

FST Trimming: Ending Dictionary Redundancy in Apertium

Francis Tyers⁰ Matthew Marting¹ Kevin Unhammer²

⁰UiT Norgga árkatalaš universitehta
Romssa, Norga
ftyers@prompsit.com

¹St. David's School
Raleigh, NC.
0

²Kaldera språkteknologi
Stavanger, Noreg
unhammer+apertium@mm.st

27th May 2014

Introduction and
background

FSTs in the Apertium pipeline
The Problem: Redundant data
A Solution: Intersection

Implementation of
lt-trim

Preprocessing the bilingual
dictionary
Prefixing the bilingual
dictionary
Moving uninflected lemma
parts
Intersection
lt-trim in use

Ending Dictionary
Redundancy

Conclusion

Acknowledgements

Outline of talk

Introduction and background

Implementation of `lt-trim`

Ending Dictionary Redundancy

Conclusion

`lt-trim`

Francis Tyers,
Matthew Marting,
Kevin Unhammer

Introduction and
background

FSTs in the Apertium pipeline

The Problem: Redundant data

A Solution: Intersection

Implementation of
`lt-trim`

Preprocessing the bilingual
dictionary

Prefixing the bilingual
dictionary

Moving uninflected lemma
parts

Intersection

`lt-trim` in use

Ending Dictionary
Redundancy

Conclusion

Acknowledgements

Introduction and background

lt-trim

Francis Tyers,
Matthew Marting,
Kevin Unhammer

Introduction and background

FSTs in the Apertium pipeline
The Problem: Redundant data
A Solution: Intersection

Implementation of lt-trim

Preprocessing the bilingual dictionary
Prefixing the bilingual dictionary
Moving uninflected lemma parts
Intersection
lt-trim in use

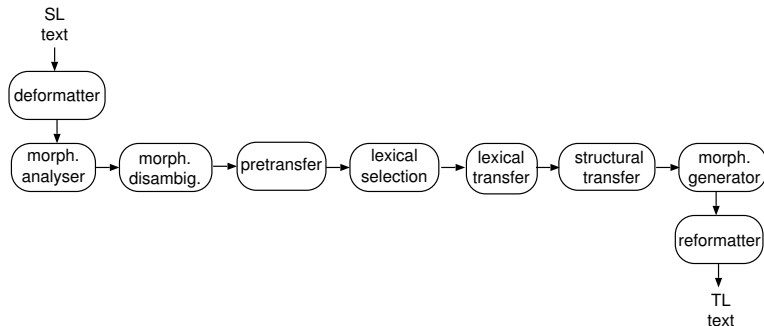
Ending Dictionary Redundancy

Conclusion

Acknowledgements

- ▶ Apertium: Free/Open Source, Rule-based Machine Translation platform
- ▶ Uses Ittoolbox Finite State Transducers for:
 - ▶ morph. analysis: 'fishes' to
`fish<n><pl>/fish<vblex><pres>`
 - ▶ lex. transfer: `fish<n><pl>` to `fisk<n><m><pl><DD>`
 - ▶ morph. generation: `fisk<n><m><pl><def>` to 'fiskane'

Apertium pipeline architecture



lt-trim

Francis Tyers,
Matthew Marting,
Kevin Unhammer

Introduction and
background

FST's in the Apertium pipeline

The Problem: Redundant data

A Solution: Intersection

Implementation of
lt-trim

Preprocessing the bilingual
dictionary

Prefixing the bilingual
dictionary

Moving uninflected lemma
parts

Intersection

lt-trim in use

Ending Dictionary
Redundancy

Conclusion

Acknowledgements

Multiword support

Ittoolbox FST's support a variety of multiwords

An Ittoolbox “lexical unit” is one token, and can be:

- ▶ simple non-multi-words: ‘fish’
- ▶ simple space-separated words: ‘hairy frogfish’ as a single token
- ▶ multiwords with **inner inflection**: ‘takes out’, analysed as `take<vblex><pri><p3><sg># out`, converted to `take# out<vblex><pri><p3><sg>` before lexical transfer

