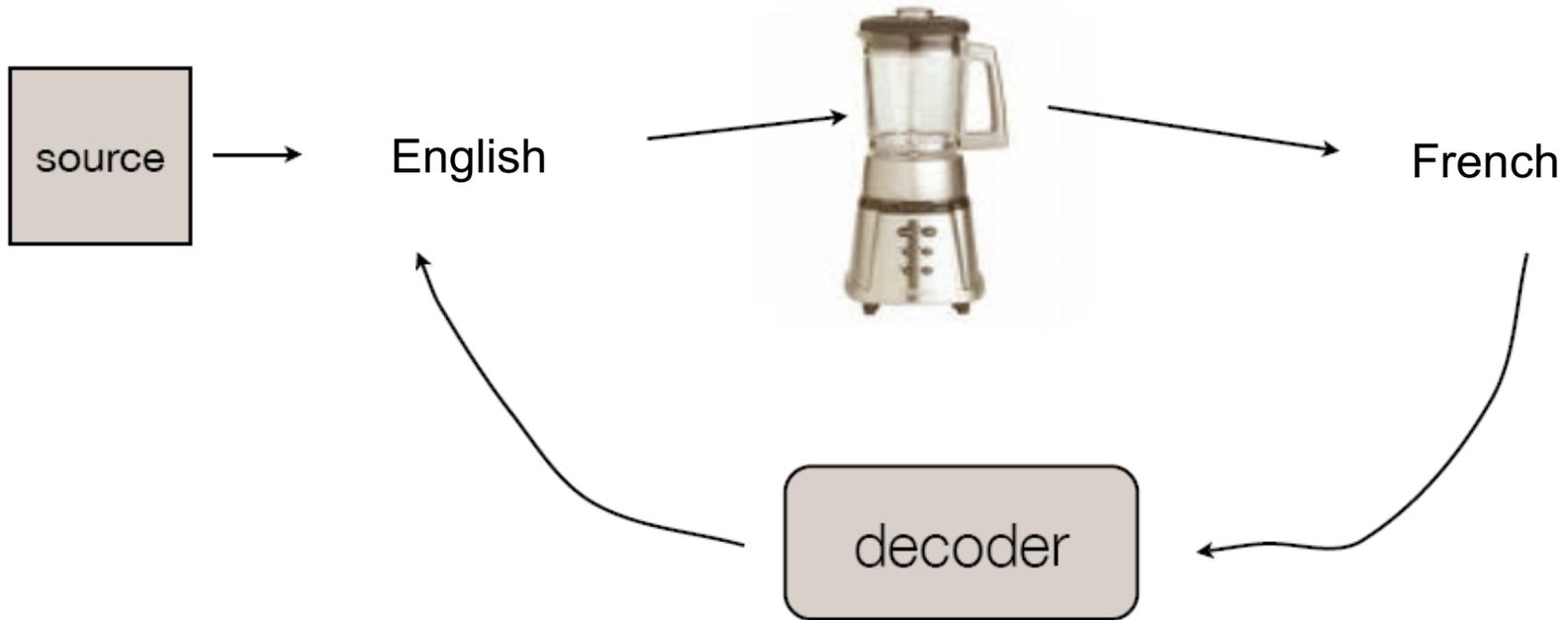


Discriminative Training

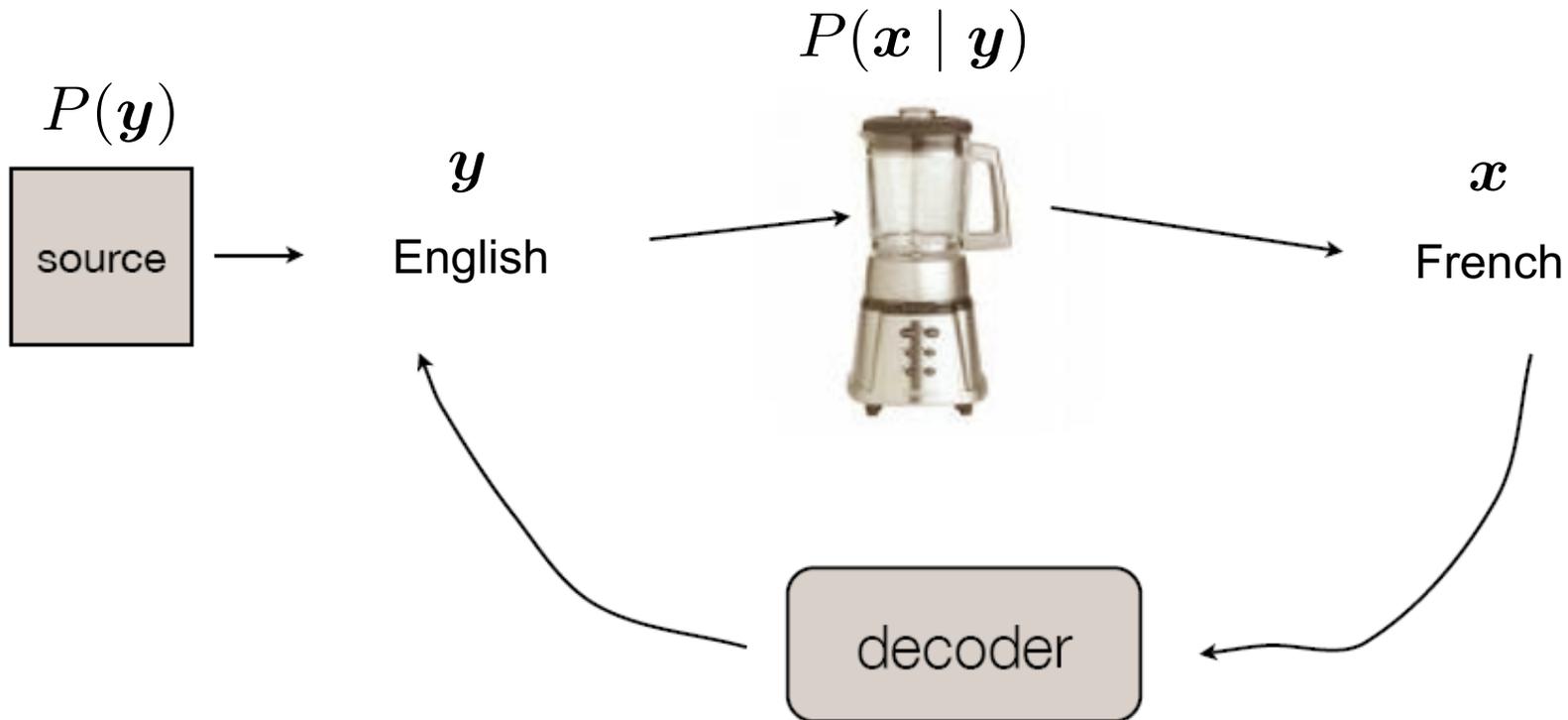
Kevin Gimpel



Noisy Channel Model



Noisy Channel Model for Translating French (x) to English (y)



$$y^* = \operatorname{argmax}_y P(y | x)$$

$$= \operatorname{argmax}_y \frac{P(x | y)P(y)}{P(x)}$$

$$= \operatorname{argmax}_y P(x | y)P(y)$$

Noisy Channel

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} P(\mathbf{x} | \mathbf{y}) P(\mathbf{y})$$

Noisy Channel

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} P(\mathbf{x} | \mathbf{y}) P(\mathbf{y})$$

assumes we have the right model, and that we estimate it perfectly

Noisy Channel

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} P(\mathbf{x} | \mathbf{y}) P(\mathbf{y})$$

assumes we have the right model, and that we estimate it perfectly

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} P(\mathbf{x} | \mathbf{y})^\alpha P(\mathbf{y})^\beta$$

Noisy Channel

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} P(\mathbf{x} | \mathbf{y}) P(\mathbf{y})$$

assumes we have the right model, and that we estimate it perfectly

$$\begin{aligned} \mathbf{y}^* &= \operatorname{argmax}_{\mathbf{y}} P(\mathbf{x} | \mathbf{y})^\alpha P(\mathbf{y})^\beta \\ &= \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y}) \end{aligned}$$

extra parameters to tune, can tune to
optimize BLEU (or whatever metric you want)

“tuning”

Noisy Channel \rightarrow Linear Model?

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y})$$

since we're not using idealized decoding rule anymore,
why not add more feature functions?

Noisy Channel \rightarrow Linear Model?

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y})$$

since we're not using idealized decoding rule anymore,
why not add more feature functions?

“word count feature”:

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y}) + \boxed{\gamma |\mathbf{y}|}$$

Noisy Channel \rightarrow Linear Model?

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y})$$

since we're not using idealized decoding rule anymore,
why not add more feature functions?

“word count feature”:

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y}) + \boxed{\gamma |\mathbf{y}|}$$

“reverse translation model feature”:

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y}) + \gamma |\mathbf{y}| + \boxed{\delta \log P(\mathbf{y} | \mathbf{x})}$$

Noisy Channel \rightarrow Linear Model?

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y})$$

since we're not using idealized decoding rule anymore,
why not add more feature functions?

“word count feature”:

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y}) + \boxed{\gamma |\mathbf{y}|}$$

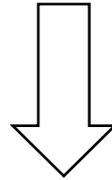
“reverse translation model feature”:

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y}) + \gamma |\mathbf{y}| + \boxed{\delta \log P(\mathbf{y} | \mathbf{x})}$$

but if we keep adding features, tuning gets harder...

Noisy Channel \rightarrow Linear Model

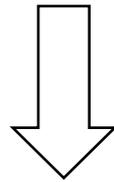
$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y})$$



generalize to a linear model

Noisy Channel \rightarrow Linear Model

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y})$$



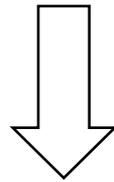
generalize to a linear model

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \sum_{i=1}^K \theta_i f_i(\mathbf{x}, \mathbf{y})$$

“feature weights” *“feature functions”*

Noisy Channel \rightarrow Linear Model

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \alpha \log P(\mathbf{x} | \mathbf{y}) + \beta \log P(\mathbf{y})$$



generalize to a linear model

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \sum_{i=1}^K \theta_i f_i(\mathbf{x}, \mathbf{y})$$

we will write it like this:

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})$$

“feature weight vector”

“feature function vector”

Log-Linear Models

$$\mathbf{y}^* = \underset{\mathbf{y}}{\operatorname{argmax}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})$$

if you want a probability distribution over translations of a given source sentence, exponentiate and normalize:

Log-Linear Models

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})$$

if you want a probability distribution over translations of a given source sentence, exponentiate and normalize:

$$P(\mathbf{y} \mid \mathbf{x}) = \frac{\exp\{\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})\}}{\sum_{\mathbf{y}'} \exp\{\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}')\}}$$

this is a (conditional) “log-linear” model

More General Formulation: Derivations

The diagram illustrates the derivation of the optimal translation and source sentence. It features a central equation: $\langle \mathbf{y}^*, \mathbf{h}^* \rangle = \operatorname{argmax}_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_x} \boldsymbol{\theta}^\top \mathbf{f}(x, \mathbf{y}, \mathbf{h})$. Four blue arrows point from descriptive labels to the components of the equation: 'feature weight vector' points to $\boldsymbol{\theta}^\top$, 'feature function vector' points to $\mathbf{f}(x, \mathbf{y}, \mathbf{h})$, 'source sentence' points to x , and 'translation' points to $\langle \mathbf{y}, \mathbf{h} \rangle$.

feature weight vector

feature function vector

$\langle \mathbf{y}^*, \mathbf{h}^* \rangle = \operatorname{argmax}_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_x} \boldsymbol{\theta}^\top \mathbf{f}(x, \mathbf{y}, \mathbf{h})$

source sentence *translation*

More General Formulation: Derivations

feature weight vector

feature function vector

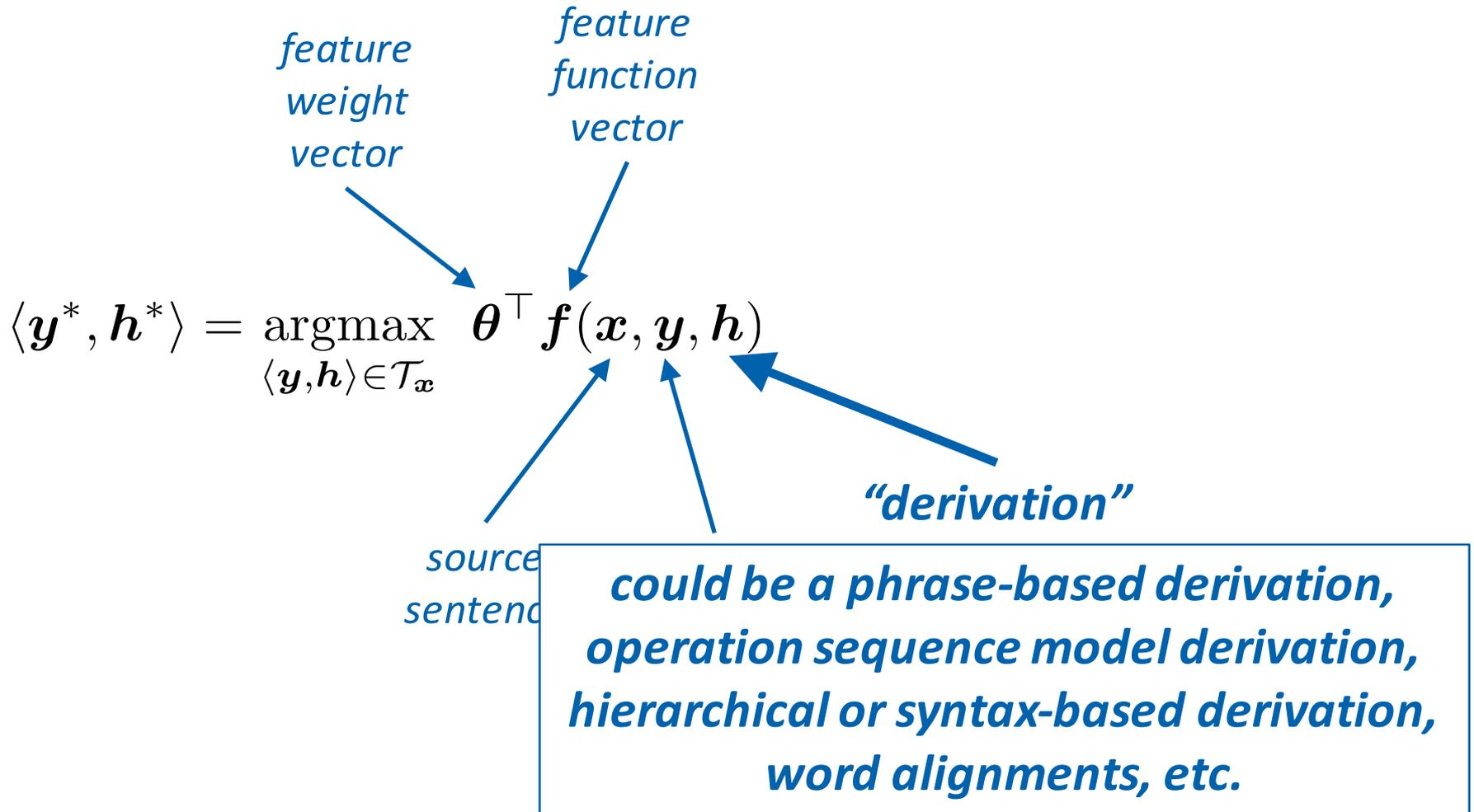
$\langle \mathbf{y}^*, \mathbf{h}^* \rangle = \operatorname{argmax}_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_x} \boldsymbol{\theta}^\top \mathbf{f}(x, \mathbf{y}, \mathbf{h})$

source sentence

translation

"derivation"

More General Formulation: Derivations



More General Formulation: Derivations

The diagram illustrates the general formulation of the maximum margin problem. It features the equation $\langle \mathbf{y}^*, \mathbf{h}^* \rangle = \operatorname{argmax}_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_x} \boldsymbol{\theta}^\top \mathbf{f}(x, \mathbf{y}, \mathbf{h})$. Annotations include: "feature weight vector" pointing to $\boldsymbol{\theta}$; "feature function vector" pointing to \mathbf{f} ; "source sentence" pointing to x ; and "derivation" pointing to \mathbf{h} .

$$\langle \mathbf{y}^*, \mathbf{h}^* \rangle = \operatorname{argmax}_{\langle \mathbf{y}, \mathbf{h} \rangle \in \mathcal{T}_x} \boldsymbol{\theta}^\top \mathbf{f}(x, \mathbf{y}, \mathbf{h})$$

feature weight vector

feature function vector

source sentence

"derivation"

We won't talk much more about derivations in this lecture, but they're (almost) always there

Overview

This lecture is about algorithms for choosing the feature weights θ

Since we have a linear model, we can use supervised machine learning

But MT differs from typical supervised tasks (as we'll see), so MT training procedures differ too

We'll start with a way to visualize training
for machine translation

African
National
Congress

非国大

opposition

反对

sanction

制裁

Zimbabwe

津巴布韦

African
National
Congress

opposition

sanction

Zimbabwe

非国大

反对

制裁

津巴布韦

Gold standard:

African National Congress opposes
sanctions against Zimbabwe

African
National
Congress

opposition

sanction

Zimbabwe

非国大

反对

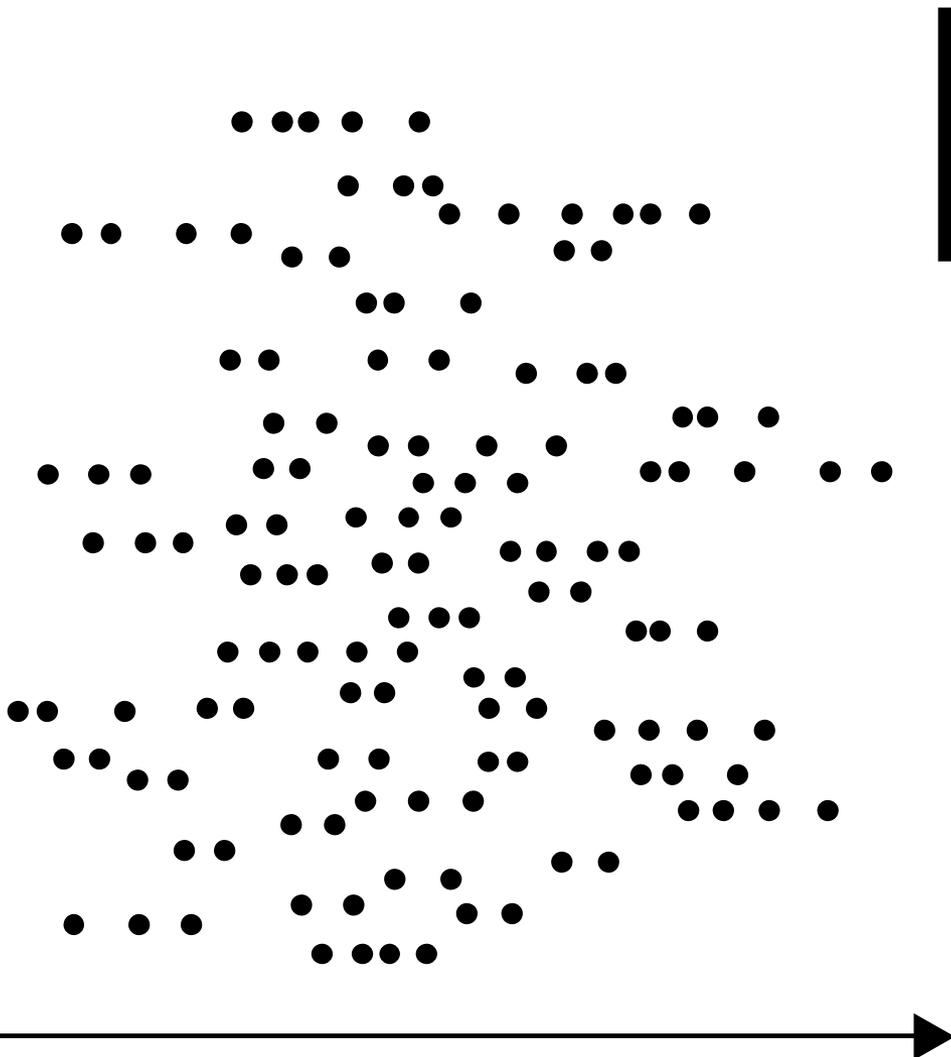
制裁

津巴布韦

Gold standard:

African National Congress opposes
sanctions against Zimbabwe

BLEU
score



**each point is a
translation**

model score

African
National
Congress

opposition

sanction

Zimbabwe

非国大

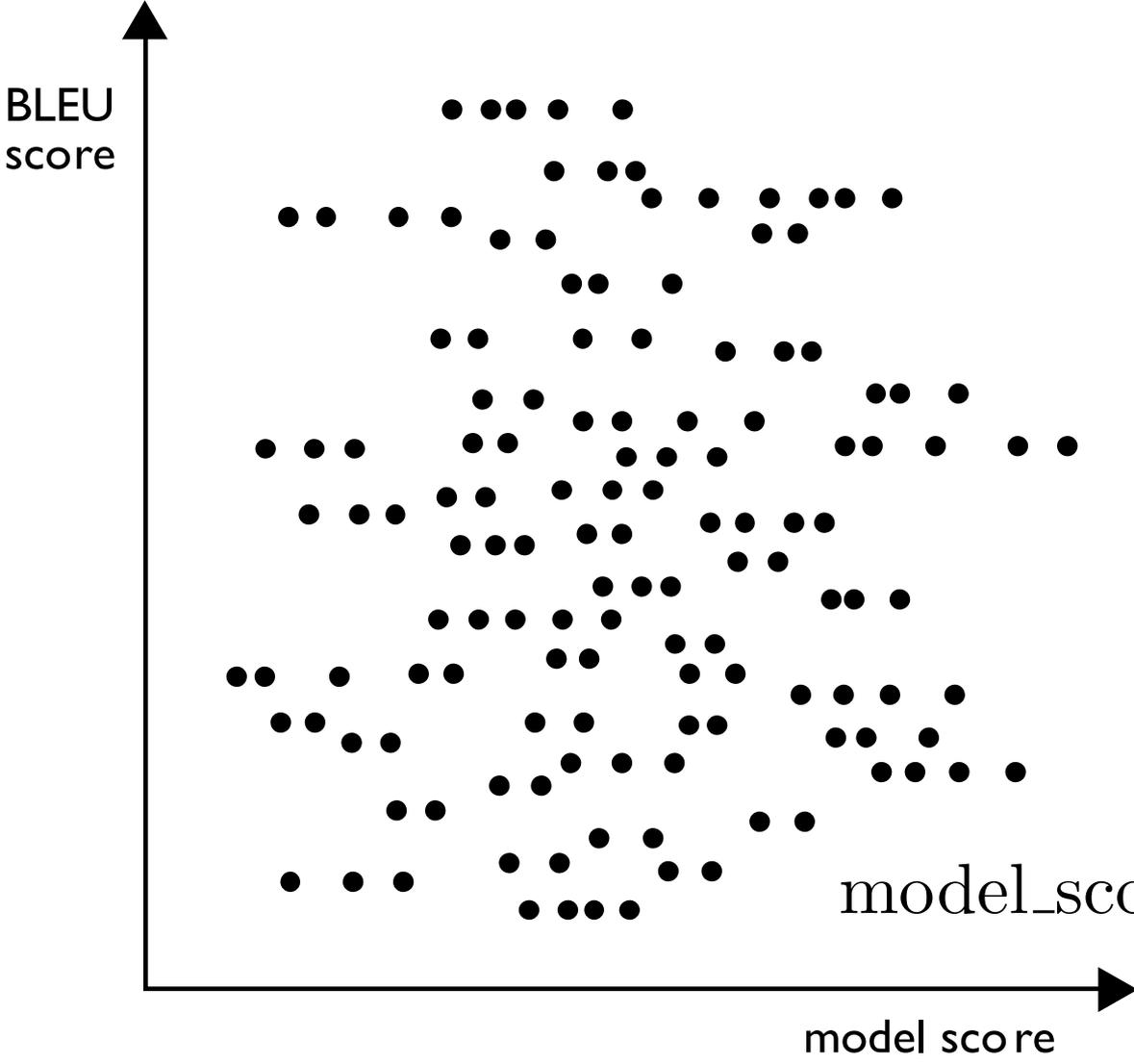
反对

制裁

津巴布韦

Gold standard:

African National Congress opposes
sanctions against Zimbabwe



$$\text{model_score}(x, y) = \theta^\top f(x, y)$$

African
National
Congress

opposition

sanction

Zimbabwe

非国大

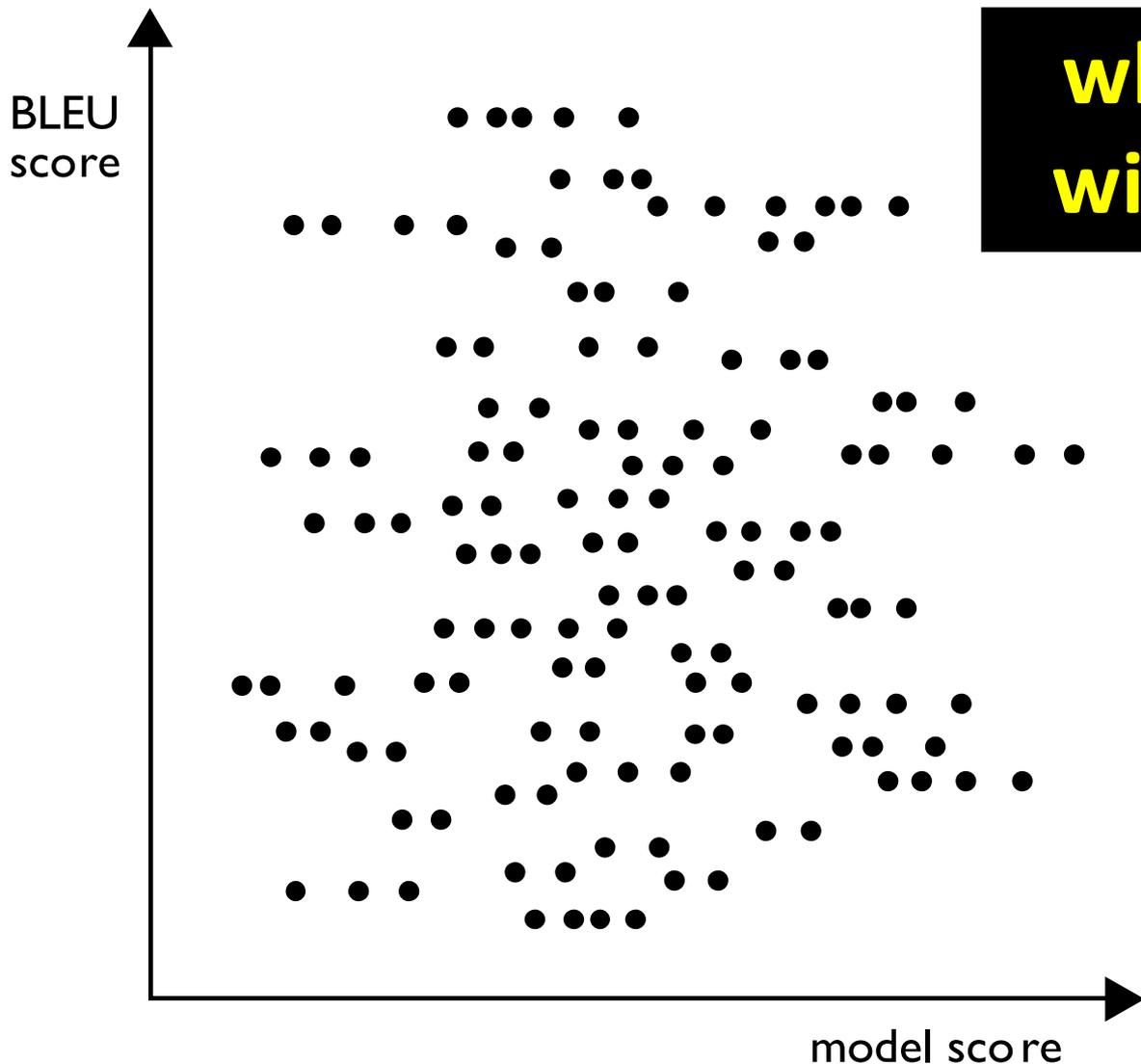
反对

制裁

津巴布韦

Gold standard:

African National Congress opposes
sanctions against Zimbabwe



**which translation
will be predicted?**

African
National
Congress

opposition

sanction

Zimbabwe

非国大

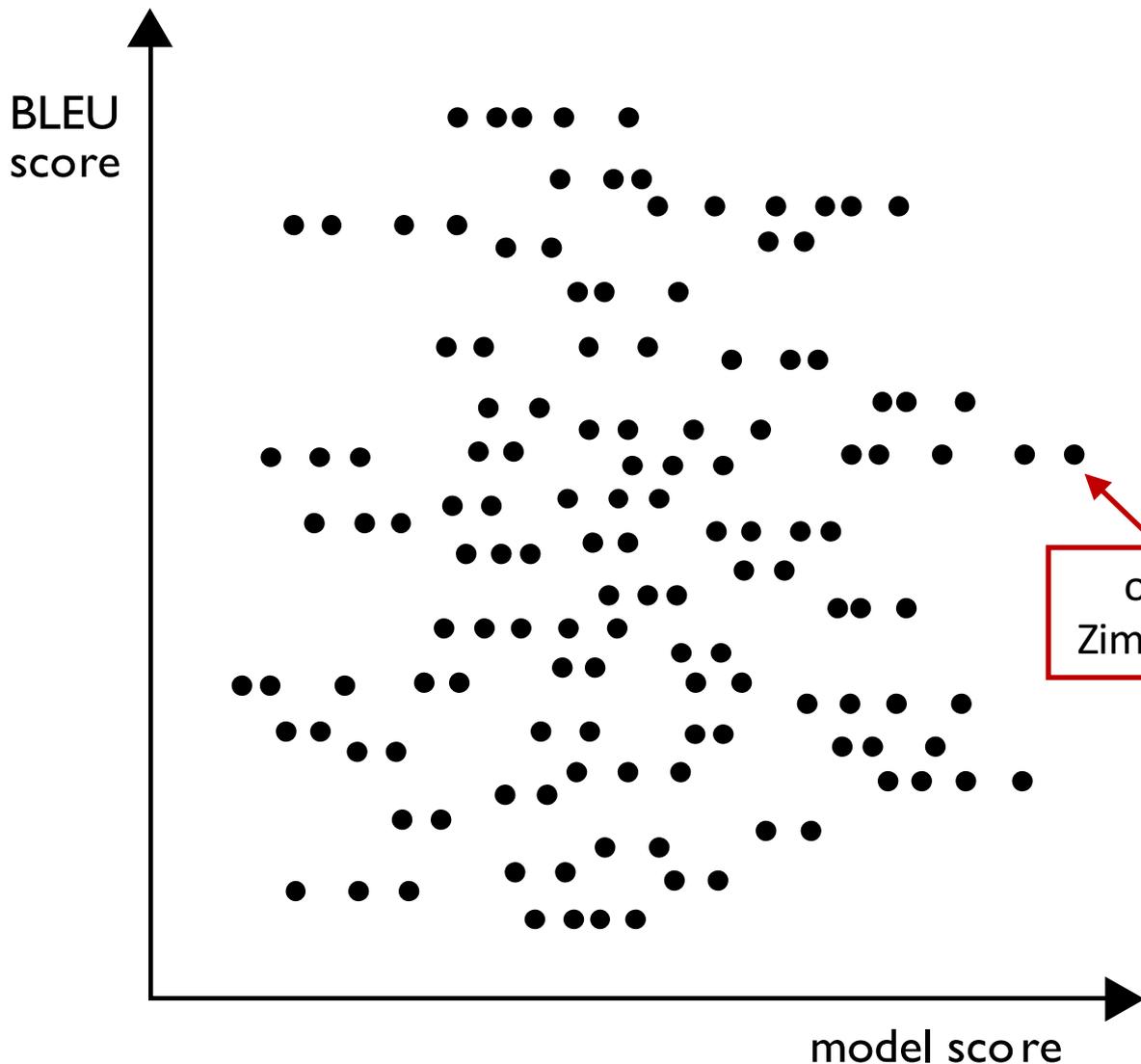
反对

制裁

津巴布韦

Gold standard:

African National Congress opposes
sanctions against Zimbabwe



predicted translation

opposition to sanctions against
Zimbabwe African National Congress

African National Congress

opposition

sanction

Zimbabwe

Gold standard:

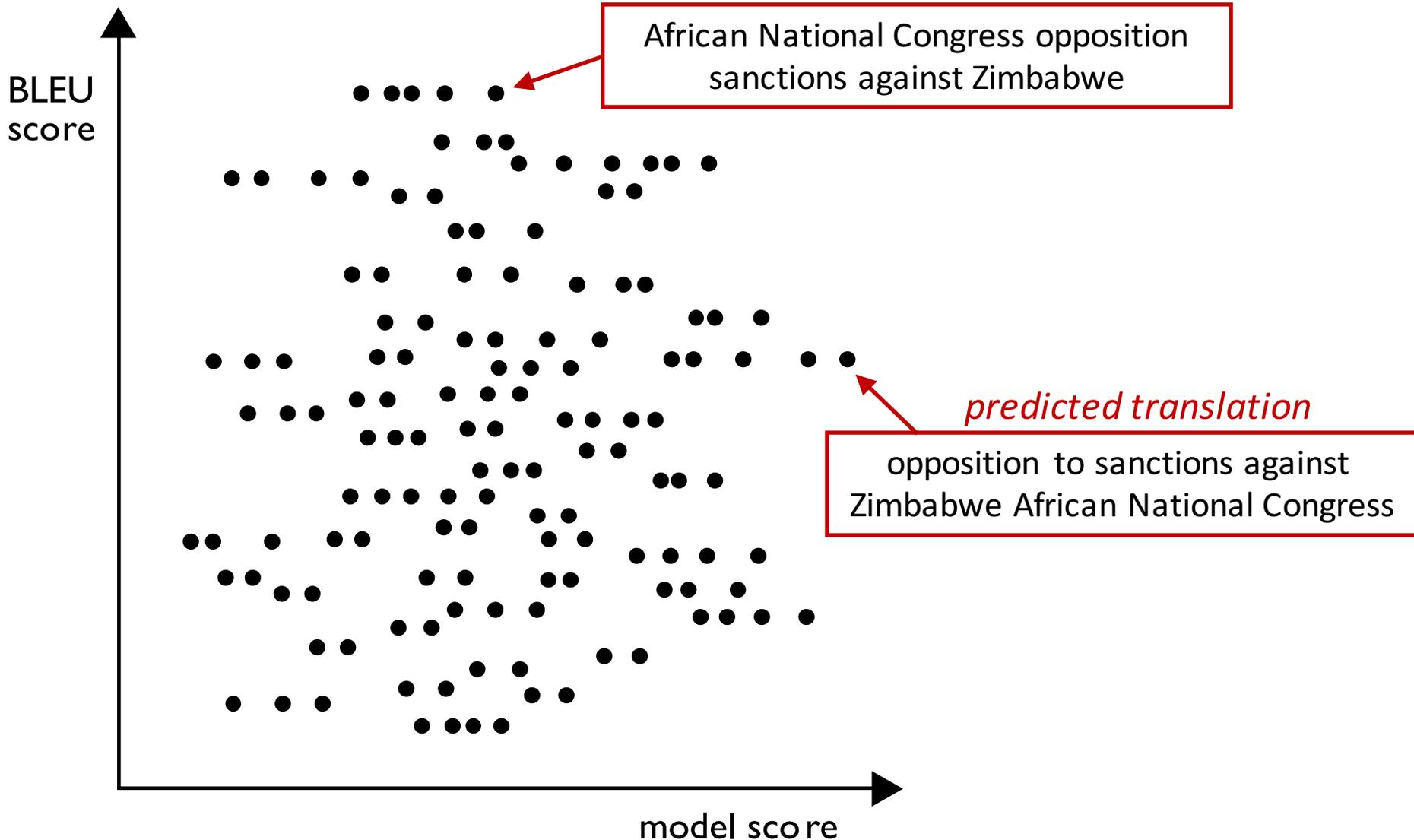
African National Congress opposes sanctions against Zimbabwe

非国大

反对

制裁

津巴布韦



African National Congress

opposition

sanction

Zimbabwe

Gold standard:

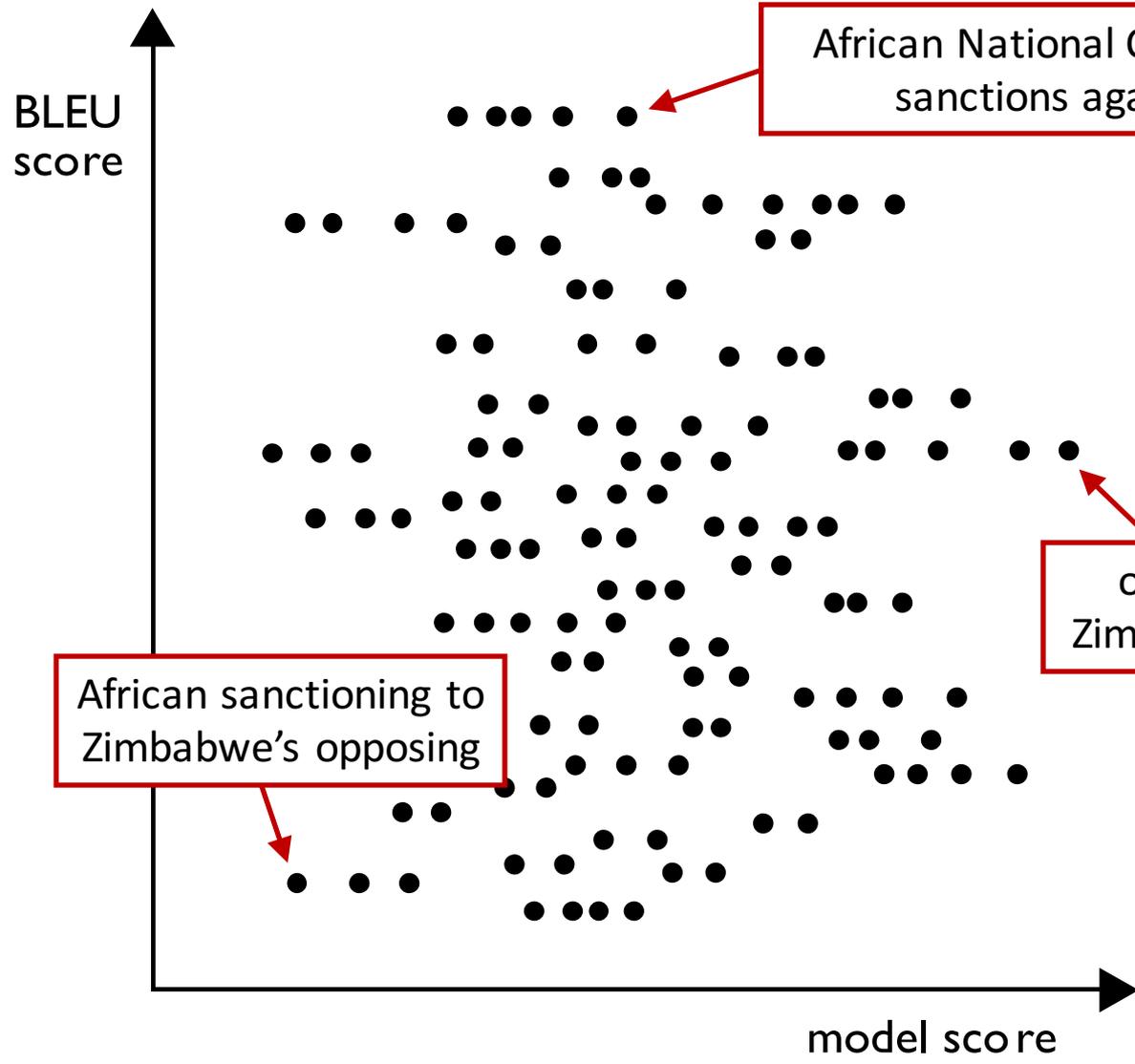
African National Congress opposes sanctions against Zimbabwe

非国大

反对

制裁

津巴布韦



African
National
Congress

opposition

sanction

Zimbabwe

非国大

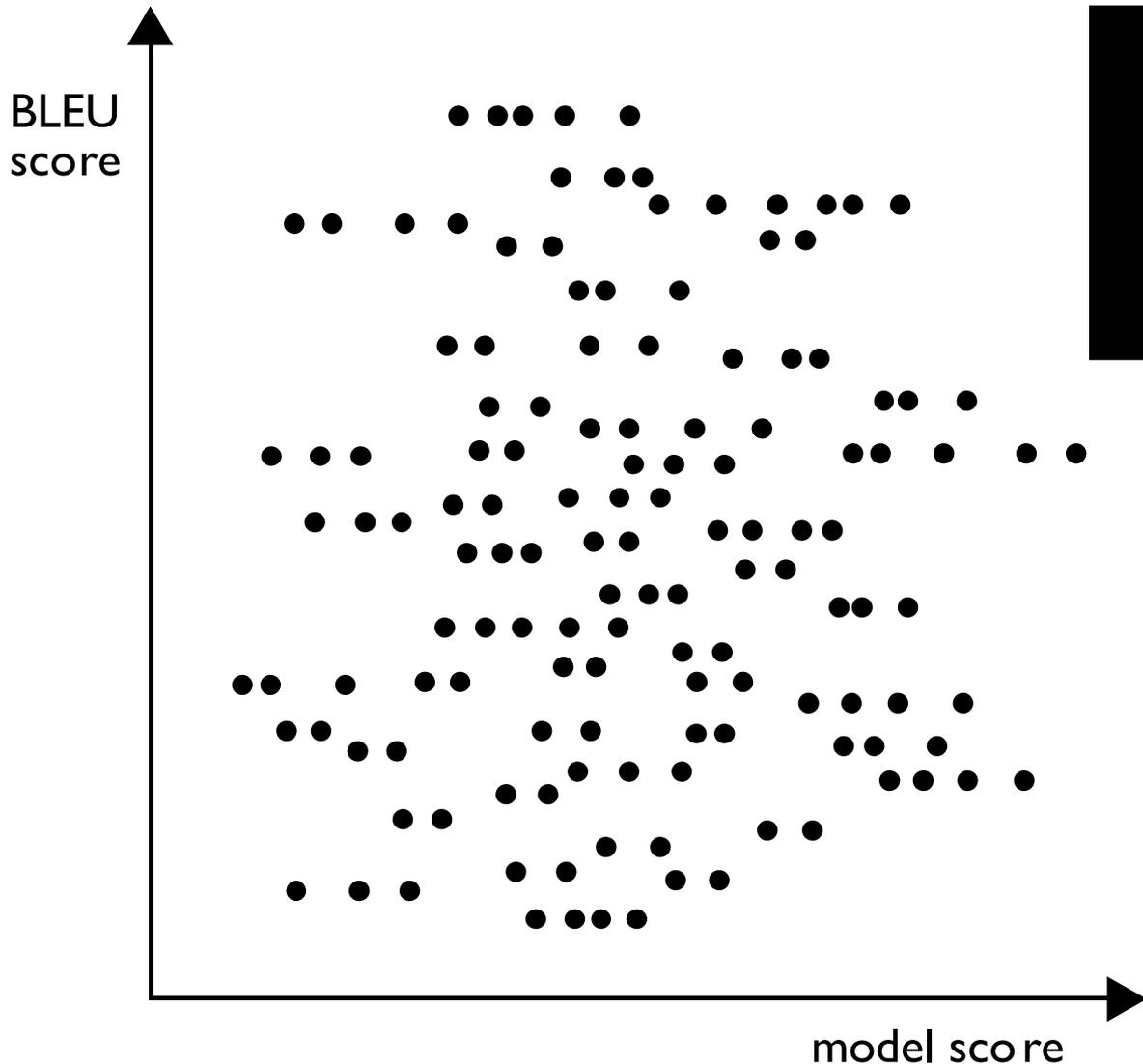
反对

制裁

津巴布韦

Gold standard:

African National Congress opposes
sanctions against Zimbabwe



**training moves
translations in
this plot**

African
National
Congress

opposition

sanction

Zimbabwe

非国大

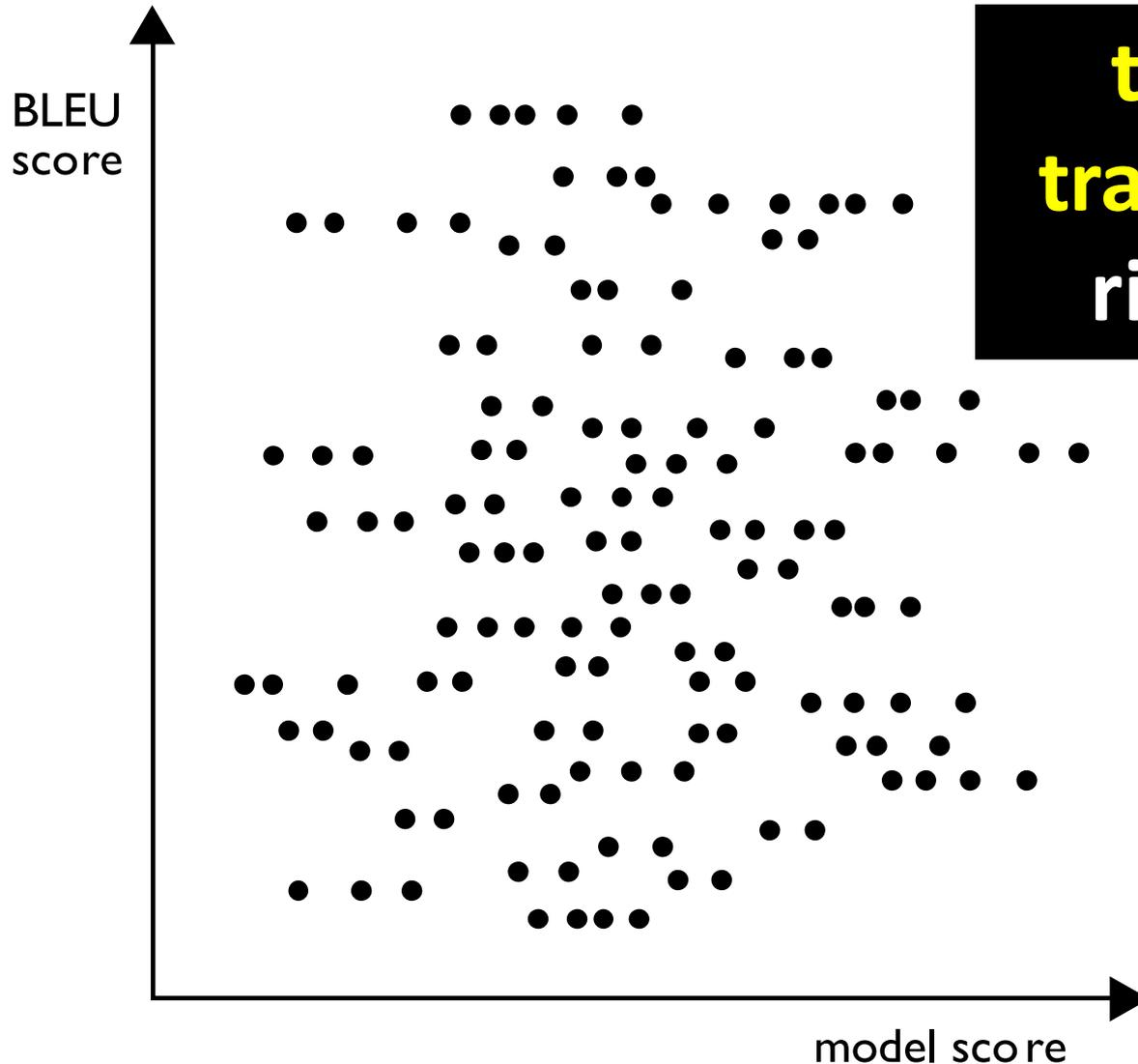
反对

制裁

津巴布韦

Gold standard:

African National Congress opposes
sanctions against Zimbabwe



training moves
translations left or
right in this plot

African
National
Congress

opposition

sanction

Zimbabwe

非国大

反对

制裁

津巴布韦

Gold standard:

African National Congress opposes
sanctions against Zimbabwe

BLEU
score

“ideal” model?

model score

African
National
Congress

opposition

sanction

Zimbabwe

非国大

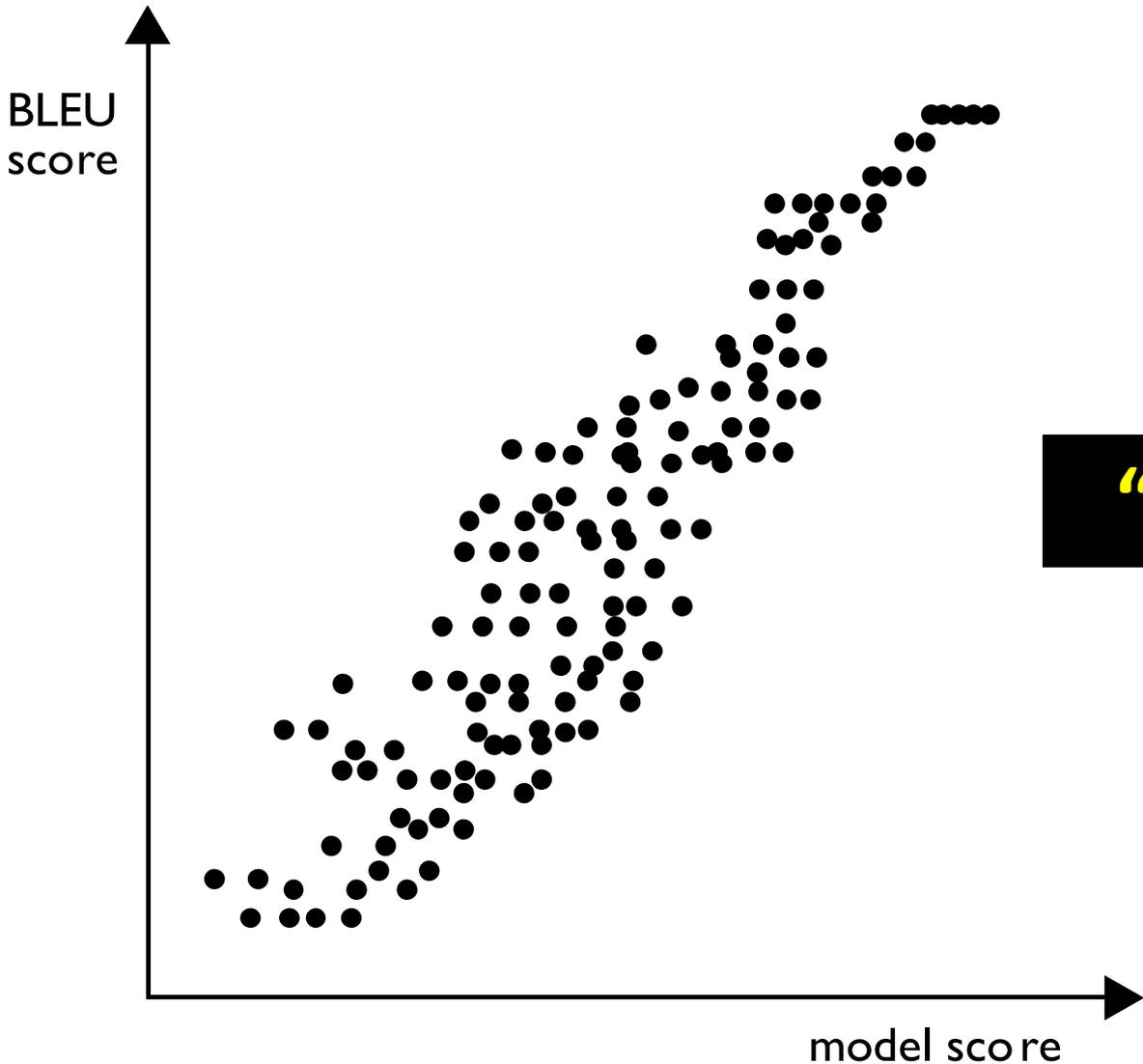
反对

制裁

津巴布韦

Gold standard:

African National Congress opposes
sanctions against Zimbabwe



“ideal” model?

