the  $^{133}\text{Xe}$  activity  
179 concentration data collected at the IMS noble gas sampler at Radionuclide Station RN33 (DEX33) at  
180 mount Schauinsland, Germany for the challenge. This sampling station is located at 1205 m above sea  
181 level on a mountain in the Black Forest. Surrounding low-level terrain ranges in elevation from 200 to  
182 600 m. The SPALAX<sup>TM</sup> system (Fontaine et al., 2004) at this station uses a sample collection period of  
183 24 hours. The time tag for each sample is the beginning of the sample collection period and the reported  
184 concentration is an average value decay-corrected to the beginning of the sample collection period. The

185 measured data at DEX33 and their uncertainties (one sigma) are shown in Fig. 2. The uncertainties range  
186 from 2.3% of the largest measured value to approximately 40% of the smallest values.

## 187 **5. Model comparison results**

188 Thirteen participants submitted 26 solutions containing modeled concentrations of  $^{133}\text{Xe}$  at the sampler  
189 (DEX33) in Germany on the time periods used by the sampler. A plot of modeled concentrations for all  
190 26 submissions and the concentrations at the sampler (black dots connected by a dotted line) is provided  
191 in Fig. 3. One submission had two predicted concentration values larger than  $100 \text{ mBq m}^{-3}$ , but the upper  
192 limit on this plot partially obscures that fact. Some of the values were zero, thus they cannot be  
193 represented on a log plot and the lines for adjacent nonzero values give the appearance of discontinuous  
194 data. However, the data were discrete values for each day and the lines on this plot are provided to aid in  
195 tracing of the time sequence of individual submissions.

196 The measured concentrations show five peaks separated in time and most modeled concentrations also  
197 show five peaks separated in time. There are three time periods (Nov. 17-19, Nov. 26-27 and Dec. 8-9)  
198 where most or all of the modeled concentrations are smaller than the measured concentrations. Data  
199 collected at DEX33 when IRE was not operating (Saey et al., 2010a) show that approximately 90% of the  
200 historical samples have concentrations above  $0.1 \text{ mBq m}^{-3}$ . Thus, it is reasonable to expect detectable  
201 background concentrations of  $^{133}\text{Xe}$  at this sampler from other sources even when the wind is blowing  
202 releases from IRE in a different direction.

203 Although the measured concentrations are influenced by releases from IRE, the highest concentrations in  
204 the plume often bypassed the sampling station during the time period shown in Fig. 3. The sample  
205 collection period of the first sample from DEX33 used in this study starts only 6 h after the first IRE  
206 release data, but it is 15 h before the first large release. Earlier simulations suggest that releases from IRE  
207 in the previous 3 d move to the northeast and almost all of the plume bypasses the sampler. An example

208 modeled  $^{133}\text{Xe}$  plume using the HYSPLIT computer code and GDAS data (3 h temporal resolution, 1°  
209 spatial resolution) corresponding to the time of the sample with collection start at 0600 UTC on  
210 November 14 is shown in Fig. 4. The plume is truncated on the south in Fig. 4 to minimize the output file  
211 size. This particular model run slightly underestimates the sampler concentration for this time period but  
212 it still illustrates the sharp gradients on the edges of the main body of the plume. As a consequence,  
213 relatively small discrepancies in the direction of movement between the modeled plume and the real  
214 plumes can lead to large concentration discrepancies at sampling locations.

### 215 **5.1 Statistical performance measures**

216 The values of the individual statistics and the ranking parameter are provided in Table 3 for every  
217 submitted solution. The entries in the table are sorted by descending rank. The best values for the  
218 individual performance measure are highlighted in bold text. The mean square error (MSE) between the  
219 modeled and predicted values is also provided because it is used in the ensemble calculation in the next  
220 section.

221 The only difference between Mau1 and Mau3 is that Mau3 used  $4 \times 10^7$  particles while Mau1 used  $3 \times 10^6$   
222 particles. The accuracy of predictions improved significantly using more particles. The submission with  
223 the largest rank (Sch) used background source estimates (average releases from other medical isotope  
224 production facilities and nuclear power plants) in addition to the releases from IRE in the calculation.  
225 This submission illustrates the effect additional sources can have on the KS statistic, because it is highly  
226 influenced by the additional sources (fewer predicted concentrations are near zero). The F5 statistic is  
227 influenced by the additional sources to a lesser extent.

### 228 **5.2 Ensemble performance measures**

229 Rather than comparing the results of individual models, one can attempt to combine them in an optimal  
230 way to provide a better prediction. A number of researchers (Kolczynski et al., 2009; Solazzo and

231 Galmarini, 2015) have started using ensembles of the individual models in an effort to produce better  
232 modeled concentrations. One of the justifications for using ensembles is to overcome the high sensitivity  
233 to the direction of plume movement illustrated in Fig. 4.

234 An ensemble reduction technique based on minimizing the mean square error between the measured and  
235 predicted concentrations is now available (Stein et al., 2015) in the HYSPLIT suite of codes. Using this  
236 approach, we calculate the average of all possible model combinations composed by increasing the  
237 number of ensemble members from 1 to 25 and estimate their MSE. The combination with the minimum  
238 MSE is then selected. In other words, we combine the 25 model outputs in 300 pairs, 2300 trios, etc., and  
239 determine which combination provides the minimum MSE. Fig. 5 shows the minimum MSE obtained as a  
240 function of the number of submissions in the reduced ensembles. The curve has a minimum at two  
241 ensemble members. In addition, the best ensembles with two, three or four members all have lower MSE  
242 than the single best model. This means that including more than about four members in the ensemble will  
243 produce a less accurate result.

244 The MSE of an average of several submissions used to select the ensemble members is different than the  
245 performance measures shown in Table 3. The ensemble of four members yields an average value that has  
246  $KS=0.42$ ,  $R=0.98$ ,  $FB=-0.25$ ,  $F5=0.61$ ,  $Rank = 3.03$ ,  $NMSE=0.81$  and  $MSE=2.74$ . As a comparison, the  
247 ensemble with only two members (Hof4 and Mau3) has  $KS=0.42$ ,  $R=0.97$ ,  $FB=0.01$ ,  $F5=0.58$ ,  
248  $Rank=3.10$ ,  $NMSE=0.31$  and  $MSE=1.34$ . The rank for the two member ensemble is better than the rank  
249 of the best submission and the rank of the four member ensemble is about equal to the rank of the best  
250 submission. The correlation ( $R$ ) of the four member ensemble is higher than for the single best  
251 submission, but the fractional bias ( $FB$ ) is worse. The modeled  $^{133}\text{Xe}$  concentrations for the ensemble  
252 members and the ensemble average for the minimum MSE ensemble of four members is provided in Fig.  
253 6. Two of the ensemble members used releases varying every 15 min while the other two used sources  
254 varying every 3 hr. These four models use four different meteorological data sets and two different

255 computer codes, implying independence between the four ensemble members. Independence among  
256 ensemble members is a necessary but not sufficient condition for building accurate ensembles  
257 (Kioutsioukis and Galmarini, 2014).

258 This study, and historical sampling data from DEX33 when IRE was not operating (Saey et al., 2010a),  
259 suggests that the largest sample values are heavily dominated by releases from IRE. A comparison of  
260 measured and predicted concentrations are provided in Table 4 for the five largest sampled values for the  
261 submissions that scored the highest on individual statistical performance measures. The ensemble with  
262 four members is also included for comparison. The percentage values are the relative difference of the  
263 predicted and measured concentrations, and a negative value means the predicted value is smaller than the  
264 measured value. The Hof2 submission had a high correlation (0.97) between the sampled and measured  
265 concentrations, but also a large fractional bias. Some of the submissions predicted the largest  
266 concentrations to within 15%. The submission (Sau) did not have the best score on any specific statistical  
267 measure, but it was one of the four members of the minimum MSE ensemble and it has the smallest  
268 maximum relative error on the five largest measured concentrations.

