



STATISTICAL PHRASE-BASED SPEECH TRANSLATION

Lambert Mathias and William Byrne

Center for Language and Speech Processing, Johns Hopkins University, Baltimore, MD 21218, USA
Dept. of Engineering, Cambridge University, Trumpington Street, Cambridge, CB2 1PZ, USA

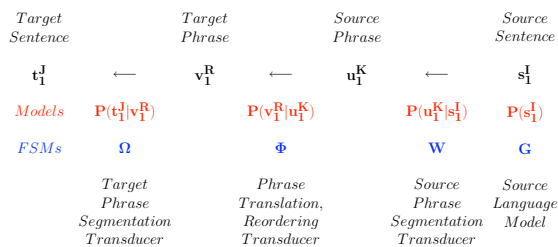


Introduction

- Objective: Tight integration of Automatic Speech Recognition (ASR) and phrase-based Statistical Machine Translation (SMT) systems.
 - ASR word lattices as input to the SMT system
 - Lattices encode a larger search space
 - Exploit sub-sentential information
- Modeling issues:
 - Propagation of ASR information to the SMT component
 - Correct disfluencies, hesitations, spontaneous speech effects...
 - Efficient phrase extraction
- Previous work:
 - (E.Matusov et al 2005) reported translation gains using word lattices
 - (Bertoldi et al 2005) used confusion networks for integration with the SMT system
- We present: Generative source-channel model of speech to text translation.
 - Tight coupling of the ASR and SMT system using word lattices
 - Implemented using weighted finite state machines (WFSMs)
- So what's new?
 - Conditional models vs joint models of target-source generation
 - Unified modeling framework
 - No need for extensive reformulation of underlying ASR and SMT models
 - Simpler decoding and estimation procedures
 - Lattice translation is a direct extension of our text translation systems

Generative Model of Text Translation

- Noisy channel model for text translation



- Translation system translates phrase sequences and not word sequences
 - Target phrase sequence acceptor $Q = \text{project_input}[\Omega \circ T]$
 - T is an acceptor for target word sequence t_1^J
 - A phrase sequence is a sequence of words to be translated
 - Different phrase sequences lead to different translations
 - Phrase sequence acceptor is unweighted

- The final translation is given by

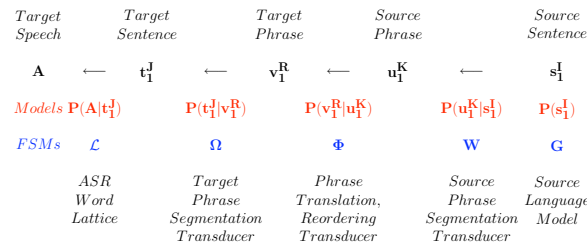
$$\hat{s}_1^I = \text{argmax}_{s_1^I} \{ \max_{v_1^R, u_1^K} P(t_1^J, v_1^R, u_1^K, s_1^I) \}$$

- Implemented as a best-path search through the translation FSM \mathcal{T}

$$\mathcal{T} = G \circ W \circ \Phi \circ Q$$

Generative Model of Speech Translation

- Noisy channel model for speech translation



- The ASR word lattice is one of the components in the noisy channel model
- The translation system translates the lattice of phrase sequences
 - Q is obtained by applying Ω to the ASR lattice \mathcal{L}
 - Unlike the text translation case, Q is a weighted acceptor with scores from the acoustic word lattice attached to it.
- The final translation is given by

$$\hat{s}_1^I = \text{argmax}_{s_1^I} \left\{ \max_{t_1^J, v_1^R, u_1^K} P(A, t_1^J, v_1^R, u_1^K, s_1^I) \right\}$$

- Implemented as a best-path search through the translation FSM \mathcal{T}

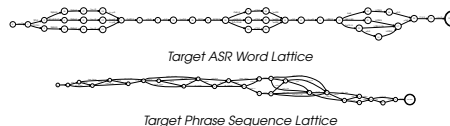
$$\mathcal{T} = G \circ W \circ \Phi \circ Q$$

Transforming ASR Word Lattices into Phrase Lattices

- Reformulate speech translation as a modeling problem
 - Efficient extraction of phrases from ASR word lattice
- ASR lattice pre-processing
 - Map unspoken tokens to NULL
 - Standard FSM operations: epsilon removal, determinization and minimization
- Controlling ambiguity in the ASR word lattice
 - Path based likelihood pruning of ASR word lattices
 - Extracting phrases under the posterior distribution

$$Q = P(v_1^R | A) = \frac{\sum_{t_1^J} P(v_1^R | t_1^J) P(A | t_1^J) P(t_1^J)}{P(A)}$$

- What about the target LM $P(t_1^J)$?
 - Shown to improve translation
 - Doesn't show up in the noisy channel model
 - Need to reformulate the SMT models to properly incorporate the target LM.
- Phrase extraction from ASR word lattice
 - Lattice subsequences of limited length
 - GRM Library tools for counting subsequences in a WFSM



Speech to Text Translation Performance

- 2005 TC-STAR Mandarin Broadcast News Translation Task

- 6 Mandarin news broadcasts, manually segmented and transcribed into Chinese sentences (525 for DEV set, 495 for EVAL set)
- 2 English reference translations per sentence for translation.
- Mandarin ASR lattices generated by LIMSI, corresponding to 231 segments for DEV set and 181 segments for EVAL set.

- Translation system: similar to 2005 JHU/CU NIST Chinese-English SMT system.

- Document level BLEU for evaluating translation performance.

- Mandarin Broadcast News Translation Performance

	Mandarin Source	DEV	EVAL
Monotone	Ref. Transcription	16.1	18.8
Phrase Order	ASR 1-Best	14.8	13.6
	ASR lattice	15.0	13.8
MJ-1 VT Phrase Reordering	Ref. Transcription	16.1	19.3
	ASR 1-Best	15.0	13.8
	ASR lattice	15.1	14.0

- How many new foreign phrases does speech translation introduce?

	DEV transcriptions	ASR Lattice
#Chinese phrases	44744	58395

- How many new foreign phrases were found in bitext?

	DEV transcriptions	ASR Lattice
#Chinese phrases	11617	12983
# English phrases	59589	60574

- Extracting phrases from ASR lattice doesn't result in a larger phrase inventory
 - Disfluencies in speech - different modeling approaches for phrase extraction
 - Mismatch in tokenization for the Chinese ASR lattices

Conclusion

- Presented a modeling framework for statistical speech-to-text translation

- Extension of the phrase-based TTM text translation model
- Tight coupling of the ASR and SMT subsystems using lattices both as
 - * Statistical models
 - * WFSM based implementation
- Demonstrated feasibility of the above approach.

- Future Work

- Initial formulation and implementation has weaknesses
 - * Proper integration of the target language model
 - * Phrase extraction under the posterior distribution
 - * Improved pruning strategies for word lattices
 - * Improved phrase coverage
- Integrated development of the component ASR and SMT systems.