

Machine Translation Zoo

Tree-to-tree transfer and Discriminative learning

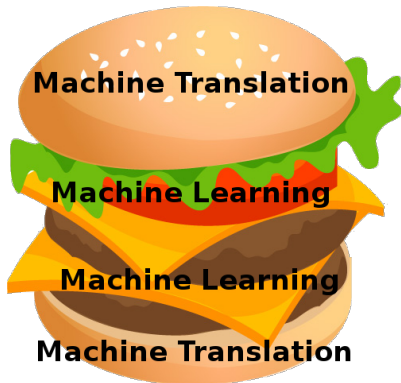
Martin Popel

ÚFAL (Institute of Formal and Applied Linguistics)
Charles University in Prague

May 5th 2013, Seminar of Formal Linguistics, Prague

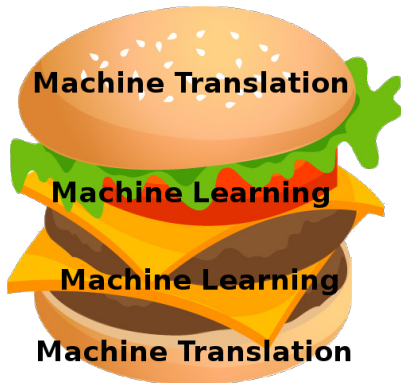
Today's Menu

- 1 MT Intro
 - Taxonomy
 - Hybrids
- 2 Online Learning
 - Perceptron
 - Structured Prediction
- 3 Guided Learning
- 4 Back to MT
 - Easy-First Decoding in MT
 - Guided Learning in MT



Today's Menu

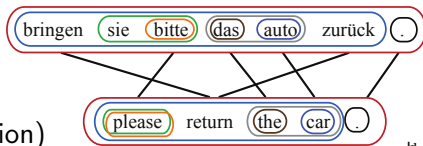
- 1 MT Intro
 - Taxonomy
 - Hybrids
- 2 Online Learning
 - Perceptron
 - Structured Prediction
- 3 Guided Learning
- 4 Back to MT
 - Easy-First Decoding in MT
 - Guided Learning in MT



Phrase-based MT (Moses)

Training

- word-alignment
(GIZA++ & symmetrization)
- phrase extraction
- tune parameters (MERT)



Decoding

- get all matching *rules*
- find one *derivation*
with a maximum score (beam search)



TectoMT

Training

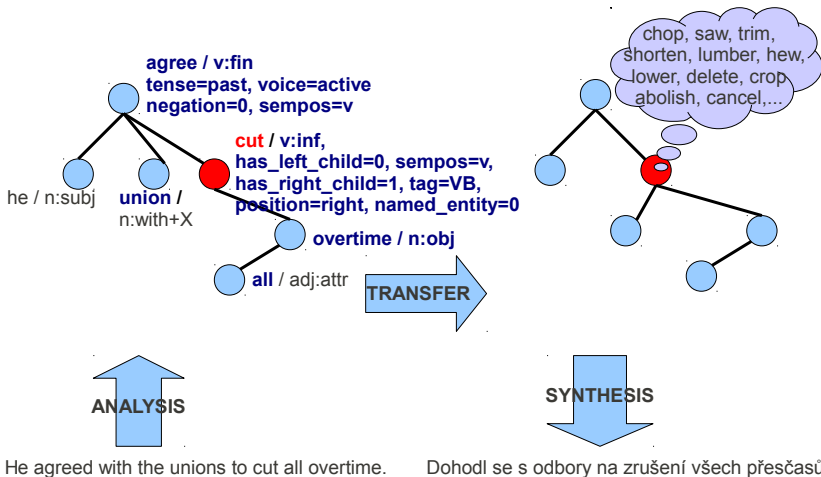
- analyze CzEng to t-layer
- t-node alignment
- learn one MaxEnt model for each source lemma and formeme

Decoding

- get all translation variants for each lemma and formeme
- find a labeling with a maximum score (HMTM)



TectoMT – MaxEnt Model



Machine Translation Taxonomy

- Level of transfer:
- Base translation unit (BTU):
- Extract more segmentations in training?
- Try (search) more segmentations in decoding?
- Use more segmentations in the output translation?
- What is the context X in $P(BTU_{target} | BTU_{source}, X)$?

Machine Translation Taxonomy

- Level of transfer: **surface**
- Base translation unit (BTU):
word
- Extract more segmentations in training? **no**
- Try (search) more segmentations in decoding? **no**
- Use more segmentations in the output translation? **no**
- What is the context X in $P(BTU_{target} | BTU_{source}, X)$?
Considering just Translation Model: **nothing**



(Brown et al., 1993)
word-based

Machine Translation Taxonomy

- Level of transfer: **surface**
- Base translation unit (BTU):
word, **phrase**
- Extract more segmentations in training? **yes**
- Try (search) more segmentations in decoding? **yes**
- Use more segmentations in the output translation? **no**
- What is the context X in $P(BTU_{target} | BTU_{source}, X)$?
Considering just Translation Model: **nothing**



(Brown et al., 1993)
word-based



(Koehn et al., 2003)
phrase-based

Machine Translation Taxonomy

- Level of transfer: **surface**
- Base translation unit (BTU):
word, phrase, **phrase with gaps**
- Extract more segmentations in training? **yes**
- Try (search) more segmentations in decoding? **yes**
- Use more segmentations in the output translation? **no**
- What is the context X in $P(BTU_{target} | BTU_{source}, X)$?
Considering just Translation Model: **nothing**



(Brown et al., 1993)
word-based



(Koehn et al., 2003)
phrase-based



(Chiang, 2005)
hierarchical

Machine Translation Taxonomy

- Level of transfer: surface, **shallow syntax**
- Base translation unit (BTU):
word, phrase, phrase with gaps, **treelet**
- Extract more segmentations in training? **no**
- Try (search) more segmentations in decoding? **no**
- Use more segmentations in the output translation? **no**
- What is the context X in $P(BTU_{target} | BTU_{source}, X)$?
Considering just Translation Model: **neighboring treelets**



Machine Translation Taxonomy

- Level of transfer: surface, shallow syntax, **tectogrammatical**
- Base translation unit (BTU):
word, phrase, phrase with gaps, treelet, **node**
- Extract more segmentations in training? **no**
- Try (search) more segmentations in decoding? **no**
- Use more segmentations in the output translation? **no**
- What is the context X in $P(BTU_{target} | BTU_{source}, X)$?
Considering just Translation Model: **neighboring nodes**



(Koehn et al., 2003)
phrase-based



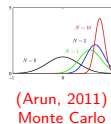
(Chiang, 2005)
hierarchical



(Mareček et al., 2010)
TectoMT

Machine Translation Taxonomy

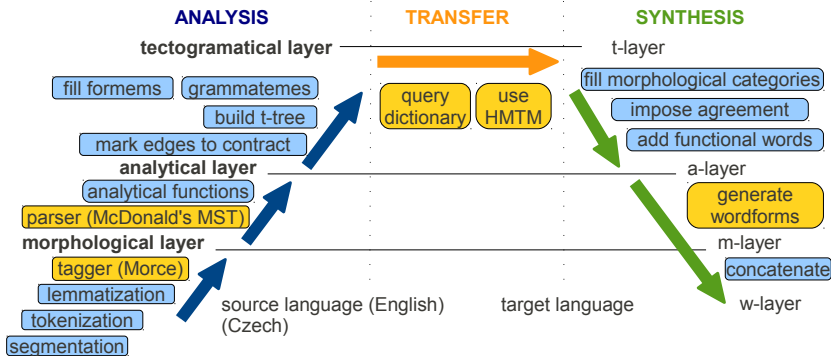
- Level of transfer: **surface**, shallow syntax, tectogrammatical
- Base translation unit (BTU):
word, **phrase**, phrase with gaps, treelet, node
- Extract more segmentations in training? **yes**
- Try (search) more segmentations in decoding? **yes**
- Use more segmentations in the output translation? **yes**
- What is the context X in $P(BTU_{target} | BTU_{source}, X)$?
Considering just Translation Model: **nothing**



Hybrids: TectoMoses

Linearize source t-trees (two factors: lemma and formeme), translate with Moses, project dependencies and use TectoMT synthesis.

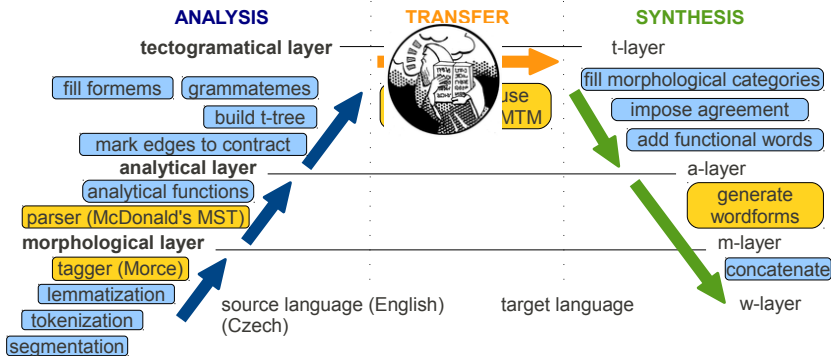
rule based & statistical blocks



Hybrids: TectoMoses

Linearize source t-trees (two factors: lemma and formeme), translate with Moses, project dependencies and use TectoMT synthesis.

