

# Consensus Decoding Approaches to Statistical Machine Translation

12/14/2011

Nan DUAN

Advisor: Ming ZHOU and Mu LI  
School of Computer Science and Technology  
Tianjin University

# Outline

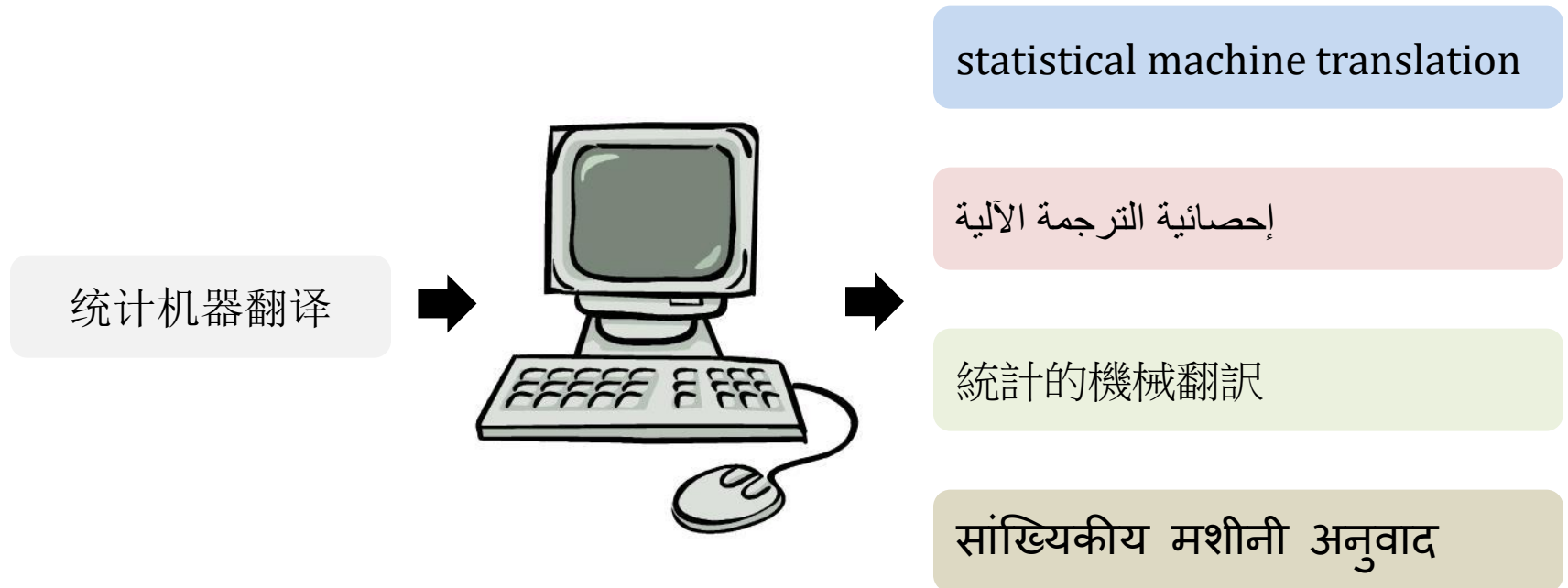
- Consensus Decoding to SMT
- Collaborative Decoding
- Mixture Model-based MBR Decoding
- Hypothesis Mixture Decoding
- Conclusion

# Outline

- Consensus Decoding to SMT
- Collaborative Decoding
- Mixture Model-based MBR Decoding
- Hypothesis Mixture Decoding
- Conclusion

# “Statistical” Machine Translation

- **MT**: translation from one natural language into another natural language, using computer
- **SMT**: an MT paradigm based on *statistical model* and analysis of *bilingual data*



# Modeling

- Source-Channel Model

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e \{P(e|f)\} \\ &= \operatorname{argmax}_e \{P(e) \cdot P(f|e)\}\end{aligned}$$

Decoding  
Algorithm

Language  
Model

Translation  
Model

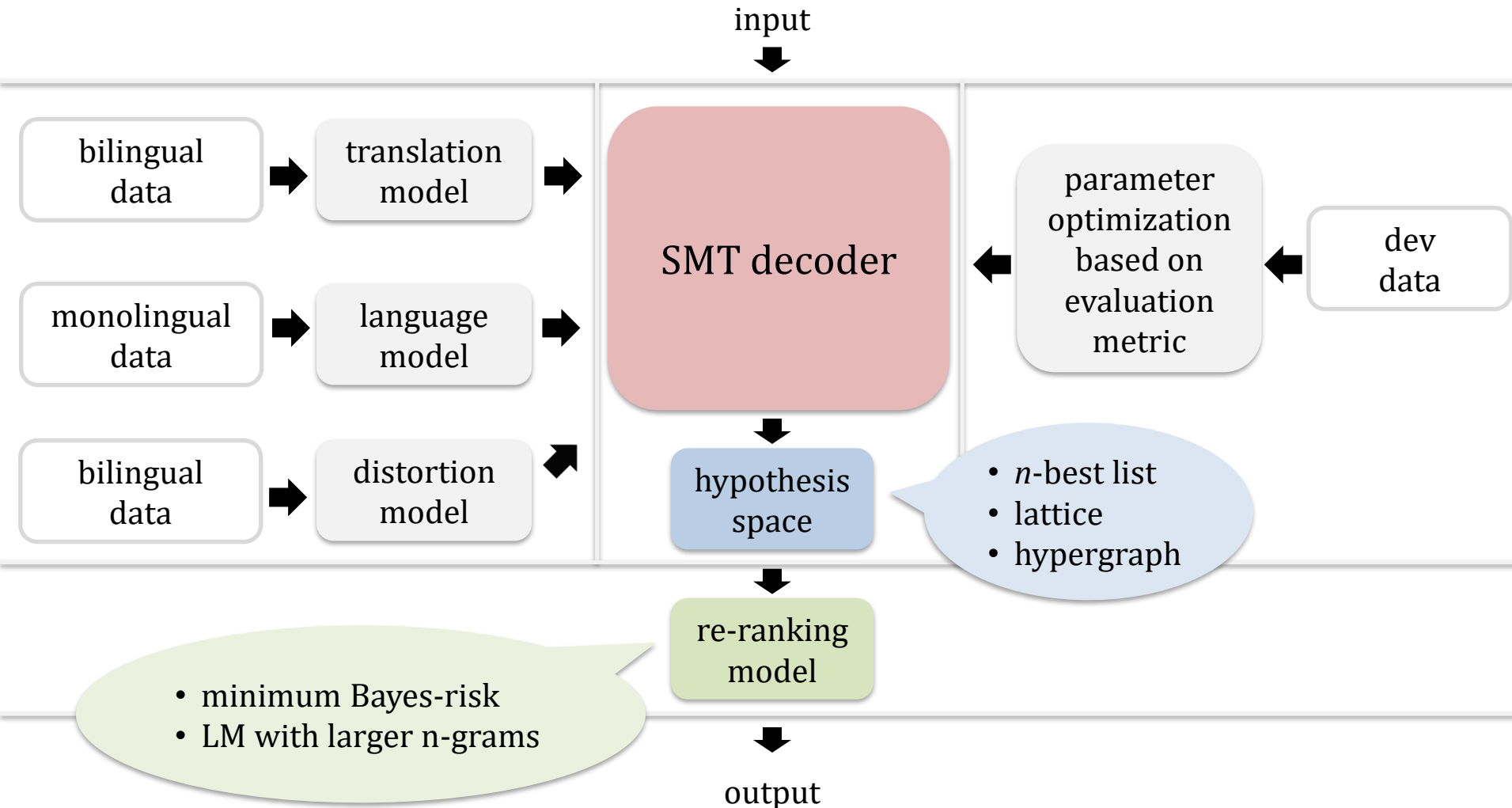
- Maximum Entropy Model

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e \frac{\exp[\sum_{m=1}^M \lambda_m \cdot h_m(e, f)]}{\sum_{e'} \exp[\sum_{m=1}^M \lambda_m \cdot h_m(e', f)]} \\ &= \operatorname{argmax}_e \{ \sum_{m=1}^M \lambda_m \cdot h_m(e, f) \}\end{aligned}$$

Parameter  
Optimization

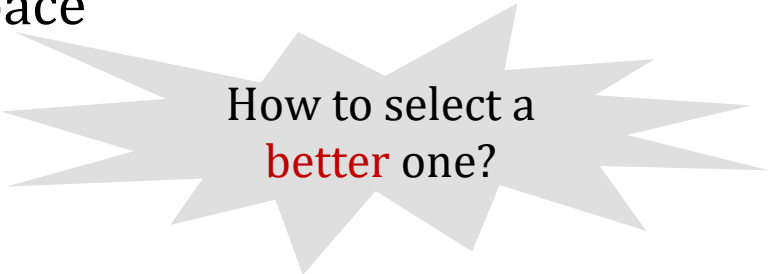
Feature  
Selection

# Generic System Overview



# Existing Problems

- Single SMT model
  - 1-best translation is usually not the REAL-BEST one
  - better ones exist in the hypothesis space



How to select a **better** one?

- Multiple SMT models
  - different pros and cons
  - phrase-based models
    - **good** at handling local reorderings
    - **bad** at dealing with long-distance dependencies
  - syntax-based models
    - **good** at handling long-distance dependencies
    - **bad** at phrase coverage



How to **combine** **merits** of different SMT models?

# Minimum Bayes-Risk Decoding

- Seek the translation with *the least expected loss* under a probability distribution

$$\begin{aligned}\hat{e} &= \underset{e' \in H(f)}{\operatorname{argmin}} \sum_{e \in H(f)} L(e, e') P(e/f) \\ &= \underset{e' \in H(f)}{\operatorname{argmax}} \sum_{e \in H(f)} G(e, e') P(e/f)\end{aligned}$$

Diagram illustrating the Minimum Bayes-Risk Decoding formula with callouts:

- loss function**:  $L(e, e')$
- hypothesis space**:  $H(f)$
- gain function**:  $G(e, e')$
- probability distribution**:  $P(e/f)$

*MBR Decoding on Hypergraph*  
[Kumar et al., 2009]

*MBR Decoding on Lattice*  
[Tromble et al., 2008]

*Using BLEU as the loss function*  
[Ehling et al., 2007]

*MBR Decoding on N-best*  
[Kumar and Byrne, 2004]

Related work



# Word-Level System Combination

- Given outputs generated by multiple SMT systems
  - 1) align multiple hypotheses to build a confusion network
  - 2) output the path with the highest score as final translation

I like eating chocolate icecream.



我喜欢巧克力冰激凌。  
 我喜欢吃巧克力冰激凌。  
 我爱吃巧克力冰激凌。  
 我爱巧克力冰激凌。



我	喜欢	~	巧克力	冰激凌	。
我	喜欢	吃	巧克力	~	。
我	爱	吃	巧克力	冰淇淋	。
我	爱	~	巧克力	冰激凌	。

confidence  
score

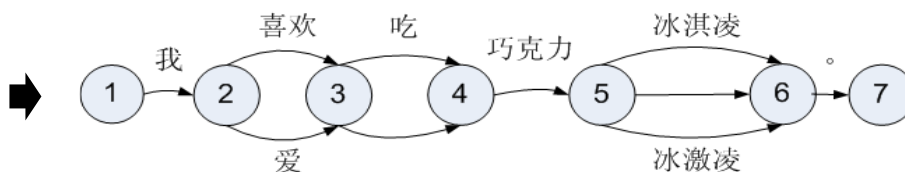
null  
penalty

LM  
score

length  
penalty

$$\hat{e} = \underset{e \in H(f)}{\operatorname{argmax}} \{ \alpha P(e) + \beta N_{\text{null}}(e) + \chi \log P_{LM}(e) + \sigma N_{\text{words}}(e) \}$$

我爱吃巧克力冰激凌。



# Consensus Decoding

- Techniques that can *re-rank* or *re-produce* better translations by making use of *consensus statistics* computed between hypotheses
  1. **system construction**
    - single system
    - multiple systems
  2. **first-pass decoding**
    - n-best list
    - hypergraph/lattice
  3. **search space re-construction**
    - simple union
    - more complicated structure
  4. **second-pass decoding**
    - re-rank existing translations
    - re-produce unseen translations

# Outline

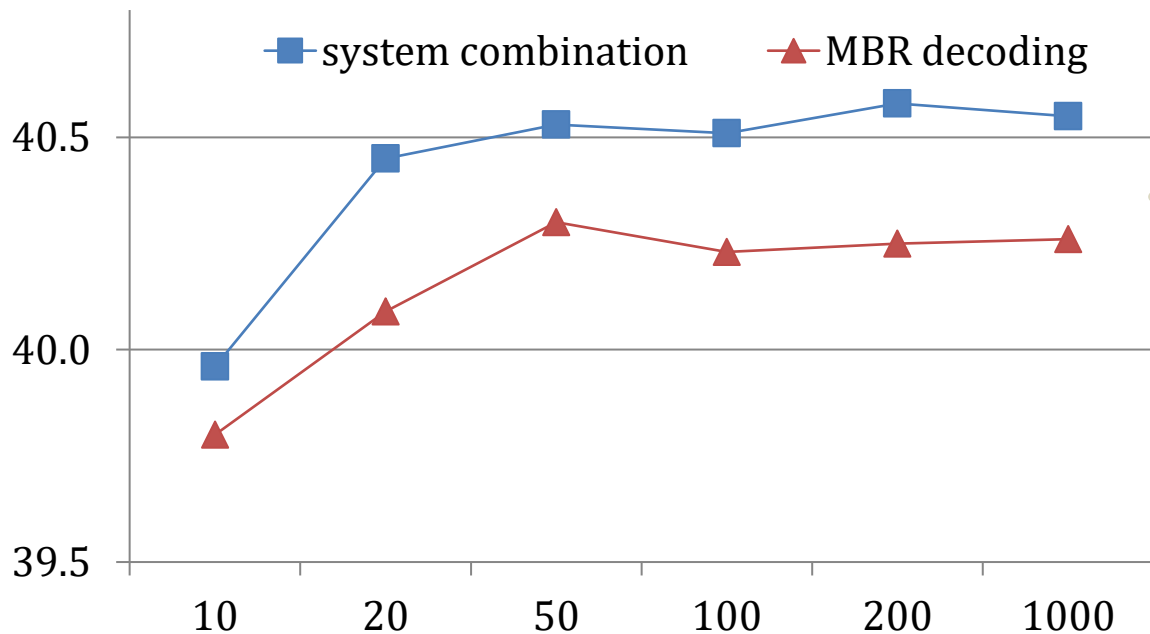
- Consensus Decoding to SMT
- Collaborative Decoding
- Mixture Model-based MBR Decoding
- Hypothesis Mixture Decoding
- Conclusion

