
FBK @ IWSLT 2007

N. Bertoldi, M. Cettolo, R. Cattoni, M. Federico
FBK - Fondazione B. Kessler, Trento, Italy

Trento, 15 October 2007

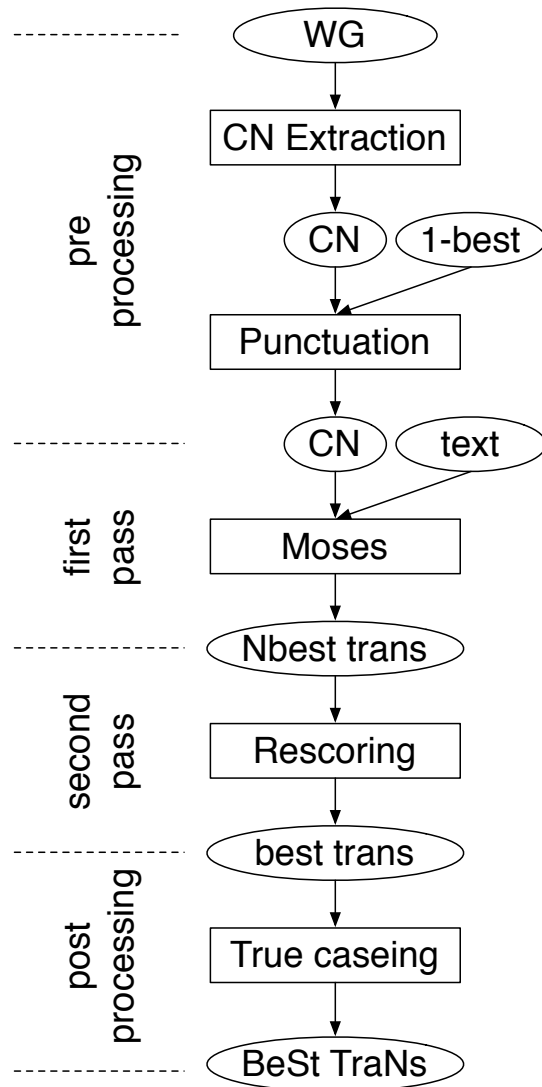
HERMES
Cross-Language Information Processing



- system architecture
- confusion network
- punctuation insertion
- improvement of lexicon
- use of multiple lexicons and language models
- system evaluation

