

# Towards a General and Extensible Phrase-Extraction Algorithm



Wang Ling      Tiago Luís      João Graça      Luísa Coheur      Isabel Trancoso  
L<sup>2</sup>F, INESC-ID/IST   L<sup>2</sup>F, INESC-ID/IST   L<sup>2</sup>F, INESC-ID/IST   L<sup>2</sup>F, INESC-ID/IST   L<sup>2</sup>F, INESC-ID/IST



{wang.ling,tiago.luis,joao.graca,luisa.coheur,isabel.trancoso}@l2f.inesc-id.pt

## MOTIVATION

### Phrase Extraction

- Plays a key role in enriching translation quality
- However:
  - ◇ there is not a general algorithm that allows this process to be easily extended
  - ◇ a great deal of time is spent reimplementing the same generic algorithm

### Goals

- Create a general algorithm for phrase extraction that is easily extensible
- Identify key control points that define the behavior of the algorithm
- Provide modular interface to allow users to extend these control points

1

## ALGORITHM - GENERAL PHRASE EXTRACTION

```
Require: Bilingual Corpus
Require: MaximumPhraseSize - max
for each sentence pair (s, t) in Corpus do
  extractedPhrasePairs = extractPhrasePairs(s, t, max)
  for each phrase pair p in extractedPhrasePairs do
    phraseTable.add(p)
  end for
end for
computeGlobalPhraseStats
pruneGlobalPhraseStats
savePhraseTable
```

### General Phrase Extraction

- Extract all phrase pairs from a given bilingual corpora

### Key Control Points

- ComputeGlobalPhraseStats - Merges extracted pairs and calculates their features
- PruneGlobalPhraseStats - Eliminates phrase pairs

3

## ALGORITHM - EXTRACT PHRASE PAIRS

```
Require: Bilingual sentence s
fl = s.foreignLen
sl = s.sourceLen
extractedPhrasePairs = {}
for fp = 0; fp ≤ fl; fp ++ do
  for fd = 1; fd ≤ maxDuration; fd ++ do
    if ForeignPhraseAcceptor.accept(s, fp, fd) then
      for sp = 0; sp ≤ sl; sp ++ do
        for sd = 1; sd ≤ maxDuration; sd ++ do
          if SourcePhraseAcceptor.accept(s, sp, sd) then
            PhrasePair p = phrase pair from s from (fp, sp) to (fd, sd)
            LocalPhrasePairFeaturesCreator.addFeatures(p)
            if PhrasePairAcceptor.accept(p) then
              extractedPhrasePairs.add(p)
            end if
          end if
        end for
      end for
    end if
  end for
end for
return extractedPhrasePairs
```

### Extract Phrase Pairs

- Extract all phrase pairs from a given sentence pair
- Invoked in the context of the General Phrase Extraction

### Key Control Points

- SourcePhraseAcceptor - Decide if the source phrase is a good translation unit
- ForeignPhraseAcceptor - Decide if the foreign phrase is a good translation unit
- LocalPhrasePairFeaturesCreator - Calculates the features for the extracted phrase pair
- PhrasePairAcceptor - Decides whether the pair phrase is a suitable candidate for extraction

2

## IMPLEMENTED EXTENSIONS

### Moses pipeline / Baseline

- SizeBasedSourceAcceptor and SizeBasedTargetAcceptor
- KohenAcceptor
- ProbabilityFeatureCreator
- LexicalWeightingFeatureCreator

### Weighted Alignments

- WeightedAlignmentScorer
- WeightedAlignmentAcceptor

### Punctuation Based Filtering

- NoPunctuationAcceptor
- NoTerminalPunctuationAcceptor

4

## EXPERIMENTS

| Method              | BTEC (Fr-En) | DIALOG (Cn-En) |
|---------------------|--------------|----------------|
| Baseline / Moses    | 62.46        | 41.42          |
| Weighted Alignments | 63.07        | 42.15          |
| NoPunct             | 62.75        | 41.20          |
| NoTerminalPunct     | <b>63.41</b> | <b>42.28</b>   |

- Using weighted alignments and NoTerminalPunct acceptor improved the results

5

## CONCLUSIONS

- Introduced a general algorithm for phrase extraction
- Easily extensible by adding implementations of control points
- Available at <https://code.google.com/p/geppetto>

6