# Collaborative Decoding

*ACL, 2009*

# Motivation

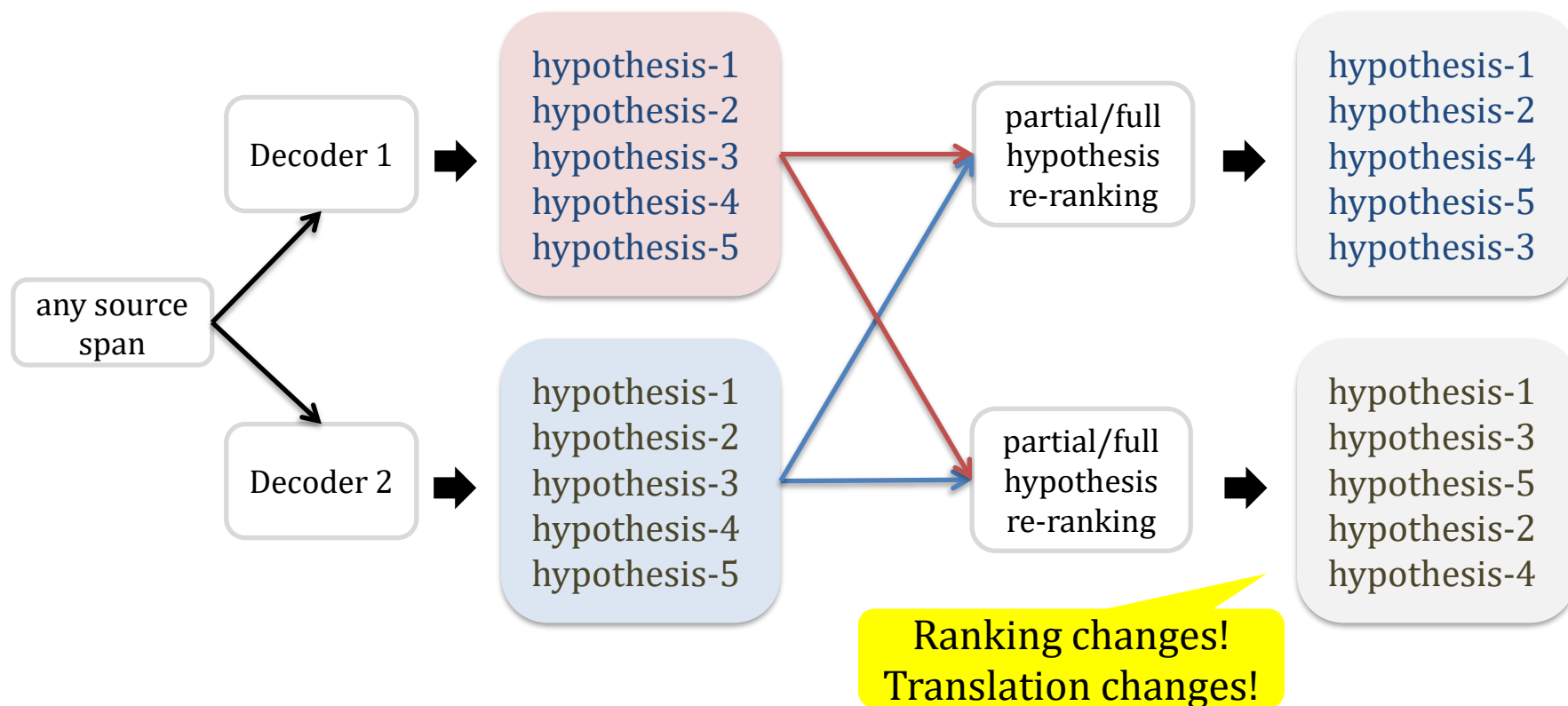
- MBR decoding & system combination
  - take **n-best translations** as input
    - present *a small portion* of the whole search space
    - some useful partial translations are *pruned* during decoding



no significant  
gains using  
larger n-best

# Collaborative Decoding

- Multiple SMT decoders work collaboratively
  - hypothesis scoring/re-ranking in decoding directly
  - explore translations beyond n-best lists of full translations



# Co-decoding Model

- Baseline model

- original SMT model

$$\Phi_m(f, e) = \sum_{i=1}^N \lambda_{m,i} \cdot h_{m,i}(f, e)$$

model specified  
features

- Augmented model

- consensus-based model

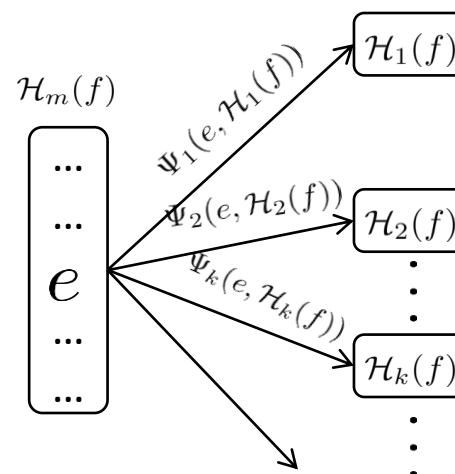
$$\sum_{k, k \neq m} \Psi_k(e, H_k(f)) = \sum_{k, k \neq m} \sum_l \lambda_{k,l} \cdot h_{k,l}(e, H_k(f))$$

n-gram -based  
features

- Co-decoding model

- baseline model + augmented models

$$F_m(e) = \Phi_m(f, e) + \sum_{k, k \neq m} \Psi_k(e, H_k(f))$$



# Co-decoding Feature

$$h_{k,l}(e, H_k(f)) = \sum_{e' \in H_k(f)} \underbrace{P(e' | H_k(f)) G_l(e, e')}_{\text{Co-decoding Feature}}$$

## Posterior Probability

- alpha is the scaling-factor that controls the shape of the probability distribution

$$P(e' | H_k(f)) = \frac{\exp\{\alpha \cdot F_k(e')\}}{\sum_{e'' \in H_k(f)} \exp\{\alpha \cdot F_k(e'')\}}$$

## n-gram Measure Functions

- n-gram **agreement** measure

$$G_n^+(e, e') = \sum_{i=1}^{|e|-n+1} \delta(e_i^{i+n-1}, e'_i)$$

- n-gram **disagreement** measure

$$G_n^-(e, e') = \sum_{i=1}^{|e|-n+1} (1 - \delta(e_i^{i+n-1}, e'_i))$$



# Experiments

- Training
  - bilingual data (5.1M sentence pairs)
    - all data available for the NIST 2008 constrained track of C-to-E MT task
  - monolingual data
    - Xinhua portion of LDC English Gigaword V3.0
- Evaluation
  - development data
    - NIST 2003 (919 sentences)
  - test data
    - NIST 2005 (1,082 sentences) and NIST 2008 (1,357 sentences)
- Metric
  - case insensitive NIST BLEU

# Experiments

- Baseline system

- Hiero
  - (Chiang, 2005)
- BTG
  - (Xiong et al., 2006)
- DepHiero
  - (Shen et al., 2008)

- Combination system

- more than one SMT model are used in co-decoding, so further combination is straightforward
  - word-level
    - (Rosti et al., 2007)
  - sentence-level
    - (Hildebrand and Vogel, 2008)

# Overall Comparison

	NIST 2005	NIST 2008
	before co-decoding/ <b>after co-decoding</b>	
Hiero	38.66%/ <b>40.08%</b>	27.67%/ <b>29.19%</b>
BTG	38.06%/ <b>39.93%</b>	27.25%/ <b>29.14%</b>
DepHiero	39.50%/ <b>40.32%</b>	28.75%/ <b>29.68%</b>
Further Combination		
Word-Level Combination	40.45%/ <b>40.85%</b>	29.52%/ <b>30.35%</b>
Sentence-Level Combination	40.09%/ <b>40.50%</b>	29.02%/ <b>29.71%</b>

# Summary

- Advantages
  - re-rank both partial and full translations
  - improve all involved SMT systems
- Future work
  - more systems included

# Mixture Model-based MBR Decoding

*COLING, 2010*

# Motivation

MBR decoding improves over max-derivation decoding

$$\hat{e} = \underset{e' \in H(f)}{\operatorname{argmin}} \sum_{e \in H(f)} L(e, e') P(e|f)$$

## Limitations of MBR decoding

- It can only perform on single SMT system
- It can not adapt to multiple SMT systems

**Model distribution is the key point!**

*MBR Decoding for  
Hypergraph*

[Kumar et al., 2009]

*MBR Decoding for Lattice*

[Tromble et al., 2008]

*MBR Decoding for n-best*

[Kumar and Byrne, 2004]

...

**Different Search Spaces**

*Log-BLEU*

[Tromble et al., 2008]

*BLEU*

[Ehling et al., 2007]

*Word Error Rate (WER)*

*Position-independent WER*

[Kumar and Byrne., 2004]

...

**Different Loss Functions**

In this work, we present Mixture Model-based MBR decoding (*MMMMBR decoding*)

- integrates multiple model distributions for Bayes-risk computation;
- integrates multiple search spaces for hypothesis selection.

# Overview of MMMBR Model

- Generates multiple search spaces

- Integrates multiple search spaces as a combined search space

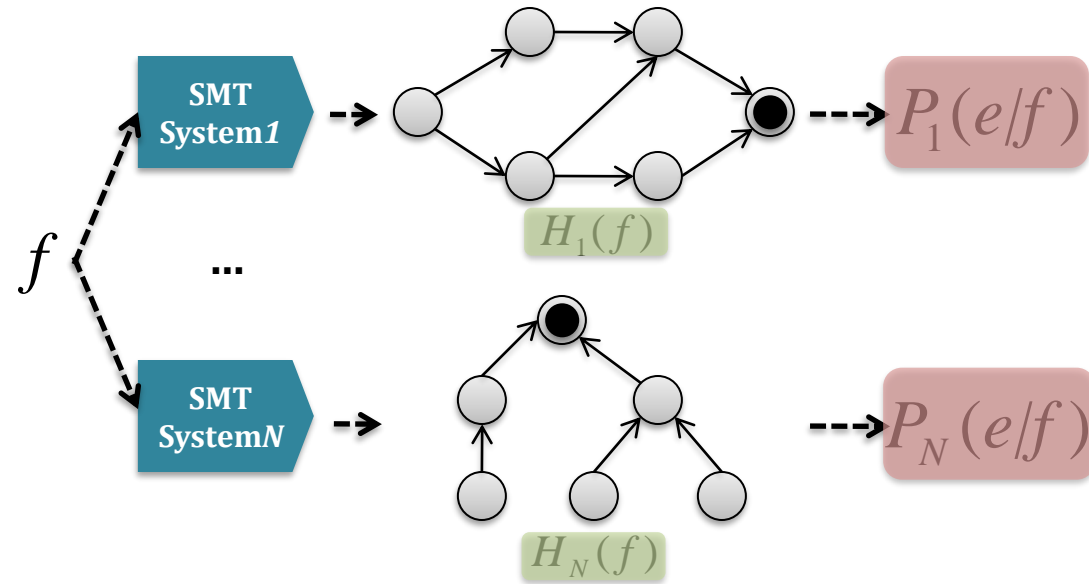
$$H(f) = \bigcup_{i=1}^N H_i(f)$$

- Integrates multiple distributions by mixture modeling

$$P(e/f) = \sum_{i=1}^N \lambda_i P_i(e/f)$$

$$\sum_{i=1}^N \lambda_i = 1 \quad 0 \leq \lambda_i \leq 1$$

- Rewrites MBR to MMMBR



$$\hat{e} = \underset{e' \in H(f)}{\operatorname{argmax}} \sum_{e \in H(f)} G(e, e') \left( \sum_{i=1}^N \lambda_i P_i(e/f) \right)$$

**Mixture Model**

# Mathematical Derivation

$$\begin{aligned}
 \hat{e} &= \underset{e' \in H}{\operatorname{argmax}} \sum_{e \in H} G(e, e') \sum_{i=1}^N \lambda_i P_i(e/f) \\
 &= \underset{e' \in H}{\operatorname{argmax}} \sum_{e \in H} \left\{ \sum_{i=1}^N \lambda_i \sum_{k=1}^N \sum_{e \in H_k} G(e, e') P_i(e/f) \right\} \\
 &= \underset{e' \in H}{\operatorname{argmax}} \sum_{e \in H} \left\{ \sum_{i=1}^N \lambda_i \sum_{k=1}^N \sum_{e \in H_k} \left\{ \theta_0/e' + \sum_{\omega} \theta_{|\omega|} \#_{\omega}(e') \delta_{\omega}(e) \right\} P_i(e/f) \right\} \\
 &= \underset{e' \in H}{\operatorname{argmax}} \sum_{e \in H} \left\{ \sum_{i=1}^N \lambda_i \sum_{k=1}^N \theta_0/e' + \sum_{\omega} \theta_{|\omega|} \#_{\omega}(e') p_i(\omega | H_i) \right\} \quad \text{MAP decoding score} \\
 &= \underset{e' \in H}{\operatorname{argmax}} \theta_0/e' + \sum_{\omega} \theta_{|\omega|} \#_{\omega}(e') \sum_i \lambda_i p_i(\omega | H_i) + \sum_k \theta_k \log C_{MAP}(e' | f, d_k)
 \end{aligned}$$

