

EGU2016-6478

Automatic Event Bulletin Built by Waveform Cross Correlation using the Global Grid of Master Events with Adjustable Templates Ivan Kitov, Dmitry Bobrov, and Mikhail Rozhkov

Abstract. We built an automatic seismic event bulletin for the whole globe using waveform cross correlation at array stations of the International Monitoring System (IMS). To detect signals and associate them into robust event hypotheses in an automatic pipeline we created a global grid (GG) of master events with a diversity of waveform templates. For the Comprehensive Nuclear-Test-Ban Treaty (CTBT), the GG provides an almost uniform distribution of monitoring capabilities and adjustable templates. For seismic areas, we select high quality signals at IMS stations from earthquakes. For test sites, signals from UNEs are best templates. Global detection and association with cross correlation

technique for research and monitoring purposes demands templates from master events outside the regions of natural seismicity and test sites. We populate aseismic areas with masters having synthetic templates calculated for predefined sets of IMS array stations. We applied various technologies to synthesize most representative signals for cross correlation and tested them using the Reviewed Event Bulletin (REB) issued by the International Data Centre (IDC). At first, we tested these global sets of master events and synthetic templates using IMS seismic data for February 13, 2013 and demonstrated excellent detection and location capability. Then, using the REB and cross correlation



Figure 1. Map of REB events found during 2013. Red circles – events found by cross correlation (i.e. two or more REB phases are within 4 s). White circles – not found events –most included 3-C stations not considered for waveform cross correlation at this stage.

An underground nuclear test can be conducted in any place on the planet, not only in seismic regions. Figure 1 presents a map of events from the Reviewed Event Bulletin (REB) as found by the International Data Centre (IDC) in 2013. Figure 2 shows the Seismic Network of the International Monitoring System (IMS) and relative positions of selected historical underground nuclear tests. The difference between geographical distribution of earthquakes and explosions is striking.

The technique of waveform cross correlation (matched filter) is naturally based

on high-quality waveform templates from a IMS arrays are then used to build training representative set of master events. seismic areas, one can use waveform templates from natural events. Similarly, one can use signals from underground nuclear tests where available. And for the rest of the world we need the best possible templates for the purpose of seismic monitoring. Here we present and test the performance of the uniformly distributed Global Grid (GG) of master events designed for the best performance of waveform cross correlation. The events in the REB with detections found by cross correlation at two and more sensors are likely among the most important.

sets for machine learning. The events built by cross correlation and not found in the REB are also used for training as noise.

Figure 3 presents the global coverage by master events. A small segment of the GG is shown in Figure 4. Each node contains a master event with templates at IMS array stations. Since the IMS network is sparse, the number of stations associated with a given master may vary from 3 to 10. The quality. sensitivity, and resolution of an array depend on many factors; its aperture and number of



Preparatory Commission for the Comprehensive Nuclear-Test Ban Treaty Organization. Provisional Technical Secretariat, Vienna International Centre, P.O. Box 1200, A-1400 Vienna, Austria. E-mail: ivan.kitov@ctbto.org

International Data Centre, CTBTO

bulletins (XSELs) experienced analysts from the IDC compared the relative performance of various templates and built reliable sets of events and detections for machine learning. In this study, we carefully compile global training sets for machine learning in order to establish statistical decision lines between reliable and unreliable event hypotheses, then apply classification the intermediate automatic cross procedures correlation bulletin based on the GG, and compile the final XSEL, which is more accurate and has lower detection threshold than the REB.

Conclusion. This work extends the research we presented at the Using the catalogue of seismic events (REB) created by the

Science and Technology 2015 conference regarding the construction of Global Grid of master events for seismic monitoring with waveform cross correlation. IDC, we have done some preliminary estimates of the portion of valid events obtained by waveform cross correlation. Machine learning was a natural choice for selection of valid events from the extremely large amount of data. As in other areas, we doubled the number of events in the REB, with all new events matching the IDC event definition criteria.

We have created a prototype of cross-correlation-based

Figure 2. A map of IMS seismic network with historical UNEs. Blue circles – primary arrays, blue triangles – primary 3-C stations. Yellow circles – auxiliary arrays, yellow triangles – auxiliary 3-C stations. Red starts – underground nuclear explosions. Only primary arrays are used for cross correlation.

> Considering the design and historical performance of the primary IMS arrays we have to distribute them over three quality groups. For a given master, we use arrays from Group 2 and 3 only when no array from Group 1 is available. The largest number of IMS stations associated with a master is 10. The smallest number is 3 as dictated by the strict IDC requirement for a valid REB event. There are some areas where 3 primary arrays are not available. Currently, we are extending the set of stations to use templates from 3-C stations to build master events for nodes in the not covered areas.