Introduction and background

FST's in the Apertium pipeline

The Problem: Redundant data

A Solution: Intersection

Implementation of lt-trim

Preprocessing the bilingual dictionary

Prefixing the bilingual dictionary

Moving uninflected lemma parts

Intersection

lt-trim in use

Ending Dictionary Redundancy

Conclusion

Acknowledgements

Multiword support

- ▶ **joined** multiwords: ‘they’ll’;

analysed as single token

`prpers<prn><subj><p3><mf><p1>+will<vaux><inf>`,

then split into two tokens

`prpers<prn><subj><p3><mf><p1>` and

`will<vaux><inf>` **before** lexical transfer

Multiword support

lt-trim

Francis Tyers,
Matthew Marting,
Kevin Unhammer

Introduction and
background

FSTs in the Apertium pipeline

The Problem: Redundant data

A Solution: Intersection

Implementation of
lt-trim

Preprocessing the bilingual
dictionary

Prefixing the bilingual
dictionary

Moving uninflected lemma
parts

Intersection

lt-trim in use

Ending Dictionary
Redundancy

Conclusion

Acknowledgements

- ▶ combinations of the 3 multiword types: 'creure-ho que', analysed as single token

creure<vblex><inf>+ho<prn><enc><p3><nt># que,

then moved and split into two tokens

creure# que<vblex><inf> and

ho<prn><enc><p3><nt> before lexical transfer

The Problem: Redundant data

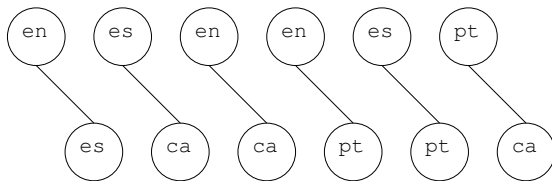


Figure: Current number of monodixes with pairs of four languages

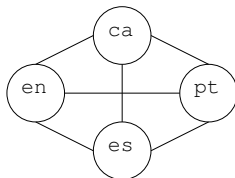


Figure: Ideal number of monodixes with four languages

A Solution: Intersection



lt-trim

Francis Tyers,
Matthew Marting,
Kevin Unhammer

Introduction and background

FSTs in the Apertium pipeline

The Problem: Redundant data

A Solution: Intersection

Implementation of lt-trim

Preprocessing the bilingual
dictionary

Prefixing the bilingual
dictionary

Moving uninflected lemma
parts

Intersection

lt-trim in use

Ending Dictionary Redundancy

Conclusion

Acknowledgements

Implementation of lt-trim

lt-trim

Francis Tyers,
Matthew Marting,
Kevin Unhammer

Introduction and background

FSTs in the Apertium pipeline

The Problem: Redundant data

A Solution: Intersection

Implementation of lt-trim

Preprocessing the bilingual
dictionary

Prefixing the bilingual
dictionary

Moving uninflected lemma
parts

Intersection

lt-trim in use

Ending Dictionary Redundancy

Conclusion

Acknowledgements



Preprocessing the bilingual dictionary

lt-trim

Francis Tyers,
Matthew Marting,
Kevin Unhammer

Introduction and background

FST's in the Apertium pipeline

The Problem: Redundant data

A Solution: Intersection

Implementation of lt-trim

Preprocessing the bilingual dictionary

Prefixing the bilingual
dictionary

Moving uninflected lemma
parts

Intersection

lt-trim in use

Ending Dictionary Redundancy

Conclusion

Acknowledgements