African
National
Congress

opposition

sanction

Zimbabwe

非国大

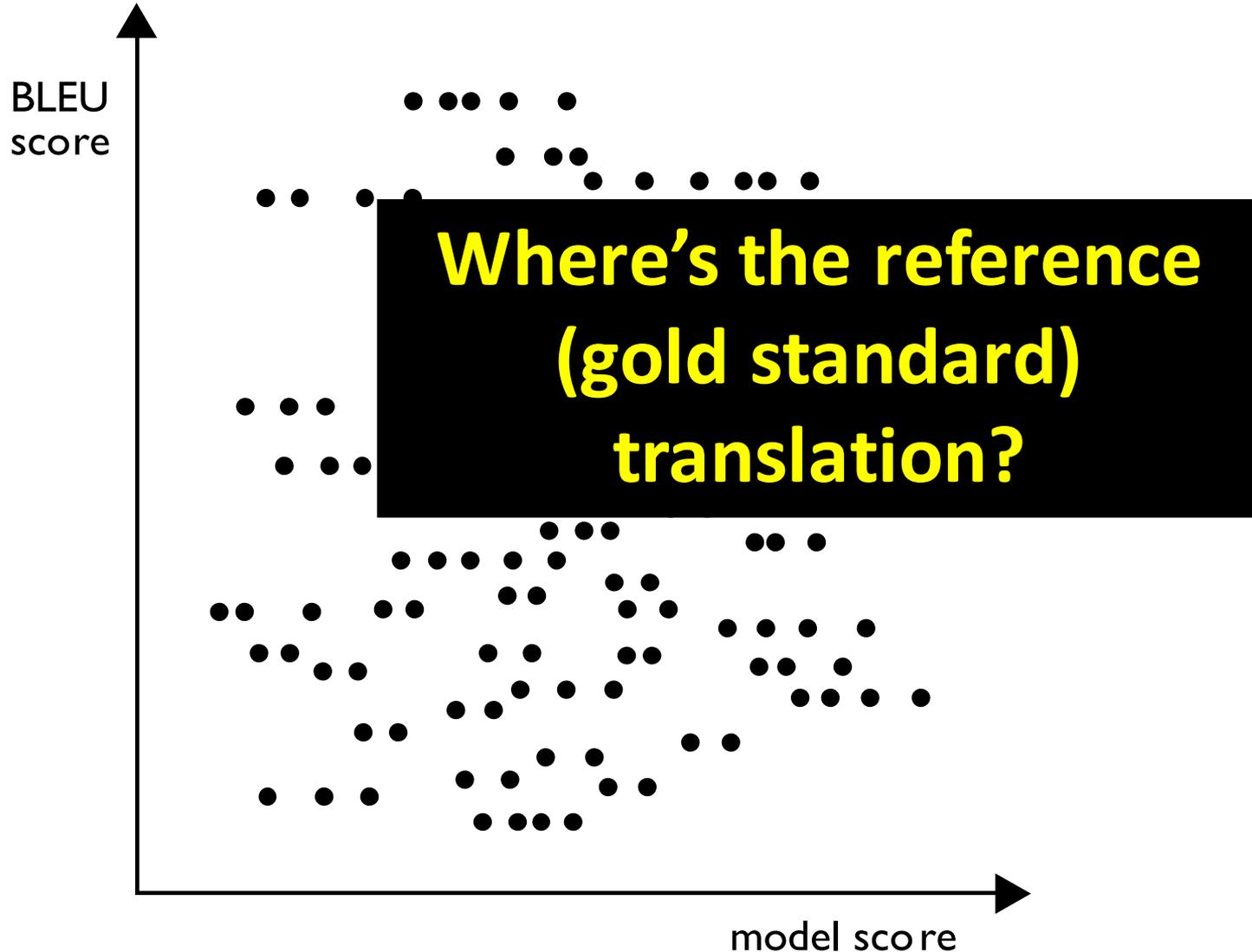
反对

制裁

津巴布韦

Gold standard:

African National Congress opposes
sanctions against Zimbabwe



African
National
Congress

opposition

sanction

Zimbabwe

非国大

反对

制裁

津巴布韦

Gold standard:

African National Congress opposes
sanctions against Zimbabwe

BLEU
score

Issue:

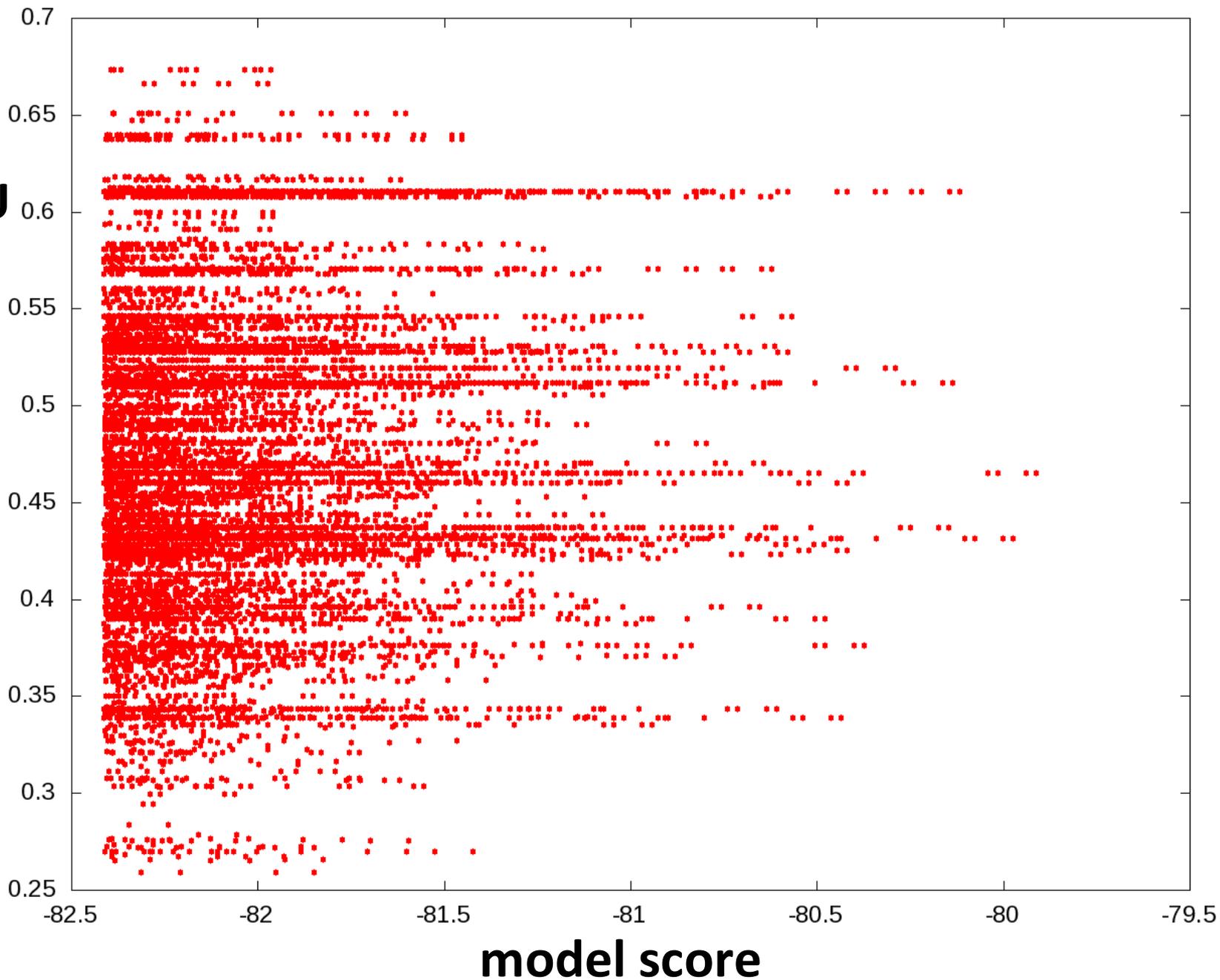
**gold standard translation is often
unreachable by the model**

Why?

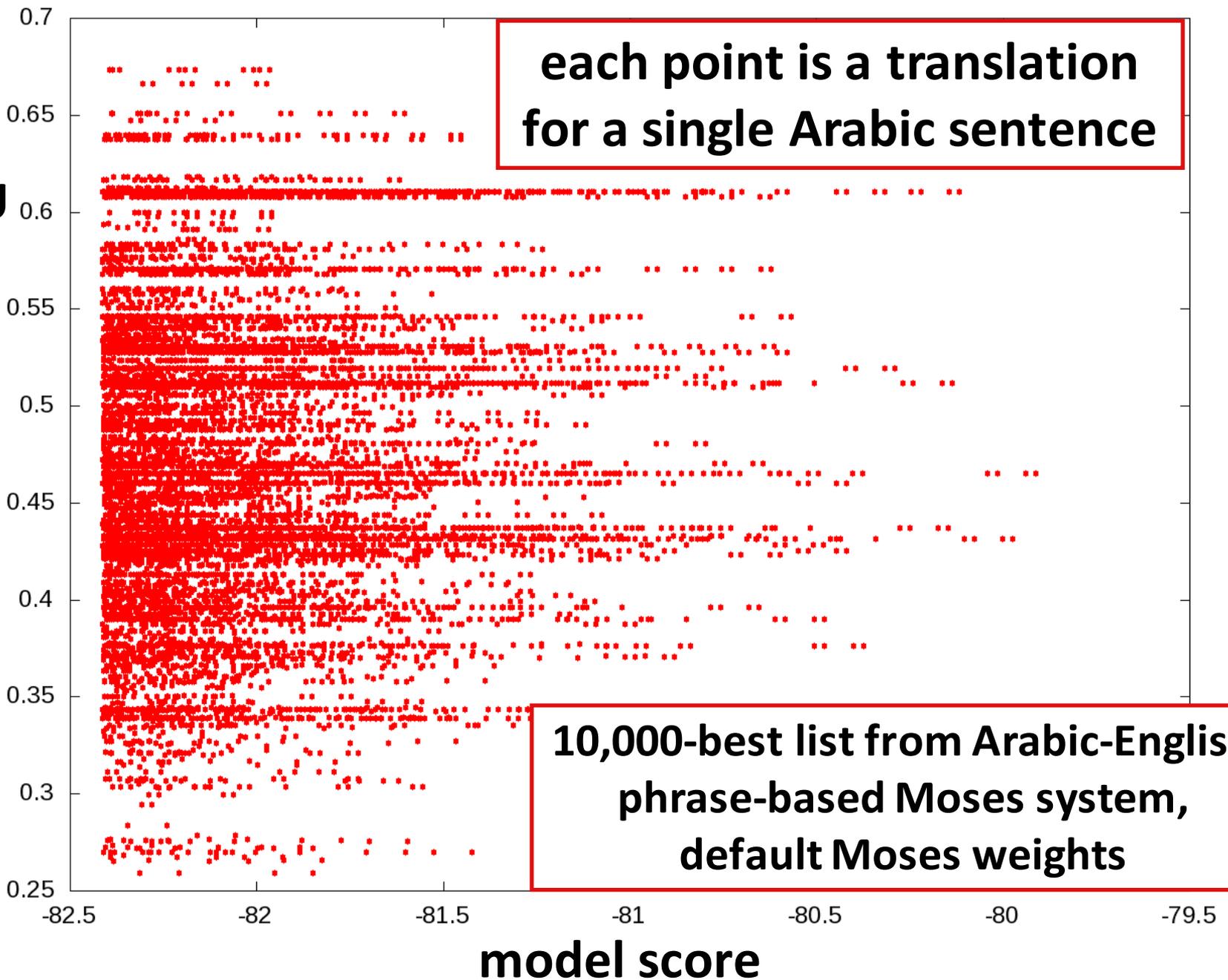
**limited translation rules,
free translations,
noisy data**

model score

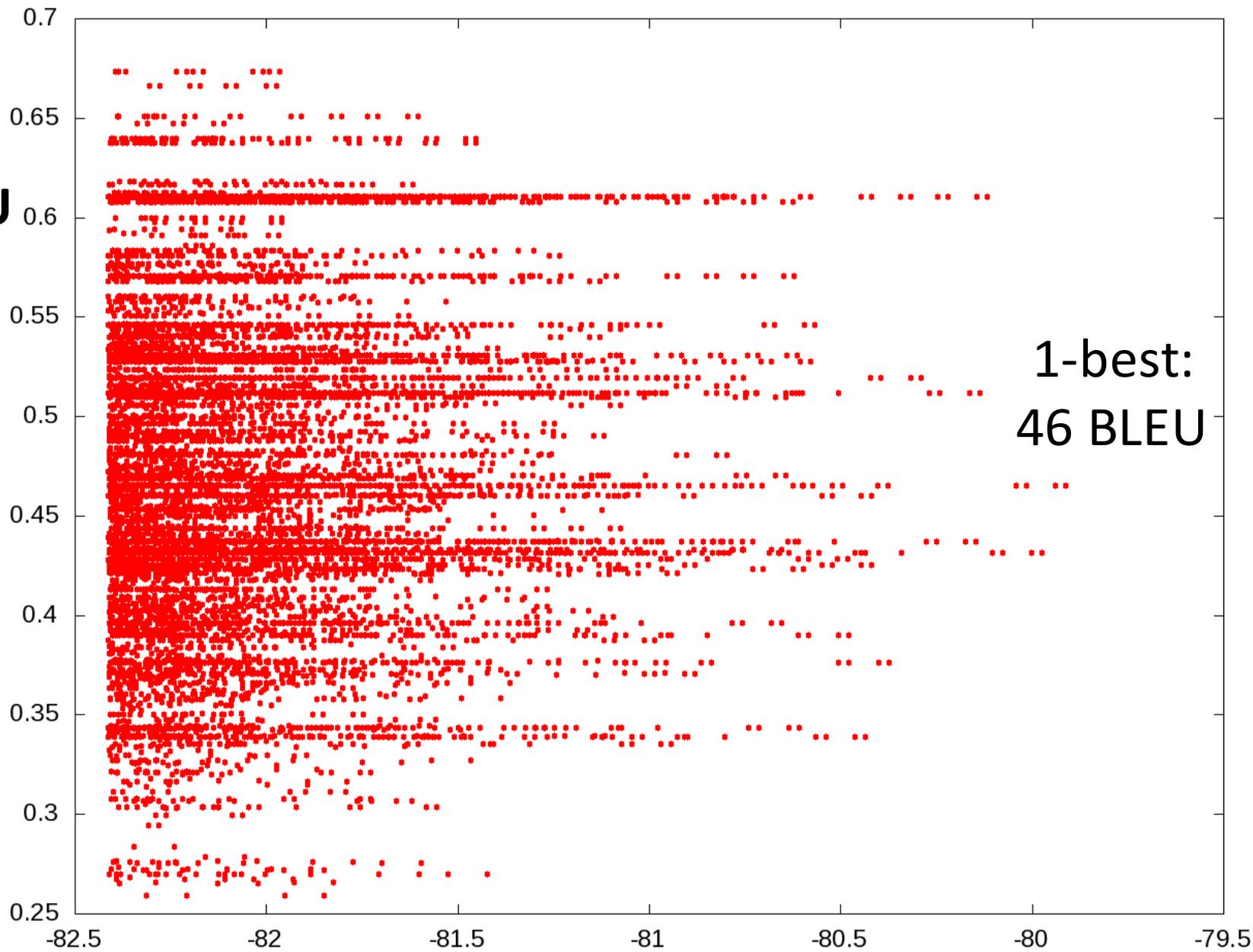
BLEU



BLEU



BLEU

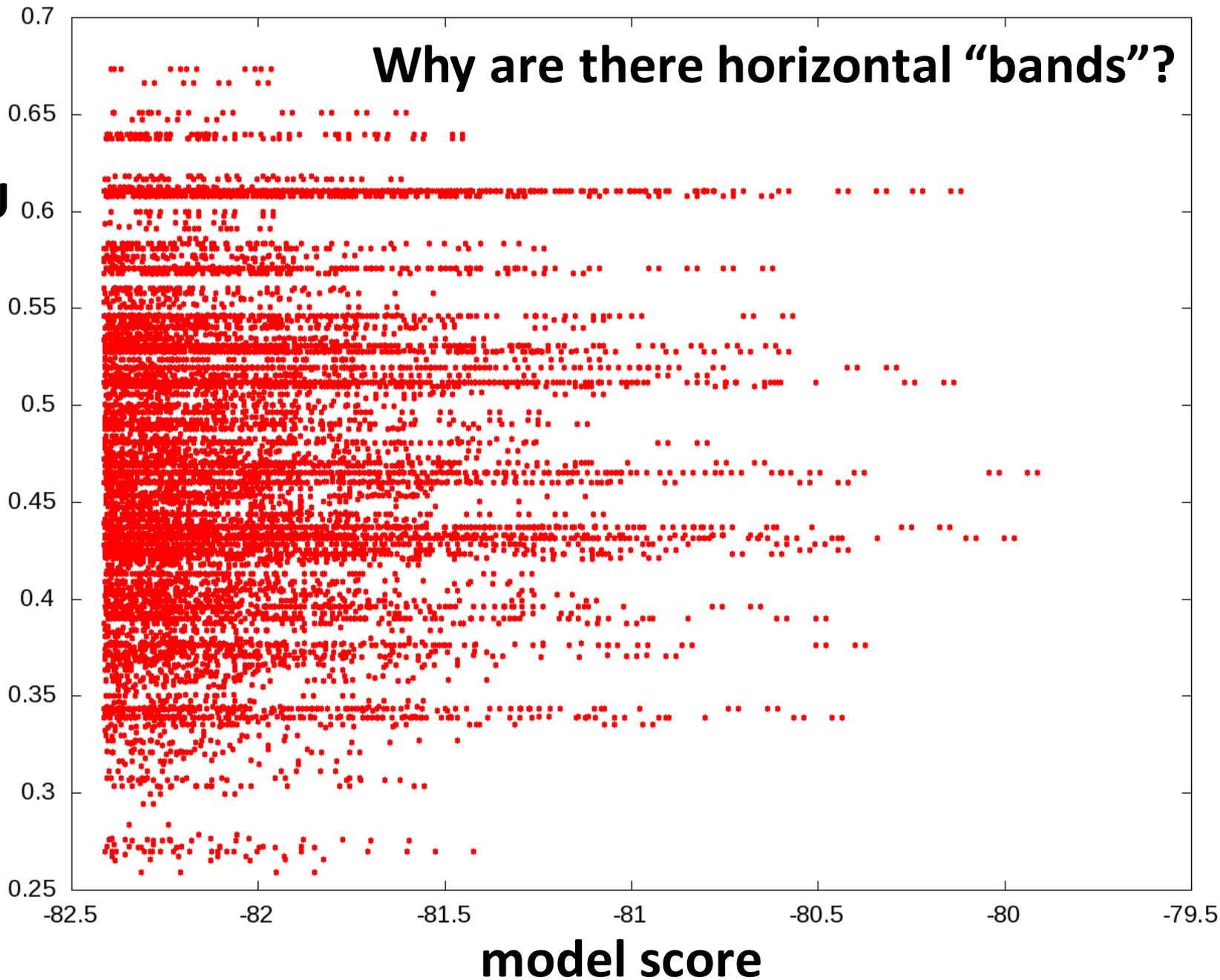


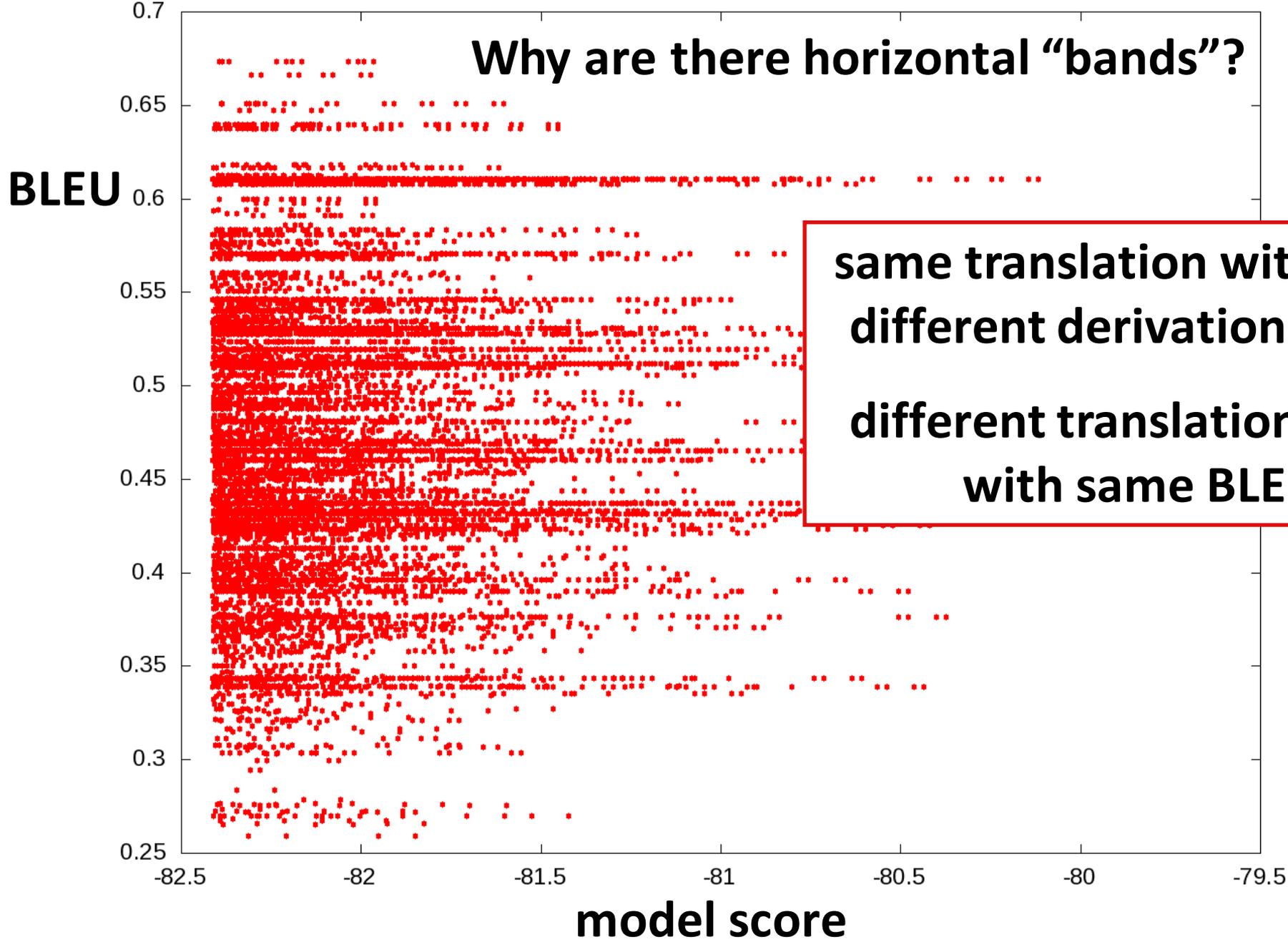
1-best:
46 BLEU

model score

Why are there horizontal “bands”?

BLEU





Roadmap

- We'll cover some of the most widely-used discriminative training algorithms for MT:
 - Minimum Error Rate Training*
 - Minimum Bayes Risk
 - Pairwise Ranking Optimization*
 - Batch MIRA*

* implemented in Moses!