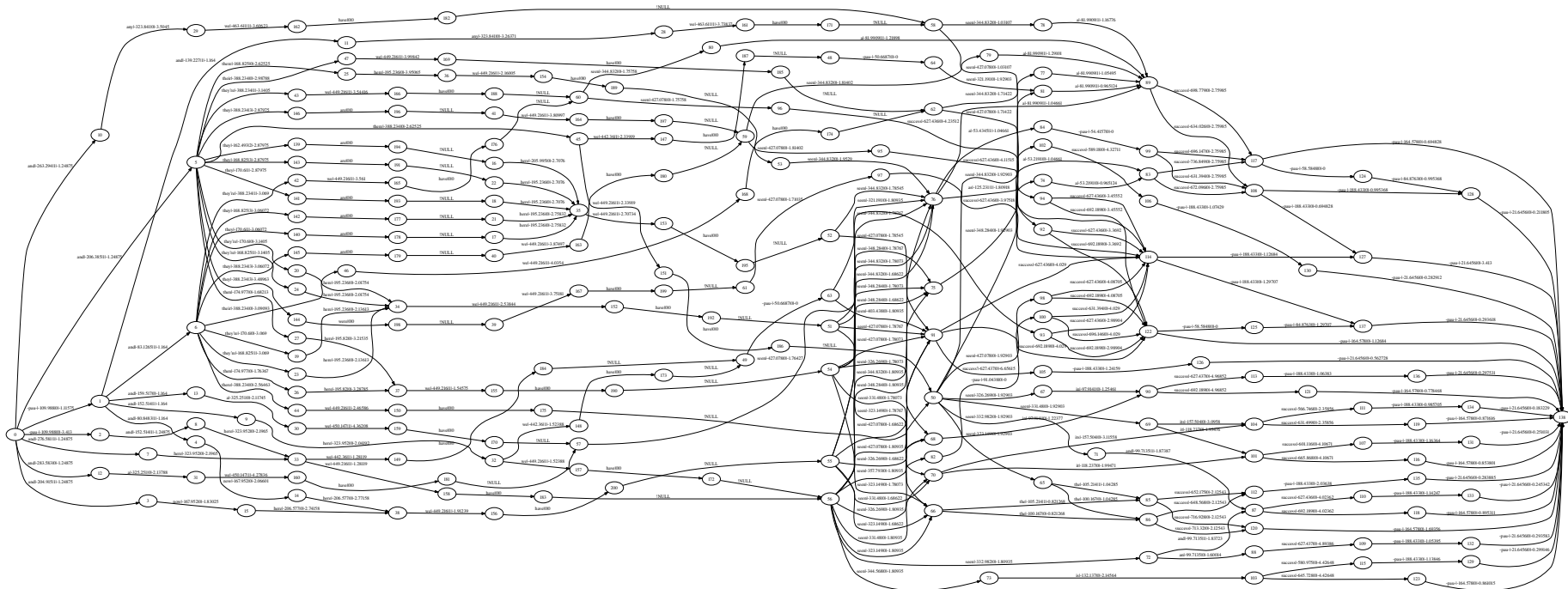
Acknowledgments

- Hermes people: Marcello, Mauro, Roldano



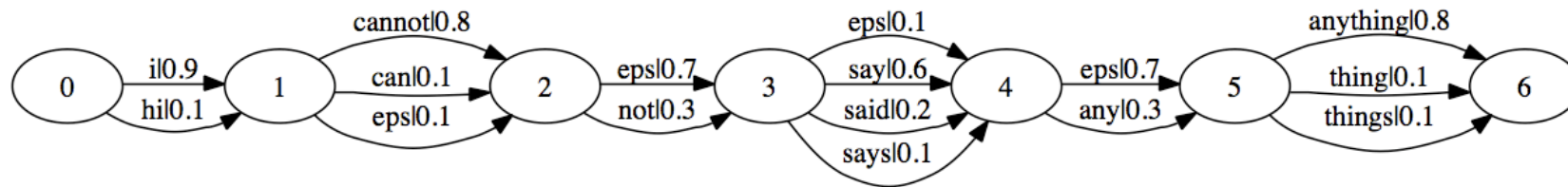
- input from speech (word-graph or 1-best) or text
- pre and post processing (optional)
 - use of the SRILM toolkit
 - **CN extraction**: lattice-tool
 - **punctuation insertion**: hidden-ngram
 - case restoring: disambig
- **Moses** is a text/CN decoder
- rescoring of N -best translations (optional)

Step 1: take the ASR *word lattice*



- arcs are labeled with *words* and *acoustic and LM scores*
- arcs have start and end *timestamps*
- any path is a *transcription hypothesis*

Step 2: approximate the word lattice into a *Confusion Network*



- a CN is a linear word graph
- arcs are labeled with *words* or with the *empty word* (ϵ -word)
- arcs are weighted with word *posterior probabilities*
- paths are a *superset* of those in the word lattice
- paths can have different lengths
- algorithm proposed by [Mangu, 2000]
 - exploit start and end timestamps of the lattice arcs
 - collapse/cluster close words
 - lattice-tool

Step 3: represent the CN as a *table*

i. ₉	cannot. ₈	€. ₇	say. ₆	€. ₇	anything. ₈
hi. ₁	can. ₁	not. ₃	said. ₂	any. ₃	thing. ₁
	€. ₁		says. ₁		things. ₁
			€. ₁		

Step 3: represent the CN as a *table*

i. ₉	cannot. ₈	€.7	say. ₆	€.7	anything. ₈
hi. ₁	can. ₁	not. ₃	said. ₂	any. ₃	thing. ₁
	€.1		says. ₁		things. ₁
			€.1		

Notes

- text is a trivial CN
- CN can be used for representing ambiguity of the input
 - transcription alternatives
 - punctuation
 - upper/lower case

The Problem

- *punctuation* improves *readability* and *comprehension* of texts
- *punctuation marks* are important clues for the translation process
- most ASR systems generate output *without* punctuation

The Problem

- *punctuation* improves *readability* and *comprehension* of texts
- *punctuation marks* are important clues for the translation process
- most ASR systems generate output *without* punctuation

Our approach [Cattoni, Interspeech 2007]

- insert punctuation as a *pre-processing* step
- exploit *multiple* hypotheses of punctuation
- use *punctuated models* (i.e. trained on texts with punctuation)
- let the decoder choose the best punctuation (and translation)

