

NIST 2017 Language Recognition Evaluation Plan

September 29, 2017

1 Introduction

The 2017 NIST language recognition evaluation (LRE17) is the 8th cycle in an on-going language recognition evaluation series that began in 1996. The objectives of the evaluation series are (1) to explore promising new ideas in language recognition, (2) to support the development of advanced technology incorporating these ideas, and (3) to measure the performance of the current state of technology. Targeting researchers working on the general problem of text-independent, speaker-independent language recognition, the evaluation is designed to focus on core technology issues and to be simple and accessible to those wishing to participate.

LRE17 will be organized in a similar manner to LRE15, focusing on differentiating closely related languages. Again *fixed* and *open* training conditions will be offered to allow cross-system comparisons and to understand the effect of additional and unconstrained amount of training data on system performance.

There are several differences between LRE17 and LRE15. In addition to conversational telephone speech (CTS) and broadcast narrow band speech (BNBS), speech extracted from videos or video speech (VS) will be used as test material. The test segments from CTS and BNBS will be extracted from longer recordings to create smaller chunks containing approximately 3s, 10s, or 30s of speech, as was in older LREs. The test segments from VS will use the entire recording. NIST will distribute to participants a small development set containing CTS, BNBS, and VS. Participants will be asked to provide score vectors representing the log-likelihood scores, rather than log-likelihood ratios. The primary metric will support equal weighting of data sources and durations.

Participation in LRE17 is open to all who find the evaluation of interest and are able to comply with the evaluation rules set forth in this plan. There is no cost to participate, but participating teams must be represented at the evaluation workshop planned for December 2017 (location TBD). Information about evaluation registration can be found on the LRE17 website¹.

2 Task Description

2.1 Task Definition

The task for LRE17 is *language detection*: given a segment of speech and a target language, automatically determine if the target language was spoken in the test segment. LRE17 has 14 target languages grouped into five language clusters as listed in Table 1.

Input to a language recognition (LR) system will be a series of test segments, and output from the system will be a series of score vectors, one vector per test segment. Each score vector is defined as a 14-dimensional vector corresponding to the 14 target languages in the order listed in Section 6.4, and representing estimated log-likelihood scores, using natural (base e) logarithms, for the corresponding languages. In terms of the conditional probabilities for the observed data (O) given a target language model (L_i), the log-likelihood score (ℓ_i) is defined as

¹<https://www.nist.gov/itl/iad/mig/nist-2017-language-recognition-evaluation>

$$\ell_i = \log(P(O|L_i)). \quad (1)$$

The likelihood function in (1) is related to the posterior probability $P(L_i|O)$ via Bayes' rule as follows

$$P(L_i|O) = \frac{P(L_i) \exp(\ell_i)}{\sum_{j=1}^{N_L} P(L_j) \exp(\ell_j)}, \quad (2)$$

where $P(L_i)$ is the *a priori* probability of the language class i , and N_L is the number of target languages.

Language Cluster	Target Languages	Language Code
Arabic	Egyptian Arabic, Iraqi Arabic, Levantine Arabic, Maghrebi Arabic	ara-arz, ara-acm, ara-apc, ara-ary
Chinese	Mandarin, Min Nan	zho-cmn, zho-nan
English	British English, General American English	eng-gbr, eng-usg
Slavic	Polish, Russian	qsl-pol, qsl-rus
Iberian	Caribbean Spanish, European Spanish, Latin American Continental Spanish, Brazilian Portuguese	spa-car, spa-eur, spa-lac, por-brz

Table 1: LRE17 target languages and language clusters

2.2 Training Conditions

The training condition is defined as the amount of data/resources used to build an LR system. The task described above can be evaluated over a *fixed* (required) or *open* (optional) training condition.

- **Fixed** – The *fixed* training condition limits the system training to the following specific data sets:
 - previous LRE data (as released in LDC2017E22)
 - Fisher corpus (LDC2004S13, LDC2004T19, LDC2005S13, LDC2005T19)
 - Switchboard corpora (LDC97S62, LDC2001S13, LDC2002S06, LDC2004S07, LDC98S75, LDC99S79)
 - LRE17 “dev” set (LDC2017E23)

The Linguistic Data Consortium (LDC) data license agreement lists the actual catalog numbers for these corpora. Participants can obtain the data from the LDC after they have signed the data license agreement. For the *fixed* training condition, only the specified speech data may be used for system training and development, including all system modules (e.g., speech activity detection) and auxiliary systems employed to build an LR system (e.g., automatic speech recognition). Publicly available non-speech audio and data (e.g., noise samples, impulse responses, filters) may be used and should be noted in the system description. Participation in the *fixed* condition is required.

Note: The use of pretrained models on data other than what is designated above (e.g., BUT Hungarian phoneme recognizer) is not allowed in this condition.