### 269 **5.3 Comparisons using grouped submissions**

270 Ranks were calculated for several different combinations of the suite of submissions in addition to the  
271 minimum MSE ensemble approach. The ranks provided in Fig. 7 are based on the seventeen submissions  
272 identified in Table 2. Except for the single submission with the highest rank, the ranks were calculated  
273 using the average of each member of the group. The average of all the submissions has a lower rank than  
274 the average from the ensemble with four members. The rank for the group of HYSPLIT models is lower  
275 than the ranks for the FLEXPART and other models. Most of the FLEXPART models used ECMWF  
276 meteorological data while most of the HYSPLIT models used GDAS data. Thus, it is not surprising that  
277 the lower ranks using the HYSPLIT model correspond to the lower ranks for GDAS data as compared to  
278 other data sets. Although the governing equations generally are time reversible, the implementations yield

279 slightly different concentration estimates depending on the time direction. The average of the forwards  
280 time runs had a slightly higher rank than the average of the backwards runs. The average of model runs  
281 using meteorological data with finer spatial resolution than 0.5° had higher rank than those using 0.5°  
282 resolution data. The average of model runs using 1.0° resolution meteorological data had a rank about  
283 equal to the average of finer resolution model runs, however, the normalized MSE for the 1.0° spatial  
284 resolution runs was 5.09 while that of the finer spatial resolution runs was 2.89. Those models that  
285 incorporated the source term on a 15-min timing basis had higher ranks than models using sources using  
286 longer source term aggregation periods.

#### 287 **5.4 Additional sources**

288 The modeling exercise was formulated to consider the hypothesis that a single larger emitter may  
289 dominate the concentrations observed at an IMS facility. However, one submission (Sch) included annual  
290 average emission rates for nuclear power plants and other medical isotope production facilities as an  
291 additional source term. The Sch results are compared to the four member ensemble average in Fig. 8. This  
292 submission suggests that the other releases are also influencing the sampler, and this result is consistent  
293 with historical data (Saey et al., 2010a). The transport runs done for submission Hof4 yielded effective  
294 atmospheric dilution factors that indicate releases from the medical isotope production facility in Chalk  
295 River, Canada, could potentially influence 18 of the 30 DEX33 samples. No Chalk River source was  
296 introduced in the Hof4 submittal even though releases from the facility seem to have influenced some of  
297 the measured data at DEX33.

#### 298 **6. Discussion**

299 The ranking and ensemble analysis in this paper suggests that combining multiple models may provide  
300 more accurate predicted concentrations than almost any single model. One ensemble selection technique  
301 was used in this paper. Further research is needed to identify optimal methods for selecting ensemble

302 members, and those methods may depend on the nature of the transport problem. Although this exercise  
303 only addressed release and transport of a nondepositing noble gas, other radionuclides of interest to the  
304 treaty monitoring community (such as  $^{137}\text{Cs}$  and  $^{131}\text{I}$ ) deposit on the ground during transport, and models  
305 that work best for predicting air concentrations may not fare as well when predicting deposition on the  
306 ground (Draxler et al., 2015).

307 Participants in this challenge predicted measured concentrations at a sampling station using only releases  
308 from one medical isotope production facility. Some of the models predicted the highest measured  
309 concentrations quite well (high rank or small MSE); however none predicted the small measured  
310 (background) concentrations very well. The one submission that included average release estimates from  
311 other nuclear facilities matched the small concentrations much better. If expected releases from future  
312 nuclear tests are small, such as releases from the 2013 test by the Democratic People's Republic of Korea  
313 (Ringbom et al., 2014), then modeling of sources from nuclear facilities with smaller releases than  
314 medical isotope production facilities may also be important.

315 The grouped model comparisons shown in Fig. 7 categorize prediction performance relative to several of  
316 the choices available to modelers. For this exercise, the ranks for submissions using FLEXPART were  
317 higher than the ranks for submissions using FLEXPART. However, most HYSPLIT runs used GDAS  
318 data while FLEXPART used other meteorological data. Interpretation of the results must recognize that  
319 most of the categories are confounded with each other. For example, all of the HYSPLIT model runs in  
320 comparisons in Fig. 7 did runs that were forwards in time. In addition, the sampler at DEX33 used a  
321 collection interval of 24 h, and 24 h may be long enough to average out some of the differences in the  
322 time resolution of the source term. The release data from IRE were provided with a time resolution of 15  
323 min. Two of the models in the four member minimum MSE ensemble used 15 min release data, but the  
324 other two aggregated releases to a 3 h basis. The average predicted concentrations for the models that



325 incorporated the source term on a 15-min timing basis had a higher rank than models using longer release  
326 periods. However, models using 3 h source averaging had a higher rank than those using 1 h averaging.

327 Other operational radioxenon samplers in the IMS use a shorter sample collection interval of 12 h  
328 (Prelovskii et al., 2007; Ringbom et al., 2003) and new generation radioxenon samplers under  
329 development (Hayes et al., 2013; Le Petit et al., 2015) can use collection periods of 6 or 8 h. These  
330 shorter collection periods may show more sensitivity to the time resolution of a highly time-variable  
331 source term than the current sampler.

332 Finally, the results of this single exercise indicate that the use of stack monitoring data to determine  
333 radionuclide concentrations at a distance of nearly 400 km can yield predicted large concentrations within  
334  $\pm 40\%$  of the measured concentrations if an ensemble is used. Individual models have a larger spread than  
335 the ensemble results. The uncertainties in the stack data do not appear to dominate the uncertainties in the  
336 modeled results. However, the uncertainty in the air flow rate in the stack is not known, so the  
337 uncertainty in the release values may be significantly larger than the 10% uncertainty in the isotope  
338 concentration data in the stack. More work will be needed to determine the achievable accuracy in other  
339 conditions, such as for larger source-receptor distances. We anticipate more exercises of this nature could  
340 help to define methods to understand the effect of emissions from fission-based medical isotope  
341 production on IMS sampling data.

## 342 **Acknowledgments**

343 Participants in the atmospheric transport modeling challenge received  $^{133}\text{Xe}$  emission data from the  
344 Institut des Radioéléments (IRE) radiopharmaceutical plant in Fleurus, Belgium. IRE granted permission  
345 to use the data for the challenge.

346 The German national authority *Bundesamt für Strahlenschutz* (BfS) granted permission to use the  $^{133}\text{Xe}$   
347 concentration data collected at the IMS noble gas sampler (DEX33) in Schauinsland, Germany for the

348 challenge. Clemens Schlosser and Verena Heidmann of Bfs manually analyzed the spectra to obtain the  
349  $^{133}\text{Xe}$  concentration data and associated error bars.

350 Some of the authors wish to acknowledge the funding support of the U.S. Department of State and the  
351 U.S. Defense Threat Reduction Agency.

