```cpp
 1 #include "wordBreakdown.h"
 2
 3 using namespace std;
 4
 5 //GLOBALS for database access
 6 sqlite3 *db;
 7 vector< vector< string > > databaseResults;
 8
 9 string deadEndDelim1 = "xxx"; //formerly "DEADBEEF"
10 string deadEndDelim2 = "fff"; //formerly "DEADerBEEF"
11
12 static int callback(void *queryterm, int nCol, char **values, char **headers){
13     int i;
14     vector<string> rowEntry;
15     //fprintf(stderr,"===Callback for query %s===\n",(char *) queryterm);
16     for(i=0; i< nCol; i++){
17         //fprintf(stderr,"%s = %s\n", headers[i], values[i] ? values[i] : "NULL");
18         rowEntry.push_back( values[i] );
19     }
20     //printf("\n");
21     databaseResults.push_back(rowEntry);
22     return 0;
23 }
24
25 vector< vector<phone> > getPhoneSeqsForSampaStrs( vector<string> sampaStrings ) {
26     vector< vector <phone> > sampaSyllPhrases;
27     for(int i = 0; i < sampaStrings.size(); i++) {
28         vector<phone> sampaSylls = parseSAMPAintoPhonemes( sampaStrings[i] );
29         sampaSyllPhrases.push_back(sampaSylls);
30      }
31     return sampaSyllPhrases;
32 }
33
34 vector< vector<phone> > getPhoneSeqsForOrthoWord( string orthoWord ) {
35     vector<string> sampaStrings = queryDBwithOrthoForSampaStrs( orthoWord );
36     return getPhoneSeqsForSampaStrs( sampaStrings );
37 }
38
```

```cpp
39  /*given an phrase of ortho words, gives all the sampa permutations
40   that it could possibly be, where each inner vector is a phonetic permutation of
41   the full phrase */
42  vector< vector<phone> > findAllPhoneSeqsForOrthoPhrase( string orthoPhrase ) {
43
44      vector<string> orthoWords = strTokOnWhitespace( orthoPhrase );
45      vector< vector<phone> > fullPhrasePhoneSeqs;
46      //cerr << "FIND ALL PERMUTATIONS" << endl;
47      for (int i = 0; i < orthoWords.size(); i++) {
48          vector<string> sampaStrings = queryDBwithOrthoForSampaStrs( orthoWords[i] );
49          vector< vector<phone> > nextWordSAMPAPhoneSeqs = getPhoneSeqsForOrthoWord( orthoWords[i] );
50          /*
51          //DEBUG
52          cerr << j <<": ";
53          for ( int k = 0; k < nextWordSAMPAPhoneSeqs[j].size(); k++ ) {
54              cerr<< "_" << nextWordSAMPAPhoneSeqs[j][k] << "_";
55          }
56          cerr << endl;
57          //END DEBUG
58          */
59          if( nextWordSAMPAPhoneSeqs.empty() ) {
60              cout<< "The word '"<<orthoWords[i]<<"was not found in our phonetic dictionary";
61              cout<< "Enter a different word now, or watch this program crash and burn: ";
62              string temp;
63              cin >> temp;
64
65              nextWordSAMPAPhoneSeqs = getPhoneSeqsForOrthoWord( temp );
66          }
67
68          //if this is ths first orthoWord
69          if( i == 0 ) {
70              for( int j = 0; j < nextWordSAMPAPhoneSeqs.size(); j++ ) {
71                  fullPhrasePhoneSeqs.push_back( nextWordSAMPAPhoneSeqs[j] );
72              }
73          } else {
74              int numFullPhrases = fullPhrasePhoneSeqs.size();
75              //cerr << "\tnumFullPhrases = "<< numFullPhrases << endl;//TODO debug
76              if ( nextWordSAMPAPhoneSeqs.size() > 1 ) {
```

```cpp
 77            //cerr << "\tnextWordSAMPAPhoneSeqs.size() = "<< nextWordSAMPAPhoneSeqs.size() << endl
 78            for(int m = 1; m < nextWordSAMPAPhoneSeqs.size(); m++) {
 79                //if there's more than one phonetic interpretation of the
 80                // ortho word to be added, then we need to create duplicates
 81                // of all existing sampaPhrase entries for each of them.