Step 1: take the input *not-punctuated CN*

i. ₉	cannot. ₈	€. ₇	say. ₆	€. ₇	anything. ₈	at. ₉	this. ₈	point. ₇	are ₁	there. ₈	€. ₈	any. ₇	comments. ₇
hi. ₁	can. ₁	not. ₃	said. ₂	any. ₃	thing. ₁	€. ₁	these. ₁	points. ₁		the. ₁	a. ₁	new. ₁	comment. ₂
	€. ₁		say. ₁		things. ₁		those. ₁	€. ₁		their. ₁	air. ₁	a. ₁	commit. ₁
			€. ₁					pint. ₁				€. ₁	

Step 2: extract the not-punctuated *consensus decoding*

i cannot say anything at this point are there any comments

Step 3: compute the *N-best* hypotheses of punctuation (with hidden-ngram)

NBEST_0	-15.270	i cannot say anything		at this point	.	are there any comments	
NBEST_1	-15.317	i cannot say anything		at this point	.	are there any comments	?
NBEST_2	-16.275	i cannot say anything		at this point		are there any comments	?
NBEST_3	-16.322	i cannot say anything		at this point	?	are there any comments	?
NBEST_4	-17.829	i cannot say anything		at this point		are there any comments	.
NBEST_5	-18.284	i cannot say anything		at this point	?	are there any comments	
NBEST_6	-18.331	i cannot say anything		at this point		are there any comments	
NBEST_7	-18.473	i cannot say anything	.	at this point		are there any comments	
NBEST_8	-18.521	i cannot say anything	.	at this point		are there any comments	?
NBEST_9	-18.834	i cannot say anything		at this point	.	are there any comments	.

Step 4: compute the *punctuating CN* with *posterior probs* of multiple marks

i ₁	cannot ₁	say ₁	anything ₁	€ ₁ .9	at ₁	this ₁	point ₁	. ₁ .7	are ₁	there ₁	any ₁	comments ₁	? ₁ .6
				. ₁ .1				€ ₁ .2					€ ₁ .3
								? ₁ .1					. ₁ .1

Step 5: *merge* the input CN and the punctuating CN

i. ₉	cannot. ₈	€. ₇	say. ₆	€. ₇	anything. ₈	at. ₉	this. ₈	point. ₇	are ₁	there. ₈	€. ₈	any. ₇	comments. ₇
hi. ₁	can. ₁	not. ₃	said. ₂	any. ₃	thing. ₁	€. ₁	these. ₁	points. ₁		the. ₁	a. ₁	new. ₁	comment. ₂
	€. ₁		say. ₁		things. ₁		those. ₁	€. ₁		their. ₁	air. ₁	a. ₁	commit. ₁
			€. ₁					pint. ₁				€. ₁	

+

i ₁	cannot ₁	say ₁	anything ₁	€. ₉	at ₁	this ₁	point ₁	. ₇	are ₁	there ₁	any ₁	comments ₁	? ₆
				. ₁				€. ₂					€. ₃
								? ₁					. ₁

Step 6: get the final *punctuated CN*

i. ₉	cannot. ₈	€.7	say. ₆	€.7	anything. ₈	€.9	at. ₉	this. ₈	point. ₇	..7	are ₁	there. ₈	€.8	any. ₇	comments. ₇	?. ₆
hi. ₁	can. ₁	not. ₃	said. ₂	any. ₃	thing. ₁	..1	€.1	these. ₁	points. ₁	€.2		the. ₁	a. ₁	new. ₁	comment. ₂	€.3
	€.1		say. ₁		things. ₁			those. ₁	€.1	?. ₁		their. ₁	air. ₁	a. ₁	commit. ₁	..1
			€.1						pint. ₁					€.1		

Step 6: get the final *punctuated CN*

i. ₉	cannot. ₈	€.7	say. ₆	€.7	anything. ₈	€.9	at. ₉	this. ₈	point. ₇	..7	are ₁	there. ₈	€.8	any. ₇	comments. ₇	?. ₆
hi. ₁	can. ₁	not. ₃	said. ₂	any. ₃	thing. ₁	..1	€.1	these. ₁	points. ₁	€.2		the. ₁	a. ₁	new. ₁	comment. ₂	€.3
	€.1		say. ₁		things. ₁			those. ₁	€.1	?. ₁		their. ₁	air. ₁	a. ₁	commit. ₁	..1
			€.1						pint. ₁					€.1		

