

MLSALT5: Speech and Language Processing Applications

Practical: Automatic MT Evaluation

FEYNMAN LIANG

1. 1-BEST KEYWORD SPOTTING

1.1. Word based KWS.

1.1.1. *Implementation description.* We implement a word-based KWS system in the `KWSIndex` class (Listing 2). Our KWS system is constructed in the `KWSIndex#apply` method, which takes a CTM file, reads each line, and parses it into a `CTMEntry` (Listing 1). Each `CTMEntry` contains the same information present in a line in a CTM file (e.g. kw file, start time, duration, token). In addition, we add some additional information such as the previous token and end times (if applicable) to each `CTMEntry` during parsing in order to allow us to identify `CTMEntry`s which are consecutive within the 1-best output.

These `CTMEntry`s are then added to an inverted index for quick lookup by token. We represent the inverted index using a hashmap which is keyed on the entry's token. We chose to use a hashmap because it provides insertions, updates, deletions, and lookups all in $O(1)$ amortized time and has amortized space complexity $O(N)$ where N is the number of entries.

At query time, the `KWSIndex.kws` method is called with a path to `queries.xml`. Each query string is passed to `KWSIndex.get`, where it is split on whitespaces to yield a sequence of individual words (i.e. tokens), and each token is looked up in the inverted index. This yields a sequence of collections of `CTMEntry`s, each collection of `CTMEntry`s representing the possible hits for that position in the query string.

We can view this as a lattice where the time axis is the position within the query string and the collection of hits at each position nodes. Since the hits for a query sequence should be consecutive in time, we need to find paths through the lattice such that each node within the path appears consecutively within the CTM file. This is where the `prevToken` and `prevEndTime` fields of the `CTMEntry` class come in. Using this information, we can filter all possible paths through this lattice to just those whose sequences of `CTMEntry`s appear consecutively within the CTM file. During this construction of paths, we also enforce the 0.5 second rule (i.e. the end time of each word in a path must be within 0.5 seconds of the next word's start time). For each path, we set the score to be equal to the average of all the CTM entry scores in the path because this yielded identical performance to the reference results provided in `/usr/local/teach/MLSALT5/Practical/scoring`.

After finding all valid paths, we are left with a collection of hits for the query sequence. These hits are used to instantiate a `QueryResults` class, whose `QueryResults.toXML` method is used to write the results in the required KWS output format.

1.1.2. *Results for word-based KWS.* We ran our word-based KWS system on `reference.ctm` and `decode.ctm` and ensured that our results (Table 2.1) matched the provided results in `/usr/local/teach/MLSALT5/Practical/scoring`.

System	IV	OOV	All
Word (<code>reference.ctm</code>)	1.000	1.000	1.000
Word (<code>decode.ctm</code>)	0.402	0.000	0.320

TABLE 1. TWV for word-based KWS on 1-best reference

Table 2.1 shows that our KWS system achieves 1.0 TWV on `reference.ctm`. This is expected because we built the KWS index on the same reference we scored against.

When we run the system on the 1-best decoding output of a word-based ASR system (`decode.ctm`), we obtain significantly worse results. These errors are due to errors from the ASR system causing incorrect

hits to be added to the KWS index as well as other entries to be omitted, resulting in worse performance. Also, notice that our out-of-vocabulary performance is 0.000, meaning that we were not able to find any hits for any out of vocabulary words. Since a word based system is keyed on words in its vocabulary, any OOV words will result in a miss. Hence, this 0.000 TWV on OOV words is expected not only for this system, but for any word-based KWS system.

1.2. Morphological Decomposition. The motivation for using morph decomposition is to better handle OOV queries. Although a particular query word may not be in the vocabulary, by decomposing both the query as well as the keys in the KWS inverted index into morphs, there is a chance that the decomposed query’s morphs will match with some decomposed entry’s morphs and produce a hit.

1.2.1. Implementation details. The `MorphDecompose` class (Listing 4) takes a morph dictionary for both the CTM entries (`morph.dct`) as well as the queries (`morph.kwslist.dct`) and provides the methods `MorphDecompose.decomposeQuery` and `MorphDecompose.decomposeEntry` for decomposing query sequences and `CTMEntry`s respectively into their corresponding morph sequences. When decomposing `CTMEntry`s, we break up each entry into multiple entries corresponding to the individual morphs. We split the duration of a word equally amongst its corresponding morphs. We chose score of each morph entry to be equal to the score of the original entry so that the score assigned to a hit sequence (see ??) remains the same before and after morph decomposition. In addition, we update the `prevToken` and `prevEndTime` fields to correctly reflect the decomposition of the previous consecutive entry (if applicable).

The `KWSIndexMorph` class (Listing 5) uses `MorphDecompose` to automatically decompose any queries submitted as well as the the inverted index keys (if required, e.g. if using a word-based rather than morph-based ASR decoding output) into morphs before performing the query. After decomposing the query (and the index inf applicable), `KWSIndexMorph.kws` performs an identical procedure to `KWSIndex.kws` to construct a lattice and find valid paths satisfying the 0.5 seconds rule.

Note that the morph-based system output (`decode-morph.xml`) is no longer a one-best list. For example, on

Hence, the entries in the CTM file are no longer disjoint and it does not make sense to enforce contiguity of the results for a query sequence. Hence, in `KWSIndexMorph` we remove the contiguity requirement when generating hit sequences (i.e. paths through the lattice of hits for each query sequence position).

1.2.2. Results. We first ran our morph-based KWS system on the output of a word-based ASR system (`decode.ctm`). This required `KWSIndexMorph` to perform morph decomposition of the inverted index as initially after parsing the CTM file each `CTMEntry`’s token was a word rather than a morph. We also tried our system using the output of a morph-based ASR system (`decode-morph.ctm`). In this case, morph decomposition of the inverted index was not required as the tokens for each `CTMEntry` were already morphs.

