# TTIC 31190:
# Natural Language Processing

Kevin Gimpel
Winter 2016

## Lecture 15:
## Introduction
## to Machine Translation

# Announcements

- Assignment 3 due Monday
- email me to sign up for your (10-minute) class presentation on 3/3 or 3/8

# Roadmap

- classification
- words
- lexical semantics
- language modeling
- sequence labeling
- neural network methods in NLP
- syntax and syntactic parsing
- computational semantics
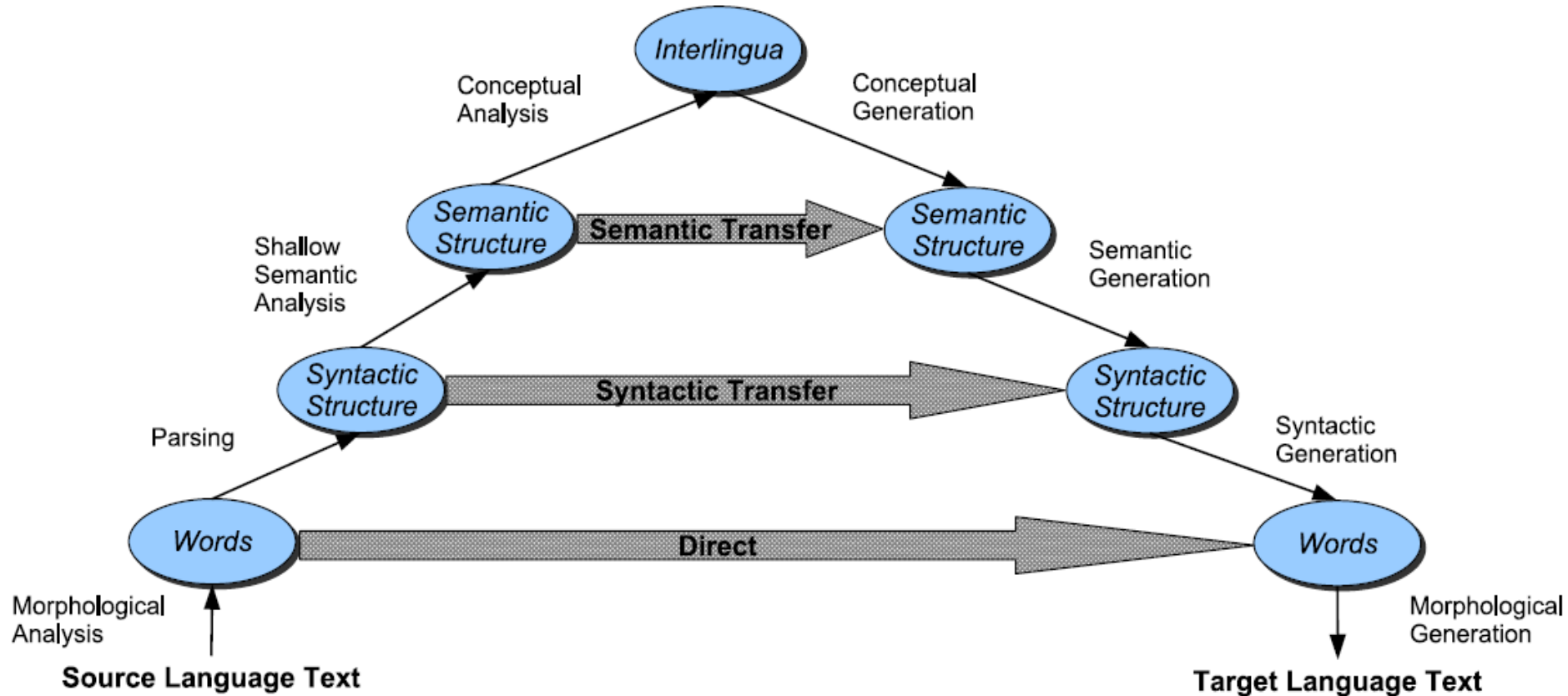- machine translation
- other NLP applications

# People rely on machine translation!

# People rely on machine translation!

# Approaches to Machine Translation: The Vauquois Triangle

# Interlingua Example

$$
\begin{bmatrix}
\text{EVENT} & \text{SLAPPING} \\
\text{AGENT} & \textsc{Mary} \\
\text{TENSE} & \text{PAST} \\
\text{POLARITY} & \text{NEGATIVE} \\
\text{THEME} & \begin{bmatrix}
\text{WITCH} \\
\text{DEFINITENESS} & \text{DEF} \\
\text{ATTRIBUTES} & \begin{bmatrix} \text{HAS-COLOR} & \text{GREEN} \end{bmatrix}
\end{bmatrix}
\end{bmatrix}
$$

Interlingual representation of *Mary did not slap the green witch*.

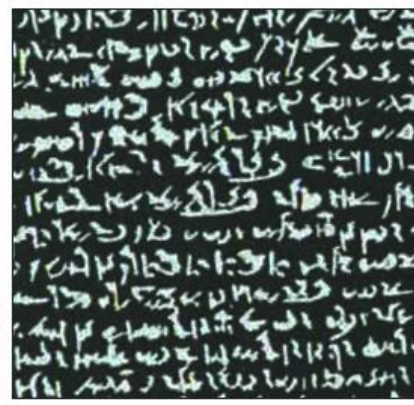# Classification Framework for Machine Translation

**inference**: solve $\operatorname{argmax}$
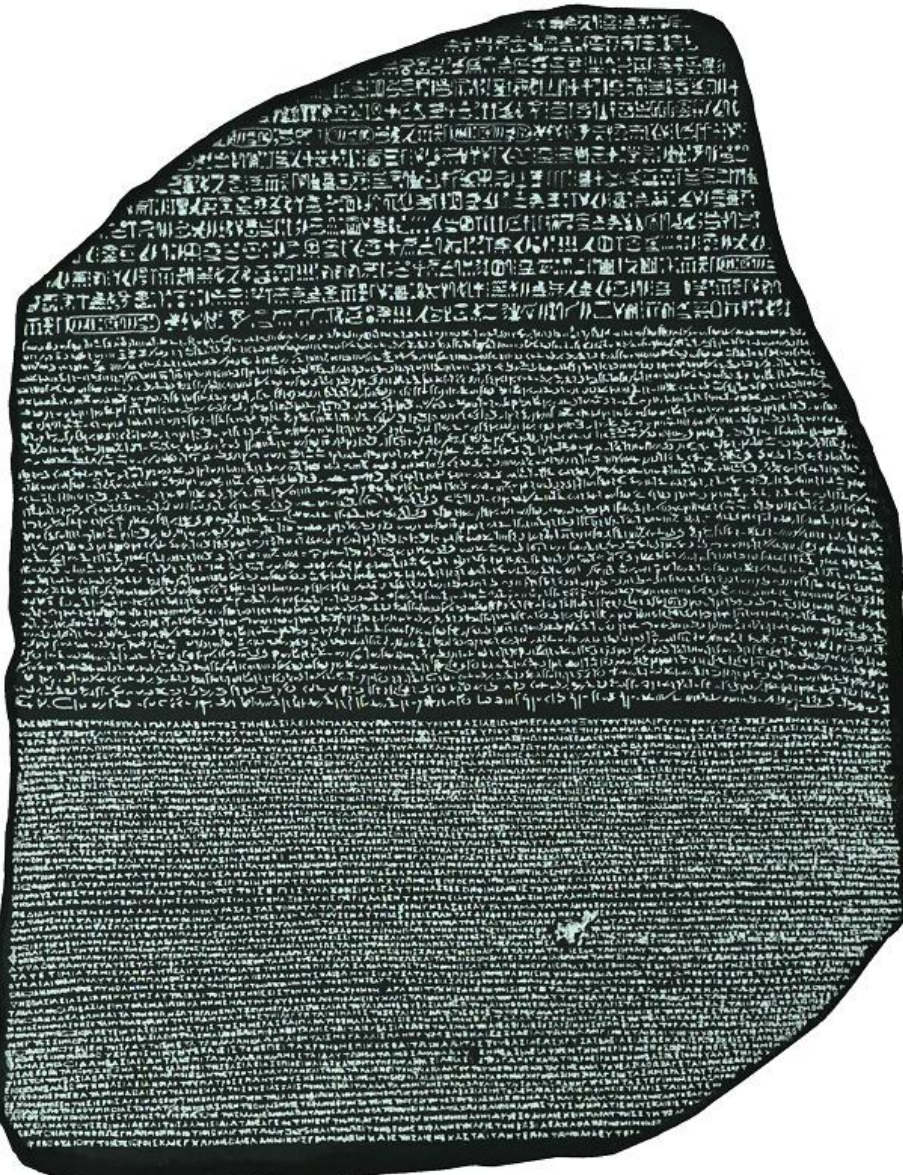
**modeling**: define $\operatorname{score}$ function

$$\operatorname{classify}(\boldsymbol{x}, \boldsymbol{\theta}) = \underset{\boldsymbol{y}}{\operatorname{argmax}} \ \operatorname{score}(\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{\theta})$$

**learning**: choose $\boldsymbol{\theta}$

- modern systems are **data-driven**

- first we need data!

# Data?

# Data?

## 碟頭飯 RICE PLATE

| 揚州炒飯 | Yang Chow Fried Rice | 7.95 |
| 咸魚雞粒炒飯 | Salted Fish w/ Chicken Fried Rice | 8.95 |
| 油雞飯 | Soy Chicken Rice | 5.95 |
| 滑雞菜遠飯 | Chicken with Vegetable on Rice | 5.95 |
| 粟米雞扒飯 | Chicken with Cream Corn on Rice | 5.95 |
| 豉椒雞球飯 | Chicken W/ Black Bean Sauce | 5.95 |
| 涼瓜牛肉飯 | Beef with Bitter Melon on Rice | 5.95 |
| 菜遠牛肉飯 | Beef with Vegetable on Rice | 5.95 |
| 牛腩飯 | Beef Stew on Rice | 6.95 |
| 滑蛋牛肉飯 | Beef with Egg on Rice | 5.95 |
| 滑蛋蝦仁飯 | Shrimp with Egg on Rice | 6.95 |
| 鮮蝦菜遠飯 | Shrimp with Vegetable on Rice | 6.95 |
| 魚片菜遠飯 | Fish with Vegetable on Rice | 6.95 |
| 咖哩尤魚飯 | Curry Squid on Rice | 6.95 |
| 滑蛋叉燒飯 | BBQ Pork with Egg on Rice | 5.95 |
| 肉片豆腐飯 | Pork with Tofu on Rice | 5.95 |

## 粥品 CONGEE

| 白粥 | Plain Congee | 2.50 |
| 皮蛋肉片粥 | Preserve Egg w/ Pork Congee | 5.50 |
| 生滾牛肉粥 | Beef Congee | 5.50 |
| 魚片粥 | Fish Congee | 5.95 |
| 滑雞粥 | Chicken Congee | 5.50 |

## Chinese Menu

## Kings Garden

### 球記-皇家園

Authentic Chinese Food

TEL: (614) 793-2234

7726 Sawmill Rd.

Dublin, Ohio 43017

(Old Sawmill Sq. Shopping Center)

OPEN HOUR

| Mon | Close |
| Tues – Sat | 11:00 am to 10:00 pm |
| Sun | 11:00 am to 9:00 pm |

Catering available.

# Data?

302 云南芫爆松茸
Sauteed trichdoma matsutake with coriander ar
蘑菇之王，素有 "海有鲟鱼子，陆地上的松茸"，含人
细嫩，香味浓溢

303 白油爆鸡枞
Stir-fried wikipedia
肉质细嫩，洁白如玉，或炒或蒸、串汤作菜，清香四
云南皱椒鸡枞
Stir-fried wikipedia with pimientos

304 香油鸡枞蒸水蛋
Steam eggs with wikipedia

|  | | Savory potato wedges | ¥ | 15 / 例 |
| 濃湯 | | Gream of pumpkin soup | ¥ | 15 / 例 |
| 薩角 | | India samosa | ¥ | 25 / 例 |
| 6. 意式火腿面包棒 | | Italian ham bread | ¥ | 15 / 例 |
| 7. 田園沙拉 | | Garden salad | ¥ | 15 / 例 |
| 8. 三明治 | | Sand wiches | ¥ | 15 / 和 |
| （培根/薩拉米/吞拿魚/火腿） | | (Bacon/Salami/Tuna/Ham) | | |
| 9. 香烤魷魚圈 | | BBQ wikipedia | ¥ | 20 / 例 |
| 10.串烤牛小排 | | BBQ beef and vegetables | ¥ | 20 / 例 |
| 11.水牛城香辣鷄翅 | | Kookaburra wings | ¥ | 25 / 6 |
| 12.德國烤腸 | | German BBQ Sausage | ¥ | 30 / 例 |
| 13.香蒜面包 | | Garlic butter bread | ¥ | 10 / 3 |

# Data?



Also:
- news articles
- company websites
- laws & patents
- subtitles

# Parallel Data

- **parallel data**: bilingual data that is naturally aligned at some level

- usually aligned at the document level

- sentence-level alignments are generated automatically

  - how might you design an algorithm for this?

  - it can be done well without dictionaries!

  - can throw out sentences that don't align with anything

# Learning from Parallel Sentences

| Chickasaw | English |
|---|---|
| 1. Ofi 'at kowi 'ã lhiyohli | 1. The dog chases the cat |
| 2. Kowi 'at ofi 'ã lhiyohli | 2. The cat chases the dog |
| 3. Ofi 'at shoha | 3. The dog stinks |

# Learning from Parallel Sentences

## Chickasaw

1. Ofi 'at kowi 'ã lhiyohli
2. Kowi 'at ofi 'ã lhiyohli
3. Ofi 'at shoha

## English

1. The dog chases the cat
2. The cat chases the dog
3. The dog stinks

# Machine Translation Evaluation

- human judgments are ideal, but expensive
  - what other problems are there with human judgments?