 82                for( int n = 0; n < numFullPhrases; n++){
 83                    vector< phone > copyOfFullPhraseN( fullPhrasePhoneSeqs[n] );
 84                    fullPhrasePhoneSeqs.push_back( copyOfFullPhraseN );
 85                }
 86            }
 87        }
 88        for( int m = 0; m < fullPhrasePhoneSeqs.size(); m++){
 89            int phrsToAppendNdx = m / numFullPhrases;
 90            vector<phone> phraseToAppend( nextWordSAMPAPhoneSeqs[phrsToAppendNdx] );
 91            fullPhrasePhoneSeqs[m].insert( fullPhrasePhoneSeqs[m].end(),
 92                                    phraseToAppend.begin(),
 93                                    phraseToAppend.end() );
 94        }
 95            //DEBUG
 96            for ( int e = 0; e < fullPhrasePhoneSeqs.size(); e++ ) {
 97                cerr << e <<"***sampa phrase after append  ";
 98                for ( int f = 0; f < fullPhrasePhoneSeqs[e].size(); f++ ) {
 99                    cerr<< "_" << fullPhrasePhoneSeqs[e][f] << "_";
100                }
101                cerr << endl;
102            }
103            //END DEBUG
104        }
105        /*
106        //DEBUG
107        cerr << j <<"++SAMPA+PHRASES++  ";
108        for ( int k = 0; k < fullPhrasePhoneSeqs[j].size(); k++ ) {
109            cerr<< "-" << fullPhrasePhoneSeqs[j][k] << "-";
110        }
111        cerr << endl;
112        //END DEBUG
113    */
114
```

```
115        }
116            // assert(0);
117
118 /*
119     vector<string> misheard;
120     for (int i = 0; i < fullPhrasePhoneSeqs.size(); i++){
121         //misheard.push_back( interpretPhrase( fullPhrasePhoneSeqs[i] ) )
122     }
123
124     for (int i = 0; i < misheard.size(); i++) {
125         string s = misheard[i];
126         DDDDDDDDDDDEBUG(s);
127     }
128
129
130     return misheard;
131 */
132
133
134
135     return fullPhrasePhoneSeqs;
136 }
137
138
139 /*given an phrase of ortho words, gives all the sampa permutations
140  that it could possibly be, where each inner vector represents all the
141  possible phonetic interpretations for each phoneme. For example:
142  Given: a nice
143  Returns:   vector[0]: { e, @, A }
144             vector[1]: { n }
145             vector[2]: { aI, i }
146             vector[3]: { s }
147              */
148  vector< set<phone> > findPhoneTreeForOrthoPhrase( string orthoPhrase ) {
149     vector< set<phone> > phoneTree;
150     cerr <<"findPhoneTreeForOrthoPhrase is broken! ("<<orthoPhrase<<")"<<endl;
151     assert(0);
152     //using set because it doesn't allow for duplicates
```

```cpp
153     vector<string> orthoWords = strTokOnWhitespace( orthoPhrase );
154     vector< vector<phone> > oldFullPhrasePhoneSeqs = findAllPhoneSeqsForOrthoPhrase( orthoPhrase );
155     vector< vector<phone> > fullPhrasePhoneSeqs;
156
157     if( oldFullPhrasePhoneSeqs.size() > 0 )   {
158         for( int i = 0; i < oldFullPhrasePhoneSeqs.size(); i++) {
159             //strip all the phoneSeqs of emph values
160             fullPhrasePhoneSeqs.push_back( getNoEmphsPhoneVect( oldFullPhrasePhoneSeqs.at(i) ) );
161         }
162
163         for(int i = 0; i < fullPhrasePhoneSeqs[0].size(); i++) {
164
165             set<phone> dummySet;
166
167             for( int j = 0; j < fullPhrasePhoneSeqs.size(); j++) {
168                 //I assume that all phone seqs will be equal length. If not, assert.