- **Open** – The *open* training condition removes the limitations of the *fixed* condition. In addition to the data listed in the *fixed* condition, participants can use any additional data including proprietary data and data that are not publicly available. The inclusion of non-publicly available data is new for LRE17. Please note that any additional data used must be adequately described by providing enough details in the system description.

LDC will also make available selected data from the IARPA Babel Program to be used in the *open* training condition. Participation in this condition is optional but strongly encouraged to demonstrate the gains that can be achieved with unconstrained amounts of data.

3 Performance Measurement

3.1 Primary Metric

Pair-wise LR performance will be computed for all target-language/non-target-language pairs (L_T, L_N). This will be done in terms of false-reject (missed detection) and false alarm (FA) probabilities, which will be computed separately for each target language and each target/non-target language pair, respectively. The miss and false alarm probabilities will then be combined using a linear cost function according to an application-motivated cost model, defined as

$$C(L_T, L_N) = C_{Miss} \times P_{Target} \times P_{Miss}(L_T) + C_{FA} \times (1 - P_{Target}) \times P_{FA}(L_T, L_N), \quad (3)$$

where L_T and L_N are target and non-target languages, respectively. Here, C_{Miss} (cost of a missed detection), C_{FA} (cost of a spurious detection), and P_{Target} (*a priori* probability of the specified target language) are the application model parameters and defined to have the following values:

Parameter ID	C_{Miss}	C_{FA}	P_{Target}
1	1	1	0.5
2	1	1	0.1

Table 2: LRE17 cost parameters

Note that the first set of parameter values are those historically used in NIST LREs and provide equal weighting to miss and false alarm errors, while the second set of parameters are not similarly balanced. Therefore, to improve the interpretability of the cost function, it will be normalized by $C_{Default}$, which is defined as the best cost that could be obtained without processing the input data (i.e., by either always accepting or always rejecting the segment language as matching the target language, whichever gives the lower cost) as follows

$$C_{Norm}(L_T, L_N) = C(L_T, L_N) / C_{Default}, \quad (4)$$

Here, the default cost for both sets of parameters defined in Table 2 is set to $C_{Default} = C_{Miss} \times P_{Target}$. Rewriting the cost model in (3) by combining all of the application model parameters yields

$$C_{Norm}(L_T, L_N) = P_{Miss}(L_T) + \beta \times P_{FA}(L_T, L_N), \quad (5)$$

where β is defined as:

$$\beta = \frac{C_{FA} \times (1 - P_{Target})}{C_{Miss} \times P_{Target}}.$$

Actual detection costs will be computed by applying detection thresholds of $\log(\beta)$ to log-likelihood *ratios* derived from the log-likelihoods output by the system².

²Log-likelihood ratios will be computed as the difference between the target language log-likelihood and the sum of the log-likelihoods of the non-target languages, i.e., $LLR(L_i) = -\log \left[\frac{1}{N_L - 1} \sum_{j \neq i} \exp(\ell_j - \ell_i) \right]$.

In addition to the performance numbers computed for each target/non-target language pair, an average cost performance for each system will be computed as

$$C_{avg}(\beta) = \frac{1}{N_L} \left\{ \sum_{L_T} P_{Miss}(L_T) + \frac{1}{N_L - 1} \left[\beta \times \sum_{L_T} \sum_{L_N} P_{FA}(L_T, L_N) \right] \right\}, \quad (6)$$

where N_L is the number of target languages. The primary metric for LRE17 will be the average cost performance defined in (6), computed using the two application model parameters given in Table 2, that are then averaged:

$$C_{primary} = \frac{C_{avg}(\beta_1) + C_{avg}(\beta_2)}{2}. \quad (7)$$

Unlike in previous LREs, in LRE17 the evaluation data will be divided into partitions based on the data source, i.e., MLS14 and VAST, for each language, resulting in a total 28 partitions (2×14). In other words, for each language, the counts for each corpus (MLS14 and VAST) will be equalized. C_{avg} will be calculated for each partition, and the final result is the average of all the partitions' C_{avg} 's. The average of basic C_{avg} scores for the two set of parameters defined in Table 2 will serve as the primary metric to measure a system performance. Also, the minimum detection cost, $minC_{avg}$, will be computed by using the detection thresholds that minimize the detection cost. Note that for minimum cost calculations, the counts for each condition set will be equalized before pooling and cost calculation (i.e., minimum cost will be computed using a single threshold not one per condition set).

NIST will make available the script that calculates the primary metric.

