

# A Systematic Exploration of Diversity in Machine Translation

Kevin Gimpel

Toyota Technological Institute at Chicago

Dhruv Batra  
Virginia Tech

Chris Dyer  
Carnegie Mellon University

Gregory Shakhnarovich  
Toyota Technological Institute at Chicago

## Motivation

In structured prediction, M-best lists are not diverse  
Packed representations partially address this, but add complexity  
Can we inject **diversity** into M-best lists for machine translation?



## Generating diverse translations

$$\langle \hat{y}, \hat{h} \rangle = \operatorname{argmax}_{\langle y, h \rangle \in \mathcal{J}_x} \mathbf{w}^\top \phi(x, y, h)$$

predicted translation    parameters    feature vector  
latent derivation    Input sentence

We use the framework of Batra et al. (2012):

$$\begin{aligned} \text{best: } \langle y^1, h^1 \rangle &= \operatorname{argmax}_{\langle y, h \rangle \in \mathcal{J}_x} \mathbf{w}^\top \phi(x, y, h) \\ \text{diverse 2}^{\text{nd}}\text{-best: } \langle y^2, h^2 \rangle &= \operatorname{argmax}_{\langle y, h \rangle \in \mathcal{J}_x} \mathbf{w}^\top \phi(x, y, h) + \lambda_1 \Delta(y^1, y) \\ &\dots \\ \text{diverse } m\text{th}\text{-best: } \langle y^m, h^m \rangle &= \operatorname{argmax}_{\langle y, h \rangle \in \mathcal{J}_x} \mathbf{w}^\top \phi(x, y, h) + \sum_{j=1}^{m-1} \lambda_j \Delta(y^j, y) \end{aligned}$$

“dissimilarity weight”  
“dissimilarity function”  
compute dissimilarity to all previous translations

Following Batra et al., we set  $\lambda = \lambda_j, \forall j$   
 $\lambda$  is tuned via grid search for each task

## Dissimilarity functions for MT

Count all  $n$ -gram matches in the two translations:

$$\Delta_n(y, y') = - \sum_{i=1}^{|y|-q} \sum_{j=1}^{|y'|-q} \mathbb{I}[y_{i:i+q} = y'_{j:j+q}]$$

$q = n - 1$   
Iverson bracket

$\lambda$  and  $n$  are tuned to maximize oracle BLEU of diverse lists

Can be implemented as another language model; no change to decoder is required

## Experimental setup

| language pair     | model                     | # sentence pairs |
|-------------------|---------------------------|------------------|
| Arabic → English  | phrase-based              | 4.3M             |
| Chinese → English | hierarchical phrase-based | 300K             |
| German → English  | phrase-based              | 1.9M             |

Moses used for all language pairs  
5-gram LMs with an extra 600M tokens from Gigaword