- decomposing Bayes-risk computation to each local search space
- using *log-BLEU* (Tromble et al., 2008) as the similarity function
- using *Algorithm 4* (Kumar et al., 2009) for  $n$ -gram posterior probability computation



# Model Training

- A two-pass training procedure

$$\begin{aligned}\hat{e} &= \underset{e' \in H}{\operatorname{argmax}} \left[ \theta_0 / e' + \sum_{\omega} \theta_{|\omega|} \#_{\omega}(e') \sum_i \lambda_i p_i(\omega | H_i) + \sum_k \theta_k \log C_{MAP}(e' | f, d_k) \right] \\ &= \underset{e' \in H}{\operatorname{argmax}} \sum_i \lambda_i \left\{ \theta_0 / e' + \sum_{\omega} \theta_{|\omega|} \#_{\omega}(e') p_i(\omega | H_i) + \sum_k \theta_k \log C_{MAP}(e' | f, d_k) \right\}\end{aligned}$$

## Step1: for MBR and MAP

- Fix system weights in mixture model
- Run Mert to tune MBR and MAP parameters

## Step2: for Mixture Model

- Fix parameters in MBR and MAP models
- Run Mert to tune system weights

**This procedure can be iteratively processed!**

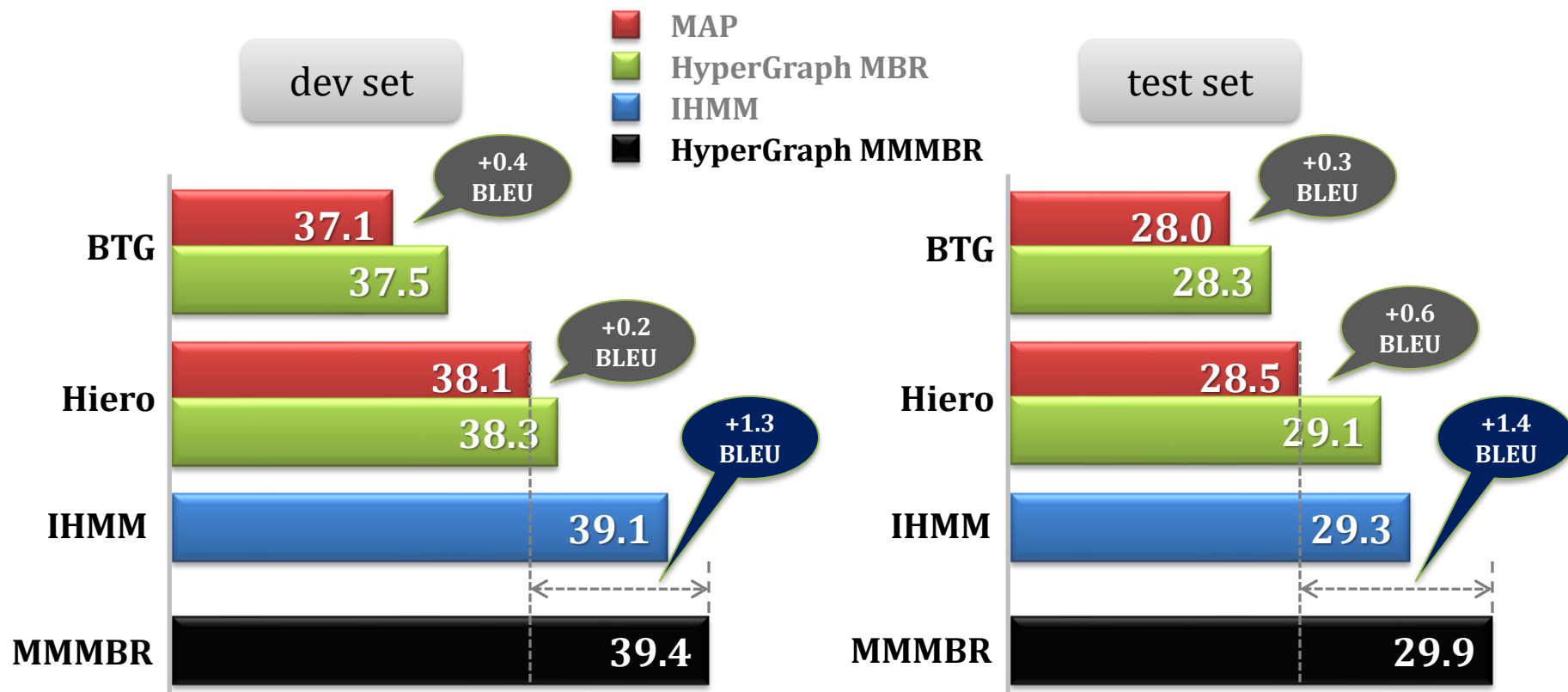
# Experiments

- Training
  - bilingual data (5.1M sentence pairs)
    - all data available for the NIST 2008 constrained track of C-to-E MT task
  - monolingual data
    - Xinhua portion of LDC English Gigaword V3.0
- Evaluation
  - development data
    - newswire portion of the NIST 2006 (616 sentences)
  - test data
    - NIST 2008 (1,357 sentences)
- Metric
  - case insensitive NIST BLEU

# Experiments




- Baseline system
  - Hiero
    - (Chiang, 2005)
  - BTG
    - (Xiong et al., 2006)
- Comparison technique
  - word-level system combination
    - (Li et al., 2009)

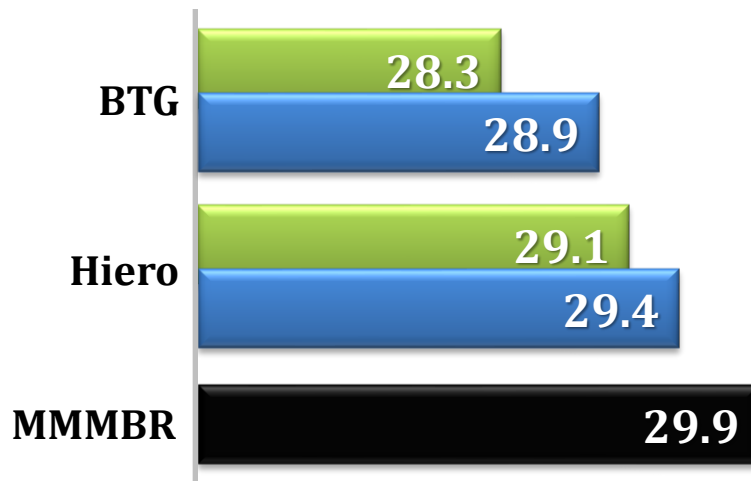
# Overall Comparison



- MBR based on single system improves, but few (+0.2~+0.6 BLEU points)
- MMMBR based on multiple systems improves significantly (+1.3~+1.4 BLEU points)
- MMMBR is also comparable to word-level system combination (+0.3~+0.6 BLEU points)

# Impacts of P(...) and H(...)

-  HyperGraph MBR
-  HyperGraph MBR based on Mixture Model
-  HyperGraph MMMBR



$$\hat{e} = \underset{e' \in H_i}{\operatorname{argmax}} \sum_{e \in H_i} G(e, e') P_i(e/f)$$

$$\hat{e} = \underset{e' \in H_i}{\operatorname{argmax}} \sum_{e \in H_i} G(e, e') \sum_{i=1}^N \lambda_i P_i(e/f)$$

$$\hat{e} = \underset{e' \in H}{\operatorname{argmax}} \sum_{e \in H} G(e, e') \sum_{i=1}^N \lambda_i P_i(e/f)$$

Using mixture  
model helps!

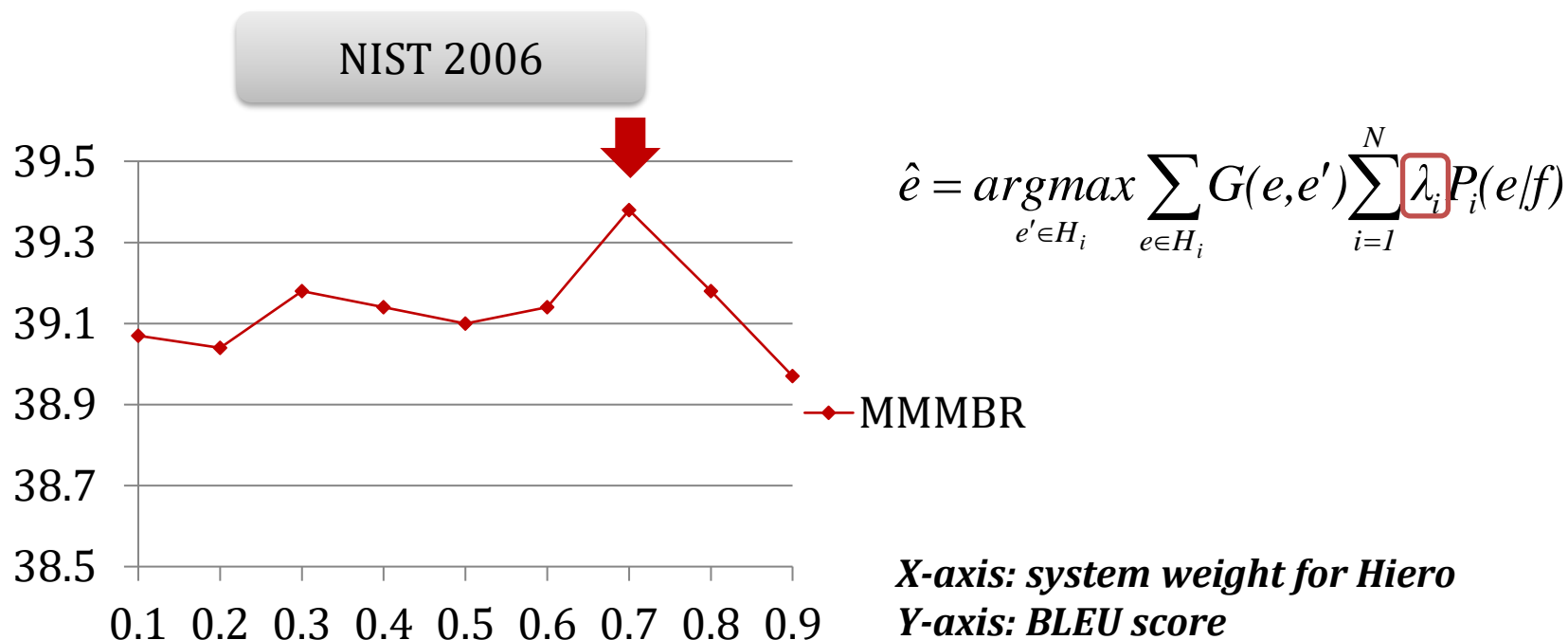
+

Enlarging search  
space helps!

=

best

# Impacts of System Weight



The performance is stable after the 1<sup>st</sup> round of two-pass training has been finished.

# Summary

- Advantages
  - significant gains achieved
    - better than MBR decoding
    - comparable to word-level system combination
  - large hypothesis spaces explored
  - flexible and tunable system weights
- Future work
  - more systems included
    - syntax-based models

# Hypothesis Mixture Decoding

*ACL, 2011*



# Motivation

## Various SMT Models

phrase-based, hierarchical phrase-based, syntax-based...



### System Combination



n-best input



confusion-network decoding



output

### MBR Decoding



n-best/hypergraph input



re-ranking



output

How to  
combine these  
two merits?