## 352 **Appendix**

353 In the following descriptions, let  $P$  denote predicted concentrations,  $M$  denote measured concentrations,  
354 an overbar denote an average over the data set, and  $i$  denote an index of the  $N$  sample values. The  
355 fractional bias (FB) is measure of the bias between measured and predicted values. The FB is normalized  
356 to the range -2 to 2 and positive values indicate predictions are larger than measured values. Small  
357 concentrations attributable to releases from facilities other than IRE have a small effect on this  
358 performance measure. The fractional bias is defined as:

$$359 \quad FB = 2 \frac{(\bar{P} - \bar{M})}{(\bar{P} + \bar{M})} \quad (1)$$

360 The correlation coefficient  $R$  is used to represent the linear relationship between measured and predicted  
361 values where the summation is taken over all samples. Possible values for  $R$  range from -1 to 1. The  
362 correlation coefficient is calculated from:

$$363 \quad R = \frac{\sum(M_i - \bar{M})(P_i - \bar{P})}{\sqrt{\sum(M_i - \bar{M})^2 (P_i - \bar{P})^2}} \quad (2)$$

364 The fraction of predicted values within a certain factor of the measured value is often used in model  
365 comparisons. This statistic can be heavily influenced if some modeled values are near zero while nuisance  
366 sources cause the measured values to be at or just above a detection limit. We define the factor of five  
367 (F5) statistic as the fraction of sample values that satisfy:

$$368 \quad 0.2 \leq \frac{P_i}{M_i} \leq 5.0 \quad (3)$$

369 The Kolmogorov–Smirnov (KS) statistic (Stephens, 1970) quantifies the differences between the  
370 distribution of unpaired measured and predicted values. The values are considered as samples from two  
371 different statistical distributions and KS is defined as the maximum difference between two cumulative  
372 distributions when  $M_k=P_k$ , where

$$373 \quad KS = \text{Max}|D(M_k) - D(P_k)|. \quad (4)$$

374 In this case, D is the cumulative distribution of the measured and predicted concentrations over the range  
375 of k values such that D is the probability that the concentration will not exceed  $M_k$  or  $P_k$ . It measures the  
376 ability of the model to reproduce the measured concentration distribution regardless of when or where it  
377 occurred. The maximum difference between any two distributions cannot be more than 100%. This  
378 statistic can be heavily influenced if some modeled values are near zero while nuisance sources cause the  
379 measured values to be at or just above a detection limit.

380 The normalized mean square error (NMSE) is a measure of the difference between paired measured and  
381 predicted values. The normalized mean square error is calculated from:

$$382 \quad NMSE = \frac{MSE}{\overline{M \cdot P}} \quad (5)$$

383 where MSE is the mean square error defined as:

$$384 \quad MSE = \frac{1}{N} \sum (M_i - P_i)^2 \quad (6)$$

## 385 **References**

- 386 Becker, A., Wotawa, G., Ringbom, A., Saey, P.R.J., 2010. Backtracking of Noble Gas Measurements  
387 Taken in the Aftermath of the Announced October 2006 Event in North Korea by Means of PTS  
388 Methods in Nuclear Source Estimation and Reconstruction. Pure Appl. Geophys. 167(4), 581-599.  
389 doi:10.1007/s00024-009-0025-0
- 390 Bowyer, T.W., Kephart, R., Eslinger, P.W., Friese, J.I., Miley, H.S., Saey, P.R.J., 2013. Maximum  
391 reasonable radionuclide releases from medical isotope production facilities and their effect on  
392 monitoring nuclear explosions. J. Environ. Radioact. 115(1), 192-200.  
393 doi:10.1016/j.jenvrad.2012.07.018
- 394 Buehner, M., McTaggart-Cowan, R., Beaulne, A., Charette, C., Garand, L., Heilliette, S., Lapalme, E.,  
395 Laroche, S., Macpherson, S.R., Morneau, J., Zadra, A., 2015. Implementation of Deterministic  
396 Weather Forecasting Systems Based on Ensemble–Variational Data Assimilation at Environment

397 Canada. Part I: The Global System. *Monthly Weather Review* 143(7), 2532-2559.  
398 doi:10.1175/MWR-D-14-00354.1

399 Buehner, M., Morneau, J., Charette, C., 2013. Four-dimensional ensemble-variational data assimilation  
400 for global deterministic weather prediction. *Nonlin. Processes Geophys.* 20(5), 669-682.  
401 doi:10.5194/npg-20-669-2013

402 Chang, J.C., Hanna, S.R., 2004. Air quality model performance evaluation. *Meteorology and  
403 Atmospheric Physics* 87(1-3), 167-196. doi:10.1007/s00703-003-0070-7

404 Charron, M., Polavarapu, S., Buehner, M., Vaillancourt, P.A., Charette, C., Roch, M., Morneau, J.,  
405 Garand, L., Aparicio, J.M., MacPherson, S., Pellerin, S., St-James, J., Heilliette, S., 2012. The  
406 Stratospheric Extension of the Canadian Global Deterministic Medium-Range Weather Forecasting  
407 System and Its Impact on Tropospheric Forecasts. *Monthly Weather Review* 140(6), 1924-1944.  
408 doi:10.1175/MWR-D-11-00097.1

409 Comprehensive Nuclear-Test-Ban Treaty. 1996. *Text of the Comprehensive Nuclear-Test-Ban Treaty*.  
410 United Nations Office for Disarmament Affairs (UNODA), Status of Multilateral Arms Regulation  
411 and Disarmament Agreements, CTBT. Accessed on September 20, 2012, at  
412 <http://www.ctbto.org/the-treaty/treaty-text/>

413 Côté, J., Desmarais, J.-G., Gravel, S., Méthot, A., Patoine, A., Roch, M., Staniforth, A., 1998. The  
414 Operational CMC-MRB Global Environmental Multiscale (GEM) Model. Part II: Results. *Monthly  
415 Weather Review* 126(6), 1397-1418. doi:10.1175/1520-0493(1998)126<1397:TOCMGE>2.0.CO;2

416 CTBTO. 2014. *Verification Regime*. Accessed on October 13, 2014, at  
417 [http://www.ctbto.org/verification-regime/monitoring-technologies-how-they-work/radionuclide-  
418 monitoring/page-5/](http://www.ctbto.org/verification-regime/monitoring-technologies-how-they-work/radionuclide-monitoring/page-5/)

419 D'Amours, R., Malo, A., Flesch, T., Wilson, J., Gauthier, J.-P., Servranckx, R., 2015. The Canadian  
420 Meteorological Centre's Atmospheric Transport and Dispersion Modelling Suite. *Atmosphere-  
421 Ocean* 53(2), 176-199. doi:10.1080/07055900.2014.1000260

422 D'Amours, R., Malo, A., Servranckx, R., Bensimon, D., Trudel, S., Gauthier-Bilodeau, J.P., 2010.  
423 Application of the atmospheric lagrangian particle dispersion model MLDP0 to the 2008 eruptions  
424 of Okmok and Kasatochi volcanoes. *Journal of Geophysical Research-Atmospheres* 115(D00L11),  
425 1-11. doi:10.1029/2009JD013602

426 Déqué, M., Dreveton, C., Braun, A., Cariolle, D., 1994. The ARPEGE/IFS atmosphere model: a  
427 contribution to the French community climate modelling. *Climate Dynamics* 10(4-5), 249-266.  
428 doi:10.1007/BF00208992

429 Déqué, M., Piedelievre, J.P., 1995. High resolution climate simulation over Europe. *Climate Dynamics*  
430 11(6), 321-339. doi:10.1007/BF00215735