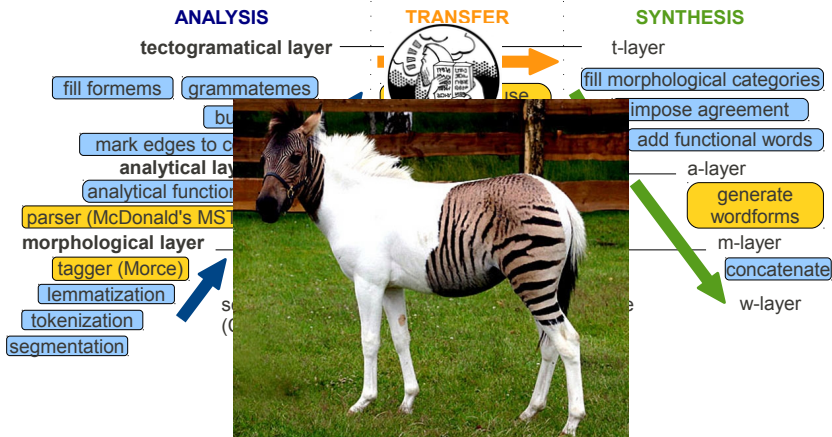
rule based & statistical blocks



Hybrids: TectoMoses

Linearize source t-trees (two factors: lemma and formeme), translate with Moses, project dependencies and use TectoMT synthesis.

rule based & statistical blocks



Hybrids: PhraseFix

Done for WMT 2013 by Petra Galuščáková:

- Post-edit TectoMT output using Moses
- trained on cs-tectomt → cs-reference (whole CzEng).
- How to post-edit only when confident?
 - filter phrase table
 - add “confidence” feature for MERT
 - improve alignment (monolingual)
 - boost phrase table (e.g. with identities)

Future work:

- use also source (English) sentences ⇒ multi-source translation
- project only content words (using TectoMT)
- factored translation with non-synchronous (overlapping) factors

Hybrids: PhraseFix

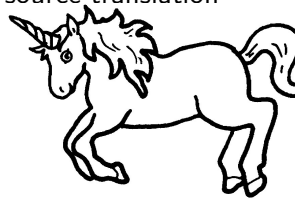
Done for WMT 2013 by Petra Galuščáková

- Post-edit TectoMT output using Moses
- trained on cs-tectomt → cs-reference (1)
- How to post-edit only when confident?
 - filter phrase table
 - add “confidence” feature for MERT
 - improve alignment (monolingual)
 - boost phrase table (e.g. with identities)



Future work:

- use also source (English) sentences ⇒ multi-source translation
- project only content words (using TectoMT)
- factored translation with non-synchronous (overlapping) factors



Even More Hybrids: DepFix, AddToTrain, Chimera

DepFix (Rosa et al., 2012)

- post-edit SMT using syntactic analysis and rules
- exploit also the source sentences, robust parsing

AddToTrain (Bojar, Galuščáková)

- translate monolingual news (or WMT devsets) with TectoMT
- add this to Moses parallel training data

Chimera

- post-edit AddToTrain output with DepFix
- sent to WMT 2013 in attempt to beat Google

Even More Hybrids: DepFix, AddToTrain, Chimera

DepFix (Rosa et al., 2012)

- post-edit SMT using syntactic analysis and rules
- exploit also the source sentences, robust parsing

AddToTrain (Bojar, Galuščáková)

- translate monolingual news (or WMT devsets)
- add this to Moses parallel training data

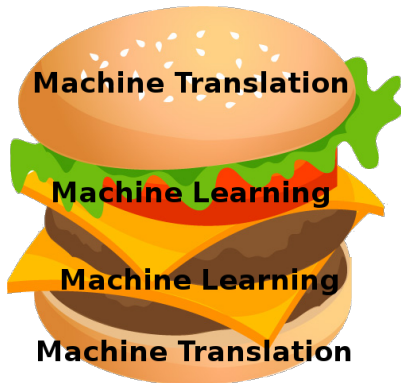
Chimera

- post-edit AddToTrain output with DepFix
- sent to WMT 2013 in attempt to beat Google



Today's Menu

- 1 MT Intro
 - Taxonomy
 - Hybrids
- 2 Online Learning
 - Perceptron
 - Structured Prediction
- 3 Guided Learning
- 4 Back to MT
 - Easy-First Decoding in MT
 - Guided Learning in MT



General Algorithm for Online Learning

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

General Algorithm for Online Learning

$\mathbf{w} := 0$

```
while ( $\mathbf{x}, y_{gold}$ ) := get_new_data()  
     $y_{pred} := \text{prediction}(\mathbf{w}, \mathbf{x})$   
     $\mathbf{w} += \text{update}(\mathbf{x}, y_{gold}, y_{pred})$ 
```

Output: \mathbf{w}

initialize all weights to zero

General Algorithm for Online Learning

```
w := 0
while (x, ygold) := get_new_data()
    ypred := prediction(w, x)
    w += update(x, ygold, ypred)
```

Output: **w**

initialize all weights to zero
for each instance (observation)

General Algorithm for Online Learning

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

initialize all weights to zero
for each instance (observation)

1. get its features **x**

General Algorithm for Online Learning

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

initialize all weights to zero
for each instance (observation)

1. get its features **x**
2. do the prediction y_{pred}

General Algorithm for Online Learning

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

initialize all weights to zero
for each instance (observation)

1. get its features **x**
2. do the prediction y_{pred}
3. get the correct label y_{gold}

General Algorithm for Online Learning

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

initialize all weights to zero
for each instance (observation)

1. get its features **x**
2. do the prediction y_{pred}
3. get the correct label y_{gold}
4. update the weights

General Algorithm for Online Learning

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

initialize all weights to zero
for each instance (observation)

1. get its features **x**
2. do the prediction y_{pred}
3. get the correct label y_{gold}
4. update the weights

General Algorithm for Online Learning

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

initialize all weights to zero
for each instance (observation)

1. get its features **x**
2. do the prediction y_{pred}
3. get the correct label y_{gold}
4. update the weights

Definition: **conservative** online learning

no error \Rightarrow no update

i.e., if $y_{pred} = y_{gold}$ then $\text{update}(\mathbf{x}, y_{gold}, y_{pred}) = 0$

General Algorithm for Online Learning

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
    //prediction(w, x) =  $y_{gold}$ 
Output: w
```

initialize all weights to zero
for each instance (observation)

1. get its features \mathbf{x}
2. do the prediction y_{pred}
3. get the correct label y_{gold}
4. update the weights

Definition: conservative online learning

no error \Rightarrow no update

i.e., if $y_{pred} = y_{gold}$ then $\text{update}(\mathbf{x}, y_{gold}, y_{pred}) = 0$

Definition: aggressive online learning

after the update, the instance would be classified correctly

Perceptron

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

Binary Perceptron

$$\text{prediction}(\mathbf{w}, \mathbf{x}) \stackrel{\text{def}}{=} [\mathbf{w} \cdot \mathbf{x} > 0]$$

$$\text{update}(\mathbf{x}, y_{gold}, y_{pred}) \stackrel{\text{def}}{=} \alpha(y_{gold} - y_{pred}) \cdot \mathbf{x}$$

Perceptron

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
   $y_{pred}$  := prediction(w, x)
  w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

dot product (similarity score)
of weights and features

$$\mathbf{w} \cdot \mathbf{x} = \sum_i w_i x_i$$

Binary Perceptron

$$\text{prediction}(\mathbf{w}, \mathbf{x}) \stackrel{\text{def}}{=} [\mathbf{w} \cdot \mathbf{x} > 0]$$

$$\text{update}(\mathbf{x}, y_{gold}, y_{pred}) \stackrel{\text{def}}{=} \alpha(y_{gold} - y_{pred}) \cdot \mathbf{x}$$

Perceptron

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

dot product (similarity score)
of weights and features

$$\mathbf{w} \cdot \mathbf{x} = \sum_i w_i x_i$$

Iverson bracket

$$[P] = \begin{cases} 1 & \text{if } P \text{ is true;} \\ 0 & \text{otherwise.} \end{cases}$$

$$\text{prediction}(\mathbf{w}, \mathbf{x}) \stackrel{\text{def}}{=} [\mathbf{w} \cdot \mathbf{x} > 0]$$

$$\text{update}(\mathbf{x}, y_{gold}, y_{pred}) \stackrel{\text{def}}{=} \alpha (y_{gold} - y_{pred}) \cdot \mathbf{x}$$

Binary Perceptron

$$[\mathbf{w} \cdot \mathbf{x} > 0]$$

$$\alpha (y_{gold} - y_{pred}) \cdot \mathbf{x}$$

Perceptron

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

dot product (similarity score)
of weights and features

$$\mathbf{w} \cdot \mathbf{x} = \sum_i w_i x_i$$

Iverson bracket

$$[P] = \begin{cases} 1 & \text{if } P \text{ is true;} \\ 0 & \text{otherwise.} \end{cases}$$

Binary Perceptron

$$\text{prediction}(\mathbf{w}, \mathbf{x}) \stackrel{\text{def}}{=} [\mathbf{w} \cdot \mathbf{x} > 0]$$

$$\text{update}(\mathbf{x}, y_{gold}, y_{pred}) \stackrel{\text{def}}{=} \alpha (y_{gold} - y_{pred}) \cdot \mathbf{x}$$

learning rate (step size) $\alpha > 0$

Perceptron

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

prediction(**w**, **x**) $\stackrel{\text{def}}{=}$

Binary Perceptron

$[\mathbf{w} \cdot \mathbf{x} > 0]$

update(**x**, y_{gold} , y_{pred}) $\stackrel{\text{def}}{=}$

$\alpha(y_{gold} - y_{pred}) \cdot \mathbf{x}$

Multi-class Perceptron

$\arg \max_y \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, y)$

$\alpha(\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred}))$

learning rate (step size) $\alpha > 0$

Perceptron

