

Accepted Manuscript

Reducing process delays for real-time earthquake parameter estimation – An application of KD tree to large databases for Earthquake Early Warning

Lucy Yin, Jennifer Andrews, Thomas Heaton

PII: S0098-3004(17)30596-4

DOI: [10.1016/j.cageo.2018.01.001](https://doi.org/10.1016/j.cageo.2018.01.001)

Reference: CAGEO 4074

To appear in: *Computers and Geosciences*

Received Date: 29 May 2017

Revised Date: 2 December 2017

Accepted Date: 10 January 2018

Please cite this article as: Yin, L., Andrews, J., Heaton, T., Reducing process delays for real-time earthquake parameter estimation – An application of KD tree to large databases for Earthquake Early Warning, *Computers and Geosciences* (2018), doi: 10.1016/j.cageo.2018.01.001.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



1 **Reducing process delays for real-time earthquake parameter estimation – an**
2 **application of KD tree to large databases for Earthquake Early Warning**
3

4
5 Lucy Yin¹, Jennifer Andrews², Thomas Heaton^{1,2}
6

7 ¹Department of Civil and Mechanical Engineering
8 California Institute of Technology
9 Pasadena, USA
10

11 ²Division of Geological and Planetary Sciences
12 California Institute of Technology
13 Pasadena, USA
14

15
16
17 Corresponding Author: Lucy Yin
18
19

20 Department of Civil and Mechanical Engineering
21 California Institute of Technology
22 1200 E California Blvd
23 Pasadena, CA USA 91106
24 Tel: 1-626-841-9702
25 Email: lyin@caltech.edu
26

27 Abstract

28 Earthquake parameter estimations using nearest neighbor searching among a
29 large database of observations can lead to reliable prediction results. However, in
30 the real-time application of Earthquake Early Warning (EEW) systems, the accurate
31 prediction using a large database is penalized by a significant delay in the
32 processing time. We propose to use a multidimensional binary search tree (KD tree)
33 data structure to organize large seismic databases to reduce the processing time in
34 nearest neighbor search for predictions. We evaluated the performance of KD tree
35 on the Gutenberg Algorithm, a database-searching algorithm for EEW. We
36 constructed an offline test to predict peak ground motions using a database with
37 feature sets of waveform filter-bank characteristics, and compare the results with
38 the observed seismic parameters. We concluded that large database provides more
39 accurate predictions of the ground motion information, such as peak ground
40 acceleration, velocity, and displacement (PGA, PGV, PGD), than source parameters,
41 such as hypocenter distance. Application of the KD tree search to organize the
42 database reduced the average searching process by 85% time cost of the exhaustive
43 method, allowing the method to be feasible for real-time implementation. The
44 algorithm is straightforward and the results will reduce the overall time of warning
45 delivery for EEW.

46 Highlights

- 47 • Presented a multidimensional binary search (KD) tree database structure for
48 seismic data

- 49 • Reduced average searching time by 85% for real-time seismology
50 predictions
- 51 • Suggested to directly predict ground motion information for Earthquake
52 Early Warning due to accuracy
- 53 • Evaluated pros and cons of modeling approach and big data search approach
54 for real-time seismology

55 **Introduction**

56 Due to the advancement of information technology in the past few decades,
57 Earthquake Early Warning (EEW) systems are able to analyze ground motions in
58 real-time and provide alerts before the onset of the destructive wave at specific
59 facilities (Heaton, 1985) (Allen et al., 2003). EEW is based on the principle that the
60 damaging earthquake ground motion propagates more slowly than electronic
61 information, so warnings can be successfully delivered immediately after detecting
62 the first earthquake signals at a seismic station (Cua, 2005). The speed of the more
63 damaging S-waves from earthquakes is about 3.5km/s, whereas electrically
64 transmitted signals from the seismic network sensors travel at about 3.0×10^5 km/s.

65 EEW is most beneficial for earthquakes causing a significant level of ground
66 shaking, so the alert speed is critical to provide a warning to the most strongly
67 affected areas close to the epicenter. Additionally, for high-cost user actions (such as
68 halting industrial processes), the accuracy of ground motion predictions at user
69 sites is important for the widespread adoption and use of EEW (Hoshiya, 2013). In
70 general, the conventional algorithms use trained models to estimate earthquake
71 source parameters (such as magnitude and hypocenter distance) from station

72 ground motion observations, and then apply ground motion prediction equations to
73 estimate the peak ground motion experienced at different user sites (Wu et al.,
74 2007) (Zuccolo et al., 2016) (Kuyuk et al., 2014). The predictive models tend to
75 compress the observed information into a few source parameters, which can overly
76 simplify the behavior of wave propagation through the Earth (Meier, 2017).
77 Significant error in final prediction results can be accumulated through the
78 uncertainties in the underlying models (Bose et al., 2009) (Allen et al., 2009). As a
79 result, for the purposes of a real-time EEW system, it is a challenge to create a
80 simple model that fully captures all the attributes that influence the peak ground
81 motion in a recorded waveform, such as magnitude, location, depth, soil type, local
82 site condition, directivity, source radiation.

83 Fingerprint searching and template match methods are alternative approaches to
84 EEW and have also recently been employed in other areas of seismology (Yoon et al.,
85 2015). In the fingerprint searching method, important waveform characteristics are
86 extracted from each earthquake record to form an extensive database of
87 “earthquake fingerprints”. During the occurrence of an on-going earthquake, the
88 algorithm searches among the database for the most similar “earthquake
89 fingerprints”, and then estimates the source parameters or peak ground motions of
90 the new event based on the searched records. A recently developed method, called
91 the Gutenberg Algorithm (GbA) (Meier et al., 2015), applies the fingerprint-
92 searching concept to EEW by abstracting the time-frequency amplitude information
93 of the real-time seismic signal for various filter bands to create a large-scale
94 database, and then estimates the earthquake source parameters such as magnitude

95 and hypocenter distance for on-going earthquakes. In addition, the template-
96 matching method in FinDer (Böse et al., 2012) compares observations with a
97 database of theoretical spatial ground motion patterns to estimate earthquake
98 source parameters and peak ground shaking at various sites. Another example of
99 real-time earthquake monitoring algorithm using search engine method was
100 developed for the estimation of source-focal mechanism (Zhang et al. 2014).
101 Although the predictable variable in the example above are different, all of the
102 methods share the common approach of searching among a pre-processed database.