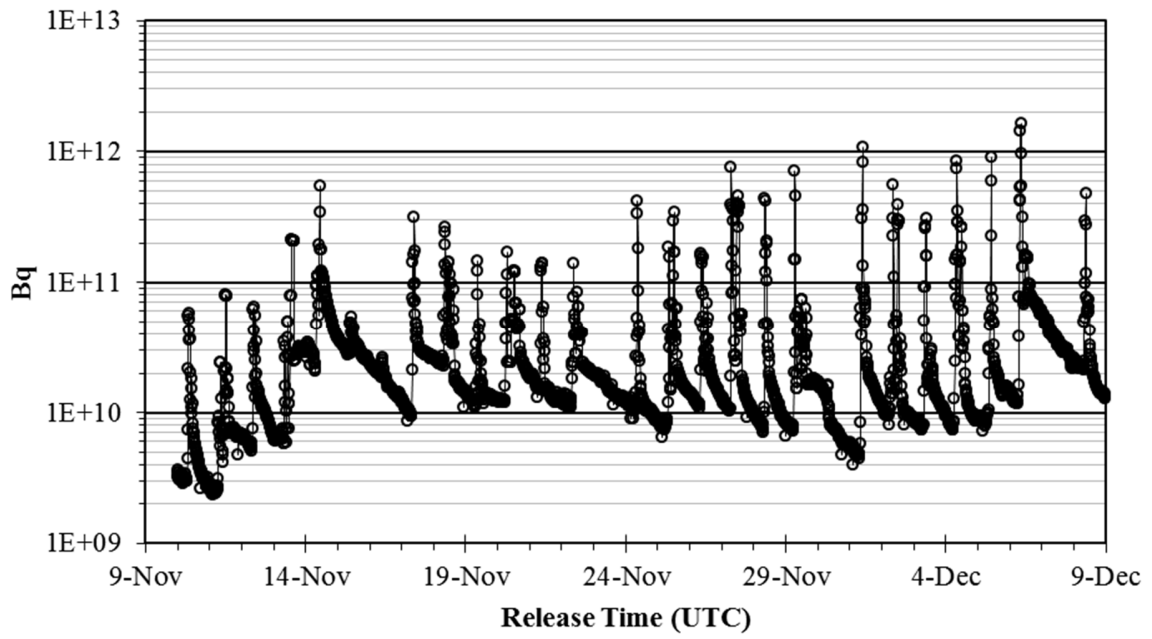
431 Done, J., Davis, C.A., Weisman, M., 2004. The next generation of NWP: explicit forecasts of convection  
432 using the weather research and forecasting (WRF) model. *Atmospheric Science Letters* 5(6), 110-  
433 117. doi:10.1002/asl.72

434 Draxler, R., Arnold, D., Chino, M., Galmarini, S., Hort, M., Jones, A., Leadbetter, S., Malo, A., Maurer,  
435 C., Rolph, G., Saito, K., Servranckx, R., Shimbori, T., Solazzo, E., Wotawa, G., 2015. World  
436 Meteorological Organization's model simulations of the radionuclide dispersion and deposition from  
437 the Fukushima Daiichi nuclear power plant accident. *J. Environ. Radioact.* 139(0), 172-184.  
438 doi:10.1016/j.jenvrad.2013.09.014

- 439 Draxler, R.R., 2006. The Use of Global and Mesoscale Meteorological Model Data to Predict the  
440 Transport and Dispersion of Tracer Plumes over Washington, D.C. *Weather Forecast* 21(3), 383-  
441 394. doi:10.1175/WAF926.1
- 442 Draxler, R.R., Hess, G.D., 1998. An overview of the HYSPLIT\_4 modeling system of trajectories,  
443 dispersion, and deposition. *Aust. Meteorol. Mag.* 47, 295-308
- 444 Draxler, R.R., Hess, G.D., 2010. Description of the HYSPLIT\_4 Modeling System, ARL-224, Air  
445 Resources Laboratory, National Oceanic and Atmospheric Administration (NOAA), Silver Springs,  
446 Maryland
- 447 Eslinger, P.W., Friese, J.I., Lowrey, J.D., McIntyre, J.I., Miley, H.S., Schrom, B.T., 2014. Estimates of  
448 radionuclides released from Southern Hemisphere medical isotope production facilities using  
449 measured air concentrations and atmospheric transport modeling. *J. Environ. Radioact.* 135(2014),  
450 94-99. doi:10.1016/j.jenvrad.2014.04.006
- 451 Ferber, G.J., Heffter, J.L., Draxler, R.R., Lagomarsino, R.J., Thomas, F.L., Deitz, R.N., Benkovitz, C.M.,  
452 1986. Cross-Appalachian Tracer Experiment (CAPTEX -83) final report., NOAA Tech. Memo. ERL  
453 ARL-142, NOAA/Air Resources Laboratory, Silver Spring, Maryland
- 454 Fontaine, J.P., Pointurier, F., Blanchard, X., Taffary, T., 2004. Atmospheric xenon radioactive isotope  
455 monitoring. *J. Environ. Radioact.* 72(1-2), 129-135. doi:10.1016/S0265-931X(03)00194-2
- 456 Gudiksen, P.H., Ferber, G.J., Fowler, M.M., Eberhard, W.L., Fosberg, M.A., Knuth, W.R., 1984. Field  
457 studies of transport and dispersion of atmospheric tracers in nocturnal drainage flows. *Atmospheric  
458 Environment* (1967) 18(4), 713-731. doi:http://dx.doi.org/10.1016/0004-6981(84)90257-9
- 459 Hayes, J.C., Ely, J.H., Haas, D.A., Harper, W.W., Heimbigner, T.R., Hubbard, C.W., Humble, P.H.,  
460 Madison, J.C., Morris, S.J., Panisko, M.E., Ripplinger, M.D., Stewart, T.L., 2013. Requirements for  
461 Xenon International, PNNL-22227 Rev.1, Pacific Northwest National Laboratory, Richland,  
462 Washington. doi:10.2172/1122330
- 463 Hoffman, I., Ungar, K., Bean, M., Yi, J., Servranckx, R., Zaganescu, C., Ek, N., Blanchard, X., Le Petit,  
464 G., Brachet, G., Achim, P., Taffary, T., 2009. Changes in Radionuclide Observations in Canada and  
465 Europe during Medical Isotope Production Facility Shut Down in 2008. *Journal of Radioanalytical  
466 and Nuclear Chemistry* 282, 767-772. doi:10.1007/s10967-009-0235-z
- 467 Kalinowski, M.B., Becker, A., Saey, P.R.J., Tuma, M.P., Wotawa, G., 2008. The Complexity of CTBT  
468 Verification. Taking Noble Gas Monitoring as an Example. *Complexity* 14(1), 89-99.  
469 doi:10.1002/cplx.20228
- 470 Kalinowski, M.B., Tuma, M.P., 2009. Global radionuclide emission inventory based on nuclear power  
471 reactor reports. *J. Environ. Radioact.* 100(1), 58-70. doi:10.1016/j.jenvrad.2008.10.015
- 472 Kanamitsu, M., Alpert, J.C., Campana, K.A., Caplan, P.M., Deaven, D.G., Iredell, M., Katz, B., Pan,  
473 H.L., Sela, J., White, G.H., 1991. Recent Changes Implemented into the Global Forecast System at  
474 NMC. *Weather Forecast* 6(3), 425-435. doi:10.1175/1520-0434(1991)006<0425:RCIITG>2.0.CO;2
- 475 Kioutsioukis, I., Galmarini, S., 2014. De praeceptis ferendis: good practice in multi-model ensembles.  
476 *Atmos. Chem. Phys.* 14(21), 11791-11815. doi:10.5194/acp-14-11791-2014
- 477 Kolczynski, W.C., Stauffer, D.R., Haupt, S.E., Deng, A., 2009. Ensemble Variance Calibration for  
478 Representing Meteorological Uncertainty for Atmospheric Transport and Dispersion Modeling.  
479 *Journal of Applied Meteorology and Climatology* 48(10), 2001-2021.  
480 doi:10.1175/2009JAMC2059.1