```

w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
  
```

Output: **w**

Special case:

multi-prototype features

$$\mathbf{f}(\mathbf{x}, y) \stackrel{\text{def}}{=} \begin{aligned} & [y = \text{class}_1] \cdot \mathbf{x}, \\ & [y = \text{class}_2] \cdot \mathbf{x}, \\ & \dots \\ & [y = \text{class}_C] \cdot \mathbf{x} \end{aligned}$$

$$\text{prediction}(\mathbf{w}, \mathbf{x}) \stackrel{\text{def}}{=}$$

Binary Perceptron

$$[\mathbf{w} \cdot \mathbf{x} > 0]$$

$$\text{update}(\mathbf{x}, y_{gold}, y_{pred}) \stackrel{\text{def}}{=}$$

$$\alpha(y_{gold} - y_{pred}) \cdot \mathbf{x}$$

Multi-class Perceptron

$$\arg \max_y \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, y)$$

$$\alpha(\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred}))$$

$$\mathbf{w} := \mathbf{w} + \alpha \mathbf{f}(\mathbf{x}, y_{gold}) - \alpha \mathbf{f}(\mathbf{x}, y_{pred})$$

Perceptron

```
w := 0
while (x,  $y_{gold}$ ) := get_new_data()
     $y_{pred}$  := prediction(w, x)
    w += update(x,  $y_{gold}$ ,  $y_{pred}$ )
```

Output: **w**

General case:

any *label-dependent* features, e.g.

$f_{101}(\mathbf{x}, y) \stackrel{\text{def}}{=} [(y=\text{NNP} \text{ or } y=\text{NNPS})$
and \mathbf{x} capitalized]

prediction(**w**, **x**) $\stackrel{\text{def}}{=} [$

Binary Perceptron

$\mathbf{w} \cdot \mathbf{x} > 0]$

update(**x**, y_{gold} , y_{pred}) $\stackrel{\text{def}}{=} \alpha$

$(y_{gold} - y_{pred}) \cdot \mathbf{x}$

Multi-class Perceptron

$\arg \max_y \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, y)$

$\alpha (\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred}))$

Structured Prediction

- the number of possible labels is huge
- labels \mathbf{y} have a structure (graph, tree, sequence, . . .)
- usually can be decomposed (factorized) into subproblems
- local features
 - $f_i(\mathbf{x}, \mathbf{y}, j)$ can use whole \mathbf{x} , but only such y_k where k is “near” j
 - $f_{101}(\mathbf{x}, \mathbf{y}, j) \stackrel{\text{def}}{=} [(y_j = \text{NNP} \text{ or } y_j = \text{NNPS}) \text{ and word } x_j \text{ capitalized}]$
 - $f_{102}(\mathbf{x}, \mathbf{y}, j) \stackrel{\text{def}}{=} [y_j = \text{NNP} \text{ and } y_{j-1} = \text{NNP} \text{ and } |\mathbf{x}| \leq 6]$
- global features
 - $F_i(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \sum_j f_i(\mathbf{x}, \mathbf{y}, j)$