System	IV	OOV	All
Morph (<code>decode.ctm</code>)	0.378	0.016	0.304
Morph (<code>decode-morph.ctm</code>)	0.373	0.065	0.310

TABLE 2. Best possible TWVs for word-based vs morph-based systems

Table 2.1 shows our results. Notice that unlike word-based systems, the OOV performance of morph-based systems is no longer zero. This is because queries that are OOV may decompose into morphs which happen to be present in the KWS index, allowing a hit to be found. However, this only allows us to handle OOV words with IV morph decompositions, which represent a limited subset of all possible OOV words. Accordingly, while OOV performance is non-zero it is far from the performance for IV words.

Also, notice that the TWV score for IV words dropped significantly, and that despite the increase in OOV TWV our overall TWV has also dropped. To understand this, recall that TWV is computed as

$$TWV = 1 - (P_{miss} + \beta P_{FA}) \quad (1.1)$$

where $\beta = 999.9$ for Babel. Notice that both the number of true positives as well as the number of false positives affect the TWV score. Since multiple words may share a common morph decomposition and our KWS system no longer requires adjacent entries in a hit sequence appear contiguously in the CTM file, the number of hits has increased from $\#Sys=725$ to $\#Sys \in \{1230, 1308\}$ (Table 3), resulting in more true positives ($\#CorDet=405$ increases to be $\in \{420, 415\}$) but also more false alarms ($\#FA=320$

System	#Sys	#CorDet	#FA	#Miss
Word (<code>decode.ctm</code>)	725	405	320	558
Morph (<code>decode.ctm</code>)	1308	420	888	543
Morph (<code>decode-morph.ctm</code>)	1230	415	815	548

TABLE 3. Detailed performance metrics for word vs morph based KWS

increases to be $\in \{888, 815\}$). The increase in number of false alarms P_{FA} offsets any reduction in P_{miss} and hence the overall TWV decreases.

1.3. Score Normalization. In the definition of TWV (Equation 1.1), the penalty for a miss depends on the frequency of a word and is more costly for rare words (i.e. less total number of occurrences) [WM14]. This suggests that TWV may be improved by boosting the scores of rare words and suppressing the scores of frequent words.

Sum-to-one (STO) normalization is one method of achieving this effect by adjusting the score s_{ki} for the i th hit for keyword k to be

$$\hat{s}_{ki} = \frac{s_{ki}^\gamma}{\sum_{j \in \{\text{hits for } k\}} s_{kj}^\gamma} \quad (1.2)$$

$\gamma > 0$ is a tunable parameter, which in our experiments was fixed to $\gamma = 1$.

We implement STO normalization in `QueryResult` (Listing 3), where it is applied in the class constructor and used when writing the KWS query results to `xml`.

System	Threshold	IV	OOV	All
Word (<code>decode.ctm</code>)	0.167	0.402	0.000	0.320
Morph (<code>decode.ctm</code>)	0.205	0.378	0.016	0.304
Morph (<code>decode-morph.ctm</code>)	0.301	0.373	0.065	0.310
Word STO (<code>decode.ctm</code>)	0.029	0.402	0.000	0.320
Morph STO (<code>decode.ctm</code>)	0.056	0.391	0.015	0.314
Morph STO (<code>decode-morph.ctm</code>)	0.047	0.389	0.065	0.323

TABLE 4. Best possible TWVs for systems with vs without sum-to-one (STO) score normalization

Table 4 compares the effects of STO normalization on our word-based, morph-based using `decode.ctm`, and morph-based using `decode-morph.ctm` KWS systems. Notice that in all cases the best threshold decreases significantly, consistent with the observation that raw scores tend to overapproximate the true posterior probability and are better adjusted by score normalization [WM14].

In general, we see that STO usually helps improve the over all TWV. For both the morph-based KWS systems, gains of 0.010 and 0.013 in overall TWV were observed for the `decode.ctm` and `decode-morph.ctm` systems respectively. Interestingly, the performance of the word-based system did not change after applying STO.

Table 5 compares more detailed metrics of STO vs unnormalized KWS systems. In general, we see that STO leads to little difference in performance when applied to any of the systems independently.

System	#Sys	#CorDet	#FA	#Miss
Word (<code>decode.ctm</code>)	725	405	320	558
Morph (<code>decode.ctm</code>)	1308	420	888	543
Morph (<code>decode-morph.ctm</code>)	1230	415	815	548
Word STO (<code>decode.ctm</code>)	725	405	320	557
Morph STO (<code>decode.ctm</code>)	1308	420	888	543
Morph STO (<code>decode-morph.ctm</code>)	1229	415	814	548

TABLE 5. Extension of Table 3 with results after STO score normalization

2. SYSTEM COMBINATION

In this section, we investigate system combinations involving three outputs with unnormalized scores from IBM WFST software (a morph-based system “Morph” and two word-based systems “Word” and Word Sys2).

2.1. STO normalization on IBM WFST outputs. We first consider the provided outputs in isolation. Table 6 investigate how STO affects each of the individual systems. Here, the increase in TWV after STO normalization is dramatic, with overall TWV gains ranging from 0.061 to 0.160. One possible explanation for this could be that the IBM WFST software yields scores which vary significantly depending on keywords queried. Hence, the effects of STO normalization vary significantly across keywords: common hits are assigned lower scores and the scores of rarer hits are increased.

