

Online Learning for CAT applications

Nicolò Cesa-Bianchi Gabriele Reverberi

Università degli Studi di Milano



S M A R T

Statistical Multilingual Analysis
for Retrieval and Translation

Interactive Machine Translation

- Learning/optimization techniques are used to tune the parameters of SMT systems
- Online learning adjusts parameters **incrementally**
[Lian et al., 2006; Arun and Koehn, 2007; Tillman and Zhang, 2008]
- Especially useful when the system **interacts** with the user



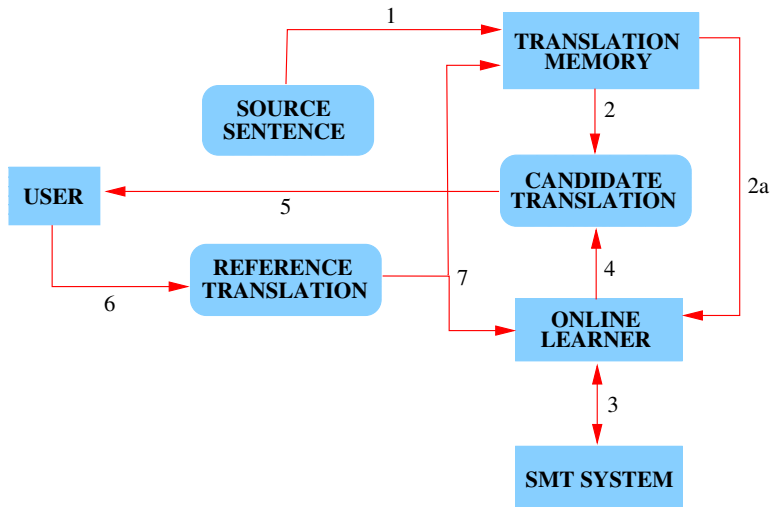
Computer Assisted Translation (CAT)

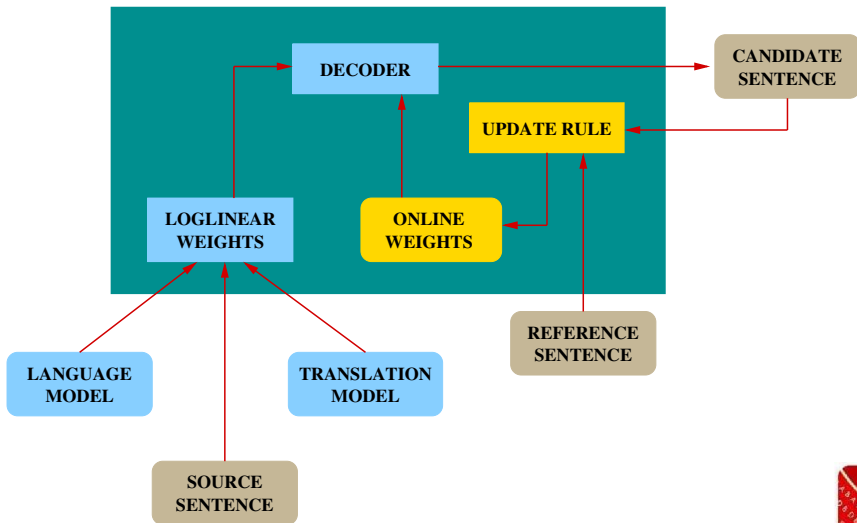
The image displays two screenshots of Computer Assisted Translation (CAT) software. The top screenshot shows the 'SMART Translation' window. It features a 'Workbench' menu and a toolbar. The main text area contains an English source text: 'The purpose of this report is to establish the scenarios which will be evaluated in the three case studies within the SMART project and to detail the requirements of the case studies towards the technical work packages (both in terms of required functionality and integration related issues)'. Below this is a yellow-highlighted Slovenian target text: 'Prevajalni spomin je sestavljen iz delov besedila iz izvornega jezika in njihovih prevodov v enega ali več ciljnih jezikov'. Further down, there is a green bar indicating a '100% match from Translation Memory'. Below this, the 'Source' and 'Target' fields are populated with the respective text segments.

The bottom screenshot shows the 'SDL Trados Translator's Workbench' window, configured for 'project - English (United Kingdom) -> Slovenian'. It displays the same source text in the top pane and the translated target text in the bottom pane. The target text is: 'Prevajalni spomin je sestavljen iz delov besedila iz izvornega jezika in njihovih prevodov v enega ali več ciljnih jezikov.'. At the bottom of the window, a status bar shows a match percentage of '100%' and 'Exact Match'.