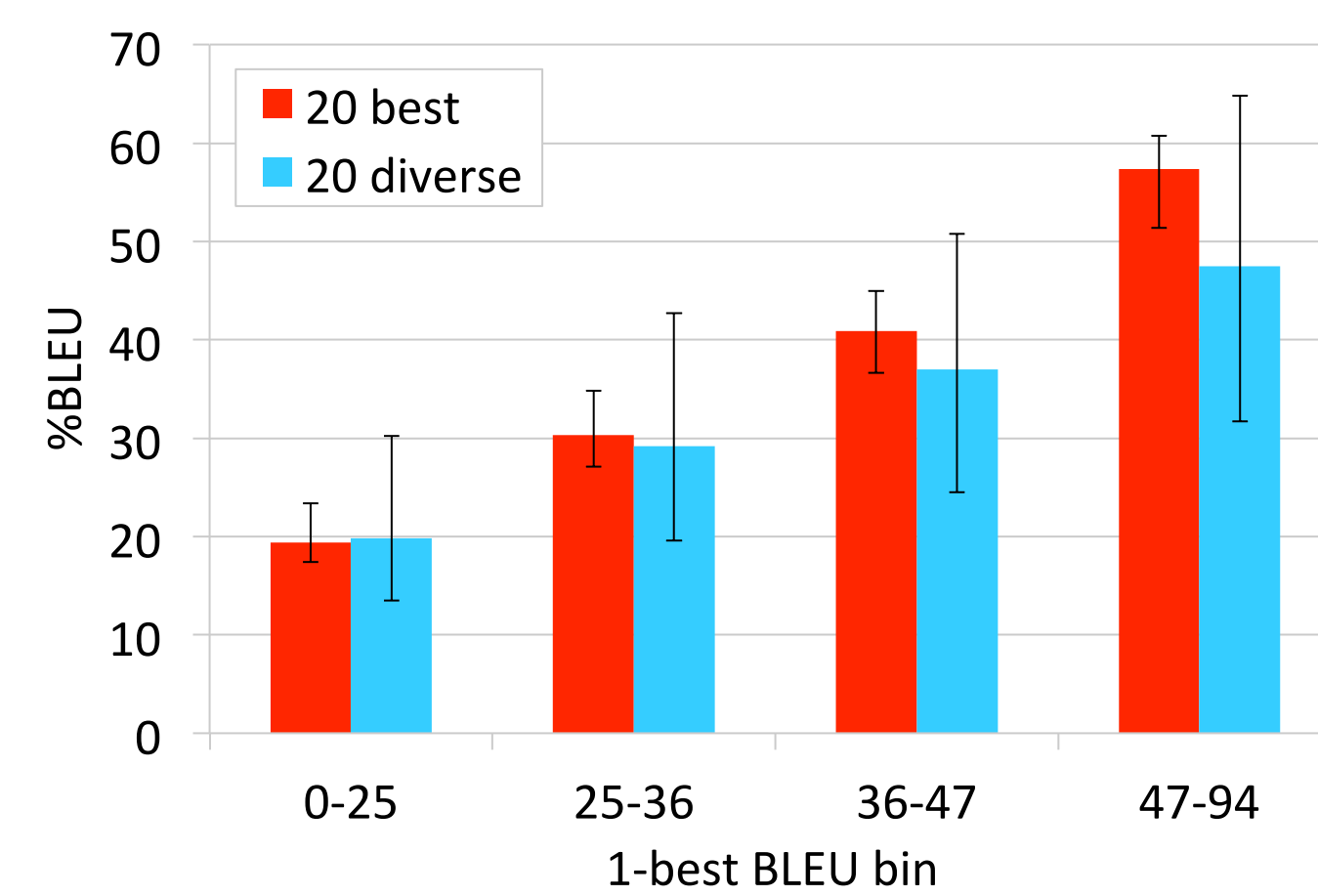
## Examples of diverse translations

|                     | reference:   | The government wants to limit the torture of "witches," a brochure was released |
|---------------------|--|---|
| 1-best translation: | The government wants the torture of 'witch' and gave out a booklet   |   |
| unique 5-best list: | The government wants the torture of "witch" and gave out a booklet<br>The government wants the torture of 'witch' and gave out a brochure<br>The government wants the torture of 'witch' and gave out a leaflet<br>The government wants the torture of "witch" and gave out a brochure   |   |
| diverse list:       | The government wants to <b>stop torture</b> of "witch" and issued a leaflet issued<br>The government wants to <b>"stop the torture</b> of" witches and gave out a brochure<br>The government intends to the torture of "witchcraft" and were issued a leaflet<br>The government is <b>the torture of "witches" stamp out</b> and gave a brochure |   |

|                     |   |
|---------------------|---|
| references:         | Kanazawa indicated as Akihito did not bleed much during the operation, the doctors used his own blood for transfusion. Kanazawa said Akihito did not lose a lot of blood during the surgery. Actually, doctors had only used the blood drawn from him for transfusions.   |
| 1-best:             | Kanazawa College said Akihito bleeding during surgery, doctors, not only in his own blood transfusion.  |
| unique 5-best list: | Kanazawa College said Akihito bleeding during surgery, doctors, not only with his own blood transfusion.<br>Kanazawa College, said Akihito bleeding during surgery, doctors, not only in his own blood transfusion.<br>Kanazawa College said that Akihito bleeding during surgery, doctors, not only in his own blood transfusion.<br>Kanazawa College said that Akihito bleeding during surgery, doctors, not only with his own blood transfusion.                                 |
| diverse list:       | Kanazawa College, said <b>Akihito did not hemorrhage during surgery</b> , the doctor <b>only used his own blood</b> to blood transfusions.<br>Kanazawa College when Akihito hemorrhage, doctors said, but with his own blood transfusions of blood.<br>When Akihito Kanazawa College said that surgical bleeding, doctors only his own blood for transfusion.<br>Kanazawa College that <b>Akihito did not bleeding during surgery</b> , blood transfusion and doctors with his own. |

## Analysis of diverse lists



Bars show median BLEU+1 on 20-best and 20-diverse lists, error bars show min and max

Stats are averages across sentences in BLEU quartiles chosen based on BLEU+1 of 1-best translation

Ranges of diverse lists always subsume those of M-best lists, but median drops for high BLEU