Each node or master event is responsible for a circular footprint of ~125 km in radius. The distance between nodes is approximately 140 km. Therefore, the GG covers the whole earth without any blind areas. It is important that the zones of responsibility of neighboring nodes intersect. Figure 5 presents the detailed design of the nodes. For all ($\sim 23,500$) nodes, cross correlation with templates at each associated station is

Figure 5. Local Association and Conflict Resolution 1. Five circles with ~25 km increment in radius. 2. 91 nodes for origin time calculation (green points) 3. All hypotheses at the outer circle (red points) are rejected since they have to be created by neighbouring masters.

where travel times to the stations relevant for the given master event are predicted (Figure 5); the spacing between circles is 25 km and the outer circle has radius of 125 km. Because the outer circle (red points) of a given master intersects with internal circles of the neighboring masters we reject all final event hypotheses built on this outer circle. Such hypotheses must belong to the neighboring events. Other conflicts between event hypotheses belonging to neighboring events are resolved origin times for arrivals on the associated stations.

The current study includes three and conflict resolution, the best events stages. At first, we conducted cross are saved in an XSEL. At this stage, the correlation of the templates based on principal goal is to find more arrivals at the first PCA component in Figure 7 IMS array stations than listed in the with the waveforms likely containing REB for future use in machine signals from all seismic events learning. We define an REB arrival as included in the 2013 REB. The set of found when there is a cross correlation successful detections of the REB detection within 4 s. events (2 and more arrivals for one There are 16.318 REB events event) builds a part of training set. having 2 and more arrivals found by Secondly, we built complete cross cross correlation. At the same time, we correlation bulletins (XSELs) for find many arrivals at other stations of February 13 and 14, 2013. Finally, we the same master, which were not applied various machine learning detected by standard procedures. methods to the training set containing Figure 11 shows the numbers of found valid detections (i.e. in the XSEL REB detections and added detections events matched in the REB) and false for 22 stations. detections (i.e. those in the XSEL For a synthetic template, events not matched in the REB) and theoretical time delays between classified all events hypotheses in the channels and the absence of empirical XSELs for Feb. 13 and 14. The SASCs may lead to poor azimuth and resulting XSEL is considered as a slowness estimates as obtained by the FK-analysis of cross correlation traces. substitute to the REB.

2013, we used 30,513, which have 2 used the closest masters to detect and more primary IMS arrays. Only signals near those existing or expected master events within 1400 km from the in the REB. The number of found and REB events were used to populate the added arrivals has increased. Figure 12 set of valid events. Time windows demonstrates that there are many more included the segment $-5 \min t_0 + 5 \min t_0$ detections with slowness and azimuth from the expected arrival times on residuals beyond the predefined limits. stations related to each master. The The best stations like WRA have one year.

hypotheses and, after local association be use in local association and then in . event.





Global Grid monitoring system, which has been tested at the IDC during the past years and for this study was populated with one synthetic template obtained by the PCA (Principal Component Analysis) of real waveforms from a hundred of underground nuclear tests distributed over the world. One template waveform was replicated over all stations and individual channels. Time delays between individual channels of an array were calculated for master event-station theoretical slownesses.

To use machine learning at the global level it is necessary to create a training dataset populated with valid events created by waveform cross correlation as well as false events created by

same method. Here we use the REB to build a cross correlation event list and process only data at relevant stations. To reduced calculations, we have processed only time intervals around known REB events. Two days, February 13 and 14, were processed continuously and a large number of false events together with associated detections were build.

Overall, we have demonstrated a significant increase in the number of detected arrivals when cross correlation is used. The obtained events and arrivals were used for training of various classification algorithms applied in the framework of continuous processing with cross correlation. The Global CC Grid technique gives an opportunity for

WAVEFORM TEMPLATES, CROSS CORRELATION, EVENT BUILDING



Figure 6. Real seismograms with varying distance, depth, and UNE source functions which were used in PCA.



Figure 7. The PCA 1st component is used as a universal template for all stations and channels

Despite its original design aimed at seismic monitoring of nuclear tests, waveform templates for the GG can be selected from a broader set of known as a signal dimensionality reduction.

The dimensionality reduction allows finding a template best describing the set of records of nuclear explosions conducted in wide ranges of epicentral distance, rock types, yield, and depth of burial. Some examples of such waveforms are presented in Figure 6. These real seismograms have different sampling rate and were recorded by sensors with different responses. Therefore, they were all reduced to one sensor type and sampling rate.

Principal Component Analysis (PCA) performed with Singular Value Decomposition (SVD) is used as a tool to build templates. Figure 7 demonstrates the first principal component obtained by the PCA. This waveform was used to build synthetic templates for all channels and all arrays. The only difference between the templates was in the sampling rate which varies from 20 Hz to 80 Hz for IMS array stations. For a given array, time delays between channels were calculated from theoretical travel time curves according to master/station positions.



ons of cross correlation coefficient between naster templates and real waveforms from the 2013 DPRK for rray stations AKASG and WRA. Almost all masters have detections at AKASG and WRA. These detections detections at other 8 IMS stations are used to build event nypotheses.