Prefixing the bilingual dictionary



lt-trim

Francis Tyers,
Matthew Marting,
Kevin Unhammer

Introduction and background

FST's in the Apertium pipeline

The Problem: Redundant data

A Solution: Intersection

Implementation of lt-trim

Preprocessing the bilingual
dictionary

Prefixing the bilingual dictionary

Moving uninflected lemma
parts

Intersection

lt-trim in use

Ending Dictionary Redundancy

Conclusion

Acknowledgements

Moving uninflected lemma parts

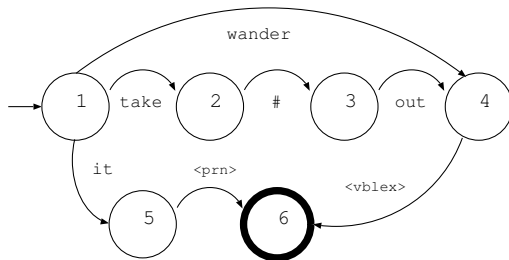
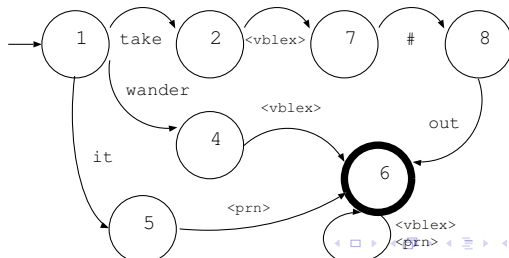


Figure: Input bilingual FST (letter transitions compressed to single arcs)



Intersection

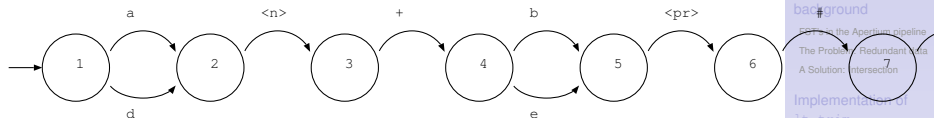


Figure: Input monolingual FST

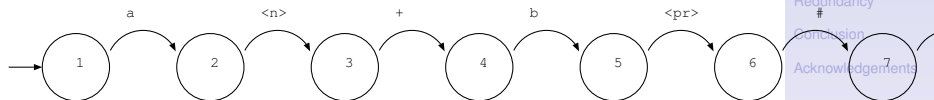


Figure: Trimmed monolingual FST

Introduction and
background

FSTs in the Apertium pipeline
The Problem: Redundant data
A Solution: Intersection

Implementation of
lt-trim

Preprocessing the bilingual
dictionary

Prefixing the bilingual
dictionary

Moving uninflected lemma
parts

Intersection

lt-trim in use

Ending Dictionary
Redundancy

Conclusion

Acknowledgements

lt-trim in use



lt-trim

Francis Tyers,
Matthew Marting,
Kevin Unhammer

Introduction and background

FST's in the Apertium pipeline

The Problem: Redundant data

A Solution: Intersection

Implementation of lt-trim

Preprocessing the bilingual
dictionary

Prefixing the bilingual
dictionary

Moving uninflected lemma
parts

Intersection

lt-trim in use

Ending Dictionary Redundancy

Conclusion

Acknowledgements

Ending Dictionary Redundancy

lt-trim

Francis Tyers,
Matthew Marting,
Kevin Unhammer

Introduction and background

FST's in the Apertium pipeline

The Problem: Redundant data

A Solution: Intersection

Implementation of lt-trim

Preprocessing the bilingual
dictionary

Prefixing the bilingual
dictionary

Moving uninflected lemma
parts

Intersection

lt-trim in use

Ending Dictionary Redundancy

Conclusion

Acknowledgements



Conclusion



lt-trim

Francis Tyers,
Matthew Marting,
Kevin Unhammer

Introduction and background

FST's in the Apertium pipeline
The Problem: Redundant data
A Solution: Intersection

Implementation of lt-trim

Preprocessing the bilingual
dictionary
Prefixing the bilingual
dictionary
Moving uninflected lemma
parts
Intersection
lt-trim in use

Ending Dictionary Redundancy

Conclusion

Acknowledgements

Acknowledgements



lt-trim

Francis Tyers,
Matthew Marting,
Kevin Unhammer

Introduction and background

FST's in the Apertium pipeline
The Problem: Redundant data
A Solution: Intersection

Implementation of lt-trim

Preprocessing the bilingual
dictionary
Prefixing the bilingual
dictionary
Moving uninflected lemma
parts
Intersection
lt-trim in use

Ending Dictionary Redundancy

Conclusion

Acknowledgements