**Pros**

new hypotheses can be generated

**Cons**

few candidates can be leveraged

**Pros**

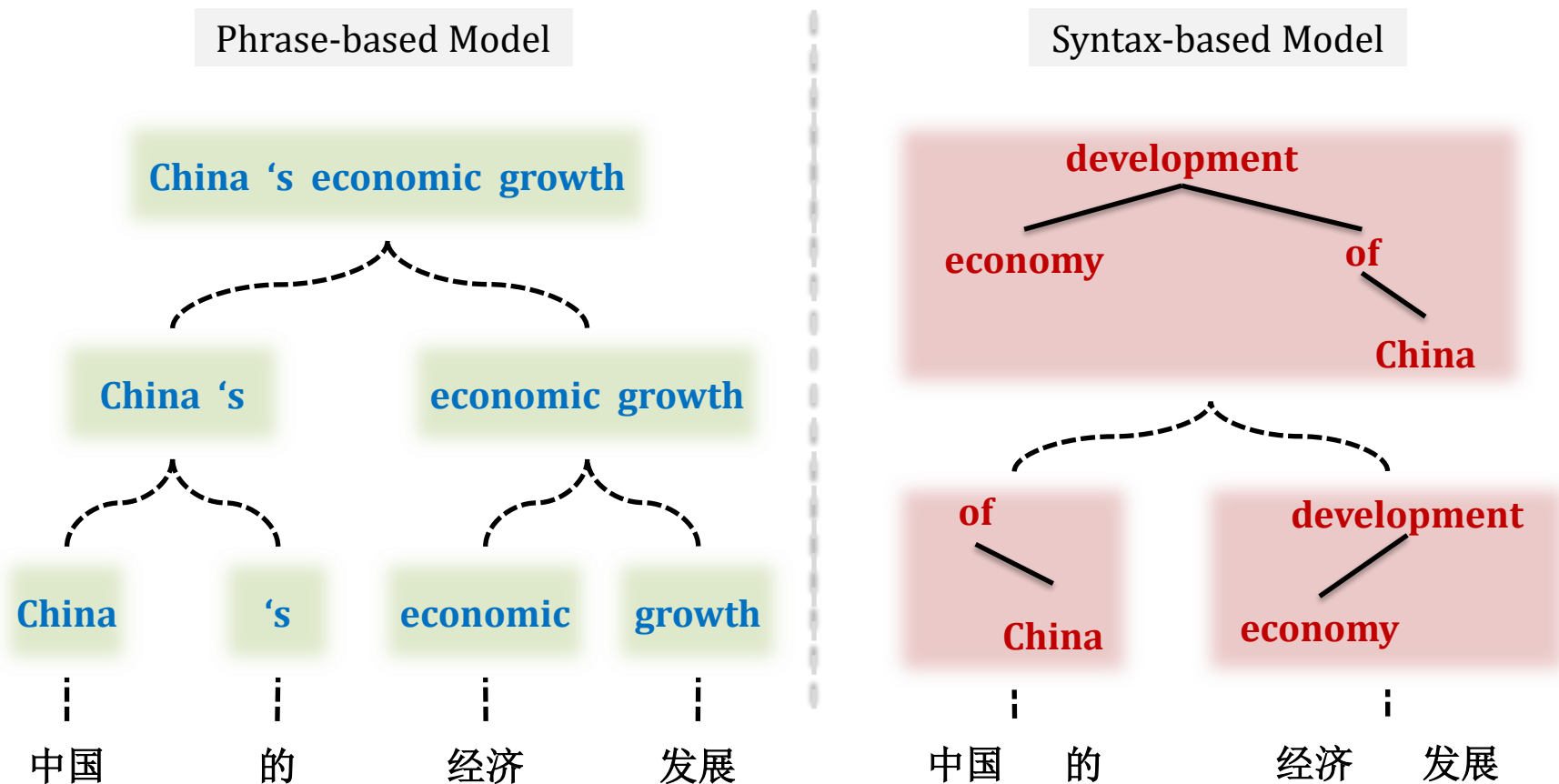
more candidates can be leveraged

**Cons**

no new hypothesis can be generated

# Hypothesis Mixture Decoding

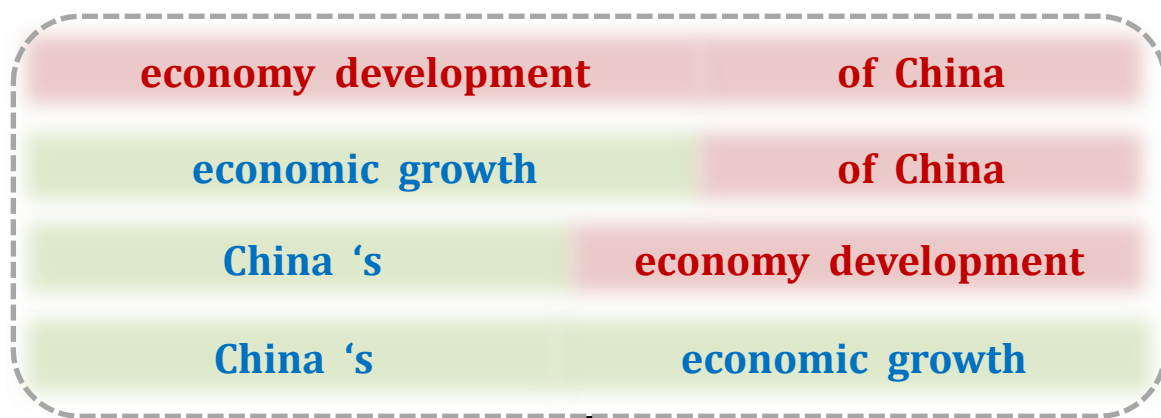
- Step-1: independent decoding



# Hypothesis Mixture Decoding

- Step-2: hypothesis mixture decoding (HMD)

Mixture Hypothesis Space



translations beyond existing hypothesis spaces can be generated



partial hypotheses generated by arbitrary SMT models can be used

中国 的

经济 发展

# Models and Features

$$\hat{e} = \underset{e \in H(f)}{\operatorname{arg\,max}} \sum_i \lambda_i \cdot h_i(e, f)$$

- Linear combination of two sets of features
  - **consensus-based features**
    - n-gram posteriors
    - length posteriors
  - **general features**
    - LM score
    - word penalty
    - count of lexicon pair
    - reordering score
    - count of new generated n-grams

# Two Decoding Algorithms

## HM Decoding Algorithm

```
1:  for each component model  $M_n$  do
2:      output the hypothesis space  $H_n(f)$  for the input
3:  end for
4:  for  $l = 1$  to  $|f| - 1$  do
5:      for all  $i, j$  s.t.  $j - i = l$  do
6:           $H(f_i^j) = \{nil\}$ 
7:          Hypothesis Re-Construction
8:          for each hypothesis  $e \in \bigcup_{n=1}^N H_n(f_i^j)$  do
9:              compute HMD features for  $e$ 
10:             add  $e$  to  $H(f_i^j)$ 
11:         end for
12:     end for
13: end for
14: return  $\hat{e} \in H(f)$  with the maximum score
```

## (I)BTG-based Hypothesis Re-Construction

```
1:  for all  $k$  s.t.  $i \leq k < j$  do
2:      for  $e_1 \in H(f_i^k)$  and  $e_2 \in H(f_{k+1}^j)$  do
3:          add  $e = Comb_{[\cdot]}(e_1, e_2)$  to  $H(f_i^j)$ 
4:          add  $e = Comb_{(\cdot)}(e_1, e_2)$  to  $H(f_i^j)$ 
5:      end for
6:  end for
```

## (II)SCFG-based Hypothesis Re-Construction

```
1:  for each rule  $r \in R$  that matches  $f_i^j$  do
2:      for  $e_1 \in H(r_{\#1})$  and  $e_2 \in H(r_{\#2})$  do
3:          add  $e = Comb_r(e_1, e_2)$  to  $H(f_i^j)$ 
4:      end for
5:  end for
```

# Experiments

- Training
  - bilingual data (5.1M sentence pairs)
    - all data available for the NIST 2008 constrained track of C-to-E MT task
  - monolingual data
    - Xinhua portion of LDC English Gigaword V3.0
- Evaluation
  - development data
    - NIST 2004 (1,788 sentences)
  - test data
    - NIST 2005 (1,082 sentences)
    - newswire portion of NIST 2006 (616 sentences) / NIST 2008 (691 sentences)
- Metric
  - case insensitive NIST BLEU

# Experiments

- **Baseline System**
  - DHPB: string-to-dependency system
    - (Shen et al., 2008)
  - PB: phrase-based system
    - (Xiong et al., 2006)
- **Comparison Technique**
  - Comb: word-level system combination
    - (Li et al., 2009)
  - CD: collaborative decoding
    - (Li et al., 2009)
  - MMMBR: mixture model-based MBR decoding
    - (Duan et al., 2010)

# Overall Comparison

		NIST 2004	NIST 2005	NIST 2006	NIST 2008
Baseline	DHPB	39.90%	39.76%	35.00%	30.43%
	PB	38.93%	38.21%	33.59%	29.62%
System Combination	Comb	41.14%	40.70%	36.04%	31.16%
Collaborative Decoding	CD-DHPB	40.81%	40.56%	35.73%	30.87%
	CD-PB	40.39%	40.34%	35.20%	30.39%
	CD-Comb	41.27%	41.02%	36.37%	31.54%
MMMBR Decoding	MMMBR	41.19%	40.96%	36.30%	31.43%
HM Decoding	HMD-BTG	41.24%	41.26%	36.76%	31.69%
	HMD-SCFG	41.31%	41.19%	36.63%	31.52%
	HMD-Comb	41.74%	41.53%	37.11%	32.06%



# Summary

- Advantages
  - HM decoding combines merits of word-level system combination and MBR decoding
    - large hypothesis spaces
    - capability to produce unseen translations
- Future work
  - explore richer consensus information
  - explore better reordering models for HM decoding

# Outline

- Consensus Decoding to SMT
- Collaborative Decoding
- Mixture Model-based MBR Decoding
- Hypothesis Mixture Decoding
- **Conclusion**

# Conclusion

- Consensus decoding helps SMT, by
  - *re-ranking* existing translations
    - MBR decoding
    - mixture model-based MBR decoding
    - model combination
    - collaborative decoding
    - ...
  - *re-producing* new translations
    - word-level system combination
    - hypothesis mixture decoding
    - ...
- Future work
  - explore richer consensus information
  - adapt to other NLP tasks

# My Research at MSRA

(since 2006.12-now)

- C-E and J-E SMT Projects
  - EngKoo
  - Bing Translator
  - Prototype SMT systems for research
- MT Competitions
  - NIST 2008 (international)
  - CWMT 2008 (domestic)
  - CWMT 2009 (domestic)
  - IWSLT 2010 (international)
    - joint work with MSR Redmond
- Publications
  - 10 international conference papers
    - 3 ACL, 2 COLING, 1 EMNLP and etc.


# APPENDIX


- More contents/details after...

# Feature Subspace-based System Combination

*EMNLP, 2009*

# Motivation

- Two ways to build multiple SMT systems for system combination
  - using models based on different paradigms
    - phrase-based
    - hierarchical phrase-based
    - syntax-based
    - ...

Challenging,  
even for “old birds”
  - using models based on different training methods
    - different word-aligners
    - different LMs
    - different word segmentations
    - ...

Time-consuming,  
sometimes not applicable
- Is there any *light-weight* method for system building?

# Feature Subspace-based System Combination (1)

- Ensemble Generation using Single SMT System

- most SMT systems are based on linear models

- feature space:  $\Omega = \{h_m(e,f) | m = 1, \dots, M\}$

- feature subspace:  $\Omega' \subset \Omega$

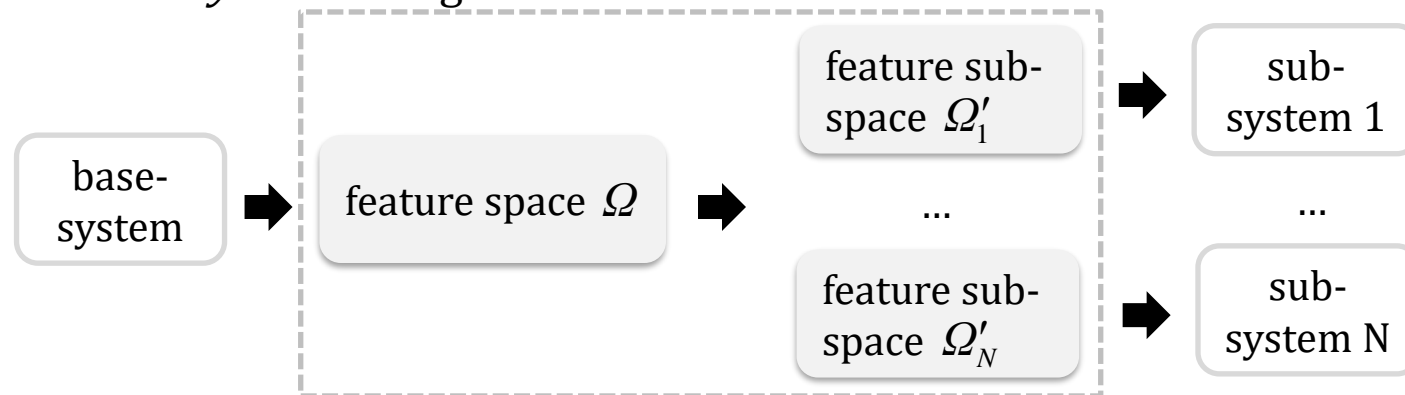
$$\hat{e} = \underset{e \in H(f)}{\operatorname{argmax}} \sum_{m=1}^M \lambda_m \cdot h_m(e,f)$$

a set of features

- each feature subspace maps to a new SMT system

- *base-system*: using all features

- *sub-system*: using a sub-set of all features





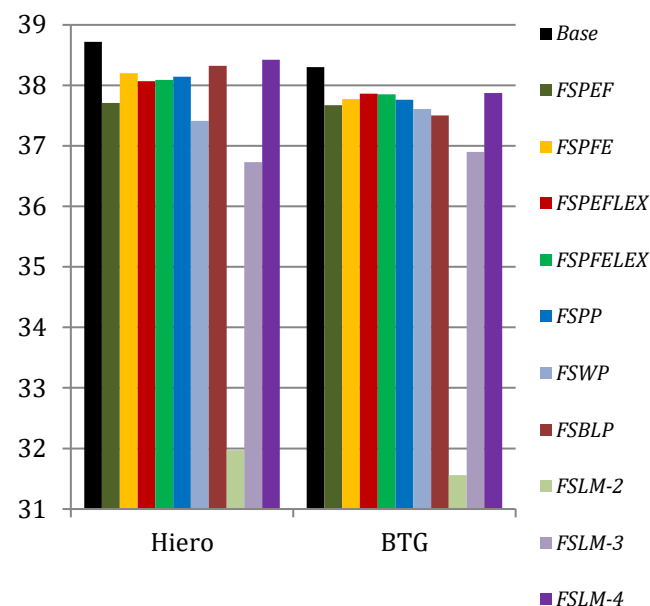
# Principle of Feature Selection

- Not all feature subspaces are helpful

- too many candidates:  $|\Omega'| = 2^{|\Omega|}$
- poor performances for most of them

- Our Feature Selection Principle

- for non-LM features
  - remove one feature each time
  - use the remainders for system construction
- for LM feature
  - lower its n-gram order
  - use it with the remainders for system construction



*Now's SMT  
systems CANNOT  
live without LM!*

# Feature Subspace-based System Combination (2)

- Sentence-Level System Combination
  - re-ranking the union of hypotheses generated by all systems

$$H(f) = \bigcup_{i=1}^N H_i(f)$$

$$\hat{e} = \underset{e \in H(f)}{\operatorname{argmax}} \{ \lambda_{LM} h_{LM}(e) + \lambda_l h_l(e) + \psi(e, H(f)) \}$$

simply union all hypotheses generated by base-system and sub-systems without bias