 - F_{101} . . . number of capitalized words with tag NNP or NNPS
 - F_{102} . . . number of NNP followed by NNP
or 0 if the sentence is longer than six words
 - We can define also features that cannot be decomposed

Structured Prediction using Online Learning

① local approach

- update after each local decision
- output of previous decisions used in local features
- e.g. Structured Perceptron (Collins, 2002)
- $y_{pred} = \arg \max_y \sum_i w_i f_i(\mathbf{x}, y_j, y_{j-1}, \dots)$

② global approach

- generate n-best list (lattice) of outputs \mathbf{y} for the whole \mathbf{x}
- compute global features, do update for each \mathbf{x} (sentence)
- we are re-ranking the n-best list
- e.g. MIRA (Crammer and Singer, 2003)
- $\mathbf{y}_{pred} = \arg \max_{\mathbf{y}} \sum_i w_i F_i(\mathbf{x}, \mathbf{y})$

Margin-based Online Learning

Definitions

- $score(y) = \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, y)$
- $margin(y) = score(y_{gold}) - score(y)$
 - $margin > 0 \Rightarrow$ no error
 - $|margin| \sim$ confidence
- $hinge_loss(y) = \max(0, 1 - margin(y))$

Online Prediction and Update

$$y_{pred} \stackrel{\text{def}}{=} \arg \max \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, y)$$
$$\mathbf{w} += \alpha (\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred}))$$

Margin-based Online Learning

Definitions

- $score(y) = \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, y)$
- $margin(y) = score(y_{gold}) - score(y)$
 - $margin > 0 \Rightarrow$ no error
 - $|margin| \sim$ confidence
- $hinge_loss(y) = \max(0, 1 - margin(y))$

Online Prediction and Update

$$y_{pred} \stackrel{\text{def}}{=} \arg \max \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, y)$$
$$\mathbf{w} += \alpha (\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred}))$$

Margin-based Online Learning

Definitions

- $score(y) = \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, y)$
- $margin(y) = score(y_{gold}) - score(y)$
 - $margin > 0 \Rightarrow$ no error
 - $|margin| \sim$ confidence
- $hinge_loss(y) = \max(0, 1 - margin(y))$

Online Prediction and Update

$$y_{pred} \stackrel{\text{def}}{=} \arg \max \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, y)$$
$$\mathbf{w} += \alpha (\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred}))$$