103 One of the most important factors required of search algorithms is that the
104 searched database needs to be sufficiently large in order to cover a wide range of
105 potential earthquakes. In other words, if similar data to the target query are not
106 included in the database, the searched result could be significantly off from the true
107 value. As an example, the records in the databases should represent the natural
108 distribution of earthquake occurrence as described by the Gutenberg-Richter
109 relationship (Gutenberg and Richter, 1944); there should be many more small
110 events than large ones because small size earthquakes occur more often than large
111 earthquakes, so the search should returns reflect real earthquake likelihoods. Of
112 course, the best strategy is to include all worldwide earthquakes recorded over a
113 long period of time. While increasing the database promises to improve estimation
114 accuracy, the trade-off is that the processing time of searching among a large
115 database increases significantly due to the rise in comparison operations. A simple
116 search of the Advanced National Seismic System (ANSS) Composite Catalog
117 (<http://www.quake.geo.berkeley.edu/anss/>) reveals that 2090 shallow crustal

118 earthquakes (depth <30km) over magnitude 2 occurred in California during 2015.
 119 Similar results are also indicated with searches on USGS/ComCat, Southern
 120 California Earthquake Data Center and other similar earthquake databases. If one
 121 wants to include all records from the network over years for all the earthquake
 122 events worldwide, the size of the database scales exponentially (Yu, 2016). As a
 123 result, the processing delay of the real-time search will significantly increase
 124 because the time required to query databases sequentially is proportional to the
 125 size of the database. While advances have been made in the development of such
 126 algorithms in EEW, very little attention has been paid to optimizing the processing
 127 time of large databases.

128 Database searching is often an application of the Nearest Neighbor (NN) search
 129 problem with the Euclidean metric. The problem is commonly encountered in many
 130 computational techniques such as event detection, pattern recognition, and data
 131 analysis (Bhatia, 2010). In general, we seek for a point in the database that
 132 minimizes the Euclidean distance to the target point (sometimes referred as the
 133 least square distance). The problem states that: for the target point $x =$
 134 $[x_{(1)}, \dots, x_{(d)}]$ and the i^{th} training point in the database $y_i = [y_{i(1)}, \dots, y_{i(d)}]$, we define
 135 the distance between x and y_i to be

$$d(x, y_i) = \left(\sum_{k=1}^d (x_{(k)} - y_{i(k)})^2 \right)^{1/2} \quad [1]$$

136 NN searches for the \hat{y} with the closest distance to the target point, mathematically
 137 represented as $\hat{y} = \operatorname{argmin}_{y_i} (d(x, y_i))$. In most cases, the k-Nearest-Neighbor (k-
 138 NN) search method is applied by finding the k closest training points to the target

139 point; this method provides a more robust estimation that avoids outliers in the
140 database. The corresponding parameters associated with the \hat{y} are used to classify
141 or estimate the parameters of interest for the target point.

142 In this study, we use a data structure, multidimensional binary search tree (KD
143 tree) NN searching concept, to organize the seismic data, and evaluate the reduction
144 of NN searching time for large datasets. KD-tree is a binary tree data structure that
145 links the relative position of all the data points, so data with similar patterns cluster,
146 thereby allowing the search procedure to become faster (Bentley J. L., 1975).
147 Although it requires initial effort to construct the tree data structure, the searching
148 process is quick. The goal is to introduce the concept of data structures in EEW to
149 minimize the processing time for waveform record searching without loss of
150 accuracy, and thereby earthquake alerts can be delivered to the sites of interest
151 much earlier. The effectiveness of fast alerts is especially valuable in the proximity
152 of the epicenter where the strongest damage occurs very quickly after event onset.
153 In this study, we describe a searching procedure that uses the KD tree NN search
154 method that identify the EEW fingerprints characterized by the Gutenberg
155 Algorithm. We 1) evaluate the influence of database size on the prediction accuracy
156 of the earthquake source parameters (magnitude and hypocenter distance) and
157 peak ground motion parameters (PGA, PGV, PGD), 2) estimate the processing
158 efficiency of the KD tree searching for databases with different sizes and extrapolate
159 the future performance by scaling to larger data sets. The KD tree is well-established
160 NN searching algorithm that has been implemented in a wide range of engineering
161 and database applications (Bentley J. , 1979). Although the method parallel

162 processing can reduce real-time latencies, the cost of allocating additional resources
163 could be unfeasible in long term. Only by efficiently design the computational
164 algorithms to optimize the processing time can EEW start to adopt the databases for
165 real-time seismology applications, and the fingerprint searching algorithms with big
166 data reveal their full practical potential.

167 **Data**

168 Theoretical analysis of the KD tree searching shows the performance complexity
169 being $O(\log N)$ verses $O(N)$ for the linear sequential search, where N is the number
170 of data points in the database (Friedman et al., 1977). Although the theoretical
171 average search time of KD tree is much shorter than the linear sequential search, the
172 performance varies depending on the distribution of the data. Our goal is to
173 determine the searching efficiency of the KD tree method for our GbA seismic
174 database. We ran a series of offline tests on the earthquake filterbank database to
175 mimic potential performance of EEW using true seismic records. The dataset used is
176 pre-processed by (Meier et al., 2015) for the GbA. The database consists of 182,805
177 near-site records with 9 feature dimensions in each record. Each of the feature
178 dimensions represents the peak ground velocity in octave-wide frequency bands for
179 a given ground motion record with a fixed time window. The frequency bands used
180 in GbA features are shown in Table 1. GbA creates such a dataset table for every
181 half-second increment in time after the P-wave arrival. In general, EEW tends to
182 consider at least 3 to 4 sec data after the P-wave arrival for the trade-off of accuracy
183 and time delay. For the purpose of this investigation we selected the database for a
184 10-sec time window because the predictions are stabilize with more data collection.

185 The collected earthquakes cover a large range of magnitude, spanning from M 2.0 to
 186 M 8.0, compiled from shallow crustal earthquakes collected from Japan, Southern
 187 California, and Next Generation Attenuation-West 1 (Chiou and Youngs, 2008).

188

Feature Dimension No.	Frequency Band (Hz)
1	0.09375 – 0.1875
2	0.1875 – 0.375
3	0.375 – 0.75
4	0.75 – 1.5
5	1.5 – 3
6	3 – 6
7	6 – 12
8	12 – 24
9	24 – 48

189 Table 1. Frequency bands for feature input in Gutenberg Algorithm. The GbA
 190 database consists of 9 feature dimensions. Each feature takes the observed peak
 191 ground velocity in the given frequency band.