169                 if( fullPhrasePhoneSeqs[j].size() != fullPhrasePhoneSeqs[0].size() ){
170                     cerr << "0-size="<< fullPhrasePhoneSeqs[0].size()<< endl;
171                     cerr << j<<"-size="<< fullPhrasePhoneSeqs[j].size()<< endl;
172                     cerr << j<<"="<< phoneVectToString( fullPhrasePhoneSeqs[j] )<< endl;
173
174                     assert(0);
175                 }
176                 //put the i-th phone from each fullPhrasePhoneSeq into the i-th phoneTree set.
177                 int deleteme1 = dummySet.size();
178                 phone toAdd = fullPhrasePhoneSeqs[j][i];
179                 cerr << "dummySet.size()="<<deleteme1<<",  toAdd="<<toAdd<<endl;
180                 dummySet.insert( toAdd );
181             }
182             phoneTree.push_back( dummySet );
183         }
184     }
185
186
187
188
189     //DEBUG
190     for ( int j = 0; j < phoneTree.size(); j++) {
```

```cpp
191            cerr << j <<"++PHONE+TREE++  ";
192            vector<phone> temp( phoneTree[j].begin(), phoneTree[j].end() );
193
194            for ( int k = 0; k < temp.size(); k++ ) {
195
196                cerr<< "-" << temp.at(k) << "-";
197            }
198            cerr << endl;
199        }
200        //END DEBUG
201
202        return phoneTree;
203 }
204
205
206 /*given an orthoPhrase, returns all possible orthoPhrases it could be misheard as*/
207 vector<string> discoverOronymsForPhrase( string origOrthoPhrase ) {
208     vector<string> orthoMisheardAsPhrases;
209     vector<vector<phone> > allPhoneSeqsOfOrigPhrase = findAllPhoneSeqsForOrthoPhrase( origOrthoPhrase
210
211     int numUniquePhoneticInterpretations = allPhoneSeqsOfOrigPhrase.size();
212     for(int i = 0; i < numUniquePhoneticInterpretations; i++) {
213        vector<phone> curPhoneSeqWithEmph( allPhoneSeqsOfOrigPhrase.at(i) );
214        string strOfCurPhoneSeq = phoneVectToString( curPhoneSeqWithEmph );
215
216        cerr << "Phonetic interpretation "<<i<<" ("<< strOfCurPhoneSeq <<")"<<endl;
217
218        //remove emphasis marking for easier lookups
219        vector<phone> curPhoneSeq = getNoEmphsPhoneVect( curPhoneSeqWithEmph );
220
221        //vector<string> altOrthoPhrases = interpretPhrase( curPhoneSeq );
222        vector<string> altOrthoPhrases = findOrthoStrsForPhoneSeq( curPhoneSeq );
223
224        cerr << "exits findOrthoStrsForPhoneSeq"<<endl;
225        for( int j = 0; j < altOrthoPhrases.size(); j++) {
226            string altOrthoPhrase = altOrthoPhrases.at(j);
227            cerr << i <<"~~>" << altOrthoPhrase << endl;
228
```

```cpp
229            //ensure it contains no deadEndDelims so we only add fully valid strings
230            if ( altOrthoPhrase.find( deadEndDelim1 ) == string::npos
231                && altOrthoPhrase.find( deadEndDelim2 ) == string::npos ) {
232                orthoMisheardAsPhrases.push_back( altOrthoPhrase );
233            }
234        }
235    }
236    //deduplicate orthoMisheardAsPhrases by putting in a set and back again
237    cerr << "DEDUPLICATION TIME!" <<endl;
238    set<string> tempSetForDeduplication( orthoMisheardAsPhrases.begin(), orthoMisheardAsPhrases.end()
239    orthoMisheardAsPhrases.assign( tempSetForDeduplication.begin(), tempSetForDeduplication.end() );
240    return orthoMisheardAsPhrases;
241 }
242
243 /*This function does the phoneme-tree-traversal thing for oronyms
244    returns orthographic phrases (I *think* each string is a full phrase...)*/
245 vector<string> interpretPhrase( vector<phone> sampaPhraseOrig ) {
246    vector<phone> sampaPhrase = getNoEmphsPhoneVect(sampaPhraseOrig);
247    vector<string> misheardOrthoPhrases;
248    assert(0);
249    /*
250    cerr << "INTERPRET PHRASE for " << phoneVectToString(sampaPhrase) << endl;
251    if( sampaPhrase.size() == 0 ) {
252        misheardOrthoPhrases.push_back("");
253        cerr << "<<<<<<phraseSize == 0, so returning all of the phrases" <<endl;
254        return misheardOrthoPhrases;
255    }
256
257    string sampaStr = "";
258    vector <phone> usedPhones;
259    cerr << "sampaPhrase.size() ="<<sampaPhrase.size()<<endl;
260    for (int i = 0; i < sampaPhrase.size(); i++) {
261        cerr << "i = "<<i<<"; sampaPhrase[i]= phone p = "<<sampaPhrase[i]<<endl;
262        phone p = sampaPhrase[i];
263        if( strcmp( "\"", p.c_str() ) == 0) {
264            assert(0);
265            continue; //TODO incorporate someday, but ignore emphases for now.