System	IV	OOV	All
Morph	0.430	0.089	0.360
Word	0.501	0.000	0.399
Word Sys2	0.507	0.000	0.403
Morph STO	0.556	0.367	0.520
Word STO	0.579	0.000	0.460
Word Sys2 STO	0.585	0.000	0.465

TABLE 6. Performance of individual IBM WFST KWS system outputs

The results in Table 6 show that the Morph STO system is the best we have considered so far, achieving 0.520 overall TWV. It’s strong performance is largely due to its strong OOV TWV of 0.367. Notice that just like in and , the word-based system achieves OOV TWV of 0.000 while the morph-based systems are able to achieve non-nil OOV TWV. Again, this is due to the fact that different (e.g. OOV) word sequences may decompose to the same morph sequence (e.g. which could happen to be IV), allowing OOV queries to be handled.

2.2. Combining IBM WFST outputs and the effects of score normalization. To combine outputs, we provide the `ResultsCombiner` class (Listing 6), who’s `ResultsCombiner#combine` method takes two `QueryResults` (possibly from different KWS systems) and returns the union of the two `QueryResults`.

One design decision was whether to apply STO before combining different outputs or over the combined output. We compare the performance of applying no STO normalization, STO just before combining, STO just after combining, and STO both before and after in Table 7. Our results show that the best performing system performed STO only before combining the outputs, but that the difference in performance is minimal. Applying STO before combining results can be justified by two factors. Firstly, scores from different KWS systems can vary so STO normalization before combination ensures combined scores are on the same scale. Secondly, we do not apply STO after combining because doing so would penalize multiple systems returning duplicate hits, which would increase the sum of the score for that keyword and hence result in all scores being scaled by a smaller factor. Due to these justifications, we apply STO only before output combination for remaining experiments.

STO	IV	OOV	All
None	0.382	0.087	0.322
After Combination	0.488	0.368	0.464
Before Combination	0.491	0.365	0.465
Both Before/After	0.489	0.358	0.462

TABLE 7. Effect of STO before and/or after combination on combined system Morph+Word+Sys2

The next item we investigated was which system combination could yield the best performance. Accordingly, we investigated all $\binom{3}{2} + \binom{3}{3}$ possible combinations, both with and without STO, in Table 8.

The overall trends in Table 8 show that STO before combination yields an improvement in both IV and OOV TWV across all combinations, but particularly for OOV performance of combinations involving

System	IV	OOV	All
Morph+Word+Sys2	0.382	0.087	0.322
Morph+Word	0.413	0.094	0.347
Morph+Sys2	0.427	0.094	0.359
Word+Sys2	0.447	0.000	0.355
Morph STO+Word STO+Sys2 STO	0.491	0.365	0.465
Morph STO+Word STO	0.530	0.363	0.496
Morph STO+Sys2 STO	0.539	0.362	0.502
Word STO+Sys2 STO	0.544	0.000	0.432

TABLE 8. Performance of combinations of IBM WFST KWS systems

morph-based systems. Combing with Table 6, we see a general trend that combining more sources tends to boost OOV performance while hurting IV performance. Intuitively, this makes sense. Adding additional outputs will increase the number of total hits, which for well-covered IV queries leads to an increase in P_{fa} (degrading IV TWV) where as for poorly-covered OOV queries tends to decrease P_{miss} (improving OOV TWV).

The best results overall were obtained from using just Morph STO without combination. This suggests that the word-based results do not add anything on top of the morph-based results when using the combination techniques we investigated.

2.3. Impact of query length. We also consider how the length of a query affects performance by segmenting the output of our best KWS system (IBM WFST Morph with STO). We wrote the `LengthMap` (??) to read `queries.xml` and generate a custom `.map` file mapping queries to their lengths. Note that since the KWS system is morph based, we defined the length of a query to be the number of morphs in a query. `scripts/termselect.sh` is then used with this `.map` file to produce TWV scores segmented by query length.

Figure 1a shows the distribution of queries over their lengths and Figure 1b shows the TWV segmented by query length. Notice that the TWV starts low at queries with a single morph. This can be explained by the large number of hits any single morph will generate, yielding a high P_{fa} which causes a low TWV (Equation 1.1). With 2 – 4 morph long queries, the morph sequence is unique enough to keep P_{fa} low without increasing P_{miss} , explaining the high TWVs of around 0.55 for these lengths. After the number of morphs in a query exceeds 4, TWV rapidly degrades. This trend can be explained by noting that a longer sequence of morphs has less probability of occurring, hence it is less likely to be covered by our KWS system. This means that P_{miss} is larger for these longer queries, leading to lower TWV.