192

193 **KD Tree and Method**

194 KD Tree

195 KD tree is a binary tree structure that stores the finite set of database points with
 196 k-dimensional feature space. In our case, we have 9 variables corresponding to 9-
 197 dimensions. The method involves two steps. First, we construct the tree to organize
 198 the information in the database. Then, the NN algorithm is applied on the KD tree to
 199 search to the most similar point to the target record during an on-going earthquake.

200 In KD tree implementation, a point in the database is also called a node in the tree.

- 201 • Construction of KD-tree

202 The construction of the KD-tree is a recursive process. Starting with the root of
 203 the tree, the first feature dimension (frequency band: 0.09375 – 0.1875Hz) is chosen
 204 as the splitting hyperplane. All nodes are ordered with respect to the value in this

205 feature dimension, and the node with the median value is inserted into the root of
206 the tree. All nodes with coordinates less than the median in the splitting hyperplane
207 create the left subtree, and the nodes with coordinates larger than the median in the
208 splitting hyperplane create the right subtree. All the feature dimensions rotate in
209 becoming the splitting hyperplane to create the next level of subtrees.

210 • Nearest Neighbor Search in KD-tree

211 Starting with the root node of the tree, the nearest distance is initialized to be the
212 distance between the target node to the root. Then recursively move down to the
213 next level in the tree, and checks if the splitting hyperplane intersects with the
214 hypersphere centered at the target record with a radius of the current nearest
215 distance. If the node falls outside of the hypersphere created by the current nearest
216 node (indicating the point is further to the target node than the current nearest
217 node), then this node and any extended child nodes further away can be eliminated
218 from the investigation. The process is repeated, recursively moving down to the
219 next level in the tree until reaching the leaves of the tree. The searching time is
220 reduced since large subsets of the database are not visited. Therefore, the average
221 searching time in a KD tree is significantly lower, especially when the size of the
222 database is large.

223 To better visualize the concept, Figure (1) demonstrates a KD tree structure for a
224 2-dimensional featured database with 10 earthquake records described by the peak
225 velocity and acceleration at initial 3 sec after triggering a station. The goal is to
226 predict magnitude of the new event based on the velocity and acceleration recorded
227 at the first 3 sec of the p-wave. We start the search process of the nearest neighbor

228 of the target data (the yellow star) with node E, which is the root of the tree. The
229 radius of the initial hypersphere is set between the target data and node E. In a 2D
230 feature space, the hypersphere is simply a circle. Since the left branch (link between
231 node A and E) does not cross the hypersphere, indicating all the nodes in the left
232 subtree (node C, A, B, D) can be eliminated from the search because their Euclidean
233 distance to the target point is clearly further than node E. This eliminates the
234 computational effort of going through almost half of the database at the first step.
235 Since the target node is closest to node H, the magnitude associated with node H
236 ($M=4.0$) is the prediction result for the target node.

237

238

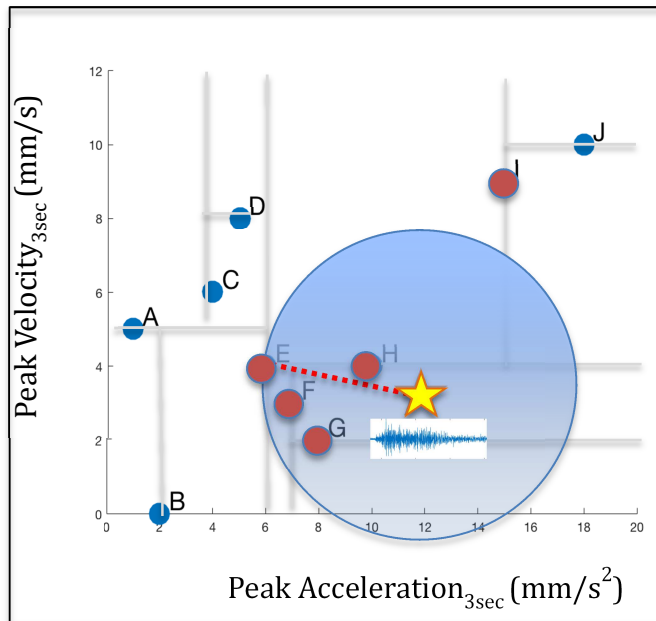
239

240

241

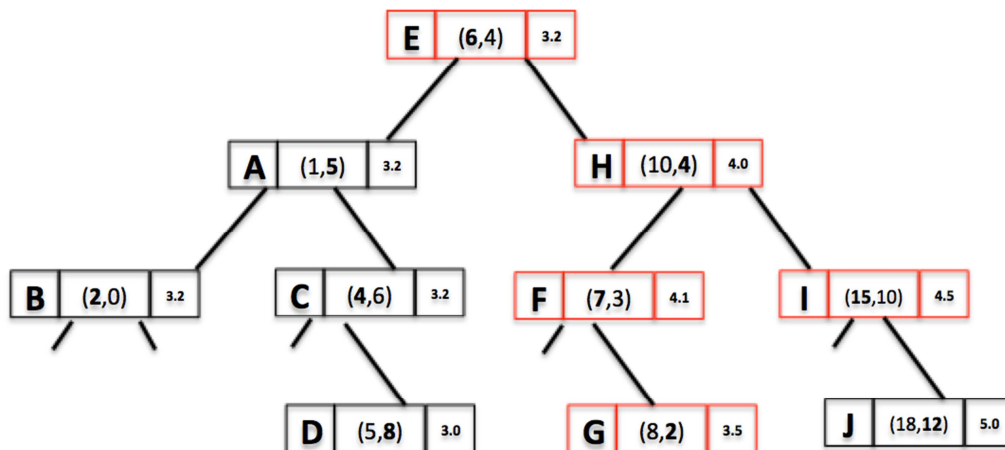
242

243 a)



244
245

b)



246
247

248 Figure 1: A 2-dimensional KD tree example: a) visual distribution of the database in
249 feature dimensions, b) tree structure of the database. A database of 10 earthquake
250 records (A - J) is organized using KD tree data structure (grey lines are the branches
251 of the tree). As a comparison, the linear sequential search requires going through all
252 10 records, which doubles the computation effort.