Notes

- this approach works with any speech input (1-best and CN) without punctuation and with partially punctuated input

Step 6: get the final *punctuated CN*

i. ₉	cannot. ₈	€.7	say. ₆	€.7	anything. ₈	€.9	at. ₉	this. ₈	point. ₇	..7	are ₁	there. ₈	€.8	any. ₇	comments. ₇	?. ₆
hi. ₁	can. ₁	not. ₃	said. ₂	any. ₃	thing. ₁	..1	€.1	these. ₁	points. ₁	€.2		the. ₁	a. ₁	new. ₁	comment. ₂	€.3
	€.1		say. ₁		things. ₁			those. ₁	€.1	?. ₁		their. ₁	air. ₁	a. ₁	commit. ₁	..1
			€.1						pint. ₁					€.1		

Notes

- this approach works with any speech input (1-best and CN) without punctuation and with partially punctuated input
- one system (with punctuated models) translates any input (text and speech)

Which is the better approach to add punctuation marks?

Which is the better approach to add punctuation marks?

- in the *source* as a *pre-processing* step

Which is the better approach to add punctuation marks?

- in the *source* as a *pre-processing* step
- in the *target* as a *post-processing* step
 - translate with not-punctuated models
 - add punctuation to the best translation (with hidden-ngram)

Which is the better approach to add punctuation marks?

- in the *source* as a *pre-processing* step
- in the *target* as a *post-processing* step
 - translate with not-punctuated models
 - add punctuation to the best translation (with hidden-ngram)
- evaluation
 - task: eval set 2006, TC-STAR English-to-Spanish
 - training data: FTE transcriptions of EPPS (36Mw English, 38Mw Spanish)
 - verbatim input (w/o punctuation), case-insensitive

approach	BLEU	NIST	WER	PER
target	42,23	9.72	46.12	34.38
source	44.92	9.84	42.84	31.77

Do multiple punctuation hypotheses help to improve translation quality?

Do multiple punctuation hypotheses help to improve translation quality?

- evaluation
 - verbatim (w/o punctuation)
 - case-insensitive

input	type	# punctuation hyps	BLEU	NIST	WER	PER
vrb		1	44.92	9.84	42.84	31.77
		1000	45.33	9.83	42.58	31.59

Do multiple punctuation hypotheses help to improve translation quality?

- evaluation
 - verbatim (w/o punctuation), 1-best
 - case-insensitive

input	type	# punctuation hyps	BLEU	NIST	WER	PER
vrb		1	44.92	9.84	42.84	31.77
		1000	45.33	9.83	42.58	31.59
asr	1-best	1	35.62	8.37	57.15	44.56
		1000	36.01	8.41	56.78	44.39

Do multiple punctuation hypotheses help to improve translation quality?

- evaluation
 - verbatim (w/o punctuation), 1-best, and CN
 - case-insensitive

input	type	# punctuation hyps	BLEU	NIST	WER	PER
vrb		1	44.92	9.84	42.84	31.77
		1000	45.33	9.83	42.58	31.59
asr	1-best	1	35.62	8.37	57.15	44.56
		1000	36.01	8.41	56.78	44.39
	CN	1	36.22	8.46	56.39	44.37
		1000	36.45	8.49	56.17	44.19

Create a phrase-pair lexicon

- take a case-sensitive parallel corpus
- word-align the corpus in direct and inverse directions (GIZA++)
- combine both word-alignments in one symmetric way:
 - grow-diag-final, union, and intersection
- extract phrase pairs from a symmetrized word-alignment
- add single word translation from direct alignment
- score phrase pairs according to word and phrase frequencies

Create a phrase-pair lexicon

- take a case-sensitive parallel corpus
- word-align the corpus in direct and inverse directions (GIZA++)
- combine both word-alignments in one symmetric way:
 - grow-diag-final, union, and intersection
- extract phrase pairs from a symmetrized word-alignment
- add single word translation from direct alignment
- score phrase pairs according to word and phrase frequencies

Ideas for improving the lexicon:

- use *case-insensitive* corpus for word-alignment, but case-sensitive extraction

Create a phrase-pair lexicon

- take a case-sensitive parallel corpus
- word-align the corpus in direct and inverse directions (GIZA++)
- combine both word-alignments in one symmetric way:
 - grow-diag-final, union, and intersection
- extract phrase pairs from a symmetrized word-alignment
- add single word translation from direct alignment
- score phrase pairs according to word and phrase frequencies

Ideas for improving the lexicon:

- use *case-insensitive* corpus for word-alignment, but case-sensitive extraction
- extract phrase pairs separately from more symmetrized word-alignments, concatenate them and compute their scores

How much improvement do we get?