2003: MERT

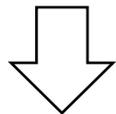
Minimum Error Rate Training in Statistical Machine Translation

Franz Josef Och

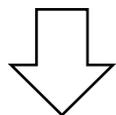
Information Sciences Institute
University of Southern California
4676 Admiralty Way, Suite 1001
Marina del Rey, CA 90292
och@isi.edu

Minimum Error Rate Training (MERT)

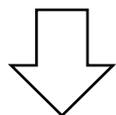
$$\min_{\theta} \text{error}(\text{predicted translations})$$



$$\min_{\theta} -\text{BLEU}(\text{references}, \text{predicted translations})$$



$$\max_{\theta} \text{BLEU}(\text{references}, \text{predicted translations})$$



$$\max_{\theta} \text{BLEU}(\text{references}, \{\text{translate}(\mathbf{x}_i, \theta)\}_{i=1}^N)$$

$$\text{translate}(\mathbf{x}, \theta) = \operatorname{argmax}_{\mathbf{y}} \theta^{\top} \mathbf{f}(\mathbf{x}, \mathbf{y})$$

Minimum Error Rate Training (MERT)

$$\max_{\theta} \text{BLEU}(\text{references}, \{\text{translate}(\mathbf{x}_i, \theta)\}_{i=1}^N)$$

optimizing this objective is intractable in general – how can we do it?

generate k-best lists of translations,
approximately optimize on k-best lists,
repeat with new parameters
(pool k-best lists across iterates)

MERT

given a k-best list,

randomly choose a feature (e.g., reverse translation model) and vary a step size δ :

$$\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) + \delta \log P(\mathbf{y} | \mathbf{x})$$

find value of δ that maximizes BLEU

MERT

more generally,

randomly choose a “search direction” ψ and

vary a scalar multiplier δ :

$$\theta^\top \mathbf{f}(x, y) + \delta \psi^\top \mathbf{f}(x, y)$$

MERT

more generally,

randomly choose a “search direction” ψ and vary a scalar multiplier δ :

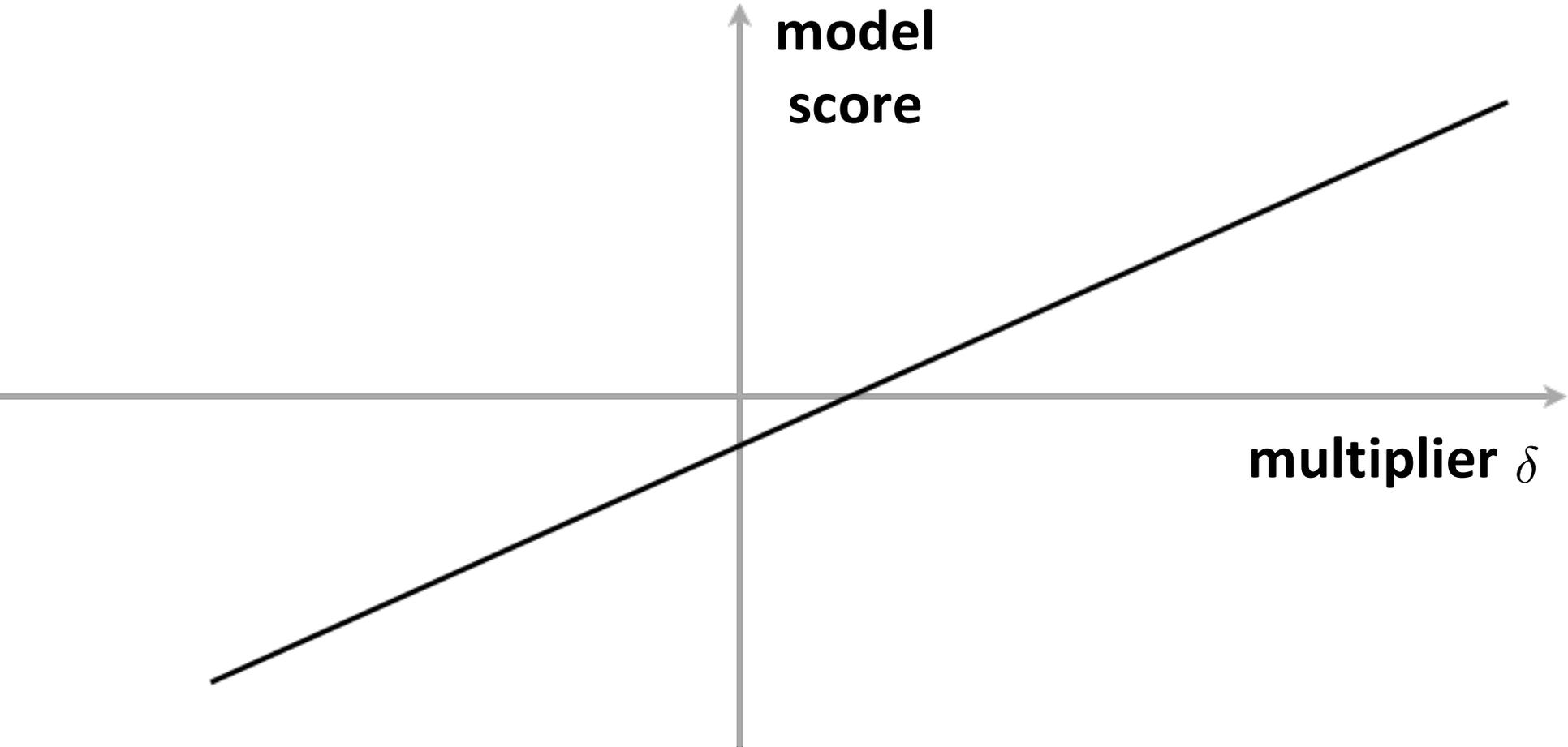
$$\theta^\top \mathbf{f}(x, y) + \delta \psi^\top \mathbf{f}(x, y)$$

single-feature special case: let ψ be a “one-hot” vector (all zeroes except a single 1)

for the reverse translation model feature one-hot vector ψ , this becomes:

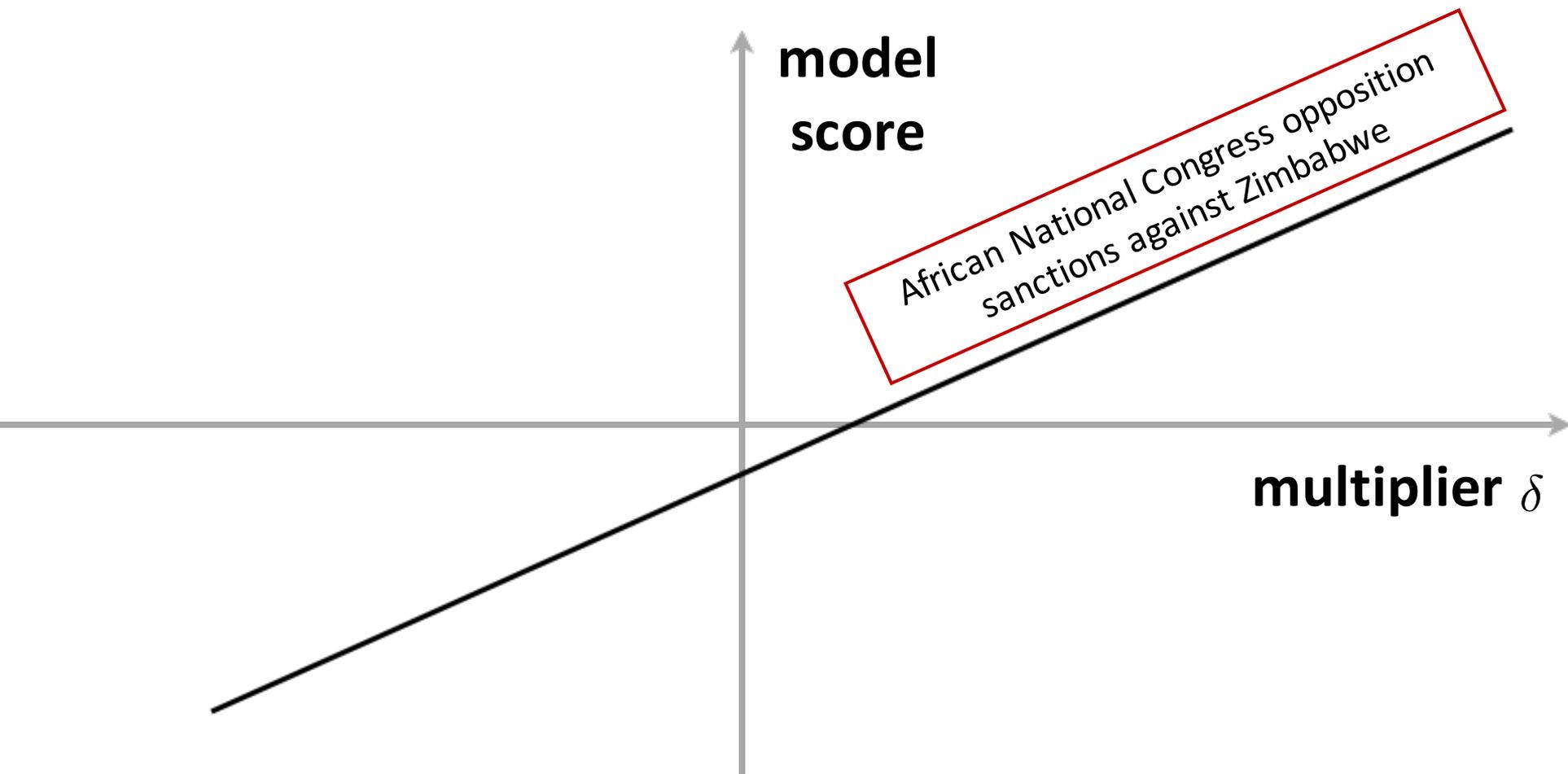
$$\theta^\top \mathbf{f}(x, y) + \delta \log P(\mathbf{y} | \mathbf{x})$$

MERT



each line is a translation
from the k-best list:

MERT



MERT

$$\theta^\top \mathbf{f}(x, y) + \delta \log P(y | x)$$

**model
score**

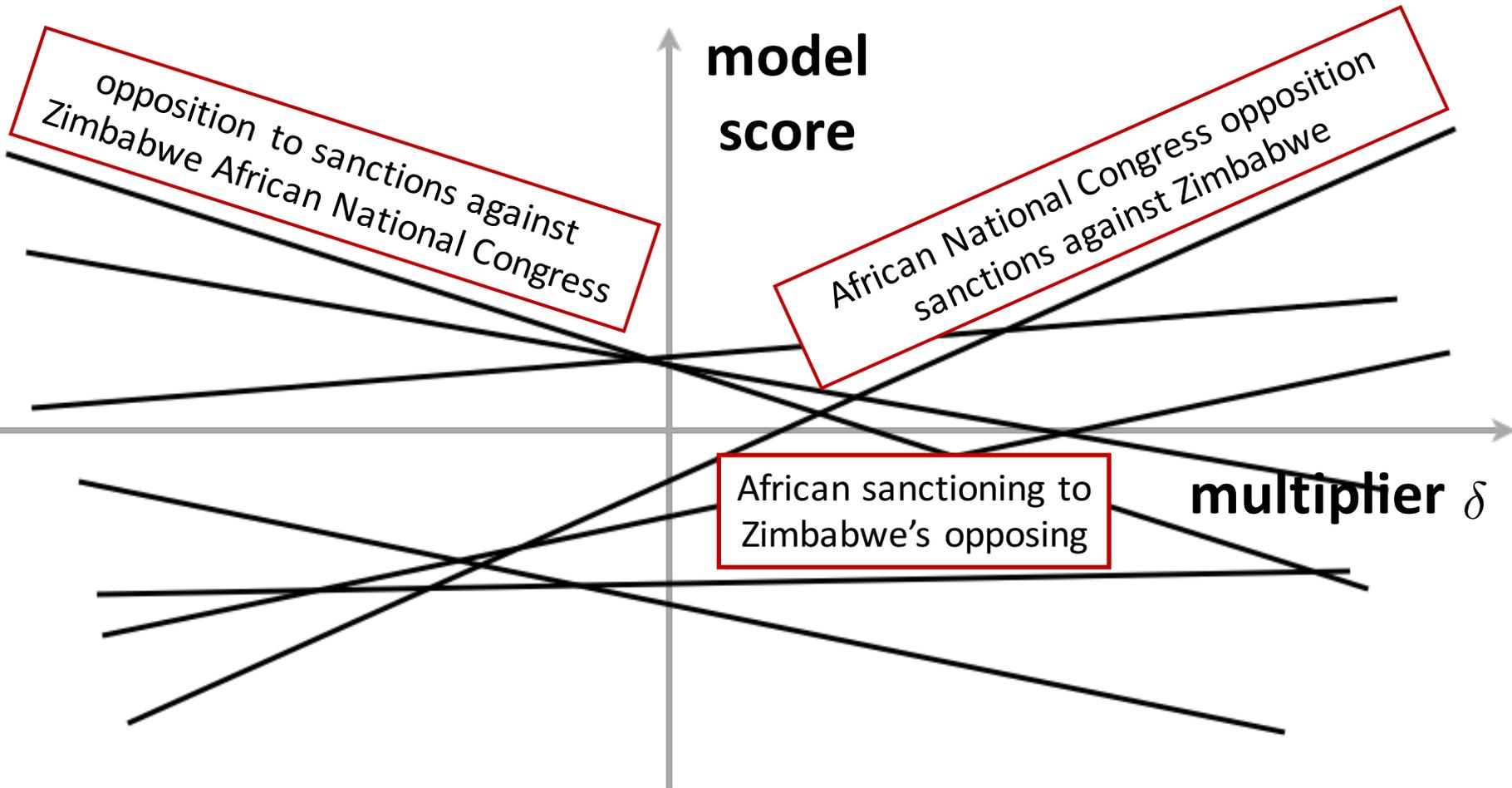
African National Congress opposition
sanctions against Zimbabwe

multiplier δ

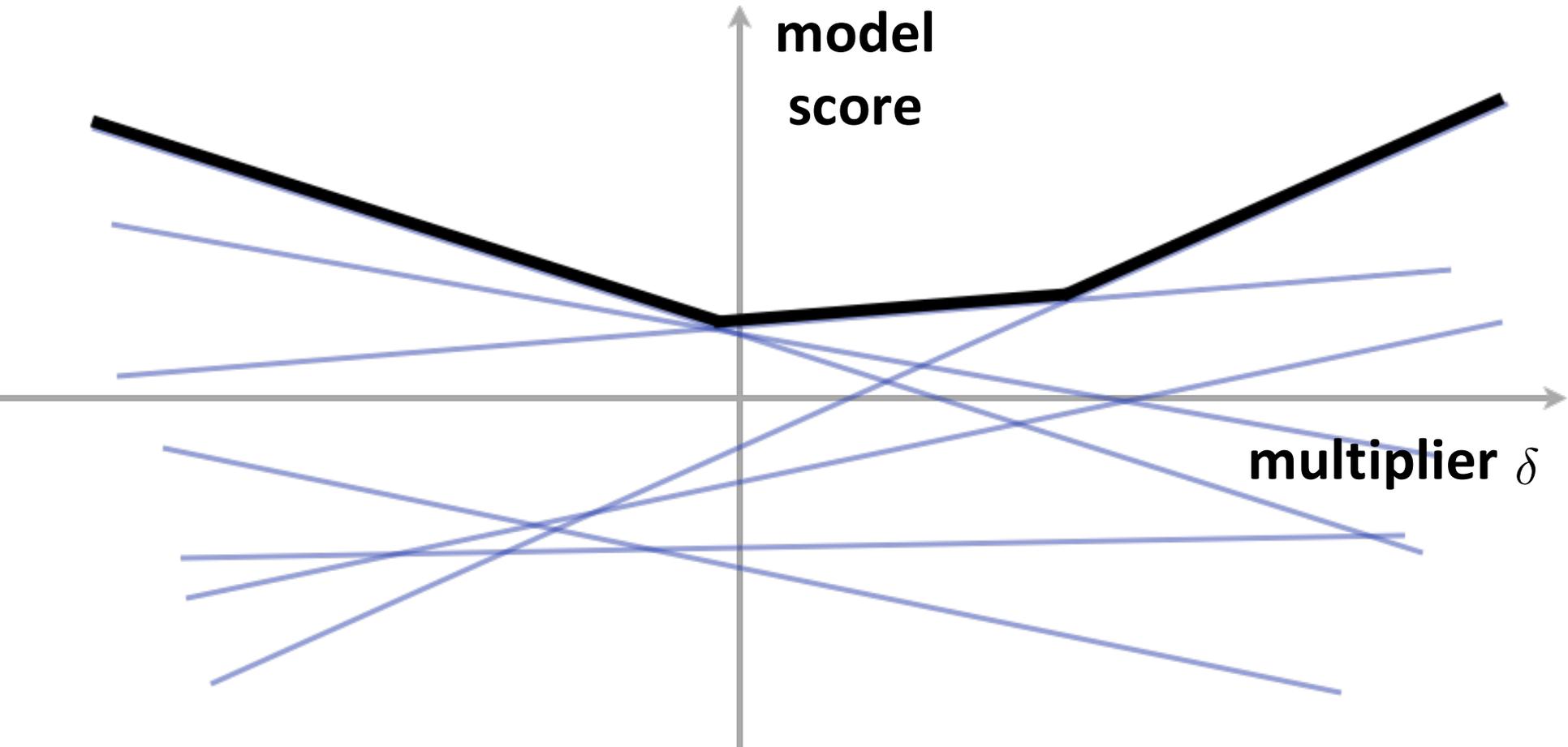
as we vary δ , the model score
of the translation changes

each line is a translation
from the k-best list:

MERT

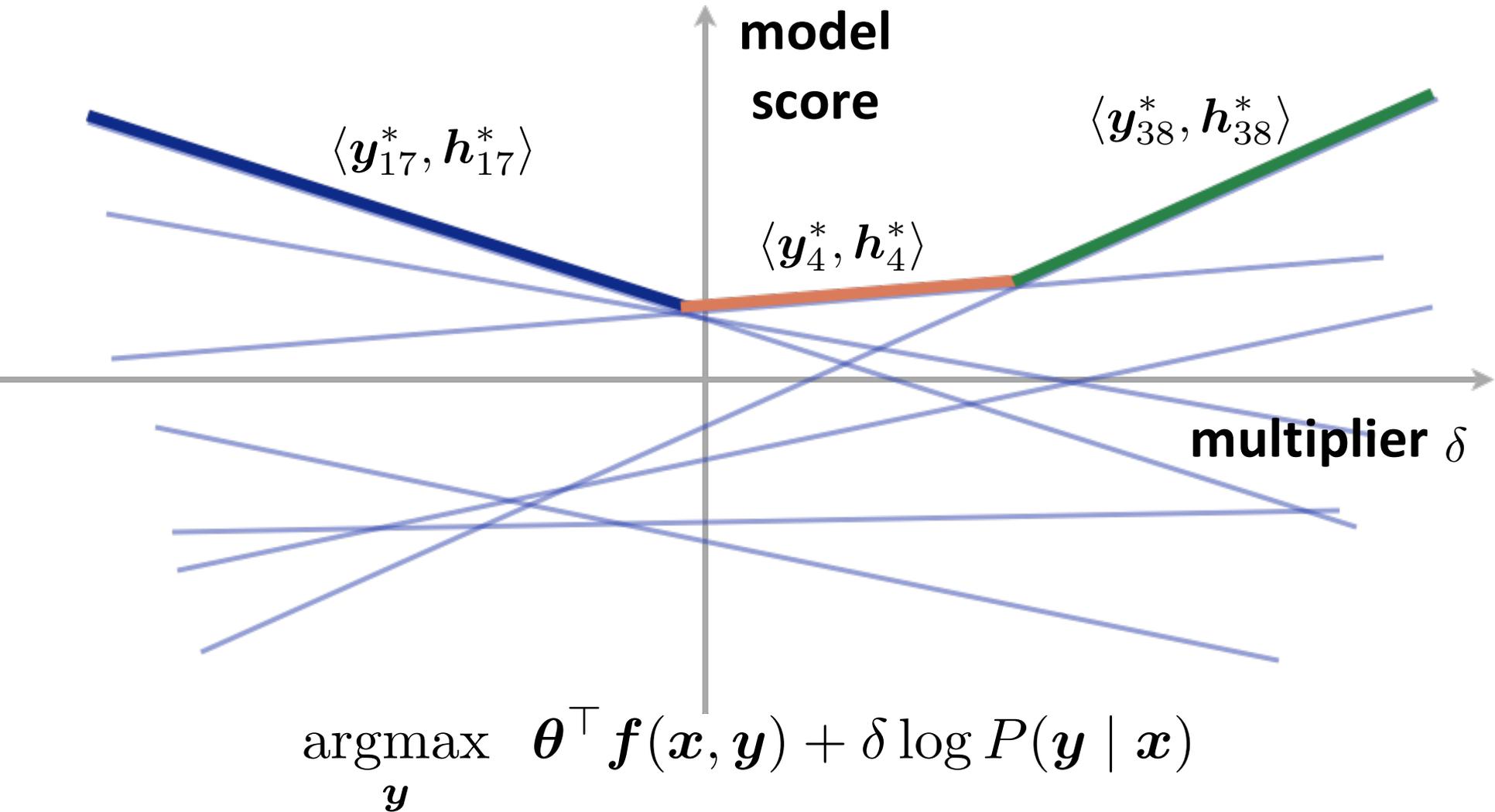


MERT

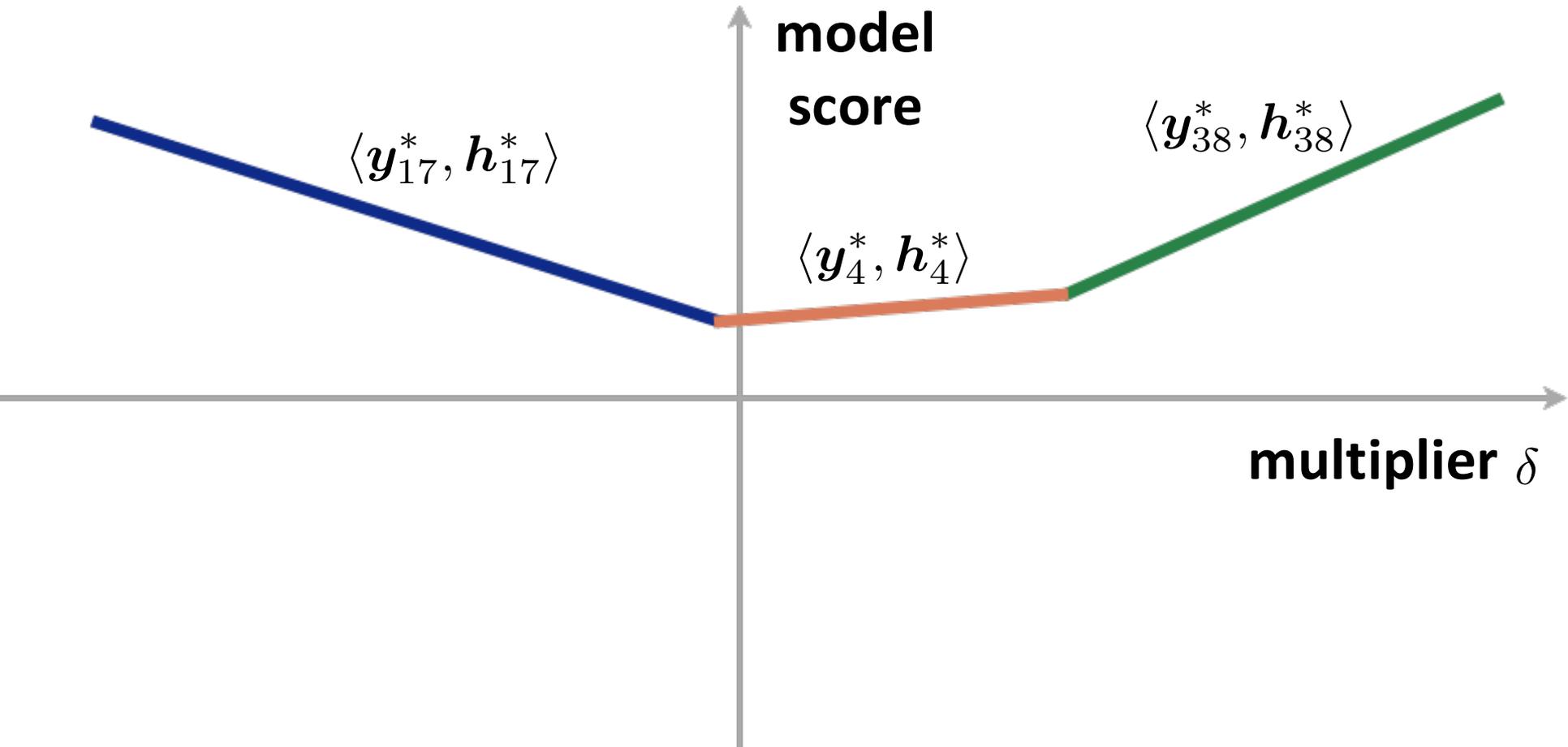


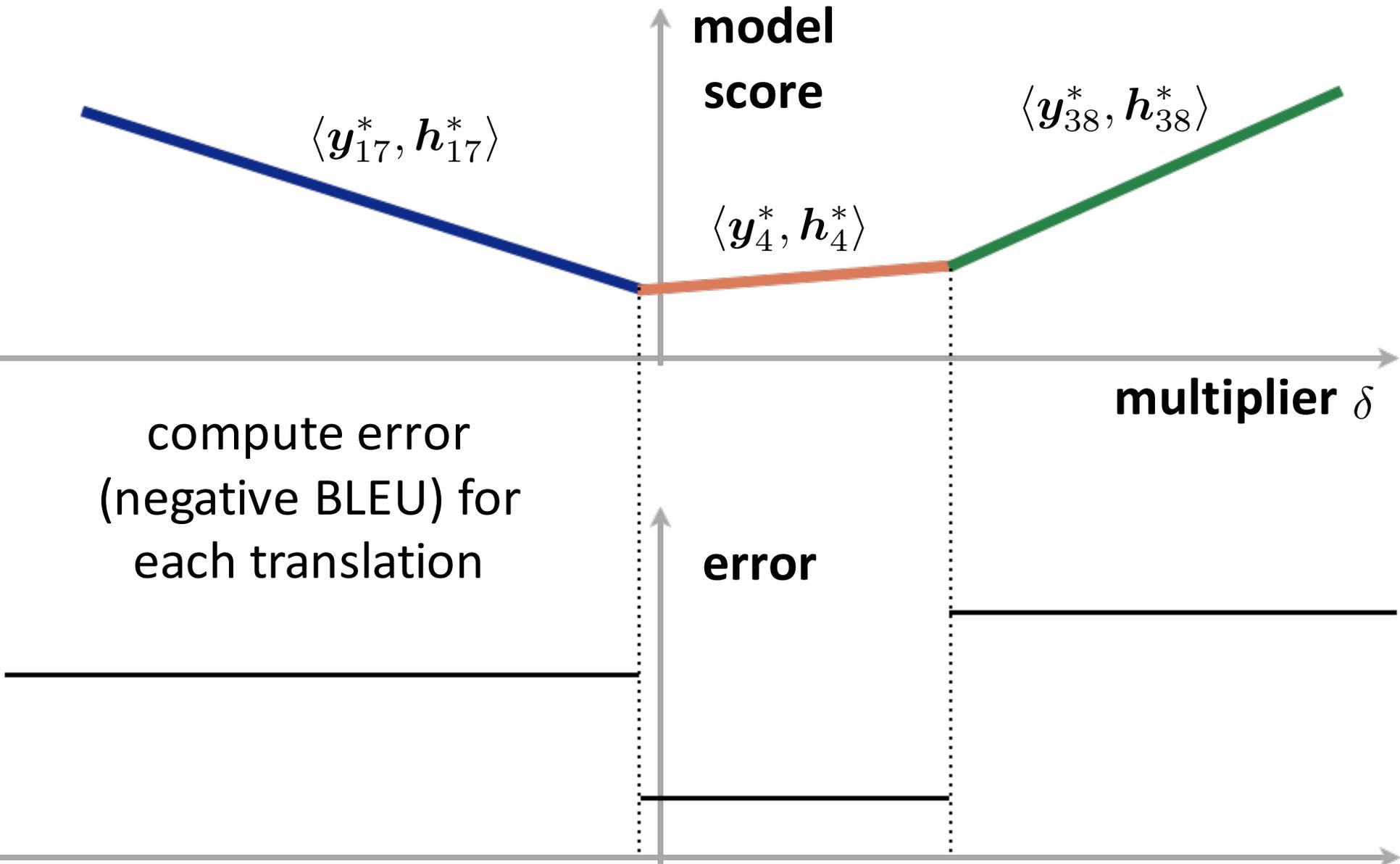
$$\max_y \theta^\top f(x, y) + \delta \log P(y | x)$$