- we need automatic evaluation metrics
  - BLEU (BiLingual Evaluation Understudy), Papineni et al. (2002)
    - compare *n*-gram overlap between system output and human-produced translation
    - correlates with human judgments surprisingly well, but only at the document level (not sentence level!)
  - other metrics do soft matching based on stemming and synonyms from WordNet
  - this is not a solved problem!

# **Statistical** Machine Translation

*One naturally wonders if the problem of translation could conceivably be treated as a problem in **cryptography**.*

*When I look at an article in Arabic, I say: "This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode."*
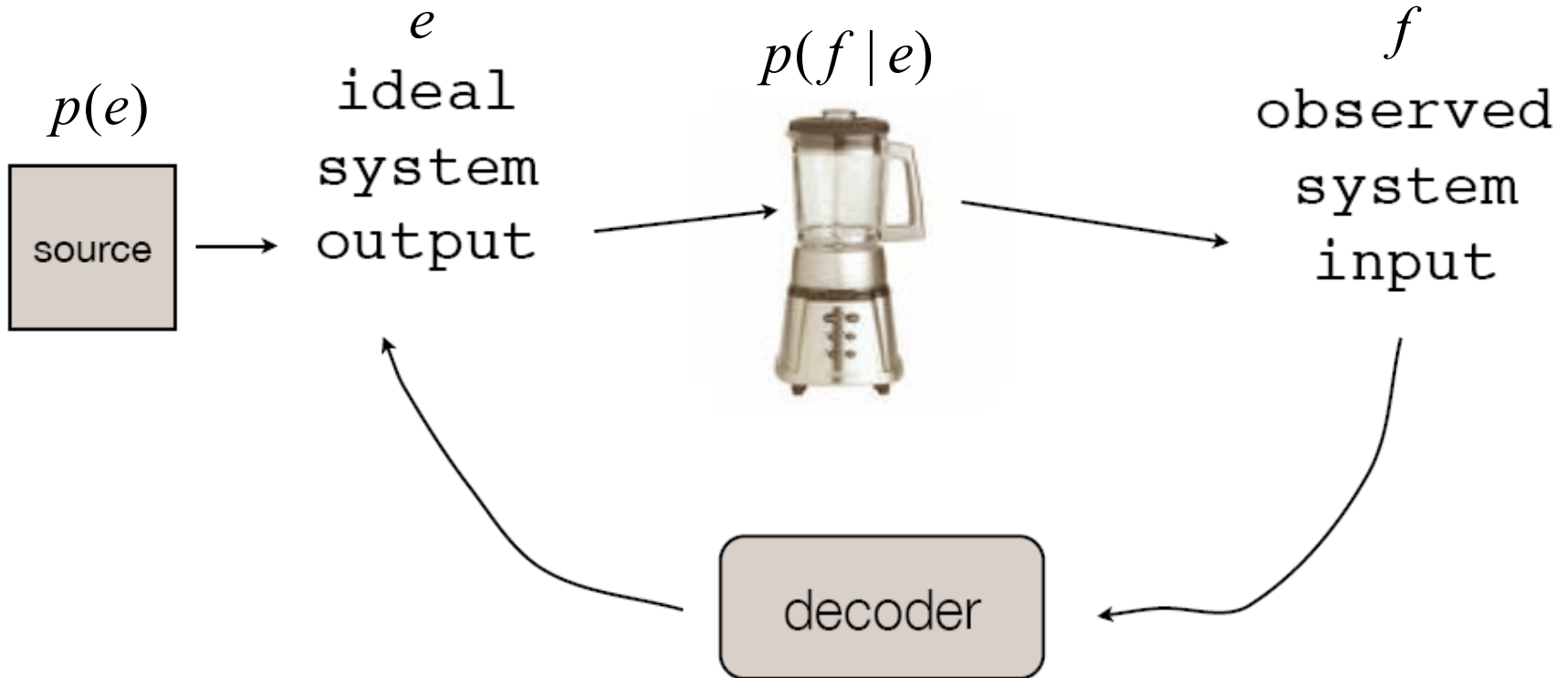
Warren Weaver, 1947

# Noisy Channel Model

# Noisy Channel Model for Translating French ( $f$ ) to English ( $e$ )

$p(e)$

$e$
ideal
system
output

$p(f \mid e)$

$f$
observed
system
input

source

decoder

$$\hat{e} = \arg\max_{e} p(e \mid f)$$

$$= \arg\max_{e} \frac{p(f \mid e) p(e)}{p(f)}$$

$$= \arg\max_{e} p(f \mid e) p(e)$$

# Modeling for the Noisy Channel

- We need to model two probability distributions: $P(e)$ and $P(f \mid e)$
  - $P(e)$ should favor fluent translations
  - $P(f \mid e)$ should favor accurate/faithful translations

# Modeling for the Noisy Channel

- We need to model two probability distributions: $P(e)$ and $P(f \mid e)$
  - $P(e)$ should favor fluent translations
  - $P(f \mid e)$ should favor accurate/faithful translations

- Let's start with $P(e)$
  - How do we compute the probability of an English sentence?
  - This is an important part of MT (e.g., Google)

# Word Alignments

And$_1$ the$_2$ program$_3$ has$_4$ been$_5$ implemented$_6$

Le$_1$ programme$_2$ a$_3$ été$_4$ mis$_5$ en$_6$ application$_7$

# Word Alignments

$$a = a_1 ... a_{|f|}$$



And$_1$     the$_2$     program$_3$     has$_4$     been$_5$     implemented$_6$

$a_1 = 2$    $a_2 = 3$    $a_3 = 4$    $a_4 = 5$    $a_5 = 6$   $a_6 = 6$   $a_7 = 6$

Le$_1$    programme$_2$    a$_3$    été$_4$    mis$_5$    en$_6$    application$_7$

- $a$ is a "hidden" variable (not part of training data)
- for each French word, it holds the index of the aligned English word (or NULL)

- remember: our goal was to model $P(f \mid e)$
- why would we introduce a hidden variable?
  - to make it "easier" to define the model
  - we often want to share certain types of information across multiple instances in our data
    - latent variables are a natural way to capture this
    - think of clustering (some of the points come from the same cluster)

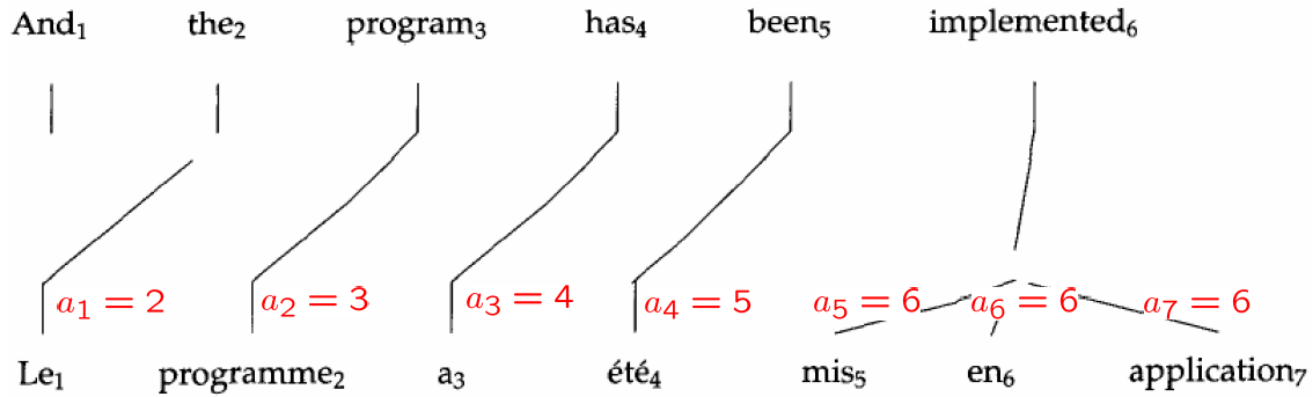# Alignments as Hidden Variables

- for simplicity, assume that each French word aligns to 1 English word (or to NULL)

- analogy to clustering:
  - each data point has 1 vote which it can distribute among all the clusters
  - here, each French word has 1 vote which it can distribute among all the English words or NULL

# Modeling Alignments: IBM Model 1

And₁     the₂     program₃     has₄     been₅     implemented₆

$a_1 = 2$    $a_2 = 3$    $a_3 = 4$    $a_4 = 5$    $a_5 = 6$   $a_6 = 6$   $a_7 = 6$

Le₁     programme₂     a₃     été₄     mis₅     en₆     application₇

$$P(f, a \mid e) = \prod_{j=1}^{|f|} P(a_j) P(f_j | e_{a_j})$$

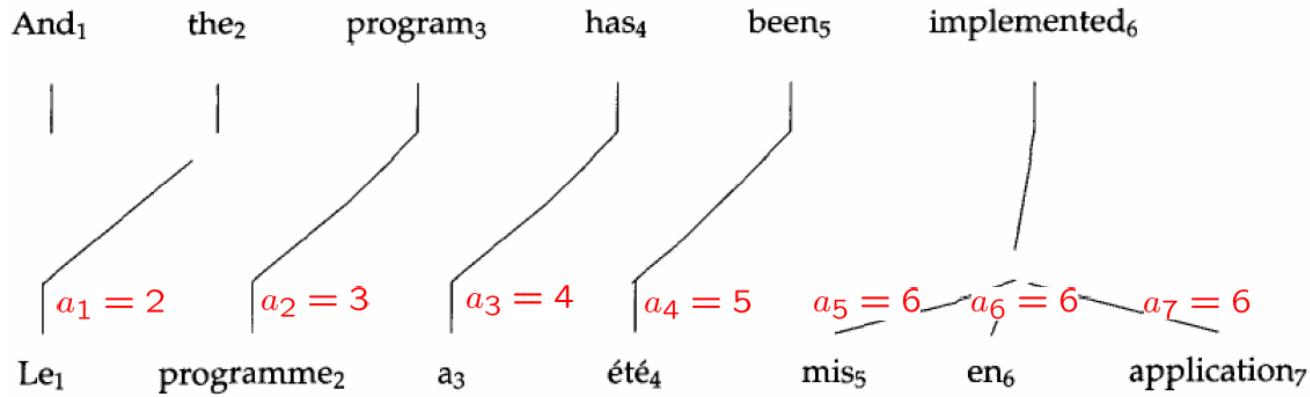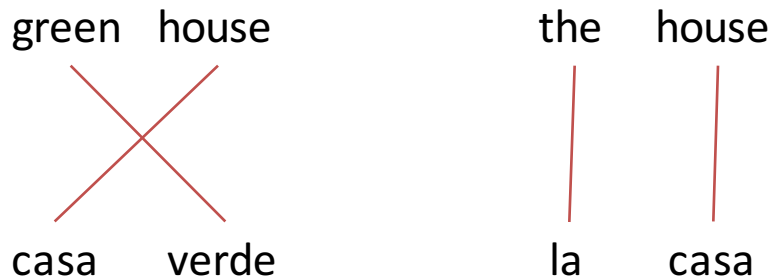$$= \prod_{j=1}^{|f|} \frac{1}{|e| + 1} P(f_j | e_{a_j})$$

# Modeling Alignments: IBM Model 1

And$_1$     the$_2$     program$_3$     has$_4$     been$_5$     implemented$_6$

$a_1 = 2$    $a_2 = 3$    $a_3 = 4$    $a_4 = 5$    $a_5 = 6$   $a_6 = 6$   $a_7 = 6$

Le$_1$     programme$_2$     a$_3$     été$_4$     mis$_5$     en$_6$     application$_7$

$$P(f, a \mid e) = \prod_{j=1}^{|f|} P(a_j) P(f_j | e_{a_j})$$

$$= \prod_{j=1}^{|f|} \frac{1}{|e| + 1} P(f_j | e_{a_j})$$

• How do we obtain $P(f \mid e)$?

# Modeling Alignments: IBM Model 1

And$_1$  the$_2$  program$_3$  has$_4$  been$_5$  implemented$_6$

$a_1 = 2$  $a_2 = 3$  $a_3 = 4$  $a_4 = 5$  $a_5 = 6$  $a_6 = 6$  $a_7 = 6$

Le$_1$  programme$_2$  a$_3$  été$_4$  mis$_5$  en$_6$  application$_7$

$$P(f, a \mid e) = \prod_{j=1}^{|f|} P(a_j) P(f_j | e_{a_j})$$

$$= \prod_{j=1}^{|f|} \frac{1}{|e| + 1} P(f_j | e_{a_j})$$

- How do we obtain $P(f \mid e)$?
- Sum over all alignments: $P(f \mid e) = \sum_a P(f, a \mid e)$

# Modeling Alignments: IBM Model 1



$$P(f, a \mid e) = \prod_{j=1}^{|f|} P(a_j) P(f_j | e_{a_j})$$

$$= \prod_{j=1}^{|f|} \frac{1}{|e| + 1} P(f_j | e_{a_j})$$

Parameters in the model,
learned using expectation maximization

# Aside: are alignments always hidden?

- certain small parallel corpora have been hand-aligned
- issues with this?
  - annotators don't agree
  - we have lots of parallel text, very little is hand-aligned
  - for some language pairs, we will never have manual alignments

- word alignment has become a fundamental part of MT, and we need unsupervised learning to solve it!