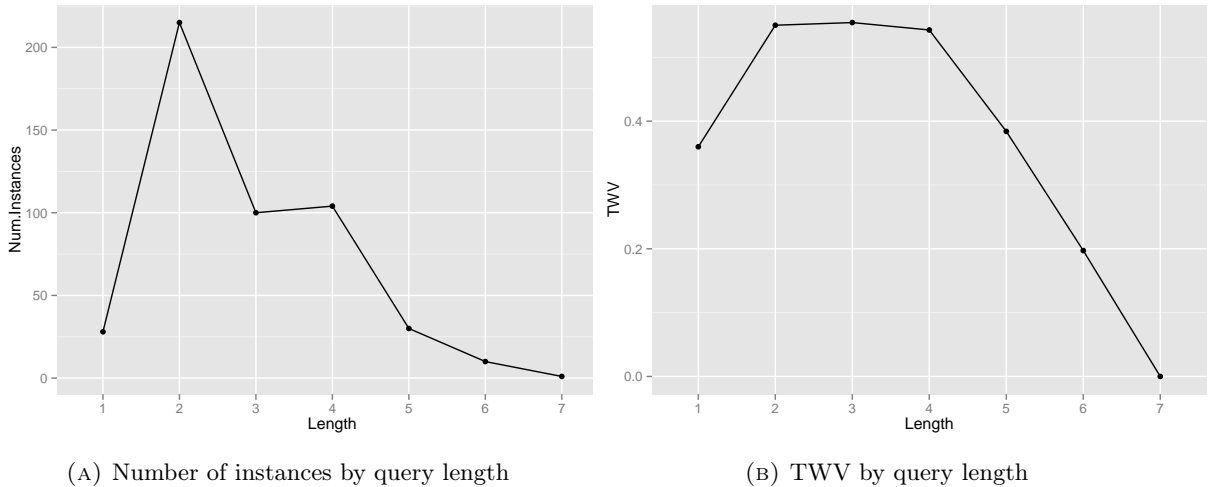


FIGURE 1. Performance of IBM Morph STO system for various query lengths (i.e. number of morphs in query)

REFERENCES

[WM14] Yun Wang and Florian Metze. An in-depth comparison of keyword specific thresholding and sum-to-one score normalization. 2014.

APPENDIX A. CODE LISTINGS

Our code is available at `/home/f1350/MLSALT5-practical/` on the MLSALT cluster. In subsection A.1 we provide listings of our implementations and in subsection A.2 we provide listings of our unit tests.

A.1. Implementations.

LISTING 1. An entry in KWS index

```
1 package com.feynmanliang.kws
2
3 import scala.xml.Elem
4
5 case class CTMEntry(
6   kwFile: String,
7   channel: Int,
8   startTime: Double,
9   duration: Double,
10  token: String,
11  prevToken: Option[String],
12  prevEndTime: Double, // for testing contiguity during phrase query
13  score: Double) extends Ordered[CTMEntry] {
14  import scala.math.Ordered.orderingToOrdered
15
16  def compare(that: CTMEntry): Int =
17    (this.kwFile, this.startTime) compare (that.kwFile, that.startTime)
18
19  def toXML(): Elem = {
20    <kw
21      file={kwFile}
22      channel={f"$channel"}
23      tbeg={f"$startTime%.2f"}
24      dur={f"$duration%.2f"}
25      score={f"$score%.6f"}
26      decision="YES" />
27  }
28 }
29
30 object CTMEntry {
31  def fromXMLNode(e: scala.xml.Node): CTMEntry = CTMEntry(
32    kwFile=(e \ "@file").text,
33    channel=(e \ "@channel").text.toInt,
34    startTime=(e \ "@tbeg").text.toDouble,
35    duration=(e \ "@dur").text.toDouble,
36    token="", // TODO: allow initializing token from XML (OK for just combining)
37    prevToken=None, // TODO: allow initializing token from XML (OK for just combining)
38    prevEndTime=0D, // TODO: allow initializing token from XML (OK for just combining)
39    score=(e \ "@score").text.toDouble)
40 }
41
42 // vim: set ts=2 sw=2 et sts=2:
```

LISTING 2. Word-based KWS system

```
1 package com.feynmanliang.kws
2
3 import java.io.{BufferedWriter, File, FileWriter}
4 import scala.io.Source
5 import scala.xml.Elem
6
7 class KWSIndex(
8   val index: Map[String, Set[CTMEntry]],
```

```

9   val scoreNormalization: String = "NONE" {
10  def get(tokens: String): Option[Set[CTMEntry]] = {
11    val res = tokens.split("\\s+")
12    .map(_._toLowerCase)
13    .map(index.get)
14    if (res.exists(_._isEmpty)) None
15    else Some(res
16      .map(_._get.map(entry => entry.copy(score=(1.0/res.size) * entry._score)))
17      .reduceLeft { (acc, x) =>
18        (for {
19          prevEntry <- acc;
20          entry <- x;
21          if (
22            entry.kwFile == prevEntry.kwFile // phrases must belong to same file
23            && prevEntry.startTime < entry.startTime && entry.startTime < (prevEntry.startTime + prevEntry.duration)
24              + 0.5
25            && prevEntry.startTime + prevEntry.duration == entry.prevEndTime) // TODO: generalize to morphs?
26          } yield {
27            prevEntry.copy(
28              duration = entry.startTime + entry.duration - prevEntry.startTime,
29              token = prevEntry.token ++ " " ++ entry.token,
30              score = prevEntry.score + entry.score
31            )
32          }).toSet
33    })
34  }
35  def kws(queryFilePath: String): QueryResult = {
36    val queryFile = scala.xml.XML.loadFile(queryFilePath)
37    val results = (queryFile \ "kw").map { kw =>
38      this._get((kw \ "kwtext").text) match {
39        case None => (kw \ "@kwid").text -> Set()
40        case Some(hits) => (kw \ "@kwid").text -> hits
41      }
42    }.toMap.asInstanceOf[Map[String, Set[CTMEntry]]]
43
44    new QueryResult(queryFilePath.split("/").last, results, scoreNormalization)
45  }
46
47
48  def write(ctmPath: String): Unit = {
49    def entry2line(entry: CTMEntry): String = {
50      s"${entry.kwFile} ${entry.channel} ${entry.startTime} ${entry.duration} " +
51      s"${entry.token} ${entry.score}"
52    }
53    val file = new File(ctmPath)
54    val bw = new BufferedWriter(new FileWriter(file))
55    index.values.toList.flatten.sorted.foreach { entry =>
56      bw.write(entry2line(entry))
57    }
58    bw.close()
59  }
60 }
61
62 object KWSIndex {
63   def apply(ctmPath: String, scoreNorm: String = "NONE"): KWSIndex = {
64     def line2entry(line: String, prevToken: Option[String], prevEndTime: Double): CTMEntry = {
65       val items = line.split("\\s+")
66       val startTime = items(2).toDouble
67       val entry = CTMEntry(
68         kwFile = items(0),
69         channel = items(1).toInt,
70         startTime = startTime,
71         duration = items(3).toDouble,
72         token = items(4).toLowerCase,
73         prevToken = prevToken,
74         prevEndTime = prevEndTime,
75         score = items(5).toDouble)
76       entry
77     }