Margin-based Online Learning

Definitions

- $score(y) = \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, y)$
- $margin(y) = score(y_{gold}) - score(y)$
 - $margin > 0 \Rightarrow$ no error
 - $|margin| \sim$ confidence
- $hinge_loss(y) = \max(0, 1 - margin(y))$

Online Prediction and Update

$$y_{pred} \stackrel{\text{def}}{=} \arg \max \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, y)$$
$$\mathbf{w} += \alpha (\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred}))$$

Perceptron

$$\alpha_{\text{Perc}} \stackrel{\text{def}}{=} 1 \text{ (or any fixed value } > 0 \text{)}$$

Passive Aggressive (PA)

$$\alpha_{\text{PA}} \stackrel{\text{def}}{=} \frac{hinge_loss(y_{pred})}{\|\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred})\|^2}$$

Passive Aggressive I

$$\alpha_{\text{PA-I}} \stackrel{\text{def}}{=} \min \{ C, \alpha_{\text{PA}} \}$$

Passive Aggressive II

$$\alpha_{\text{PA-II}} \stackrel{\text{def}}{=} \frac{hinge_loss(y_{pred})}{\|\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred})\|^2 + \frac{1}{2C}}$$

Margin-based Online Learning

Definitions



Online Prediction and Update

$$y_{pred} \stackrel{\text{def}}{=} \arg \max_{\mathbf{w}} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, y)$$

$$\mathbf{w} += \alpha (\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred}))$$

Perceptron

$$\alpha_{\text{Perc}} \stackrel{\text{def}}{=} 1 \text{ (or any fixed value } > 0 \text{)}$$

Passive Aggressive (PA)

$$\alpha_{\text{PA}} \stackrel{\text{def}}{=} \frac{\text{hinge_loss}(y_{pred})}{\|\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred})\|^2}$$

Passive Aggressive I

$$\alpha_{\text{PA-I}} \stackrel{\text{def}}{=} \min \{ C, \alpha_{\text{PA}} \}$$

Passive Aggressive II

$$\alpha_{\text{PA-II}} \stackrel{\text{def}}{=} \frac{\text{hinge_loss}(y_{pred})}{\|\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred})\|^2 + \frac{1}{2C}}$$

Cost-sensitive Online Learning

Definitions

- $cost(y)$ = external error metric (non-negative)
e.g. 1 - similarity of y and y_{gold}
- $hinge_loss(y) = \max(0, cost(y) - margin(y))$

Hope and Fear

- $\mathbf{w} += \alpha(\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred}))$
- min-cost $y_{hope} \stackrel{\text{def}}{=} \arg \max_y -cost(y)$