# IBM Model 1 Example

- Consider a training set of two sentence pairs:

green    house                the    house

casa    verde                la    casa

Initial Parameter Estimates:

| t(casa\|green) | = $\frac{1}{3}$ | t(verde\|green) | = $\frac{1}{3}$ | t(la\|green) | = $\frac{1}{3}$ |
|---|---|---|---|---|---|
| t(casa\|house) | = $\frac{1}{3}$ | t(verde\|house) | = $\frac{1}{3}$ | t(la\|house) | = $\frac{1}{3}$ |
| t(casa\|the) | = $\frac{1}{3}$ | t(verde\|the) | = $\frac{1}{3}$ | t(la\|the) | = $\frac{1}{3}$ |

$t(f \mid e)$

= probability of translating $e$ into $f$

After 1 iteration of EM:

| t(casa\|green) | = $\frac{1/2}{1} = \frac{1}{2}$ | t(verde\|green) | = $\frac{1/2}{1} = \frac{1}{2}$ | t(la\|green) | = $\frac{0}{1} = 0$ |
|---|---|---|---|---|---|
| t(casa\|house) | = $\frac{1}{2} = \frac{1}{2}$ | t(verde\|house) | = $\frac{1/2}{2} = \frac{1}{4}$ | t(la\|house) | = $\frac{1/2}{2} = \frac{1}{4}$ |
| t(casa\|the) | = $\frac{1/2}{1} = \frac{1}{2}$ | t(verde\|the) | = $\frac{0}{1} = 0$ | t(la\|the) | = $\frac{1/2}{1} = \frac{1}{2}$ |

# IBM Model 1

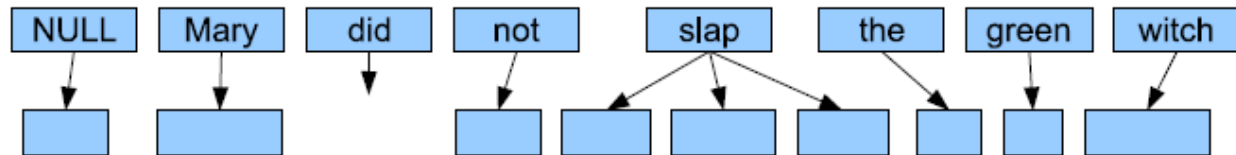$$P(f, a \mid e) = \prod_{j=1}^{|f|} \frac{1}{|e| + 1} P(f_j | e_{a_j})$$

# IBM Model 2

$$P(f, a \mid e) = \prod_{j=1}^{|f|} P(a_j \mid j, |f|, |e|) P(f_j | e_{a_j})$$
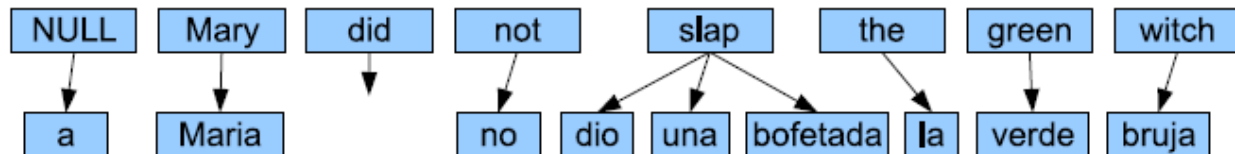
# IBM Model 3

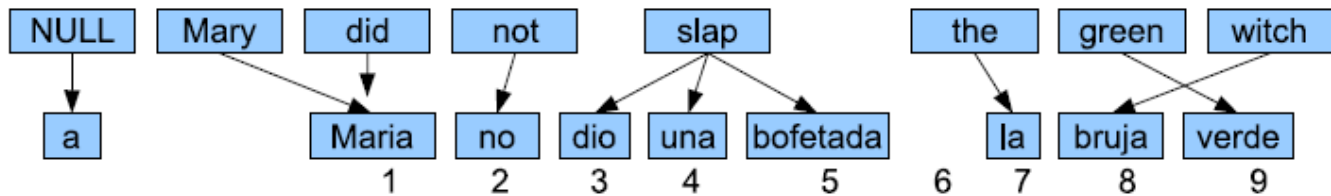Step 1: Choose fertility for each English word

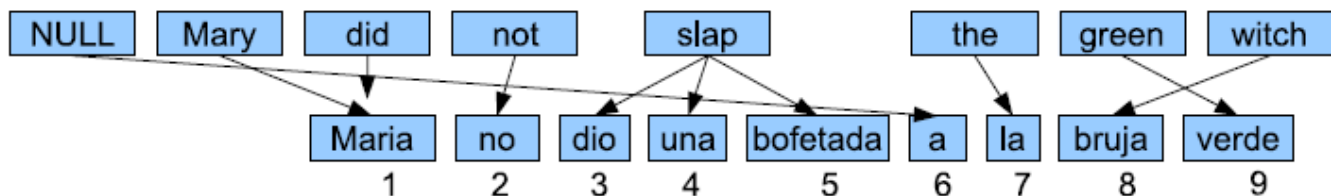| NULL | Mary | did | not | slap | the | green | witch |

Step 2: Choose fertility for NULL

| NULL | Mary | did | not | slap | the | green | witch |

Step 3: Create Spanish words by translating aligned English word

| NULL | Mary | did | not | slap | the | green | witch |

a  Maria  no  dio  una  bofetada  la  verde  bruja

Step 4: Move the Spanish words into final slots

| NULL | Mary | did | not | slap | the | green | witch |

a  Maria  no  dio  una  bofetada  la  bruja  verde
1    2    3    4    5            6  7    8     9

Step 4: Move spurious Spanish words into unclaimed slots

| NULL | Mary | did | not | slap | the | green | witch |

Maria  no  dio  una  bofetada  a  la  bruja  verde
1    2    3    4    5          6  7   8     9

# Moving to Phrases



NULL Auf diese Frage habe ich leider keine Antwort bekommen

I did not unfortunately receive an answer to this question

# Moving to Phrases

Not necessarily syntactic phrases

| Auf diese Frage | habe ich | leider | keine | Antwort bekommen |
|---|---|---|---|---|

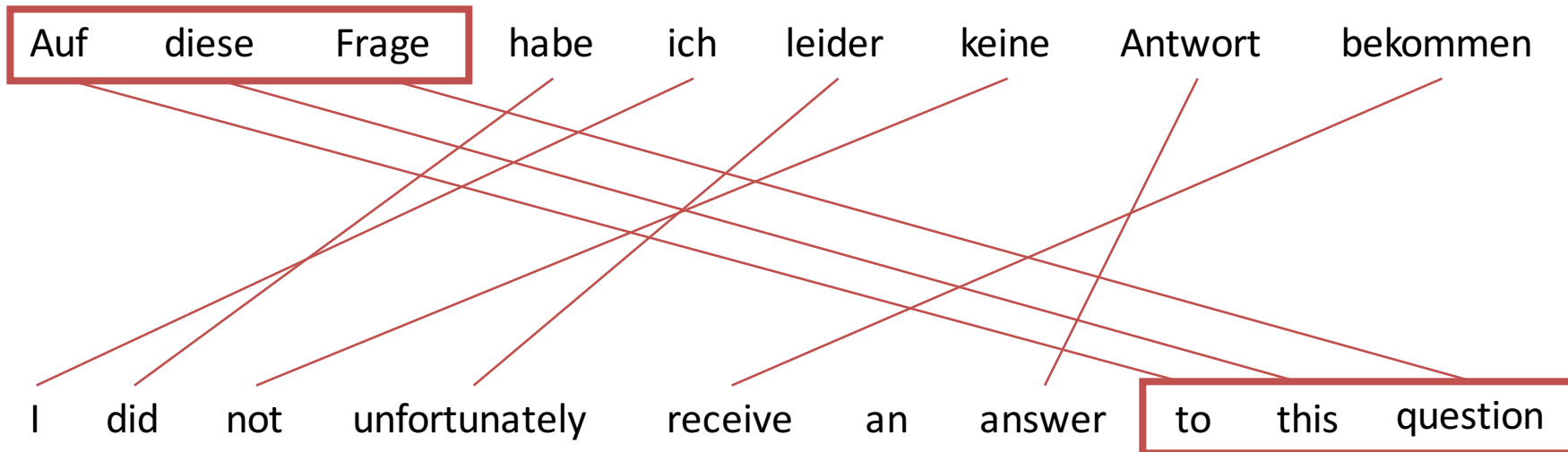| I did | not | unfortunately | receive an answer | to this question |
|---|---|---|---|---|

# "Phrase-Based" Translation

- Relies on a **phrase table**

  – massive bilingual phrase dictionary, with probabilities

- To build:

  – Find the best word alignment for each sentence pair

  – Extract all phrase pairs **consistent** with the word alignment

  – Compute probabilities using relative frequency estimation
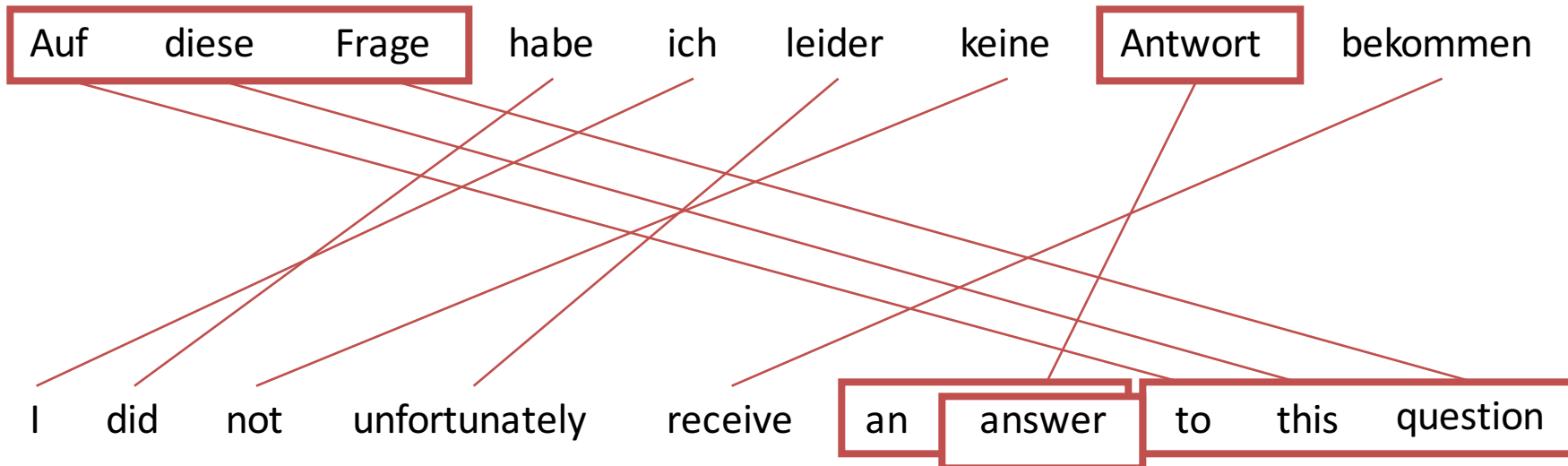
# Phrase-Based Translation

- ## Relies on a **phrase table**

  - massive bilingual phrase dictionary, with probabilities

- ## To build:

  - Find the best word alignment for each sentence pair
  - Extract all phrase pairs **consistent** with the word alignment
  - Compute probabilities using relative frequency estimation

Auf    diese    Frage    habe    ich    leider    keine    Antwort    bekommen

I    did    not    unfortunately    receive    an    answer    to    this    question
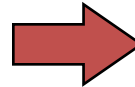
# Phrase-Based Translation

- Relies on a **phrase table**
  - massive bilingual phrase dictionary, with probabilities

- To build:
  - Find the best word alignment for each sentence pair
  - Extract all phrase pairs **consistent** with the word alignment
  - Compute probabilities

| Auf diese Frage | to this question | 1.0 |
| --- | --- | --- |

# Phrase-Based Translation

- ## Relies on a **phrase table**
  - massive bilingual phrase dictionary, with probabilities

- ## To build:
  - Find the best word alignment for each sentence pair
  - Extract all phrase pairs **consistent** with the word alignment
  - Compute probabilities

| Auf diese Frage | to this question | 1.0 |
| Antwort | an answer | 1.0 |
| Antwort | answer | 1.0 |
| | ... | |

Auf    diese    Frage    habe    ich    leider    keine    Antwort    bekommen

I    did    not    unfortunately    receive    an    answer    to    this    question

# Phrase-Based Translation

- ## Relies on a **phrase table**
  - massive bilingual phrase dictionary, with probabilities

- ## To build:
  - Find the best word alignment for each sentence pair
  - Extract all phrase pairs **consistent** with the word alignment
  - Compute probabilities using relative frequency estimation:

$$p(e \mid f) = \frac{count(e, f)}{\sum_{e'} count(e', f)}$$

| German | English | Count |
|---|---|---|
| Auf diese Frage | to this question | 1.0 |
| Antwort | an answer | 1.0 |
| Antwort | answer | 1.0 |
| … | | |

| German | English | P( e \| f ) |
|---|---|---|
| Auf diese Frage | to this question | 1.0 |
| Antwort | an answer | 0.5 |
| Antwort | answer | 0.5 |
| … | | |

# Adding Syntax: Synchronous Context-Free Grammars

**CFG**

$$NP \longrightarrow DT\ NPB$$
$$NPB \longrightarrow JJ\ NPB$$
$$NPB \longrightarrow \mathbf{NN}$$
$$DT \longrightarrow the$$
$$JJ \longrightarrow strong$$
$$JJ \longrightarrow north$$
$$NN \longrightarrow wind$$