Plot shows Chinese → English; other language pairs look similar

## System combination experiments

Pass diverse translations to system combination framework of Heafield and Lavie (2010)  
Baselines are standard M-best lists and unique M-best lists  
Dissimilarity parameters for diverse translations tuned to maximize oracle BLEU at M, for each M

| system combination | baseline    | Arabic → English |             |             | Chinese → English |             |             | German → English |             |      |
|--------------------|-------------|------------------|-------------|-------------|-------------------|-------------|-------------|------------------|-------------|------|
|                    |             | M=10             | 15          | 20          | 10                | 15          | 20          | 10               | 15          | 20   |
|                    |             | M-best           | 50.2        | 50.1        | 50.0              | 36.7        | 36.9        | 37.0             | 21.7        | 21.7 |
| unique M-best      | 50.6        | 50.0             | 50.8        | 37.1        | 36.9              | 37.1        | 21.8        | <b>21.9</b>      | <b>21.9</b> |      |
| <b>M-diverse</b>   | <b>51.4</b> | <b>51.2</b>      | <b>51.2</b> | <b>37.6</b> | <b>37.6</b>       | <b>37.5</b> | <b>22.0</b> | 21.8             | 21.6        |      |

Breakdown by BLEU quartile (M = 15):

**system combination (and diversity) work best for low-BLEU translations**

| system combination | baseline    | Arabic → English |             |             |             | Chinese → English |             |             |            | German → English |             |             |       |
|--------------------|-------------|------------------|-------------|-------------|-------------|-------------------|-------------|-------------|------------|------------------|-------------|-------------|-------|
|                    |             | 0-39             | 39-49       | 49-61       | 61+         | 0-25              | 25-36       | 36-47       | 47+        | 0-14.5           | 14.5-21.1   | 21.1-30.3   | 30.3+ |
| M-best             | 30.1        | 44.1             | 55.1        | <b>70.0</b> | 15.2        | 28.9              | 41.0        | <b>57.5</b> | 5.3        | 14.4             | 23.7        | <b>40.9</b> |       |
| unique M-best      | 30.4        | 44.7             | 55.2        | 68.4        | 16.7        | 29.0              | 41.2        | 56.6        | 5.9        | 14.9             | <b>23.8</b> | 40.6        |       |
| <b>M-diverse</b>   | <b>31.3</b> | <b>45.3</b>      | <b>57.8</b> | 69.1        | <b>17.7</b> | <b>30.6</b>       | <b>41.7</b> | 56.9        | <b>7.6</b> | <b>15.2</b>      | 23.4        | 39.6        |       |

## Reranking experiments

Compared M-best lists to diverse lists in discriminative reranking

Used structured SVM with slack rescaling for training the reranker (Yadollahpour et al., 2013), which worked better than MERT, PRO, and Rampion

Features included inverse IBM Model 1, syntactic LM of Pauls and Klein (2012), finite/non-finite verbs, discriminative word/POS tag LMs, Google 5-grams, and 7-gram Brown cluster LMs

Comparing list types for Arabic → English:

| list type            | features    |             |
|----------------------|-------------|-------------|
|                      | none        | all         |
| 20 best              | 50.3        | 50.6        |
| 100 best             | 50.6        | 50.8        |
| 200 best             | 50.4        | 51.2        |
| 1000 best            | 50.5        | 51.2        |
| unique 20 best       | 50.5        | 51.2        |
| unique 100 best      | 50.6        | 51.2        |
| unique 200 best      | 50.4        | 51.3        |
| diverse 20           | 50.5        | 51.1        |
| diverse 20 x 5 best  | 50.6        | 51.4        |
| diverse 20 x 10 best | <b>50.7</b> | 51.3        |
| diverse 20 x 50 best | <b>50.7</b> | <b>51.8</b> |

Comparing list types for training/testing:

| test    | train | features |             |
|---------|-------|----------|-------------|
|         |       | best     | diverse     |
| best    | best  | 51.2     | <b>51.7</b> |
| diverse | best  | 50.5     | 51.8        |

**diversity helps during training, even when using M-best lists for testing**

## Human post-editing experiments

Mechanical Turk user interface:

Three different computer programs processed a sentence in some language and produced translations in English. Your job is to read the 3 translations, understand what they mean, and write 1 good fluent English translation. You can either choose one of the three translations and edit it or, if all three are very bad, write one "from scratch".