MERT

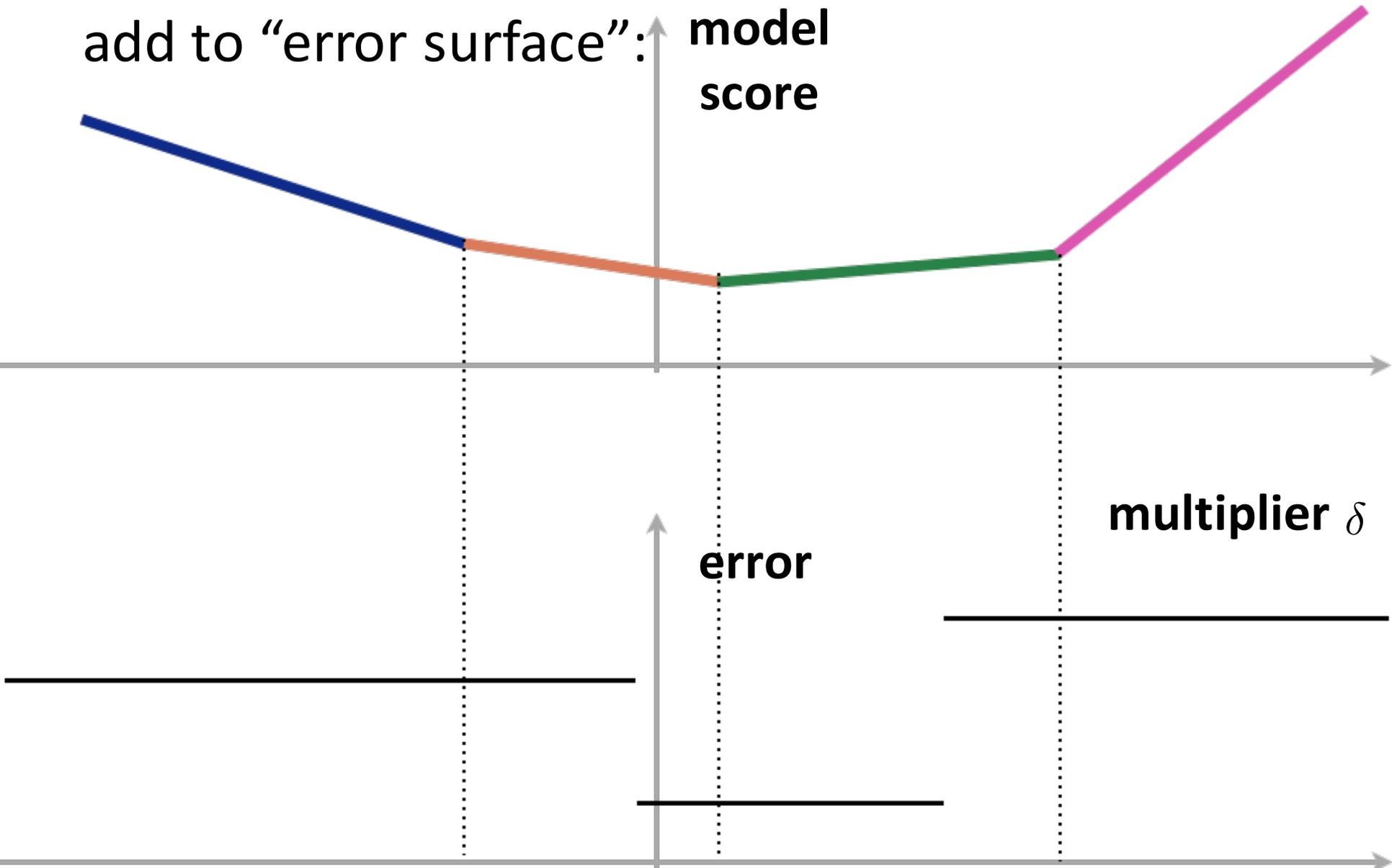


MERT

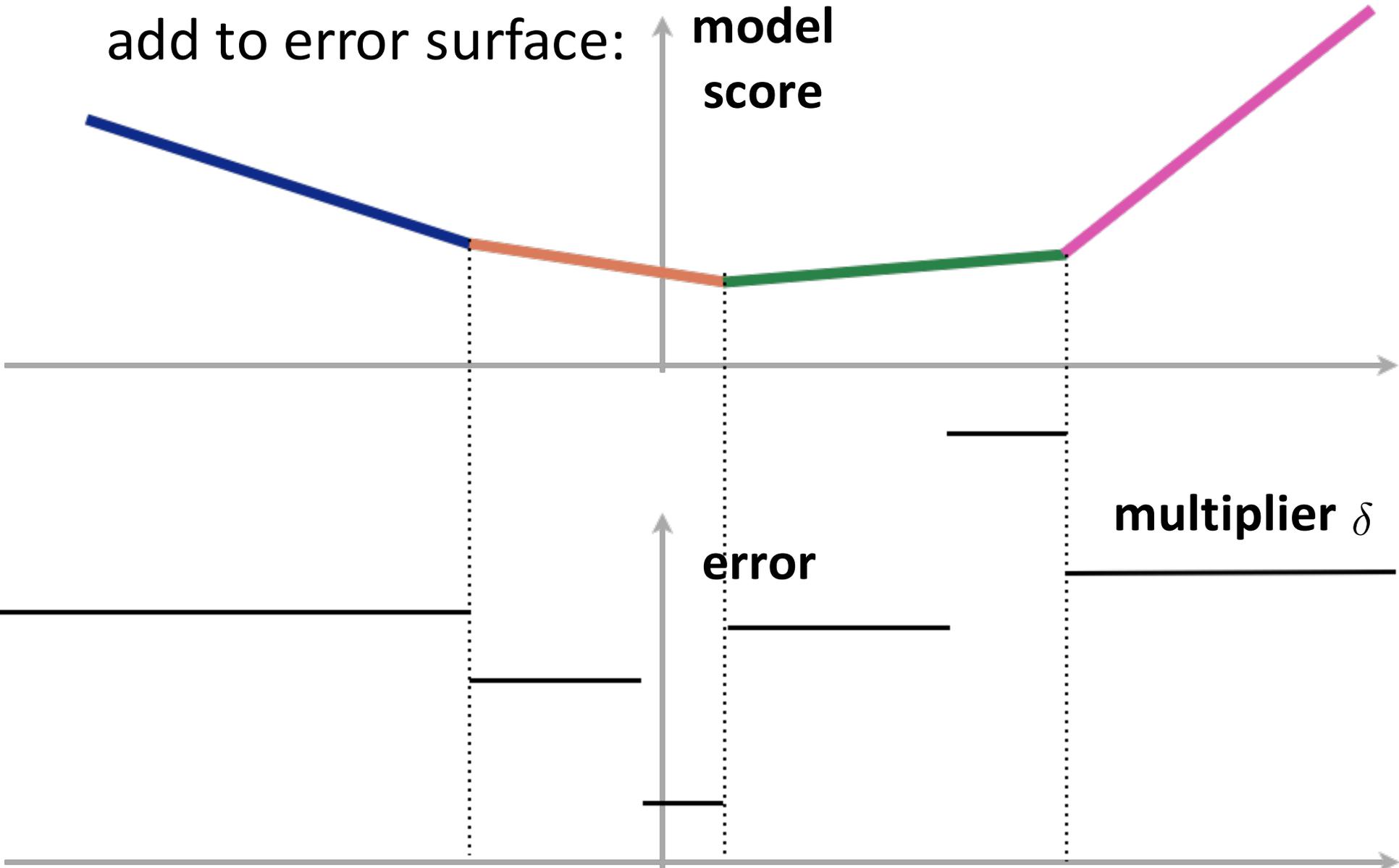


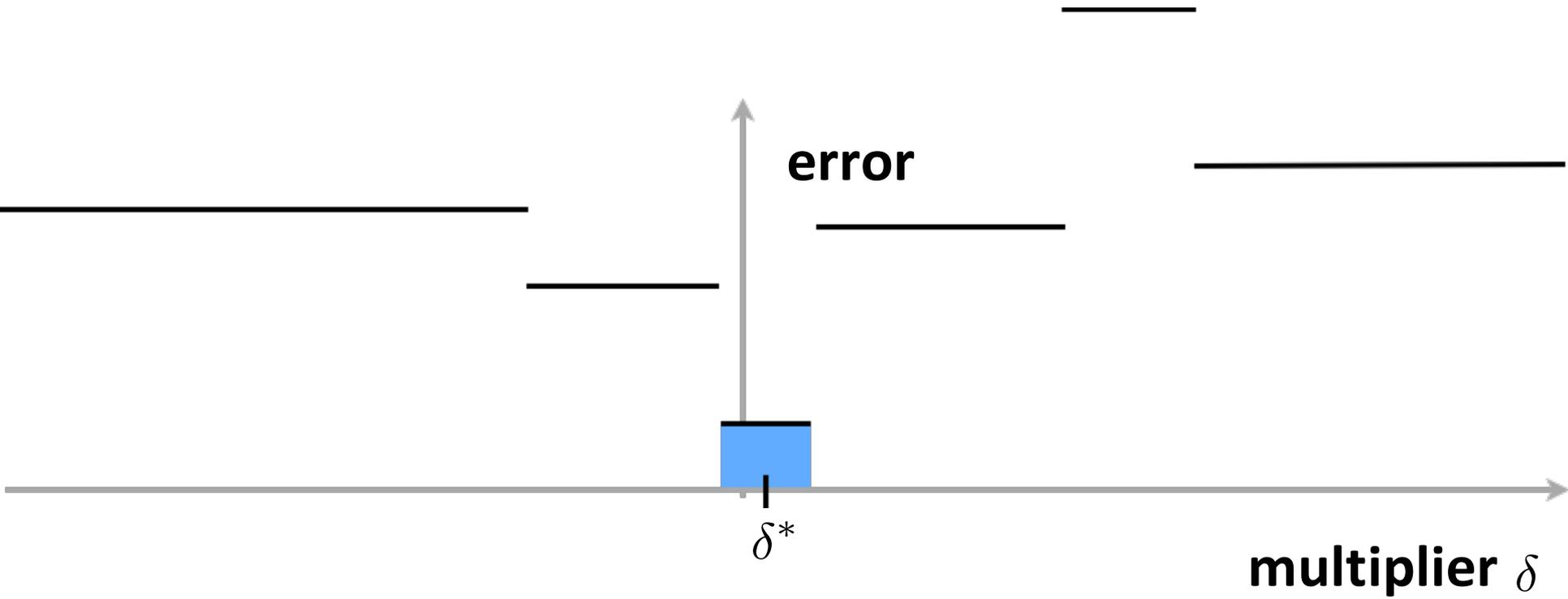


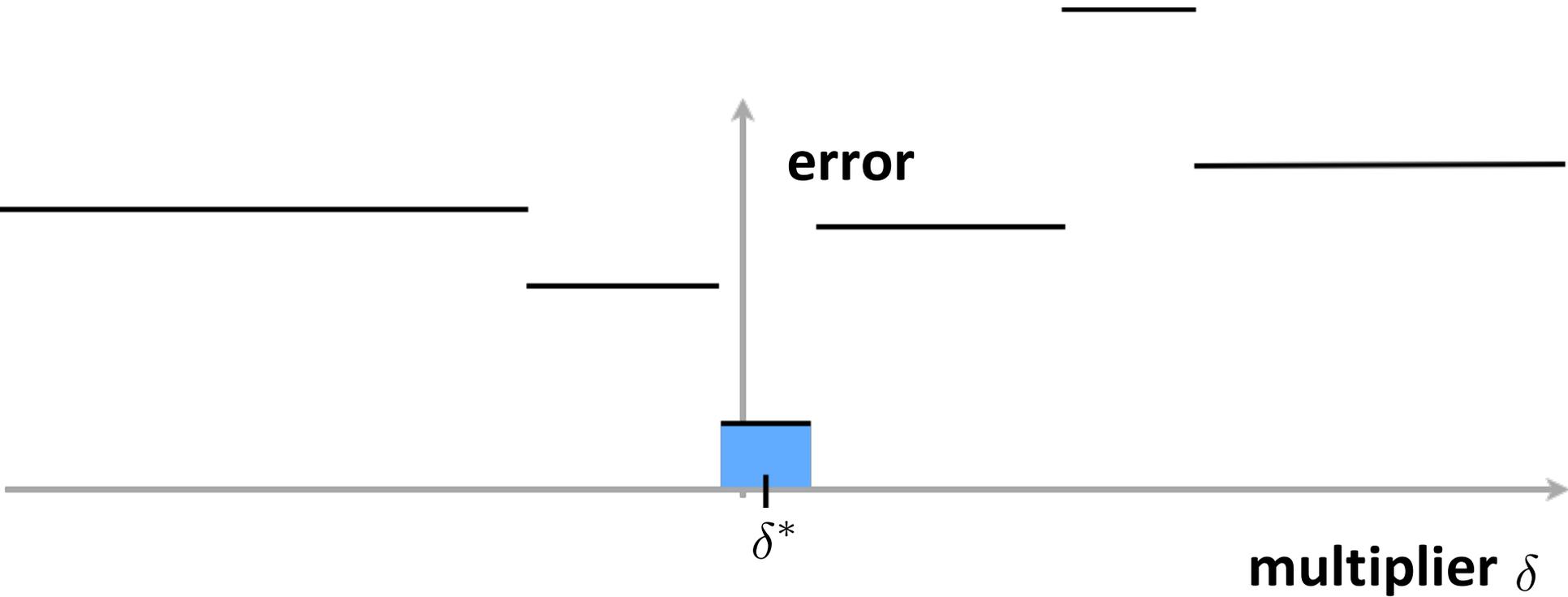
repeat for each sentence,
add to "error surface":



repeat for each sentence,
add to error surface:







$$\theta^{\text{new}} = \theta + \delta^* \psi$$

Error Surface (Smith and Eisner, 2006)

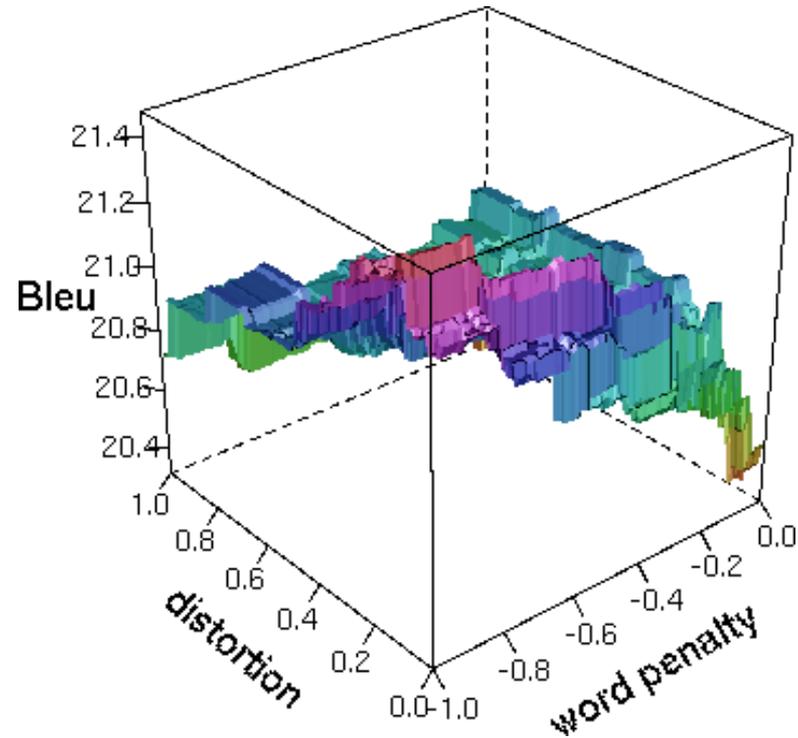
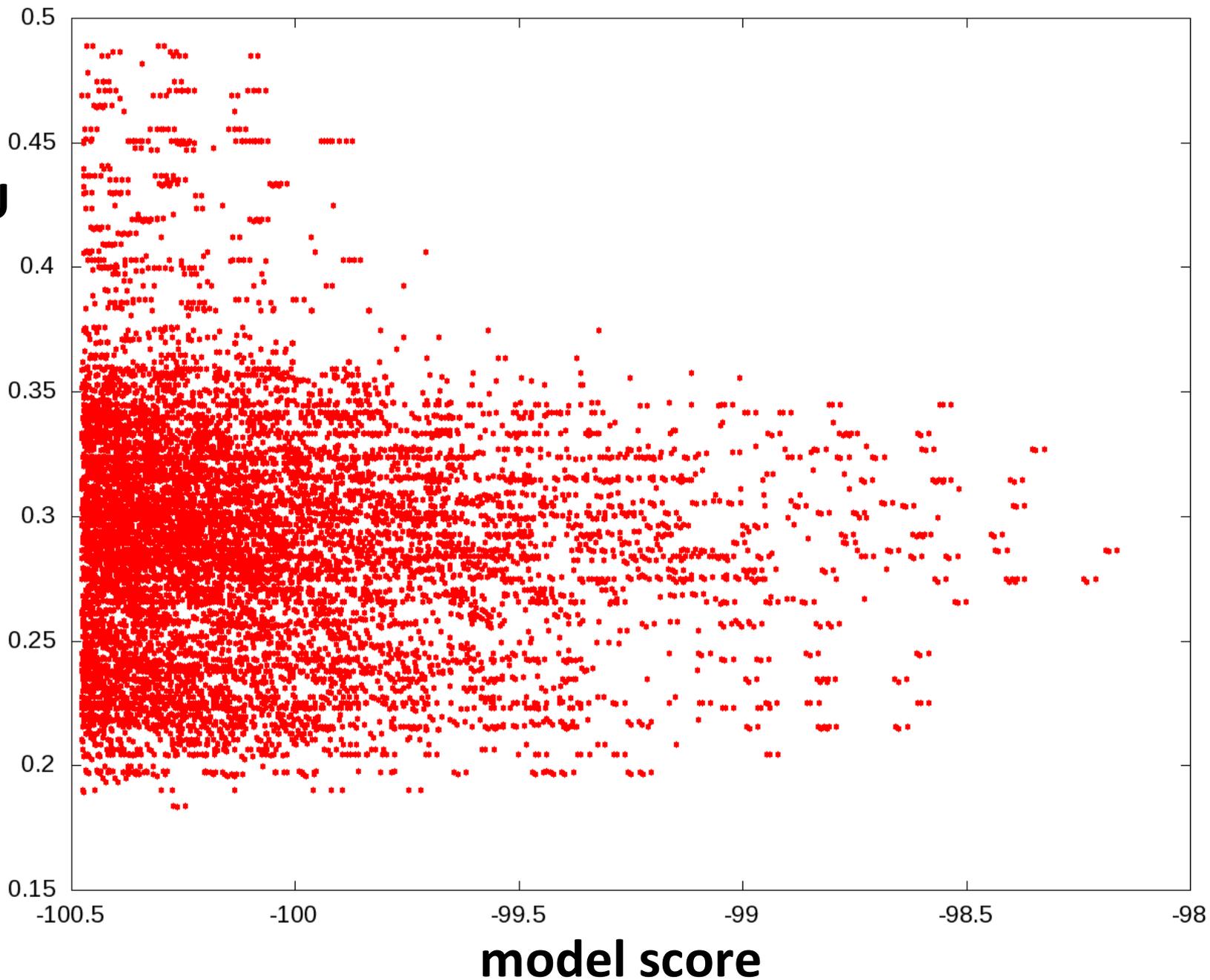


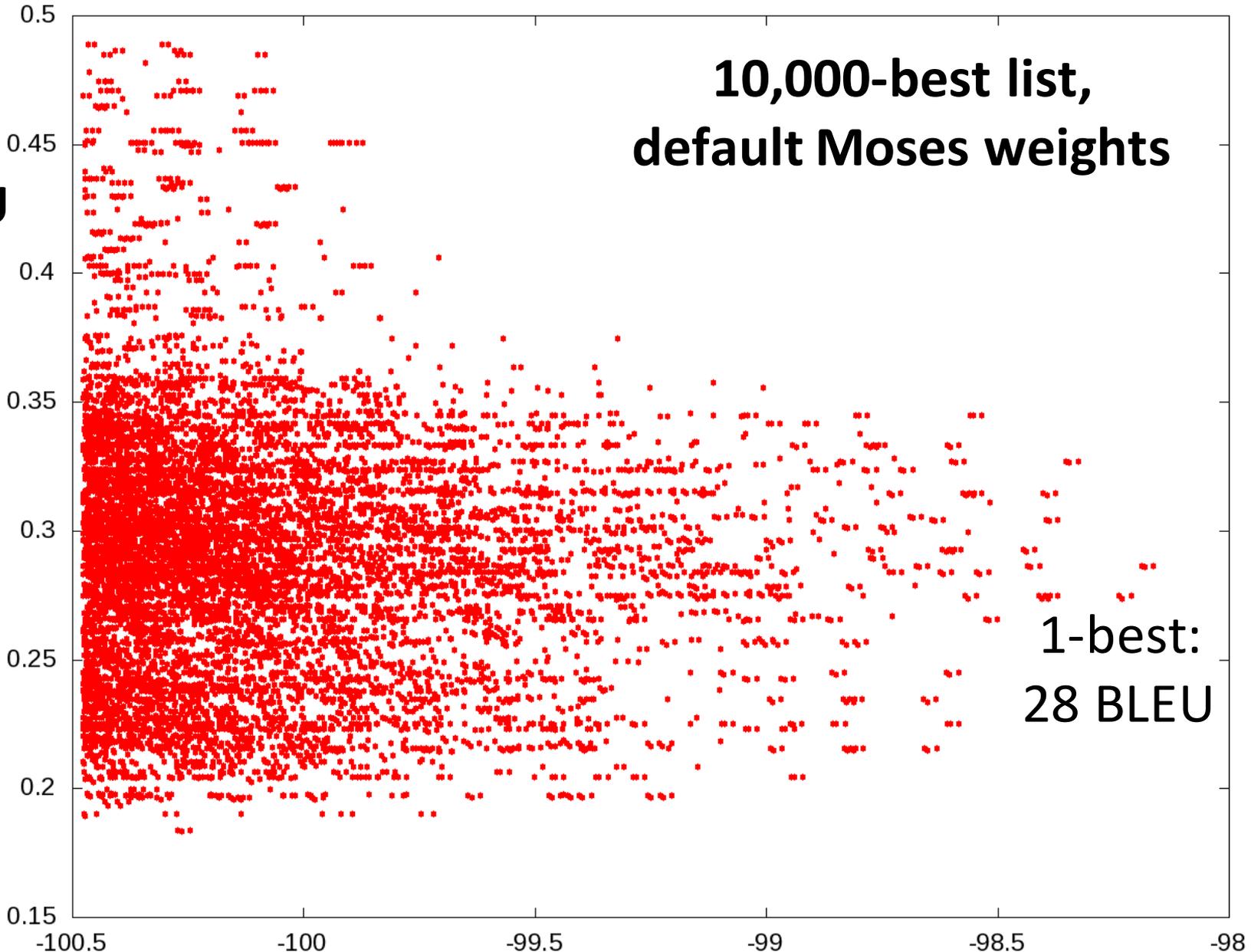
Figure 1: The loss surface for a machine translation system: while other parameters are held constant, we vary the weights on the distortion and word penalty features. Note the piecewise constant regions with several local maxima.