CAT meets online learning





Experimental setup (based on Portage SMT system)

Feature set for online weights

A new feature is created for **each phrasetable entry**



Experimental setup (based on Portage SMT system)

Feature set for online weights

A new feature is created for **each phrasetable entry**

Phase 1 – **offline mode**

- Building of phrasetable on a training corpus
- Tuning of loglinear weights on a development corpus

→ This gives the **baseline** system



Experimental setup (based on Portage SMT system)

Feature set for online weights

A new feature is created for **each phrasetable entry**

Phase 1 – **offline mode**

- Building of phrasetable on a training corpus
- Tuning of loglinear weights on a development corpus

→ This gives the **baseline** system

Phase 2 – **online mode**

Online weights are adapted during CAT process



Adaptive decoding — basic definitions

- $f(x_t, y)$ is the vector of **phrasetable feature values** when considering y as candidate translation for the source sentence x_t
- The vector w contains the decoder **online weights**
- The decoder builds a N -best list Y_t of candidate translations y by ranking them according to **margin**

$$w^T f(x_t, y)$$



Adaptive decoding — basic definitions

- $f(x_t, y)$ is the vector of **phrasetable feature values** when considering y as candidate translation for the source sentence x_t
- The vector w contains the decoder **online weights**
- The decoder builds a **N-best** list Y_t of candidate translations y by ranking them according to **margin**

$$w^T f(x_t, y)$$

- The **1-best** translation is $\hat{y}_t = \operatorname{argmax}_{y \in Y_t} w^T f(x_t, y)$
- The **pseudo-target translation** is $y_t^* = \operatorname{argmax}_{y \in Y_t} \text{BLEU}(y_t, y)$



Adaptive decoding — basic algorithmic framework

Recall:

Decoder ranks translations \mathbf{y} according to $\mathbf{w}^\top \mathbf{f}(x_t, \mathbf{y})$

- **Margin difference** for weight \mathbf{w} when \mathbf{y} is chosen instead of \mathbf{y}^*

$$\text{MARGIN}_t(\mathbf{y}^*, \mathbf{y}) = \mathbf{w}^\top (\mathbf{f}(x_t, \mathbf{y}^*) - \mathbf{f}(x_t, \mathbf{y}))$$

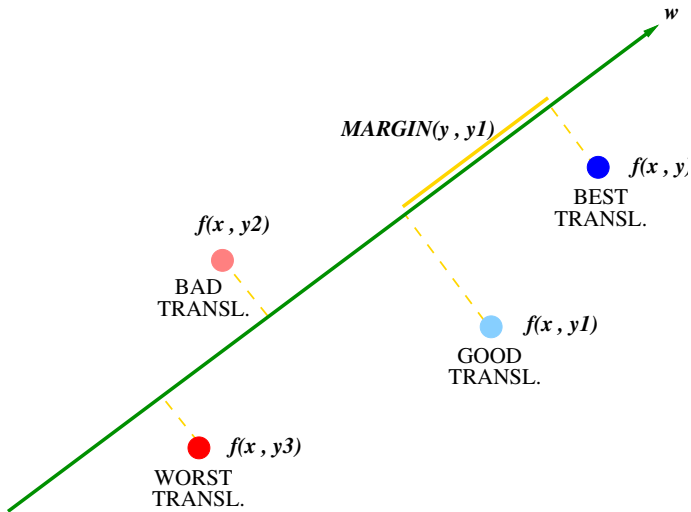
- **Linear constraints** the learner tries to enforce at each step t

$$\text{MARGIN}_t(\mathbf{y}^*, \mathbf{y}) \geq \text{BLEU}(\mathbf{y}_t, \mathbf{y}_t^*) - \text{BLEU}(\mathbf{y}_t, \mathbf{y}) \quad \forall \mathbf{y} \in Y_t$$

- Constraints are approximately enforced by projecting current \mathbf{w} onto (some of the) hyperplanes defined by constraints



Cost-sensitive margin condition



Update of parameters

Recall:

\mathbf{y} = reference translation

\mathbf{y}^* = pseudo-target translation (highest BLEU in N-best)

$\hat{\mathbf{y}}$ = guessed translation (1-best)