3.2 Alternative Metric

In addition to the cost metric C_{avg} described above, an alternative information theoretic metric will also be used to calculate the performance of an LR system. The multiclass cross-entropy metric H_{mce} measures the information the LR system provides through the log-likelihood scores and is defined as follows³

$$H_{mce} = - \sum_{i=1}^{N_L} \frac{P(L_i)}{\|S_i\|} \sum_{t \in S_i} \log P(L_i | O_t), \quad (8)$$

where S_i is the subset of indices for segments of target language i , $\|S_i\|$ is the number of segments of target language i .

For a *do-nothing* default system, the multiclass cross-entropy is given by

$$H_{max} = - \sum_{i=1}^{N_L} P(L_i) \log P(L_i). \quad (9)$$

If $H_{mce} \geq H_{max}$ for an LR system, then it does not improve upon the default *do-nothing* system. To facilitate the interpretation of the cross-entropy or mutual information, a normalized version of H_{mce} is calculated as *confidence* score which is defined as⁴

$$Confidence = 1 - \frac{H_{mce}}{H_{max}}. \quad (10)$$

Given that the cross-entropy is non-negative, a perfect LR system achieves a confidence score of 1 (i.e., it has zero confusion), while a totally confused system can achieve a confidence score of zero (or less).

³L.J. Rodriguez-Fuentes *et al.* "The Albayzin 2012 Language Recognition Evaluation," in *Proc. INTERSPEECH*, September 2013, pp. 1497-1501.

⁴J. Fiscus *et al.* "2000 NIST evaluation of conversational speech recognition over the telephone: English and Mandarin performance results," in *Proc. Speech Transcription Workshop*, 2000.

It is worth noting that the confidence metric in (10) is being considered for use as the primary metric in future LREs (after LRE17).

4 Development and Test Data Description

The data collected by the LDC as part of the MLS14 and VAST corpora will be used to compile the LRE17 test set. A small dataset extracted from these two corpora will also be distributed for system development. Test segments from the MLS14 corpus will be 8-bit (μ -law) SPHERE files sampled at 8kHz, while recordings from the VAST corpus will be 16-bit FLAC files sampled at 44kHz.

All data will be distributed by LDC. Please refer to Section 2.2 for more information about the training data conditions and what is allowable for each condition.

4.1 Data Organization

The development and test sets follow a similar directory structure:

```
<base_directory>/
  README.txt
  data/
    dev/
    eval/
  docs/
  metadata/ (in dev set only)
```

4.2 Trial File

The trial file named `lre17_{dev|eval}_trials.tsv` and located in the `docs/` directory is composed of a header and a list of test segments:

```
segmentid<NEWLINE>
<segmentid><NEWLINE>
...
```

For example:

```
segmentid
1001.lre17
1002.lre17
1003.lre17
```

5 Evaluation Rules and Requirements

LRE17 is conducted as an open evaluation where the test data is sent to the participants who process the data locally and submit the output of their systems to NIST for scoring. As such, the participants have agreed to process the data in accordance with the following rules:

- The participants agree that for each evaluation test segment the information available to the system is limited to that segment only (along with the training data); scores for a particular test segment must be computed without benefit from any information that might be derived from other test segments.
- The participants agree not to probe the test segments via manual/human means such as listening to the data or producing the transcript of the speech during the evaluation period and before all submissions are made.

- The participants are allowed to use information available in the audio file header.

In addition to the above data processing rules, participants agree to comply with the following general requirements:

- The participants agree to follow the submission requirements. See Section 6.4.
- The participants agree to have one or more representatives at the evaluation workshop to present a meaningful description of their system(s). Evaluation participants failing to do so will be excluded from future evaluation participation.
- The participants agree to the guidelines governing the publication of the results:
 - Participants are free to publish results for their own system but must not publicly compare their results with other participants (ranking, score differences, etc.) without explicit written consent from the other participants.
 - While participants may report their own results, participants may not make advertising claims about winning the evaluation or claim NIST endorsement of their system(s). The following language in the U.S. Code of Federal Regulations (14 C.F.R. § 200.113) shall be respected⁵: *NIST does not approve, recommend, or endorse any proprietary product or proprietary material. No reference shall be made to NIST, or to reports or results furnished by NIST in any advertising or sales promotion which would indicate or imply that NIST approves, recommends, or endorses any proprietary product or proprietary material, or which has as its purpose an intent to cause directly or indirectly the advertised product to be used or purchased because of NIST test reports or results.*
 - At the conclusion of the evaluation NIST generates a report summarizing the system results for conditions of interest, but these results/charts do not contain the participant names of the systems involved. Participants may publish or otherwise disseminate these charts, unaltered and with appropriate reference to their source.
 - The report that NIST creates should not be construed or represented as endorsements for any participant's system or commercial product, or as official findings on the part of NIST or the U.S. Government.

6 Evaluation Protocol

To facilitate efficient information exchange between the participants and NIST, all evaluation activities are conducted over a web-interface.