BLEU



BLEU

**10,000-best list,
default Moses weights**



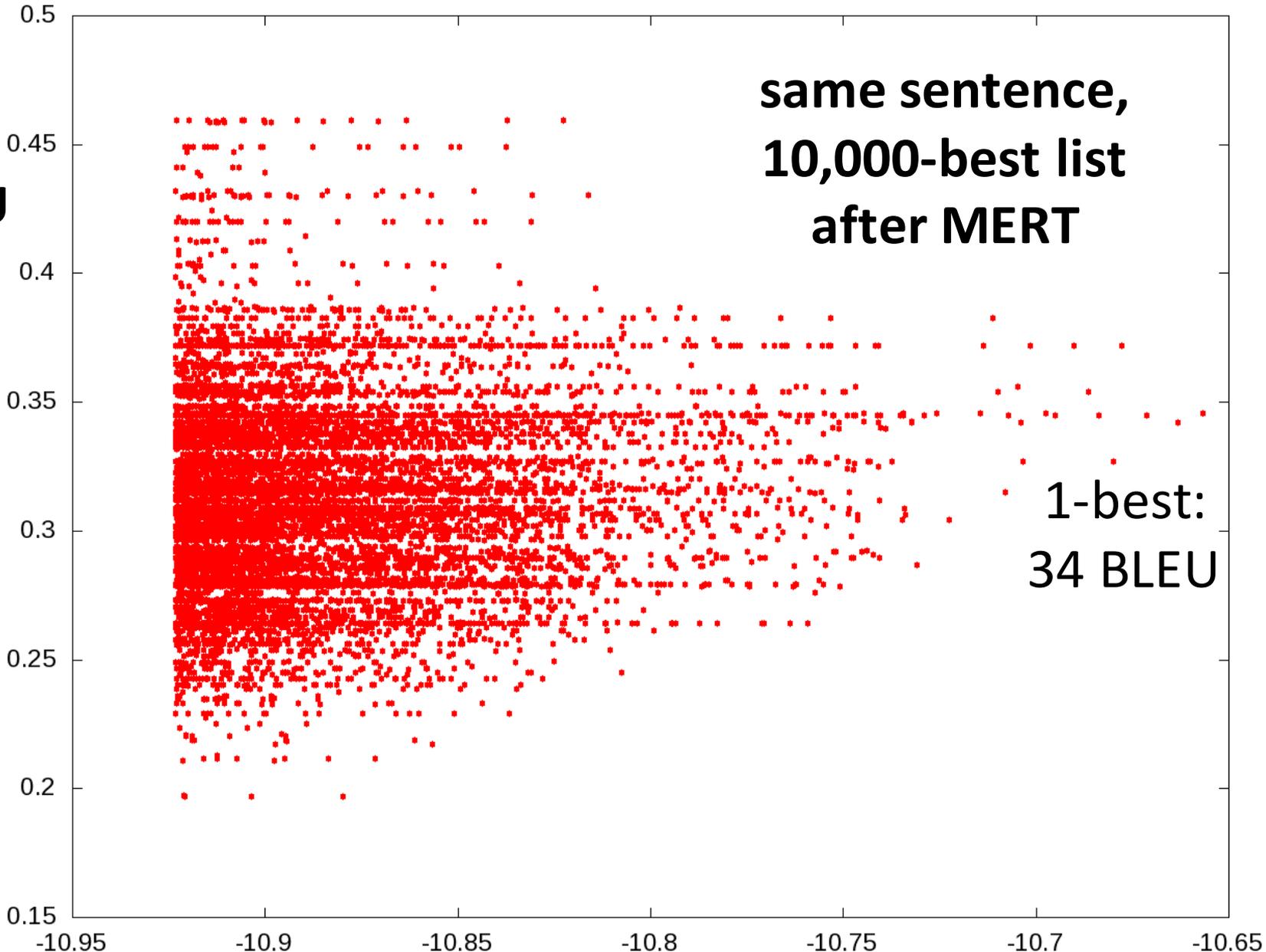
**1-best:
28 BLEU**

model score

BLEU

**same sentence,
10,000-best list
after MERT**

1-best:
34 BLEU

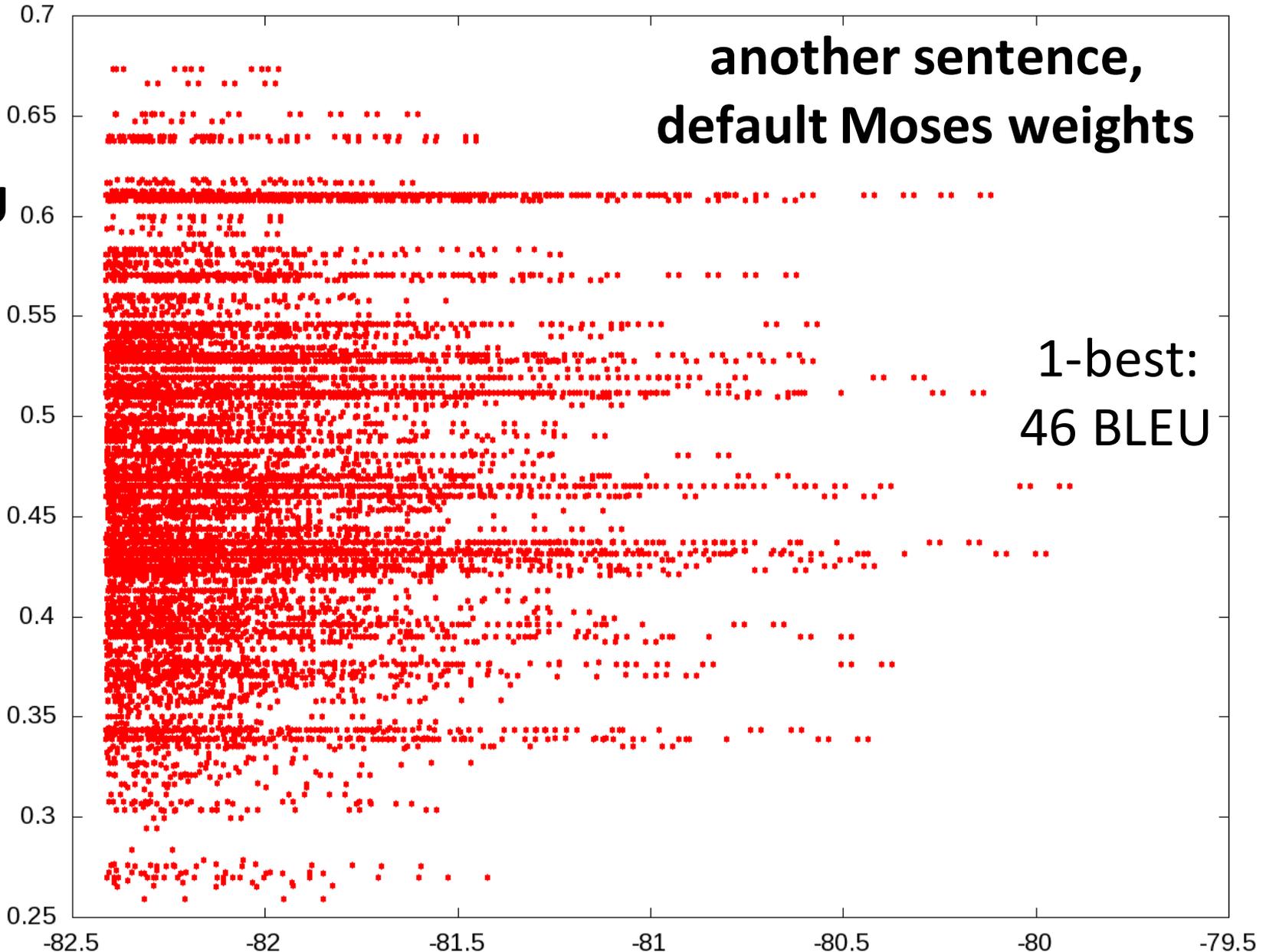


model score

BLEU

**another sentence,
default Moses weights**

1-best:
46 BLEU

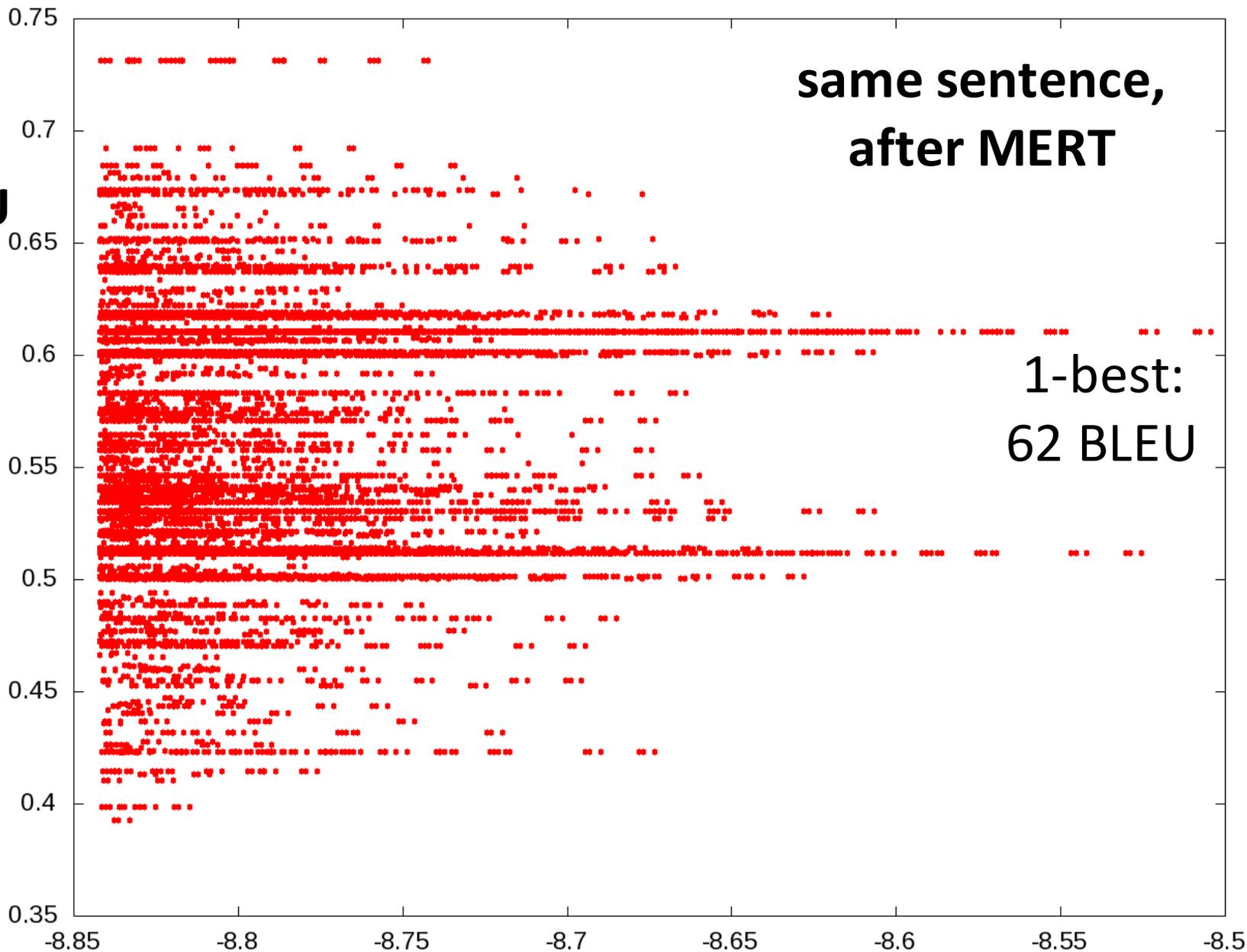


model score

BLEU

**same sentence,
after MERT**

**1-best:
62 BLEU**



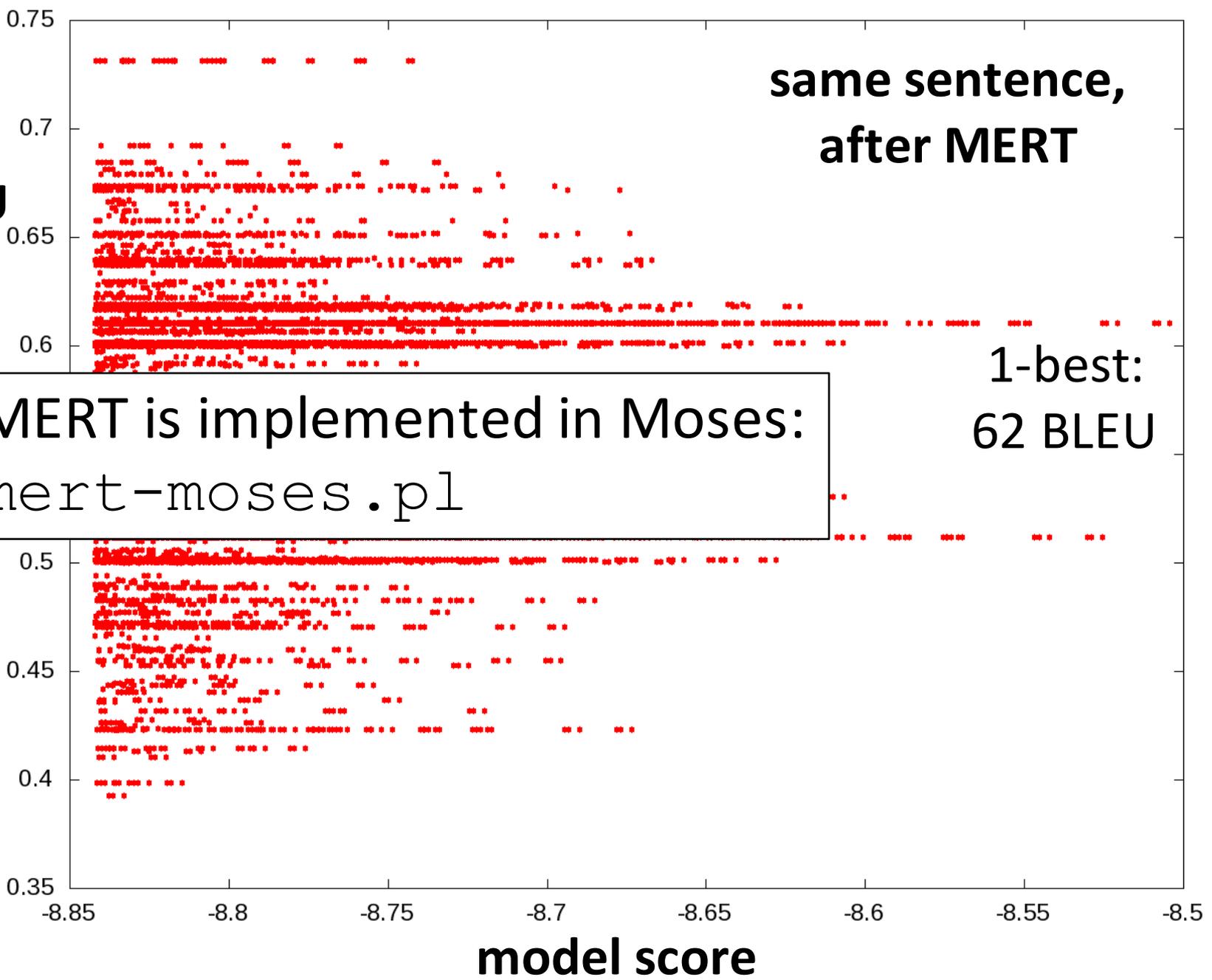
model score

BLEU

**same sentence,
after MERT**

**1-best:
62 BLEU**

MERT is implemented in Moses:
`mert-moses.pl`



$$\max_{\theta} \text{BLEU}(\text{references}, \{\text{translate}(\mathbf{x}_i, \theta)\}_{i=1}^N)$$

Problems with the MERT objective function?

Discontinuous & non-convex \rightarrow optimization relies on randomized search

No regularization \rightarrow frequently overfits to tuning set

As a result, MERT is only effective for very small models (<40 parameters)

Smoothing the Objective

$$\max_{\boldsymbol{\theta}} \text{BLEU} \left(\text{references}, \left\{ \underset{\mathbf{y}}{\operatorname{argmax}} \boldsymbol{\theta}^{\top} \mathbf{f}(\mathbf{x}^{(i)}, \mathbf{y}) \right\}_{i=1}^N \right)$$



first, convert to sentence-level objective:

$$\max_{\boldsymbol{\theta}} \sum_{i=1}^N \text{BLEU}^{+1} \left(\text{references}^{(i)}, \underset{\mathbf{y}}{\operatorname{argmax}} \boldsymbol{\theta}^{\top} \mathbf{f}(\mathbf{x}^{(i)}, \mathbf{y}) \right)$$



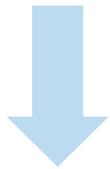
then, smooth using expectation under log-linear distribution:

$$\max_{\boldsymbol{\theta}} \sum_{i=1}^N \sum_{\mathbf{y}} \text{BLEU}^{+1} \left(\text{references}^{(i)}, \mathbf{y} \right) \frac{\exp\{\boldsymbol{\theta}^{\top} \mathbf{f}(\mathbf{x}^{(i)}, \mathbf{y})\}^{\alpha}}{\sum_{\mathbf{y}'} \exp\{\boldsymbol{\theta}^{\top} \mathbf{f}(\mathbf{x}^{(i)}, \mathbf{y}')\}^{\alpha}}$$

alpha = “smoothness factor”

Smoothing the Objective

$$\max_{\theta} \text{BLEU} \left(\text{references}, \left\{ \underset{\mathbf{y}}{\operatorname{argmax}} \theta^{\top} \mathbf{f}(\mathbf{x}^{(i)}, \mathbf{y}) \right\}_{i=1}^N \right)$$



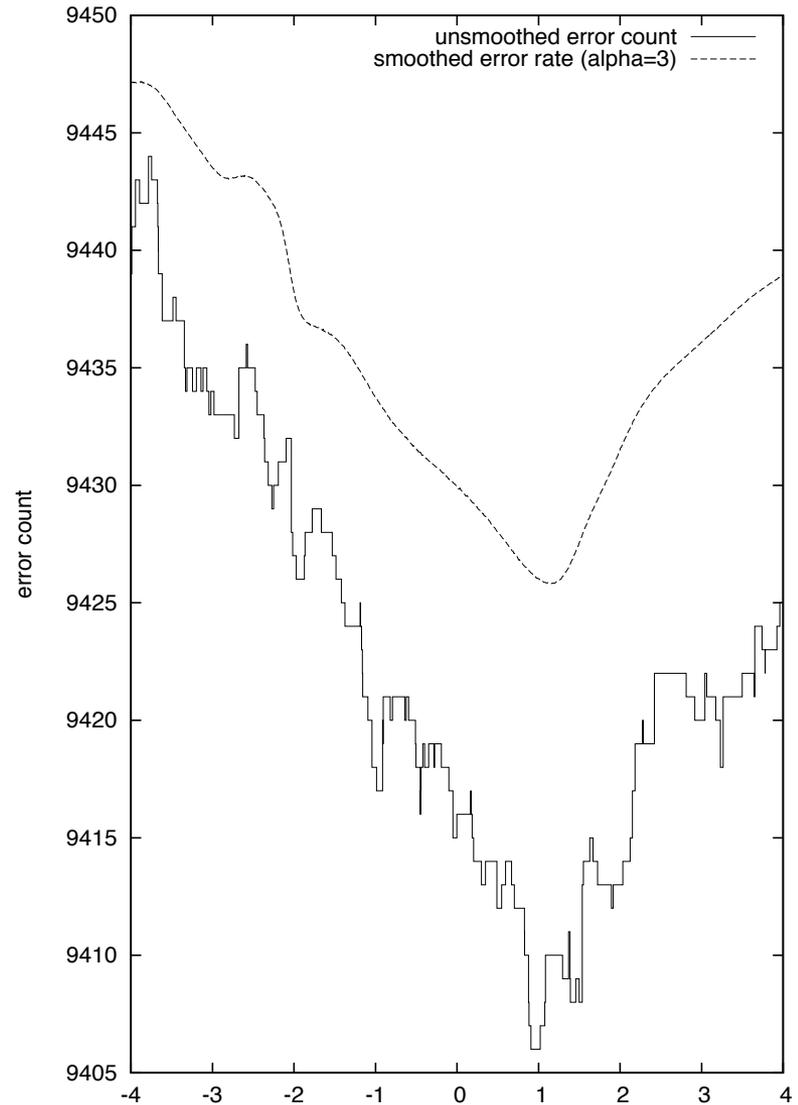
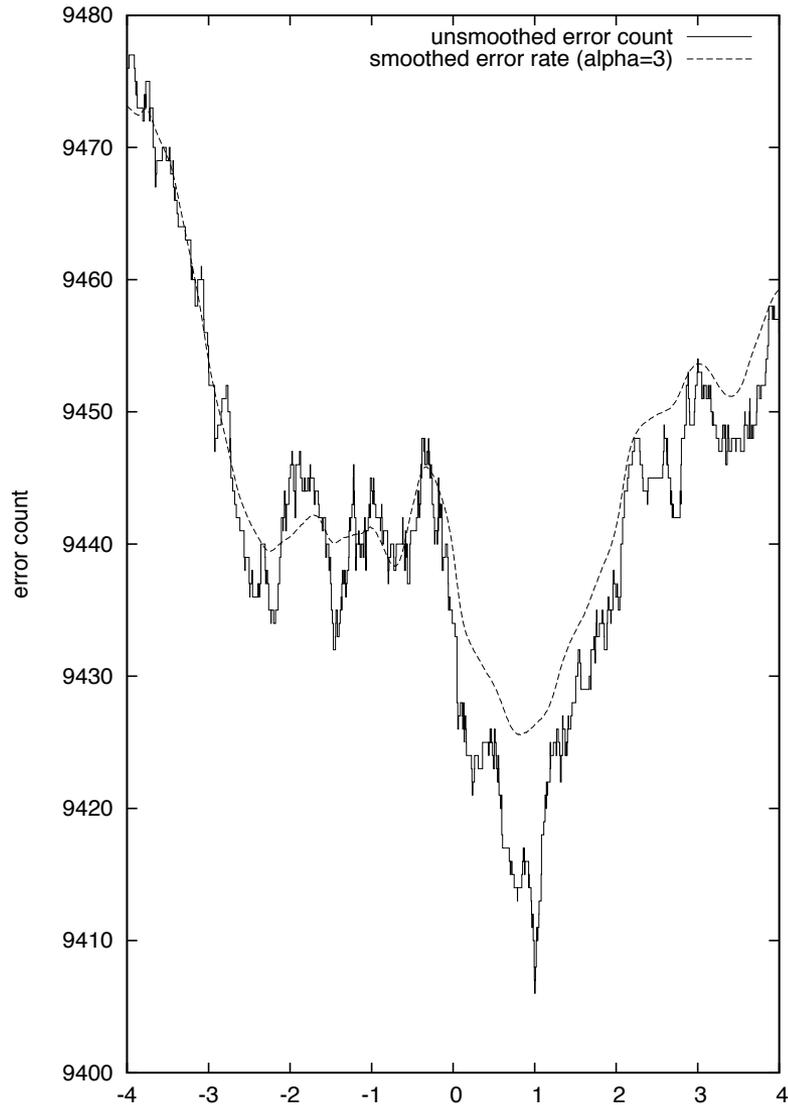
first, convert to sentence-level objective:

**This is often
called “risk” or
“Bayes risk”**

$$\max_{\theta} \sum_{i=1}^N \text{BLEU} \left(\text{references}^{(i)}, \underset{\mathbf{y}}{\operatorname{argmax}} \theta^{\top} \mathbf{f}(\mathbf{x}^{(i)}, \mathbf{y}) \right)$$
$$\max_{\theta} \sum_{i=1}^N \sum_{\mathbf{y}} \text{BLEU}^{+1} \left(\text{references}^{(i)}, \mathbf{y} \right) \frac{\exp\{\theta^{\top} \mathbf{f}(\mathbf{x}^{(i)}, \mathbf{y})\}^{\alpha}}{\sum_{\mathbf{y}'} \exp\{\theta^{\top} \mathbf{f}(\mathbf{x}^{(i)}, \mathbf{y}')\}^{\alpha}}$$

alpha = “smoothness factor”

Smoothing Error Surfaces (Och, 2003)



Many other researchers tried to improve MERT:

Regularization and Search for MERT (Cer et al., 2008)

Random Restarts in MERT for MT (Moore & Quirk, 2008)

Stabilizing MERT (Foster & Kuhn, 2009)

Issues remain:

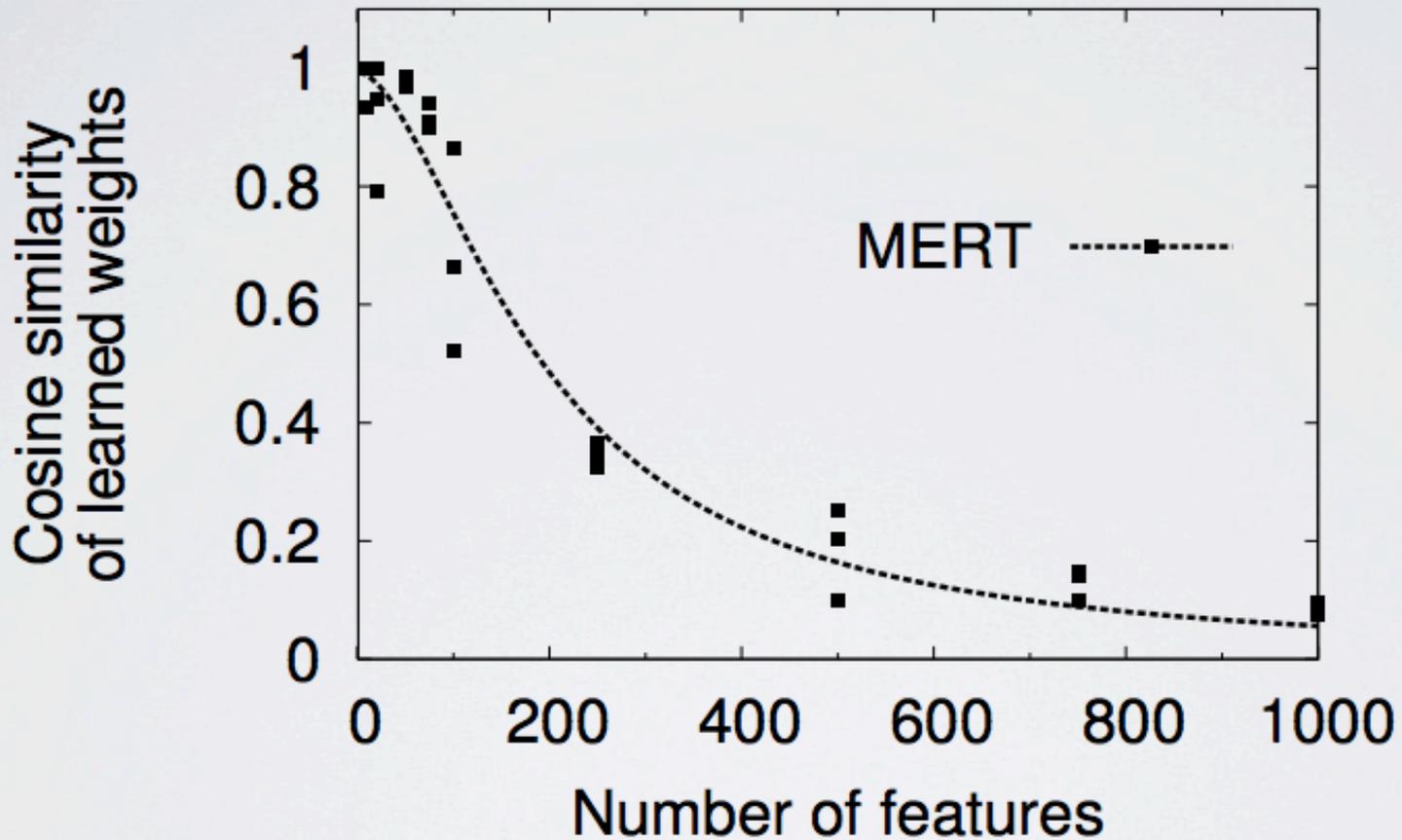
Better Hypothesis Testing for Statistical MT: Controlling for Optimizer Instability (Clark et al., 2011)

They suggest running MERT 3-5 times due to its instability



MERT *doesn't scale*

Synthetic weight learning of MERT

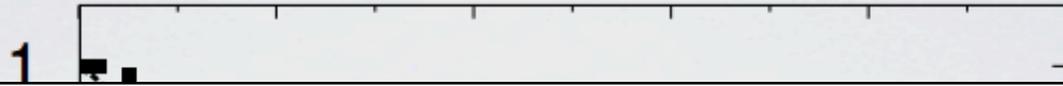


The synthetic experiment in ideal conditions validates what has long been accepted as truth