\mathbf{w} = current value of online weights

Enforce margin difference between pseudo-target \mathbf{y}^* and 1-best $\hat{\mathbf{y}}$

$$\min_{\mathbf{w}', \xi} \|\mathbf{w} - \mathbf{w}'\|^2 + C \xi$$

such that $\text{MARGIN}(\mathbf{y}^*, \hat{\mathbf{y}}) \geq (\text{BLEU}(\mathbf{y}, \mathbf{y}^*) - \text{BLEU}(\mathbf{y}, \hat{\mathbf{y}})) - \xi$

Passive-aggressive update

[Crammer et al., 2006]

$$\mathbf{w} \leftarrow \mathbf{w} + \eta_t (\text{BLEU}(\mathbf{y}_t, \mathbf{y}_t^*) - \text{BLEU}(\mathbf{y}_t, \hat{\mathbf{y}}_t))$$

Theoretical guarantees

For any sequence $(x_1, y_1), (x_2, y_2), \dots$ of source/reference pairs

- If there exists choice \mathbf{u} for the parameters that satisfies all constraints at each step, then

$$\sum_t \text{BLEU}(y_t, \hat{y}_t) \geq \sum_t \text{BLEU}(y_t, y_t^*) - \|\mathbf{u}\|^2$$

- If no such \mathbf{u} exists, then $\sum_t \text{BLEU}(y_t, \hat{y}_t)$ is at least

$$\sum_t \text{BLEU}(y_t, y_t^*) - \inf_{\mathbf{u}} \left(1 + \frac{1}{C} \right) \left(\|\mathbf{u}\|^2 + C \sum_t H_t(\mathbf{u}) \right)$$

- C is the aggressiveness parameter associated with the constraints
- $H_t(\mathbf{u})$ measures by how much the margin of \mathbf{u} fails the **worst constraint** at time t

Performance measure

- Learning algorithms and their analysis do not require BLEU
- For robustness reasons, we train and test the system using **BLEUMIX**, an average of different sentence-level measures
($1 \text{ BLEUMIX} \approx 0.65 \text{ BLEU}$)



Performance measure

- Learning algorithms and their analysis do not require BLEU
- For robustness reasons, we train and test the system using **BLEUMIX**, an average of different sentence-level measures
($1 \text{ BLEUMIX} \approx 0.65 \text{ BLEU}$)

Cumulative BLEUMIX difference

The cumulative difference in sentence-level BLEUMIX points between **online system translations** \hat{y}_t and **Portage baseline translations** y'_t with respect to the common reference translation y_t

$$\sum_{t=1}^T \left(\text{BLEUMIX}(y_t, \hat{y}_t) - \text{BLEUMIX}(y_t, y'_t) \right)$$



Experimental setup

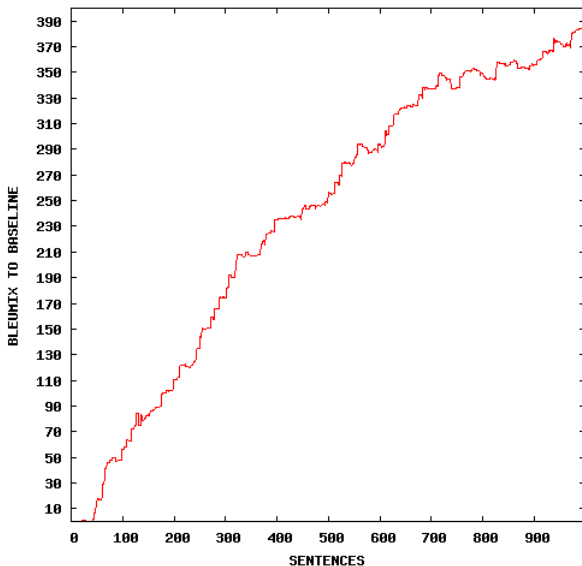
- **Corpus:** English → Spanish section of Europarl
- **Training set:** 165,000 sentences
- **Dev set:** (used to tune Portage) 6,000 sentences
- **Test set:** (used for online learning)
Five adjacent nonoverlapping blocks of 1,000 sentences each



- Online learner attempts to improve on **tuned Portage** performance by a *single run* over **1,000** sentences
 - less than **0.6%** of Portage training set!
- Learner does so by simultaneously tuning **1,7M** parameters associated with the phrasetable entries
 - about **1,700** parameters per observed sentence!
- We get an improvement of about **0.4 BLEUMIX** points **per observed sentence**