**SCFG**

$$NP \longrightarrow DT_{1}NPB_{2}\ /\ DT_{1}NPB_{2}$$
$$NPB \longrightarrow JJ_{1}NPB_{2}\ /\ JJ_{1}NPB_{2}$$
$$NPB \longrightarrow JJ_{1}NPB_{2}\ /\ NPB_{2}JJ_{1}$$
$$NPB \longrightarrow NP_{1}\ /\ NP_{1}$$
$$DT \longrightarrow the\ /\ \varepsilon$$
$$JJ \longrightarrow strong\ /\ 呼啸$$
$$JJ \longrightarrow north\ /\ 北$$
$$NN \longrightarrow wind\ /\ 风$$

# CFG

NP $\longrightarrow$ DT NPB

NPB $\longrightarrow$ JJ NPB

NPB $\longrightarrow$ **NN**

DT $\longrightarrow$ the

JJ $\longrightarrow$ strong

JJ $\longrightarrow$ north

NN $\longrightarrow$ wind

# SCFG

NP $\longrightarrow$ DT$_{[1]}$NPB$_{[2]}$ / DT$_{[1]}$NPB$_{[2]}$

NPB $\longrightarrow$ JJ$_{[1]}$NPB$_{[2]}$ / JJ$_{[1]}$NPB$_{[2]}$

NPB $\longrightarrow$ JJ$_{[1]}$NPB$_{[2]}$ / NPB$_{[2]}$JJ$_{[1]}$

NPB $\longrightarrow$ NP$_{[1]}$ / NP$_{[1]}$

DT $\longrightarrow$ the / $\varepsilon$

JJ $\longrightarrow$ strong / 呼啸

JJ $\longrightarrow$ north / 北

NN $\longrightarrow$ wind / 风

# Noisy Channel

$$\boldsymbol{y}^* = \underset{\boldsymbol{y}}{\operatorname{argmax}} \; P(\boldsymbol{x} \mid \boldsymbol{y}) \, P(\boldsymbol{y})$$

# Noisy Channel

$$\boldsymbol{y}^* = \operatorname*{argmax}_{\boldsymbol{y}} P(\boldsymbol{x} \mid \boldsymbol{y}) \, P(\boldsymbol{y})$$

predicted
translation

source
sentence

# Noisy Channel

$$\boldsymbol{y}^* = \operatorname*{argmax}_{\boldsymbol{y}} P(\boldsymbol{x} \mid \boldsymbol{y}) \, P(\boldsymbol{y})$$

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

# Noisy Channel

$$\boldsymbol{y}^* = \operatorname*{argmax}_{\boldsymbol{y}} P(\boldsymbol{x} \mid \boldsymbol{y}) \, P(\boldsymbol{y})$$

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

$$\boldsymbol{y}^* = \operatorname*{argmax}_{\boldsymbol{y}} P(\boldsymbol{x} \mid \boldsymbol{y})^{\alpha} \, P(\boldsymbol{y})^{\beta}$$

# Noisy Channel

$$y^* = \underset{y}{\operatorname{argmax}}\; P(x \mid y)\, P(y)$$

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

$$y^* = \underset{y}{\operatorname{argmax}}\; P(x \mid y)^{\alpha}\, P(y)^{\beta}$$

$$= \underset{y}{\operatorname{argmax}}\; \alpha \log P(x \mid y) + \beta \log P(y)$$

extra parameters to tune, can tune to optimize BLEU

# Noisy Channel

$$y^* = \underset{y}{\mathrm{argmax}}\; P(x \mid y)\, P(y)$$

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

$$y^* = \underset{y}{\mathrm{argmax}}\; P(x \mid y)^{\alpha}\, P(y)^{\beta}$$

$$= \underset{y}{\mathrm{argmax}}\;\; \alpha \log P(x \mid y) + \beta \log P(y)$$

extra parameters to tune, can tune to optimize BLEU

## "tuning"

# Noisy Channel → Linear Model?

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

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

# Noisy Channel → Linear Model?

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

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

"word count feature":

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

# Noisy Channel → Linear Model?

$$y^* = \underset{y}{\operatorname{argmax}} \ \alpha \log P(x \mid y) + \beta \log P(y)$$

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

"word count feature":

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

"reverse translation model feature":

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

African
National
Congress

opposition　sanction　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

African
National
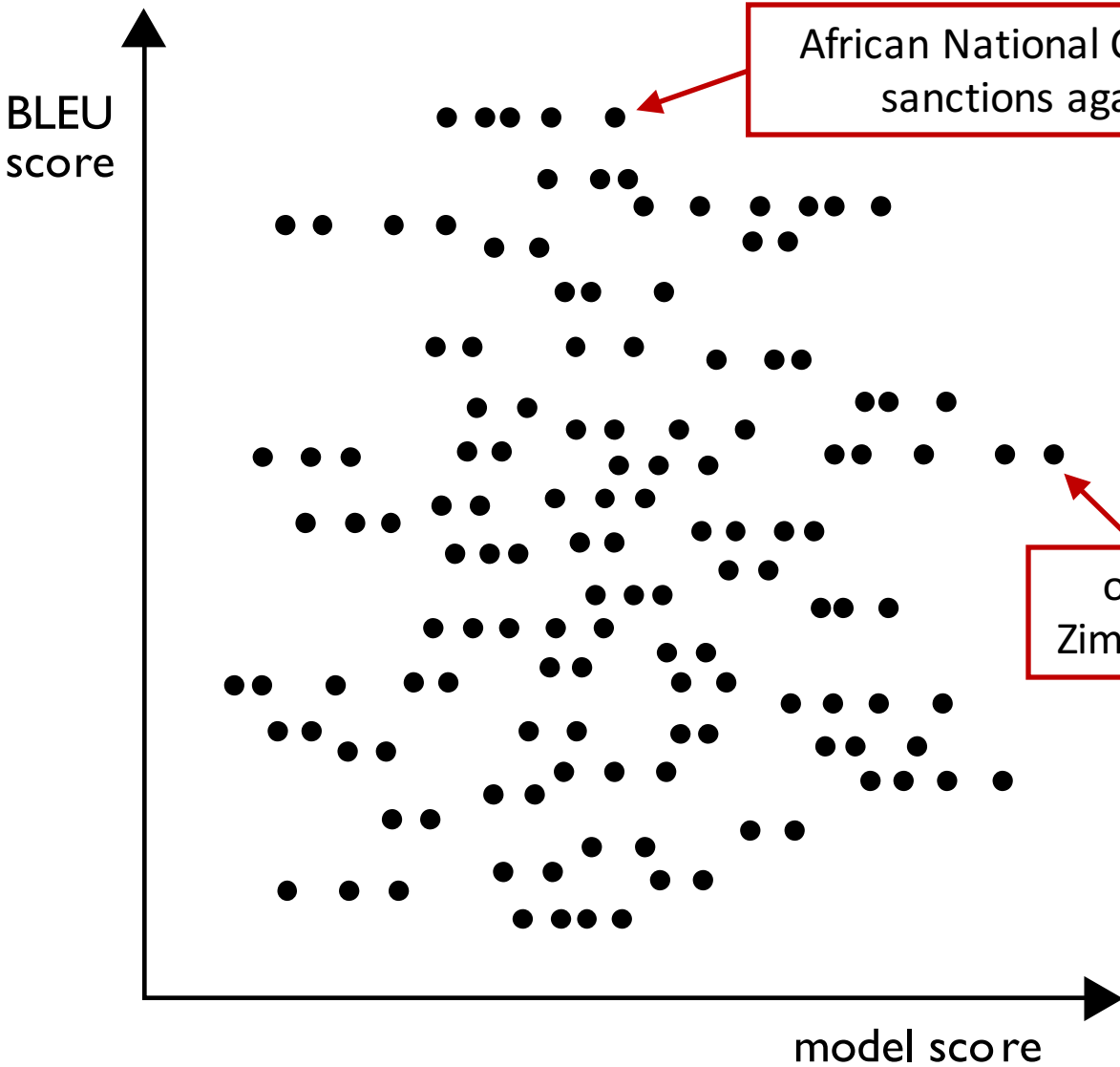Congress    opposition    sanction    Zimbabwe