253
254

The algorithm can be easily extended to k nearest neighbor (k-NN) search to find

255 k most similar points to the target point in order to give a more probabilistic

256 estimate of target parameters. It requires two modifications. First, we need to keep

257 track of all the current nearest points in an ordered queue with length k ; if the
258 queue contains fewer than k points, the subtrees on both sides need to be visited.
259 Second, instead of comparing the splitting hyperplane with the hypersphere of the
260 nearest point, we should check if the hyperplane intersects with the hypersphere of
261 the last nearest point in the queue. If they intersect, the new node is inserted into
262 the queue of k -nearest neighbors to the target point. At the end of the search, the
263 algorithm returns k points from the database that are located with minimum
264 distances to the target point.

265 Method

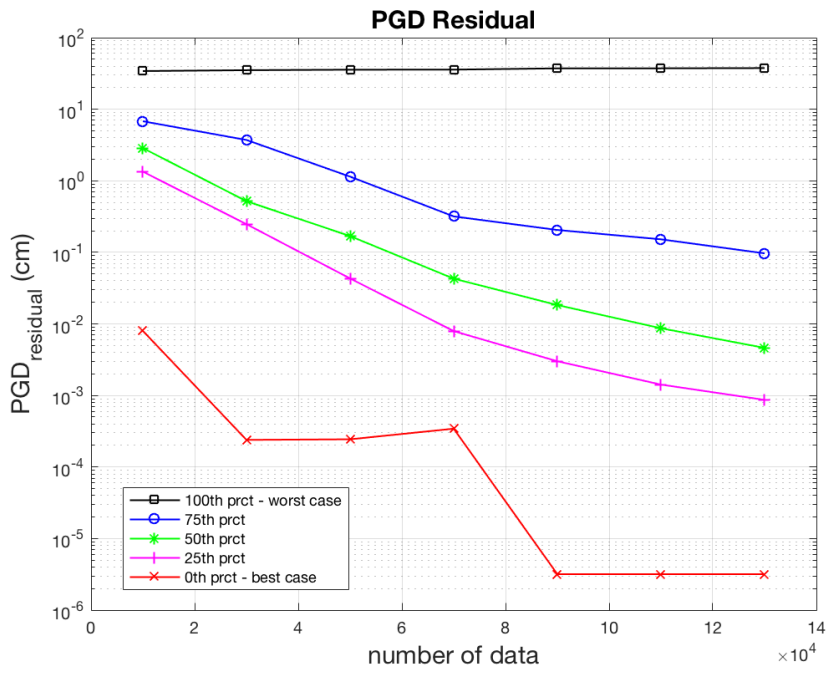
266 Since one of the ultimate goals of EEW aims to predict ground shaking, we
267 extracted 500 records from the entire database to validate the prediction of
268 earthquake source and ground motion parameters. The validation set was sampled
269 uniformly with even spacing on the Peak Ground Acceleration (PGA) of the records.
270 The reason is to cover the full spectrum of ground shaking intensity, in order to
271 mimic all circumstances that could be encountered in the future. The performance
272 of parameter estimations is evaluated with different dataset sizes. The estimated
273 seismic parameters include station-specific ground motions: Peak Ground
274 Acceleration (PGA), Peak Ground Velocity (PGV), Peak Ground Displacement (PGD),
275 and earthquake source parameters: magnitude, hypocenter distance. The procedure
276 first requires a 30-NN search in the Euclidean distance defined in Eq (1), and then a
277 prediction using the Gaussian mean of the corresponding parameters from the 30-
278 NN matched records. The value 30 is chosen to match the original model parameter

279 in the GbA. Later, we compared the searching time of the KD-tree search to the
280 Linear Sequential search, in both the CPU time and the number of operations.