- max-score $y_{fear} \stackrel{\text{def}}{=} \arg \max_y score(y)$
- cost-diminished $y_{hope} \stackrel{\text{def}}{=} \arg \max_y score(y) - cost(y)$
- cost-augmented $y_{fear} \stackrel{\text{def}}{=} \arg \max_y score(y) + cost(y)$
- max-cost $y_{fear} \stackrel{\text{def}}{=} \arg \max_y cost(y)$

Cost-sensitive Online Learning

Definitions

- $cost(y) =$ external error metric (non-negative)
e.g. 1 - similarity of y and y_{gold}
- $hinge_loss(y) = \max(0, cost(y) - margin(y))$

Hope and Fear

- $\mathbf{w} += \alpha(\mathbf{f}(\mathbf{x}, y_{gold}) - \mathbf{f}(\mathbf{x}, y_{pred}))$
- min-cost $y_{hope} \stackrel{\text{def}}{=} \arg \max_y -cost(y)$
- max-score $y_{fear} \stackrel{\text{def}}{=} \arg \max_y score(y)$
- cost-diminished $y_{hope} \stackrel{\text{def}}{=} \arg \max_y score(y) - cost(y)$
- cost-augmented $y_{fear} \stackrel{\text{def}}{=} \arg \max_y score(y) + cost(y)$
- max-cost $y_{fear} \stackrel{\text{def}}{=} \arg \max_y cost(y)$

Cost-sensitive Online Learning

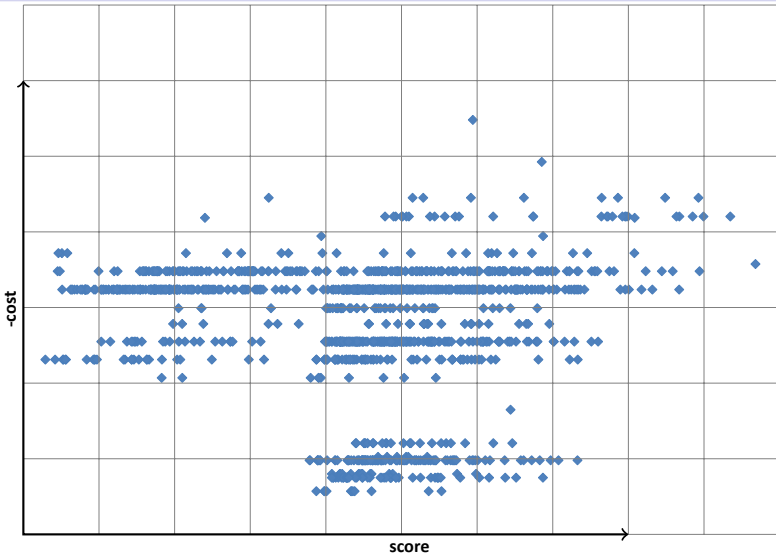
Definitions

- $cost(y) =$ external error metric (non-negative)
e.g. 1 - similarity of y and y_{gold}
- $hinge_loss(y) = \max(0, cost(y) - margin(y))$

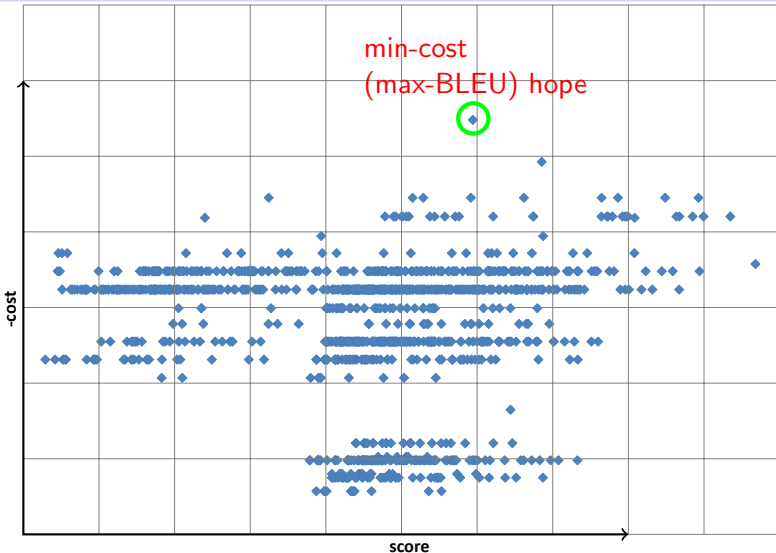
Hope and Fear

- $\mathbf{w} += \alpha(\mathbf{f}(\mathbf{x}, y_{hope}) - \mathbf{f}(\mathbf{x}, y_{fear}))$
- min-cost $y_{hope} \stackrel{\text{def}}{=} \arg \max_y -cost(y)$
- max-score $y_{fear} \stackrel{\text{def}}{=} \arg \max_y score(y)$
- cost-diminished $y_{hope} \stackrel{\text{def}}{=} \arg \max_y score(y) - cost(y)$
- cost-augmented $y_{fear} \stackrel{\text{def}}{=} \arg \max_y score(y) + cost(y)$
- max-cost $y_{fear} \stackrel{\text{def}}{=} \arg \max_y cost(y)$

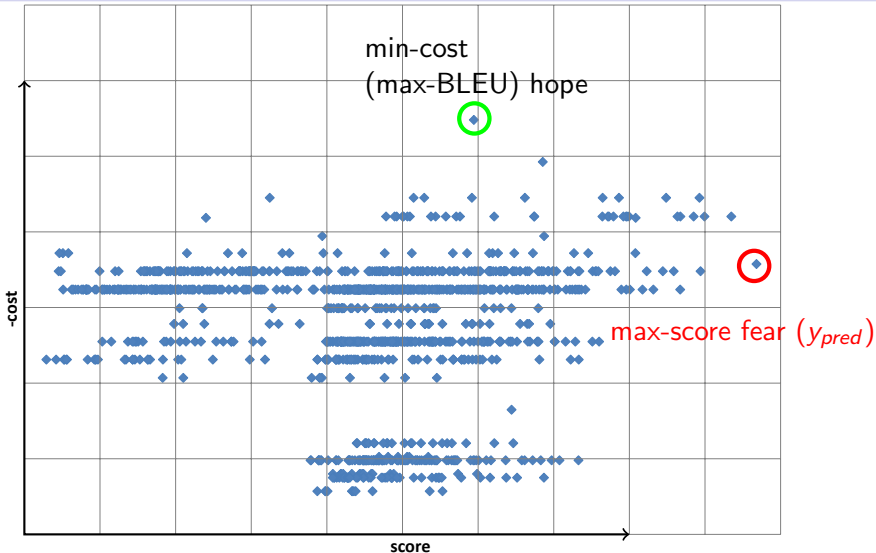
Cost-sensitive Online Learning



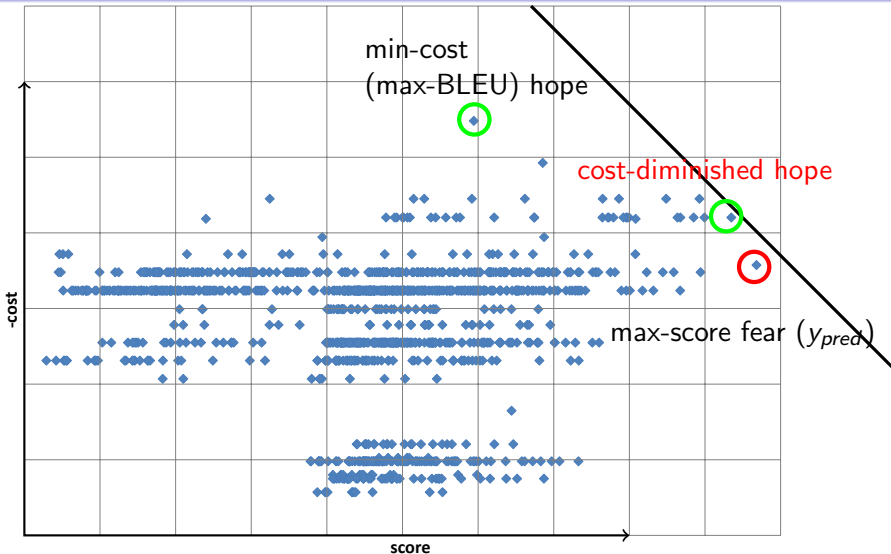
Cost-sensitive Online Learning



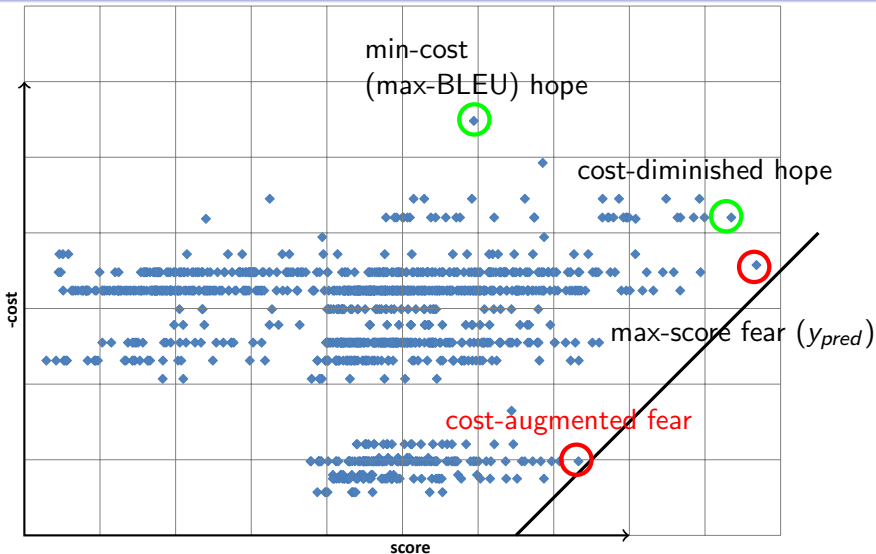
Cost-sensitive Online Learning



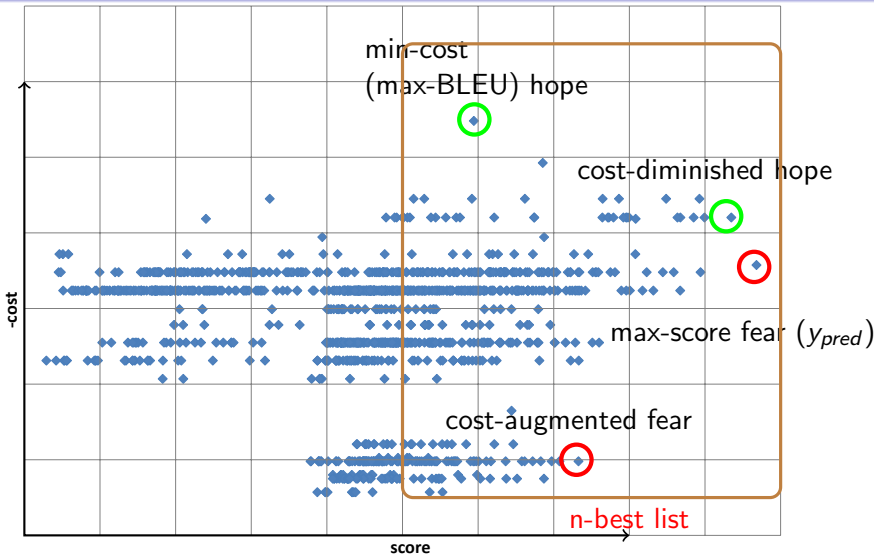
Cost-sensitive Online Learning



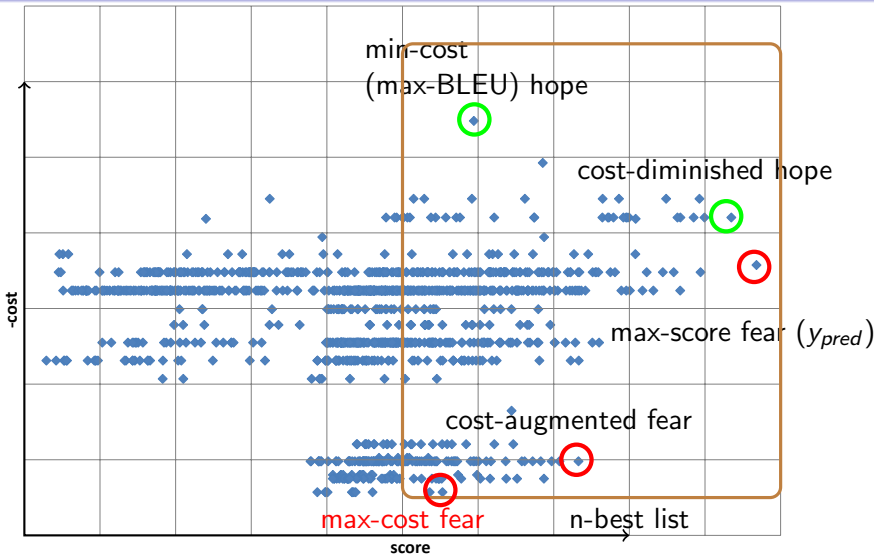
Cost-sensitive Online Learning



Cost-sensitive Online Learning



Cost-sensitive Online Learning



Application to MT

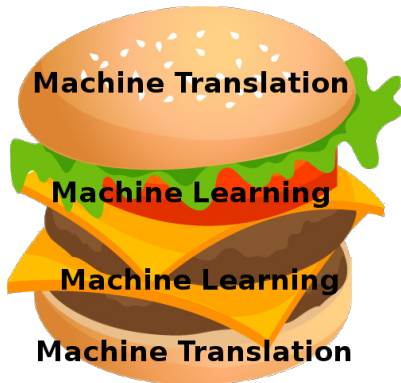
x = source sentence

y_{gold} = its reference translation

- more references sometimes available
- reference may be *unreachable*
