



Ensemble Decoding for Statistical Machine Translation

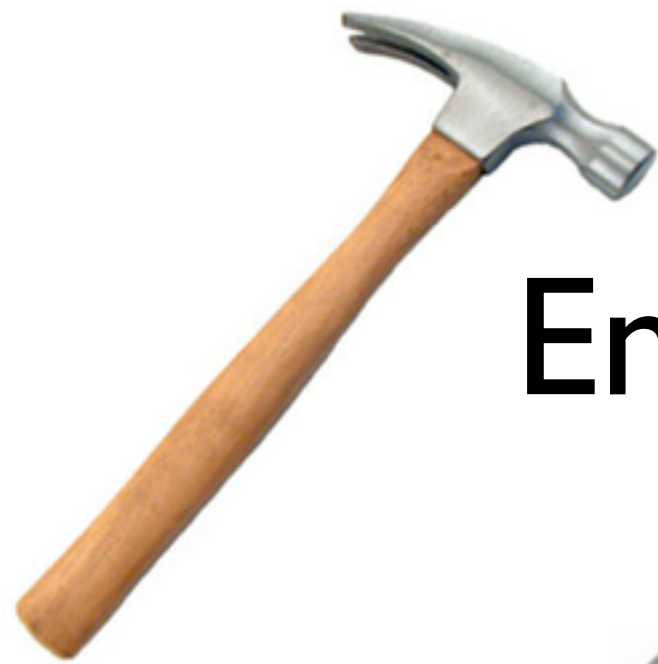
Anoop Sarkar

*Joint work with **Majid Razmara** and **Baskaran Sankaran***

SFU Natural Language Lab

Simon Fraser University, Vancouver, Canada

<http://natlang.cs.sfu.ca>



Ensemble Decoding



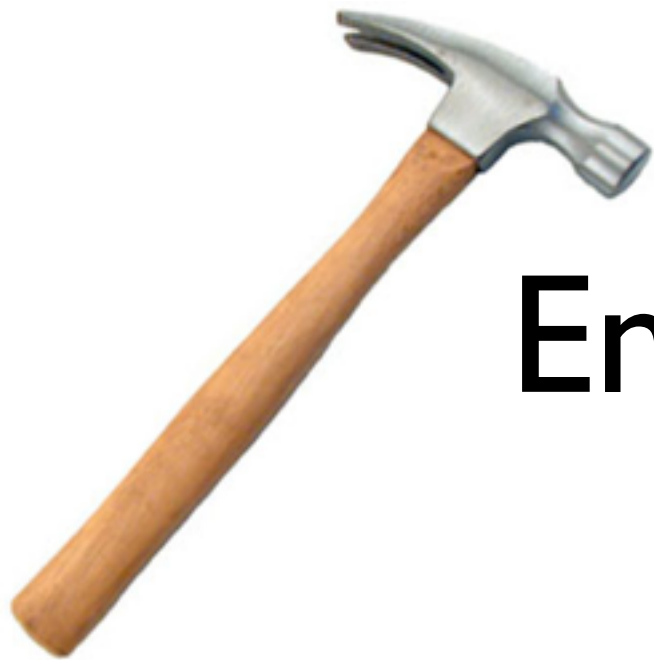
Domain Adaptation



Multi-metric optimization

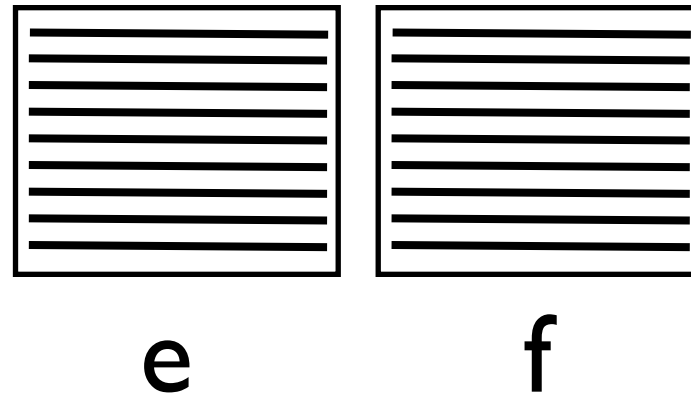
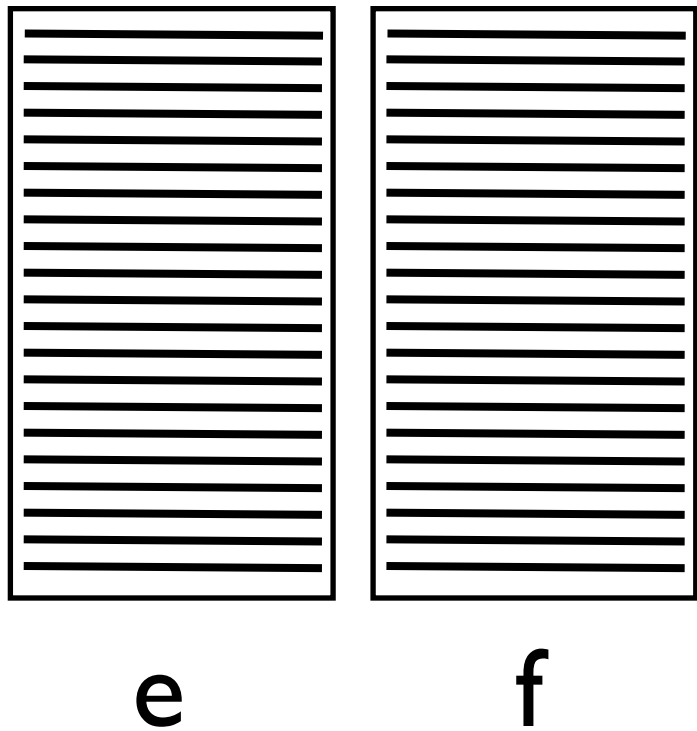


Pivot language triangulation



Ensemble Decoding

But first, Mixtures



Translation models: $m = 1 \dots M$

Log-linear mixture:
$$p(\bar{e}|\bar{f}) \propto \exp \left(\sum_m^M \lambda_m \log p_m(\bar{e}|\bar{f}) \right)$$

- Each (phrase-table) component in the usual discriminative SMT model is a mixture.
- The mixture weights are tuned on a dev set.

Linear Mixtures

$$p(\bar{e}|\bar{f}) = \sum_m^M \lambda_m p_m(\bar{e}|\bar{f})$$

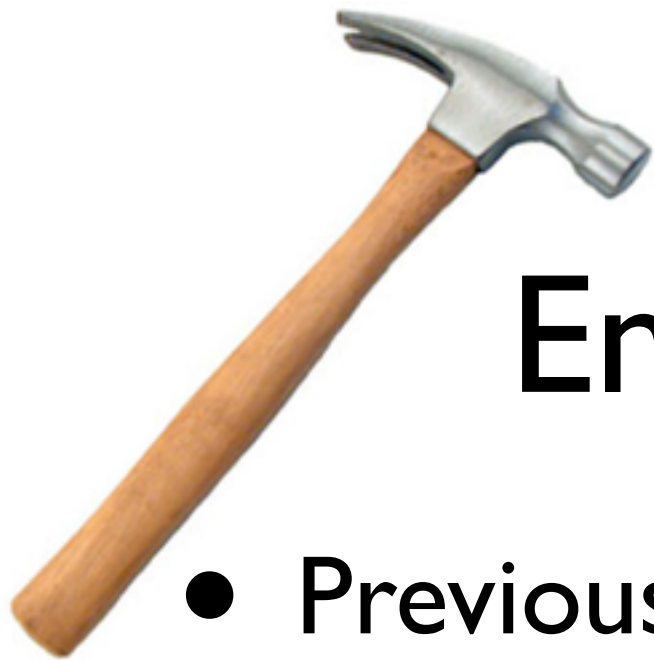
- Extract joint phrase pair distribution $p(e, f)$
- Find the weights that minimize the cross-entropy of the mixture $p(e | f)$ with respect to $p(e, f)$

$$\hat{\lambda} = \operatorname{argmax}_{\lambda} \sum_{\bar{e}, \bar{f}} \tilde{p}(\bar{e}, \bar{f}) \log \sum_m^M \lambda_m p_m(\bar{e}|\bar{f})$$

Linear Mixtures

$$\hat{\lambda} = \operatorname{argmax}_{\lambda} \sum_{\bar{e}, \bar{f}} \tilde{p}(\bar{e}, \bar{f}) \log \sum_m^M \lambda_m p_m(\bar{e} | \bar{f})$$

- Train the weights on the dev set using any optimization technique (L-BFGS).
- Linear mixtures are used as feature functions in standard discriminative SMT.
- State of the art for domain adaptation in SMT (Foster et al, EMNLP 2010).



Ensemble Decoding

- Previous mixtures of translation models were pre-processing steps.
- This work: Explore mixtures of translation models in the decoder.
- On the fly combination of models in Hiero

0.5 yu X_1 you X_2 /
have X_2 with X_1



1.5 Beihan /
North Korea



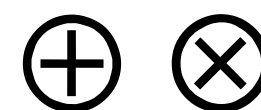
3.0 bangjiao /
diplomatic relations



1.5 yu X_1 you X_2 /
with X_1 have X_2



Semi-ring



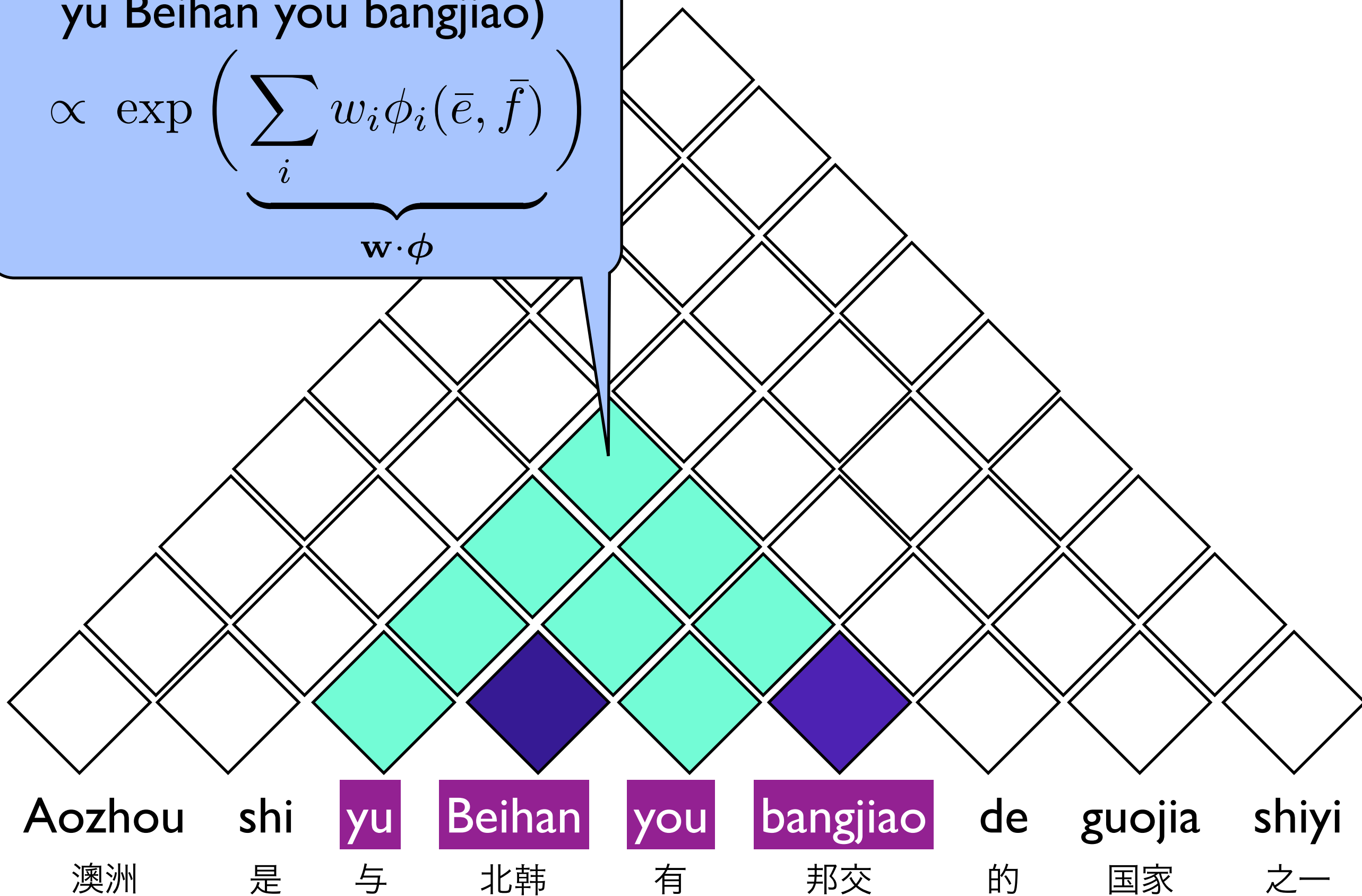
$(\mathbb{R} \cup \{\infty\}, \min, +, \infty, 0)$

1. have dipl. relns. with N. K = 4.0
2. with N.K. have dipl. relns. = 6.0
• $\min(4.0, 6.0)$ option #1 wins

Aozhou shi yu Beihan you bangjiao de guojia shiyi
澳洲 是 与 北韩 有 邦交 的 国家 之一

P(have dipl. relns. with N.K |
yu Beihan you bangjiao)

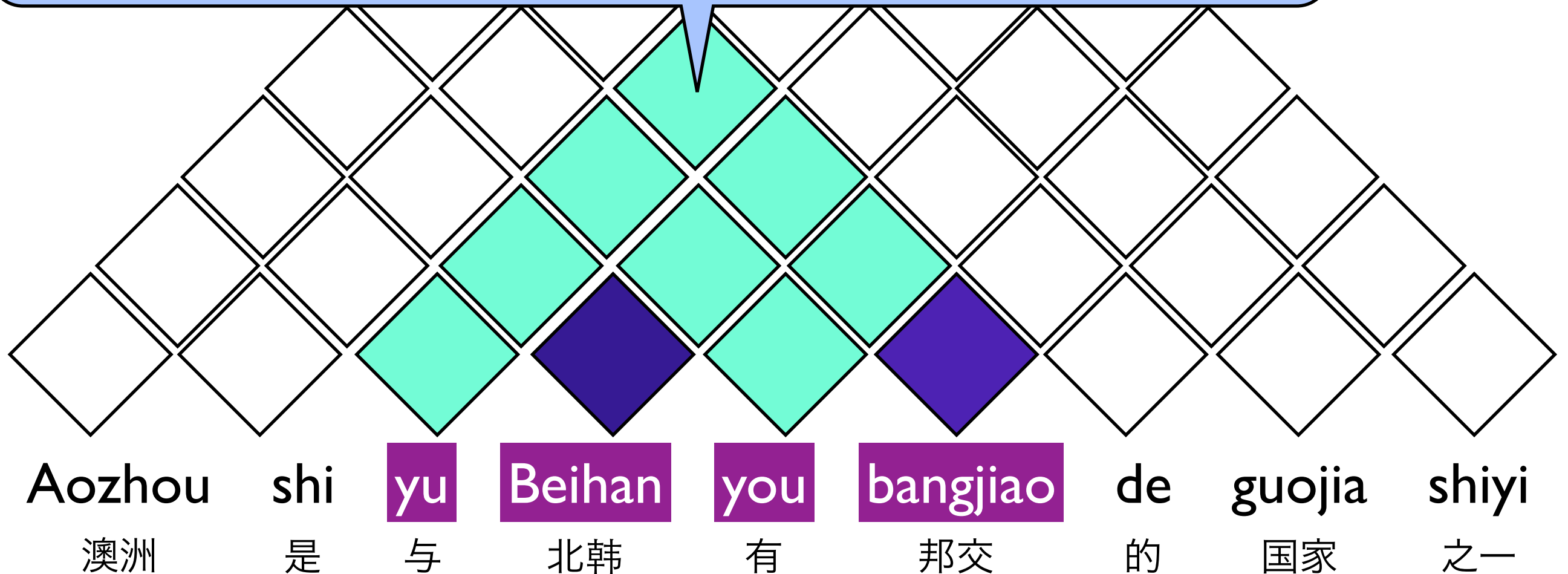
$$\propto \exp \left(\underbrace{\sum_i w_i \phi_i(\bar{e}, \bar{f})}_{\mathbf{w} \cdot \boldsymbol{\phi}} \right)$$



Ensemble Decoding

$P(\text{have dipl. relns. with N.K.} \mid \text{yu Beihan you bangjiao})$

$$\propto \exp \left(\underbrace{\mathbf{w}_1 \cdot \phi_1}_{1^{st} \text{ model}} \otimes \underbrace{\mathbf{w}_2 \cdot \phi_2}_{2^{nd} \text{ model}} \otimes \dots \right)$$



Ensemble Operations

Weighted Sum (wsum)

$$p(\bar{e} | \bar{f}) \propto \sum_m^M \lambda_m \exp(\mathbf{w}_m \cdot \phi_m)$$

- Ensemble score is the weighted sum of individual model scores
- m is each component model, a total of M components in the ensemble.

fr
en

une autre maladie métabolique héréditaire

or other disease hereditary metabolic

or another hereditary metabolic disease

m1	4.5
m2	10.5
ens	7.25

m1	16.5
m2	3.5
ens	10

Weighted Max (wmax)

$$p(\bar{e} \mid \bar{f}) \propto \max_m \left(\lambda_m \exp \left(\mathbf{w}_m \cdot \phi_m \right) \right)$$

- Ensemble score is the weighted max of all the model scores
- The n-best list can contain entries from different models

fr	une autre maladie métabolique héréditaire	m1	4.5
		m2	10.5
		ens	4.5
en	or other disease hereditary metabolic	m1	16.5
		m2	3.5
	or another hereditary metabolic disease	ens	3.5

Model Switching (Switch)

$$p(\bar{e} \mid \bar{f}) = \sum_m^M \delta(\bar{f}, m) p_m(\bar{e} \mid \bar{f})$$

- Switch models in each CKY cell. Possibly picking a different model from the ensemble.
- The n-best list can contain entries from only one model.

$$\delta(\bar{f}, m) = \begin{cases} 1, & m = \operatorname{argmax}_{n \in M} \psi(\bar{f}, n) \\ 0, & \text{otherwise} \end{cases}$$

Model Switching (Switch)

$$\delta(\bar{f}, m) = \begin{cases} 1, & m = \operatorname{argmax}_{n \in M} \psi(\bar{f}, n) \\ 0, & \text{otherwise} \end{cases}$$

For each cell the model that has the highest weighted score wins:

$$\psi(\bar{f}, n) = \lambda_n \max_{\bar{e}} (\mathbf{w}_n \cdot \phi_n(\bar{e}, \bar{\mathbf{f}}))$$

For each cell, the model with highest weighted sum of scores wins:

$$\psi(\bar{f}, n) = \lambda_n \sum_{\bar{e}} \exp (\mathbf{w}_n \cdot \phi_n(\bar{e}, \bar{\mathbf{f}}))$$

Product (prod)

$$p(\bar{e} | \bar{f}) \propto \exp \left(\sum_m^M \lambda_m (\mathbf{w}_m \cdot \phi_m) \right)$$

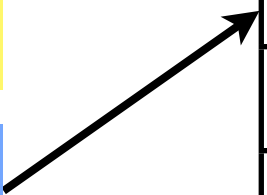
- Compute the product of all the probabilities in the ensemble (sum of log-probs).
- A Logarithmic Opinion Pool (LOP).
- LOPs work best when the ensemble is used to down-vote a highly confident but incorrect candidate.

fr

une autre maladie métabolique héréditaire

en

or other disease hereditary metabolic

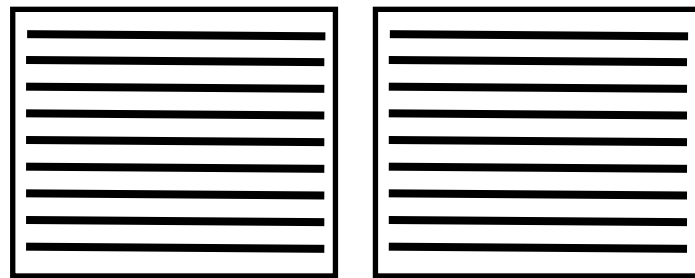


m1	1.5
m2	10.5
m3	12.0
ens	24.0



Domain Adaptation

Domain Adaptation



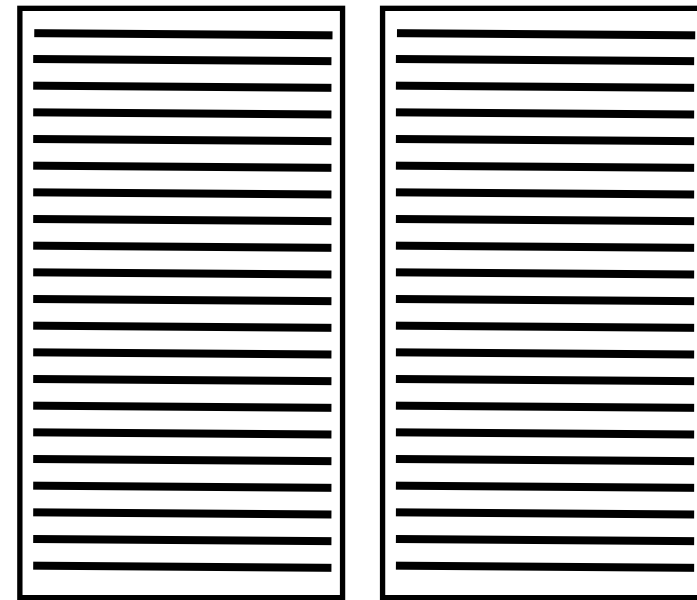
e

f

IN-domain

EMEA (Medical)

Train	11770
Dev	1533
Test	1522



e

f

OUT-of-domain

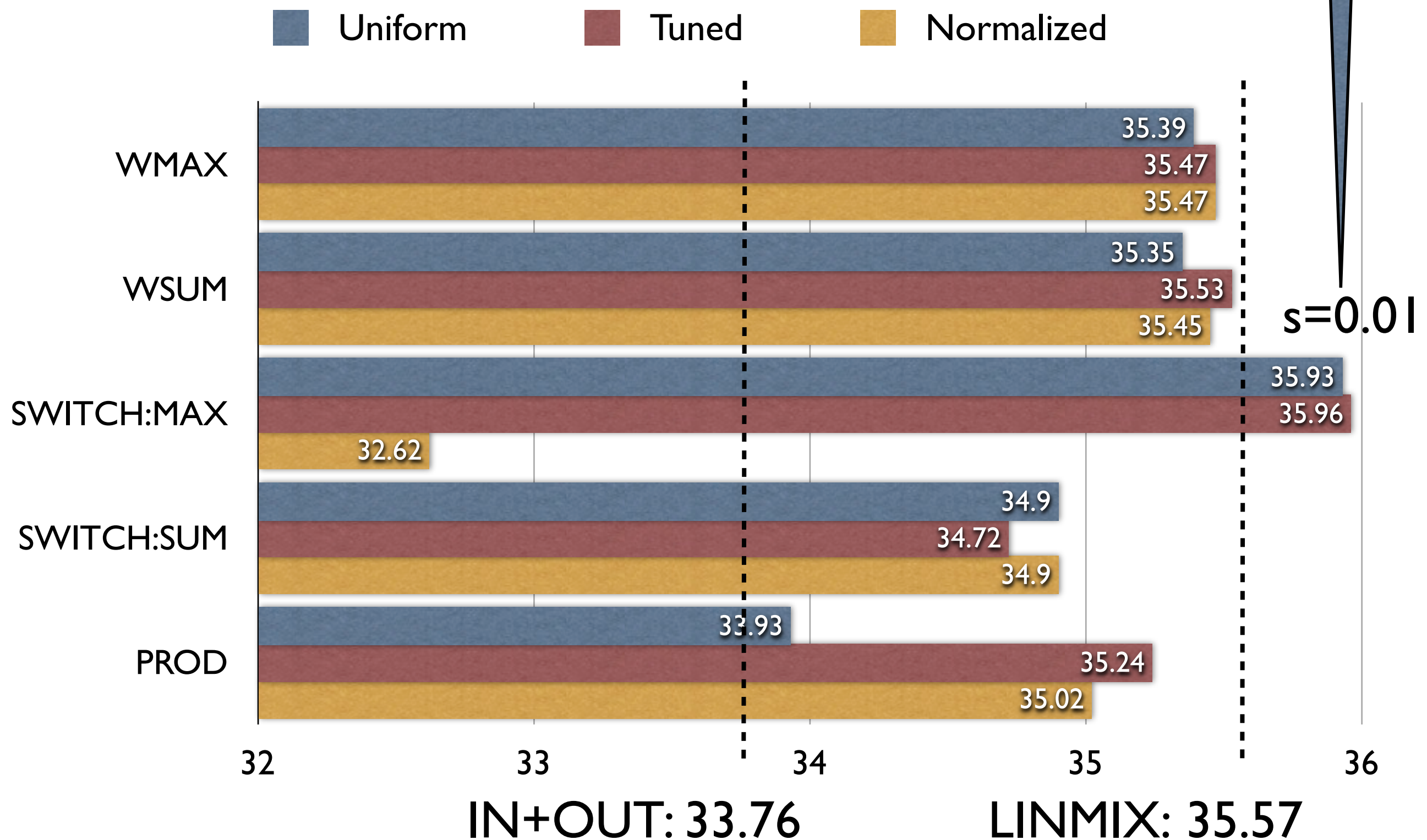
EuroParl (fr-en)

Train	1.3M
-------	------

Domain Adaptation

- Scaling the model scores using ensemble weights:
- Find the appropriate model scores that can participate in an ensemble.
- We use CONDOR (Vanden Berghen and Bersini, 2005) which uses Powell's algorithm and no gradient information.
- Component weights for each mixture operation is tuned on the dev set.

Domain Adaptation Results



Example

SOURCE	aménorrhée , menstruations irrégulières	
REF	amenorrhoea , irregular menstruation	
IN	amenorrhoea	, menstruations irrégulières
OUT	aménorrhée ,	irregular menstruation
ENSEMBLE	amenorrhoea	, irregular menstruation

Example

SOURCE	le traitement par naglazyme doit être supervisé par un médecin ayant l' expérience de la prise en charge des patients atteints de mps vi ou d' une autre maladie métabolique héréditaire .
REF	naglazyme treatment should be supervised by a physician experienced in the management of patients with mps vi or other inherited metabolic diseases .
IN	naglazyme treatment should be supervisé by a doctor the with in the management of patients with mps vi or other hereditary metabolic disease .
OUT	naglazyme 's treatment must be supervised by a doctor with the experience of the care of patients with mps vi. or another disease hereditary metabolic .
ENSEMBLE	naglazyme treatment should be supervised by a physician experienced in the management of patients with mps vi or other hereditary metabolic disease .



Multi-metric optimization

Joint work with Baskaran Sankaran and Kevin Duh

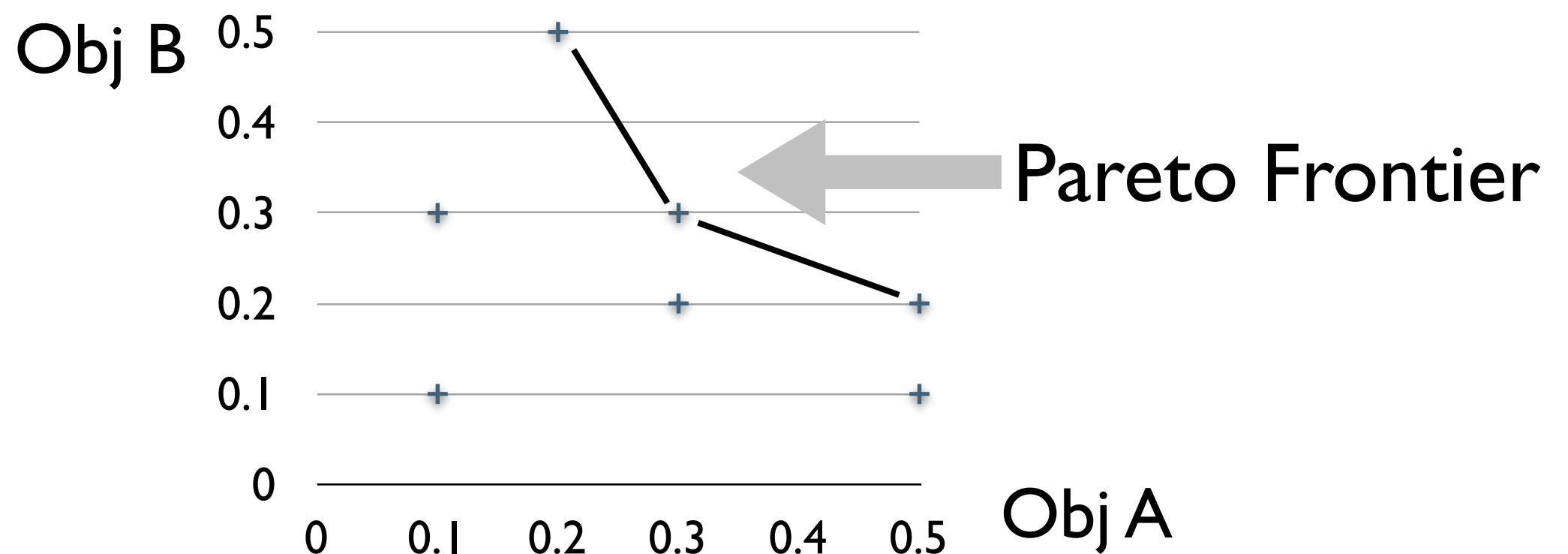
Multi-metric optimization

- Quite a few proposals for MT evaluation.
- In this talk, the focus is on BLEU, RIBES, TER, METEOR.
- Can other metrics be useful as a loss function for training SMT systems. (be useful how?)
- Most systems tune towards BLEU and test on BLEU. Can other metrics provide a second opinion?

Multi-objective optimization

$$\max_w (F_1(w), F_2(w), \dots, F_k(w))$$

- Find one w that simultaneously optimizes k objectives.
- A well formed notion of optimality wrt multiple objectives: Pareto optimality



Finding Pareto points

- Duh et al (ACL 2012) give an algorithm called PMO-PRO that finds Pareto optimal points as part of the tuning step.
- PRO (May and Hopkins, EMNLP 2011) show a pairwise ranking classifier can be used to train an SMT log-linear model.
- PMO-PRO puts Pareto points as positive examples and low scoring non-Pareto points as negative examples.
- This can be used to find Pareto points in the dev set.

Using the Pareto Points

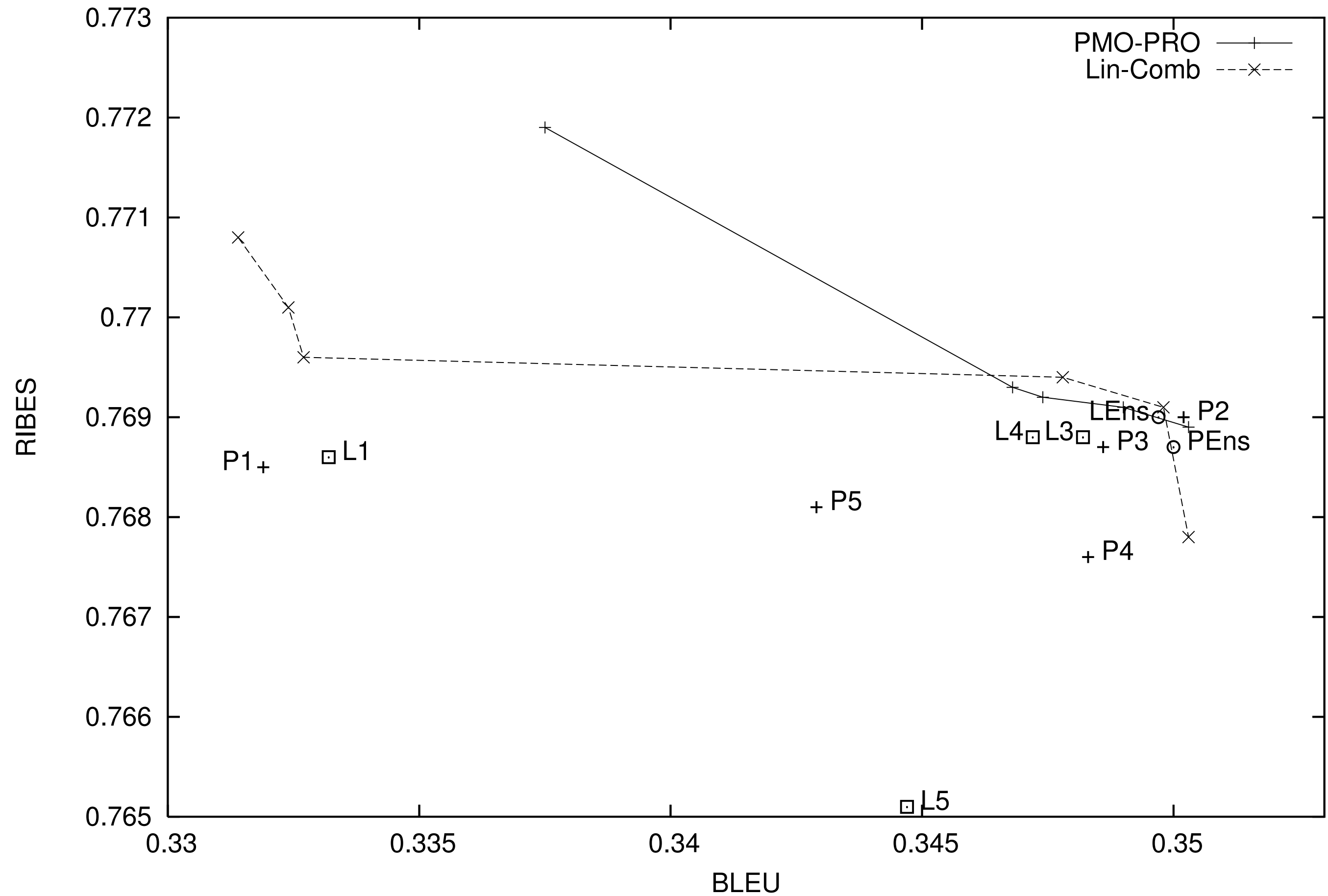
- Each Pareto point in the dev set is a weight vector that produced that point.
- **PMO-Ensemble**: Each of these weight vectors is a model and we can simply combine them using an ensemble model.
- **Union**: Take the union of “good” points wrt multiple objectives as positive examples and vice versa for negative examples. Simpler version of PMO-PRO.

Using the Pareto Points

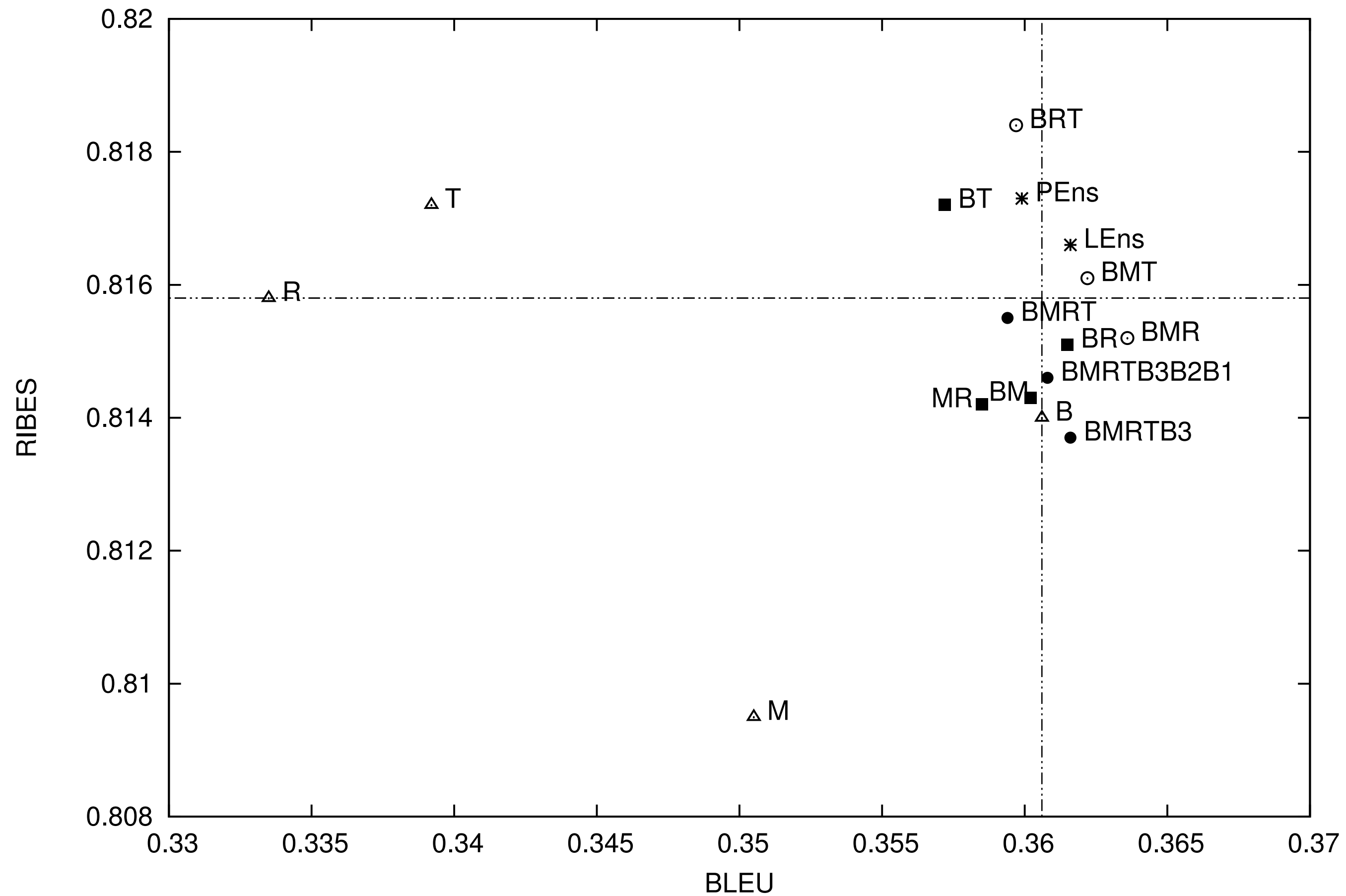
$$\max_w (F_1(w), F_2(w), \dots, F_k(w))$$

- **Ensemble Tuning:**
 - For each $F_i(w)$ perform error rate tuning (using PRO) to obtain the best w_i according to F_i in each iteration of tuning.
 - When decoding the dev set for the next search for w use an ensemble model with the same features but weights: w_1, \dots, w_k
 - Tune the ensemble model hyperparameters using PMO-PRO to get Pareto points in the ensemble.
 - Repeat.

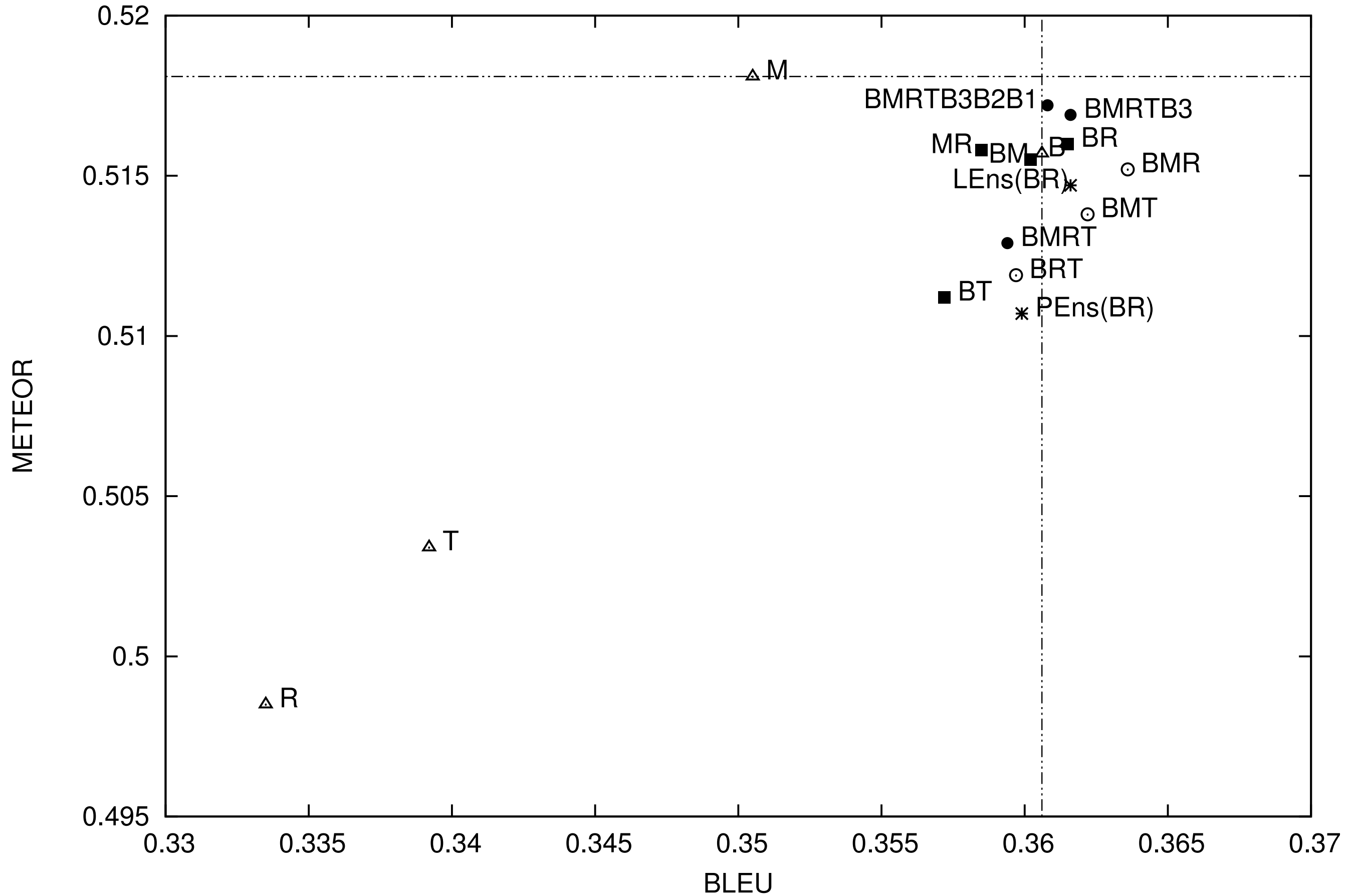
Ar-En: MTA-devset (redecode)



Ar-En: MTA-testset

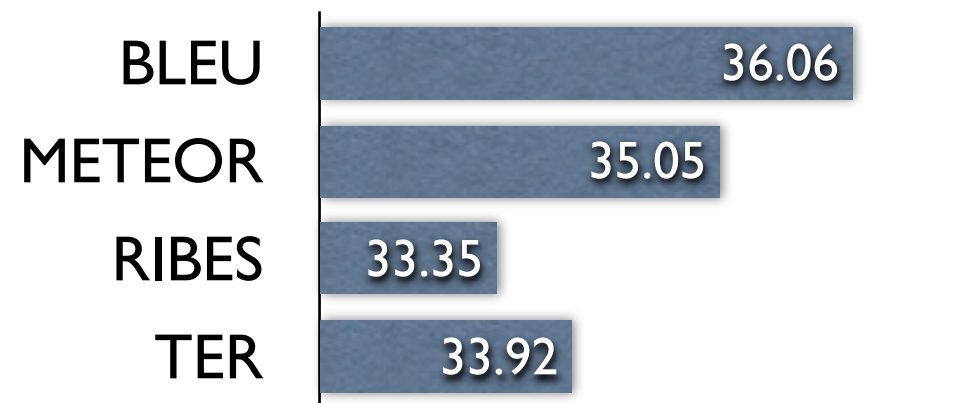


Ar-En: MTA-testset

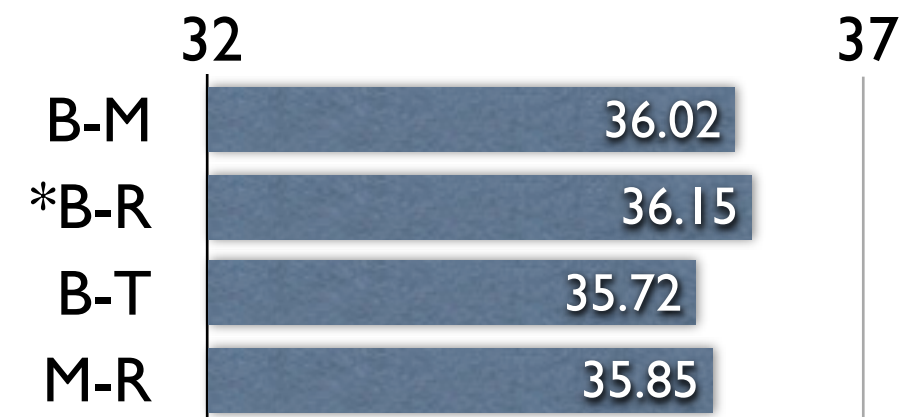


Can multi-metric tuning help a single metric?

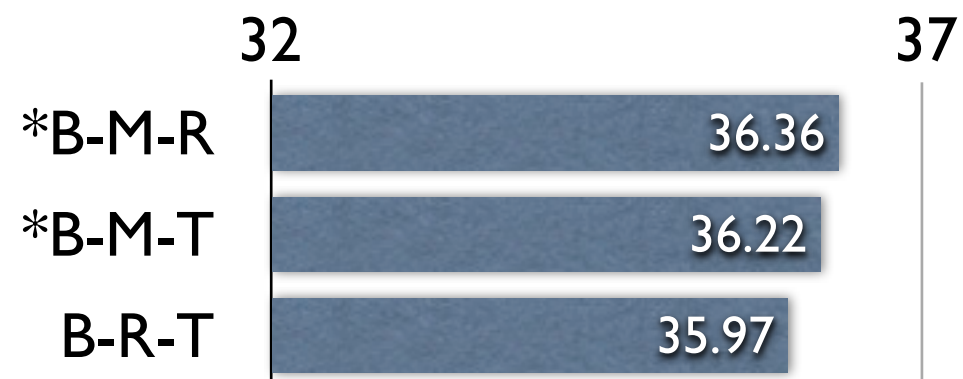
Single Objective



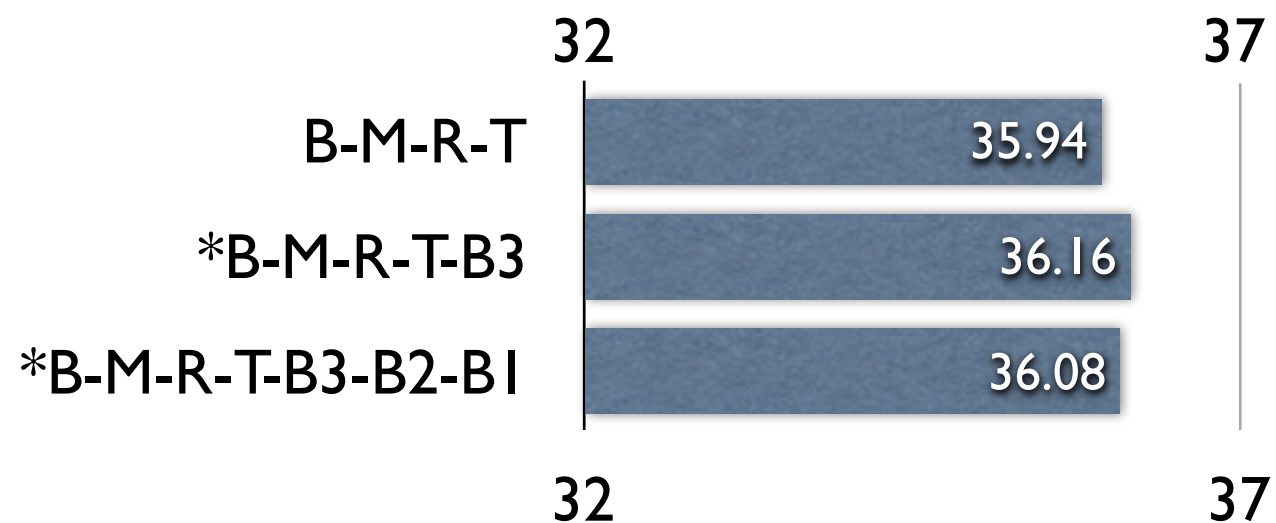
Ensemble Tuning: 2 Metrics



Ensemble Tuning: 3 Metrics



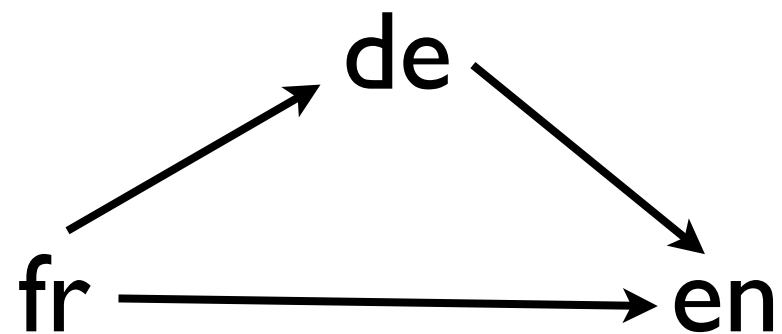
>3 Metrics



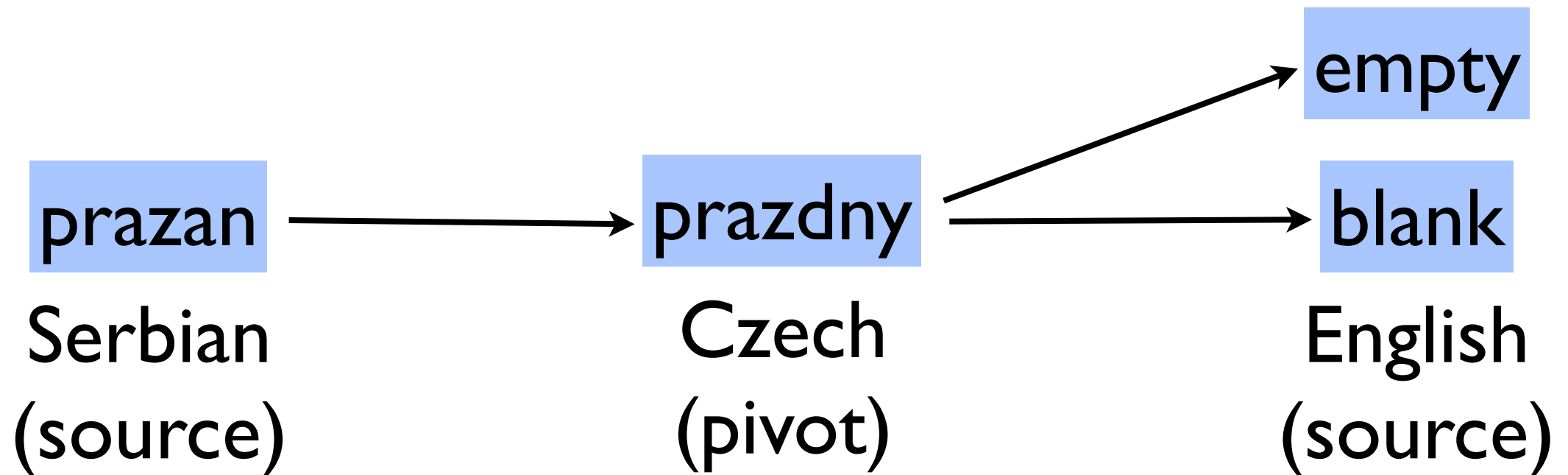
■ BLEU



Pivot language triangulation



Triangulation



- **Direct:** From source to target using available data
- Phrase-based triangulation (Cohn and Lapata, ACL 2007)

Phrase-based Triangulation

- Create a triangulated phrase table using source, pivot and target language data.

$$\begin{aligned} p(e|f) &= \sum_i p(e, i|f) \\ &= \sum_i p(e|i, f)p(i|f) \\ &\approx \sum_i p(e|i)p(i|f) \end{aligned}$$

- **Mixture:** Interpolate the triangulated model with the direct source to target model

Ensemble-based Triangulation

- Typically, one pivot language does not provide an improvement.
- More pivot languages used, the better.
- Ensemble-based Triangulation: an ensemble of different pivot models.
- Each one goes from source to pivot_k to target for pivot_1 to pivot_M
- **Ensemble:** Finally add the direct source to target model to the ensemble as well.

Experiment

- Compare Direct, Mixture and Ensemble
- Use EuroParl (en, fr, de, es, it)
- Each source language is translated to a target language through 3 pivot languages.
- For example, en to fr goes through de, es, it
- 10K sentence pairs (as in Cohn and Lapata, ACL 2007)
- To be done: 700K EuroParl corpus.

Ensemble Triangulation

	en	es	fr	de
en	-	+0.2	+1.0	-0.09
es	+0.24	-	+1.06	+0.38
fr	+0.9	+0.93	-	+0.03
de	+0.75	-0.48	+0.06	-

Comparison of Ensemble Model and Mixture model
(Cohn and Lapata, ACL 2007)

Summary

- Ensemble models combine translation models during SMT decoding.
 - Allows more dynamic combination methods.
 - Do not need to be tuned (with uniform weights)
- Applied to:
 - Domain adaptation
 - Multi-metric optimization
 - Pivot language triangulation

Direct Mixture Ensemble