A weighted combination of n-gram consensus features

- n-gram *agreement* feature:  $h_n^+(e, H(f)) = \sum_{e' \in H(f)} G_n(e, e')$
- n-gram *disagreement* feature:  $h_n^-(e, H(f)) = \sum_{e' \in H(f)} (|e| - n + 1 - G_n(e, e'))$
- n-gram consensus:  $G_n(e, e') = \sum_{i=1}^{|e|-n+1} \delta(e_i^{i+n-1}, e')$

# Experiments

- Training
  - bilingual data (5.1M sentence pairs)
    - all data available for the NIST 2008 constrained track of C-to-E MT task
  - monolingual data
    - Xinhua portion of LDC English Gigaword V3.0
- Evaluation
  - development data
    - NIST 2003 (919 sentences)
  - test data
    - NIST 2004 (1,788 sentences) and NIST 2005 (1,082 sentences)
- Metric
  - case insensitive NIST BLEU

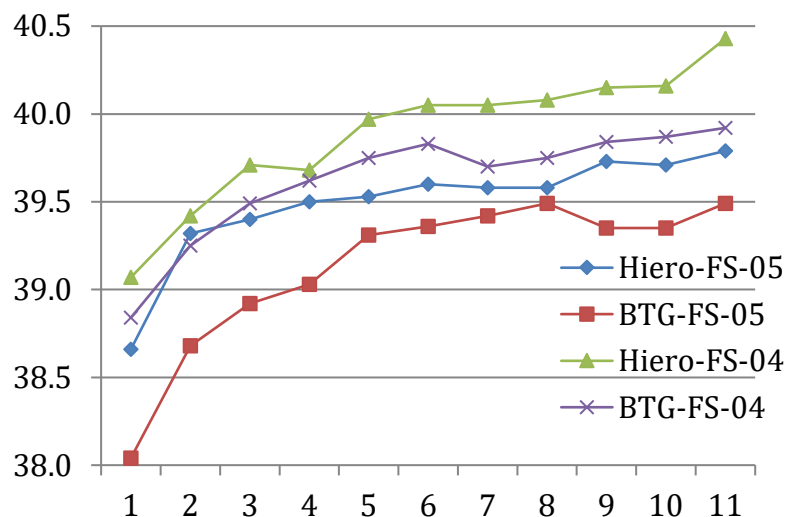
# Experiments

- Base-System
  - Hiero (Chiang, 2005)
  - BTG (Xiong et al., 2006)
- Sub-Systems
  - removing one of the non-LM features
    - *PEF/PFE*: source-to-target/target-to-source translation probability
    - *PEFLEX/PFELEX*: source-to-target/target-to-source lexical weight
    - *PP*: phrase penalty
    - *WP*: word penalty
    - *BLP*: bi-lexicon counting
  - lowering the order of an 5-gram LM
    - *LM-4/LM-3/LM-2*
  - (*7*+*3*)=10 sub-systems for each base-system

# Overall Comparison

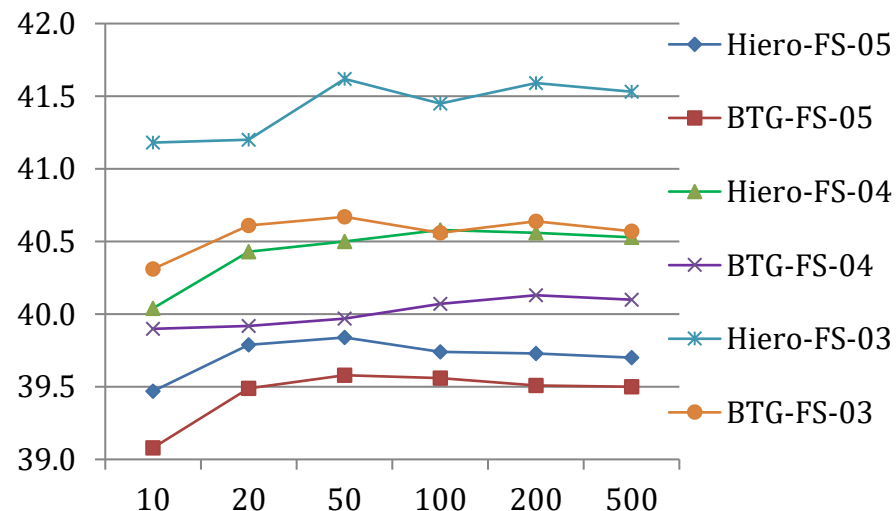
		NIST 2004	NIST 2005
Hiero	Base-system <small>220-best</small>	39.07%	38.72%
	Combination (1 base-system+ 10 sub-systems) <small>20-best for each system</small>	40.43% <b>(+1.36%)</b>	39.79% <b>(+1.07%)</b>
BTG	Base-system <small>220-best</small>	38.84%	38.30%
	Combination (1 base-system+ 10 sub-systems) <small>20-best for each system</small>	39.92% <b>(+1.08%)</b>	39.49% <b>(+1.19%)</b>
Hiero + BTG	Combination (Hiero + BTG) <small>220-best for each system</small>	39.98%	39.43%
	Combination (2 base-systems+ 20 sub-systems) <small>20-best for each system</small>	40.96% <b>(+0.98%)</b>	39.49% <b>(+1.06%)</b>

# Impacts of Different Settings



Impacts of different numbers of sub-systems for system combination

- X-axis: number of sub-systems
- Y-axis: BLEU score



Impacts of different sizes of n-best lists for system combination

- X-axis: n-best size for each system
- Y-axis: BLEU score

# Contributions of Sub-Systems

	NIST 2004	NIST 2005
Hiero	69.71%	69.69%
BTG	59.07%	58.54%

Ratio of unique hypotheses from sub-systems

	NIST 2004	NIST 2005
Hiero	44.63%	46.12%
BTG	47.54%	44.73%

Ratio of final translations from sub-systems

Oracle Performance		Hiero	BTG
		BLEU/TER	BLEU/TER
NIST 2004	Base-System <small>220-best</small>	49.68%/0.6411	49.50%/0.6349
	1 Base-System + 10 Sub-Systems <small>20-best for each system</small>	51.05%/0.6089	50.53%/0.6056
NIST 2005	Base-System <small>220-best</small>	48.89%/0.5946	48.37%/0.5944
	1 Base-System + 10 Sub-Systems <small>20-best for each system</small>	50.69%/0.5695	49.81%/0.5684

# Summary

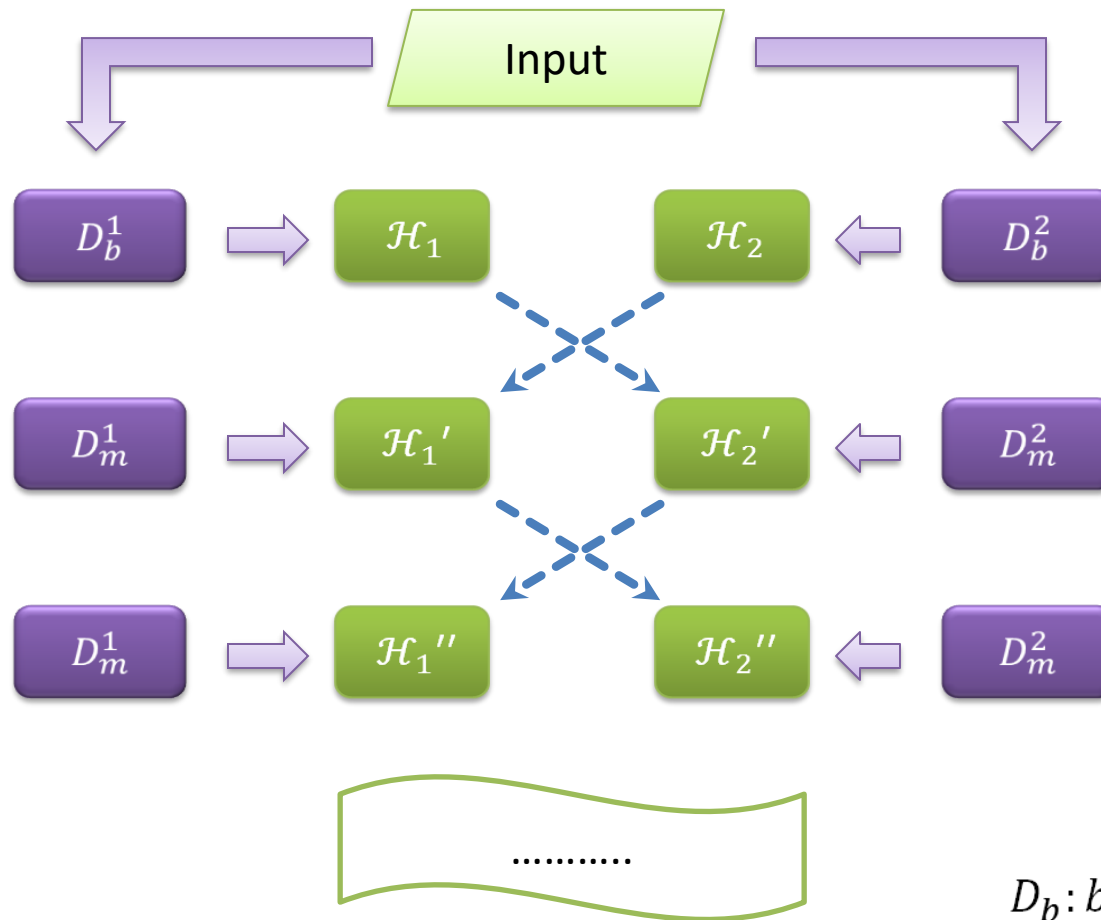
- Pros
  - simple and effective way for multiple system construction
  - useful to all SMT systems based on the linear model
- Cons
  - feature selection principle is simple



# Collaborative Decoding

*ACL, 2009*

# Decoder coordination: Iterative Decoding



$D_b$ : baseline decoder  
 $D_m$ : member decoder

# Impacts of Different Settings

NIST 2005	Hiero	BTG
Baseline	38.66%	38.04%
+agree -disagree	39.36%	39.02%
-agree +disagree	39.12%	38.67%
+agree +disagree	39.68%	39.61%

Impact of different consensus features

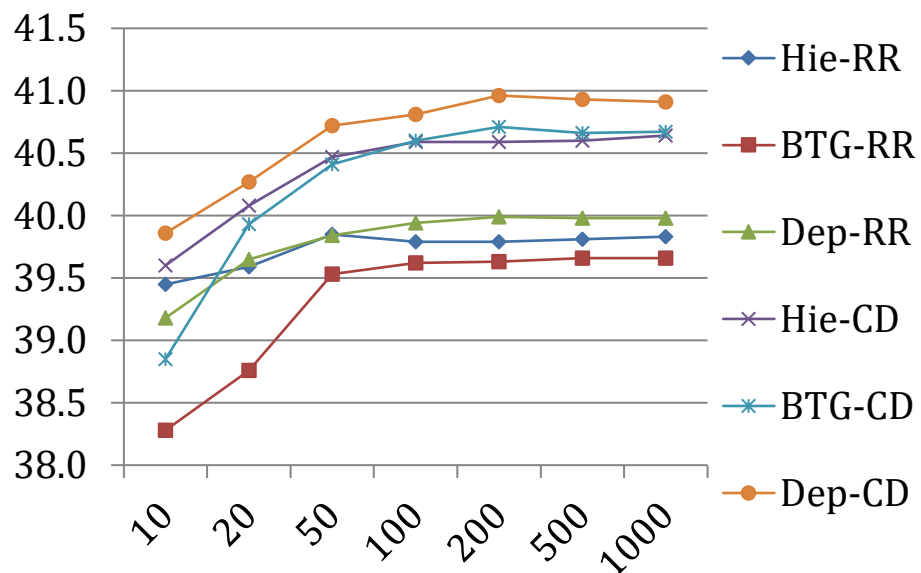
- Both n-gram agreement features and disagreement features are useful

TER before CD TER after CD	NIST 2005	NIST 2008
Hiero-BTG	0.3190 <b>0.2204</b>	0.4016 <b>0.2686</b>
Hiero-DepHiero	0.3252 <b>0.1840</b>	0.4176 <b>0.2469</b>
BTG-DepHiero	0.3498 <b>0.2171</b>	0.4238 <b>0.2665</b>

TER scores between translations

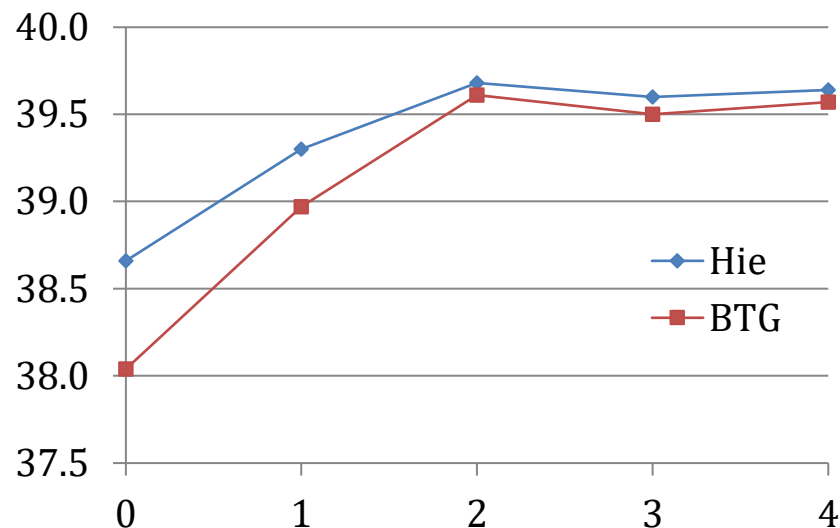
- Outputs become more similar due to the use of consensus information