```

```

78 def line2endTime(line: String): Double = {
79     val items = line.split(" ")
80     items(2).toDouble + items(3).toDouble
81 }
82 val ctmFile = Source.fromFile(ctmPath)
83 val index = ctmFile.getLines().map { line =>
84     line2entry(line, None, 0D)
85 }.foldLeft(List.empty[CTMEntry]) { case (acc, entry) =>
86     acc match {
87         case Nil =>
88             entry.copy(
89                 prevToken = None,
90                 prevEndTime = 0D) :: acc
91         case prevEntry :: _ =>
92             if (prevEntry.kwFile != entry.kwFile) {
93                 entry.copy(
94                     prevToken = None,
95                     prevEndTime = 0D) :: acc
96             } else {
97                 entry.copy(
98                     prevToken = if (prevEntry.startTime >= entry.startTime) {
99                         prevEntry.prevToken
100                     } else {
101                         Some(prevEntry.token)
102                     },
103                     prevEndTime = prevEntry.startTime + prevEntry.duration) :: acc
104             }
105     }
106 }.foldLeft(Map.empty[String, Set[CTMEntry]]) { (acc, entry) =>
107     acc + (entry.token -> (acc.getOrElse(entry.token, Set.empty[CTMEntry]) + (entry)))
108 }
109 new KWSIndex(index, scoreNorm)
110 }
111
112 def main(args: Array[String]): Unit = {
113     case class Config(
114         ctmFile: File = new File("."),
115         queryFile: File = new File("."),
116         scoreNorm: String = "NONE",
117         morphDecompose: Boolean = false,
118         out: File = new File("."))
119
120     val parser = new scopt.OptionParser[Config]("KWSIndex") {
121         head("kwsindex")
122         opt[File]('c', "ctmFile") required() valueName("<file>") action { (x, c) =>
123             c.copy(ctmFile = x) } text("ctmFile is a required file property")
124         opt[File]('q', "queryFile") required() valueName("<file>") action { (x, c) =>
125             c.copy(queryFile = x) } text("queryFile is a required file property")
126         opt[String]('s', "scoreNorm") required() valueName("<string>") action { (x, c) =>
127             c.copy(scoreNorm = x) } text("scoreNorm is a required string property")
128         opt[File]('o', "out") required() valueName("<file>") action { (x, c) =>
129             c.copy(out = x) } text("out is a required file property")
130     }
131     // parser.parse returns Option[C]
132     parser.parse(args, Config()) match {
133         case Some(config) =>
134             // do stuff
135             val index = KWSIndex(config.ctmFile.getPath(), config.scoreNorm)
136             val queryResults = index.kws(config.queryFile.getPath())
137             scala.xml.XML.save(config.out.getPath(), queryResults.toXML())
138         case None =>
139             // arguments are bad, error message will have been displayed
140             sys.error("Error parsing arguments!")
141     }
142 }
143 }
144
145 // vim: set ts=2 sw=2 et sts=2:

```


LISTING 3. Result of a KWS query

```

1 package com.feynmanliang.kws
2
3 import scala.xml.Elem
4
5 case class QueryResult(
6   val file : String,
7   private[kws] val results : Map[String, Set[CTMEntry]],
8   scoreNormalization : String = "NONE") {
9
10  val resultsNorm = scoreNormalization match {
11    case "NONE" => results
12    case "STO" => results.map { case (k,v) =>
13      val sumScore = v.map(_._score).sum
14      k -> v.map(entry => entry.copy(
15        score = entry.score / sumScore
16      ))
17    }
18    case _ => sys.error(s"Unknown score normalization, got: ${scoreNormalization}")
19  }
20
21  def toXML(): Elem = {
22    <kwslst
23      kwlist_filename="IARPA-babel202b-v1.0d_conv-dev.kwlist.xml"
24      language="swahili"
25      system_id="">
26      {for (kw <- resultsNorm.keys) yield {
27        <detected_kwlist kwid={s"$kw"} oov_count="0" search_time="0.0">
28          {for {
29            entry <- resultsNorm(kw)
30            // if (entry.score >= 0.042)
31          } yield {
32            entry.toXML()
33          }}
34        </detected_kwlist>
35      }}
36    </kwslst>
37  }
38
39 }
40
41 object QueryResult {
42   def fromXML(e: Elem, scoreNorm: String = "NONE"): QueryResult = {
43     QueryResult(
44       (e \ "@kwlist_filename").text,
45       (e \ "detected_kwlist").map { kwlist =>
46         (kwlist \ "@kwid").text -> ((kwlist \ "kw").map { kw =>
47           CTMEntry.fromXMLNode(kw)
48         }).toSet
49       }.toMap,
50       scoreNorm)
51   }
52 }
53
54 // vim: set ts=2 sw=2 et sts=2:

```

LISTING 4. Performs morphological decomposition using morph dicts

```

1 package com.feynmanliang.kws
2
3 import scala.io.Source
4
5 class MorphDecompose private (
6   qDict: Map[String, List[String]], obDict: Map[String, List[String]]) {
7   def decomposeQuery(tokens: Iterable[String]): List[List[String]] = {
8     tokens.map(decomposeQuery).toList
9   }
10  def decomposeQuery(token: String): List[String] = {