USING THE REB FOR TESTING GG PROCEDURES

In order to assess the portion of correct the hypotheses far enough (600 s) from the REB hypotheses which can be potentially built by the events. In total, the class of valid events included 1031 arrivals at 7 IMS array stations GG we have tested a smaller case of seismicity in North Atlantic region. Because of task and there were 5631 arrivals associated with the dimension, machine learning was used with the events not matching REB criteria. After training, we classified 177,520 MATLAB TreeBagger tool as the first choice. arrivals in 48,443 hypotheses. Classification Following general procedure for this algorithm allowed to build 252 new valid events application, we have compiled a dataset defined by three or more valid arrivals as containing two classes of arrivals as associated with valid and invalid (strictly according to the required by the IDC event definition criteria. IDC definition) events. To populate the set of Figure 10 displays locations for these qualified valid events and detections, we selected the events, which repeat the pattern of the Midevent hypotheses obtained by the GG, which Atlantic Ridge seismicity. To be included in the REB, all selected have close origin times to the events reported in hypotheses have to be confirmed interactively.

the REB. As in our previous studies, cross correlation Overall, we selected 258 valid event approximately doubles the amount of REB hypotheses obtained by the GG prototype for the training set. The set of invalid events was events and reduces the detection threshold b created by a random choice of $\sim 3\%$ (1859) of 0.4 magnitude units.

From 33,710 REB events built in Before building event hypotheses, we



To perform machine learning with valid false detections we used scikit-learn (0.17.1, www.scikit-learn.org) - a set of python modules for machine learning and data mining. In order optimize machine learning for monitoring purposes we have tested various applications including: SVM with linear kernel, SVM with Radial Basis Function kernel, Decision Tree, Nearest Neighbors, Randon Forest, AdaBoost, Naive Bayes, Linear Discriminant Analysis, and Quadratic Discriminat Analysis. Cross-validation tests were carried times - 2/3 for learning and 1/3 for test. To assess the output we used **Precision**, i.e. the proportion of instances predicted as positives that were correctly evaluated, Recall, i.e. the proportion of positive instances that were correctly evaluated, and F1. score, the harmonic mean of precision and recall. The harmonic mean is used instead of the arithmetic mean because the latter compensates low values for **NVAR** 0.818 0.850 precision and with high values for recall. On the 0.930 other hand, with harmonic mean we will always SONM have low values if either precision or recall is low. After a thorough investigation we have selected Random Forest algorithm which uses perturb-andcombine techniques specifically designed for trees. This means a diverse set of classifiers is created by USRK introducing randomness in the classifier construction. The prediction of the ensemble given as the averaged prediction of the individual YKA 272 ZALV 0.971 1304 classifiers. Station results are shown in the table. ZALV 224 0.821 0.833 0.827 N

We have tested two sets of parameters used for For the first set, the XSEL includes 166 REB-The second set of parameters produced 10,145 reduced set of masters (around 300 very low proportion of such detections signal classification. First set included only compatible events from the total of 285 in the XSEL events and, after application of the instead of ~25,000 in the GG) and short as they provide reliable estimates of the parameters which were associated with detections: official REB. In this XSEL, there were 153 events classification based on the results of Random Forest time window allow fast calculation for studied parameters. However, judging CC, SNR, residuals of azimuth and having 3 and more arrivals within 4 s from the training, the final XSEL included 349 REBby the distribution of CC and SNR_CC slowness, etc. In the second set, we also used the arrivals in one REB event, and 159 events with 2 compatible events. In this XSEL, 86 events had 3+ All cross correlation detections for for these detections many of them are number of arrivals associated with a given valid or and more found REB arrivals. Essentially, the XSEL arrivals in the REB, and 143 events 2+ arrivals. all masters were used to build event of relatively high quality and thus can false event, the averaged and cumulative CC for this based only on detection parameters can find only In order to achieve higher performance and XSEL reliability we have to merge both approaches. REB events.



| preparatory commission for the | comprehensive nuclear-test-ban | treaty organization

complete implementation of CC-based global detection and

location. The specific features of the P-waves from underground

nuclear tests used in this study can reduce the global detection

threshold of seismic monitoring under the CTBT by 0.4 to 0.5

units of magnitude. This corresponds to the reduction in the test

yield by a factor of 2 to 3 for any location and depth of burial

Considering the history of seismic monitoring of UNEs this is a

crucial improvement which can be also enhanced by a more

effective use of IMS array stations.



Figure 9. Finding the 2013 DPRK with the Global Grid. The number stations associated with the event hypotheses obtained in local association process. Only the hypotheses with 3 and more associated stations are shown. The best hypothesis has 9 associated stations and is the closest to the test position estimated from satellite data. All detections at the involved array stations are found using the first PCA component shown un Figure 7.



Figure 10. Locations of 252 new qualified event hypotheses, which match the REB criteria.

MACHINE LEARNING FOR THE GCCG