- we score *derivations* (which include latent variables)
one translation may have more derivations

Today's Menu

- 1 MT Intro
 - Taxonomy
 - Hybrids
- 2 Online Learning
 - Perceptron
 - Structured Prediction
- 3 Guided Learning
- 4 Back to MT
 - Easy-First Decoding in MT
 - Guided Learning in MT



Easy-First Decoding (PoS Tagging)

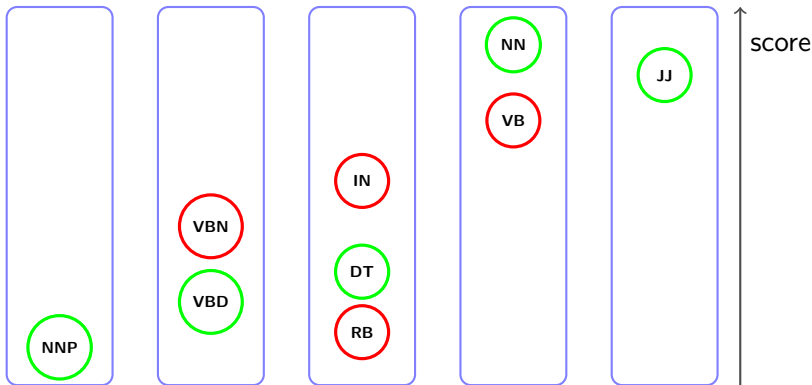
Agatha

found

that

book

interesting



Easy-First Decoding (PoS Tagging)

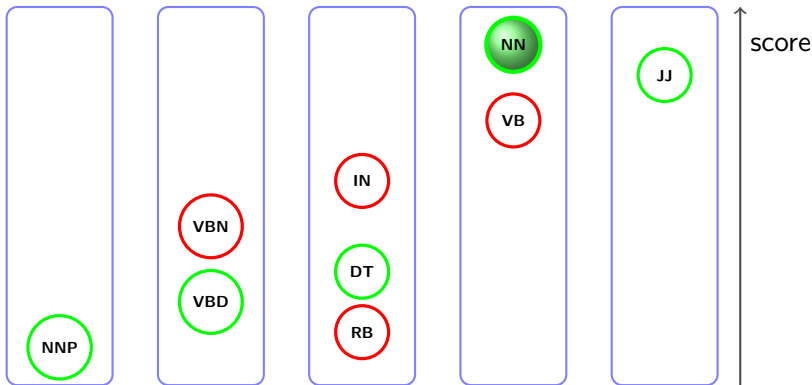
Agatha

found

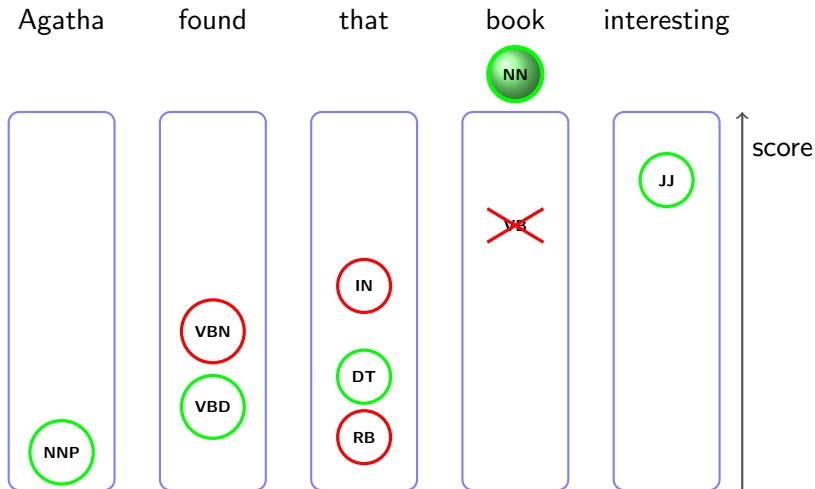
that

book

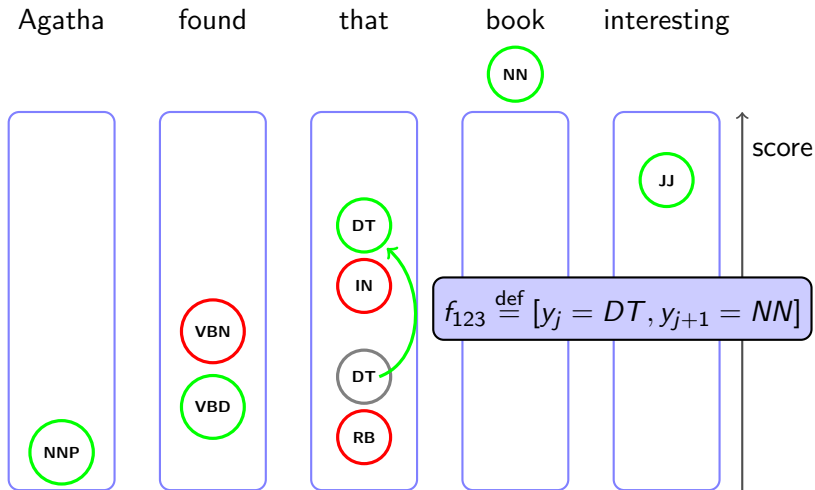
interesting



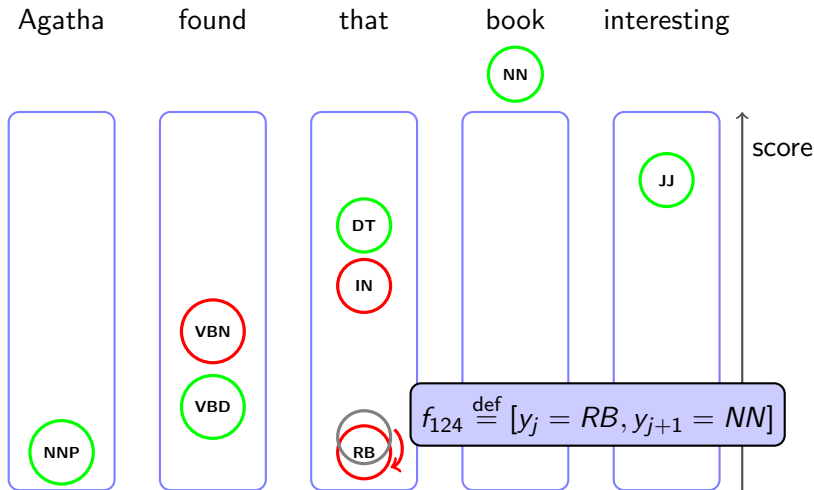
Easy-First Decoding (PoS Tagging)



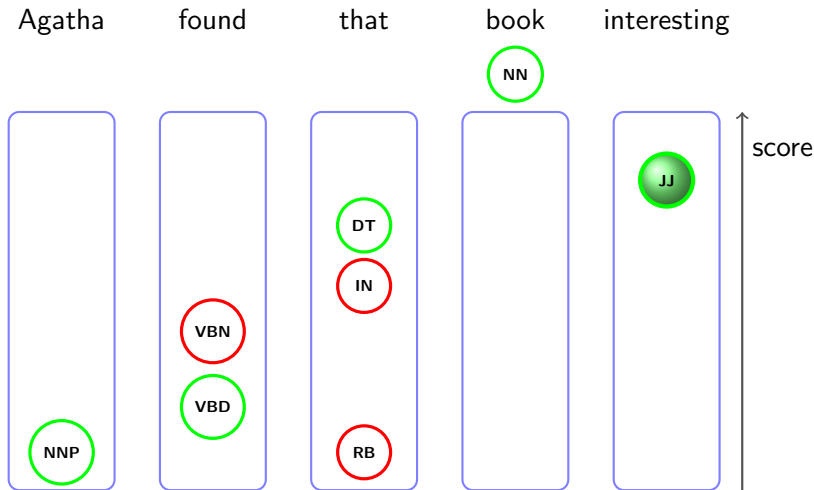
Easy-First Decoding (PoS Tagging)



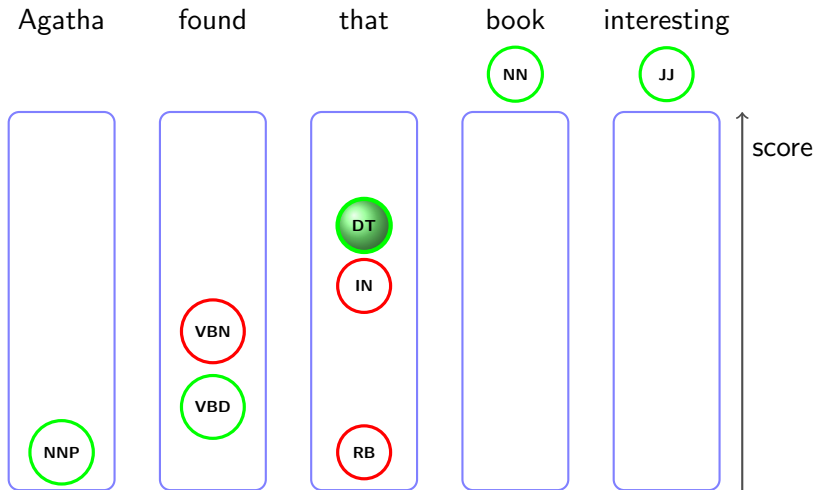
Easy-First Decoding (PoS Tagging)



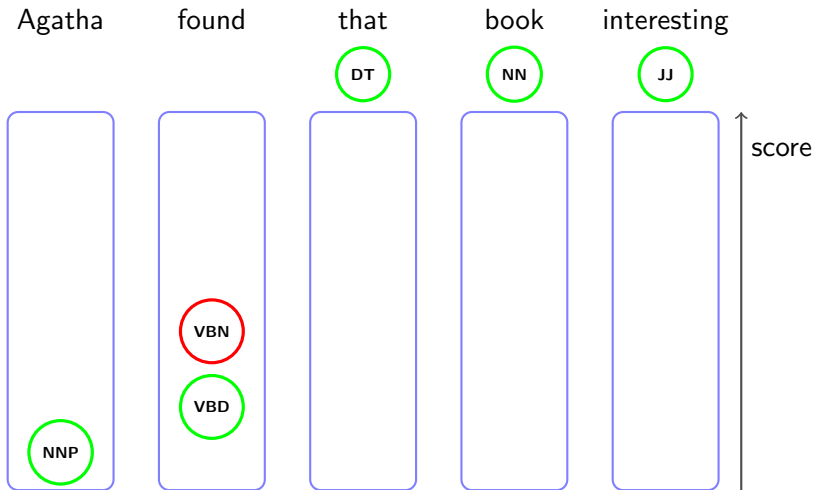
Easy-First Decoding (PoS Tagging)



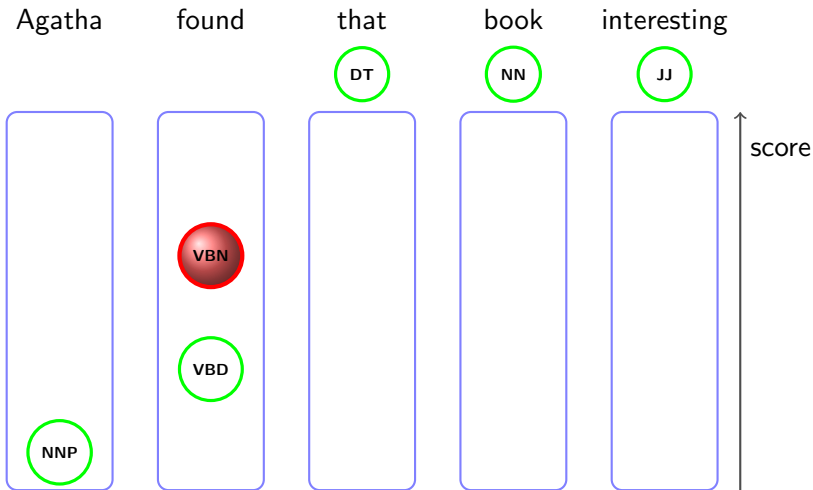
Easy-First Decoding (PoS Tagging)



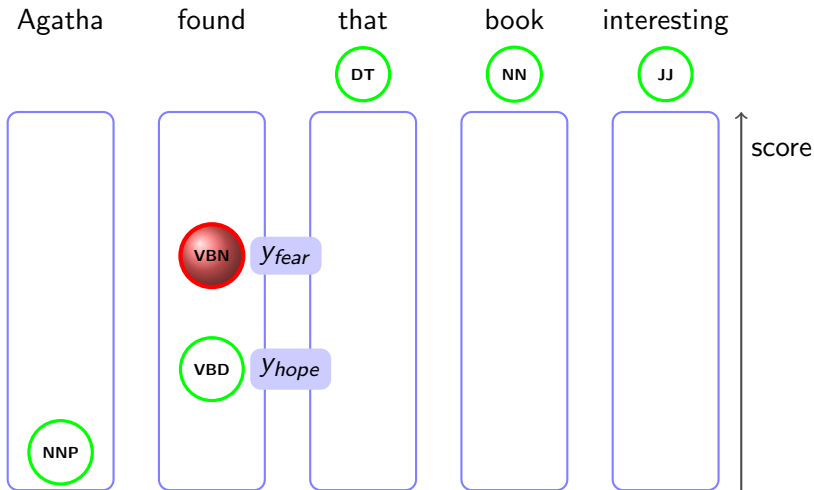
Easy-First Decoding (PoS Tagging)



Easy-First Decoding (PoS Tagging)

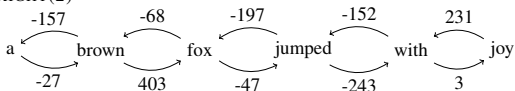


Guided Learning (PoS Tagging)



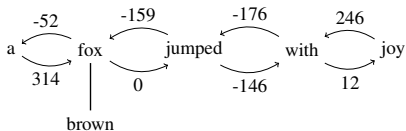
Easy-First Decoding (Dependency Parsing)

(1) ATTACHRIGHT(2)



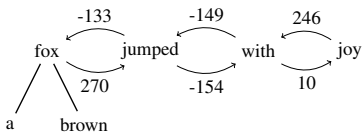
Easy-First Decoding (Dependency Parsing)

(2) ATTACHRIGHT(1)



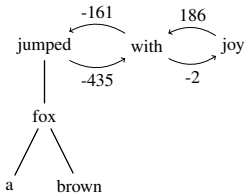
Easy-First Decoding (Dependency Parsing)

(3) ATTACHRIGHT(1)

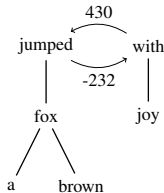


Easy-First Decoding (Dependency Parsing)

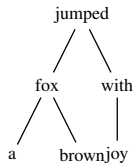
(4) ATTACHLEFT(2)



(5) ATTACHLEFT(1)

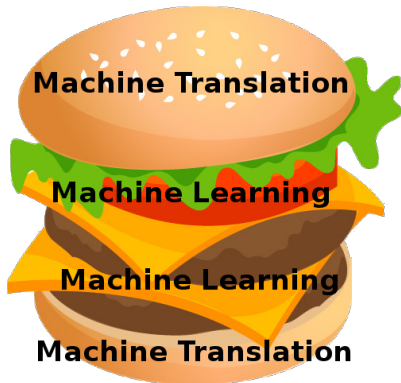


(6)

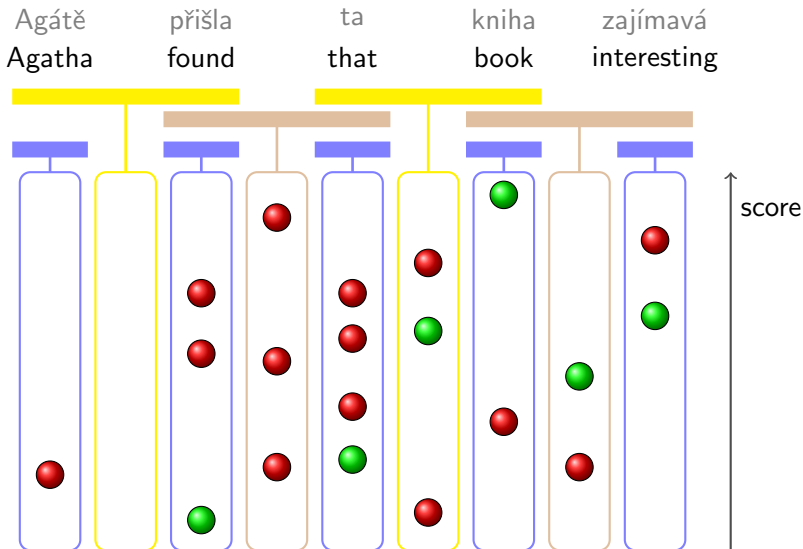


Today's Menu

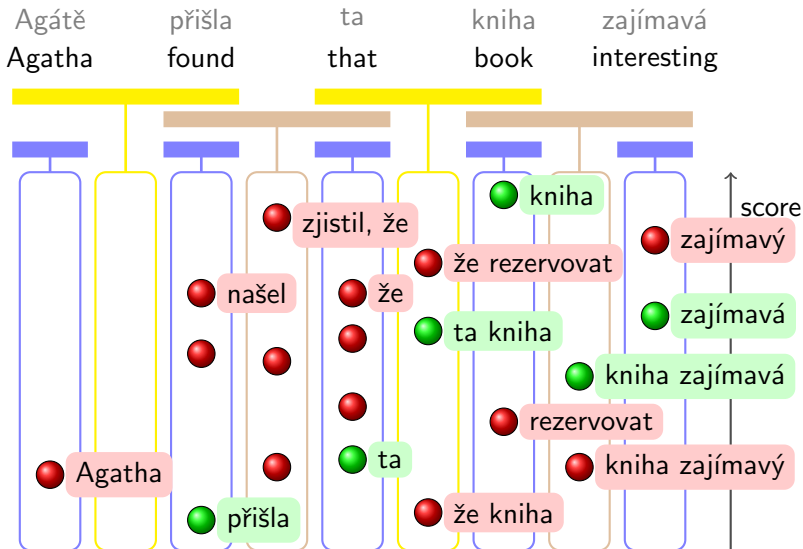
- 1 MT Intro
 - Taxonomy
 - Hybrids
- 2 Online Learning
 - Perceptron
 - Structured Prediction
- 3 Guided Learning
- 4 Back to MT
 - Easy-First Decoding in MT
 - Guided Learning in MT



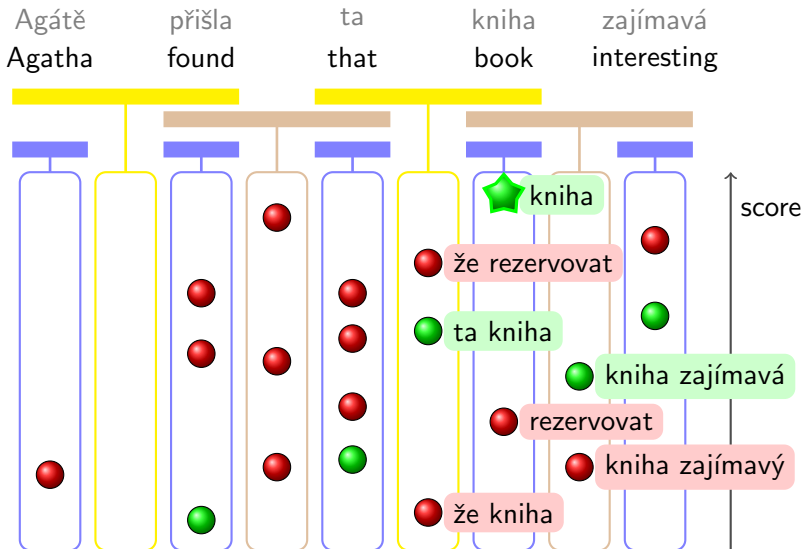
Easy-First Decoding (Phrase-Based MT)



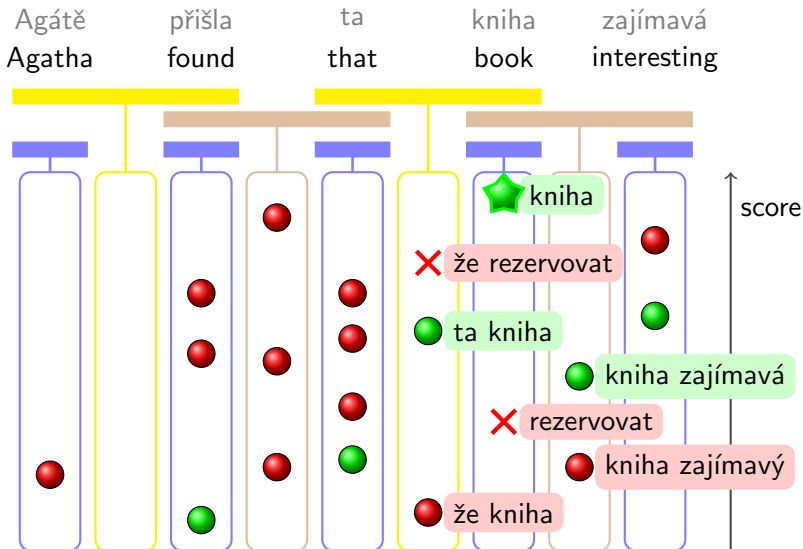
Easy-First Decoding (Phrase-Based MT)



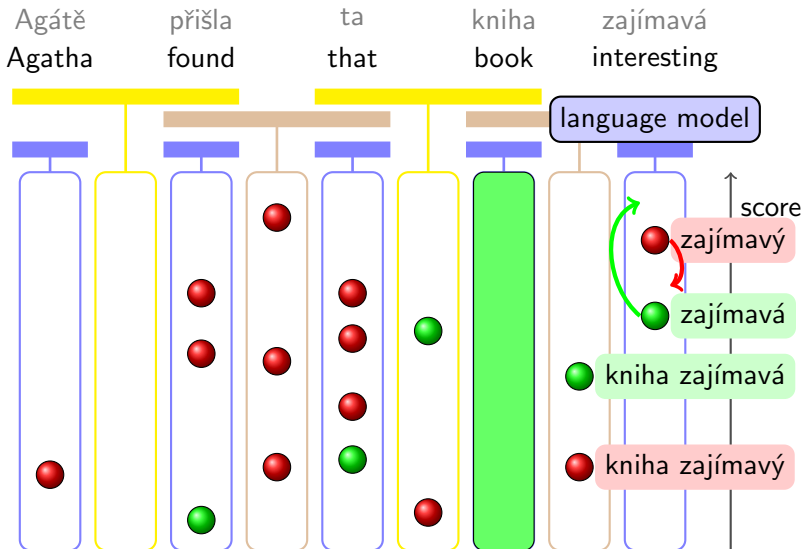
Easy-First Decoding (Phrase-Based MT)



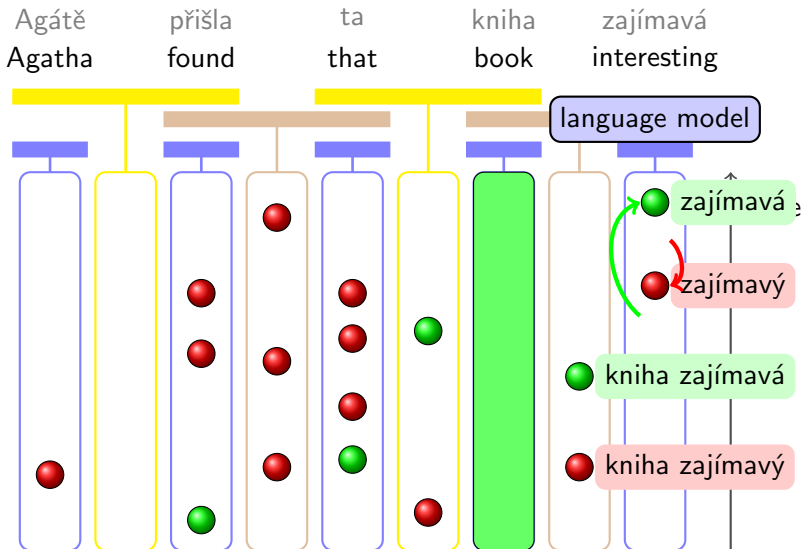
Easy-First Decoding (Phrase-Based MT)