MERT *doesn't scale*

Synthetic weight learning of MERT



Tuning as Ranking

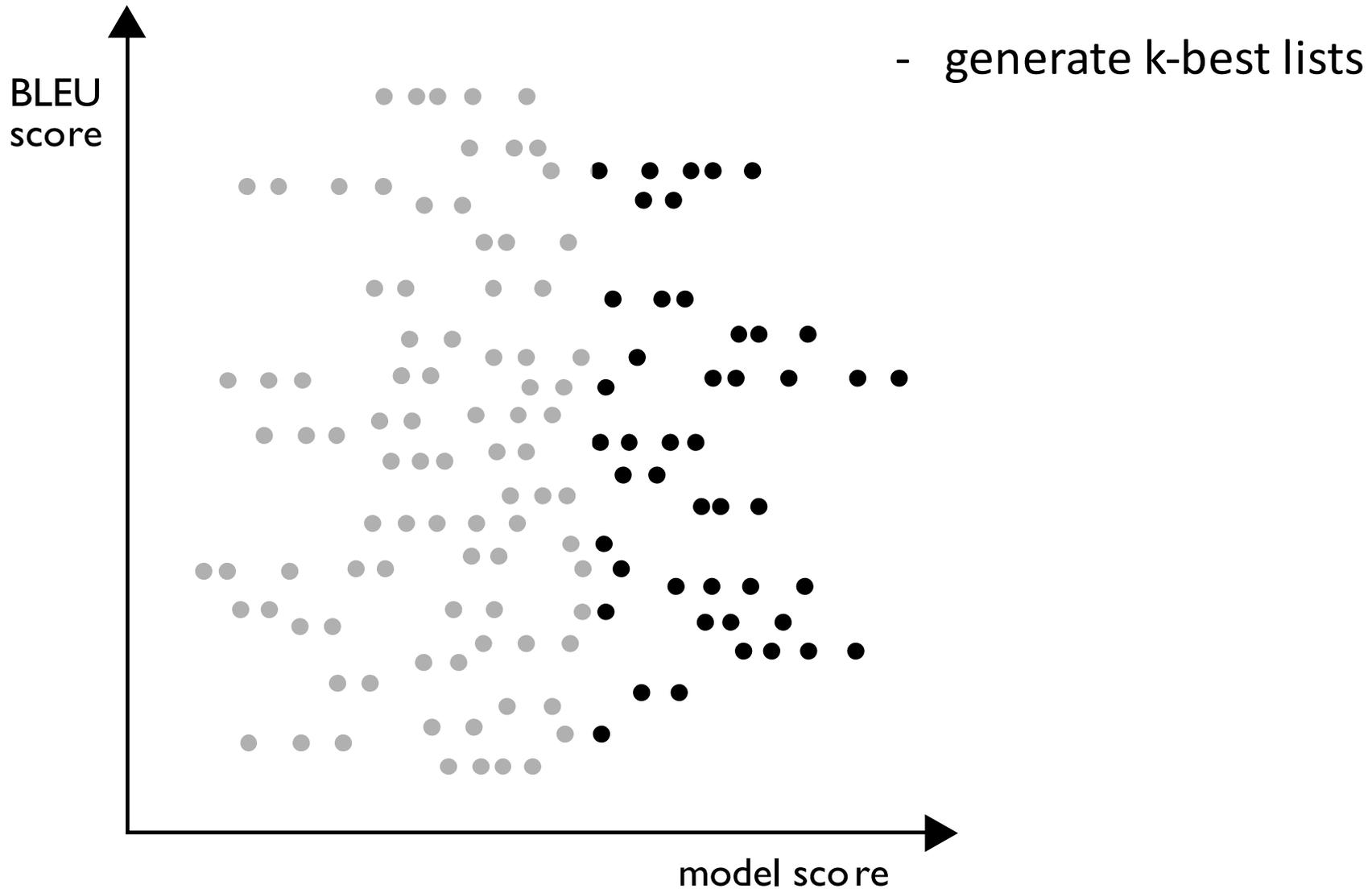
Mark Hopkins and Jonathan May
SDL Language Weaver
Los Angeles, CA 90045
{mhopkins, jmay}@sdl.com

0 200 400 600 800 1000

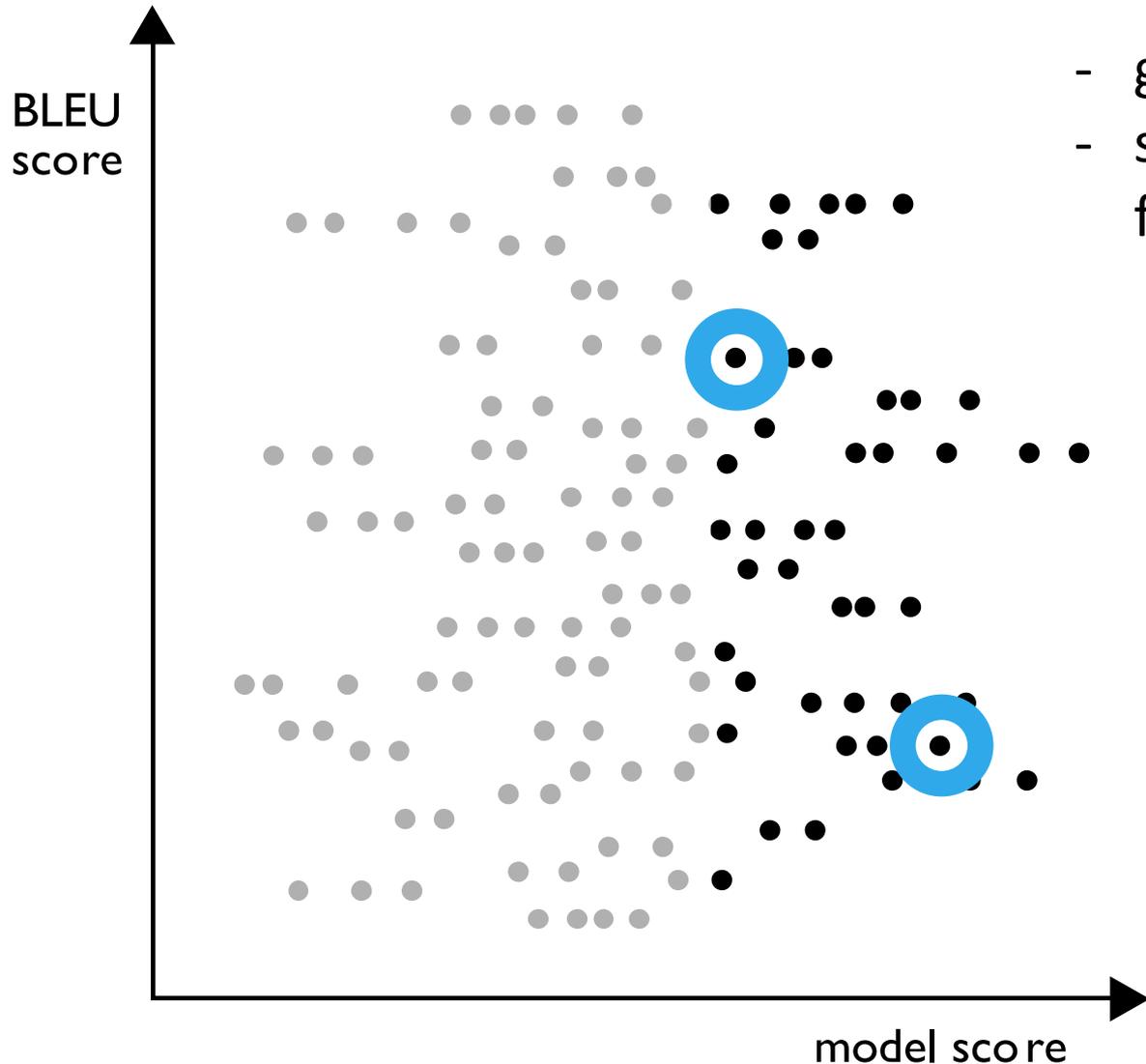
Number of features

The synthetic experiment in ideal conditions validates what has long been accepted as truth

Pairwise Ranking Optimization (Hopkins & May, 2011)



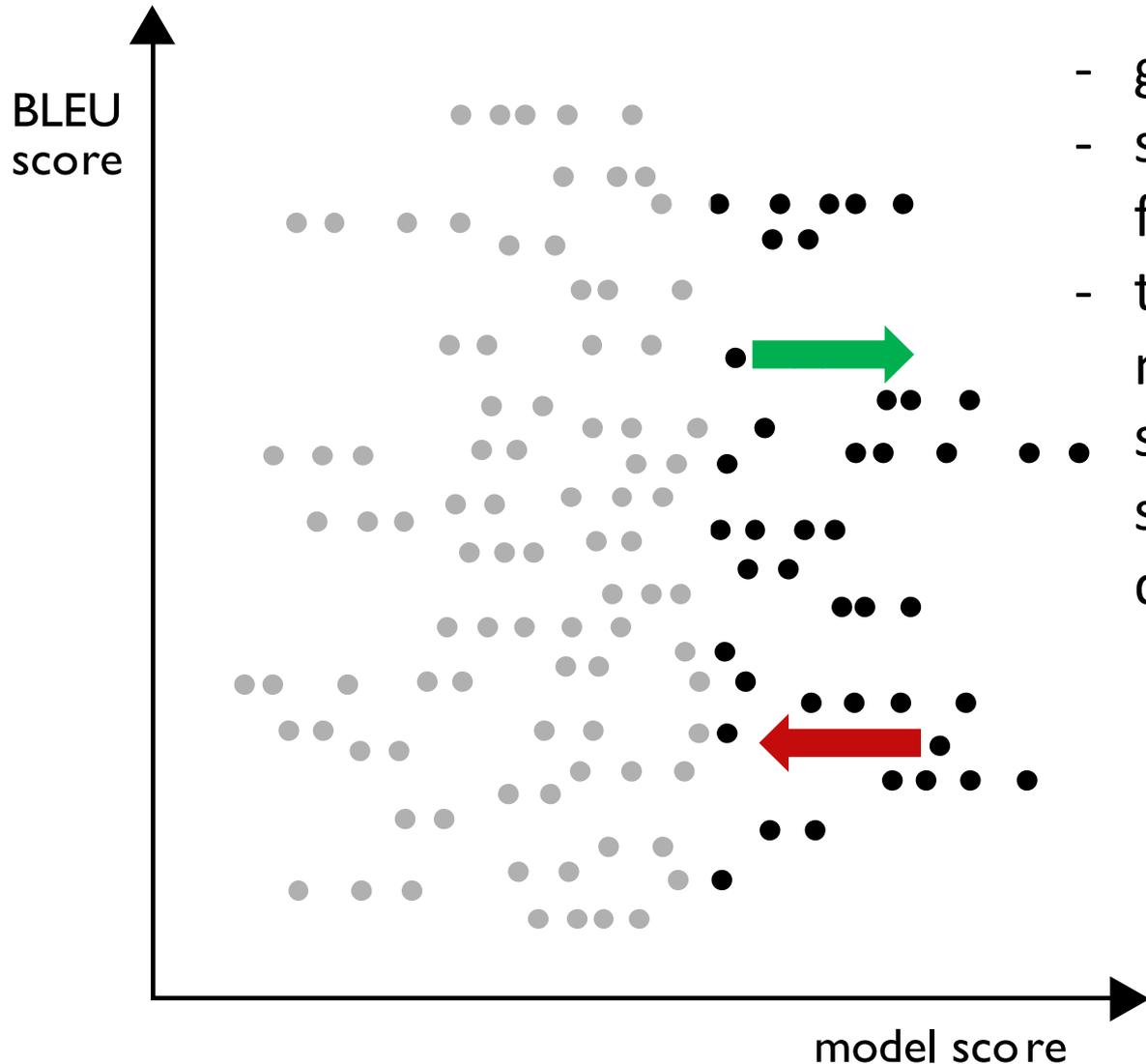
Pairwise Ranking Optimization (Hopkins & May, 2011)



- generate k-best lists
- sample translation pairs from k-best list

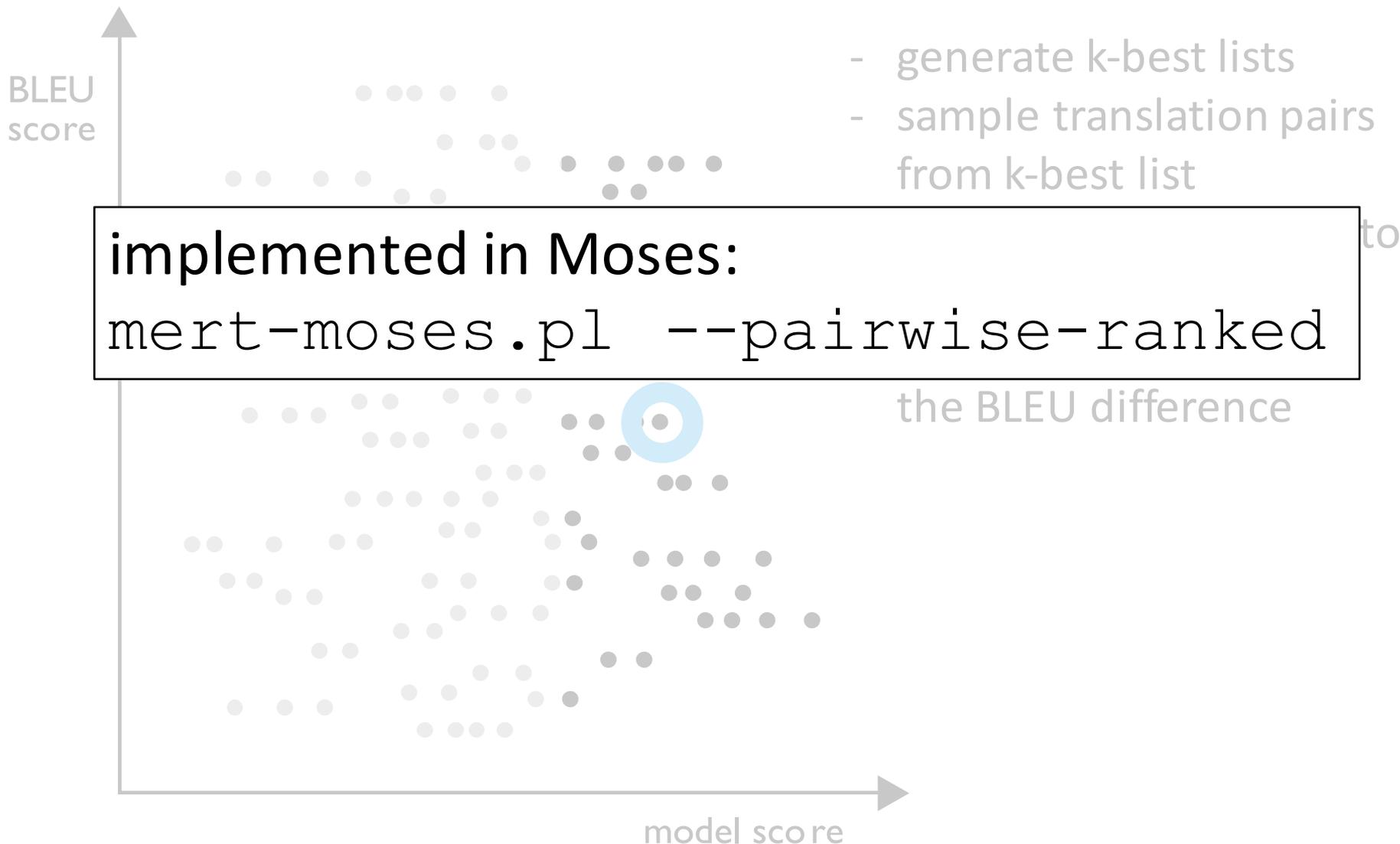
Pairwise Ranking Optimization

(Hopkins & May, 2011)



- generate k-best lists
- sample translation pairs from k-best list
- train a binary classifier to make the pairs' model score difference have same sign as the BLEU difference

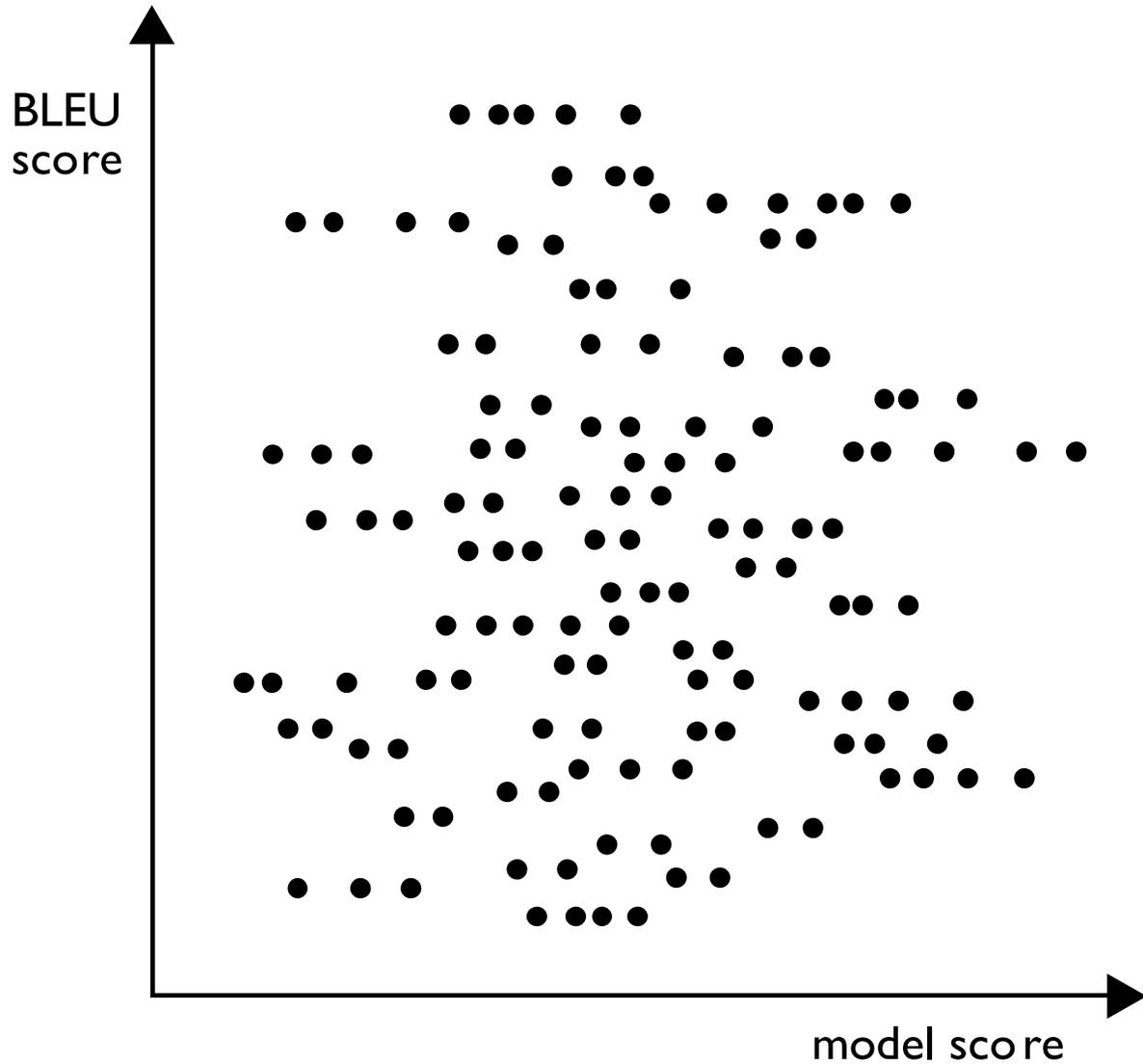
Pairwise Ranking Optimization (Hopkins & May, 2011)



How about standard machine learning algorithms
for structured prediction?

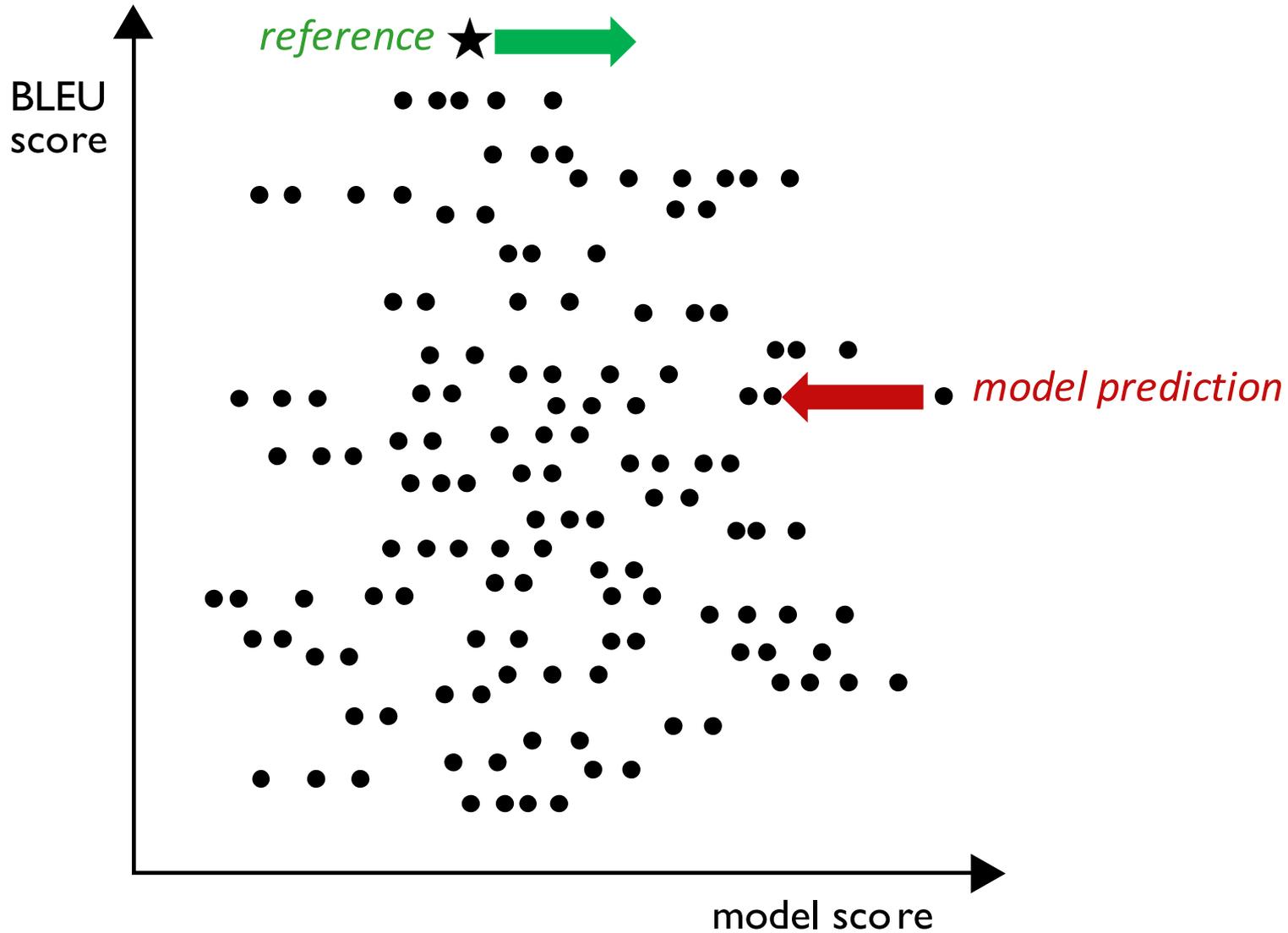
Structured Perceptron?