How much improvement do we get?

- evaluation
 - task: IWSLT Chinese-to-English, 2006 eval set
 - training data: BTEC and dev sets ('03-'05)
 - weight optimization on 2006 dev set
 - verbatim input, case-sensitive

symmetrization	text for word-alignment	# phrase pairs	BLEU	NIST
grow-diag-final	case-sensitive	496K	20.50	5.57

How much improvement do we get?

- evaluation
 - task: IWSLT Chinese-to-English, 2006 eval set
 - training data: BTEC and dev sets ('03-'05)
 - weight optimization on 2006 dev set
 - verbatim input, case-sensitive

symmetrization	text for word-alignment	# phrase pairs	BLEU	NIST
grow-diag-final	case-sensitive	496K	20.50	5.57
"	case-insensitive	507K	21.86	5.59

How much improvement do we get?

- evaluation
 - task: IWSLT Chinese-to-English, 2006 eval set
 - training data: BTEC and dev sets ('03-'05)
 - weight optimization on 2006 dev set
 - verbatim input, case-sensitive

symmetrization	text for word-alignment	# phrase pairs	BLEU	NIST
grow-diag-final	case-sensitive	496K	20.50	5.57
"	case-insensitive	507K	21.86	5.59
+union	"	507K	22.35	6.20

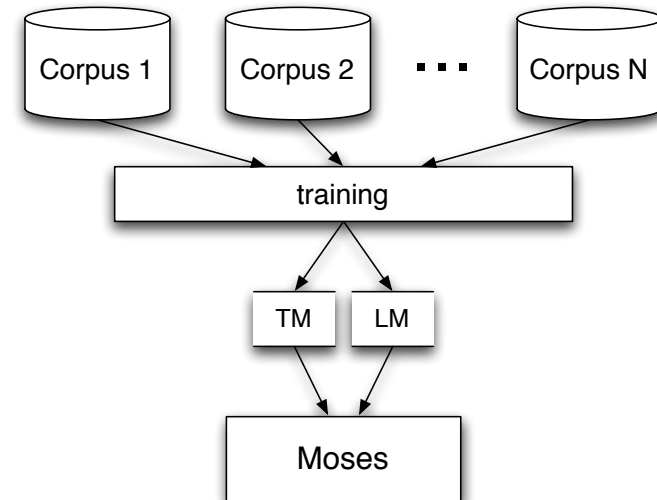
How much improvement do we get?

- evaluation
 - task: IWSLT Chinese-to-English, 2006 eval set
 - training data: BTEC and dev sets ('03-'05)
 - weight optimization on 2006 dev set
 - verbatim input, case-sensitive

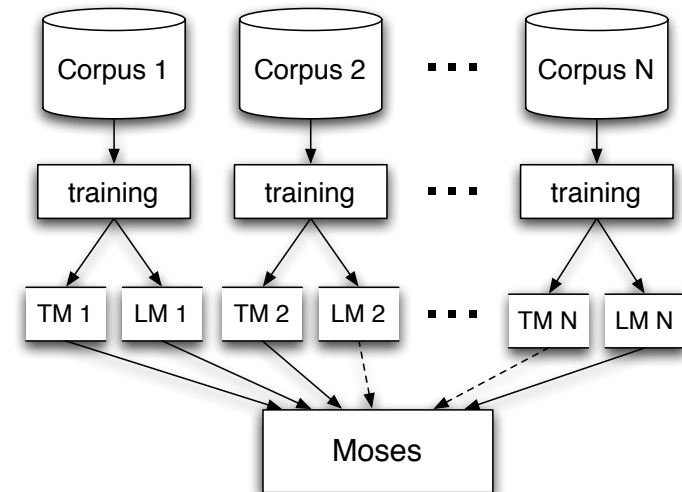
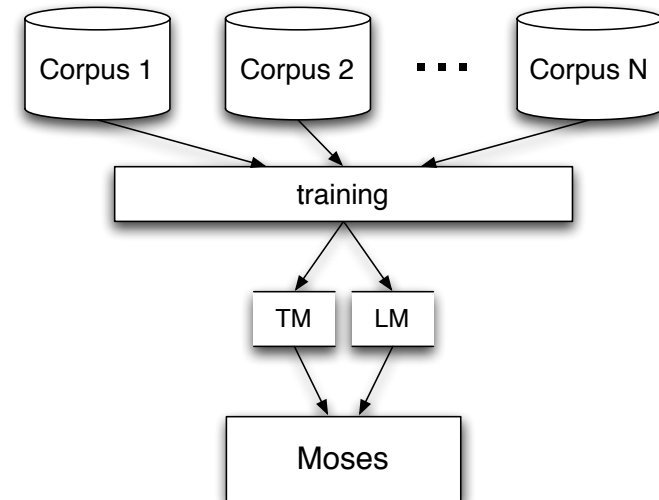
symmetrization	text for word-alignment	# phrase pairs	BLEU	NIST
grow-diag-final	case-sensitive	496K	20.50	5.57
"	case-insensitive	507K	21.86	5.59
+union	"	507K	22.35	6.20
+intersection	"	5.2M	22.71	6.31