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

predicted translation

opposition to sanctions against Zimbabwe African National Congress

model score

African National Congress    opposition    sanction    Zimbabwe

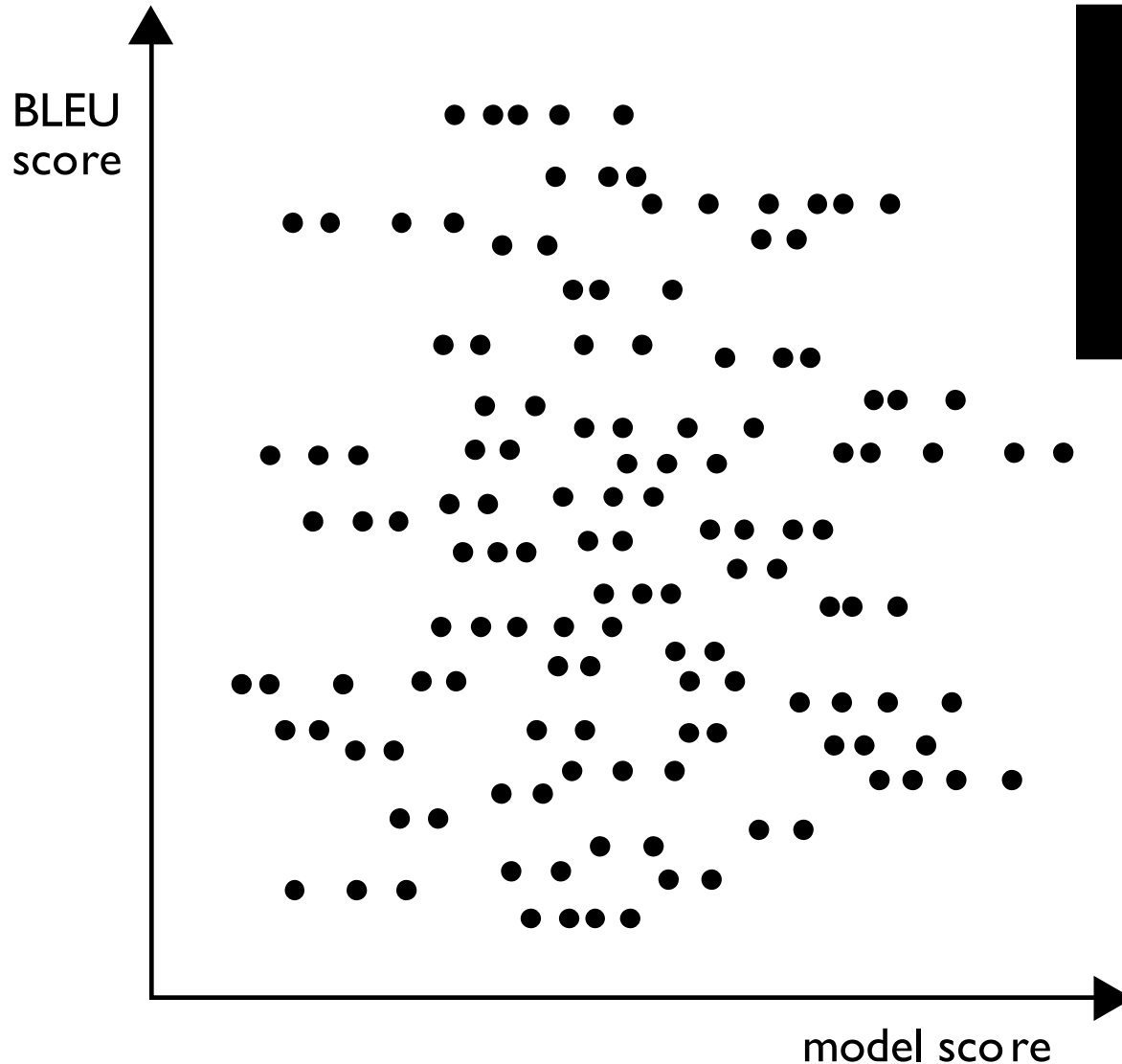非国大    反对    制裁    津巴布韦

**Gold standard:**
African National Congress opposes sanctions against Zimbabwe

BLEU score

African National Congress opposition sanctions against Zimbabwe

predicted translation

opposition to sanctions against Zimbabwe African National Congress

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

BLEU score

African National Congress opposition
sanctions against Zimbabwe

predicted translation

opposition to sanctions against
Zimbabwe African National Congress

African sanctioning to
Zimbabwe's opposing

model score

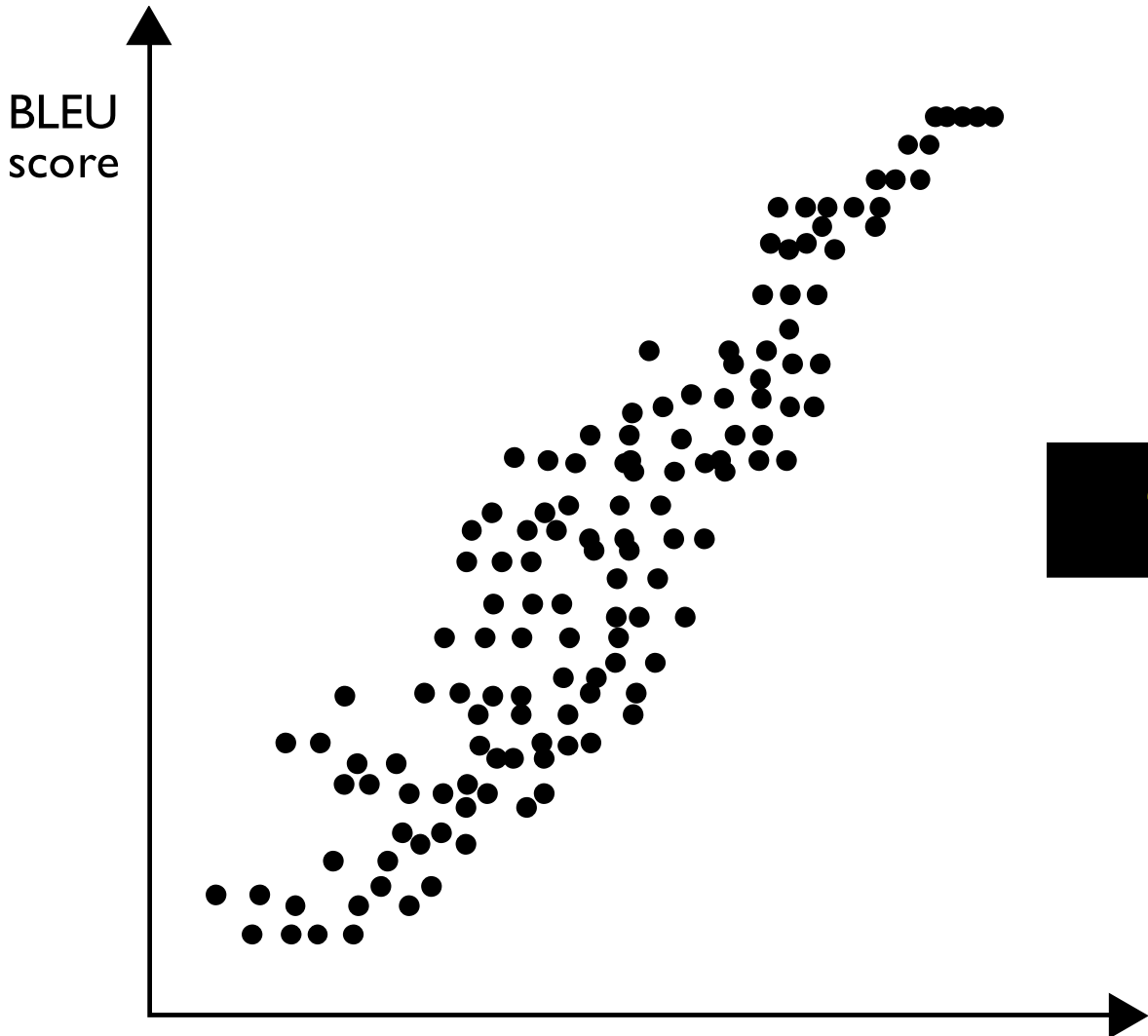African National Congress opposition sanction Zimbabwe

非国大　反对　制裁　津巴布韦

**Gold standard:**
African National Congress opposes
sanctions against Zimbabwe



BLEU score

model score

**"ideal" model**

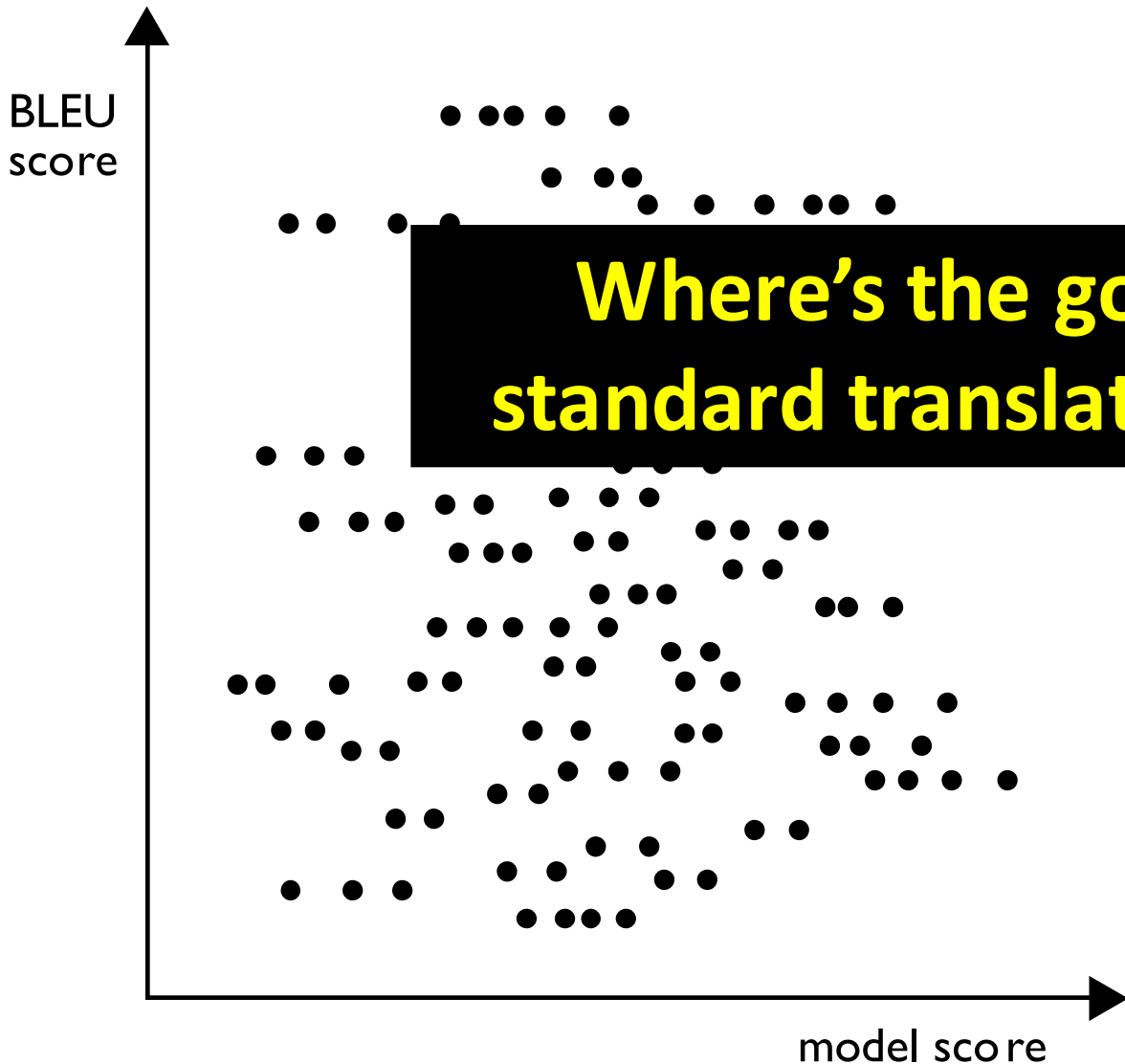African National Congress   opposition   sanction   Zimbabwe

非国大    反对    制裁    津巴布韦

**Gold standard:**
African National Congress opposes
sanctions against Zimbabwe

BLEU score

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

**Why?**

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

model score

# Free Translations

**Machine translation:**

Sharon's office said, leader of the main opposition Labor Party has admitted defeat and congratulatory telephone calls to Sharon.

**Human-generated translation:**

According to a representative of Sharon's office, the leader of the main opposition Labor Party has admitted defeat and made the obligatory congratulating telephone call to Sharon.

**Even if gold standard translation was reachable by model, we might not want to learn from it directly**

to Sharon.

**Applicable to other tasks:**
**summarization**
**image caption generation**

# Loss Functions

| name | loss | where used |
|------|------|------------|
| cost ("0-1") | $\text{cost}(y, \text{classify}(\boldsymbol{x}, \boldsymbol{\theta}))$ | intractable, but underlies "direct error minimization" |
| perceptron | $-\text{score}(\boldsymbol{x}, y, \boldsymbol{\theta}) + \max_{y' \in \mathcal{L}} \text{score}(\boldsymbol{x}, y', \boldsymbol{\theta})$ | perceptron algorithm (Rosenblatt, 1958) |
| hinge | $-\text{score}(\boldsymbol{x}, y, \boldsymbol{\theta}) + \max_{y' \in \mathcal{L}} (\text{score}(\boldsymbol{x}, y', \boldsymbol{\theta}) + \text{cost}(y, y'))$ | support vector machines, other large-margin algorithms |
| log | $-\log p_{\boldsymbol{\theta}}(y \mid \boldsymbol{x})$ $= \text{score}(\boldsymbol{x}, y, \boldsymbol{\theta}) + \log \sum_{y' \in \mathcal{L}} \exp\{\text{score}(\boldsymbol{x}, y', \boldsymbol{\theta})\}$ | logistic regression, conditional random fields, maximum entropy models |

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

| name | | where used |
|------|---|------------|
| cost ("0-1") | $\mathrm{cost}(y, \mathrm{classify}(\boldsymbol{x}, \boldsymbol{\theta}))$ | intractable, but underlies "direct error minimization" |
| perceptron | $-\mathrm{score}(\boldsymbol{x}, y, \boldsymbol{\theta}) + \max_{y' \in \mathcal{L}} \mathrm{score}(\boldsymbol{x}, y', \boldsymbol{\theta})$ | perceptron algorithm (Rosenblatt, 1958) |
| hinge | $-\mathrm{score}(\boldsymbol{x}, y, \boldsymbol{\theta}) + \max_{y' \in \mathcal{L}} (\mathrm{score}(\boldsymbol{x}, y', \boldsymbol{\theta}) + \mathrm{cost}(y, y'))$ | support vector machines, other large-margin algorithms |
| log | $-\log p_{\boldsymbol{\theta}}(y \mid \boldsymbol{x})$ $= \mathrm{score}(\boldsymbol{x}, y, \boldsymbol{\theta}) + \log \sum_{y' \in \mathcal{L}} \exp\{\mathrm{score}(\boldsymbol{x}, y', \boldsymbol{\theta})\}$ | logistic regression, conditional random fields, maximum entropy models |

**intractable, but it doesn't need to compute model score of gold standard!**

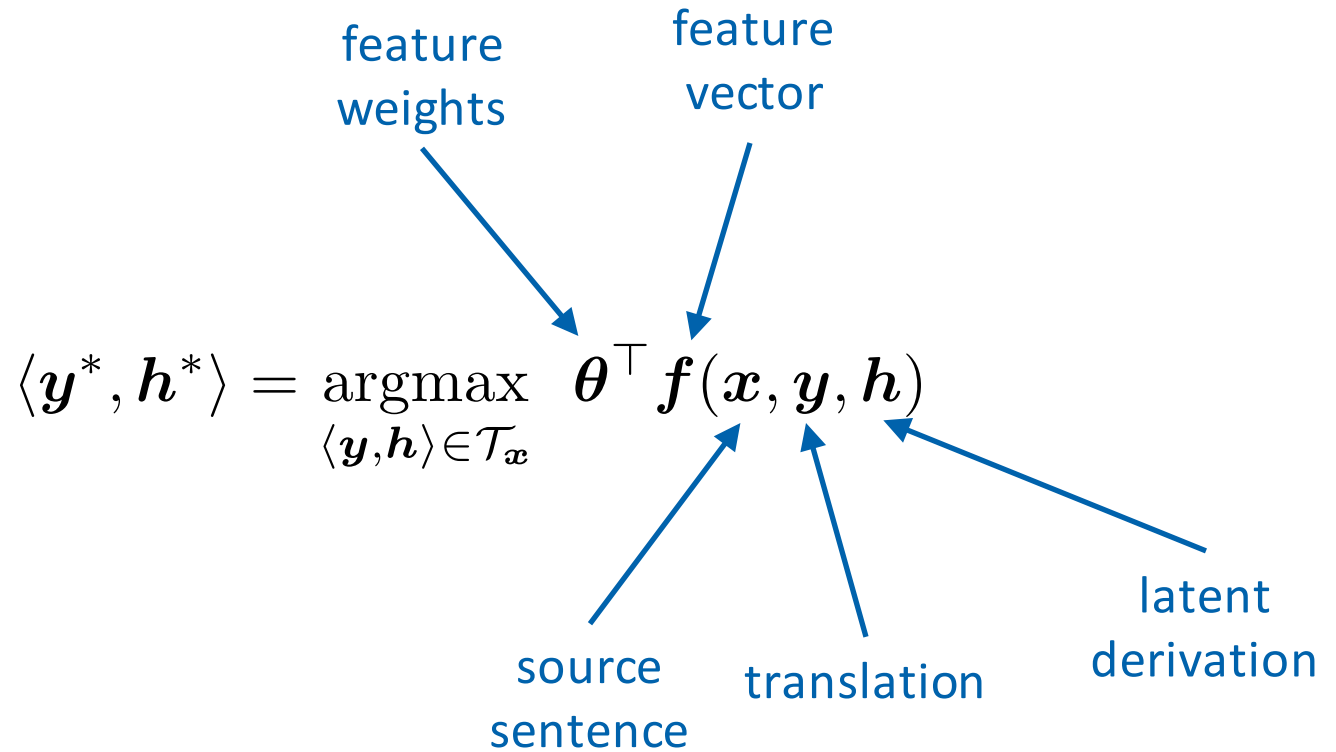| name | | where used |
|------|---|------------|
| cost ("0-1") | $\mathrm{cost}(y, \mathrm{classify}(\boldsymbol{x}, \boldsymbol{\theta}))$ | intractable, but underlies "direct error minimization" |
| perceptron | $-\mathrm{score}(\boldsymbol{x}, y, \boldsymbol{\theta}) + \max_{y' \in \mathcal{L}} \mathrm{score}(\boldsymbol{x}, y', \boldsymbol{\theta})$ | perceptron algorithm (Rosenblatt, 1958) |
| hinge | $-\mathrm{score}(\boldsymbol{x}, y, \boldsymbol{\theta}) + \max_{y' \in \mathcal{L}} (\mathrm{score}(\boldsymbol{x}, y', \boldsymbol{\theta}) + \mathrm{cost}(y, y'))$ | support vector machines, other large-margin algorithms |
| log | $-\log p_{\boldsymbol{\theta}}(y \mid \boldsymbol{x})$ $= \mathrm{score}(\boldsymbol{x}, y, \boldsymbol{\theta}) + \log \sum_{y' \in \mathcal{L}} \exp\{\mathrm{score}(\boldsymbol{x}, y', \boldsymbol{\theta})\}$ | logistic regression, conditional random fields, maximum entropy models |