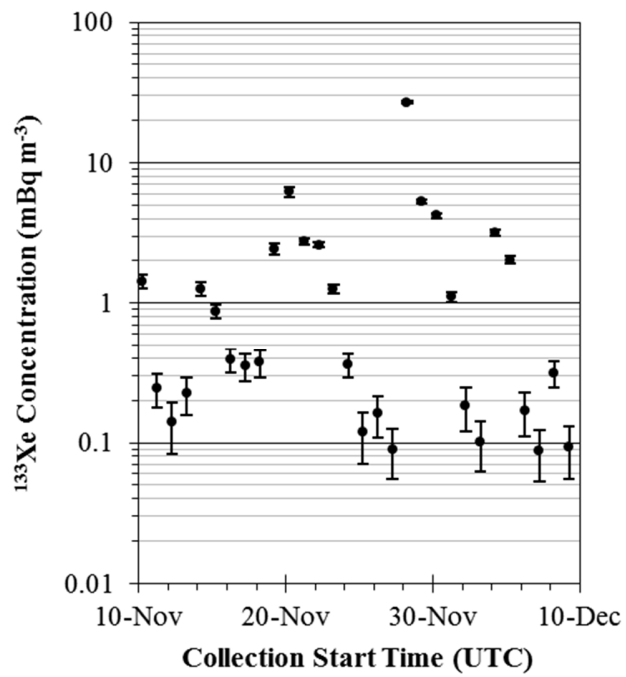
- 481 Le Petit, G., Cagniant, A., Gross, P., Douysset, G., Topin, S., Fontaine, J.P., Taffary, T., Moulin, C.,  
 482 2015. Spalax™ new generation: A sensitive and selective noble gas system for nuclear explosion  
 483 monitoring. *Appl. Radiat. Isot.* 103(0), 102-114. doi:10.1016/j.apradiso.2015.05.019
- 484 Michalakes, J., Chen, S., Dudhia, J., Hart, L., Klemp, J., Middlecoff, J., Skamarock, W., 2001.  
 485 Development of a next generation regional weather research and forecast model, in: Zwiefelhofer,  
 486 W., Kreitz, N. (Eds.), *Developments in Teracomputing: Proceedings of the Ninth ECMWF*  
 487 *Workshop on the Use of High Performance Computing in Meteorology.* World Scientific  
 488 Publishing, Singapore, pp. 269–276.
- 489 Peykov, P., Cameron, R., 2014. Medical Isotope Supply in the Future: Production Capacity and Demand  
 490 Forecast for the <sup>99</sup>Mo/<sup>99m</sup>Tc Market, 2015-2020, NEA/SEN/HLGMR(2014)2, Organisation for  
 491 Economic Co-Operation and Development, Nuclear Energy Agency, Issy-les-Moulineaux, France
- 492 Prelovskii, V.V., Kazarinov, N.M., Donets, A.Y., Popov, V.Y., Popov, I.Y., Skirda, N.V., 2007. The  
 493 ARIX-03F mobile semiautomatic facility for measuring low concentrations of radioactive xenon  
 494 isotopes in air and subsoil gas. *Instrum Exp Tech* 50(3), 393-397. doi:10.1134/S0020441207030165
- 495 Ringbom, A., Axelsson, A., Aldener, M., Auer, M., Bowyer, T.W., Fritioff, T., Hoffman, I., Khrestalev,  
 496 K., Nikkinen, M., Popov, V., Popov, Y., Ungar, K., Wotawa, G., 2014. Radioxenon detections in the  
 497 CTBT international monitoring system likely related to the announced nuclear test in North Korea  
 498 on February 12, 2013. *J. Environ. Radioact.* 128(0), 47-63. doi:10.1016/j.jenvrad.2013.10.027
- 499 Ringbom, A., Larson, T., Axelsson, A., Elmgren, K., Johansson, C., 2003. SAUNA—a system for  
 500 automatic sampling, processing, and analysis of radioactive xenon. *Nucl Instrum Meth A* 508(3),  
 501 542-553. doi:10.1016/s0168-9002(03)01657-7
- 502 Saey, P.R.J., 2009. The influence of radiopharmaceutical isotope production on the global radioxenon  
 503 background. *J. Environ. Radioact.* 100(5), 396-406. doi:10.1016/j.jenvrad.2009.01.004
- 504 Saey, P.R.J., Auer, M., Becker, A., Hoffmann, E., Nikkinen, M., Ringbom, A., Tinker, R., Schlosser, C.,  
 505 Sonck, M., 2010a. The influence on the radioxenon background during the temporary suspension of  
 506 operations of three major medical isotope production facilities in the Northern Hemisphere and  
 507 during the start-up of another facility in the Southern Hemisphere. *J. Environ. Radioact.* 101(9),  
 508 730-738. doi:10.1016/j.jenvrad.2010.04.016
- 509 Saey, P.R.J., Schlosser, C., Achim, P., Auer, M., Axelsson, A., Becker, A., Blanchard, X., Brachet, G.,  
 510 Cella, L., De Geer, L.-E., Kalinowski, M.B., Le Petit, G., Peterson, J., Popov, V., Popov, Y., et al.,  
 511 2010b. Environmental Radioxenon Levels in Europe: a Comprehensive Overview. *Journal of Pure*  
 512 *and Applied Geophysics* 167(4-5), 499-515. doi:10.1007/s00024-009-0034-z
- 513 Schöppner, M., Plastino, W., Hermanspahn, N., Hoffmann, E., Kalinowski, M., Orr, B., Tinker, R., 2013.  
 514 Atmospheric transport modelling of time resolved <sup>133</sup>Xe emissions from the isotope production  
 515 facility ANSTO, Australia. *J. Environ. Radioact.* 126(2013), 1-7. doi:10.1016/j.jenvrad.2013.07.003
- 516 Simmons, A.J., Burridge, D.M., Jarraud, M., Girard, C., Wergen, W., 1989. The ECMWF medium-range  
 517 prediction models development of the numerical formulations and the impact of increased  
 518 resolution. *Meteorology and Atmospheric Physics* 40(1-3), 28-60. doi:10.1007/BF01027467
- 519 Skamarock, W., Klemp, J.B., Dudhia, J., Gill, D.O., Barker, D., Duda, M.G., Huang, X., Wang, W., 2008.  
 520 A Description of the Advanced Research WRF Version 3, NCAR/TN-475+STR, National Center  
 521 for Atmospheric Research, Boulder, Colorado, USA. doi:10.5065/D68S4MVH
- 522 Solazzo, E., Galmarini, S., 2015. The Fukushima-137Cs deposition case study: properties of the multi-  
 523 model ensemble. *J. Environ. Radioact.* 139, 226-233. doi:10.1016/j.jenvrad.2014.02.017

- 524 Stein, A.F., Ngan, F., Draxler, R.R., Chai, T., 2015. Potential Use of Transport and Dispersion Model  
525 Ensembles for Forecasting Applications. *Weather Forecast* 30(3), 639-655. doi:10.1175/WAF-D-14-  
526 00153.1
- 527 Stephens, M.A., 1970. Use of the Kolmogorov-Smirnov, Cramér-Von Mises and related statistics without  
528 extensive tables. *Journal of the Royal Statistical Society. Series B (Methodological)*, 115-122
- 529 Stohl, A., Forster, C., Frank, A., Seibert, P., Wotawa, G., 2005. Technical note: The Lagrangian particle  
530 dispersion model FLEXPART version 6.2. *Atmos. Chem. Phys.* 5(9), 2461-2474. doi:10.5194/acp-  
531 5-2461-2005
- 532 Stohl, A., Hittenberger, M., Wotawa, G., 1998. Validation of the lagrangian particle dispersion model  
533 FLEXPART against large-scale tracer experiment data. *Atmos. Environ.* 32(24), 4245-4264.  
534 doi:10.1016/s1352-2310(98)00184-8
- 535 Tombette, M., Quentric, E., Quelo, D., Benoit, J.P., Mathieu, A., Korsakissok, I., Didier, D., 2014. C3X:  
536 a software platform for assessing the consequences of an accidental release of radioactivity into the  
537 atmosphere, Poster presented at Fourth European IRPA Congress, 23-27 June 2014, Geneva.
- 538 Wotawa, G., Becker, A., Kalinowski, M., Saey, P., Tuma, M., Zähringer, M., 2010. Computation and  
539 Analysis of the Global Distribution of the Radioxenon Isotope <sup>133</sup>Xe based on Emissions from  
540 Nuclear Power Plants and Radioisotope Production Facilities and its Relevance for the Verification  
541 of the Nuclear-Test-Ban Treaty. *Pure Appl. Geophys.* 167(4-5), 541-557. doi:10.1007/s00024-009-  
542 0033-0
- 543 Wotawa, G., De Geer, L.-E., Denier, P., Kalinowski, M., Toivonen, H., D'Amours, R., Desiato, F.,  
544 Issartel, J.-P., Langer, M., Seibert, P., Frank, A., Sloan, C., Yamazawa, H., 2003. Atmospheric  
545 transport modelling in support of CTBT verification—overview and basic concepts. *Atmos.*  
546 *Environ.* 37(18), 2529-2537. doi:10.1016/s1352-2310(03)00154-7
- 547 Zähringer, M., Becker, A., Nikkinen, M., Saey, P., Wotawa, G., 2009. CTBT radioxenon monitoring for  
548 verification: today's challenges. *J Radioanal Nucl Chem* 282(3), 737-742. doi:10.1007/s10967-009-  
549 0207-3
- 550