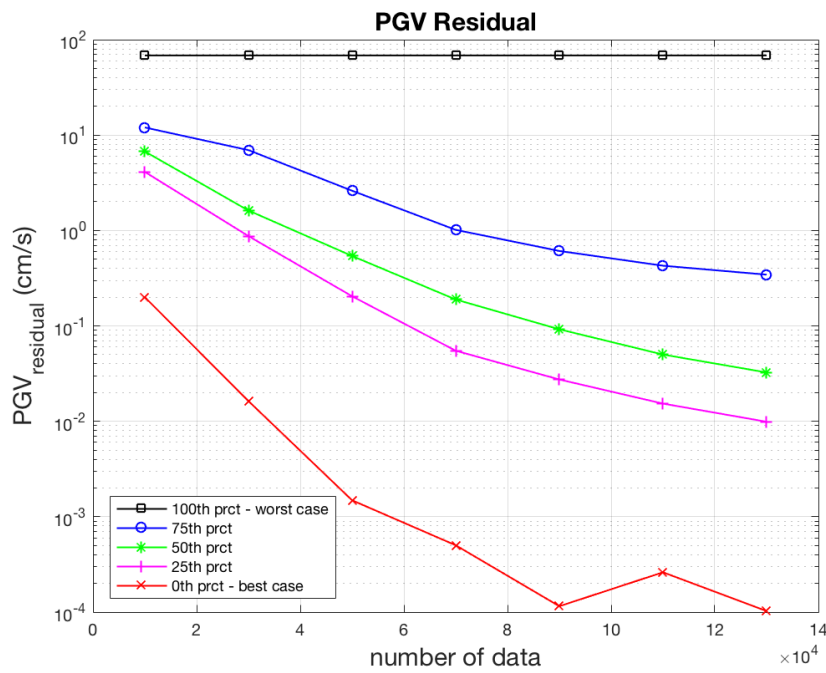
281

282 **Results**

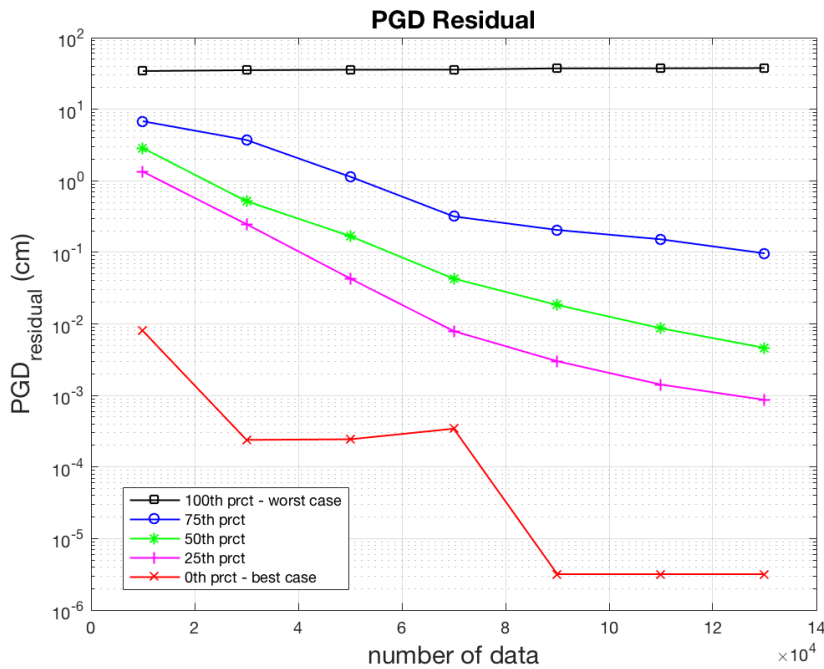
283 We computed the earthquake parameter estimation error of the validation set for
284 databases with different sizes. Figure 2 shows the 100th, 75th, 50th, 25th, and 0th
285 percentile residual errors for the estimated PGA, PGV, and PGD of the 500-validation
286 dataset, respectively. The residual error is defined as the absolute difference
287 between the true observed parameter and the predicted parameter. The 50th
288 percentile is the average residual errors; the 100th and 0th percentile indicate the
289 maximum error and minimum error, respectively. As expected, the residual error
290 decreases as the database size increases on average. The 50th percentile is not
291 flattened near the largest given database size showing that the residual errors might
292 not yet reached the global minimum; this suggests that the estimation accuracy
293 could further be improved by increasing the size of the database. Since there is
294 always a possibility of outlier data regardless how large the database gets, the
295 maximum error residuals are not affected by the size of databases as shown in the
296 100th percentile line in Figure 2. Statistically, there will always be residuals on the
297 predictions, unless the features are truly uniquely diagnostic. Of course, if a
298 sufficiently large database were compiled, the probability of encountering outliers
299 would decrease.



300 a)



301 b)



302 c)

303 Figure 2. Ground motion residuals for the 500-validation dataset with different
 304 database sizes. a) Peak Ground Acceleration, b) Peak Ground Velocity, c) Peak
 305 Ground Displacement residuals are given in absolute ground motion units. The lines
 306 show the percentile according to the legend. The 50th percentile is the average
 307 residual error; the 100th and 0th percentiles indicate the maximum and minimum
 308 errors respectively.

309

310 We also estimated the earthquake source parameters using the databases:

311 magnitude and hypocenter distance. Although ground motion parameters are more

312 useful outputs for EEW alerts, predicting source parameters is the conventional

313 approach in real-time seismology (Minson et al., 2017). Figure 3 shows that the size

314 of the database has less impact on hypocenter distance than magnitude estimation.

315 Since hypocenter distance predictions from the observed waveform are a result of

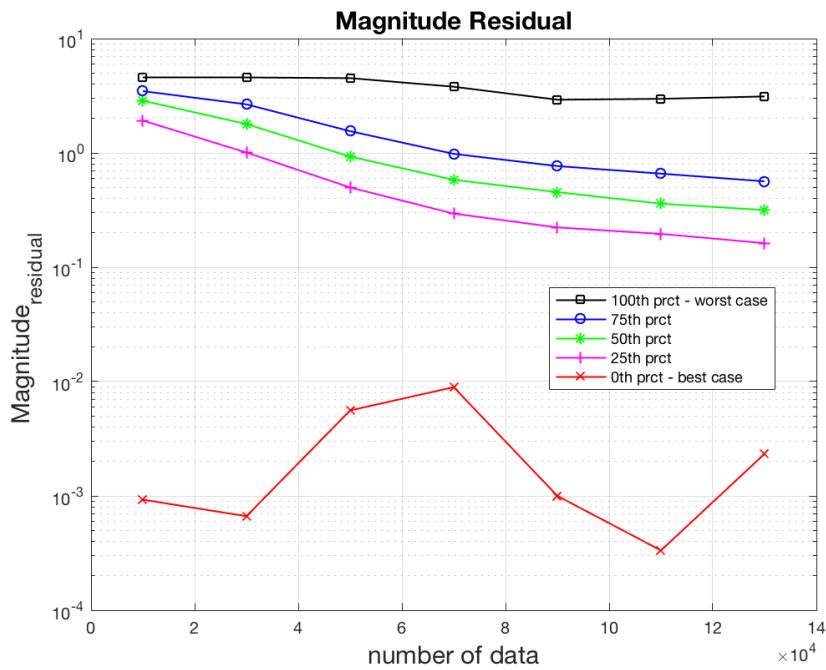
316 source energy and soil properties, the additional constraints might be necessary.

317 For example, seismicity location forecast could be introduced as prior knowledge to

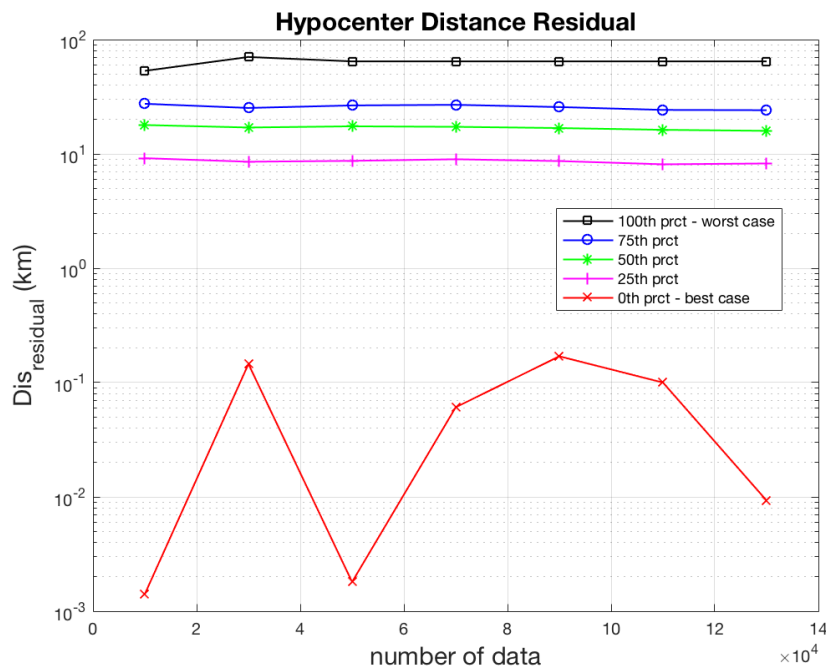
318 reduce the uncertainties in earthquake location estimation (Yin et al., 2017). This