# Impacts of Different Settings



## Co-decoding (CD) vs. Re-ranking(RR)

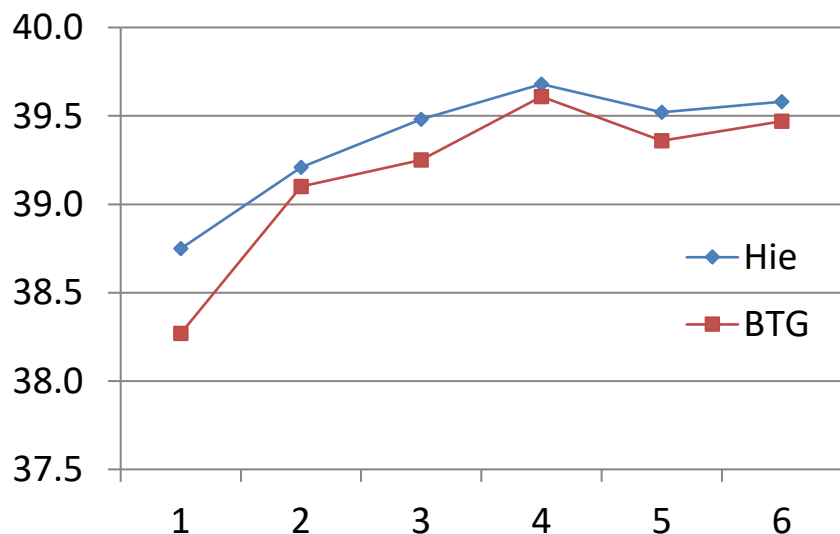
- X-axis: n-best size for each system
- Y-axis: BLEU score
- test data: NIST 2005



## Co-decoding with different iterations

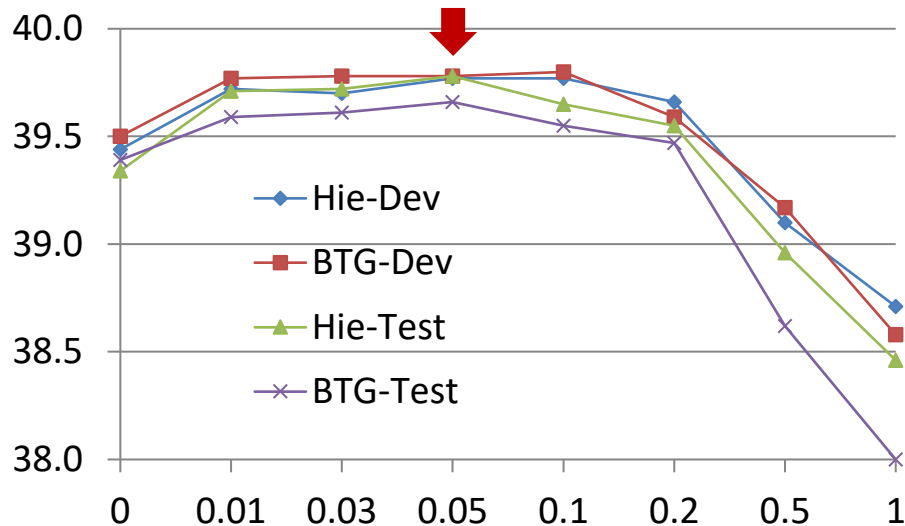
- X-axis: times of iteration
- Y-axis: BLEU score
- test data: NIST 2005

# Impacts of Different Settings



## Impact of n-gram consensus measure

- X-axis: the order of n-grams
- Y-axis: BLEU score
- NIST 2005



## Impact of scaling factor

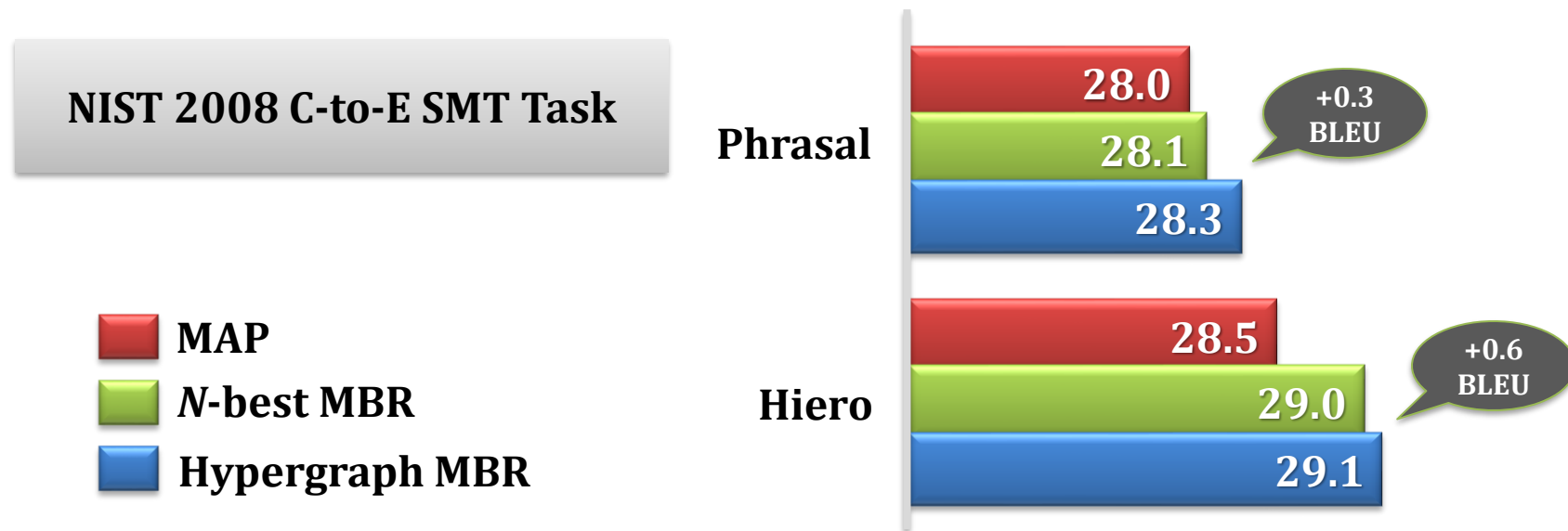
- X-axis: value of scaling factor
- Y-axis: BLEU score
- NIST 2003 (Dev) and NIST 2005 (Test)

# Mixture Model-based MBR Decoding

*COLING, 2010*

# Limitations of MBR Decoding

- Single system –based
- High correlation between hypotheses
- Relative few improvements achieved



***Extending MBR decoding to multiple SMT systems is straight-forward!***

# MMMBR Decoding Algorithm

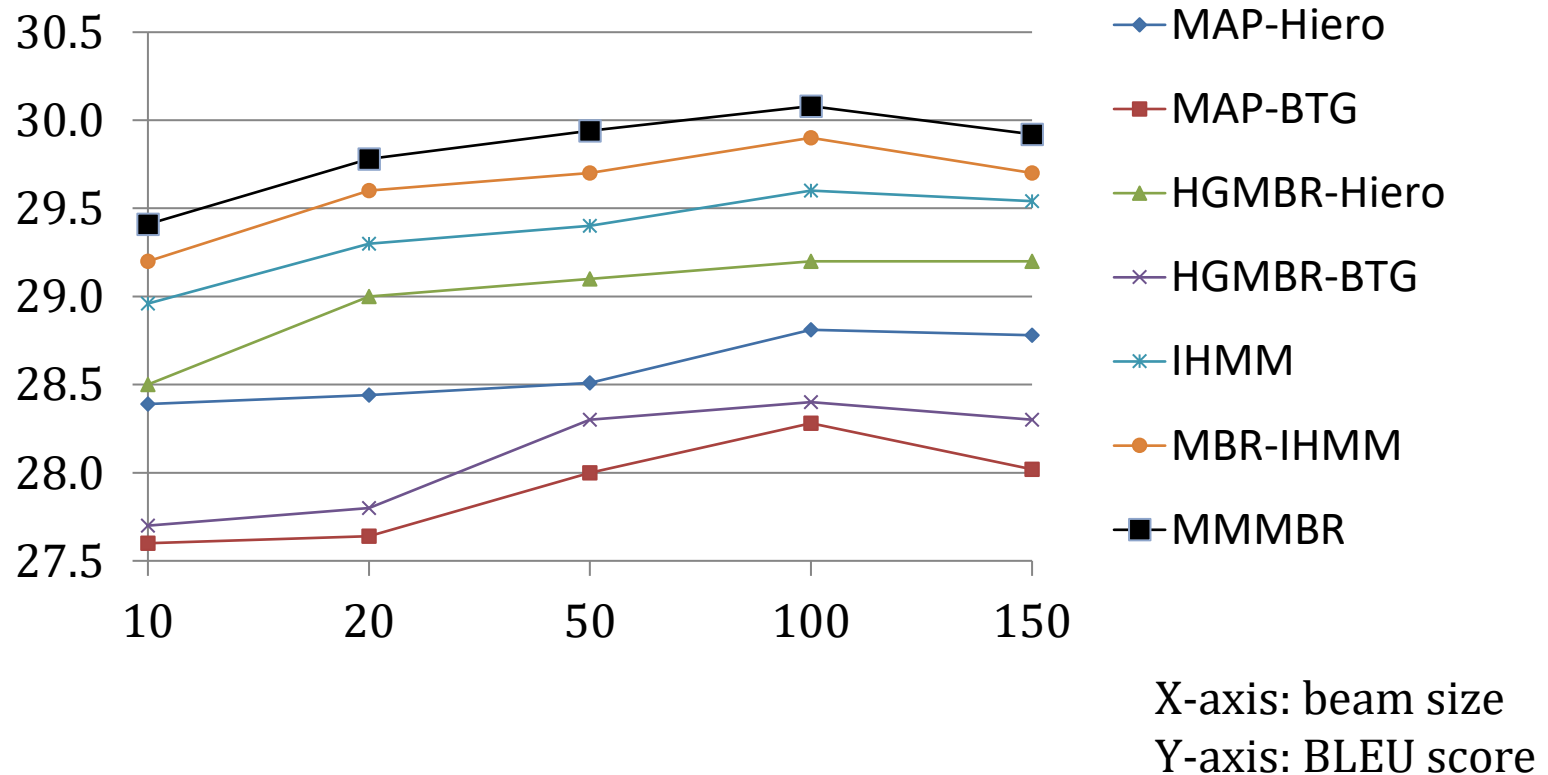
## MMMBR decoding on multiple MT systems

```
1:   for each component MT system  $d_k$  do
2:       generate the output search space  $H_k$ 
3:       compute the  $n$ -gram posterior probability set  $\{p_k(w/H_k)\}$  for  $H_k$ 
4:   end for
5:   compute the mixture  $n$ -gram posterior probability for each  $n$ -gram:
6:   for each unique  $n$ -gram  $w$  appeared in  $U_k H_k$  do
7:       for each search space  $H_i$  do
8:            $p(w) += \lambda_i p_i(w/H_i)$ 
9:       end for
10:  end for
11:  for each hyperedge  $e$  in  $U_k H_k$  do
12:      assign  $p(w)$  to the edge  $e$  for all  $w$  contained in  $e$ 
13:  end for
14:  return the best path according to the MMMBR decision rule
```



# Impacts of Beam Size

NIST 2008



# MBR Decoding on Single Model

- Decision Rule**

$$\begin{aligned}\hat{E} &= \operatorname{argmax}_{E' \in H} \sum_{E \in H} G(E, E') P(E | F) \\ &= \operatorname{argmax}_{E' \in H} \text{HGMBR}(H, E' | P)\end{aligned}$$

- Loss Function**

- $|E'|$  is hypothesis length
- $\#_{\omega}(E')$  is the number of times  $\omega$  occurs in  $E'$
- $\text{Cost}(E')$  is the model cost of the first-pass decoding
- $\theta_i$  are feature weights

$$\begin{aligned}\hat{E} &= \operatorname{argmax}_{E' \in H} \sum_{E \in H} \left\{ \theta_0 |E'| + \sum_{\omega} \theta_{|\omega|} \#_{\omega}(E') \delta_{\omega}(E') + \theta_m \text{Cost}(E') \right\} P(E | F) \\ &= \operatorname{argmax}_{E' \in H} \theta_0 |E'| + \sum_{\omega} \theta_{|\omega|} \#_{\omega}(E') p(\omega | H) + \theta_m \text{Cost}(E')\end{aligned}$$

$$p(\omega | H) = \sum_{E \in H} \delta_{\omega}(E) P(E | F)$$

# Algorithm for HGMBR

Step-1

- Sort the hypernodes topologically

Step-2

- Compute inside-outside probabilities of each hypernode

Step-3

- Compute posterior probabilities of each hyperedge

Step-4

- Compute posterior probabilities of each  $n$ -gram\*

Step-5

- Find the path with the highest score in hypergraph

# Inside-Outside Probability

- **Inside Probability**

- $w(h)$  is the weight of hyperedge  $h$ , given by  $\exp(\alpha x^T \phi(h))$ ,  $\alpha$  is the scaling factor
- $I(u)$  and  $O(u)$  refer to the incoming and outgoing hyperedges of hypernode  $u$

$$I(u) = \sum_{h \in \text{In}(u)} w(h) \left[ \prod_{v \in \text{tail}(h)} I(v) \right]$$

- **Outside Probability**

- For the root node of the hypergraph,  $O(\text{root}) = 1$
- $\text{head}(h)$  is the hypernode which  $h$  pointed to

$$O(u) = \sum_{h \in \text{Out}(u)} w(h) \left[ O(\text{head}(h)) \prod_{\substack{v \in \text{tail}(h) \\ v \neq u}} I(v) \right]$$

# Posterior Probability

- **Hyperedge Posterior Probability**

- $Z(f) = I(\text{root})$  is the inside probability of the root of the hypergraph

$$p(h \mid H) = \frac{1}{Z(f)} w(h) O(\text{head}(h)) \prod_{v \in \text{tail}(h)} I(v)$$

- ***n*-gram Posterior Probability**

- Approximate this quantity by the **sum** of posterior probabilities of edges which contribute the *n*-gram  $\omega$  and have the highest edge posterior probability relative to their predecessors on each path

$$p(\omega \mid H) = \sum_{e \in H} 1_{\omega \in e} f^*(e, \omega, H) p(e \mid H)$$