266        } else if ( strcmp( "$", p.c_str() ) == 0) {
```

```
267            assert(0);
268            continue; //TODO incorporate someday, but ignore emphases for now.
269        } else if ( strcmp( "%", p.c_str() ) == 0) {
270            assert(0);
271            continue; //TODO incorporate someday, but ignore emphases for now
272        }
273        sampaStr += p;
274        cerr<< "Sampastr = "<<sampaStr<<endl;
275        usedPhones.push_back(p);
276        vector<string> orthoMatches = queryDBwithSampaForOrthoStrs( sampaStr );
277        cerr << "orthoMatches.size() =="<<orthoMatches.size()<<endl;
278        //DEBUG
279        for(int o = 0; o < orthoMatches.size(); o++) {
280            cerr<<"++"<<orthoMatches.at(o);
281        }
282        cerr<<endl;
283        //END DEBUG
284        //if there are no exact matches
285        if ( orthoMatches.size() == 0 ) {
286            vector<string> prefixMatches = queryDBForOrthoStrsWithSampaPrefix( sampaStr );
287            //if there are no partial matches, we have a dead end, so exit
288            if( prefixMatches.size() == 0 ) {
289                misheardOrthoPhrases.push_back( deadEndDelim1 );
290                //TODO might have to delete rest of phone seq? we'll see.
291                continue;
292            } else {
293                continue; //go to next loop iter and add next phone
294            }
295            //return misheardOrthoPhrases;
296            cerr <<" OLD RETURN STATEMENT WAS HERE for if no exact matches"<< endl;
297        }
298
299        for (int j = 0; j < orthoMatches.size(); j++) {
300        cerr << "enter orthomatches loop"<<endl;
301            string orthoWord = orthoMatches[j];
302            cerr << "----"<<i<<"--orthoword--"<< orthoMatches[j] << endl;
303            vector<phone> sampaPhraseTail( sampaPhrase.begin()+j+1, sampaPhrase.end() );
304            cerr << "----"<<i<<"--sampaPhraseTail--"<< phoneVectToString(sampaPhraseTail) <<"|--"<< endl
```

```cpp
305
306            vector<string> orthoLeaves = interpretPhrase ( sampaPhraseTail );
307            if ( orthoLeaves.size() == 0 ) {
308                if( sampaPhraseTail.size() > 0 ) {
309                    cerr<< "--"<<orthoWord<<"---no leaves, has tail: "<< phoneVectToString(sampaPhraseTai
310                    //TODO RESTART TRACE AT NEXT LINE!perhaps want a continue?