6.1 Evaluation Account

Participants must sign up for an evaluation account where they can perform various activities such as registering for the evaluation, signing the data license agreement, uploading the submission and system description, and more. To sign up for an evaluation account, go to <https://lre.nist.gov>. The password must be at least 12 characters long and must contain a mix of upper and lowercase letters, numbers, and symbols. After the evaluation account is confirmed, the participant is asked to join a site or create one if it does not exist. The participant is also asked to associate his or her site to a team or create one if it does not exist. This allows multiple members with their individual accounts to perform activities on behalf of their site and/or team (e.g., making a submission) in addition to performing their own activities (e.g., requesting workshop invitation letter). Please note that the first person that creates the site or team is deemed the team representative. Site and team representatives have to approve participants who want to join his/her site/team.

⁵See <http://www.ecfr.gov/cgi-bin/ECFR?page=browse>

- A site is defined as a single organization (e.g., NIST).
- A team is defined as a group of organizations collaborating on a task (e.g., Team1 consisting of NIST and LDC).
- A participant is defined as a member or representative of a site who takes part in the evaluation (e.g., John Doe).

6.2 Evaluation Registration

One representative from the team⁶ must formally register his team to participate in the evaluation by agreeing to the terms of participation. For more information about the terms of participation, see Section 5.

6.3 Data License Agreement

One representative from each site must sign the LDC data license agreement to obtain the training data for the *fixed* training condition and Babel data for the *open* training condition.

6.4 Submission Requirements

Each team must participate in the *fixed* training condition. Teams are encouraged to participate in the *open* training condition to demonstrate the gains that can be achieved leveraging unconstrained amounts of data. There is no submission limit, but for each training condition participating teams must designate one submission as the *primary* submission that NIST can use for cross-team comparisons.

There should be one output file per training condition per system. Teams must process all test segments. Submission with missing test segments will not pass validation and will be rejected.

Each team is required to submit a system description at the designated time (see Section 7). The evaluation results are made available only after the system description report is received and confirmed to comply with guidelines described in Section 6.4.1.

The system output file is composed of a header and a set of records where each record contains a test segment given in the file list (see Section 4.2) and a 14-dimension vector of log-likelihood scores. The order of the test segments in the system output file must follow the same order as the file list. Each record is a single line containing 15 fields, separated by tab character, in the order listed below:

1. Segment ID<TAB>
2. ara-acm<TAB>
3. ara-apc<TAB>
4. ara-ary<TAB>
5. ara-arz<TAB>
6. eng-gbr<TAB>
7. eng-usg<TAB>
8. por-brz<TAB>
9. qsl-pol<TAB>
10. qsl-rus<TAB>

⁶Please note that the registration is done at the team level while the data license is done at the site level. If a team is registered, all sites in that team are registered. However, all sites in the team must sign the data license separately.

11. spa-car<TAB>
12. spa-eur<TAB>
13. spa-lac<TAB>
14. zho-cmn<TAB>
15. zho-nan<NEWLINE>

For example:

```
segmentid ara-acm ara-apc ara-ary ... zho-nan7
1001_lre17 -0.10017 -0.61518 -1.98380 ... -2.47851
1002_lre17 -0.15862 -0.35402 -0.04077 ... -0.96342
1003_lre17 -0.53162 -0.46526 -0.98556 ... -1.23140
```

There should be one output file for each training condition for each system. System outputs will be automatically validated through the online submission platform and a report will be generated and displayed in case there are any errors.

6.4.1 System Description Format

Each team is required to submit a system description. The system description must include the following items:

- a complete description of the system components, including front-end (e.g., speech activity detection, features, normalization) and back-end (e.g., background models, i-vector extractor, classifier) modules along with their configurations (i.e., filterbank configuration, dimensionality and type of the acoustic feature parameters, as well as the acoustic model and the backend model configurations),
- a complete description of the data partitions used to train the various models (as mentioned above). Teams are encouraged to report whether and how having access to the development set helped improve the performance,
- performance of the submission systems (primary and secondary) on the LRE17 development set, using the scoring software provided via the web platform (<https://lre.nist.gov>). Teams are encouraged to quantify the contribution of their major system components that they believe resulted in significant performance gains,
- a report of the CPU execution time (single threaded) and the amount of memory used to process a single trial (i.e., the time needed for processing a test segment to compute the score vector).

The system description should follow the latest IEEE ICASSP conference proceeding template.

7 Schedule

Milestone	Date
Evaluation plan published	May 2017
Registration period	May - September 2017
Training & development data available	June 15, 2017
Test data available to participants	September 20, 2017
System output due to NIST	October 20, 2017
Preliminary results released	November 01, 2017
Post evaluation workshop	December 12-13, 2017

⁷Note that the header is in lower case and output files without the header will not pass the validation step.