```

```

11     qDict.getOrElse(token, List(token))
12 }
13
14 /** Decompose a CTM entry into a list of morphemes
15     We distribute startTime/durations and scores uniformly across all morphemes.
16 */
17 def decomposeEntry(entry: CTMEntry): List[CTMEntry] = {
18     obDict.get(entry.token) match {
19         case None => List(entry)
20         case Some(morphs) => {
21             val morphDur = (entry.duration / morphs.size)
22             if (morphs.size < 2) {
23                 List(entry)
24             } else {
25                 morphs.sliding(2).zipWithIndex.flatMap { case (morphPair, i) =>
26                     if (i == 0) {
27                         List(
28                             entry.copy(
29                                 duration = morphDur,
30                                 token = morphPair(0),
31                                 prevToken = entry.prevToken.flatMap(obDict.get).map(_._last)
32                             ),
33                             entry.copy(
34                                 startTime = entry.startTime + (i+1)*morphDur,
35                                 duration = morphDur,
36                                 token = morphPair(1),
37                                 prevToken = Some(morphPair(0)),
38                                 prevEndTime = entry.startTime + (i+1)*morphDur
39                             )
40                         )
41                     } else {
42                         List(
43                             entry.copy(
44                                 startTime = entry.startTime + (i+1)*morphDur,
45                                 duration = morphDur,
46                                 token = morphPair(1),
47                                 prevToken = Some(morphPair(0)),
48                                 prevEndTime = entry.startTime + (i+1)*morphDur
49                             )
50                         )
51                     }
52                 }.toList
53             }
54         }
55     }
56
57 object MorphDecompose {
58     private def parseMorphDict(dictPath: String): Map[String, List[String]] = {
59         Source.fromFile(dictPath).getLines().map { line =>
60             val words = line.split("\\s+")
61             words.head -> words.tail.toList
62         }.toMap
63     }
64
65     def apply(qDictPath: String, obDictPath: String) = {
66         val qDict = parseMorphDict(qDictPath)
67         val obDict = parseMorphDict(obDictPath)
68         new MorphDecompose(qDict, obDict)
69     }
70 }
71
72 // vim: set ts=2 sw=2 et sts=2:

```

LISTING 5. Morph-based KWS system which automatically applies morph decomposition to the query as well as to index entries if required

```

1 package com.feynmanliang.kws
2

```

```

3 import java.io.File
4
5 class KWSIndexMorph(
6   index: Map[String, Set[CTMEntry]],
7   md: MorphDecompose,
8   scoreNorm: String = "NONE") extends KWSIndex(index, scoreNorm) {
9
10  val morphIndex = if (index.values.reduce(_+_).exists(entry => md.decomposeEntry(entry).size > 1)) {
11    // Apply entry decomposition (for decode.ctm)
12    println("Decomposing index => morphIndex")
13    index.values
14      .flatMap(_.flatMap(md.decomposeEntry))
15      .foldLeft(Map[String, Set[CTMEntry]]()) { (acc, x) =>
16        acc + (x.token -> (acc.getOrElse(x.token, Set[CTMEntry]()) + x))
17      }
18  } else {
19    // No need to decompose (for decode-morph.ctm)
20    println("index already built over morphs, setting morphIndex <- index")
21    index
22  }
23
24  override def get(tokens: String): Option[Set[CTMEntry]] = {
25    val hitsPerMorph = md.decomposeQuery(tokens.split("\\s+").map(_.toLowerCase))
26      .flatten // treat morphs same as words
27      .map(morphIndex.get)
28    if (hitsPerMorph.exists(hits => hits.isEmpty)) None
29    else Some(
30      hitsPerMorph
31        .map(hits =>
32          hits.get.map { entry =>
33            entry.copy(score=(1.0/hitsPerMorph.size) * entry.score)
34          }
35        )
36      .reduceLeft { (acc, x) =>
37        (for {
38          prevEntry <- acc;
39          entry <- x;
40          if (
41            prevEntry.kwFile == entry.kwFile
42            && prevEntry.startTime < entry.startTime && entry.startTime < (prevEntry.startTime + prevEntry.duration)
43              + 0.5)
44            //&& (entry.prevToken.isEmpty || (entry.prevToken.get == prevEntry.last.token))
45            //&& prevEntry.startTime + prevEntry.duration == entry.prevEndTime
46          } yield {
47            entry.copy(
48              startTime = prevEntry.startTime,
49              duration = entry.startTime + entry.duration - prevEntry.startTime,
50              token = prevEntry.token ++ " " ++ entry.token,
51              prevEndTime = prevEntry.prevEndTime,
52              score = prevEntry.score + entry.score
53            )
54          }).toSet
55      })
56  }
57
58  object KWSIndexMorph {
59    def apply(
60      ctmPath: String, obDictPath: String, qDictPath: String, scoreNorm: String = "NONE"): KWSIndexMorph = {
61      val index = KWSIndex(ctmPath)
62      val md = MorphDecompose(qDictPath, obDictPath)
63
64      new KWSIndexMorph(index.index, md, scoreNorm)
65    }
66
67    def main(args: Array[String]): Unit = {
68      case class Config(
69        ctmFile: File = new File("."),
70        queryFile: File = new File("."),
71        scoreNorm: String = "NONE",