# MERT, Och (2003)

**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

# Notation

feature weights

feature vector

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

source sentence

translation

latent derivation

# Minimum Error Rate Training (MERT)

$$\min_{\boldsymbol{\theta}} \quad \text{cost}\left( \{\boldsymbol{y}^{(i)}\}_{i=1}^{N}, \left\{ \operatorname*{argmax}_{\langle \boldsymbol{y}, \boldsymbol{h} \rangle \in \mathcal{T}_{\boldsymbol{x}^{(i)}}} \boldsymbol{\theta}^{\top} \boldsymbol{f}(\boldsymbol{x}^{(i)}, \boldsymbol{y}, \boldsymbol{h}) \right\}_{i=1}^{N} \right)$$

# Minimum Error Rate Training (MERT)

set of source sentences

$$\min_{\boldsymbol{\theta}} \quad \text{cost} \left( \underbrace{\{\boldsymbol{y}^{(i)}\}_{i=1}^{N}}_{\text{references}}, \underbrace{\left\{ \operatorname*{argmax}_{\langle \boldsymbol{y}, \boldsymbol{h} \rangle \in \mathcal{T}_{\boldsymbol{x}^{(i)}}} \boldsymbol{\theta}^{\top} \boldsymbol{f}(\boldsymbol{x}^{(i)}, \boldsymbol{y}, \boldsymbol{h}) \right\}_{i=1}^{N}}_{\substack{\text{decoder} \\ \text{outputs}}} \right)$$

# Minimum Error Rate Training (MERT)

"how bad are these translations?"
e.g., negative BLEU

set of source sentences

$$\min_{\boldsymbol{\theta}} \quad \text{cost} \left( \underbrace{\{\boldsymbol{y}^{(i)}\}_{i=1}^{N}}_{\text{references}}, \underbrace{\left\{ \operatorname*{argmax}_{\langle \boldsymbol{y}, \boldsymbol{h} \rangle \in \mathcal{T}_{\boldsymbol{x}^{(i)}}} \boldsymbol{\theta}^{\top} \boldsymbol{f}(\boldsymbol{x}^{(i)}, \boldsymbol{y}, \boldsymbol{h}) \right\}_{i=1}^{N}}_{\substack{\text{decoder} \\ \text{outputs}}} \right)$$

# Minimum Error Rate Training (MERT)

minimize the cost of the decoder output

intractable in general – how can we solve it?

$$\min_{\boldsymbol{\theta}} \quad \text{cost}\left(\underbrace{\{\boldsymbol{y}^{(i)}\}_{i=1}^{N}}_{\text{references}}, \underbrace{\left\{\operatorname*{argmax}_{\langle\boldsymbol{y},\boldsymbol{h}\rangle\in\mathcal{T}_{\boldsymbol{x}^{(i)}}} \boldsymbol{\theta}^{\top}\boldsymbol{f}(\boldsymbol{x}^{(i)},\boldsymbol{y},\boldsymbol{h})\right\}_{i=1}^{N}}_{\substack{\text{decoder}\\\text{outputs}}}\right)$$

# Minimum Error Rate Training (MERT)

"h

minimize the cost of the decoder output

intractable in general – how can we solve it?

$$\min_{\theta} \ \mathrm{cost}\left(\{\boldsymbol{y}^{(i)}\}_{i=1}^{N}, \left\{ \ \mathrm{argmax} \ \ \boldsymbol{\theta}^{\top}\boldsymbol{f}(\boldsymbol{x}^{(i)}, \boldsymbol{y}, \boldsymbol{h}) \right\}^{N} \right)$$

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

BLEU

each point is a translation
for the same sentence

Arabic-English,
phrase-based

model score

BLEU (y-axis)
model score (x-axis)

10,000-best list,
default Moses weights

1-best:
28 BLEU

**BLEU** (y-axis)

**model score** (x-axis)

same sentence,
10,000-best list
after MERT

1-best:
34 BLEU

another sentence,
default Moses weights

1-best:
46 BLEU

BLEU

model score
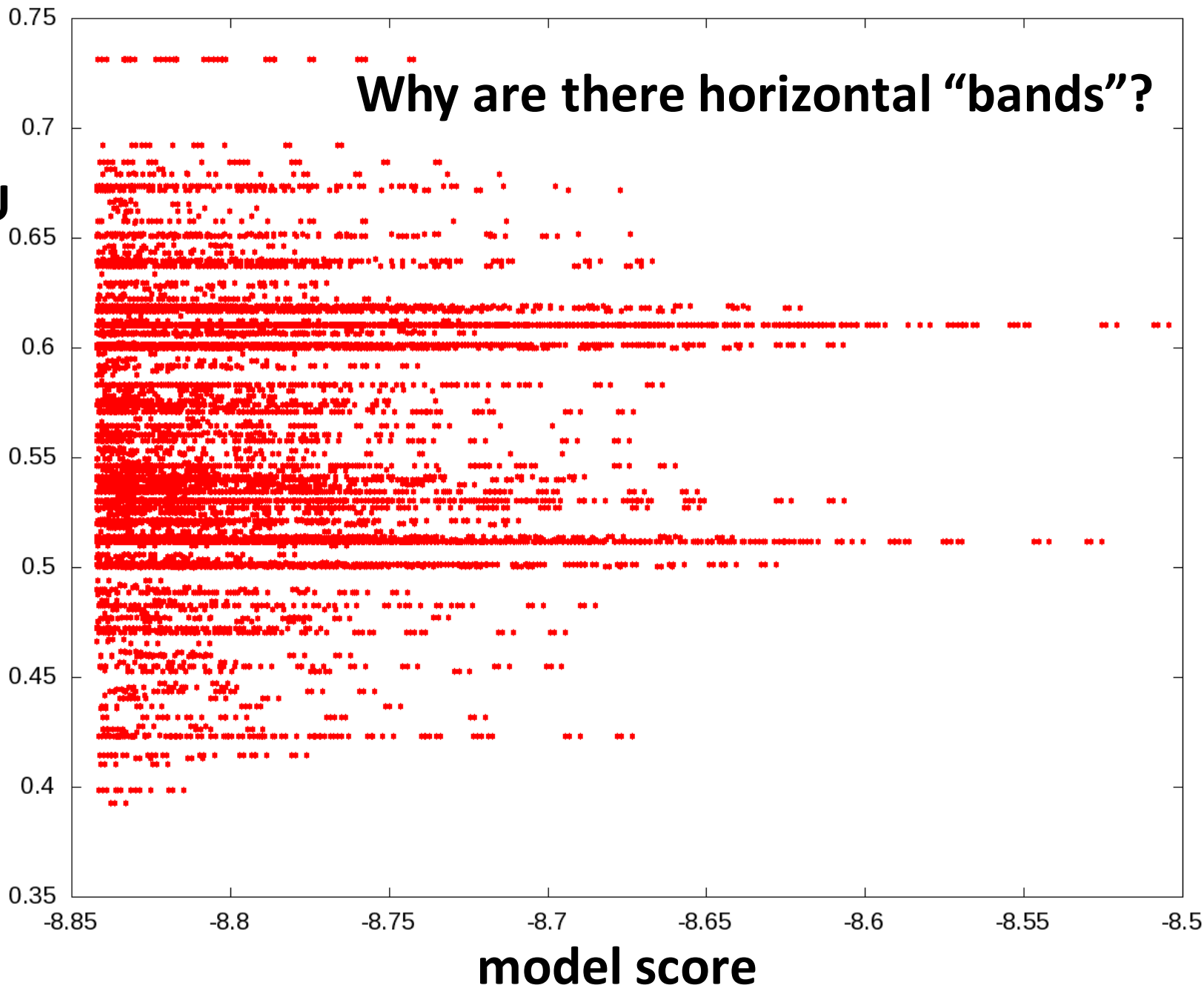
same sentence,
after MERT

1-best:
62 BLEU

BLEU

model score

$$\min_{\boldsymbol{\theta}} \quad \text{cost} \left( \underbrace{\{\boldsymbol{y}^{(i)}\}_{i=1}^{N}}_{\text{references}}, \left\{ \underbrace{\underset{\langle \boldsymbol{y}, \boldsymbol{h} \rangle \in \mathcal{T}_{\boldsymbol{x}^{(i)}}}{\text{argmax}} \quad \boldsymbol{\theta}^{\top} \boldsymbol{f}(\boldsymbol{x}^{(i)}, \boldsymbol{y}, \boldsymbol{h})}_{\substack{\text{decoder} \\ \text{outputs}}} \right\}_{i=1}^{N} \right)$$

## What are some issues with this loss function?

Discontinuous & non-convex → optimization relies on randomized search

No regularization → leads to overfitting

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

Many researchers tried to improve MERT:

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

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

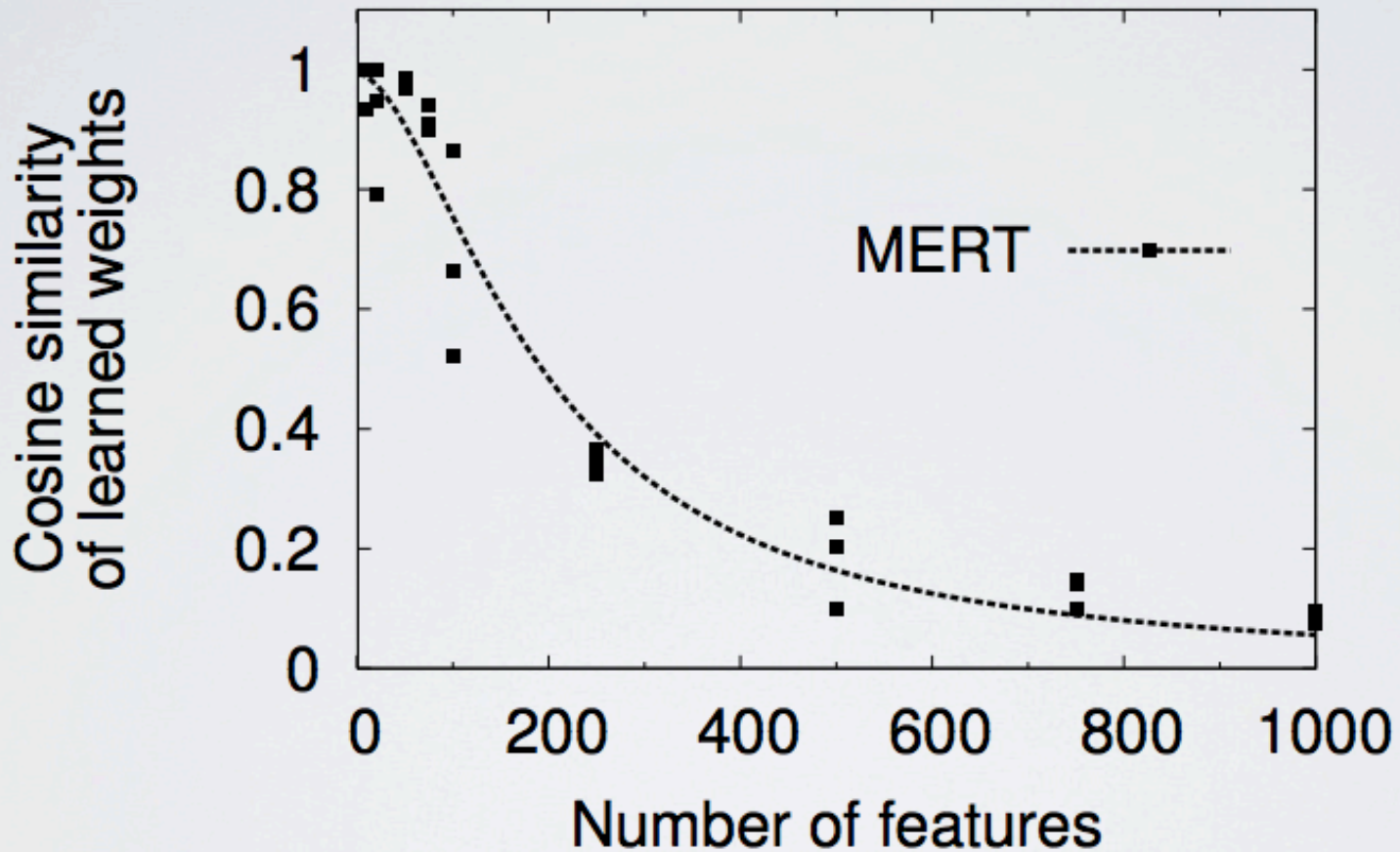*Stabilizing MERT* (Foster & Kuhn, 2009)

