

End-to-End Discriminative Training for Machine Translation

Percy Liang
Dan Klein

Alex Bouchard-Côté
Ben Taskar



UC Berkeley

Computer Science Division

Discriminative machine translation

x: le parlement adopte la résolution législative



y: parliament adopted the legislative resolution

Discriminative machine translation

x: le parlement adopte la résolution législative



y: parliament adopted the legislative resolution

Training examples

x₁: mr barn cresco wishes to make a comment .

y₁: m. barn cresco veut intervenir .

x₂: i consider it very important to continue this work .

y₂: à mon avis , il est très important de continuer dans cette voie-là .

x₃: this all augurs well , and represents real progress .

y₃: tout cela est de bon augure et constitue un réel progrès .

• • •

Discriminative machine translation

x: le parlement adopte la résolution législative



y: parliament adopted the legislative resolution

Training examples

Features

x₁: mr barn crespo wishes to make a comment .

y₁: m. barn crespo veut intervenir .

x₂: i consider it very important to continue this work .

y₂: à mon avis , il est très important de continuer dans cette voie-là .

x₃: this all augurs well , and represents real progress .

y₃: tout cela est de bon augure et constitue un réel progrès .

• • •

Discriminative machine translation

x: le parlement adopte la résolution législative



y: parliament adopted the legislative resolution

Training examples

x₁: mr barn crespo wishes to make a comment .

y₁: m. barn crespo veut intervenir .

x₂: i consider it very important to continue this work .

y₂: à mon avis , il est très important de continuer dans cette voie-là .

x₃: this all augurs well , and represents real progress .

y₃: tout cela est de bon augure et constitue un réel progrès .

• • •

Features

NN JJ → JJ NN 1

Discriminative machine translation

x: le parlement adopte la résolution législative



y: parliament adopted the legislative resolution

Training examples

x₁: mr barn crespo wishes to make a comment .

y₁: m. barn crespo veut intervenir .

x₂: i consider it very important to continue this work .

y₂: à mon avis , il est très important de continuer dans cette voie-là .

x₃: this all augurs well , and represents real progress .

y₃: tout cela est de bon augure et constitue un réel progrès .

• • •

Features

NN JJ → JJ NN 1

Lang. model -4.253

Discriminative machine translation

x: le parlement adopte la résolution législative



y: parliament adopted the legislative resolution

Training examples

*x*₁: mr barn crespo wishes to make a comment .

*y*₁: m. barn crespo veut intervenir .

*x*₂: i consider it very important to continue this work .

*y*₂: à mon avis , il est très important de continuer dans cette voie-là .

*x*₃: this all augurs well , and represents real progress .

*y*₃: tout cela est de bon augure et constitue un réel progrès .

• • •

Features

NN JJ → JJ NN 1

Lang. model -4.253

• • •

Discriminative machine translation

x: le parlement adopte la résolution législative



y: parliament adopted the legislative resolution

Training examples

*x*₁: mr barn crespo wishes to make a comment .

*y*₁: m. barn crespo veut intervenir .

*x*₂: i consider it very important to continue this work .

*y*₂: à mon avis , il est très important de continuer dans cette voie-là .

*x*₃: this all augurs well , and represents real progress .

*y*₃: tout cela est de bon augure et constitue un réel progrès .

• • •

Features

NN JJ → JJ NN

1

w
2.3

Lang. model

-4.253

1.7

• • •

Discriminative machine translation

x: le parlement adopte la résolution législative



y: parliament adopted the legislative resolution

Training examples

*x*₁: mr barn crespo wishes to make a comment .

*y*₁: m. barn crespo veut intervenir .

*x*₂: i consider it very important to continue this work .

*y*₂: à mon avis , il est très important de continuer dans cette voie-là .

*x*₃: this all augurs well , and represents real progress .

*y*₃: tout cela est de bon augure et constitue un réel progrès .

...

Features

NN JJ → JJ NN

1

w
2.3

Lang. model

-4.253

1.7

...

That's it.

Discriminative machine translation

x: le parlement adopte la résolution législative



y: parliament adopted the legislative resolution

Training examples

*x*₁: mr barn crespo wishes to make a comment .

*y*₁: m. barn crespo veut intervenir .

*x*₂: i consider it very important to continue this work .

*y*₂: à mon avis , il est très important de continuer dans cette voie-là .

*x*₃: this all augurs well , and represents real progress .

*y*₃: tout cela est de bon augure et constitue un réel progrès .

...

Features

NN JJ → JJ NN

1

Lang. model

-4.253

...

w

2.3

1.7

That's it. Well, actually...

Why is this hard?

x: le parlement adopte la résolution législative

y: parliament adopted the legislative resolution

Why is this hard?

x: le parlement adopte la résolution législative

y: parliament adopted the legislative resolution

- Correct output is ill-defined

Why is this hard?

x: le parlement adopte la résolution législative

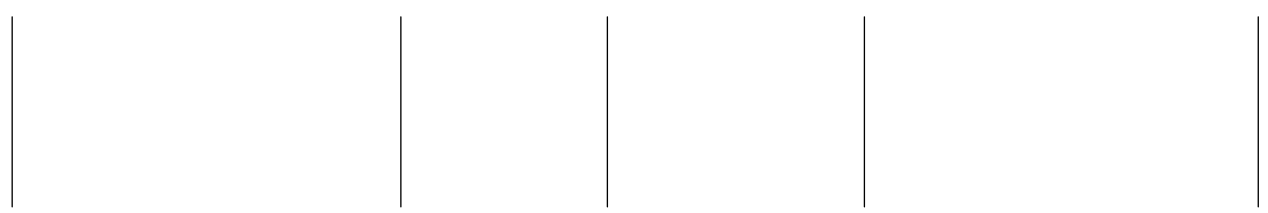
y: parliament adopted the legislative resolution

- Correct output is ill-defined
- Correspondences are missing

Why is this hard?

x: le parlement adopte la résolution législative

y: DT NN VBD DT NN JJ



- Correct output is ill-defined
- Correspondences are missing

Why is this hard?

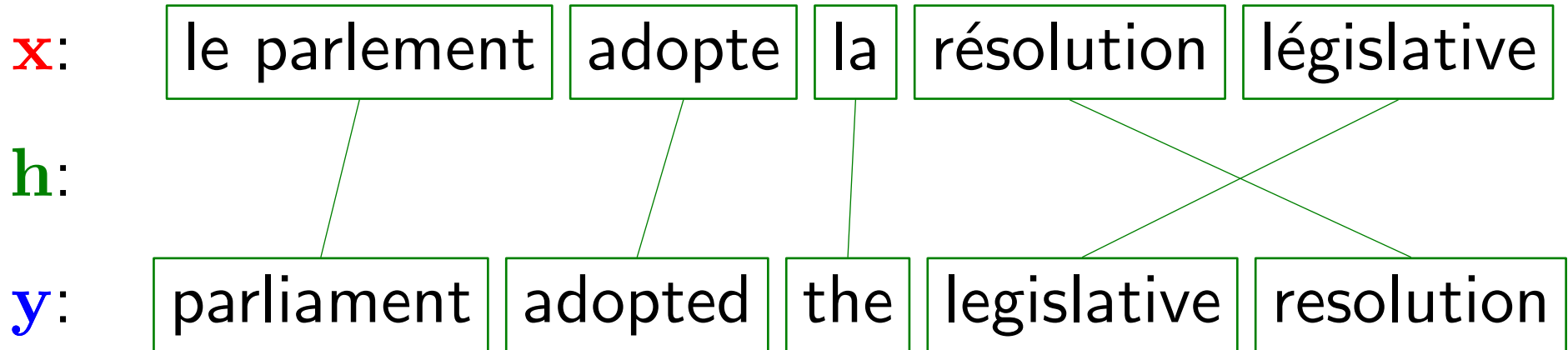
x: le parlement adopte la résolution législative

???

y: parliament adopted the legislative resolution

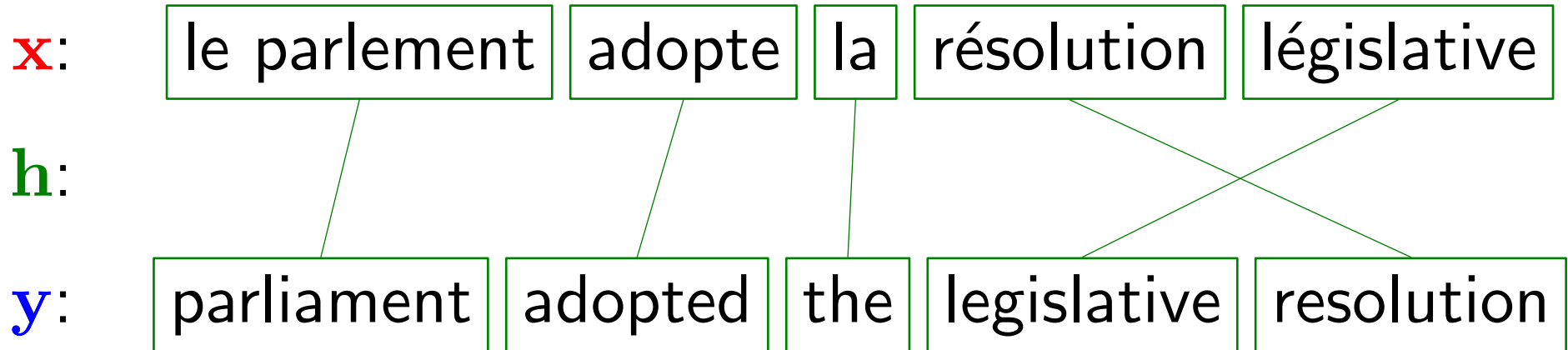
- Correct output is ill-defined
- Correspondences are missing

Why is this hard?



- Correct output is ill-defined
- Correspondences are missing

Why is this hard?



- Correct output is ill-defined
- Correspondences are missing
- Hidden correspondence is abused

Why is this hard?

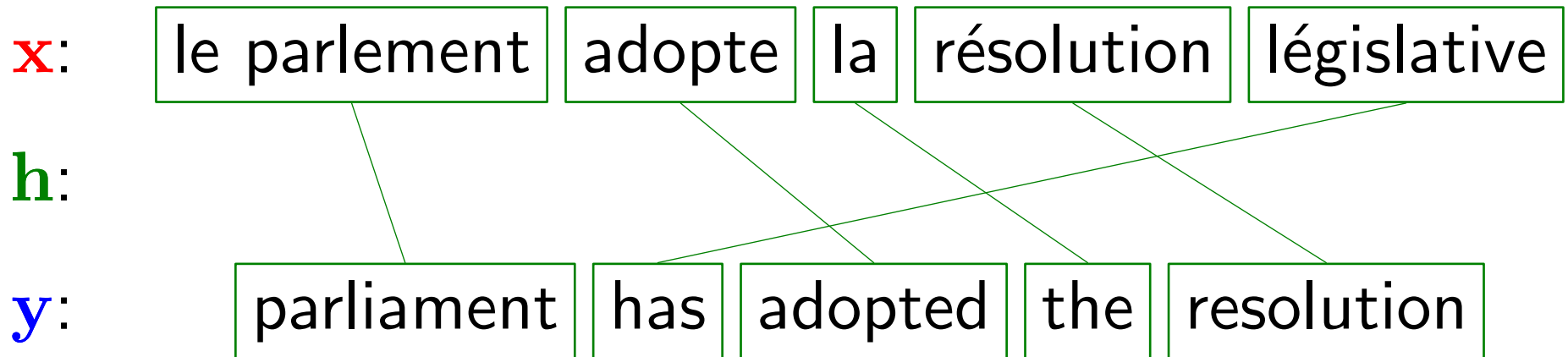
x: le parlement adopte la résolution législative

h:

y: parliament has adopted the resolution

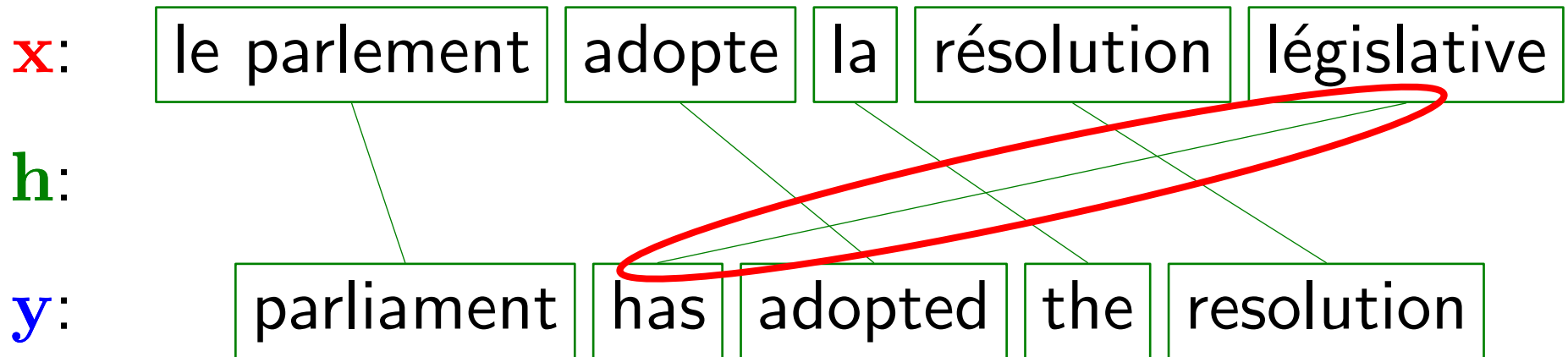
- Correct output is ill-defined
- Correspondences are missing
- Hidden correspondence is abused

Why is this hard?



- Correct output is ill-defined
- Correspondences are missing
- Hidden correspondence is abused

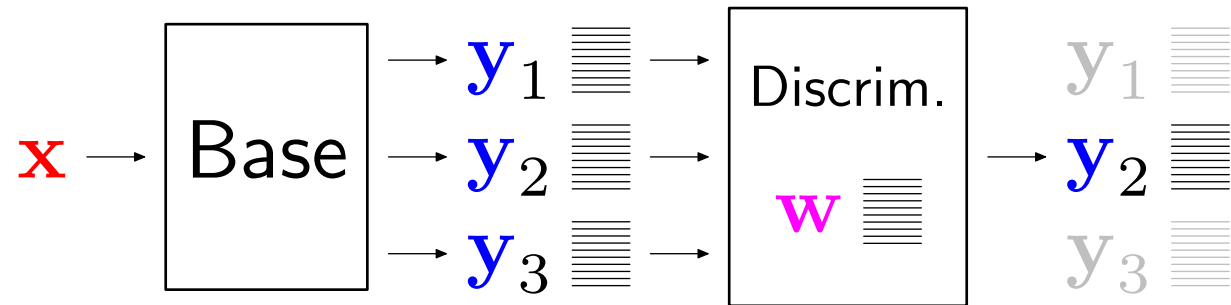
Why is this hard?



- Correct output is ill-defined
- Correspondences are missing
- Hidden correspondence is abused

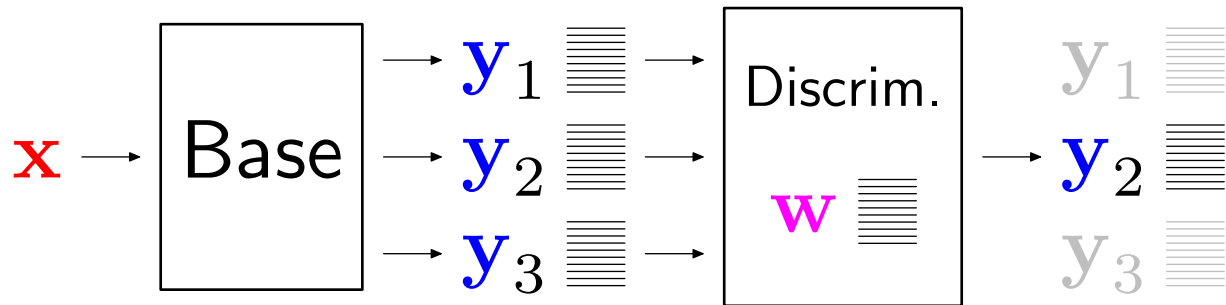
Discriminative approaches

Reranking [Shen, et al. '04; Och, et al. '04]

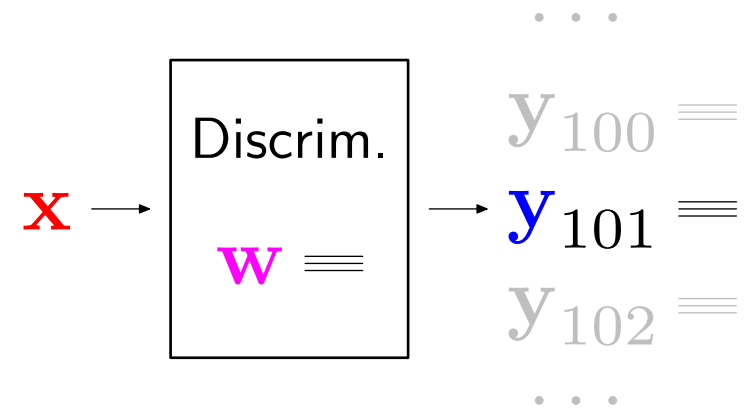


Discriminative approaches

Reranking [Shen, et al. '04; Och, et al. '04]

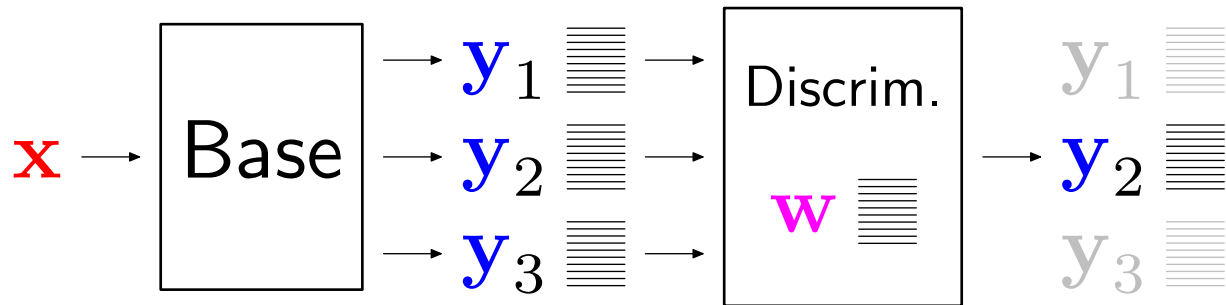


MERT [Och '03]

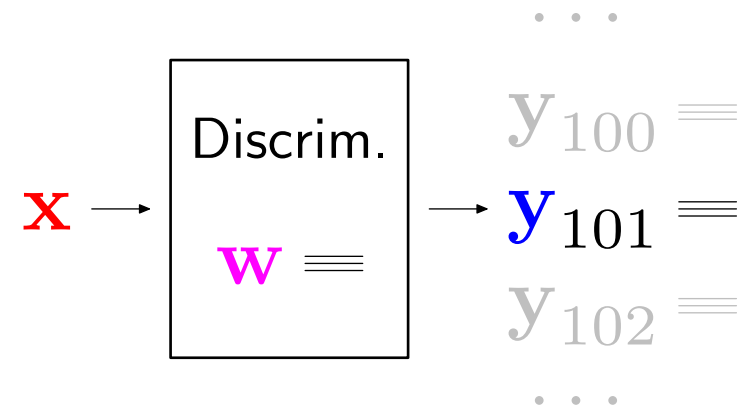


Discriminative approaches

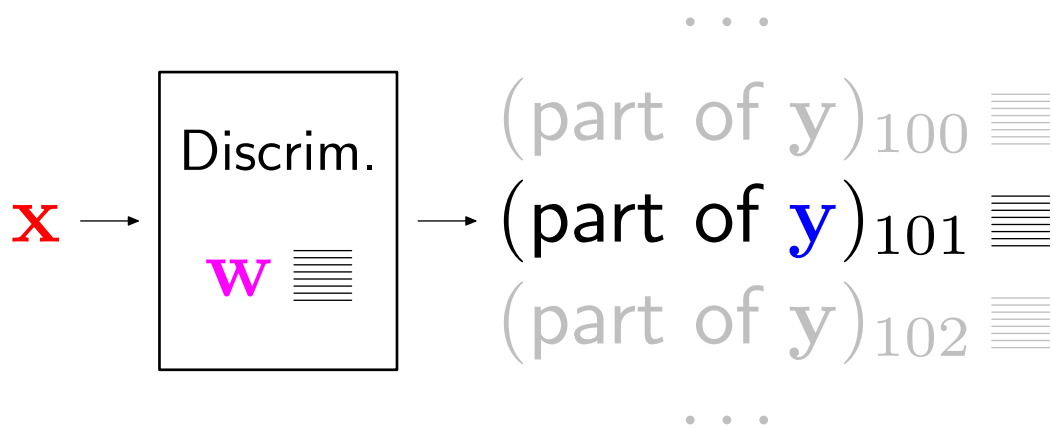
Reranking [Shen, et al. '04; Och, et al. '04]



MERT [Och '03]

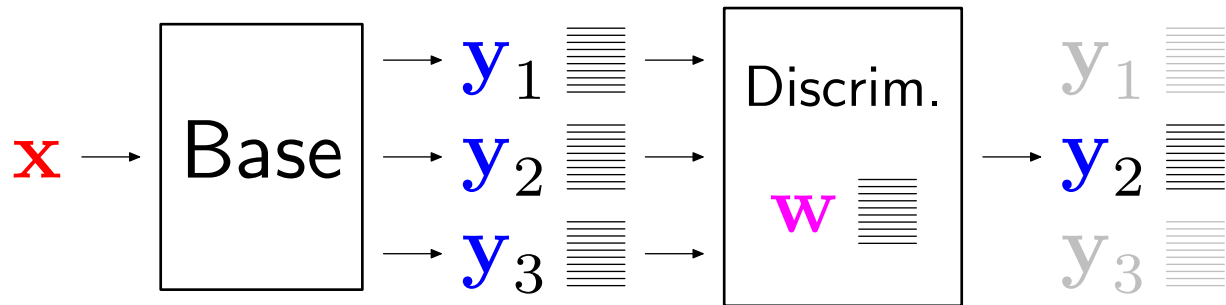


Local [Tillmann, Zhang '05]

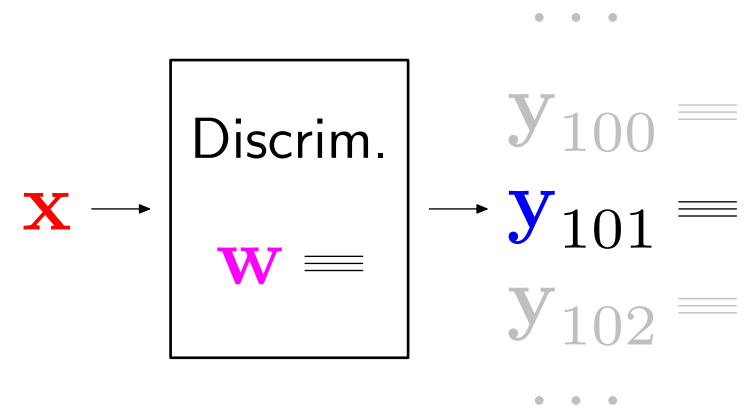


Discriminative approaches

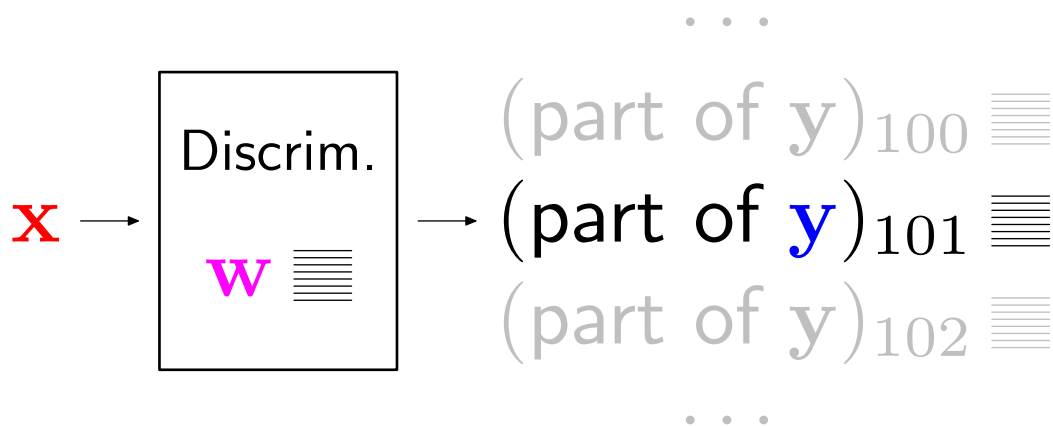
Reranking [Shen, et al. '04; Och, et al. '04]



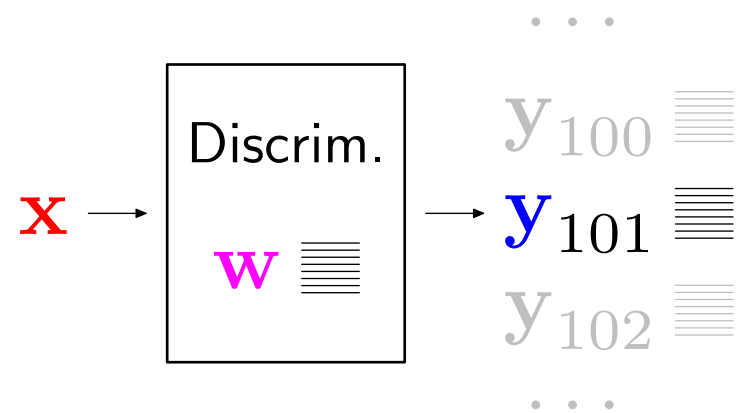
MERT [Och '03]



Local [Tillmann, Zhang '05]



Our end-to-end approach



Experimental setup

TRAIN ('99–'01)

414 sentence pairs

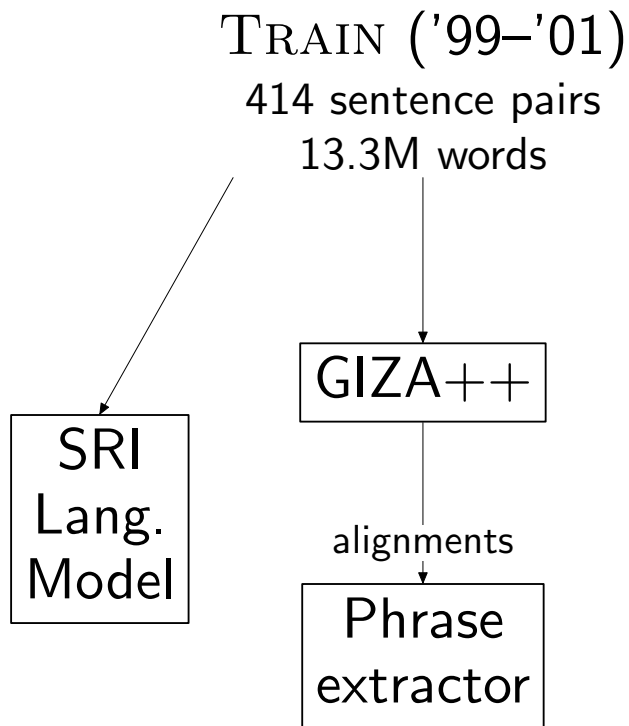
13.3M words

French-English Europarl corpus

DEV.5–15 ('02)
first 1K sentence pairs
10.4K words

TEST.5–15 ('03)
first 1K sentence pairs
10.8K words

Experimental setup

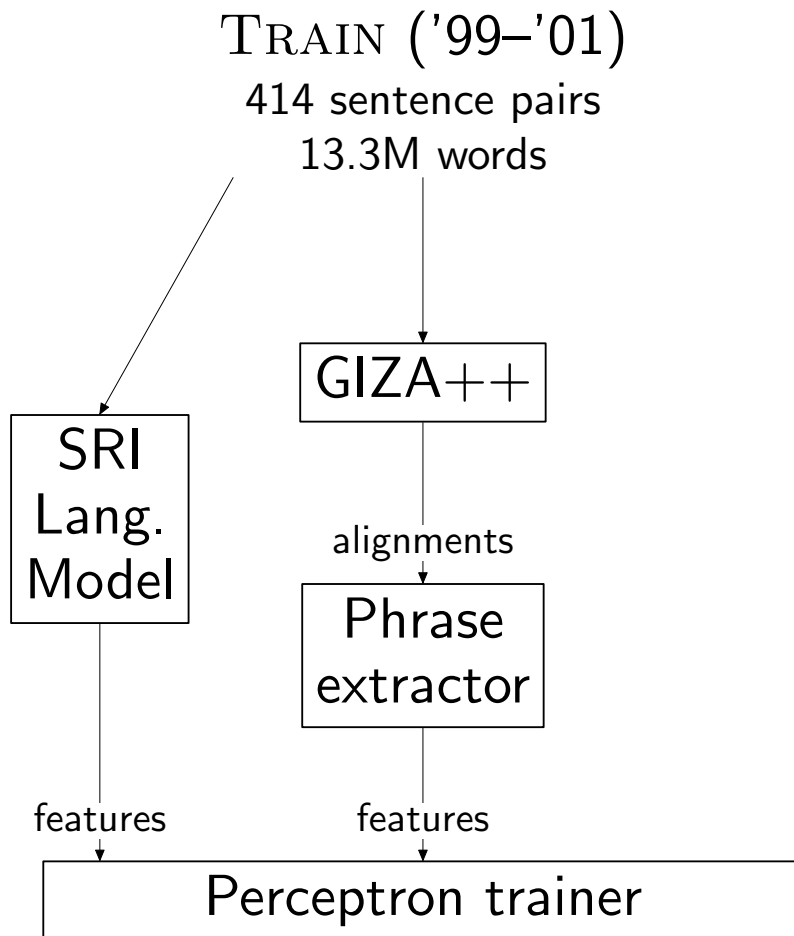


French-English Europarl corpus

DEV.5-15 ('02)
first 1K sentence pairs
10.4K words

TEST.5-15 ('03)
first 1K sentence pairs
10.8K words

Experimental setup

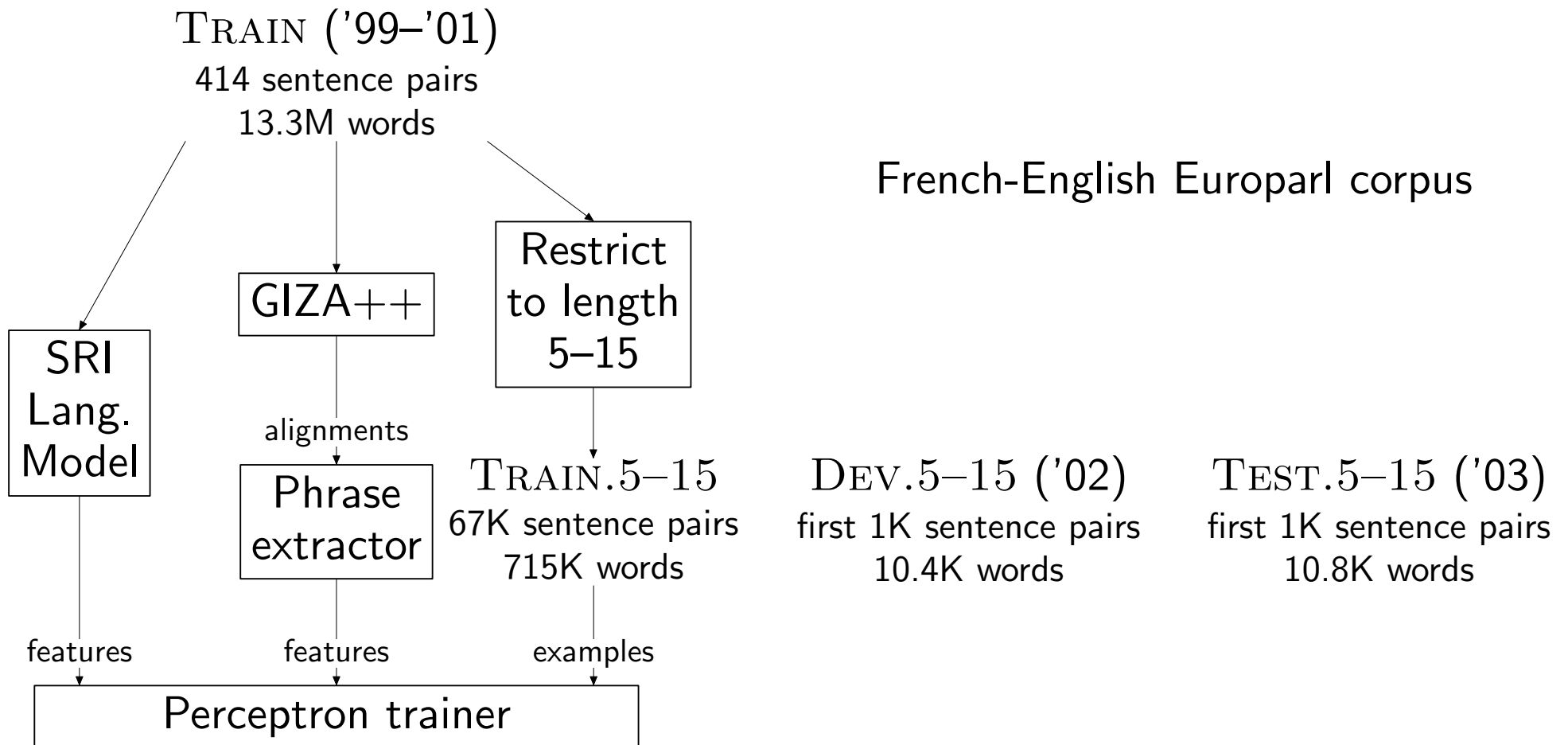


French-English Europarl corpus

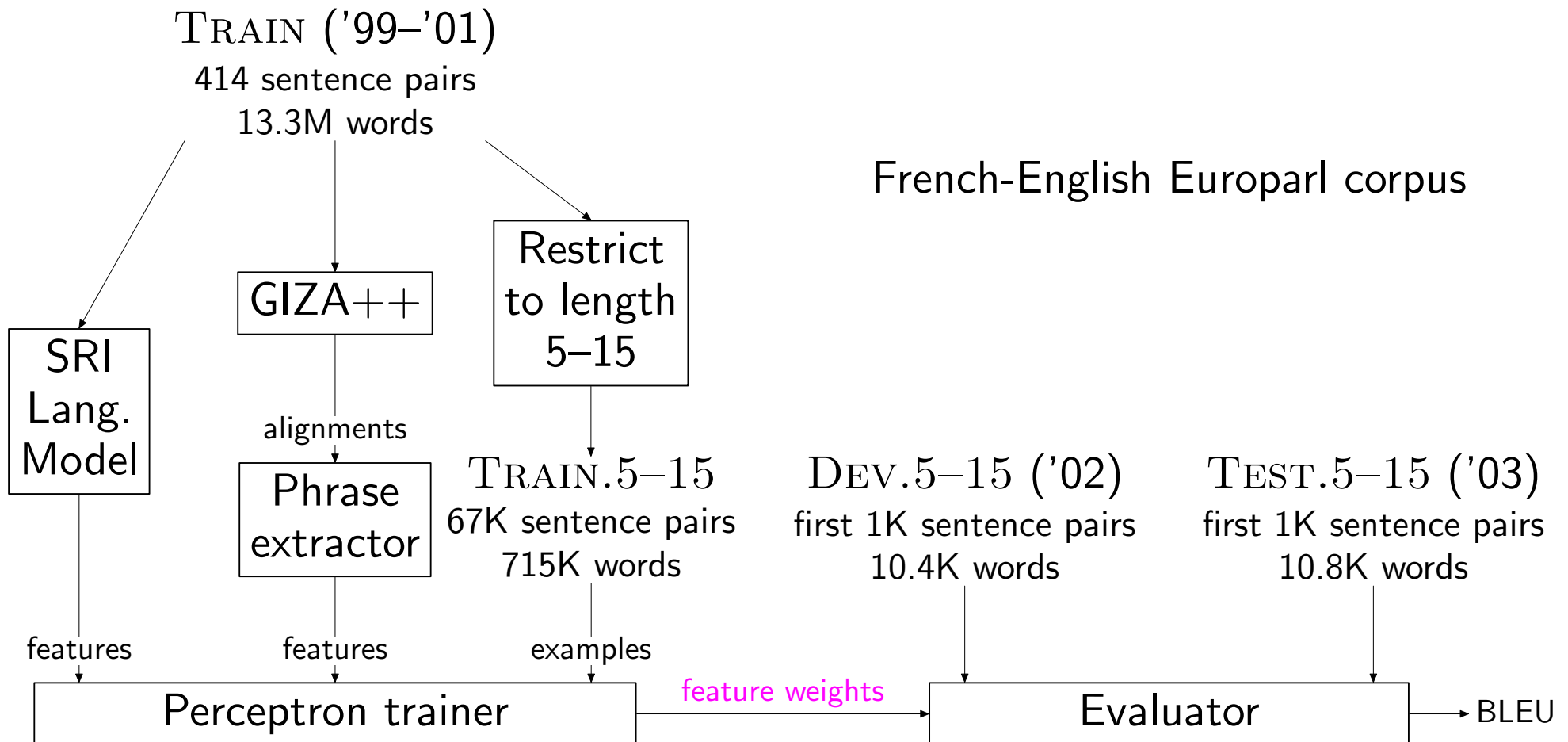
DEV.5-15 ('02)
first 1K sentence pairs
10.4K words

TEST.5-15 ('03)
first 1K sentence pairs
10.8K words

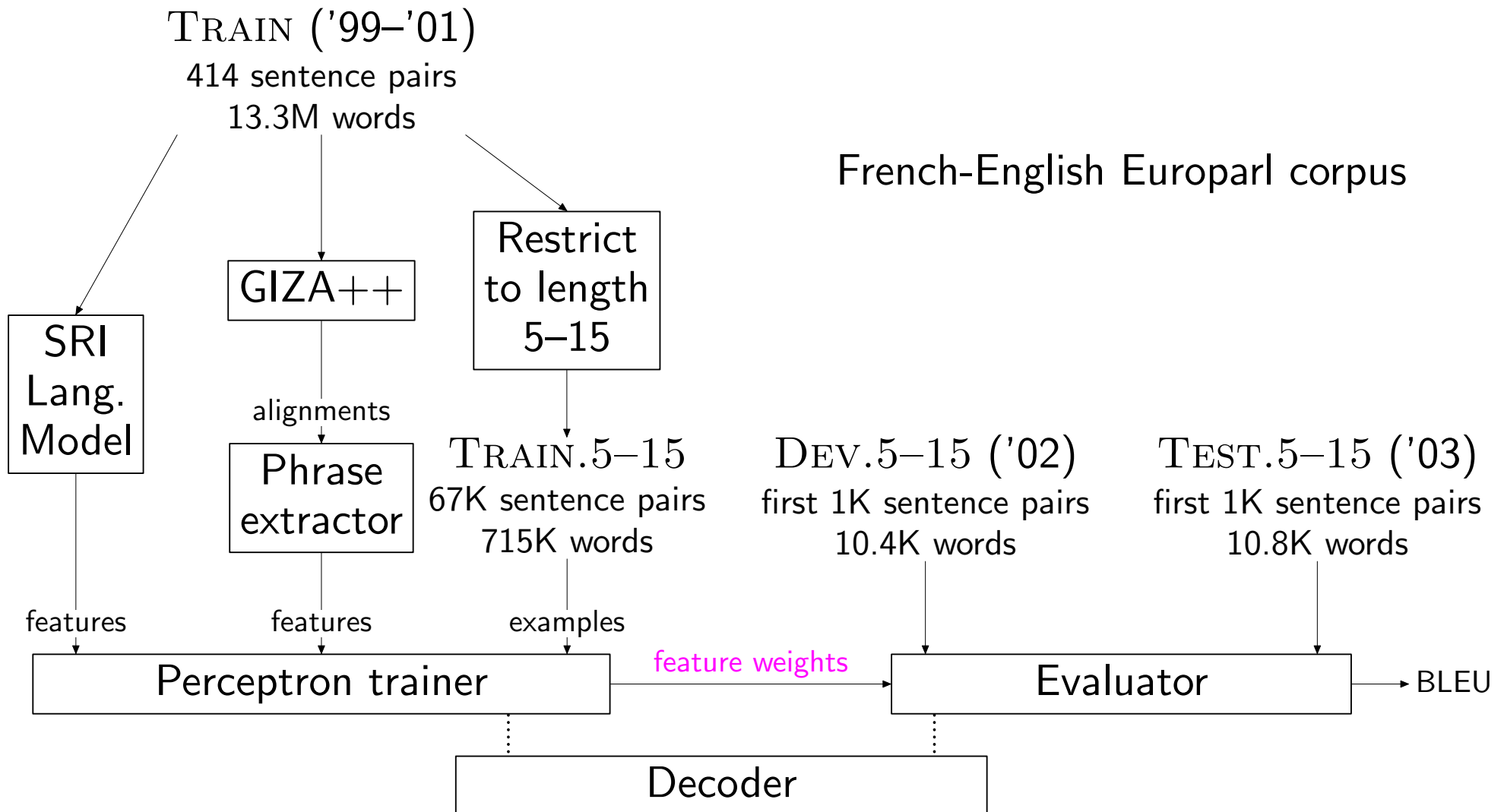
Experimental setup



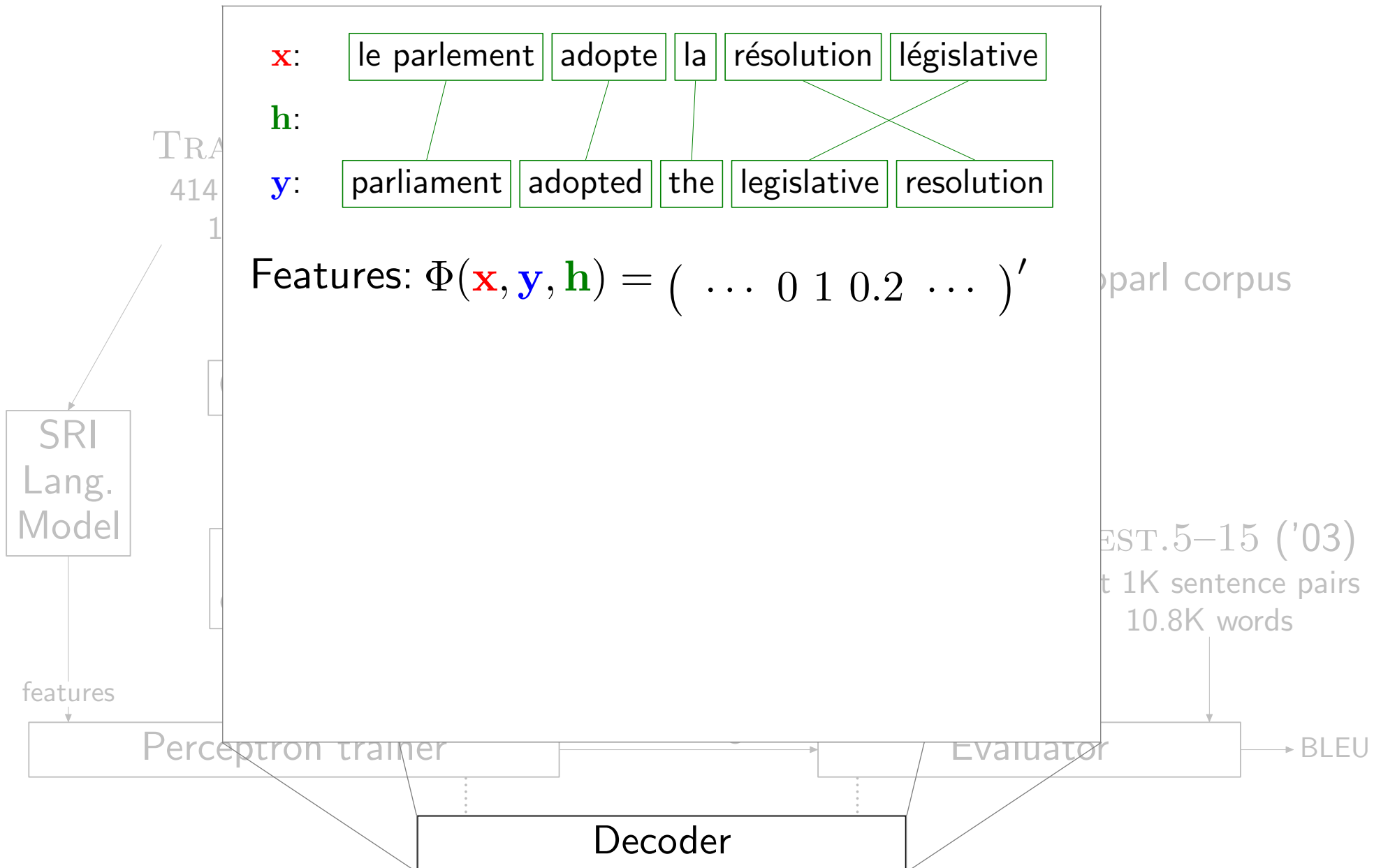
Experimental setup



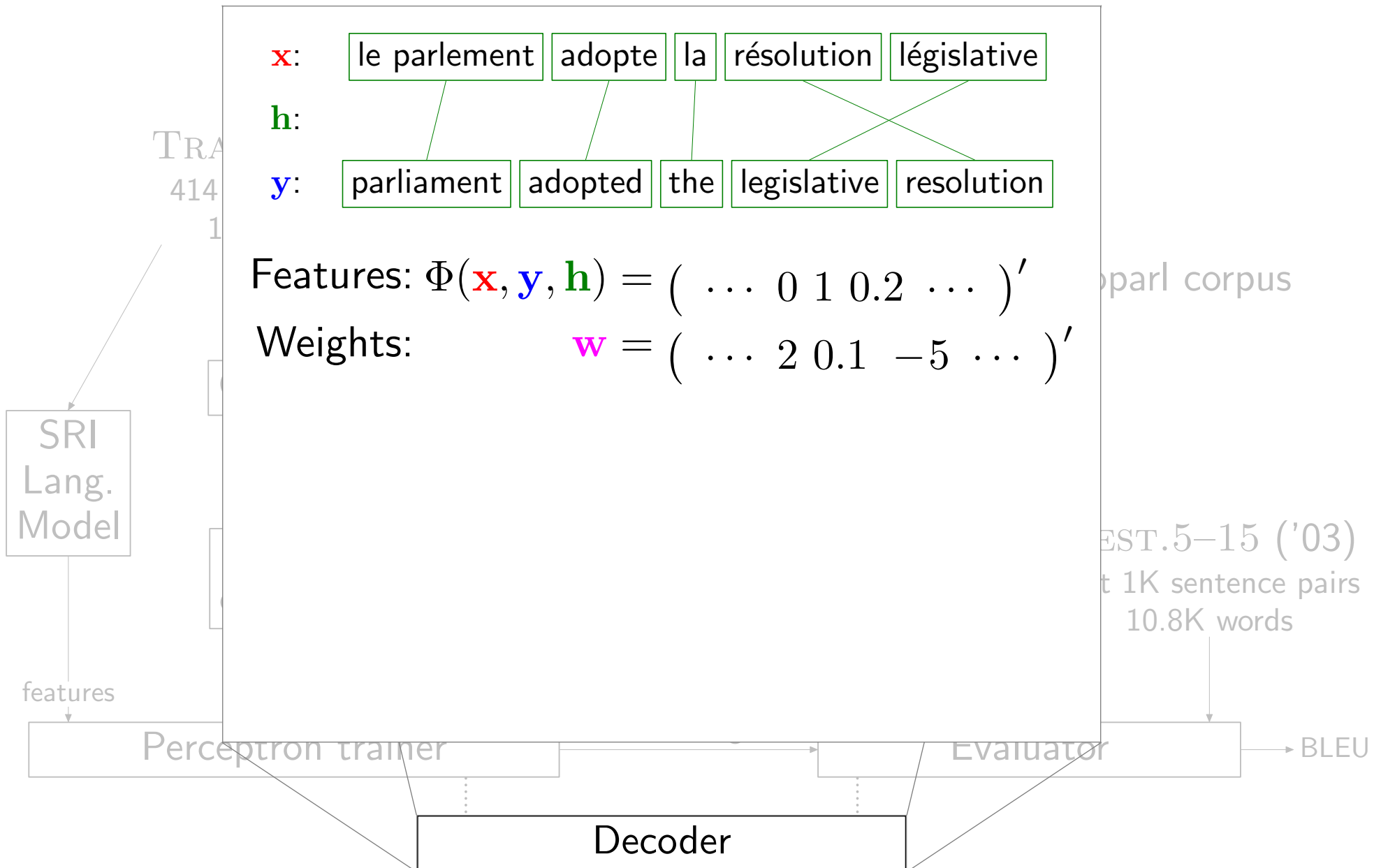
Experimental setup



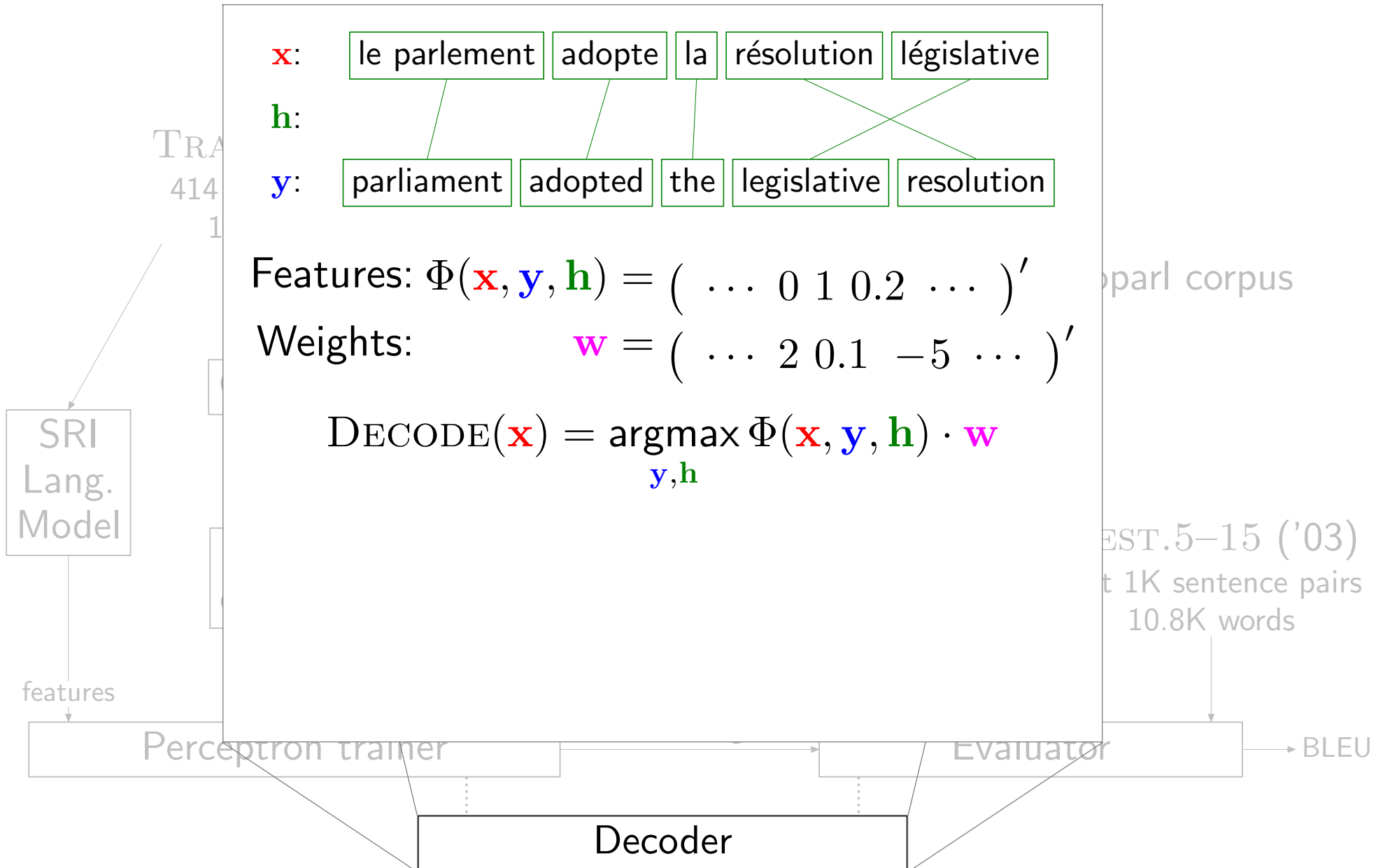
Experimental setup



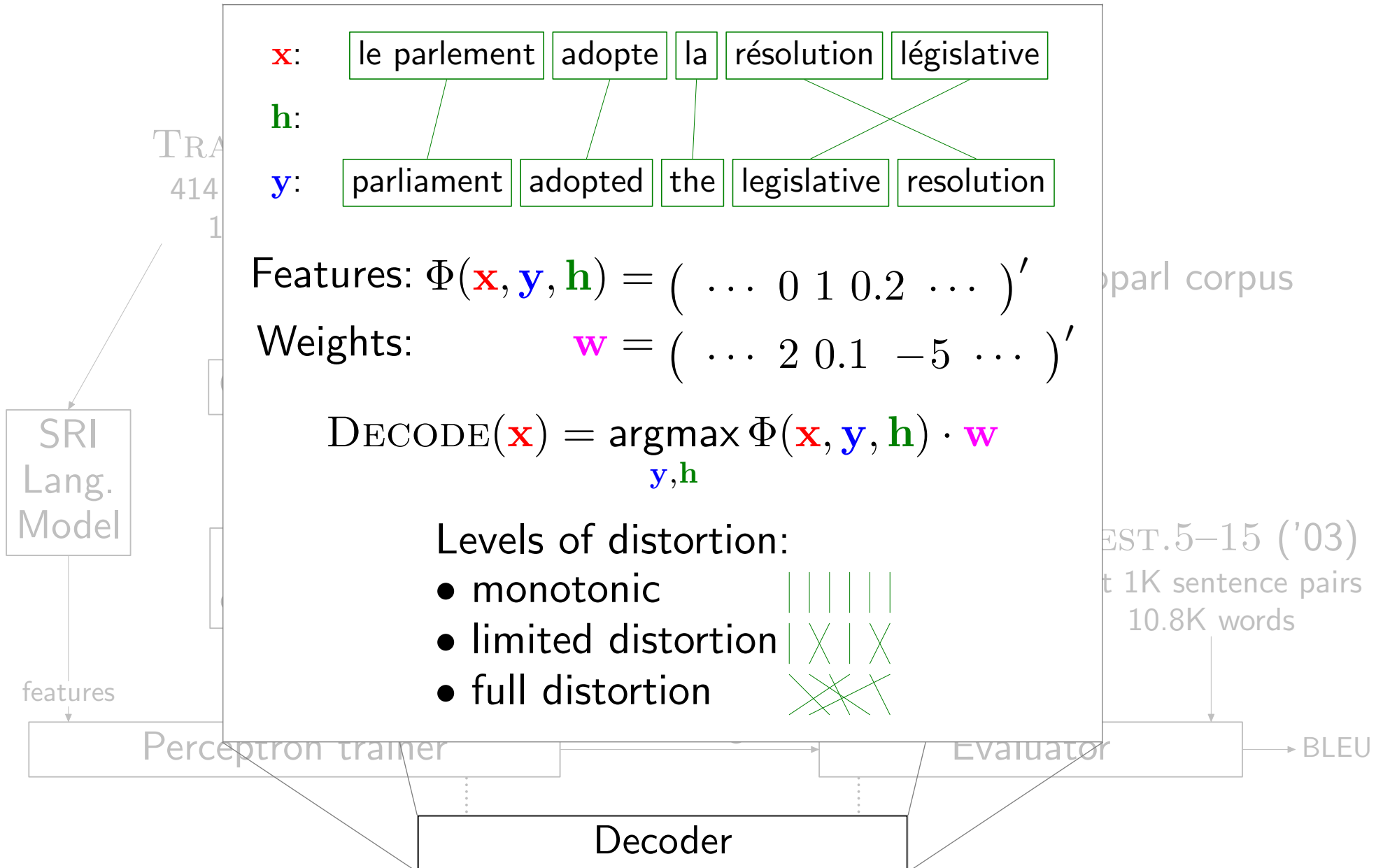
Experimental setup



Experimental setup



Experimental setup



Perceptron training

For each training example (\mathbf{x}, \mathbf{y}) : [Collins '02]

$$\begin{array}{l|l} \mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t) & \mathbf{y}_t = \mathbf{y} \\ \mathbf{w} \leftarrow \mathbf{w} - \Phi(\mathbf{x}, \mathbf{y}_p) & \mathbf{y}_p = \text{DECODE}(\mathbf{x}) \end{array}$$

Perceptron training

For each training example (\mathbf{x}, \mathbf{y}) : [Collins '02]

$$\begin{array}{l|l} \mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t) & \mathbf{y}_t = \mathbf{y} \\ -\Phi(\mathbf{x}, \mathbf{y}_p) & \mathbf{y}_p = \text{DECODE}(\mathbf{x}) \end{array}$$

$$\begin{array}{l|l} \mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) & \mathbf{y}_t, \mathbf{h}_t = ??? \\ -\Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p) & \mathbf{y}_p, \mathbf{h}_p = \text{DECODE}(\mathbf{x}) \end{array}$$

Perceptron training

For each training example (\mathbf{x}, \mathbf{y}) : [Collins '02]

$$\mathbf{w} \leftarrow \mathbf{w} \begin{array}{l} +\Phi(\mathbf{x}, \mathbf{y}_t) \\ -\Phi(\mathbf{x}, \mathbf{y}_p) \end{array} \quad \left| \begin{array}{l} \mathbf{y}_t = \mathbf{y} \\ \mathbf{y}_p = \text{DECODE}(\mathbf{x}) \end{array} \right.$$

$$\mathbf{w} \leftarrow \mathbf{w} \begin{array}{l} +\Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) \\ -\Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p) \end{array} \quad \left| \begin{array}{l} \mathbf{y}_t, \mathbf{h}_t = ??? \\ \mathbf{y}_p, \mathbf{h}_p = \text{DECODE}(\mathbf{x}) \end{array} \right.$$

How to choose $\mathbf{y}_t, \mathbf{h}_t$?

There are several choices and the choice does matter

Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)

\mathbf{x} : voté sur demande d'urgence

\mathbf{y} : vote on a request for urgent procedure

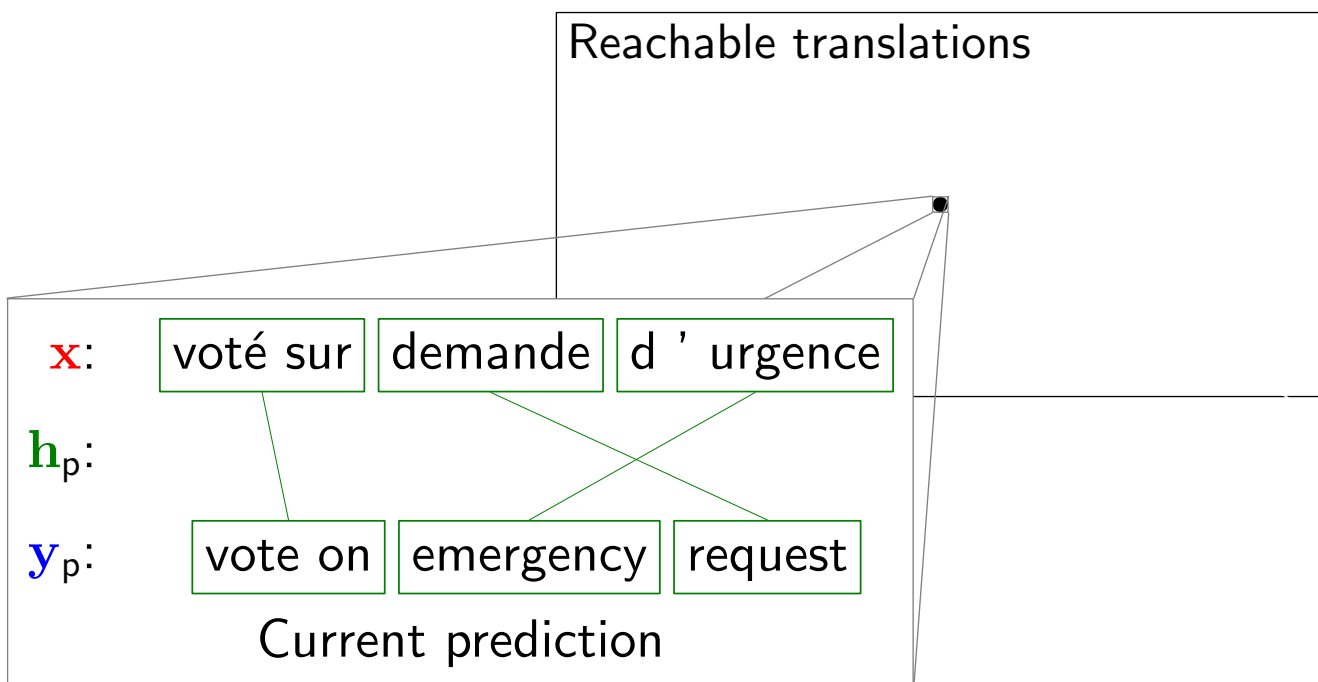
Update strategies

Training example (reference)

\mathbf{x} : voté sur demande d ' urgence

\mathbf{y} : vote on a request for urgent procedure

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$



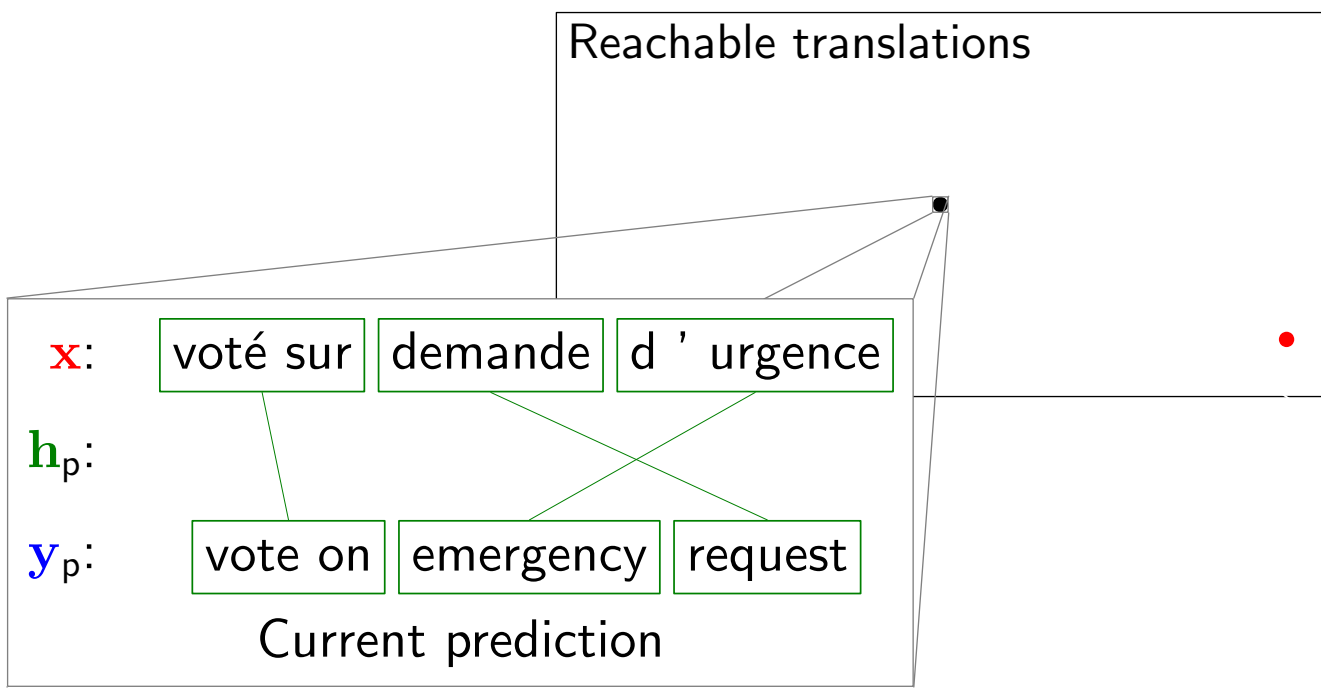
Update strategies

Training example (reference)

\mathbf{x} : voté sur demande d ' urgence

\mathbf{y} : vote on a request for urgent procedure

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$



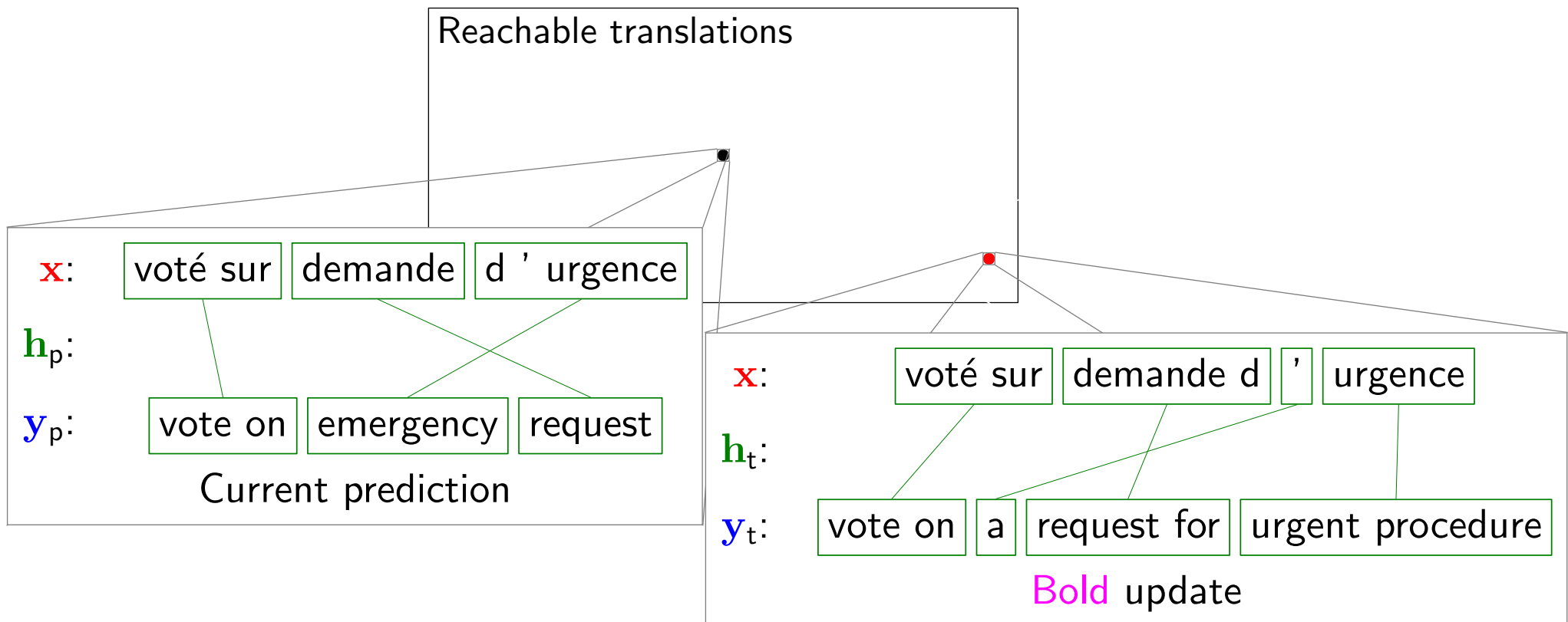
Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)

\mathbf{x} : voté sur demande d ' urgence

\mathbf{y} : vote on a request for urgent procedure



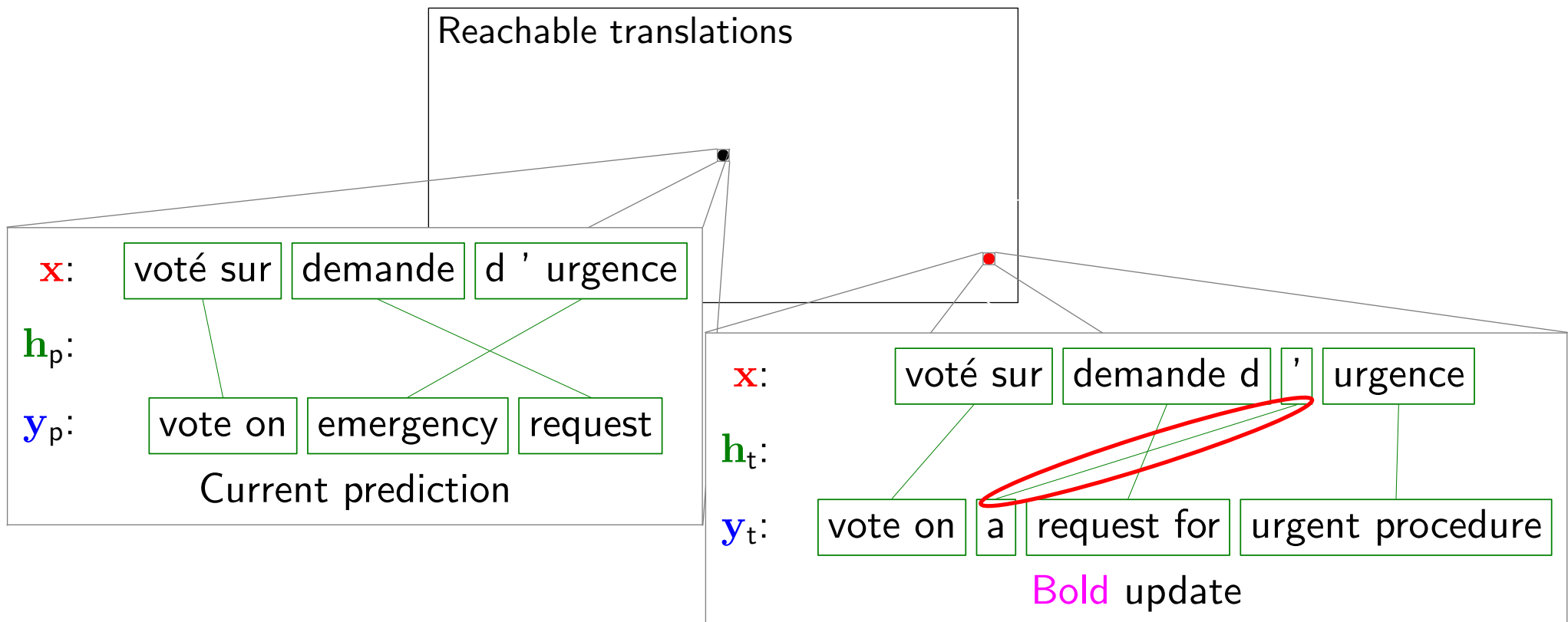
Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)