311                    misheardOrthoPhrases.push_back( orthoWord.append(  deadEndDelim1  ) );
312                }
313                //return misheardOrthoPhrases;
314                cerr <<" OLD RETURN STATEMENT WAS HERE for if no ortholeaves"<< endl;
315                continue;
316            } else {
317
318                for (int k = 0; k < orthoLeaves.size(); k++) {
319                    string orthoLeaf = orthoLeaves[k];
320                    misheardOrthoPhrases.push_back( orthoWord + orthoLeaf );
321                }
322
323            }
324        }
325    }
326
327    cerr<<"EXITING interpretPhrase"<<endl;
328    return misheardOrthoPhrases;
329    */
330 }
331
332 vector<string> findOrthoStrsForPhoneSeq( vector<phone> phoneSeq ) {
333    //cerr<<"+++findOrthoStrsForPhoneSeq, for "<< phoneVectToString( phoneSeq );
334    //cerr<<",  size = "<<phoneSeq.size()<<endl;
335    vector<phone> usedPhonesForOrtho;
336
337    vector<phone> curPhoneSeq;
338
339    vector< string > fullOrthoStrs;
340    if( phoneSeq.size() == 0 ) {
341        //cerr<<"+++ PHONE SEQ SIZE = 0, we're SUCCEEEDED! WOOOHOOO!"<<endl;
342        fullOrthoStrs.push_back("___SUCCESS!___");
```

```cpp
343        return fullOrthoStrs;
344    }
345    for ( int i = 0; i < phoneSeq.size(); i++) {
346        phone p = phoneSeq.at(i);
347        curPhoneSeq.push_back(p);
348        string curPhoneSeqStr = phoneVectToString( curPhoneSeq );
349
350        //cerr<<"+++"<<"+++p"<<i<<":"<<p<<"  full subseq = "<<curPhoneSeqStr<<"."<<endl;
351
352        /////STEP 1: EXACT MATCHES
353        //Query for exact ortho matches of the curPhoneSeq
354        vector<string> orthoInterps = queryDBwithSampaForOrthoStrs(curPhoneSeqStr);
355
356        //cerr<<"+++"<<"+++orthoInterps.size() = "<<orthoInterps.size()<<endl;
357
358        //If there is one or more exact ortho match for the phoneSeq
359        if( orthoInterps.size() > 0 ) {
360
361            //Since we're using the phones in curPhoneSeq for our ortho matches,
362            //we don't want to re-look-up those phonemes. the line below
363            // removes all the phones we used in curPhoneSeq from the phoneSeq to
364            // make newPhoneSeqTail;
365            vector<phone> newPhoneSeqTail( phoneSeq.begin() + i + 1 , phoneSeq.end() );
366            //cerr<<"+++++++++newPhoneSeqTail is "<< phoneVectToString( newPhoneSeqTail )<<endl;
367
368            //findOrthoStrings for the tail phonemes
369            vector<string> tailOrthoStrs = findOrthoStrsForPhoneSeq(newPhoneSeqTail);
370
371            //for each orthoInterpretation of the curPhoneSeq,
372            for(  int j = 0; j < orthoInterps.size(); j++ ) {
373                for ( int k = 0; k < tailOrthoStrs.size(); k++ ) {
374                    string headPlusTailOrtho = orthoInterps.at(j) + " "+ tailOrthoStrs.at(k);
375                    fullOrthoStrs.push_back( headPlusTailOrtho );
376                    cerr<<"+++"<<"+++"<<"+++"<<headPlusTailOrtho<<endl;
377                }
378            }
379        } else if ( i == phoneSeq.size() - 1 ) {
380            //then there are no phonemes left after this one.
```

```
381            //it would be stupid to check for partials if there's nothing to append
382            //so we DEADBEEF THAT SHIT
383
384            fullOrthoStrs.push_back( deadEndDelim1 );
385
386
387        } else {
388            //STEP 2: PARTIAL PREFIX MATCHES
389
390            // query for ortho matches that have the current phoneSeq as a prefix
391            vector<string> orthoPartials = queryDBForOrthoStrsWithSampaPrefix(curPhoneSeqStr);
392
393            //if there are partial matches
394            if( orthoPartials.size() > 0 ) {
395                continue;
396            } else {
397                // there are no partial matches even.
398                // How should I denote thsi?
399                fullOrthoStrs.push_back( deadEndDelim2 );
400                break;
401            }
402        }
403    }
404    return fullOrthoStrs;
405 }
406
407 vector<string> queryDBforStrings( char* sqlQuery, string queryCallback4thArg ) {
408     char* zErr;
409     /*The following line calls the callback function, passing its 4th arg as the
410      first param of the callback function.  The sqlite3_exec function
411      queries the database, then for every result that it gets, it calls the
412      callback function.*/
413     int rc = sqlite3_exec(db, sqlQuery, callback, (void*)queryCallback4thArg.c_str(), &zErr);
414     if ( rc != SQLITE_OK ) {
415         if ( zErr != NULL ) {
416             fprintf(stderr, "SQL error: %s\n", zErr);
417             sqlite3_free(zErr);
418         }
```