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*
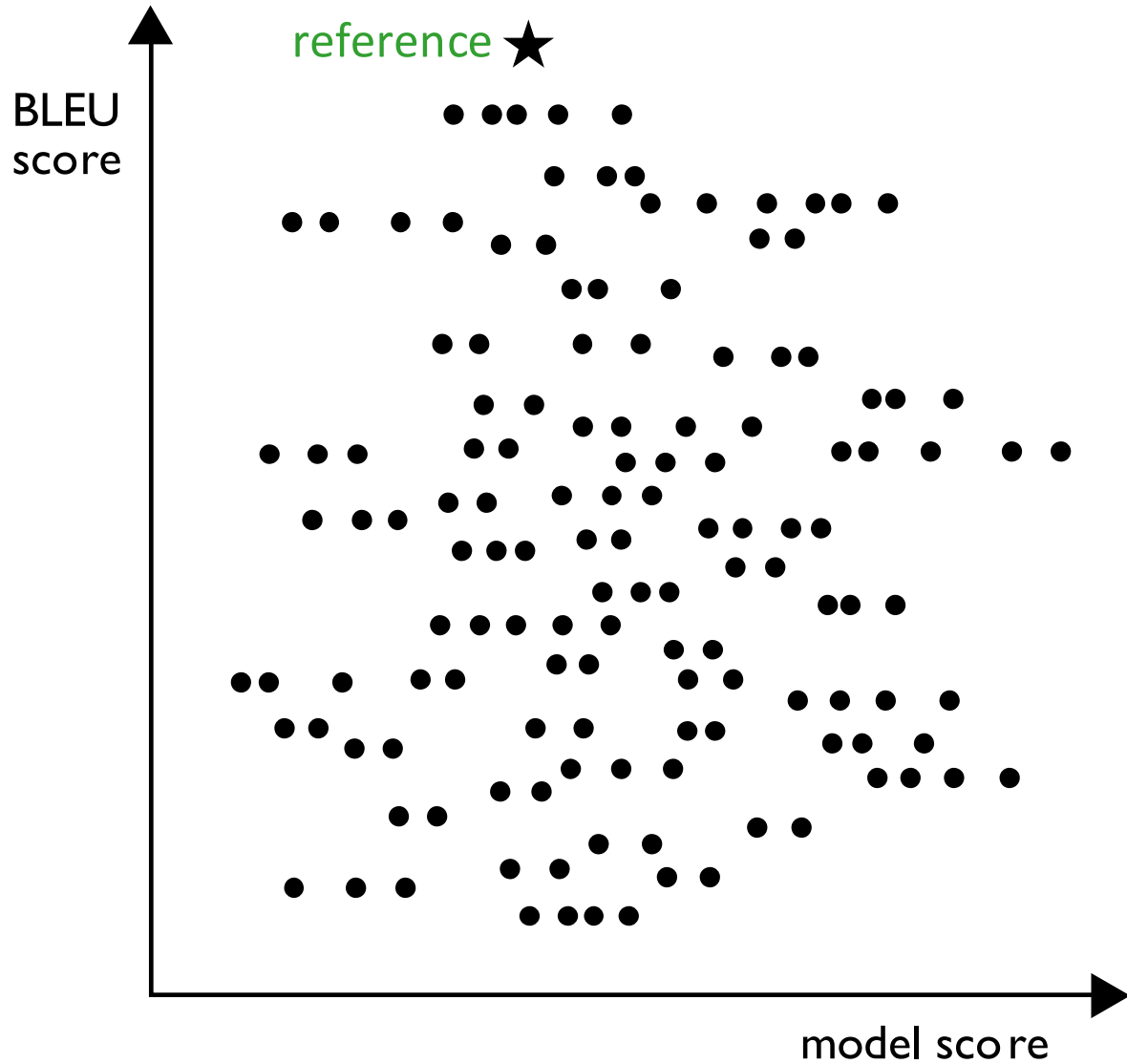
## Synthetic weight learning of MERT

**Tuning as Ranking**

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

1

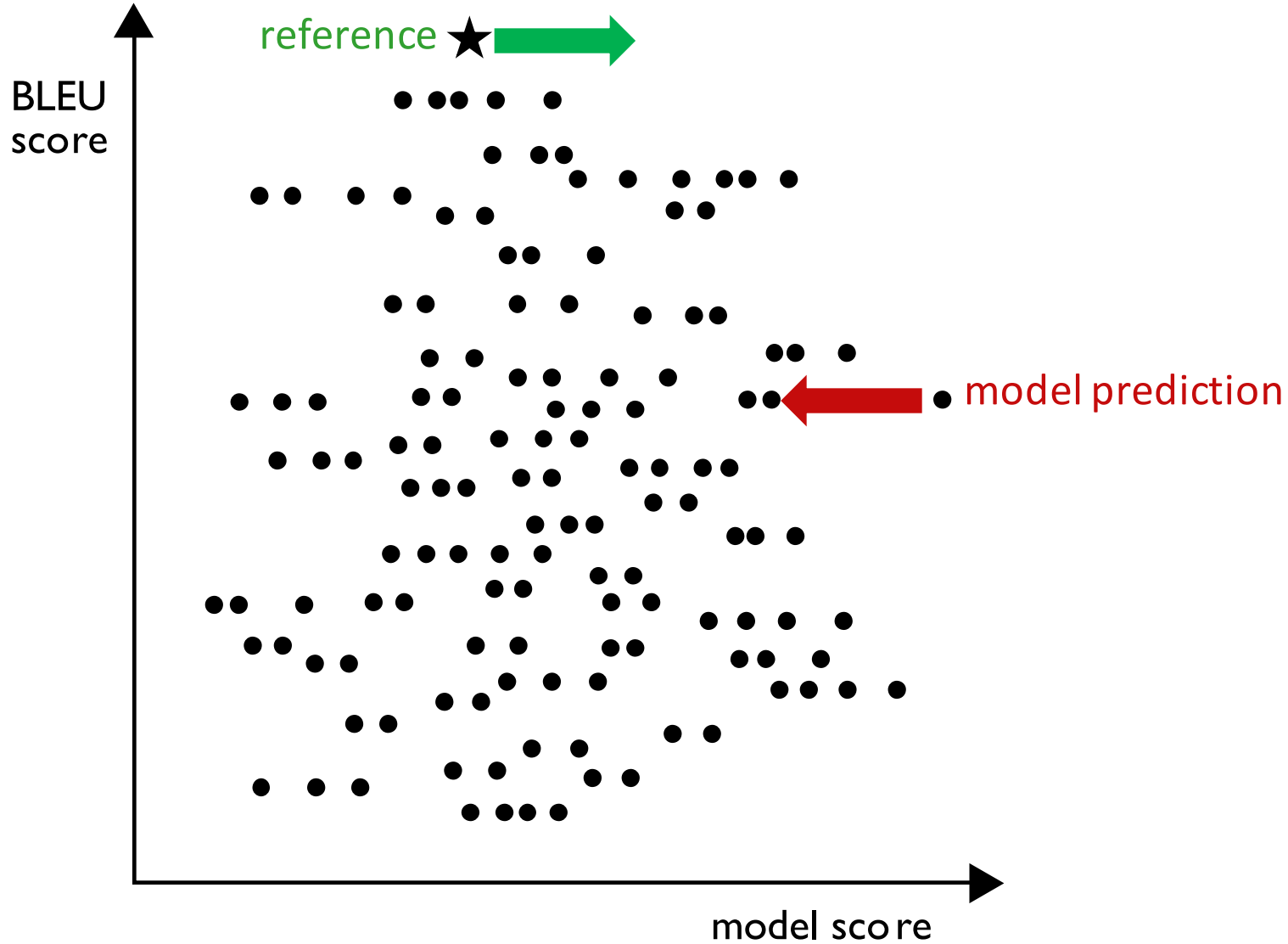0    200    400    600    800    1000

### Number of features

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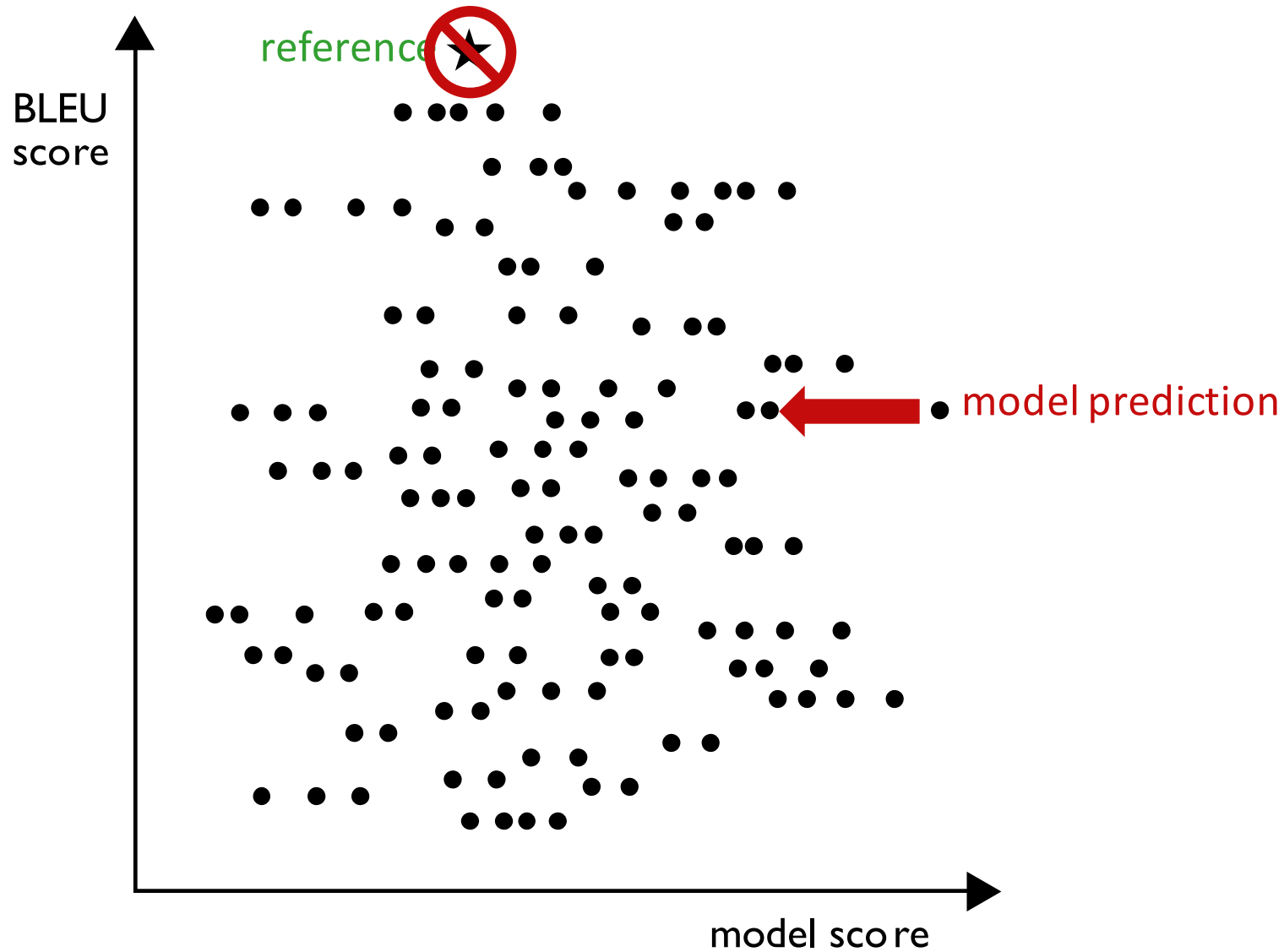validates what has long been accepted as truth
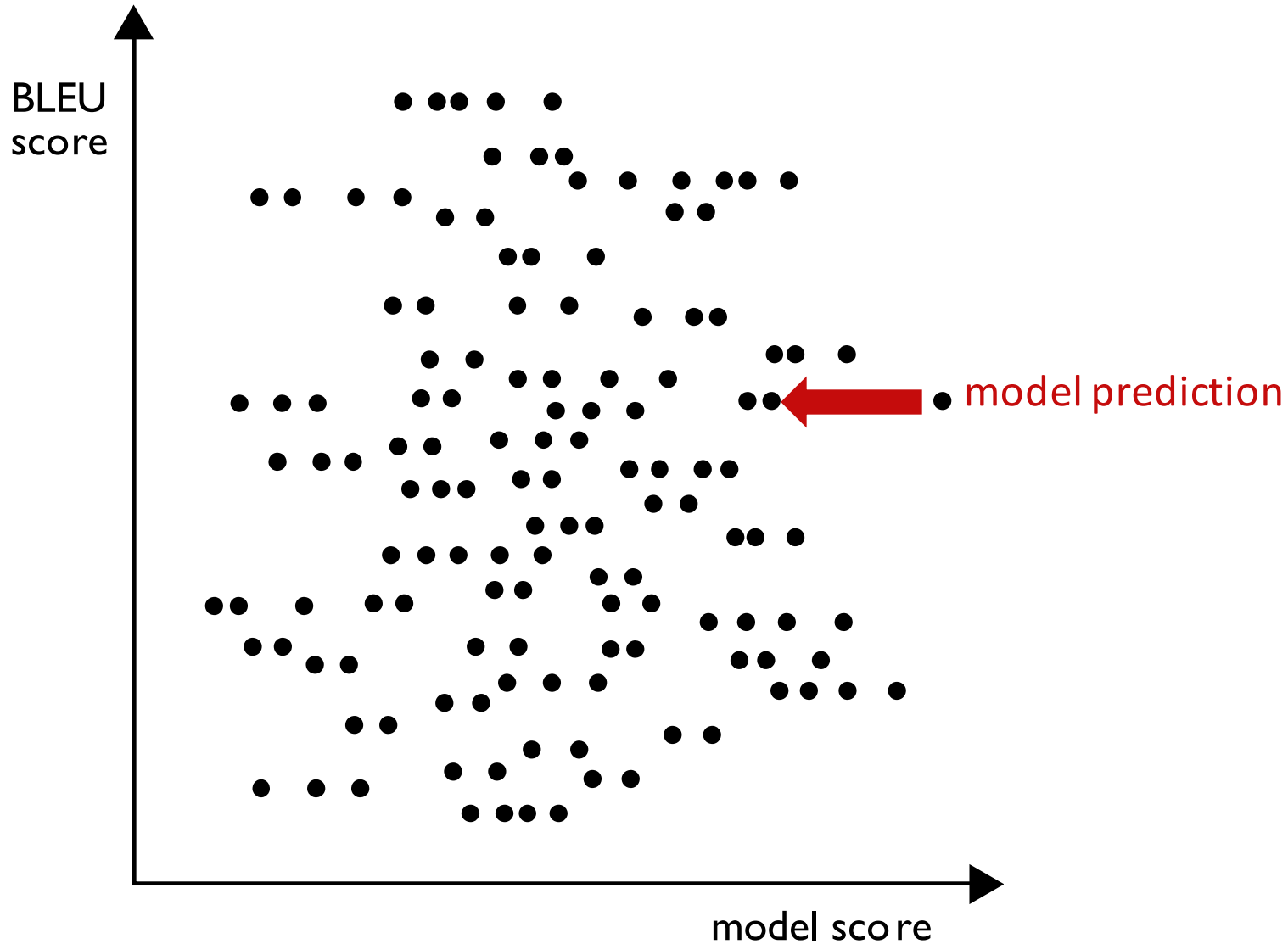
# Perceptron Loss