319 analysis implies that it is essential to select data features intelligently to

320 characterize the parameters we are aiming to predict. Frequency band features
 321 might be more suitable to predict the ground motions than source parameters, since
 322 local site effects may be implicitly being accounted for.
 323



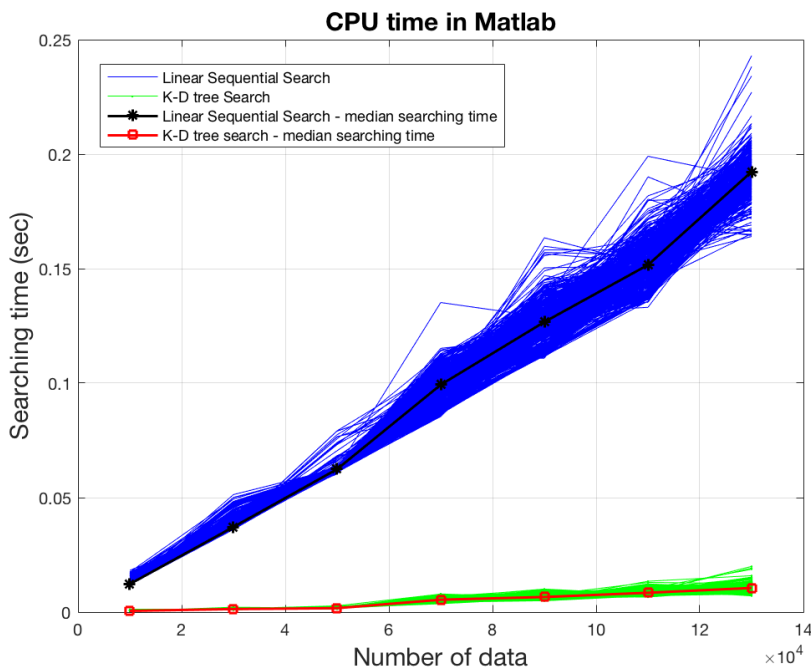
324 a)



325 b)

326 Figure 3: Source parameter residual for the 500-validation dataset with different
327 database size. a) Magnitude, b) hypocenter distance residuals are given in absolute
328 units. The lines show the percentile according to the legend. The 50th percentile is
329 the average residual error; the 100th and 0th percentiles indicate the maximum and
330 minimum errors respectively.

331
332 Through the performance analysis for databases with different sizes, we conclude
333 that large databases can help to provide more accurate ground motion estimations
334 for EEW. Next, we compare the computational time difference for the 30-NN search
335 using the KD tree methods for each validation test. The implementation is in Matlab.
336 For comparison, a Linear Sequential search method is also implemented as a base
337 case. The Matlab function follows the pseudo code concept from the Appendix with
338 optimization modules that efficiently process the data. In Figure 4, the solid lines
339 show that the average CPU search time of a database with 130, 000 points is about
340 0.2 sec for the Linear Sequential search method and 0.03 sec for the KD tree search
341 method; the significant reduction in time reduces computational effort by 85%.
342 Although the Linear Sequential search is capable of handling the real-time
343 processing with limited delay using the current size of the database, a significant
344 delay would be introduced as the database size rapidly increases in the future. The
345 dashed lines show extrapolated computational time up to double of the current
346 database size. The results anticipate that the advantages of the KD tree application
347 would be emphasized in the future as global seismic databases are growing
348 significantly (Yu, 2016).



349

350 Figure 4: CPU searching time for different database sizes using linear sequential
 351 search and KD tree search. The implementation is in Matlab.

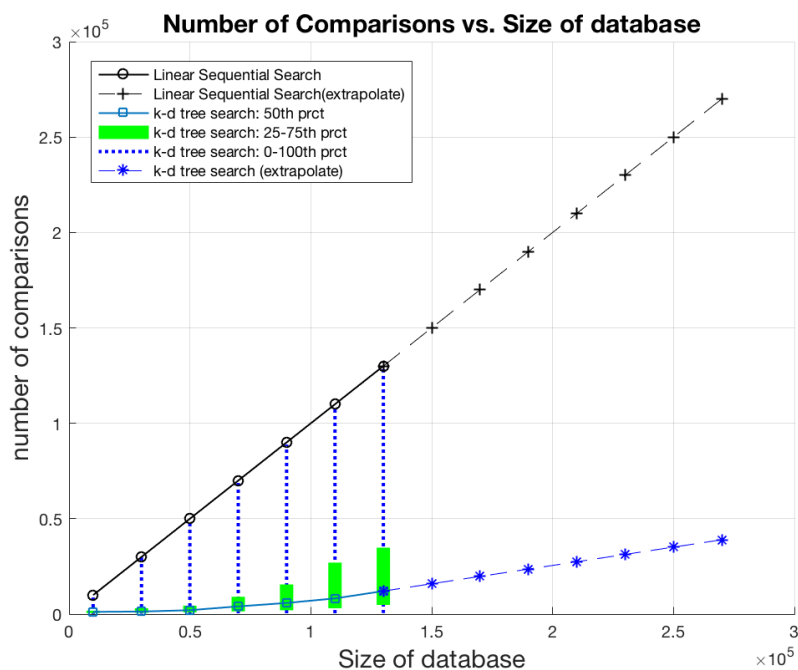
352

353

The measured operational time for the searching process varies significantly
 354 between different software languages and implementations; different optimization
 355 modules with parallelization might also bias towards one method over another.
 356 Implementations in C++ tend to be much faster than Matlab. In order to compare the
 357 true efficiency of the method across all platforms, we further compared the number
 358 of data points visited for both NN search algorithms. Since the majority of the
 359 searching time is made up by the visit to each data point to compute the Euclidean
 360 distance to the target point, the fewer data points visited ensures less time effort. In
 361 the Linear Sequential search, the operation is required for all the data in the
 362 database in a serial manner. However, in KD tree, subsections of the database can be
 363 eliminated depending on the distribution of the tree structure and location of the
 364 target point. As shown in Figure 5, the number of data points visited in the KD tree

365 for each validation varies; on average, the KD tree approach only visits about 10% of
 366 the entire database to find the closest data point to the target, confirming the
 367 performance in CPU searching time in Matlab. In the worst-case scenario, all the
 368 data points are visited, which leads to the same operational complexity as the
 369 exhaustive approach (linear sequential search).