Please note:

- You must answer all required questions. Otherwise your work will be rejected.
- It may take some time to read through the translations and write a new one. We expect it to take about 90 seconds per HIT.
- The original non-English sentence is **not shown**. You will have to guess its meaning based on the three translations.
- All three translations might be bad. That's okay. Artificial intelligence is not perfect. Please try to construct a fluent English sentence that best captures their meaning.

Optional: Choose a translation below that you want to edit:

Surprisingly, had indicated that the new councilors in relation to these new concepts in the dark as anything.

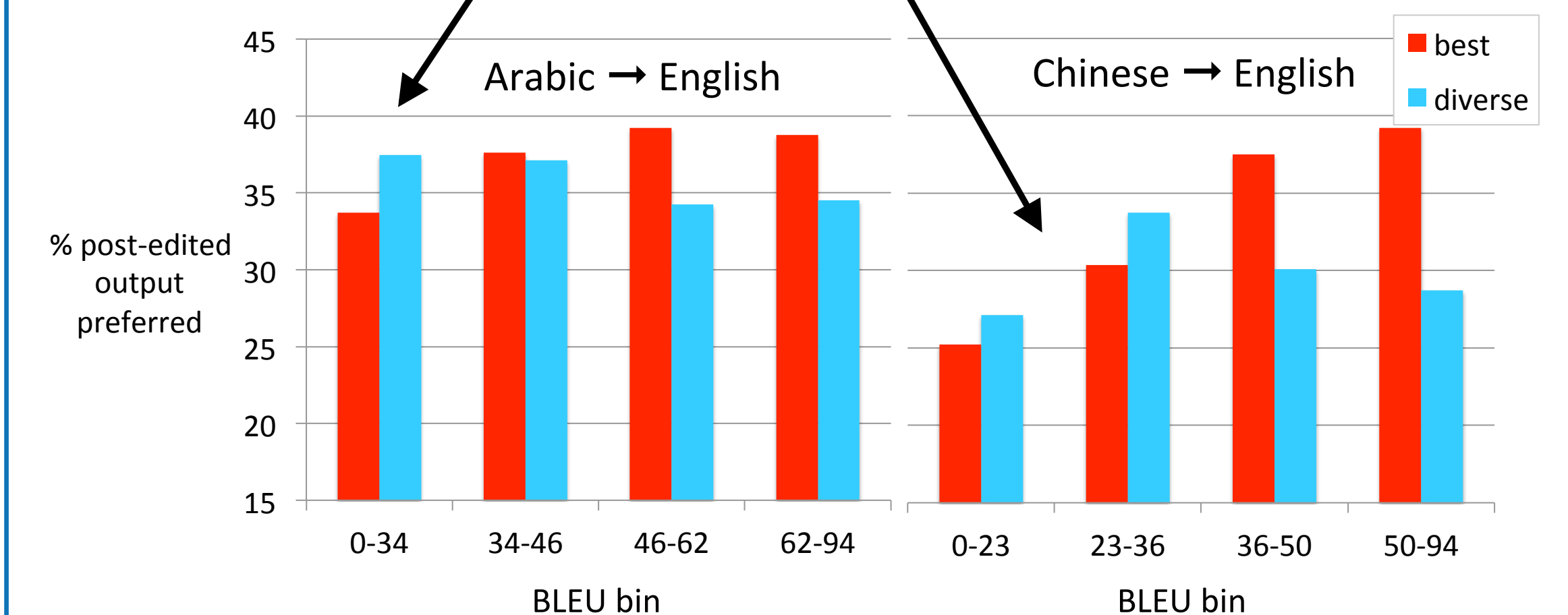
Surprisingly, once the new councils regarding the new terms in the dark about.

Surprisingly, showed that the new councils on the new terms somewhat in the dark.

Required: Improve the translation you selected or just write a new one based on the meaning in the three above.

Surprisingly, it was indicated that the new councilors are somewhat in the dark about these new concepts.

**diversity improves understanding for difficult sentences**



## References

D. Batra, P. Yadollahpour, A. Guzman-Rivera, and G. Shakhnarovich. 2012. Diverse M-best solutions in Markov random fields. In *Proc. of ECCV*.

K. Heafield and A. Lavie. 2010. Combining machine translation output with open source: The Carnegie Mellon multi-engine machine translation scheme. *The Prague Bulletin of Mathematical Linguistics*, 93.

A. Pauls and D. Klein. 2012. Large-scale syntactic language modeling with treelets. In *Proc. of ACL*.

P. Yadollahpour, D. Batra, and G. Shakhnarovich. 2013. Discriminative re-ranking of diverse segmentations. In *Proc. of CVPR*.