# Perceptron Loss

# Perceptron Loss for MT?

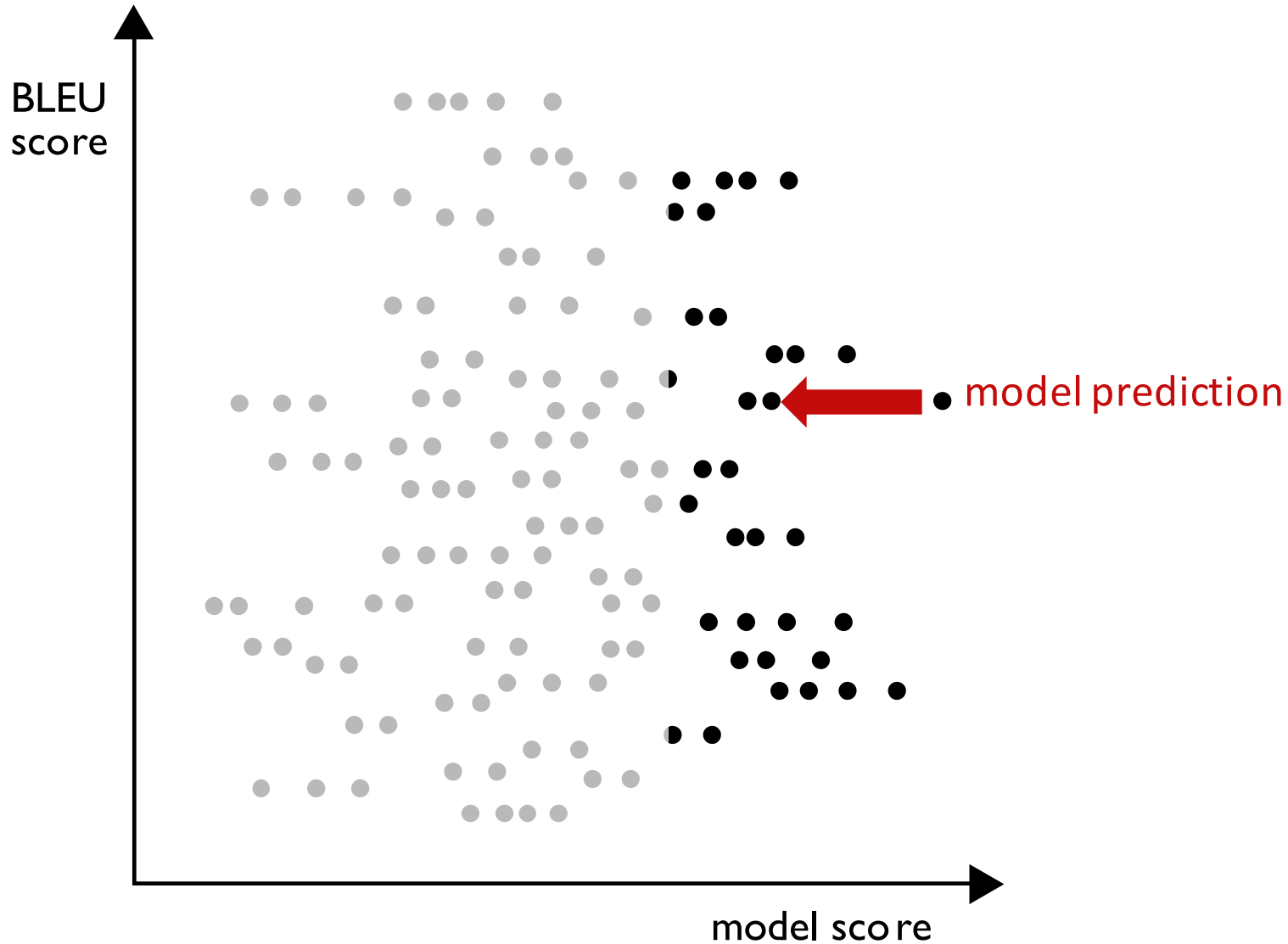# k-Best Perceptron for MT

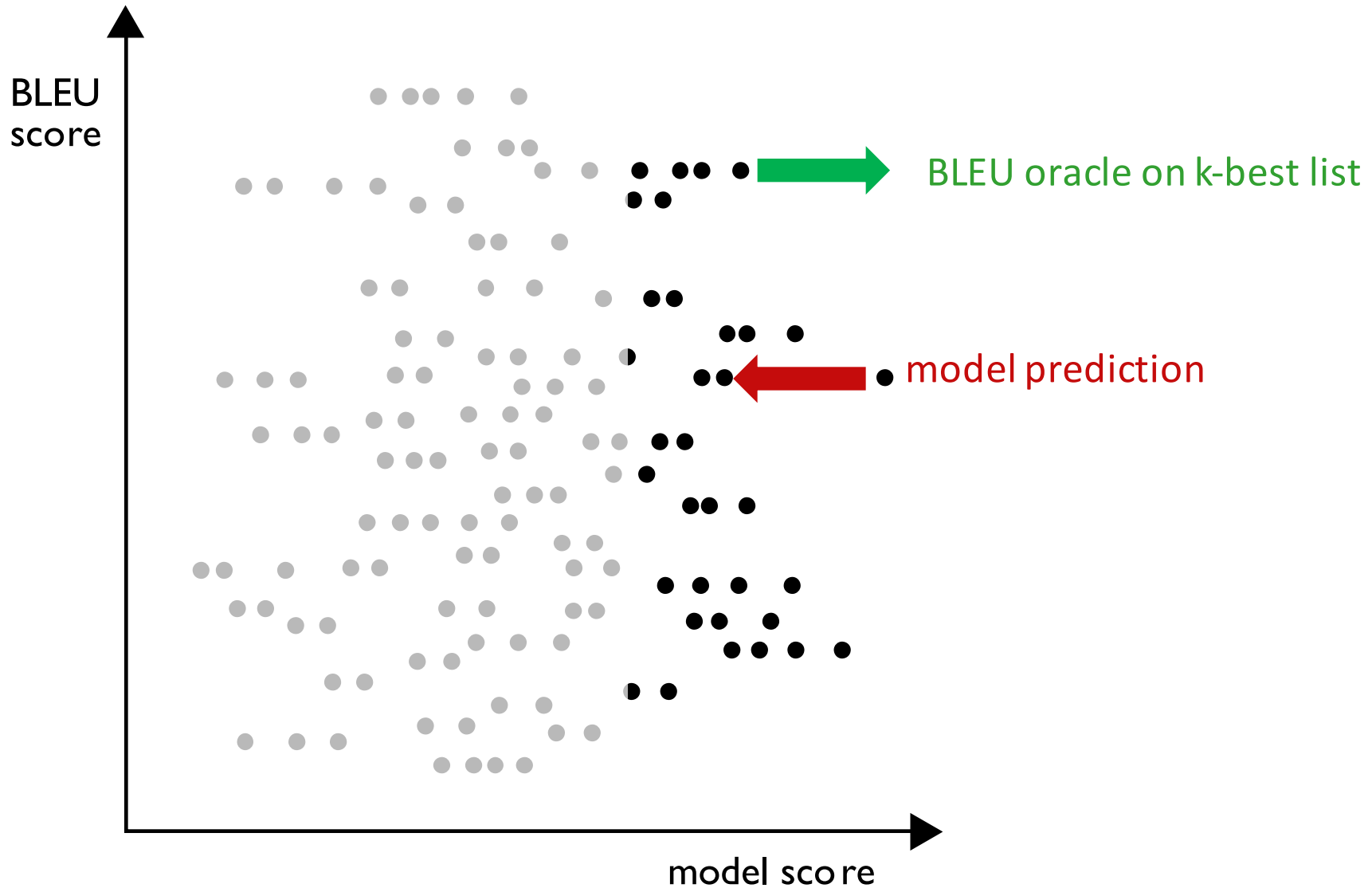(Liang et al., 2006)



model prediction

# k-Best Perceptron for MT
(Liang et al., 2006)

# k-Best Perceptron for MT
## (Liang et al., 2006)



BLEU oracle on k-best list

model prediction

# Ramp Loss Minimization

# Ramp Loss Minimization



model prediction

BLEU score

model score

# Ramp Loss Minimization



model prediction

"fear" translation

BLEU score

model score

# "Fear" Ramp Loss
## (Do et al., 2008)



BLEU score

$$\max_{\boldsymbol{y} \in \mathcal{Y}(\boldsymbol{x})} \mathrm{score}(\boldsymbol{x}^{(i)}, \boldsymbol{y})$$

model prediction

"fear" translation

$$\max_{\boldsymbol{y} \in \mathcal{Y}(\boldsymbol{x})} \Big( \mathrm{score}(\boldsymbol{x}^{(i)}, \boldsymbol{y}) + \mathrm{cost}(\boldsymbol{y}^{(i)}, \boldsymbol{y}) \Big)$$

model score

# "Fear" Ramp Loss
(Do et al., 2008)



BLEU score

$$\max_{\boldsymbol{y} \in \mathcal{Y}(\boldsymbol{x})} \operatorname{score}(\boldsymbol{x}^{(i)}, \boldsymbol{y})$$

→ model prediction

gold standard

← "fear" translation

$$\max_{\boldsymbol{y} \in \mathcal{Y}(\boldsymbol{x})} \Big( \operatorname{score}(\boldsymbol{x}^{(i)}, \boldsymbol{y}) + \operatorname{cost}(\boldsymbol{y}^{(i)}, \boldsymbol{y}) \Big)$$

model score

# "Hope" Ramp Loss

(McAllester & Keshet, 2011; Liang et al., 2006)

# "Hope" Ramp Loss

$$\max_{\boldsymbol{y} \in \mathcal{Y}(\boldsymbol{x})} \left( \mathrm{score}(\boldsymbol{x}^{(i)}, \boldsymbol{y}) - \mathrm{cost}(\boldsymbol{y}^{(i)}, \boldsymbol{y}) \right)$$

"hope" translation

model prediction

BLEU score

model score

# "Hope-Fear" Ramp Loss

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



BLEU score

→ "hope" translation

← "fear" translation

model score

# "Hope-Fear" Ramp Loss

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



BLEU score

→ "hope" translation

$$\operatorname*{argmax}_{\langle \boldsymbol{y}, \boldsymbol{h}\rangle \in \mathcal{T}_{\boldsymbol{x}^{(i)}}} \left( \boldsymbol{\theta}^\top \boldsymbol{f}(\boldsymbol{x}^{(i)}, \boldsymbol{y}, \boldsymbol{h}) - \operatorname{cost}(\boldsymbol{y}^{(i)}, \boldsymbol{y}) \right)$$

$$\operatorname*{argmax}_{\langle \boldsymbol{y}, \boldsymbol{h}\rangle \in \mathcal{T}_{\boldsymbol{x}^{(i)}}} \left( \boldsymbol{\theta}^\top \boldsymbol{f}(\boldsymbol{x}^{(i)}, \boldsymbol{y}, \boldsymbol{h}) + \operatorname{cost}(\boldsymbol{y}^{(i)}, \boldsymbol{y}) \right)$$

← "fear" translation

model score

# Experiments

averages over 8 test sets across 3 language pairs

| | Moses | Hiero |
|---|---|---|
| | %BLEU | %BLEU |
| MERT | **35.9** | **37.0** |
| Fear Ramp (away from bad) | 34.9 | 34.2 |
| Hope Ramp (toward good) | 35.2 | 36.0 |
| Hope-Fear Ramp (toward good + away from bad) | 35.7 | **37.0** |

# Pairwise Ranking Optimization

## (Hopkins & May, 2011)