- *multiple training corpora*
 - non-homogeneous data (size, domain)
 - small corpus for domain adaptation

- *multiple training corpora*
 - non-homogeneous data (size, domain)
 - small corpus for domain adaptation
- *one TM* and *one LM*
 - concatenation of all corpora
 - corpus characteristics are (too?) smoothed



- *multiple training corpora*
 - non-homogeneous data (size, domain)
 - small corpus for domain adaptation
- *one TM* and *one LM*
 - concatenation of all corpora
 - corpus characteristics are smoothed
- *multiple TMs* and *multiple LMs*
 - *advantages*
 - * more specialized models, more flexibility
 - * easy combination/selection of models
 - * effective (for TMs)
 - *drawbacks*
 - * complexity of the model



How much improvement do we get?

How much improvement do we get?

- evaluation
 - task: IWSLT Italian-to-English, second half of 2007 dev set
 - training data:
 - * baseline: BTEC, Named Entities, MultiWordNet and dev sets ('03-'06):
3.8M phrase pairs, 362K 4-grams
 - * EU Proceedings (39M phrase pairs, 16M 4-grams)
 - * Google Web 1T (336M 5-grams)
 - weight optimization on the first half of 2007 devset
 - verbatim input repunctuated with CN, case-insensitive

TM ₁ ,LM ₁	TM ₂ ,LM ₂	LM ₃	OOV	BLEU	NIST
baseline	-	-	1.68	28.70	5.76

How much improvement do we get?

- evaluation
 - task: IWSLT Italian-to-English, second half of 2007 dev set
 - training data:
 - * baseline: BTEC, Named Entities, MultiWordNet and dev sets ('03-'06):
3.8M phrase pairs, 362K 4-grams
 - * EU Proceedings (39M phrase pairs, 16M 4-grams)
 - * Google Web 1T (336M 5-grams)
 - weight optimization on the first half of 2007 devset
 - verbatim input repunctuated with CN, case-insensitive

TM ₁ ,LM ₁	TM ₂ ,LM ₂	LM ₃	OOV	BLEU	NIST
baseline	-	-	1.68	28.70	5.76
"	-	web	"	29.66	5.83

How much improvement do we get?

- evaluation
 - task: IWSLT Italian-to-English, second half of 2007 dev set
 - training data:
 - * baseline: BTEC, Named Entities, MultiWordNet and dev sets ('03-'06):
3.8M phrase pairs, 362K 4-grams
 - * EU Proceedings (39M phrase pairs, 16M 4-grams)
 - * Google Web 1T (336M 5-grams)
 - weight optimization on the first half of 2007 devset
 - verbatim input repunctuated with CN, case-insensitive

TM ₁ ,LM ₁	TM ₂ ,LM ₂	LM ₃	OOV	BLEU	NIST
baseline	-	-	1.68	28.70	5.76
"	-	web	"	29.66	5.83
"	EP	"	0.28	30.79	5.92