# Hypothesis Mixture Decoding

*ACL, 2011*

# Consensus-based Features

- n-gram posteriors based on **existing** hypothesis space

$$h_{H_n(f)}(e, f) = \sum_{\omega \in e} \#_{\omega}(e) p(\omega | H_n(f))$$

- n-gram posteriors based on **stemmed** hypothesis space

$$h_{H_n^S(f)}(e^S, f) = \sum_{\omega \in e^S} \#_{\omega}(e^S) p(\omega | H_n^S(f))$$

- n-gram posteriors based on **mixture** hypothesis space

$$h_{H(f)}(e, f) = \sum_{\omega \in e} \#_{\omega}(e) p(\omega | H(f))$$

- length posterior based on **mixture** hypothesis space

$$h_I(e, f) = \sum_{e' \in H(f), |e'|=|e|=I} p(e' | f)$$

# General Features

- LM score + word penalty
- count of lexicon pairs
- reorder features for (BTG/SCFG)-HMD respectively

$$h_{[\cdot]}(e, f) = \sum_{r \in D(e)} \delta_r([\cdot])$$

$$h_{\langle \cdot \rangle}(e, f) = \sum_{r \in D(e)} \delta_r(\langle \cdot \rangle)$$

$$h_{SCFG-Rule}(e, f) = \sum_{r \in D(e)} \delta_r(R)$$

$$h_{Glue-Rule}(e, f) = \sum_{r \in D(e)} \delta_r([\cdot])$$

- count of new generated n-grams

$$h_{New}(e, f) = \sum_{\omega \in e} \#_{\omega}(e) \bar{\delta}_{\omega}(\bigcup_{n=1}^N H_n(f))$$



# Background

# On-line SMT Engines

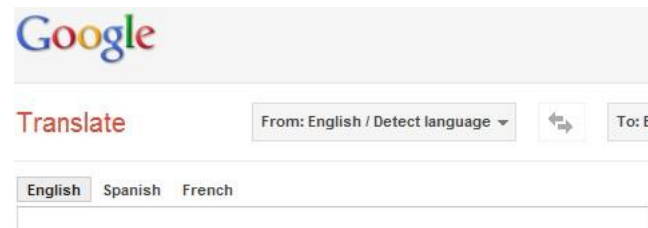
- Microsoft

- <http://dict.bing.com.cn>



- Google

- <http://translate.google.com>



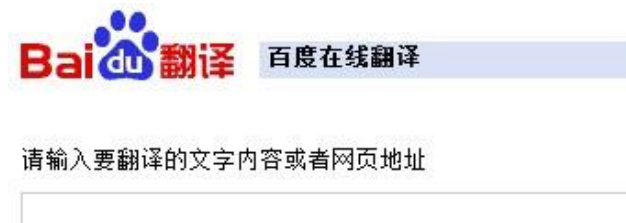
- Yahoo!