```

```

72     dict: File = new File("."),
73     kwDict: File = new File("."),
74     morphDecompose: Boolean = false,
75     out: File = new File(".")
76
77     val parser = new scopt.OptionParser[Config]("KWSIndex") {
78       head("kwsindex")
79       opt[File]('c', "ctmFile") required() valueName("<file>") action { (x, c) =>
80         c.copy(ctmFile = x) } text("ctmFile is a required file property")
81       opt[File]('q', "queryFile") required() valueName("<file>") action { (x, c) =>
82         c.copy(queryFile = x) } text("queryFile is a required file property")
83       opt[File]('d', "dict") required() valueName("<file>") action { (x, c) =>
84         c.copy(dict = x) } text("dict is a required file property")
85       opt[File]('k', "kwDict") required() valueName("<file>") action { (x, c) =>
86         c.copy(kwDict = x) } text("kwDict is a required file property")
87       opt[String]('s', "scoreNorm") required() valueName("<string>") action { (x, c) =>
88         c.copy(scoreNorm = x) } text("scoreNorm is a required string property")
89       opt[File]('o', "out") required() valueName("<file>") action { (x, c) =>
90         c.copy(out = x) } text("out is a required file property")
91     }
92     parser.parse(args, Config()) match {
93       case Some(config) =>
94         val indexMorph = KWSIndexMorph(
95           ctmPath = config.ctmFile.getPath(),
96           obDictPath = config.dict.getPath(),
97           qDictPath = config.kwDict.getPath(),
98           scoreNorm = config.scoreNorm)
99         val queryResults = indexMorph.kws(config.queryFile.getPath())
100         scala.xml.XML.save(config.out.getPath(), queryResults.toXML())
101       case None =>
102         // arguments are bad, error message will have been displayed
103         sys.error("Error parsing arguments!")
104     }
105   }
106 }
107
108
109 // vim: set ts=2 sw=2 et sts=2:

```

LISTING 6. Performs system combination by running multiple KWS systems in parallel and unioning the `QueryResults`

```

1 package com.feynmanliang.kws
2
3 import java.io.File
4 import scala.xml.Elem
5
6 object ResultCombiner {
7   def combine(qr1: QueryResult, qr2: QueryResult): QueryResult = qr1.copy(
8     scoreNormalization = "NONE", // do not renormalize resultsNorm since it is reused in accumulator
9     results = (qr1.results.keys ++ qr2.results.keys).map {
10       k => k -> (qr1.resultsNorm.getOrElse(k, Set())) ++ qr2.resultsNorm.getOrElse(k, Set())
11     }.toMap)
12
13   def main(args: Array[String]): Unit = {
14     case class Config(
15       postingLists: Seq[File] = Seq(),
16       scoreNorms: Seq[String] = Seq(),
17       finalScoreNorm: String = "NONE",
18       out: File = new File(".")
19     )
20
21     val parser = new scopt.OptionParser[Config]("ResultCombiner") {
22       head("resultCombiner")
23       opt[Seq[File]]('p', "postingLists") required() valueName("<postingList1>,...") action { (x, c) =>
24         c.copy(postingLists = x) } text("posting lists to combine")
25       opt[Seq[String]]('n', "scoreNorms") required() valueName("<scoreNorm1>,...") action { (x, c) =>
26         c.copy(scoreNorms = x) } text("score normalizations to apply before combining (same order as posting lists)")
27       opt[String]('f', "finalScoreNorm") required() valueName("<finalScoreNorm>") action { (x, c) =>

```

```

27     c.copy(finalScoreNorm = x) } text("score normalizations to apply to combined output")
28     opt[File]('o', "out") required() valueName("<file>") action { (x, c) =>
29     c.copy(out = x) } text("file to output to")
30 }
31 parser.parse(args, Config()) match {
32     case Some(config) =>
33         require(config.postingLists.size == config.scoreNorms.size)
34         val combinedQr = config.postingLists.zip(config.scoreNorms)
35             .map { case (pl, sn) =>
36             QueryResult.fromXML(scala.xml.XML.load(pl.getPath()), sn)
37             }
38             .reduceLeft(ResultCombiner.combine)
39         scala.xml.XML.save(
40             config.out.getPath(),
41             combinedQr.copy(scoreNormalization = config.finalScoreNorm).toXML())
42     case None =>
43         // arguments are bad, error message will have been displayed
44         sys.error("Error parsing arguments!")
45 }
46 }
47 }
48
49 // vim: set ts=2 sw=2 et sts=2:

```

LISTING 7. Generates a `map` file mapping keyword IDs to their lengths

```

1 package com.feynmanliang.kws
2
3 import java.io.{BufferedWriter, File, FileWriter}
4
5 object LengthMap {
6     def main(args: Array[String]): Unit = {
7         case class Config(
8             queryFile: File = new File("."),
9             lengthType: String = "WORD",
10            out: File = new File("."))
11
12         val parser = new scopt.OptionParser[Config]("LengthMap") {
13             head("lengthmap")
14             opt[File]('q', "queryFile") required() valueName("<file>") action { (x, c) =>
15             c.copy(queryFile = x) } text("queryFile is a required file property")
16             opt[String]('l', "lengthType") required() valueName("<WORD|MORPH>") action { (x, c) =>
17             c.copy(lengthType = x) } text("lengthType is a required file property")
18             opt[File]('o', "out") required() valueName("<file>") action { (x, c) =>
19             c.copy(out = x) } text("out is a required file property")
20         }
21         // parser.parse returns Option[C]
22         parser.parse(args, Config()) match {
23             case Some(config) =>
24                 val md = MorphDecompose("lib/dicts/morph.kwslst.dct", "lib/dicts/morph.dct")
25
26                 val file = new File(config.out.getPath())
27                 val bw = new BufferedWriter(new FileWriter(file))
28                 val queryFile = scala.xml.XML.loadFile(config.queryFile.getPath())
29                 (queryFile \ "kw").map { kw =>
30                     val kwld = (kw \ "@kwld").text.drop(6)
31                     val length = config.lengthType match {
32                         case "WORD" => (kw \ "kwtext").text.split("\\s+").size
33                         case "MORPH" => (kw \ "kwtext").text.split("\\s+").flatMap(md.decomposeQuery).size
34                     }
35                     s"${length} ${kwld}\n"
36                 }
37                 .foreach(bw.write)
38                 bw.close()
39
40             case None =>
41                 // arguments are bad, error message will have been displayed
42                 sys.error("Error parsing arguments!")