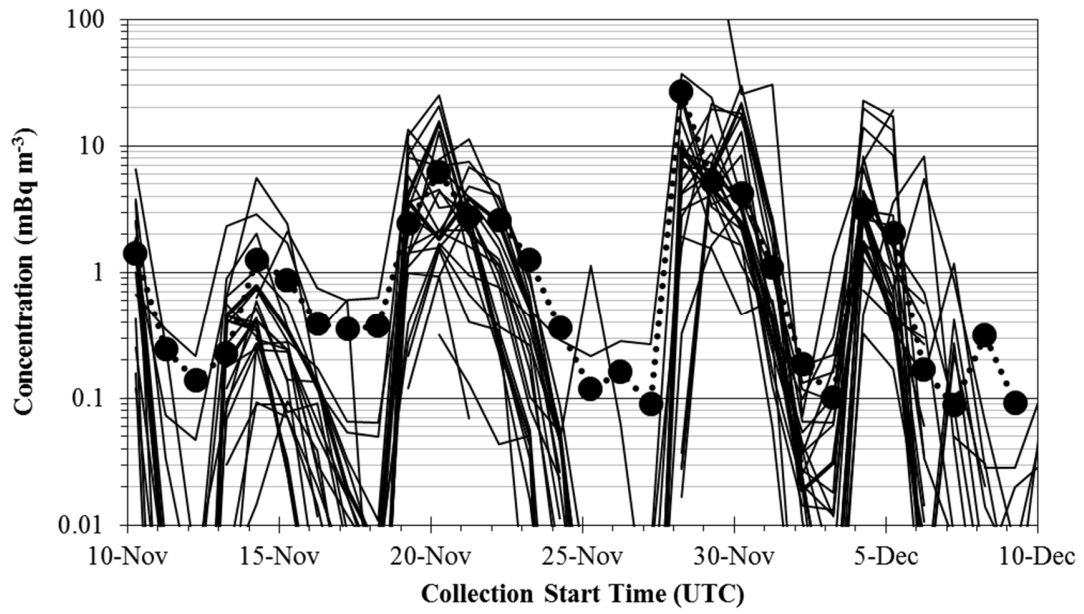


**Fig. 1.** Releases of  $^{133}\text{Xe}$  (Bq) in contiguous 15 minute intervals from the exhaust stack at the Institut des Radioéléments (IRE) radiopharmaceutical plant in Fleurus, Belgium.





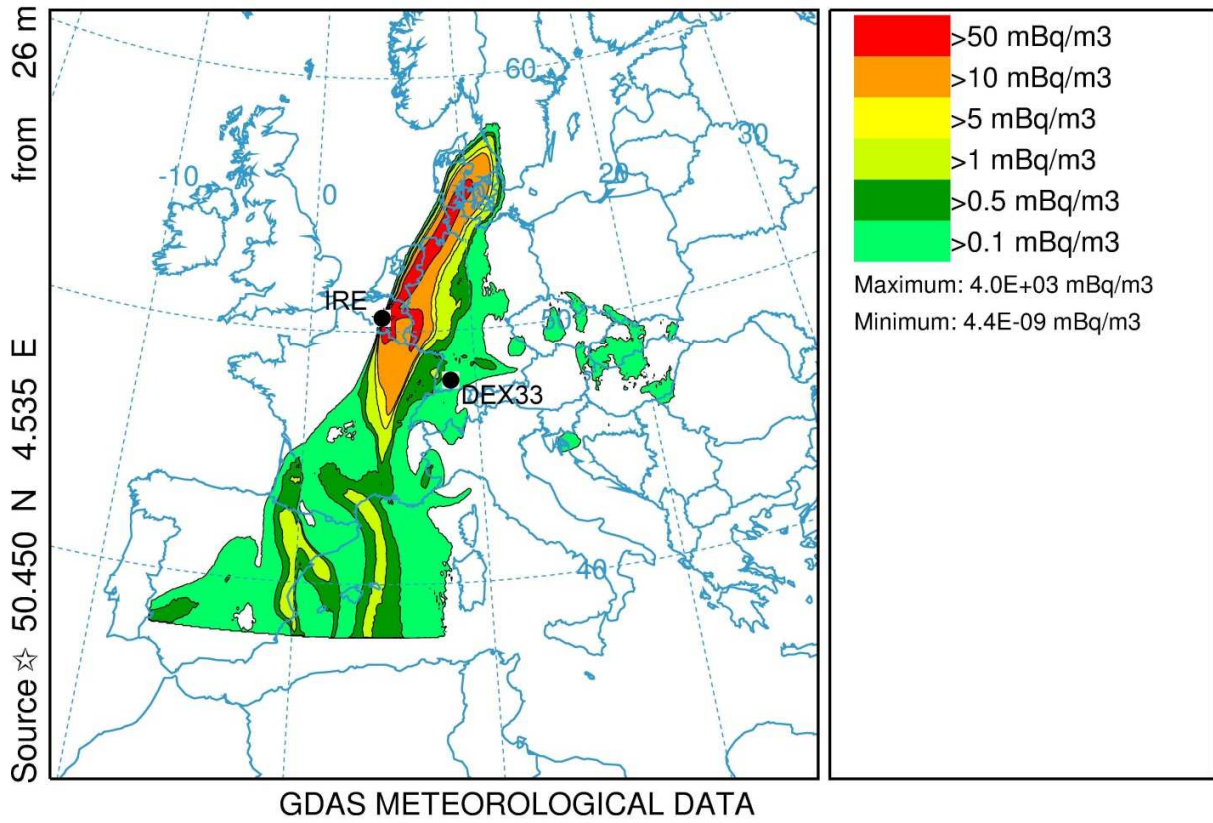
**Fig. 2.** Measured  $^{133}\text{Xe}$  activity concentrations at DEX33. The error bars represent one sigma uncertainties.



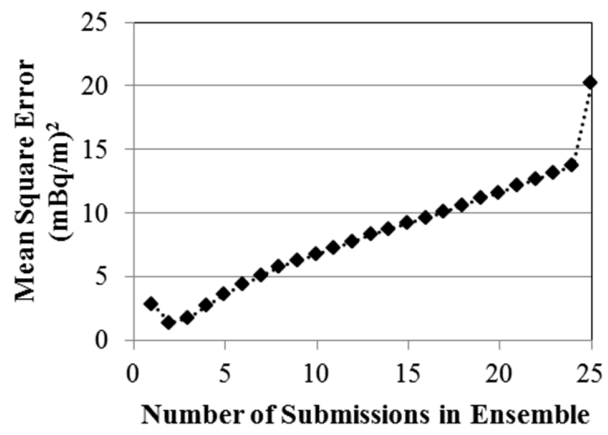
**Fig. 3.** Modeled <sup>133</sup>Xe concentrations for all submissions (solid lines) and measured concentrations at the sampler (large black dots connected by dotted lines).

### NOAA HYSPLIT MODEL

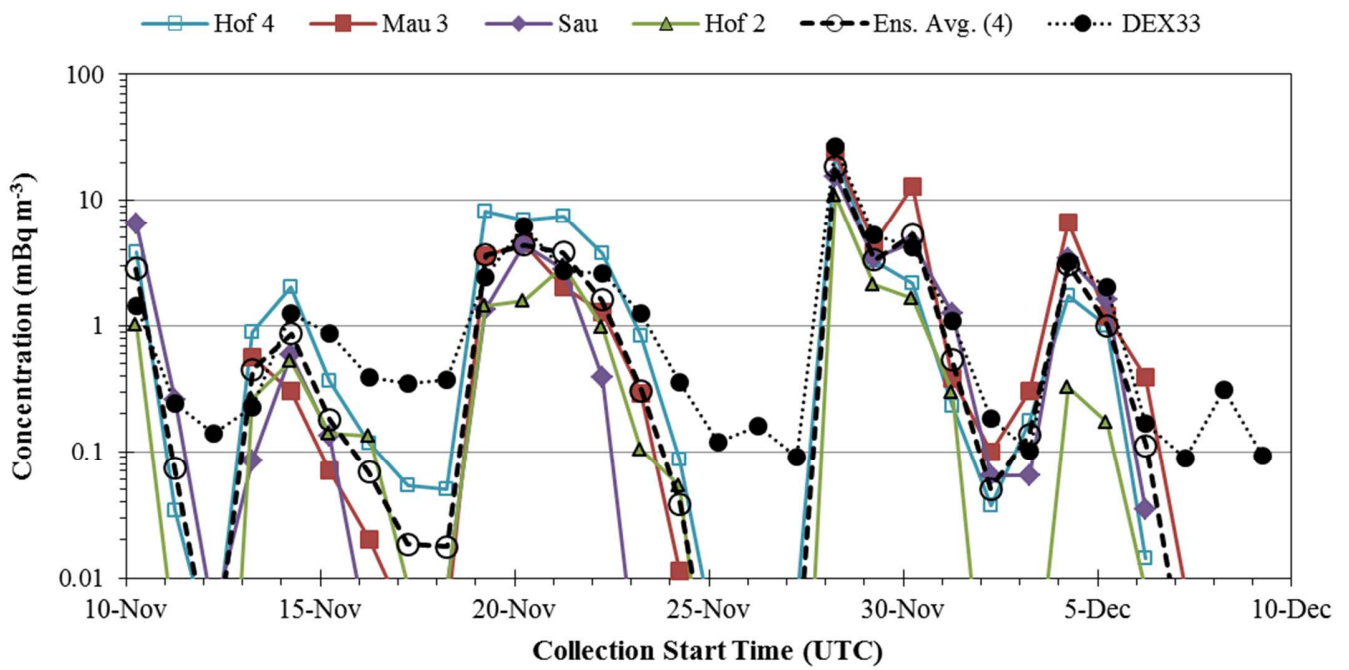
Concentration (mBq/m<sup>3</sup>) averaged between 0 m and 100 m  
Integrated from 0600 14 Nov to 0600 15 Nov 13 (UTC)  
Xeno Release started at 0000 10 Nov 13 (UTC)



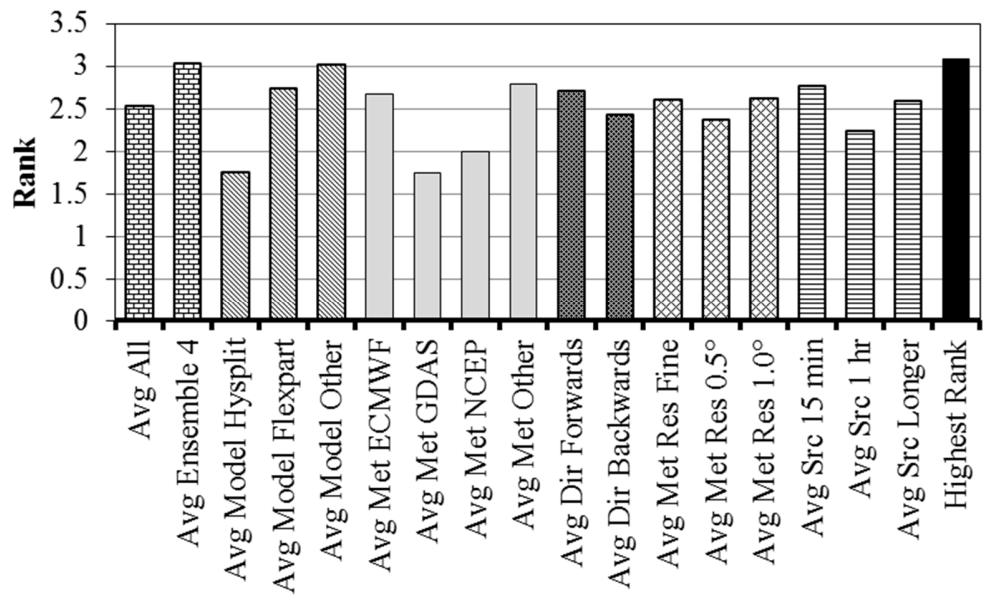
**Fig. 4.** Modeled <sup>133</sup>Xe concentrations using the HYSPLIT computer code and GDAS data corresponding to the DEX33 sample with collection start at 0600 UTC on November 14.



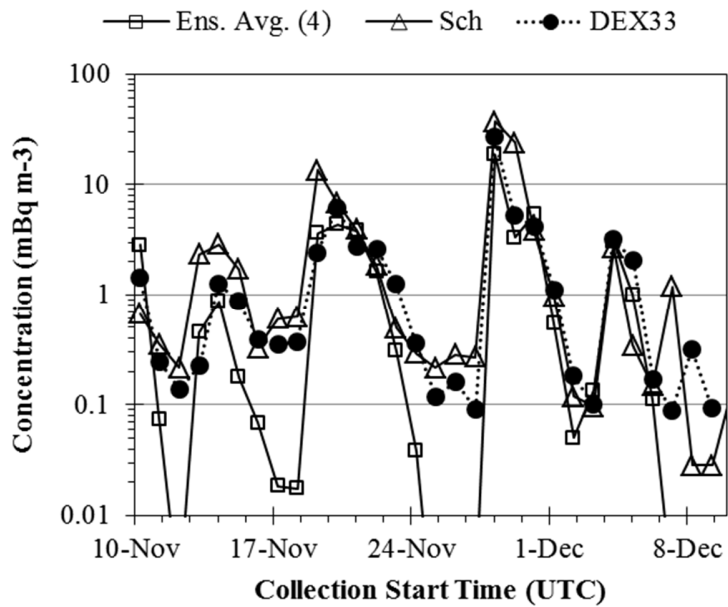
**Fig. 5.** Minimum MSE as a function of the number of submissions in the ensemble.



**Fig. 6.** Modeled  $^{133}\text{Xe}$  concentrations for the individual submissions and the ensemble average for the minimum MSE ensemble of four members.



**Fig. 7.** Rank parameters for grouped model comparisons.



**Fig. 8.** Modeled <sup>133</sup>Xe concentrations for the average of the minimum MSE ensemble of four members and a submission (Sch) that includes emissions from nuclear power plants.