- <http://babelfish.yahoo.com/>



- Baidu

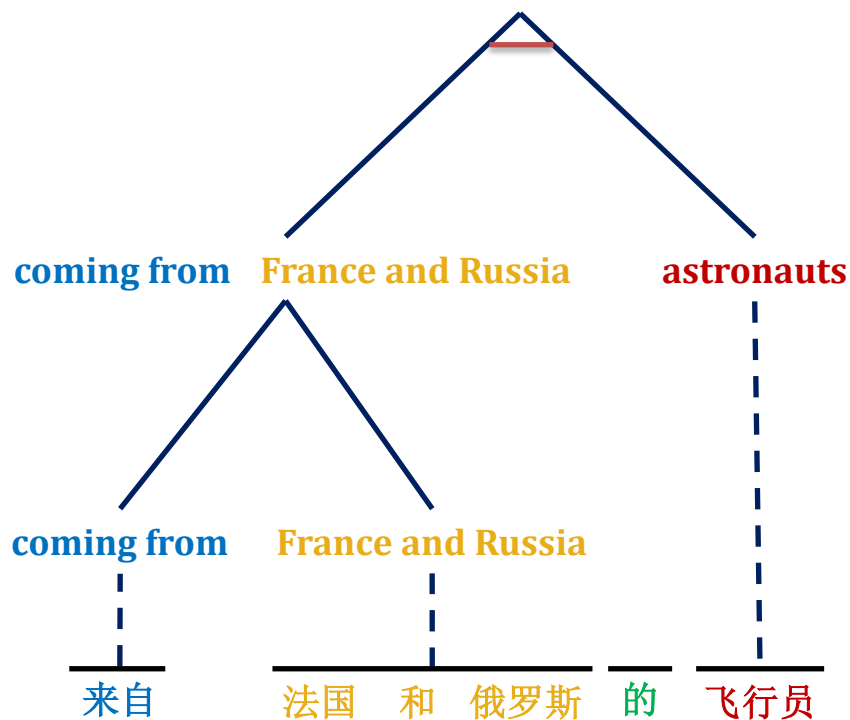
- <http://fanyi.baidu.com/>



# State-of-the-Art SMT Models

- Phrase-based Model

astronauts coming from France and Russia



$X_1 X_2 ||| X_2 X_1$

飞行员 ||| astronauts

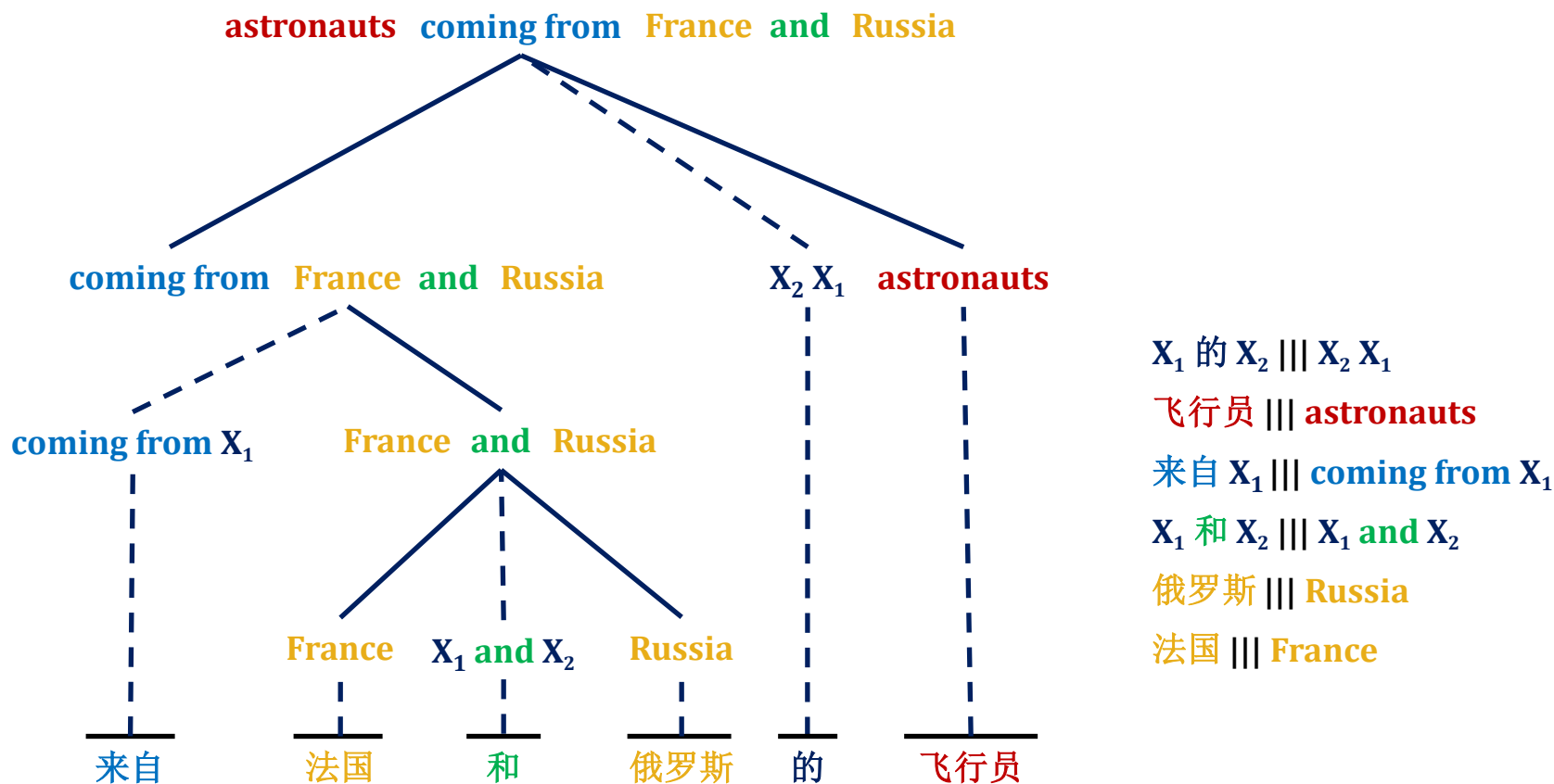
$X_1 X_2 ||| X_1 X_2$

法国 和 俄罗斯 ||| France and Russia

来自 ||| coming from

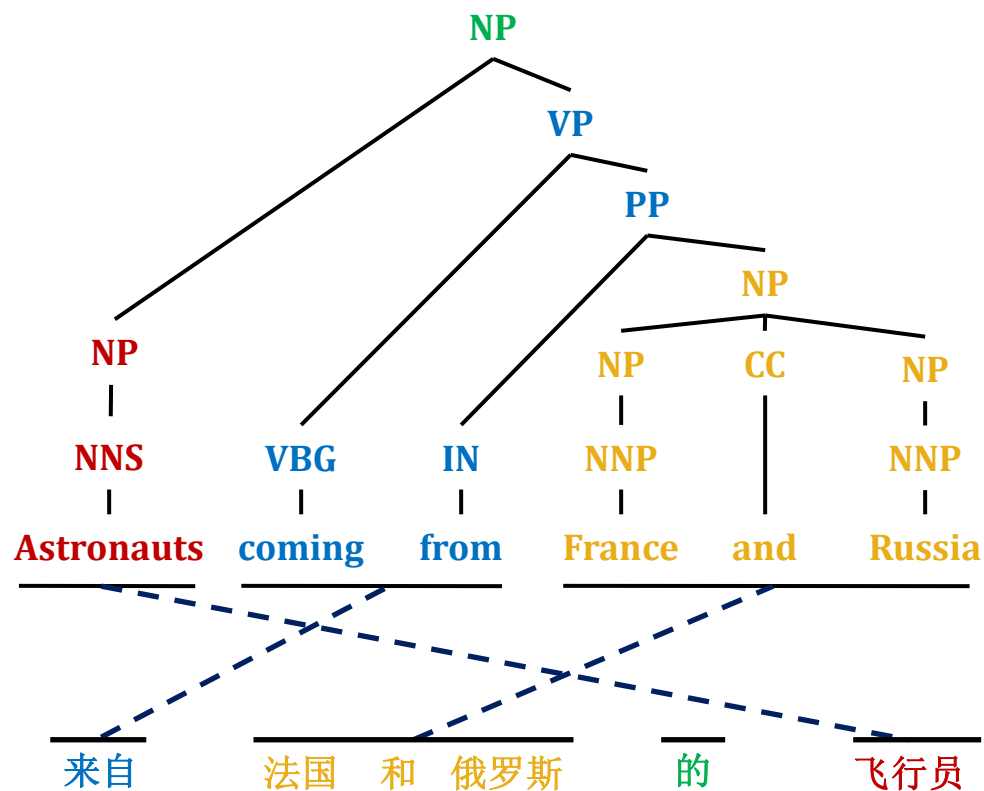
# State-of-the-Art SMT Models

- Hierarchical Phrase-based Model



# State-of-the-Art SMT Models

- Syntax-based Model



法国 和 俄罗斯 ||| NP{NNP{France}}

CC{and} NP{NNP{Russia}} ||| NP

来自  $X_1$  ||| VBG{coming} IN{from} NP: $X_1$   
||| VP

飞行员 ||| NNS{astronauts} ||| NP

$X_1$  的  $X_2$  ||| NP: $X_2$  VP: $X_1$  ||| NP

# n-best List & Hypergraph

n-best list

a cat on the mat

the mat a cat

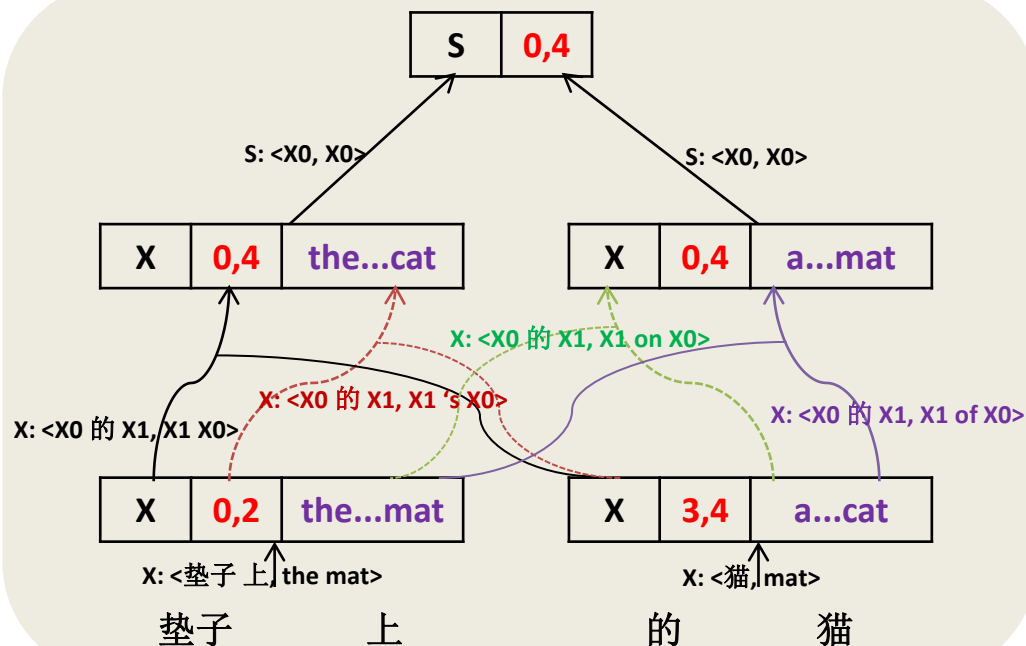
a cat of the mat

the mat 's a cat

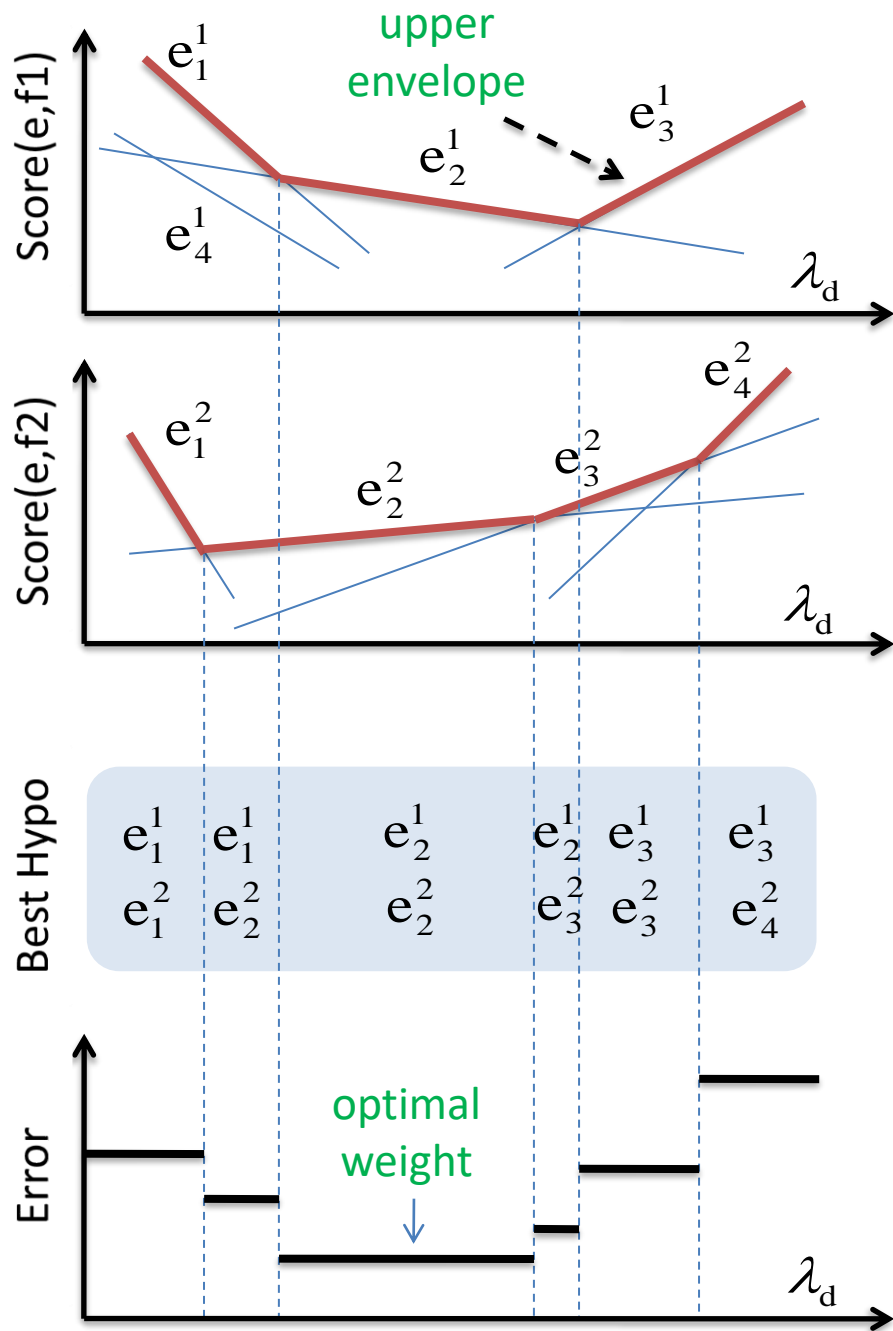


垫子 上 的 猫

Hypergraph



垫子 上 的 猫



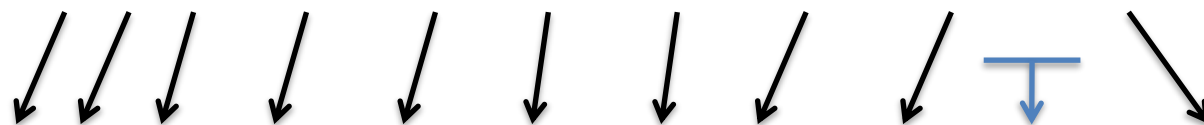
# IBM Model-4

这 7 人 包括 来自 法国 和 俄罗斯 的 宇航员



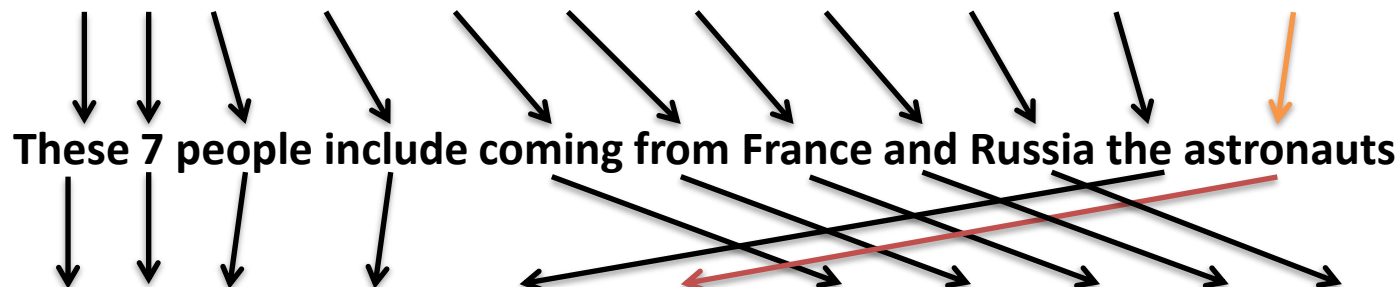
$n(2|\text{来自})$

这 7 人 包括 来自 来自 法国 和 俄罗斯 宇航员



$p\text{-null}$

这 7 人 包括 来自 来自 法国 和 俄罗斯 NULL 宇航员



$p(\text{astronauts}|\text{宇航员})$

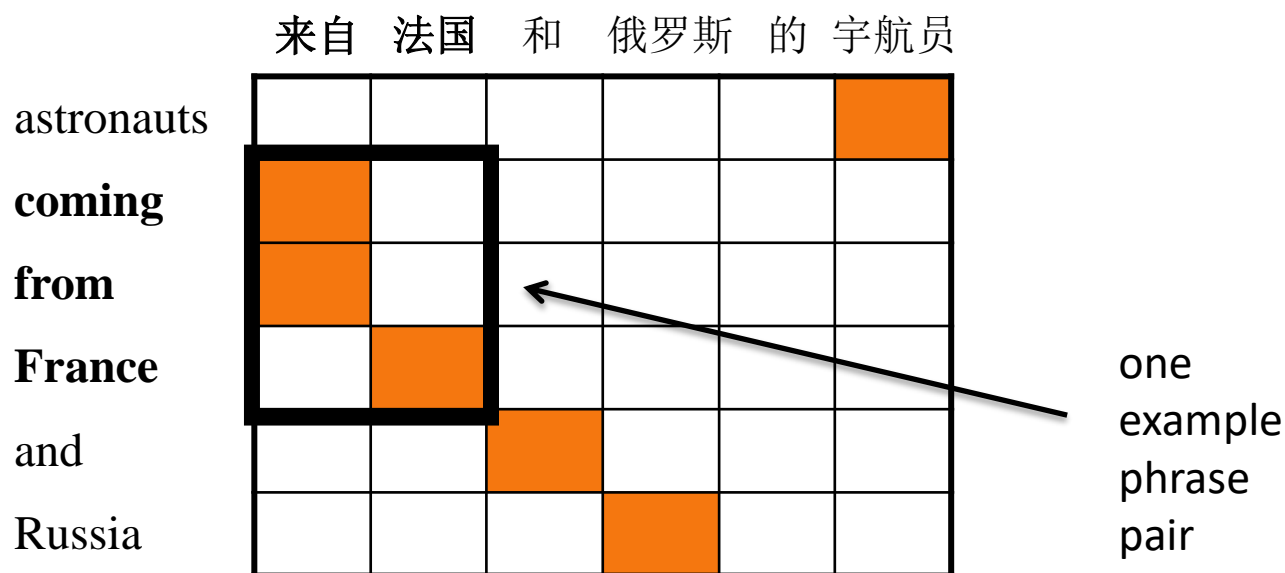
$d(5-4|\text{class}(\text{the}), \text{class}(\text{astronauts}))$

These 7 people include the astronauts coming from France and Russia

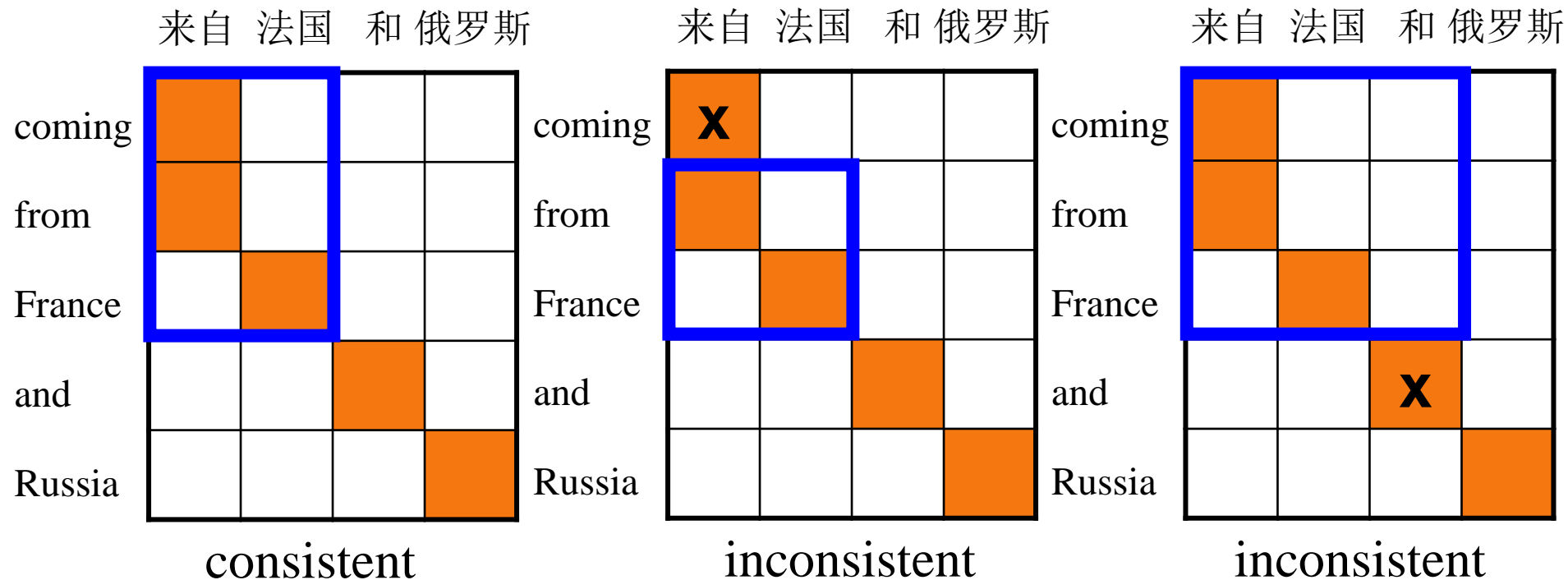


# How to Learn the Phrase Translation Table?

- Collect all phrase pairs *that are consistent with the word alignment*



# Consistent with Word Alignment



Phrase alignment must contain all alignment points for all the words in both phrases!

# Word Alignment Induced Phrases

来自 法国 和 俄罗斯 的 宇航员

astronauts						
coming						
from						
France						
and						
Russia						

(宇航员, astronauts) (来自, coming from) (法国, France) (和, and) (俄罗斯, Russia)

# Word Alignment Induced Phrases

来自 法国 和 俄罗斯 的 宇航员

astronauts						
coming						
from						
France						
and						
Russia						

(宇航员, astronauts) (来自, coming from) (法国, France) (和, and) (俄罗斯, Russia)  
(的 宇航员, astronauts) ...

# Word Alignment Induced Phrases

来自 法国 和 俄罗斯 的 宇航员

astronauts						
coming						
from						
France						
and						
Russia						

(宇航员, astronauts) (来自, coming from) (法国, France) (和, and) (俄罗斯, Russia)

(的 宇航员, astronauts) ...

(来自 法国, coming from France) (来自 法国 和, coming from France and)

(来自 法国 和 俄罗斯, coming from France and Russia) ...

# Word Alignment Induced Phrases

来自 法国 和 俄罗斯 的 宇航员

astronauts						
coming						
from						
France						
and						
Russia						

(宇航员, astronauts) (来自, coming from) (法国, France) (和, and) (俄罗斯, Russia)

(的 宇航员, astronauts) ...

(来自 法国, coming from France) (来自 法国 和, coming from France and)

(来自 法国 和 俄罗斯, coming from France and Russia) ...

(和 俄罗斯, and Russia) (法国 和 俄罗斯, France and Russia)

(法国 和 俄罗斯的, France and Russia) ...