370



371

372 Figure 5 Number of data points visited for linear sequential search and KD tree
 373 search. The dashed lines are extrapolated to estimate the performance for larger
 374 database in the future.

375

376

377 Discussion and Conclusion

378 In this study, we evaluated the viability of earthquake fingerprint searching

379 methods for EEW, using database structure to reduce searching time for large

380 databases. Specifically, we evaluated the GbA as an example of the EEW fingerprint

381 search algorithm. We found that database size is a critical factor in providing

382 reliable predictions of ground motion (PGA, PGV, PGD) and source parameters
383 (magnitude and hypocenter distance) for EEW. We also present the KD tree
384 approach to reduce the searching time, so that large database searching is feasible
385 for real-time implementations in EEW. By empirical validation, we demonstrated
386 that the searching time using KD tree can be approximately 85% less than the
387 exhaustive approach for the GbA EEW earthquake database. (Strauss et al, 2017)
388 has studied extensively on the cost-benefit effects of a warning system in the United
389 States; the study has shown that the number of injuries from earthquakes can be
390 reduced by more than 50% if EEW can provide timely and accurate alerts.

391 One of the potential applications of the database searching method is to directly
392 estimate peak ground motions from the observed ground motions for any given site
393 in real-time seismology application such as EEW; it avoids the multi-step modeling
394 errors that could be accumulated through source parameter estimation and the
395 ground motion attenuation relationship, since the final errors can lead to significant
396 uncertainties in the final shaking information. Ideally, the goal of EEW is to serve as
397 an alarm for severe ground shaking in real-time rather than source characterization.
398 The fingerprint searching methodology could also be extended to tackle other
399 challenges in EEW, such as event detection (i.e. earthquake/noise discrimination).
400 In such a problem, characteristics of additional ambient noise and teleseismic
401 records need to be incorporated in the database. This would vastly increase the
402 database size, since incorporating many different types of noise, teleseisms, regional
403 events, and calibration/maintenance signals could potentially be huge. The vision is

404 to be able to accomplish efficient searching for large databases, so that these novel
405 EEW methods are feasible in real-time in the future.

406 Although we emphasized the importance of having a large number of data, a
407 question is often raised about what should be the minimum size of database in order
408 to get reasonable accurate solutions. Assuming the standard deviation of $\log_{10}(\text{PGV})$
409 estimation of 0.309 by (Kanamori, 2007) is acceptable, the database size needed to
410 achieve this marginal error of ground motion in EEW is about 70 000 to 100 000
411 data points, as shown in Figure 2b). The (Kanamori, 2007) study focuses on two
412 EEW parameters, τ_c and P_d , that are extensively used in the existing EEW
413 algorithms, such as Onsite (Bose et al., 2009). The minimum database size calculated
414 varies with geological region, event types, predictive parameters, etc.

415

416 Creating a database for real-time seismology is not simple. In addition to the sizes
417 of databases, feature engineering also significantly affects the prediction results.
418 Selecting parameters that correlate to the predictive results requires extensive
419 scientific domain knowledge. In the observation of local earthquake records, the
420 higher frequency band features are more informative than the low frequency
421 features because the high frequency amplitude of ground motion decays rapidly
422 with distance (Hanks & McGuire, 1981) (Kong & Zhao, 2012). Although it is out of
423 the scope of this study, we plan to further investigate in the effects of using a
424 weighted Euclidean distance in the Nearest Neighbor Search to emphasize the high
425 frequency information as a significant attribute in the feature space. Continuous
426 monitoring and modifying of the features will help to improve the performance of

427 the system. As the number of features increases, the process time saved by KD tree
428 search decreases (Andoni & Indyk, 2008). For features over 20 or 30 dimensions,
429 alternative approximation to approach high dimensional searching, such as Locality
430 Sensitive Hashing, would be more appropriate [e.g. (Yoon et al., 2015)].

431 EEW is an interdisciplinary project that involves collaboration among different
432 scientific and engineering communities. The accuracy and speed of rapid
433 earthquake source parameter algorithms has significantly improved over the past
434 decade, but are potentially limited by the simplification involved in model
435 parameterization. The earthquake fingerprint searching techniques have the
436 capacity to guide the development of EEW to a new phase with the assistance of
437 better computational power and data mining techniques.

438 **Acknowledgement**

439 We would like to thank Dr Men-Andrin Meier for providing the preprocessed
440 seismic database from GbA used in this study (obtained on September 2016),
441 including Japanese waveforms data obtained from <http://www.kik.bosai.go.jp>, The
442 Next Generation Attenuation – West 1 data obtained from <http://peer.berkeley.edu>,
443 and southern California waveform from Southern California Earthquake data center
444 obtained from <http://www.scedc.caltech.edu>. All of the simulations and figures are
445 generated with Matlab 2016a. This research was supported by the Gordon and Betty
446 Moore Foundation grants 3023 and 5229, and USGS/NEHRP Cooperative agreement
447 G16AC00355. This research was also supported by the Natural Sciences and
448 Engineering Research Council of Canada's (NSERC) Postgraduate Scholarships-
449 Doctoral Program.
450

451 **Bibliography**

- 452 Allen, R., & Kanamori, H. (2003). The potential for Earthquake Early Warning in
453 Southern California. *Science*, *300*, 786-789.
454
455 Allen, R., Brown, H., Hellweg, M., Khainovski, O., & Lombard, P. (2009). Real-time
456 earthquake detection and hazard assessment by ElarmS across California. *36*, 2009.
457