(Collins, 2002)

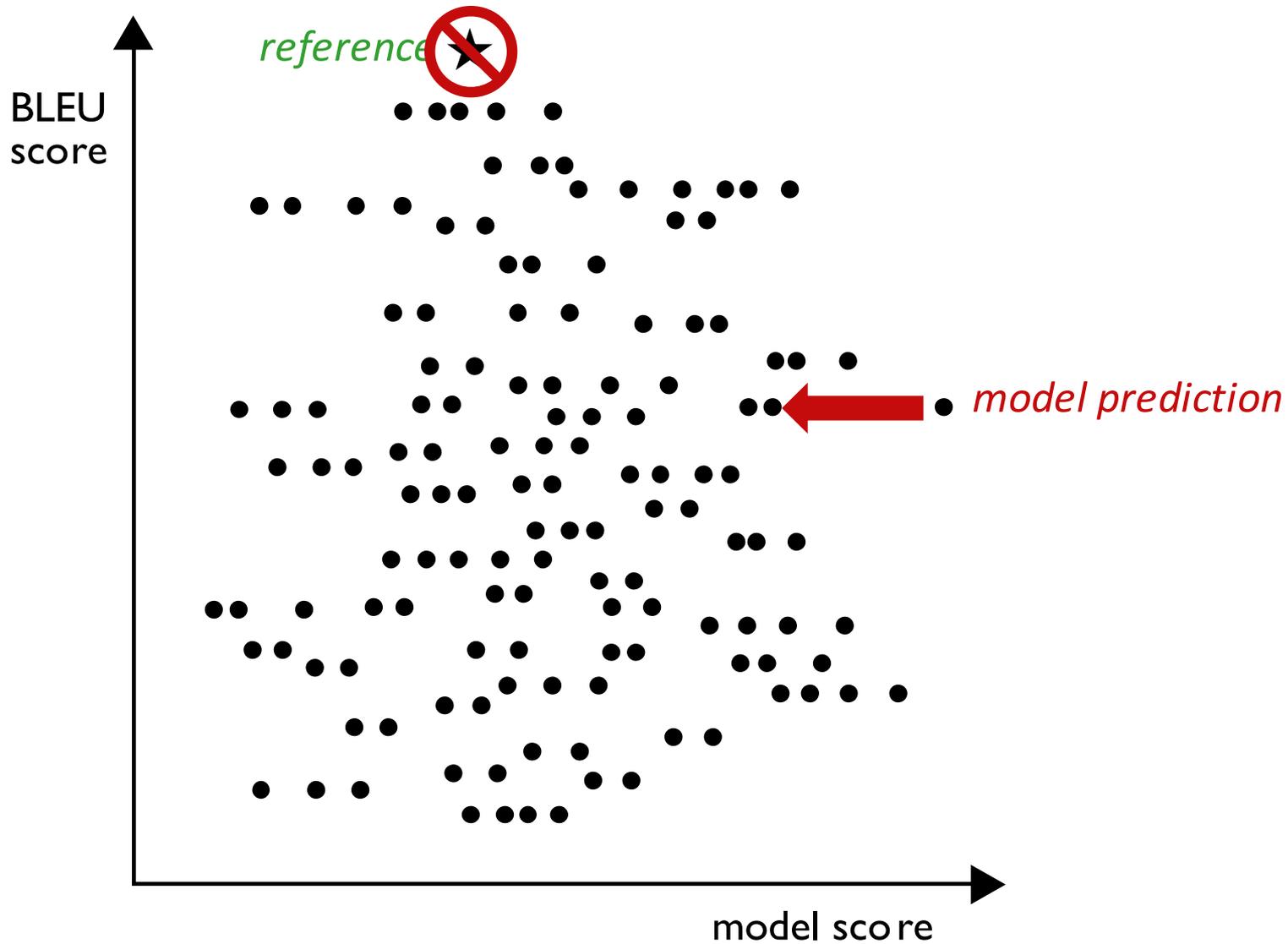


Structured Perceptron

(Collins, 2002)

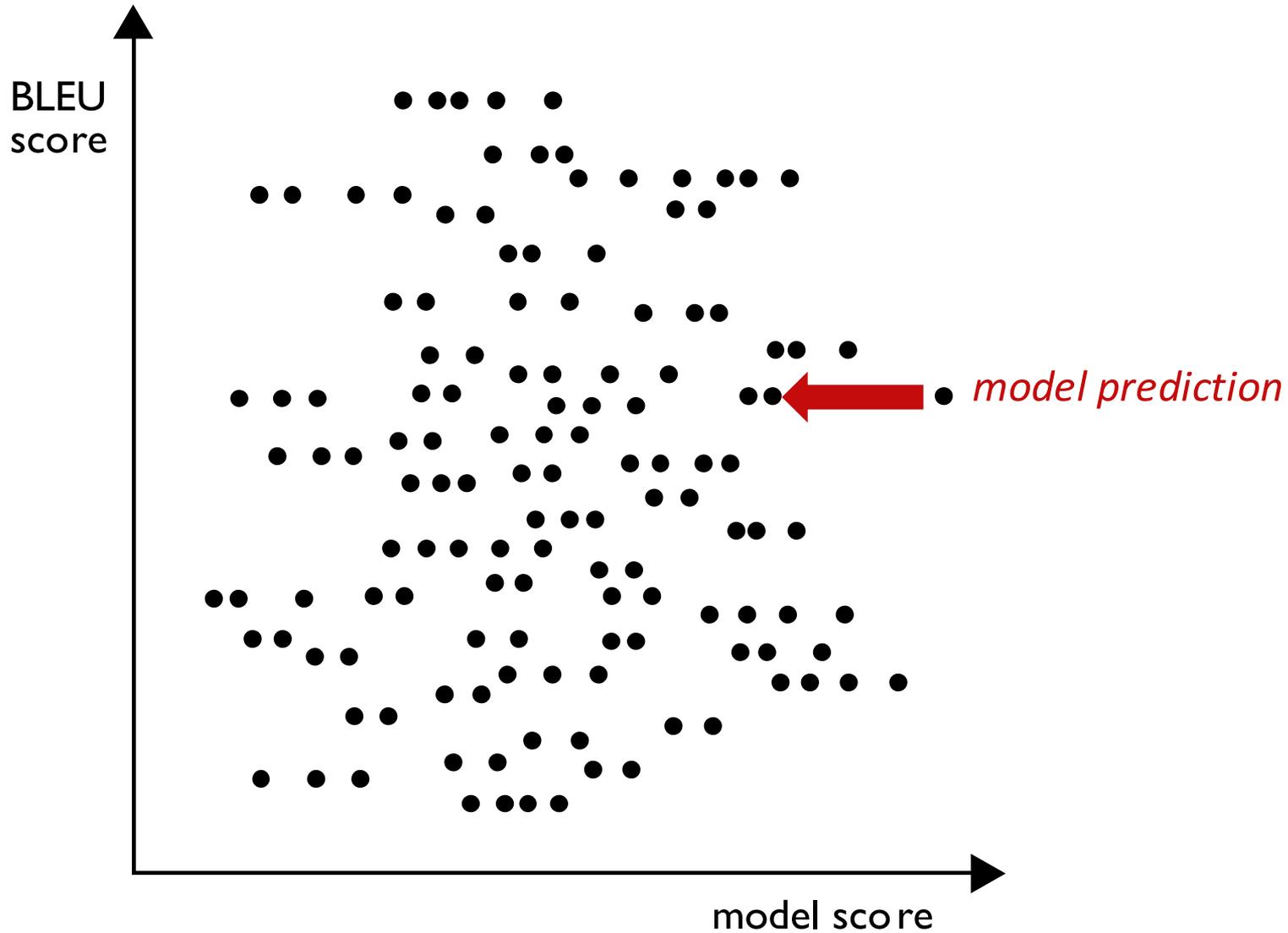


Structured Perceptron for MT?



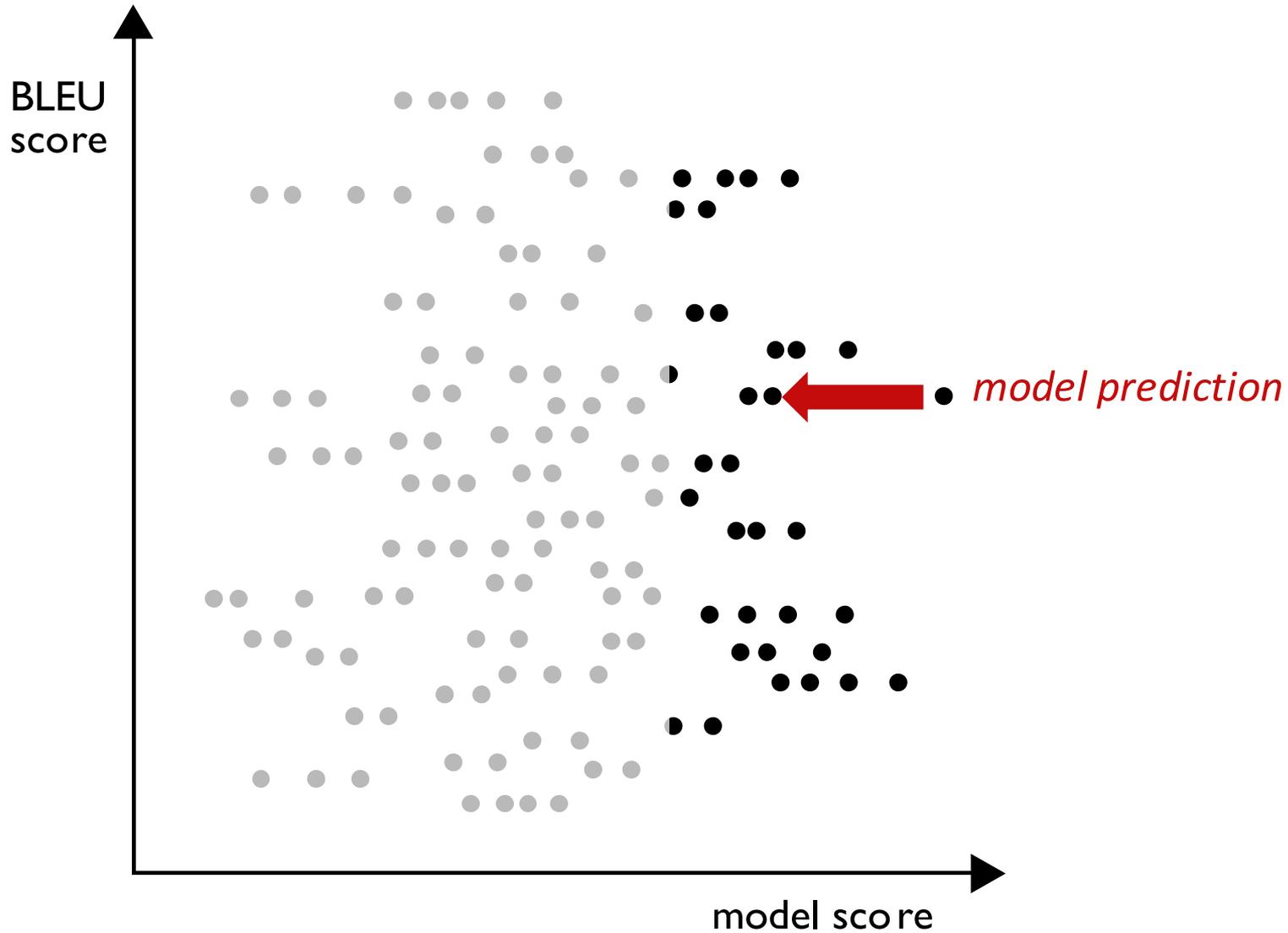
k-Best Perceptron for MT

(Liang et al., 2006)



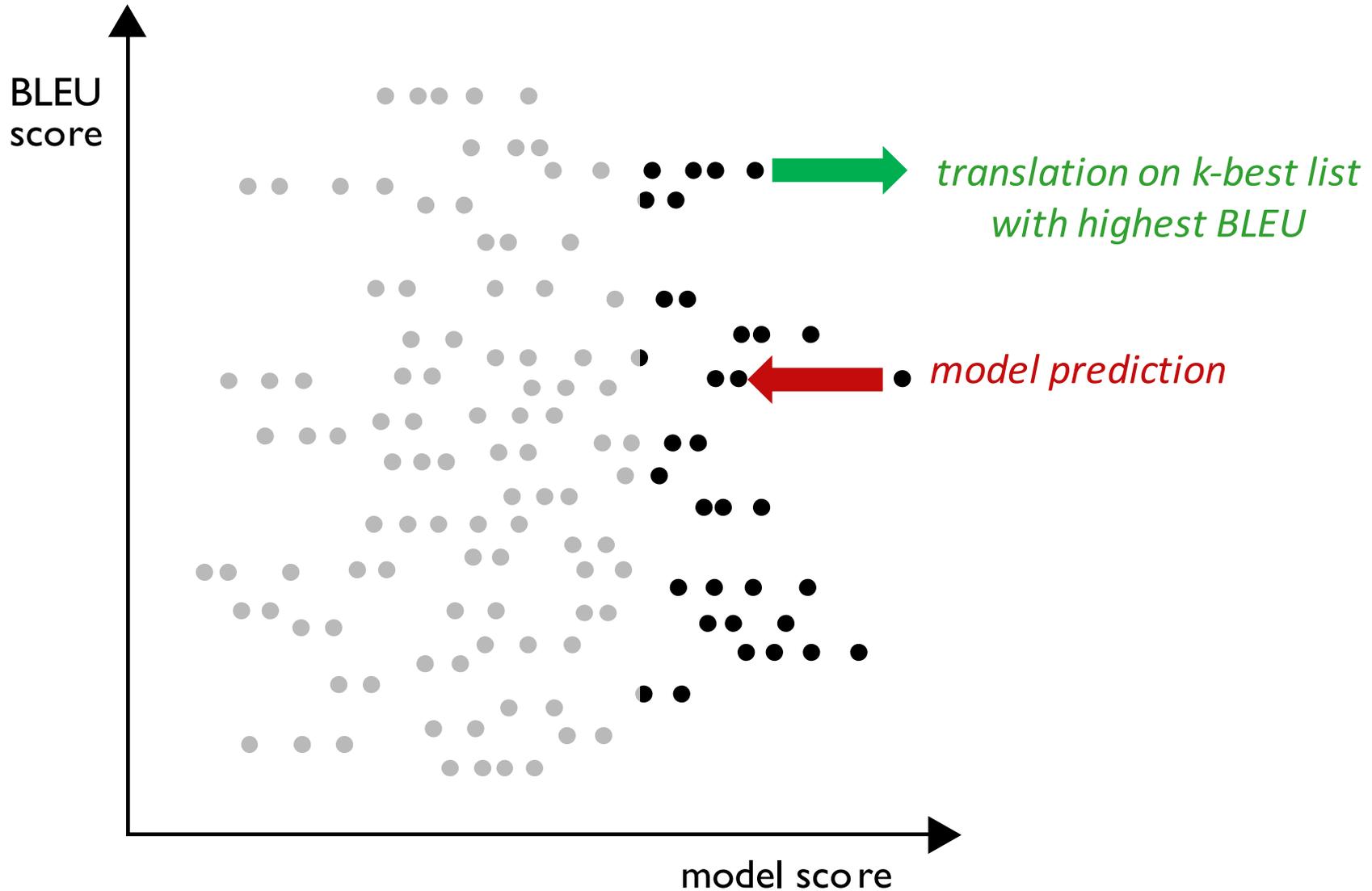
k-Best Perceptron for MT

(Liang et al., 2006)



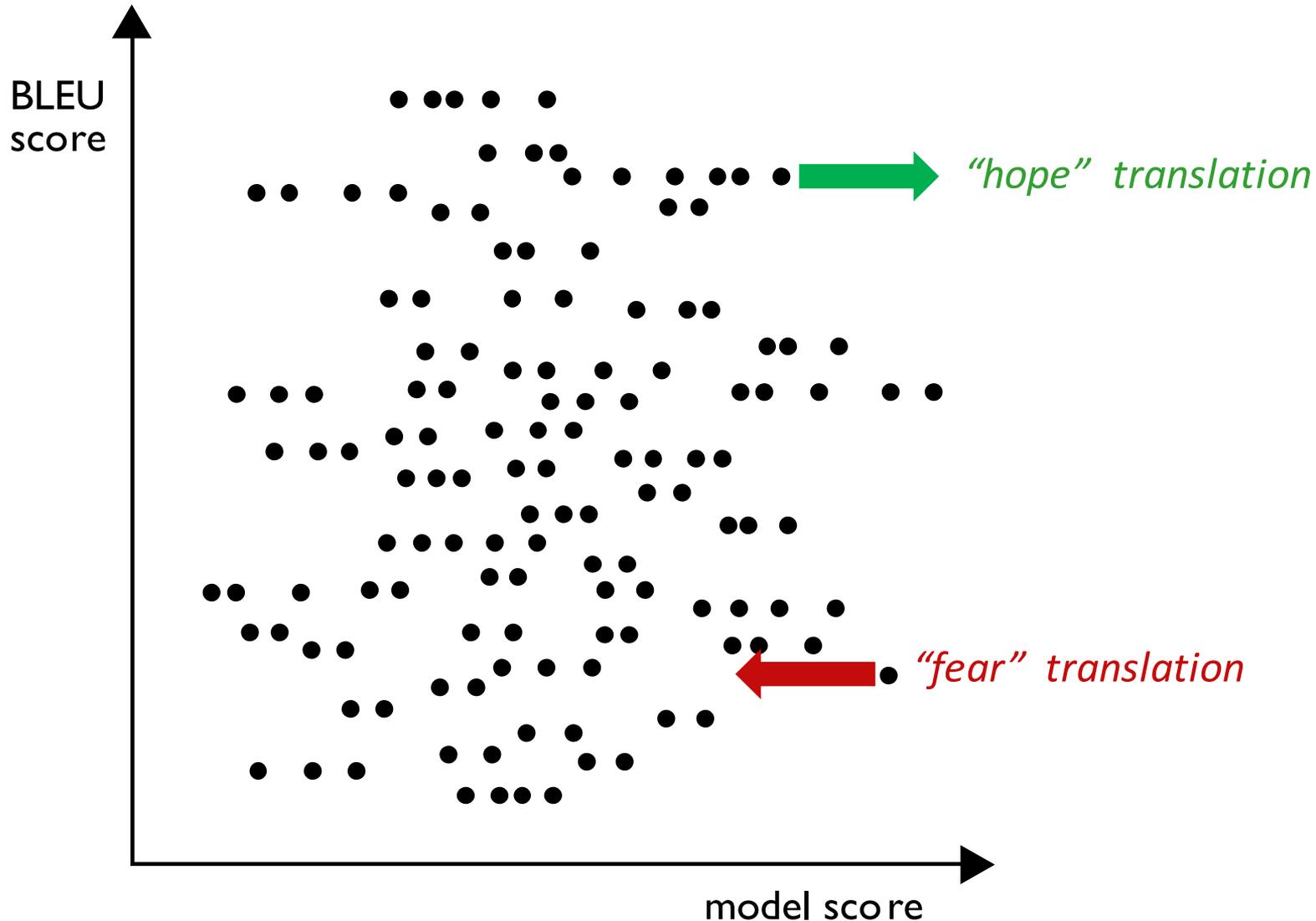
k-Best Perceptron for MT

(Liang et al., 2006)



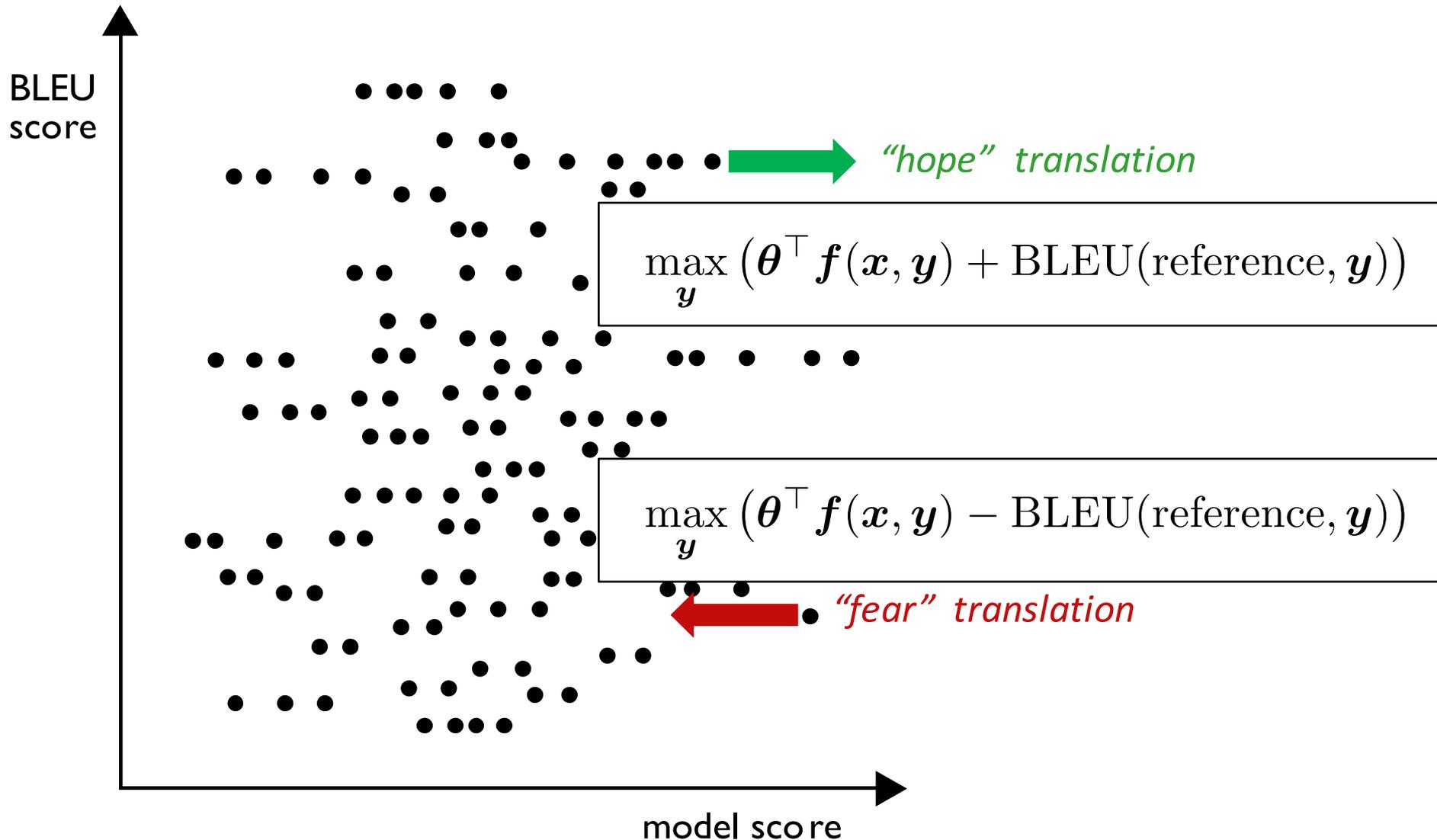
“Hope-Fear” MIRA

(Chiang et al., 2008; 2009; Cherry & Foster, 2012; Chiang, 2012)



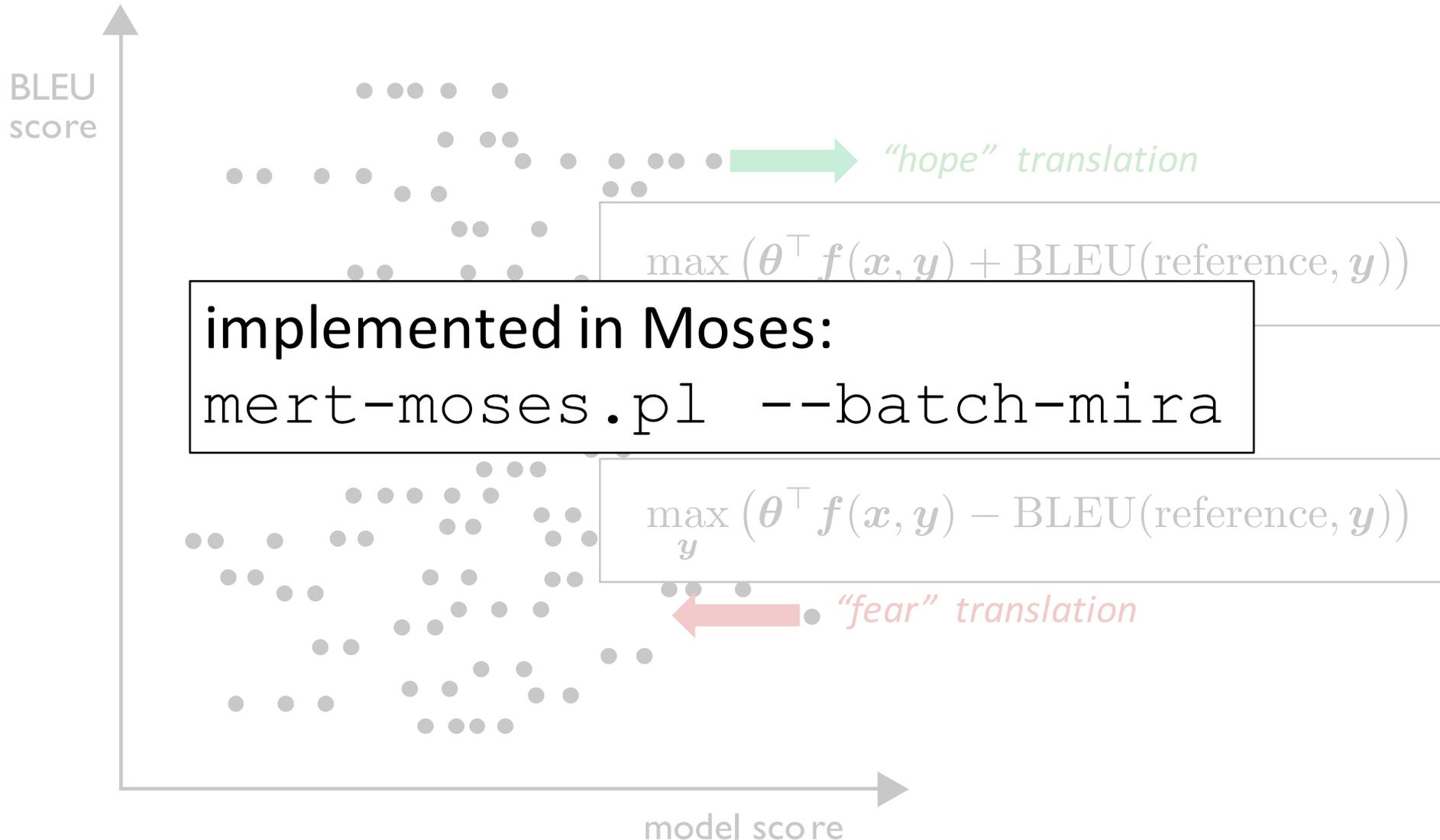
“Hope-Fear” MIRA

(Chiang et al., 2008; 2009; Cherry & Foster, 2012; Chiang, 2012)



“Hope-Fear” MIRA

(Chiang et al., 2008; 2009; Cherry & Foster, 2012; Chiang, 2012)



MT Experiments

(Gimpel, 2012)

averages over 8 test sets across 3 language pairs

	Moses %BLEU	Hiero %BLEU
MERT	35.9	37.0
PRO	35.9	36.9
Bayes Risk	35.6	36.4
Fear MIRA	34.9	34.2
Hope MIRA	35.2	36.0
Hope-Fear MIRA	35.7	37.0

Questions?

2002: conditional log-likelihood

Proceedings of the 40th Annual Meeting of the Association for
Computational Linguistics (ACL), Philadelphia, July 2002, pp. 295-302.

Discriminative Training and Maximum Entropy Models for Statistical Machine Translation

Franz Josef Och and Hermann Ney

Lehrstuhl für Informatik VI, Computer Science Department

RWTH Aachen - University of Technology

D-52056 Aachen, Germany

{och,ney}@informatik.rwth-aachen.de

2002: conditional log-likelihood

Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL), Philadelphia, July 2002, pp. 295-302.

Discriminative Training and Maximum Entropy Models for Statistical Machine Translation

Table 2: Effect of maximum entropy training for alignment template approach (WP: word penalty feature, CLM: class-based language model (five-gram), MX: conventional dictionary).

	objective criteria [%]					subjective criteria [%]	
	SER	WER	PER	mWER	BLEU	SSER	IER
Baseline($\lambda_m = 1$)	86.9	42.8	33.0	37.7	43.9	35.9	39.0
ME	81.7	40.2	28.7	34.6	49.7	32.5	34.8
ME+WP	80.5	38.6	26.9	32.4	54.1	29.9	32.2
ME+WP+CLM	78.1	38.3	26.9	32.1	55.0	29.1	30.9
ME+WP+CLM+MX	77.8	38.4	26.8	31.9	55.2	28.8	30.9



Discriminative

Hinge Loss