Weight adaptation



Oracular phrasetable adaptation


Dynamic growth of phrasetable 

Problem: on-the-fly alignment of new segments 



Oracular phrasetable adaptation

Dynamic growth of phrasetable

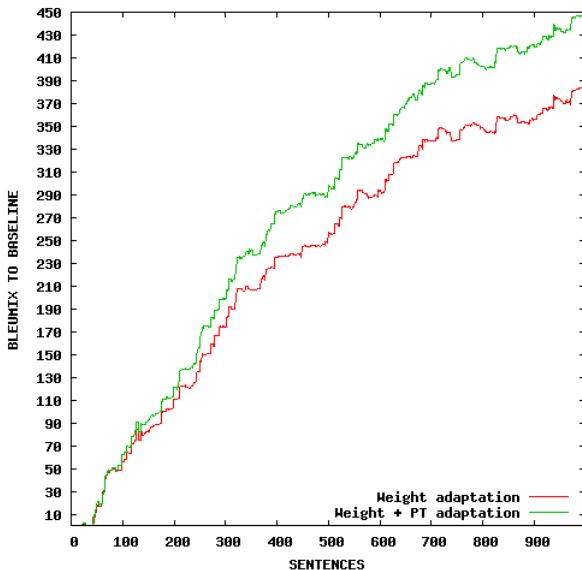
Problem: on-the-fly alignment of new segments 

Oracular PT

- Fake alignment by building an oracular PT on train + test corpora
- After translating each new sentence, the relevant segments are moved from the oracular PT to the working PT
- The weights associated with new segments are incrementally learned



Weight adaptation + PT adaptation

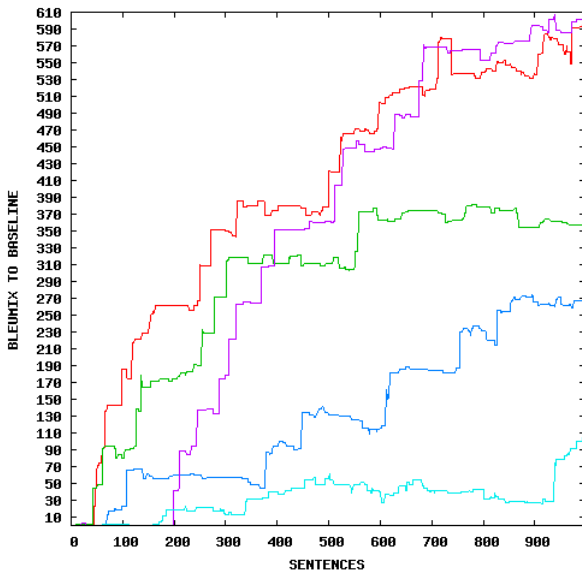


- We estimate the probability p that the performance difference **increases** when each translation in turn is obtained from a random system (adaptive or baseline)
- This is a p -value for the null hypothesis that baseline and adaptive have the same performance

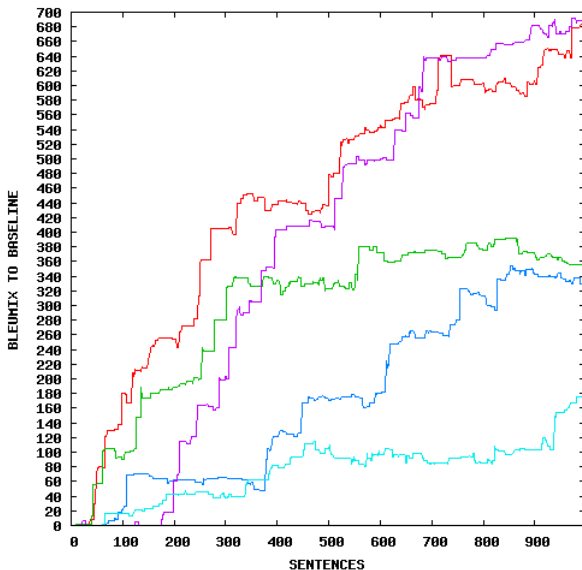
p-values				
0.01	0.28	0.33	0.18	0.45
0.01	0.40	0.20	0.13	0.41



Weight adaptation — 5 runs



Weight adaptation + PT adaptation — 5 runs



- More stable learning curves
- On-the-fly alignment to replace oracular PTT
- TM's crippling effect