- 458 Andoni, A., & Indyk, P. (2008). Near-Optimal Hasing Algorithms for Approximate
459 Nearest Neighbor in High Dimensions. *Communications of the ACM* , 51 (1), 117-122.
460
- 461 Böse, M., Heaton, T. H., & Hauksson, E. (2012). Real-Time Finite Fault Rupture
462 Detector for large earthquakes. *Geophysical Journal International* , 191, 803-812.
463
- 464 Bentley, J. L. (1975). Multidimensional binary search trees used for associative
465 searching. *Comm. ACM* , 18, 509-517.
466
- 467 Bentley, J. (1979). Multidimensional Binary Search Trees in Database Applications.
468 *IEEE Transactions on Software Engineering* , SE-5 (4), 333-340.
469
- 470 Bhatia, N. (2010). Survey of Nearest Neighbor Techniques . *International Journal of*
471 *Computer Science and Information Security* , 8 (2), 302-305.
472
- 473 Bose, M., Hauksson, E., Solanki, K., Kanamori, H., & Heaton, T. (2009). Real-time
474 testing of the on-site warning algorithm in southern California and its performance
475 during the July 29 2008 Mw 5.4 Chino Hills earthquake. *Geophysical Research Letters*
476 , 36, L00B03.
477
- 478 Cua, G. (2005). *Creating the virtual seismologist: developments in earthquake early*
479 *warning and ground motion characterization*. Ph.D. Thesis, California Institute of
480 Technology, Department of Civil Engineering.
481
- 482 Friedman, J., Bentley, J., & Finkel, R. (1977). An algorithm for finding best matches in
483 logarithmic expected time. *ACM Transactions on Mathematical Software* , 3 (3), 209-
484 226.
485
- 486 Gutenberg, B., & Richter, C. (1944). Frequency of earthquakes in California. *Bulletin*
487 *of the Seismological Society of America* , 4, 185-188.
488
489
- 490 Hanks, T., & McGuire, R. (1981). The Character of high-frequency strong ground
491 motion. *Bulletin of the Seismological Society of America* , 71 (6), 2071-2095.
492
- 493 Heaton, T. H. (1985). A Model for a Seismic Computerized Alert Network. *Science* ,
494 228 (4702), 987-990.
495
- 496 Hoshiya, M. (2013). Real-time prediction of ground motion by Kirchhoff-Fresnel
497 boundary integral equation method: Extended front detection method for
498 earthquake early warning . *J. Geophys. Res* , 118 (3), 1038-1050.
499
- 500 Kong, Q., & Zhao, M. (2012). Evaluation of Earthquake Signal Characteristics for
501 Early Warning. *Earthquake Engineering and Engineering Vibration* , 11, 435-443.
502

- 503 Kuyuk, H., & Allen, R. (2014). Designing a Network-Based Earthquake Early Warning
 504 Algorithm for California: ElarmS-2. *Bull. Seismo. Soc. Am.* (104), 162-173.
 505
- 506 Meier, M., Heaton, T., & Clinton, J. (2015). The Gutenberg Algorithm: Evolutionary
 507 Bayesian Magnitude Estimates for Earthquake Early Warning with a Filter Bank.
 508 *Bulletin of the Seismological Society of America* , 105 (5).
 509
- 510 Meier, M. (2017) How 'Good' are Real-Time Ground Motion Predictions from
 511 Earthquake Early Warning Systems? *Journal of Geophysical Research* (122), 5561–
 512 5577.
 513
- 514 Minson, S., Wu, S., Beck, J., & Heaton, T. (2017). Combining Multiple Earthquake
 515 Models in Real Time for Earthquake Early Warning. *Bulletin of the Seismological*
 516 *Society of America* .
 517
- 518 Strauss, J., Allen, R., (2017). Benefits and Costs of Earthquake Early Warning.
 519 *Seismological Research Letters* (87), 765 – 772.
 520
- 521 Wu, Y.-M., Kanamori, H., Allen, R., & Hauksson, E. (2007). Determination of
 522 earthquake early warning parameters, tc and Pd, for southern California. *Geophys.*
 523 *J.Int.*
 524
- 525 Yin, L., Meier, M., & Heaton, T. (2017). Making Earthquake Early Warning faster and
 526 more accurate using ETAS seismicity models as Bayesian Prior. *16th World*
 527 *Conference on Earthquake Engineering*.
 528
- 529 Yoon, C., O'Reilly, O., Bergen, K., & Beroza, G. (2015). Earthquake detection through
 530 computationally efficient similarity search. *Science Advances* , 1-13.
 531
- 532 Yu, E. (2016). Products and Services Available from the Southern California
 533 Earthquake Data Center (SCEDC) and the Southern California Seismic Network
 534 (SCSN). *AGU Fall Meeting*, (pp. S53A-2817).
 535
- 536 Zhang, J., Zhang H., Chen, E., Zheng, Y., Kuang, W., Zhang, X. (2014). Real-time
 537 earthquake monitoring using a search engine. *Nature Communications* 5:5664.
 538
- 539 Zuccolo, E., Bozzoni, F., & Lai, C. (2016). Regional Low-Magnitude GMPE to Estimate
 540 Spectral Accelerations for Earthquake Early Warning Applications in Southern Italy.
 541 *Seismological Research Letters* .
 542
 543
- 544 **Appendix – Pseudo Code**
 545 KD Tree Construction
 546 Function construction_kdtree(points in database, depth)
 547 Split_axis=depth mod k_dim;

```
548     Median = select median from fingerprints(split_axis)
549     Leftdatabase={fingerprints in database| fingerprints(split_axis)<median}
550     Rightdatabase={fingerprints in database| fingerprints(split_axis)>median}
551
552     %create node
553     node.location=median
554     node.left= kdtree(leftdatabase,depth+1)
555     node.right= kdtree(rightdatabase,depth+1)
556
557     Searching in KD tree
558     Function search_kdtree(target, node, nearest_dist)
559     Split_dim= split dimension at the depth of the tree
560     Hyperplane_dist=target(split_dim)-node(split_dim)
561
562     If nearest_dist > |hyperplane_dist| then
563         curr_dist:= distance between target and curr_node
564         If curr_dist<nearest_dist then
565             nearest_dist:=curr_dist
566             nearest_fingerprint=curr_node
567         search_kdtree(target, curr_node.left, nearest_dist)
568         search_kdtree(target, curr_node.right, nearest_dist)
569     else
570         if hyperplane_dist<0 then
571             search_kdtree(target, curr_node.left, nearest_dist)
572         else
573             search_kdtree(target, curr_node.right, nearest_dist)
574
```

Highlights

- Presented a multidimensional binary search (KD) tree database structure for seismic data
- Reduced average searching time by 85% for real-time seismology predictions
- Suggested to directly predict ground motion information for Earthquake Early Warning due to accuracy
- Evaluated pros and cons of modeling approach and big data search approach for real-time seismology