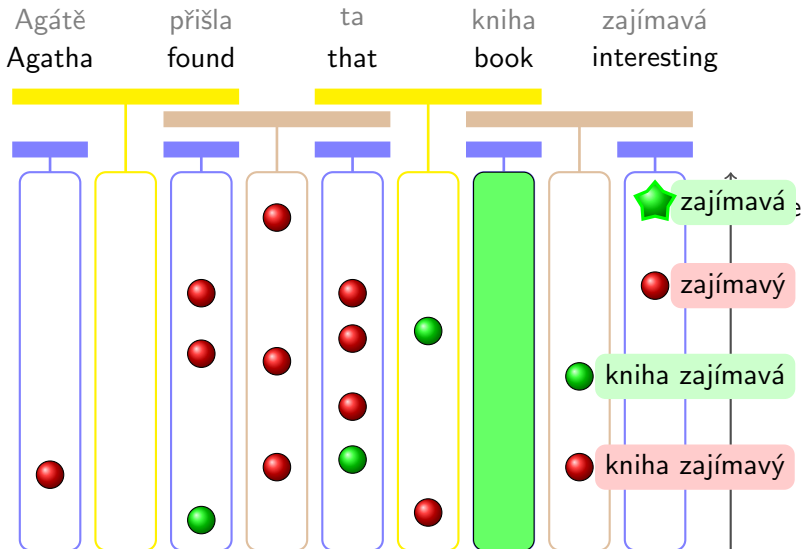
Easy-First Decoding (Phrase-Based MT)



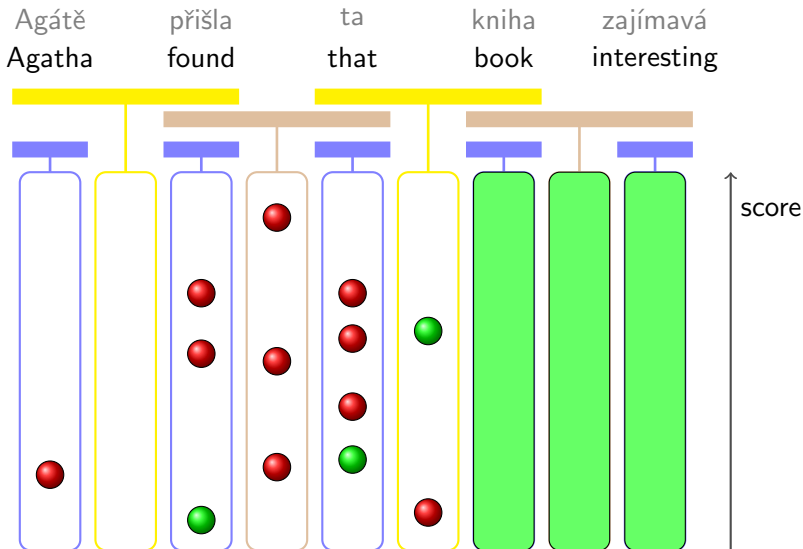
Easy-First Decoding (Phrase-Based MT)



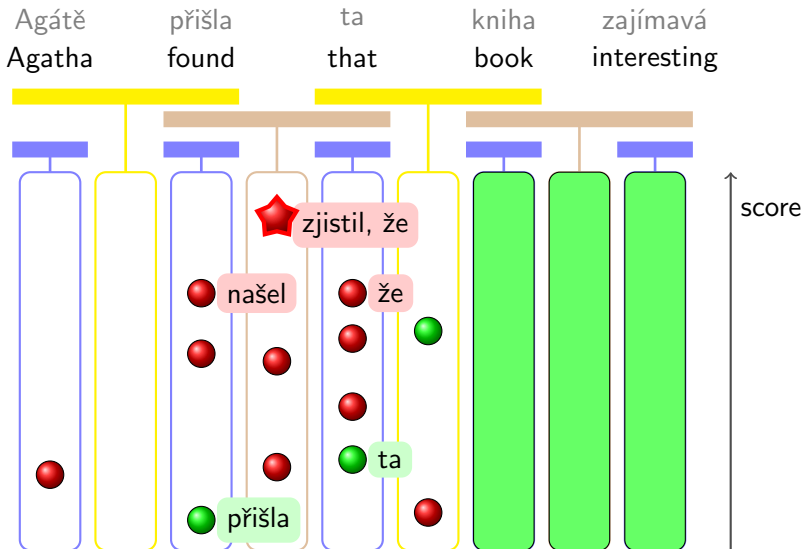
Easy-First Decoding (Phrase-Based MT)



Easy-First Decoding (Phrase-Based MT)



Easy-First Decoding (Phrase-Based MT)



Features for Guided Learning in MT

Source Segment Features

- segment size (number of words)
- entropy $P(\text{target}|\text{source}) = - \sum_i P(\text{src}, \text{trg}_i) \cdot \log P(\text{trg}_i|\text{src})$
- log count(source)
- source language model: $\log P(\text{source})$
- word identity, e.g. $f_{42} \stackrel{\text{def}}{=} [\text{src}=\text{found that}]$
- PoS identity, e.g. $f_{43} \stackrel{\text{def}}{=} [\text{src_pos}=\text{VBD IN}]$

Target-dependent Features

- $\log P(\text{trg}|\text{src})$
- target language model: $\log P(\text{target} | \text{previous segment})$
- log count(target)?
- identity, e.g. $f_{142} \stackrel{\text{def}}{=} [\text{src}=\text{found that} \ \& \ \text{trg}=\text{zjistil}]$

Features for Guided Learning in MT

Source Segment Features

- **segment size (number of words)**
- entropy $P(\text{target}|\text{source}) = - \sum_i P(\text{src}, \text{trg}_i) \cdot \log P(\text{trg}_i|\text{src})$
- log count(source)
- source language
- word identity, e.g.
- PoS identity, e.g. etc.

Combinations and Quantizations

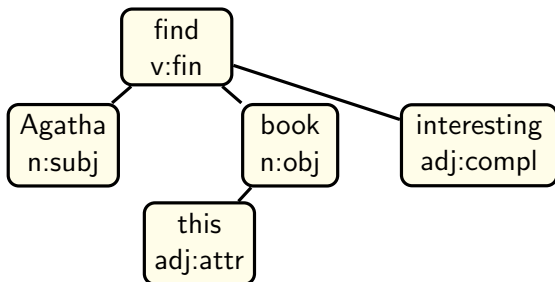
$$[\text{size}(\text{src}) = 3] \cdot \log P(\text{trg}|\text{src})$$

$$[\text{size}(\text{src}) = 3 \quad \& \quad -3 < \log P(\text{trg}|\text{src}) < -2]$$

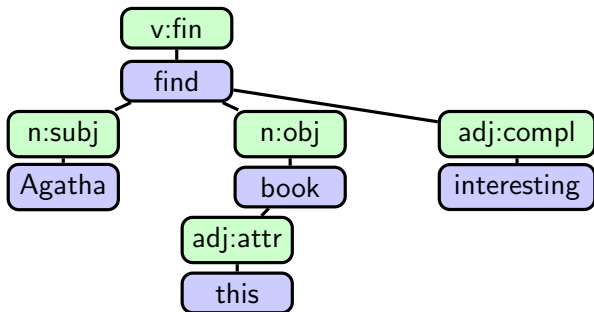
Target-dependent Features

- **log $P(\text{trg}|\text{src})$**
- target language model: log $P(\text{target} | \text{previous segment})$
- log count(target)?
- identity, e.g. $f_{142} \stackrel{\text{def}}{=} [\text{src}=\text{found that} \ \& \ \text{trg}=\text{zjistil}]$

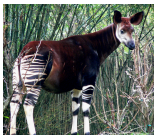
Application to Tecto Trees



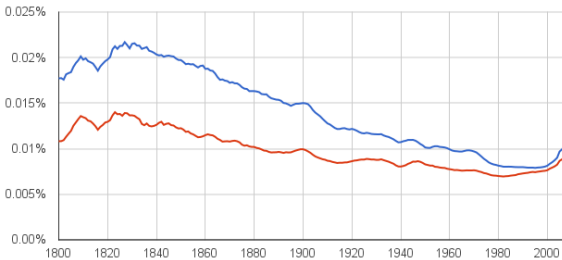
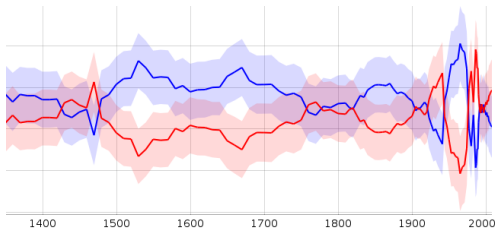
Application to Tecto Trees



What have you seen in the Zoo

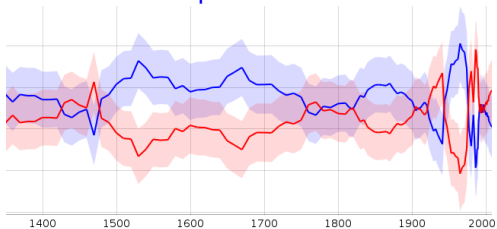


Predictions?

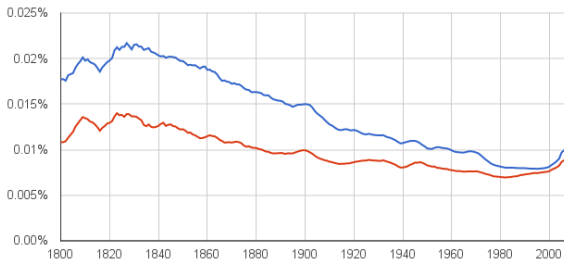


Predictions?

Hope and Fear



<http://syd.korpus.cz/>



<http://books.google.com/ngrams/>