**Table 1**

Participants in the challenge exercise

<b>ID</b>	<b>Name</b>	<b>Organization</b>
Cha	Tianfeng Chai Fong Ngan Ariel Stein Roland Draxler	National Oceanic and Atmospheric Administration (NOAA) Air Resources Laboratory, College Park, Maryland, USA
Esl	Paul W. Eslinger Ted Bowyer Brian Schrom	Pacific Northwest National Laboratory, Richland, Washington, USA
Gen	Pascal Achim Sylvia Generoso	Commissariat à l’Energie Atomique, CEA, DAM, DIF, 91297 Arpajon, France
Hay	Philip Hayes	Air Force Technical Applications Center, Patrick Air Force Base, Florida, USA
Hof	Ian Hoffman Jing Yi Kurt Ungar Alain Malo	Health Canada, Radiation Protection Bureau, Ottawa, Canada Environment Canada, Canadian Meteorological Centre, Dorval, Canada
Kij	Yuichi Kijima	Japan Atomic Energy Agency, Tokai, Ibaraki, Japan
Kry	Monika Krysta	Comprehensive Test Ban Treaty Organization (CTBTO), International Data Center, Vienna, Austria
Mau	Christian Maurer	Zentralanstalt für Meteorologie und Geodynamik, Vienna, Austria
Rob	Peter Robins Verena Heidmann	Atomic Weapons Establishment (AWE), Aldermaston, Reading, RG7 4PR, United Kingdom
Ros	Jens Ole Ross	Federal Institute for Geosciences and Natural Resources (BGR), Hannover, Germany
Sau	Olivier Saunier	French Institute for Radiation protection and Nuclear Safety, Fontenay-aux-Roses, France
Sch	Michael Schoeppner	Program on Science and Global Security, Princeton University, Princeton, New Jersey USA
Sei	Petra Seibert	University of Natural Resources and Life Sciences, Institute of Meteorology and University of Vienna, Faculty of Earth Sciences, Vienna, Austria



**Table 2**

Metadata for models used to explore the effects of common characteristics (see text for definitions of the acronyms)

ID	Code	Met. Data Source	Met. Time Resolution (h)	Met. Spatial Resolution (°)	Model Time Direction	Release Length (h)	Include
Cha	HYSPLIT	WRF	1	27/9 km	Forwards	0.25	Yes
Esl	HYSPLIT	NCEP (GDAS)	3	0.5	Forwards	1	Yes
Gen	FLEXPART	NCEP	6	0.5	Forwards	2	Yes
Hay <sup>a</sup>	WRF		Ensemble	18/6/2 km	Forwards	0.25	Yes
	HYSPLIT						
Hof 1	FLEXPART	ECMWF	3	1	Backwards	3	Yes
Hof 2	FLEXPART	NCEP	3	1	Backwards	3	Yes
Hof 4	MLDP0	CMC	6	0.5	Backwards	3	Yes
Kij	HYSPLIT	NCEP (GDAS)	3	0.5	Forwards	6	Yes
Kry 1	FLEXPART	ECMWF	3	1.0	Backwards	3	Yes
Kry 2	FLEXPART	NCEP	6	1.0	Backwards	6	Yes
Mau 2	FLEXPART	ECMWF	3	0.2	Forwards	0.25	Yes
Mau 3	FLEXPART	NCEP	3	0.5	Forwards	0.25	Yes
Rob	FLEXPART	ECMWF	3	1.0	Backwards	0.25	Yes
Ros 1	HYSPLIT	ECMWF	6	0.2	Forwards	0.25	Yes
Ros 3	HYSPLIT	NCEP (GDAS)	3	0.5	Forwards	0.25	Yes
Sau	Eulerian IdX	ARPEGE	1	0.1	Forwards	0.25	Yes
Sch	FLEXPART	NCEP	1	0.5	Backwards	3	No
Sei 1	FLEXPART	ECMWF	3 <sup>b</sup>	0.2	Backwards	1.25 <sup>c</sup>	Yes
Sei 2	FLEXPART	ECMWF	3	0.2	Backwards	1.25 <sup>c</sup>	No
Sei 3	FLEXPART	ECMWF	1	0.125	Backwards	1.25 <sup>c</sup>	No

- a. This submission was the mean of an 85 member ensemble
- b. Forecasts up to 23 hours are used
- c. Five-sample moving average in time

**Table 3**

Values of the individual statistics and the model rank parameter (Rank) for every model submission. Statistics include the Kolmogorov-Smirnov parameter (KS), Pearson correlation (R), fractional bias (FB), factor of five parameter (F5), normalized mean square error (NMSE) and the mean square error (MSE). Bold values indicate the best score on each statistic.

Model	KS	R	FB	F5	Rank	NMSE	MSE
Sch <sup>a</sup>	0.10	0.89	0.50	0.81	3.25	2.63	19.2
Hof 4	0.39	0.94	0.03	0.61	<b>3.09</b>	<b>0.63</b>	18.3
Mau 3	0.45	0.93	<b>-0.02</b>	0.52	2.92	0.81	<b>3.50</b>
Sau	0.52	0.92	-0.33	0.52	2.68	1.77	5.60
Hof 3	0.45	0.90	-0.58	0.55	2.62	4.25	36.5
Hof 1	0.45	0.75	-0.32	0.58	2.53	3.79	25.9
Hof 2	0.45	<b>0.97</b>	-0.89	0.39	2.43	5.87	25.0
Rob	<b>0.29</b>	0.35	-0.19	<b>0.68</b>	2.41	5.72	20.8
Ros 2	0.52	0.81	-0.56	0.39	2.24	4.87	11.9
Mau 1	0.58	0.79	-0.36	0.35	2.22	3.24	9.90
Ros 1	0.52	0.73	-0.56	0.45	2.18	5.42	13.3
Kry 1	0.42	0.47	-0.42	0.58	2.17	6.41	16.2
Sei 1	0.52	0.46	0.13	0.45	2.08	5.45	25.0
Gen	0.39	0.23	0.36	0.58	2.06	6.56	20.5
Esl	0.45	0.30	-0.08	0.35	1.95	7.62	41.4
Sei 2	0.55	0.43	-0.07	0.35	1.95	6.14	37.5
Kry 2	0.52	0.61	-0.67	0.35	1.87	7.40	27.3
Kij	0.45	0.17	-0.13	0.35	1.87	9.80	40.0
Sei 3	0.58	0.20	-0.03	0.35	1.80	8.89	36.6
Sei 7	0.55	0.19	-0.10	0.35	1.79	9.27	35.7
Sei 8	0.55	0.19	-0.13	0.35	1.78	9.29	59.7
Sei 9	0.58	0.19	0.28	0.32	1.64	10.3	25.5
Hay	0.65	0.71	-1.41	0.16	1.31	26.9	25.3
Cha	0.71	0.83	-1.69	0.06	1.20	62.7	23.2
Mau 2	0.58	0.59	1.75	0.23	1.12	192.	12400
Ros 3	0.55	0.18	-1.17	0.23	1.12	21.5	24.5
Average <sup>b</sup>	0.42	0.69	0.27	0.61	2.53	3.52	19.6

<sup>a</sup> This submission used other sources in addition to the releases from IRE. The statistical performance measures for this submission should not be compared directly with those of other submissions.

<sup>b</sup> The Average row is calculated by averaging all of the modeled values for each sample period and treating the averaged values as atmospheric transport model output.

**Table 4**

Comparison of measured and predicted concentrations (mBq m<sup>-3</sup>) for the five samples with the highest concentrations and the five submissions with highest values of the individual statistics. Statistics include the Pearson correlation (R), model rank (Rank), Kolmogorov-Smirnov parameter (KS) and fractional bias (FB). The Sau submission was a member of the best ensemble with four members

DEX33	Hof 2 (R)	Hof 4 (Rank)	Rob (KS)	Mau 3 (FB)	Sau (Ensemble)	Best 4 Ensemble
6.19	1.58 (-75%)	6.91 (12%)	3.26 (-47%)	4.56 (-26%)	4.38 (-29%)	4.36 (-30%)
26.8	11.1 (-59%)	23.4 (-13%)	4.18 (-84%)	24.5 (-9%)	15.4 (-42%)	18.6 (-31%)
5.28	2.11 (-60%)	3.29 (-38%)	6.21 (18%)	4.48 (-15%)	3.43 (-35%)	3.33 (-37%)

4.18	1.65 (-61%)	2.20 (-47%)	2.34 (-44%)	12.9 (208%)	4.56 (9%)	5.33 (27%)
3.17	0.32 (-90%)	1.75 (-45%)	2.82 (-11%)	6.65 (110%)	3.44 (9%)	3.04 (-4%)

---