1-best vs. Confusion Networks

1-best vs. Confusion Networks

task	input	BLEU
IE, ASR	1bst	41.51
	cn	42.29*

* primary run

- CN outperforms 1-best

1-best vs. Confusion Networks

task	input	BLEU
IE, ASR	1bst	41.51
	cn	42.29*
JE, ASR	1bst	39.46*
	cn	39.69

* primary run

- CN outperforms 1-best
- no inspection on CN for JE

Multiple TMs and LMs

Multiple TMs and LMs

task	TMs	LMs	BLEU
IE, clean	baseline	baseline	43.41
	+EP	+EP+web	44.32*

* primary run

Multiple TMs and LMs

task	TMs	LMs	BLEU
IE, clean	baseline	baseline	43.41
	+EP	+EP+web	44.32*
IE, ASR, CN	baseline	baseline	40.74
	+EP	+EP+web	41.51*

* primary run

Multiple TMs and LMs

task	TMs	LMs	BLEU
IE, clean	baseline	baseline	43.41
	+EP	+EP+web	44.32*
IE, ASR, CN	baseline	baseline	40.74
	+EP	+EP+web	41.51*
CE, clean	baseline	baseline	35.08
	"	+web	33.94
	+LDC	"	34.72*

* primary run

- additional TMs improves performance (+0.77 BLEU)
- Google Web LM severely affects performance on CE (-1.14 BLEU)

- punctuation insertion in other languages (Chinese, Japanese)
- use of *caseing* CN to for case restoring

- punctuation insertion in other languages (Chinese, Japanese)
- use of *caseing* CN to for case restoring
- automatic way of selecting corpora

- punctuation insertion in other languages (Chinese, Japanese)
- use of *caseing* CN to for case restoring
- automatic way of selecting corpora
- further inspection on the use of Google Web corpus

Thank you!

Chinese-to English

- word-alignment on ci texts, grow-diag-final + union + inter
- case sensitive models
- distortion models: distance-based and orientation-bidirectional-fe
- (stack size, translation option limit, reordering limit)=(2000,50,7)
- BTEC and dev sets ('03-'07) (TM₁: 5.9M phrase pairs, LM₁: 39K 6-grams)
LDC: (TM₂: 27M phrase pairs)
Google Web (LM₂: 336M 5-grams)
- 5 official runs

Japanese-to English

- word-alignment on ci texts, grow-diag-final + union + inter
- case sensitive models
- distortion models: distance-based and orientation-bidirectional-fe
- (stack size, translation option limit, reordering limit)=(2000,50,7)
- BTEC and dev sets ('03-'07) (TM₁: 9.1M phrase pairs, LM₁: 39K 6-grams)
Reuters: (TM₂, 176K phrase pairs)
- 6 official runs

Italian-to English

- word-alignment on ci texts, grow-diag-final + union
- case insensitive TMs and LMs and case restoring
- distortion models: distance-based
- (stack size, translation option limit, reordering limit)=(200,20,6)
- BTEC NE, MWN, dev sets ('03-'07) (TM₁: 3.8M phrase pairs, LM₁: 362K 4-grams)
EU Proceedings: (TM₂: 39M phrase pairs, LM₂: 16M 4-grams)
Google Web (LM₃: 336M 5-grams)
- rescoring with 5K-best translations
- case-restoring with a 4-gram LM
- 12 official runs

- **Toolkit for SMT:**
 - translation of both text and CN inputs
 - incremental pre-fetching of translation options
 - handling multiple lexicons and LMs
 - handling of huge LMs and LexMs (up to Giga words)
 - on-demand and on-disk access to LMs and LexMs
 - factored translation model (surface forms, lemma, POS, word classes, ...)
- **Multi-stack DP-based decoder:**
 - theories stored according to the coverage size
 - synchronous on the coverage size
- **Beam search:**
 - deletion of less promising partial translations:
 - histogram and threshold pruning
- **Distortion limit:** reduction of possible alignments
- **Lexicon pruning:** limit the amount of translation options per span

- *log-linear statistical model*
- features of the *first* pass
 - (multiple) language models
 - direct and inverted word- and phrase-based (multiple) lexicons
 - word and phrase penalties
 - reordering model: distance-based and lexicalized (CE, JE)
- (additional) features of the *second* pass (IE)
 - direct and inverse IBM Model 1 lexicon scores
 - weighted sum of n -grams relative frequencies ($n = 1, \dots, 4$) in N -best list
 - the reciprocal of the rank
 - counts of hypothesis duplicates
 - n -gram posterior probabilities in N -best list [Zens, 2006]
 - sentence length posterior probabilities [Zens, 2006]