```

```

43     }
44 }
45 }
46
47 // vim: set ts=2 sw=2 et sts=2:

```

A.2. Tests.

```

1 package com.feynmanliang.kws
2
3 import org.scalatest.FlatSpec
4
5 class KWSIndexSpec extends FlatSpec {
6   val ctmPath = "lib/ctms/reference.ctm"
7   val queryFilePath = "lib/kws/queries.xml"
8   val index = KWSIndex(ctmPath)
9   val queryResults = index.kws(queryFilePath)
10
11   "A KWSIndex" should "contain words known to be in the CTM" in {
12     assert(!index.get("halo").isEmpty)
13   }
14
15   it should "be case insensitive" in {
16     assert(!index.get("rachael").isEmpty)
17     assert(!index.get("kisumu").isEmpty)
18     assert(index.get("Wenga fans").size === 1)
19   }
20
21   it should "only contain phrase queries when words are <0.5 sec apart" in {
22     val tokens = "what she has gone"
23     assert(index.get(tokens).get.size === 1)
24     assert(index.get("pat pat pat").get.size === 1)
25   }
26
27   it should "perform KWS when given a queryFile" in {
28     val ctmPath = "lib/ctms/reference.ctm"
29     val index = KWSIndex(ctmPath)
30
31     val queryFilePath = "lib/kws/queries.xml"
32     val queryResults = index.kws(queryFilePath)
33
34     assert(queryResults.file === "queries.xml")
35     assert(!(queryResults.toXML() \ "detected_kwlist").isEmpty)
36   }
37
38   it should "perform the same as provided scoring/decode" in {
39     val ctmPath = "lib/ctms/decode.ctm"
40     val index = KWSIndex(ctmPath)
41
42     assert(index.get("that's why am").isEmpty)
43     assert(index.get("kae tukae kae").isEmpty)
44   }
45 }
46
47 // vim: set ts=2 sw=2 et sts=2:

```

```

1 package com.feynmanliang.kws
2
3 import org.scalatest.FlatSpec
4
5 class QueryResultSpec extends FlatSpec {
6   val ctmPath = "lib/ctms/reference.ctm"
7   val queryFilePath = "lib/kws/queries.xml"
8   val index = KWSIndex(ctmPath)

```

```

9   val queryResults = index.kws(queryFilePath)
10
11  // TODO: "NONE" vs "STO" score normalization tests
12  "Multiple QueryResults" should "be able to be read from XML" in {
13    val qr = QueryResult.fromXML(scala.xml.XML.loadFile("lib/kws/morph.xml"))
14    assert(qr.results.get("KW202-00001").isDefined)
15  }
16 }
17
18 // vim: set ts=2 sw=2 et sts=2:

```

```

1 package com.feynmanliang.kws
2
3 import org.scalatest.FlatSpec
4
5 class MorphDecomposeSpec extends FlatSpec {
6   val qPath = "lib/dicts/morph.kwslist.dct"
7   val obPath = "lib/dicts/morph.dct"
8   val md = MorphDecompose(qPath, obPath)
9
10  "A MorphDecompose" should "decompose a query string" in {
11    assert(md.decomposeQuery("walini") === List("wali", "ni"))
12    assert(md.decomposeQuery("very vile".split("\\s+")) === List(List("ve", "ry"), List("vi", "le")))
13  }
14
15  it should "return the query unchanged if no entry in morph dict" in {
16    assert(md.decomposeQuery("abcdefg hijkl".split("\\s+")) === List(List("abcdefg"), List("hijkl")))
17  }
18
19  it should "decompose a CTMEntry" in {
20    val entry = CTMEntry("", 1, 0.5, 0.4, "vile", None, 0.33, 0.09)
21    assert(md.decomposeEntry(entry) == List(
22      CTMEntry("", 1, 0.5, 0.2, "vi", None, 0.33, 0.09),
23      CTMEntry("", 1, 0.7, 0.2, "le", Some("vi"), 0.7, 0.09)
24    ))
25  }
26
27  it should "return the CTMEntry unchanged if no entry in morph dict" in {
28    val entry = CTMEntry("", 1, 0.5, 0.4, "abcdefg", None, 0.33, 0.09)
29    assert(md.decomposeEntry(entry) == List(entry))
30  }
31 }
32
33 // vim: set ts=2 sw=2 et sts=2:

```

```

1 package com.feynmanliang.kws
2
3 import org.scalatest.FlatSpec
4
5 class KWSIndexMorphSpec extends FlatSpec {
6   val ctmPath = "lib/ctms/decode-morph.ctm"
7   val queryFilePath = "lib/kws/queries.xml"
8   val obPath = "lib/dicts/morph.dct"
9   val qPath = "lib/dicts/morph.kwslist.dct"
10
11   val indexMorph = KWSIndexMorph(ctmPath, obPath, qPath)
12
13   "A KWSIndexMorph" should "be queryable" in {
14     val queryResults = indexMorph.kws(queryFilePath)
15     assert((queryResults.toXML() \ "detected_kwlist").length > 100)
16   }
17
18   it should "return the same results as the reference" in {
19     assert(indexMorph.get("hangesikia").isEmpty)
20   }
21 }

```

22

23

24 // vim: set ts=2 sw=2 et sts=2:
