

Structured Ramp Loss Minimization for Machine Translation

Kevin Gimpel

Noah A. Smith



LEARNING IN MACHINE TRANSLATION

Why is learning in MT different from other tasks?

References often unreachable, so **surrogate references** are used instead (e.g., BLEU-oracles on k -best lists; Och & Ney, 2002)

How are learning algorithms affected?

Loss functions are changed:

Perceptron

Latent perceptron loss:

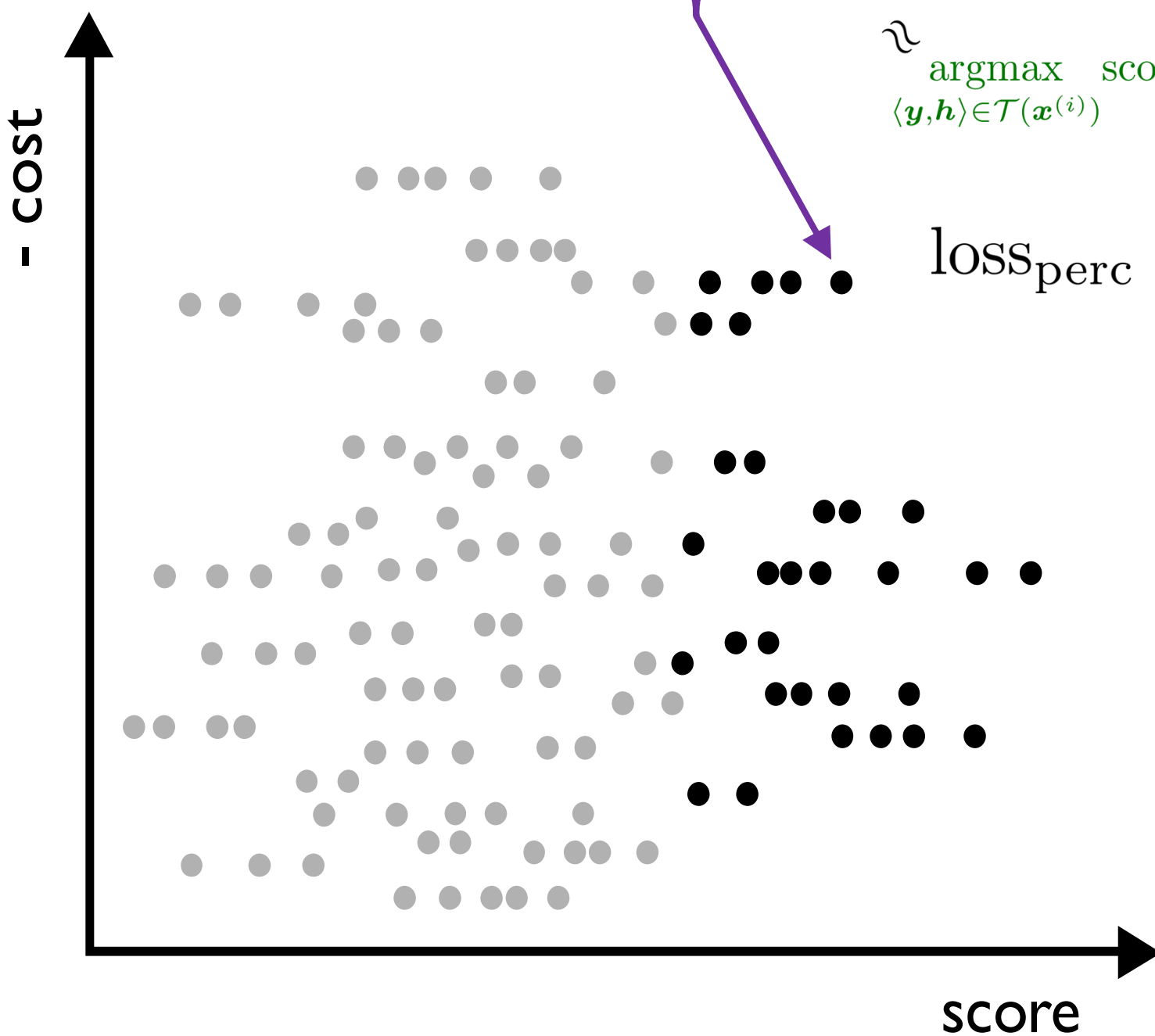
$$\text{loss}_{\text{perc}} = - \max_{\mathbf{h}: \langle \mathbf{y}^{(i)}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, \mathbf{h}) + \max_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h})$$

Latent perceptron with k -best BLEU oracle (Liang et al., 2006):

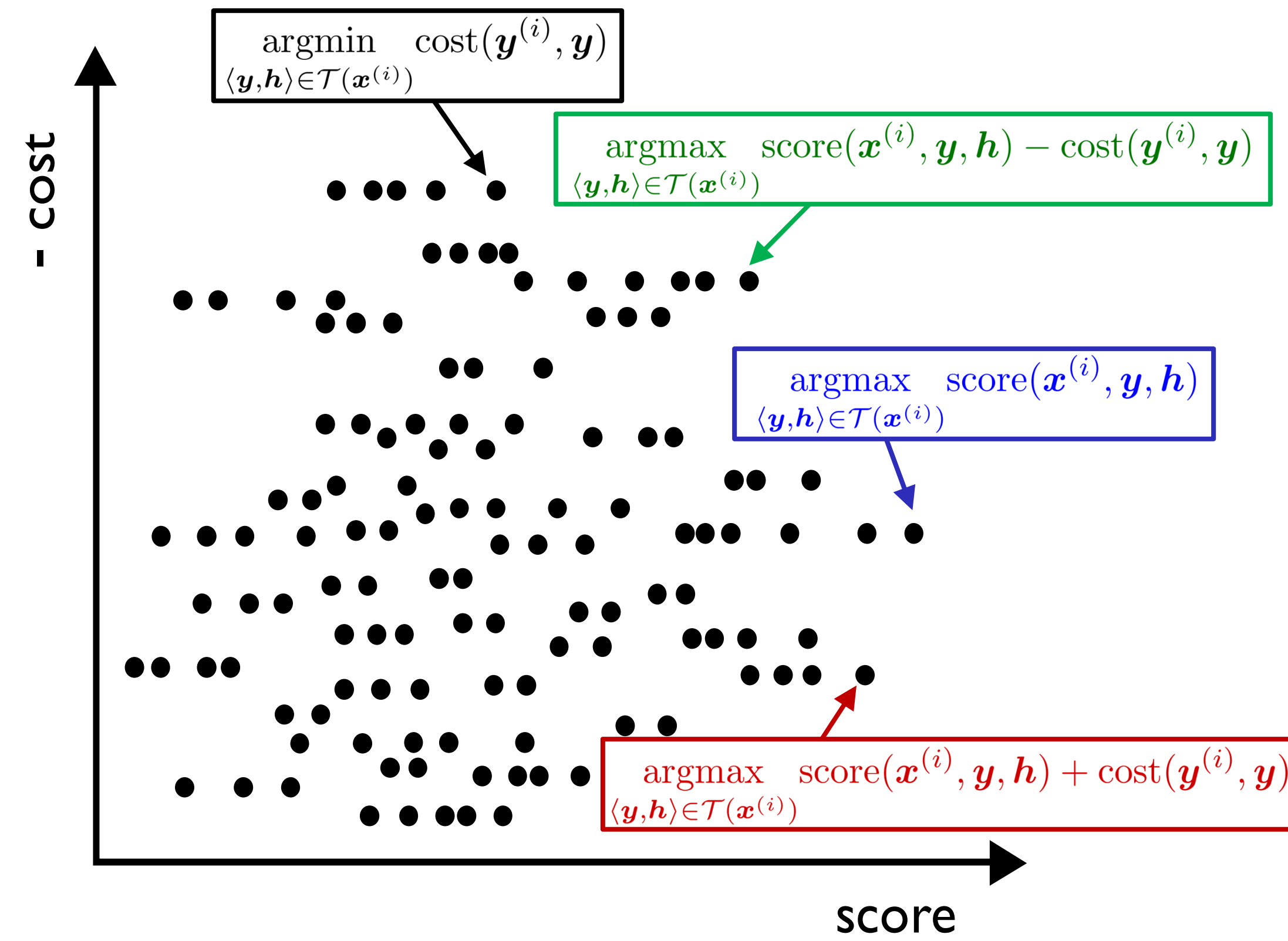
$$\text{loss}_{\text{perc } k\text{best}} = - \text{score} \left(\mathbf{x}^{(i)}, \underset{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{K}_i}{\text{argmin}} (\text{cost}(\mathbf{y}^{(i)}, \mathbf{y})) \right) + \max_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h})$$

$$\approx \underset{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}(\mathbf{x}^{(i)})}{\text{argmax}} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h}) - \text{cost}(\mathbf{y}^{(i)}, \mathbf{y})$$

$$\text{loss}_{\text{perc } k\text{best}} \approx \text{loss}_{\text{ramp } 2}$$



STRUCTURED RAMP LOSSES



$$\text{loss}_{\text{ramp } 1} = - \max_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h}) + \max_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h}) + \text{cost}(\mathbf{y}^{(i)}, \mathbf{y})$$

$$\text{loss}_{\text{ramp } 2} = - \max_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h}) - \text{cost}(\mathbf{y}^{(i)}, \mathbf{y}) + \max_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h})$$

$$\text{loss}_{\text{ramp } 3} = - \max_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h}) - \text{cost}(\mathbf{y}^{(i)}, \mathbf{y}) + \max_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h}) + \text{cost}(\mathbf{y}^{(i)}, \mathbf{y})$$

EXPERIMENTS

Small-Feature Experiments

Moses phrase-based MT system, 14 default features, default Moses initialization, 3 runs of each algorithm

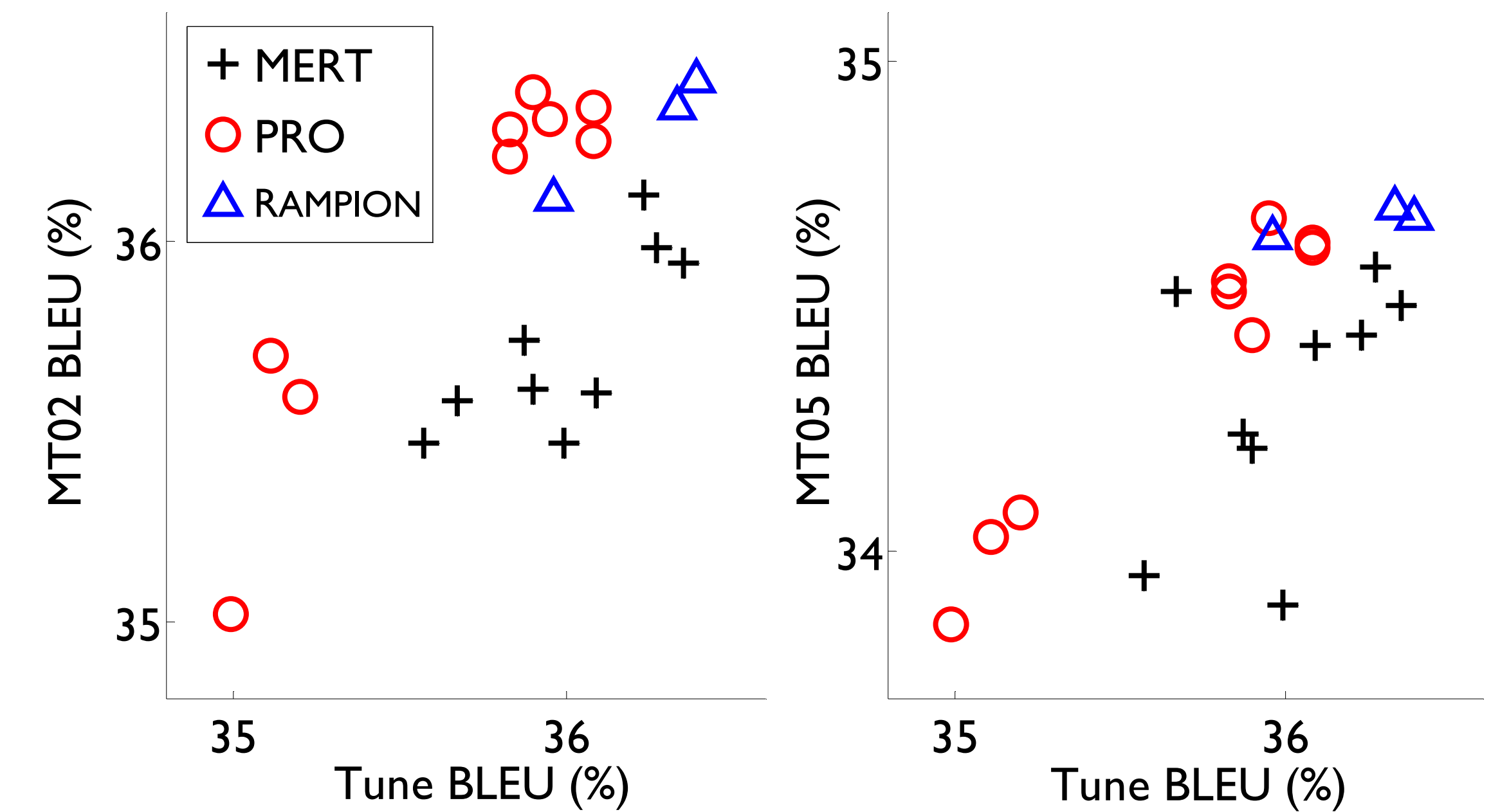
Method	Urdu-English		Chinese-English		Arabic-English			Avg
	MT08*	MT09	MT02	MT05	MT05	MT08 NW	MT08 WB	
MERT	24.5 ± 0.1	24.6 ± 0.0	35.7 ± 0.3	34.2 ± 0.2	55.0 ± 0.7	49.8 ± 0.3	32.6 ± 0.2	36.6
PRO	24.2 ± 0.1	24.2 ± 0.1	36.3 ± 0.1	34.5 ± 0.0	55.6 ± 0.3	49.6 ± 0.0	31.7 ± 0.0	36.6
RAMPION	24.5	24.6	36.4	34.7	55.5	49.8	32.1	36.8

Sensitivity Analysis

Chinese-English, 14 default Moses features

3 initializers (default Moses initialization + 2 random initializers)

3 runs of each algorithm for each initializer



Large-Feature Experiments

14 default Moses features + 7,200 additional features:

1k most frequent bilingual word pairs

200 most frequent unigrams, 1k most frequent bigrams, 1k most frequent trigrams

4k top trigger pairs, ranked by mutual information (Rosenfeld, 1996)

Method	Urdu-English			Chinese-English		
	Tune	MT08*	MT09	Tune	MT02	MT05
PRO	29.4	22.3	23.0	40.9	35.7	33.6
RAMPION	27.8	24.2	24.6	38.8	36.2	34.3

Code available for Moses: www.ark.cs.cmu.edu/MT

Also in cdec: www.cdec-decoder.org

MIRA

Latent hinge loss:

$$\text{loss}_{\text{hinge}} = - \max_{\mathbf{h}: \langle \mathbf{y}^{(i)}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, \mathbf{h}) + \max_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h}) + \text{cost}(\mathbf{y}^{(i)}, \mathbf{y})$$

MIRA for MT (Chiang et al., 2008; 2009):

$$\text{loss}_{\text{mira MT}} \approx - \max_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h}) - \text{cost}(\mathbf{y}^{(i)}, \mathbf{y}) + \max_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h}) + \text{cost}(\mathbf{y}^{(i)}, \mathbf{y})$$

$$\text{loss}_{\text{mira MT}} \approx \text{loss}_{\text{ramp } 3}$$

CRF

Latent log loss:

$$\text{loss}_{\text{log}} = - \log \sum_{\mathbf{h}: \langle \mathbf{y}^{(i)}, \mathbf{h} \rangle \in \mathcal{T}_i} \exp \{ \text{score}(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, \mathbf{h}) \} + \log \sum_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \exp \{ \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h}) \}$$

Log loss for MT (Och & Ney, 2002):

$$\text{loss}_{\text{log MT}} \approx - \max_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h}) - \text{cost}(\mathbf{y}^{(i)}, \mathbf{y}) + \log \sum_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_i} \exp \{ \text{score}(\mathbf{x}^{(i)}, \mathbf{y}, \mathbf{h}) \}$$

$$\text{loss}_{\text{log MT}} \approx \text{loss}_{\text{soft ramp } 2}$$

ALGORITHM

We optimize $\text{loss}_{\text{ramp } 3}$ using a concave-convex procedure (CCCP; Yuille & Rangarajan, 2002). CCCP is a batch optimization algorithm for a sum of a concave and a convex function.

```

Input: inputs  $\{\mathbf{x}^{(i)}\}_{i=1}^N$ , references  $\{\mathbf{y}^{(i)}\}_{i=1}^N$ , initial weights  $\theta_0$ , # iterations  $T$ , # CCCP iterations  $T'$ ,
# SSD iterations  $T''$ ,  $k$ -best list size  $k$ , step size  $\eta$ ,  $\ell_2$  regularization coefficient  $C$ 
Output: learned weights:  $\theta$ 
 $\theta \leftarrow \theta_0$ ;
for iter  $\leftarrow 1$  to  $T$  do // run Rampion for  $T$  iterations
   $\{\mathcal{K}_i\}_{i=1}^N \leftarrow \text{Decode}(\{\mathbf{x}^{(i)}\}_{i=1}^N, \theta, k)$ ; // get  $k$ -best lists
  for iter'  $\leftarrow 1$  to  $T'$  do // run CCCP (concave-convex procedure) for  $T'$  iterations
    for  $i \leftarrow 1$  to  $N$  do // impute "hope" translations for all sentences
       $(\mathbf{y}_i^+, \mathbf{h}_i^+) \leftarrow \text{argmax}_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{K}_i} \text{score}_i(\mathbf{y}, \mathbf{h}; \theta) - \text{cost}_i(\mathbf{y})$ ;
    end
    for iter''  $\leftarrow 1$  to  $T''$  do // run stochastic subgradient descent for  $T''$  epochs
      for  $i \leftarrow 1$  to  $N$  do
         $(\mathbf{y}^-, \mathbf{h}^-) \leftarrow \text{argmax}_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{K}_i} \text{score}_i(\mathbf{y}, \mathbf{h}; \theta) + \text{cost}_i(\mathbf{y})$ ; // find "fear" translation
         $\theta^- = \eta C \left( \frac{\theta - \theta^-}{N} \right)$ ; // regularize back to  $\theta_0$ , not 0
         $\theta^+ = \eta (f(\mathbf{x}^{(i)}, \mathbf{y}_i^+, \mathbf{h}_i^+) - f(\mathbf{x}^{(i)}, \mathbf{y}^-, \mathbf{h}^-))$ ; // subgradient update for loss
      end
    end
  end
end
return  $\theta$ ; // return final parameters (no averaging needed)

```



Algorithm 1. RAMPION *Campanula rapunculoides*

"A hardy biennial, cultivated for the use of its fleshy roots in salads, either boiled or in a raw state, generally the latter; the leaves are also used in winter salads" (Nicholson, 1884)

REFERENCES

- D. Chiang, Y. Marton, and P. Resnik. (2008) "Online large-margin training of syntactic and structural translation features." EMNLP.
- D. Chiang, W. Wang, and K. Knight. (2009) "11,001 new features for statistical machine translation." NAACL.
- P. Liang, A. Bouchard-Cote, D. Klein, and B. Taskar. (2006) "An end-to-end discriminative approach to machine translation." ACL.
- G. Nicholson. (1884) *The Illustrated Dictionary of Gardening, Div. VI*. London: L. Upcott Gill.
- F. Och and H. Ney. (2002) "Discriminative training and maximum entropy models for statistical machine translation." ACL.
- R. Rosenfeld. (1996) "A maximum entropy approach to adaptive statistical language modeling." *Computer, Speech, and Language*, 10(3).
- A. Yuille and A. Rangarajan. (2002) "The concave-convex procedure (CCCP)." NIPS.