\mathbf{x} : voté sur demande d ' urgence

\mathbf{y} : vote on a request for urgent procedure



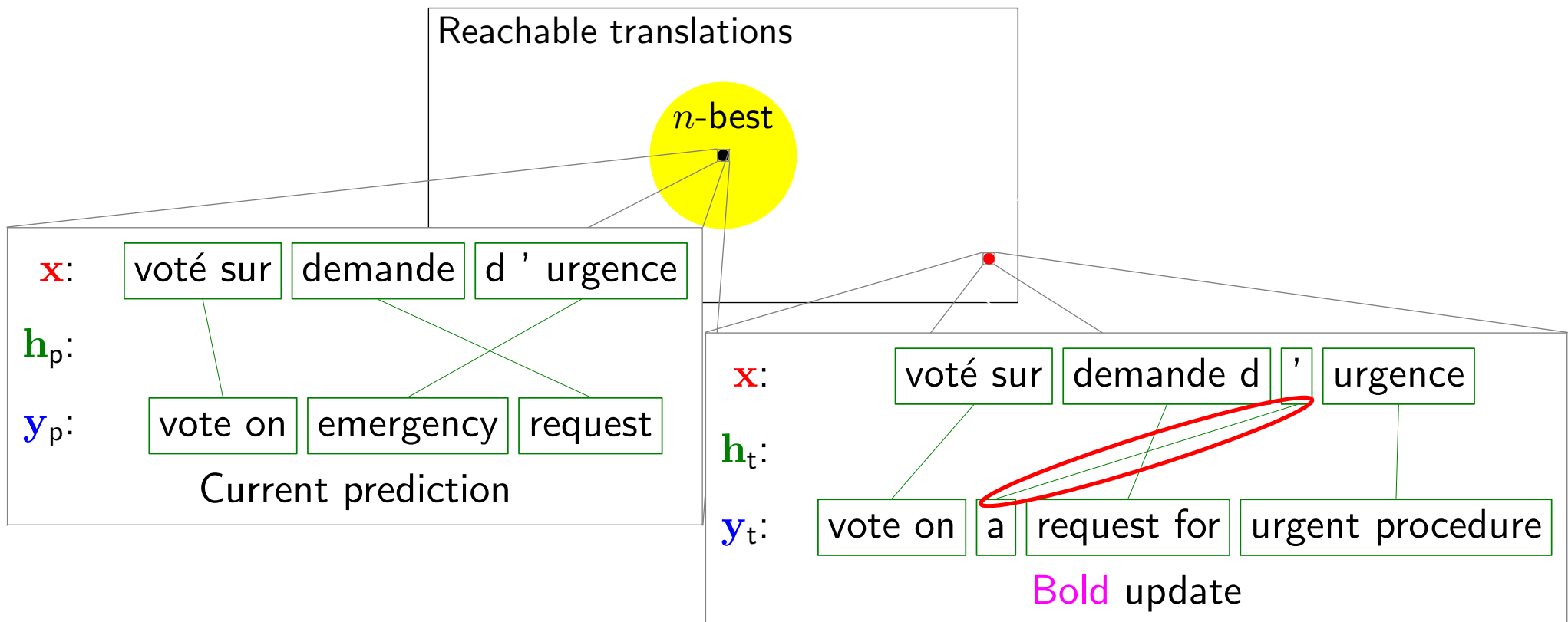
Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)

\mathbf{x} : voté sur demande d ' urgence

\mathbf{y} : vote on a request for urgent procedure

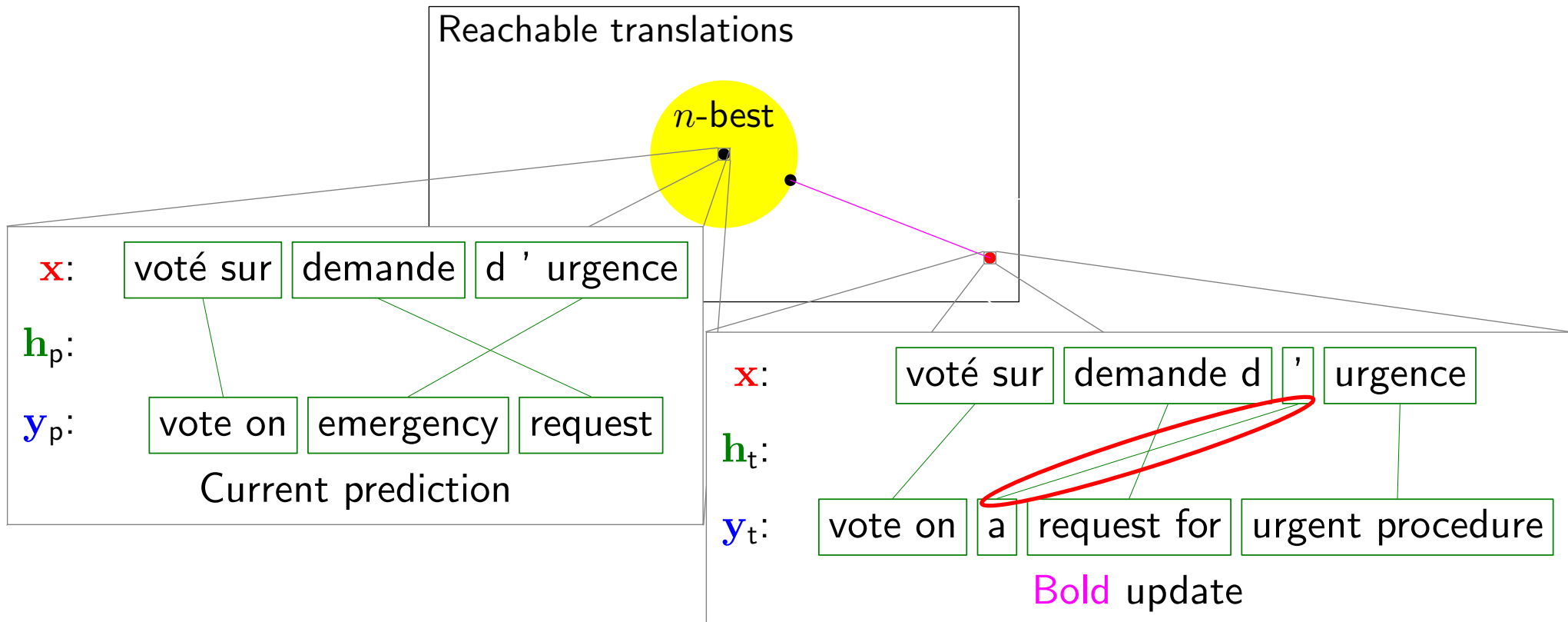


Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)

\mathbf{x} : voté sur demande d ' urgence
 \mathbf{y} : vote on a request for urgent procedure

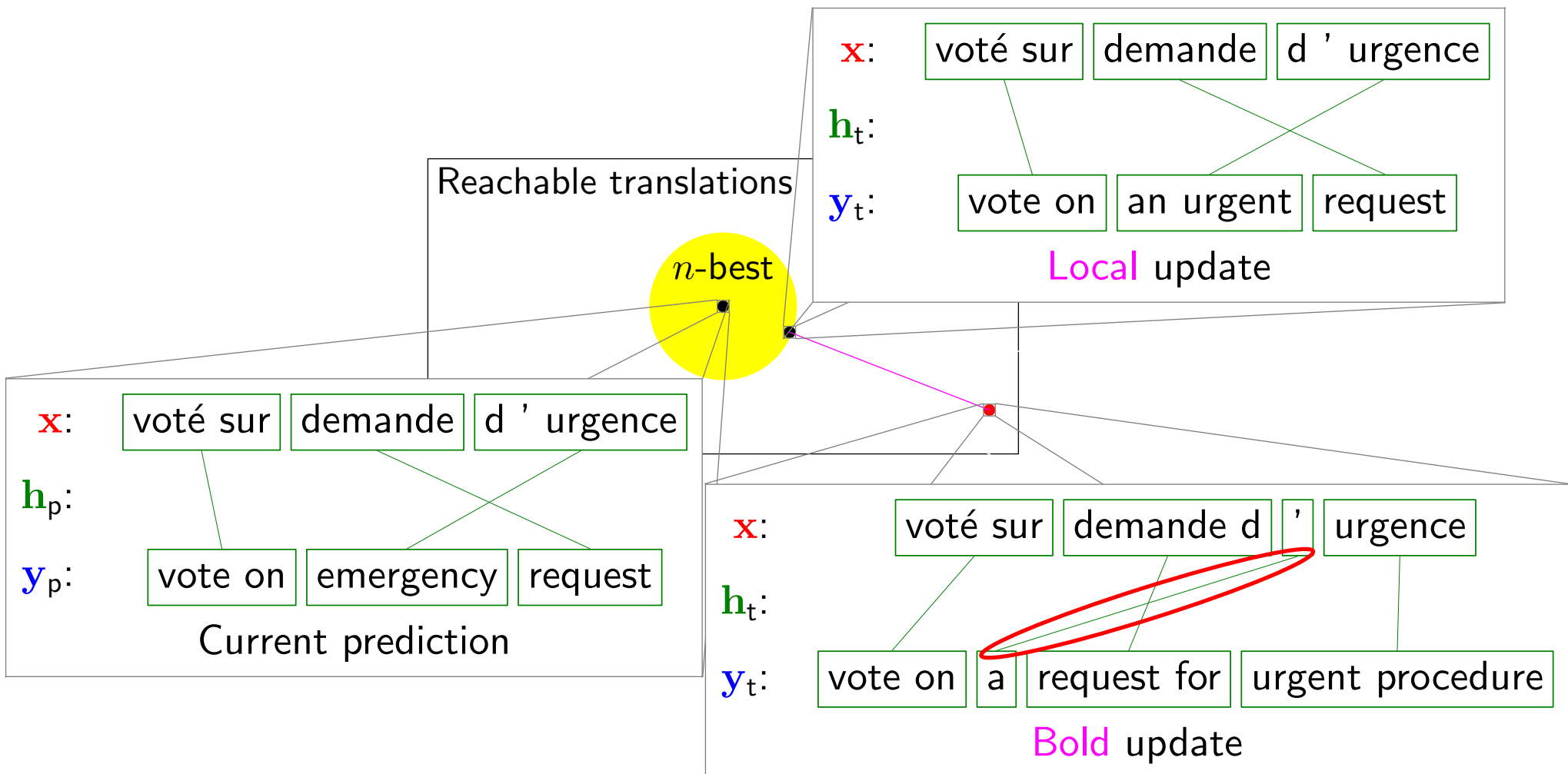


Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)

\mathbf{x} : voté sur demande d'urgence
 \mathbf{y} : vote on a request for urgent procedure



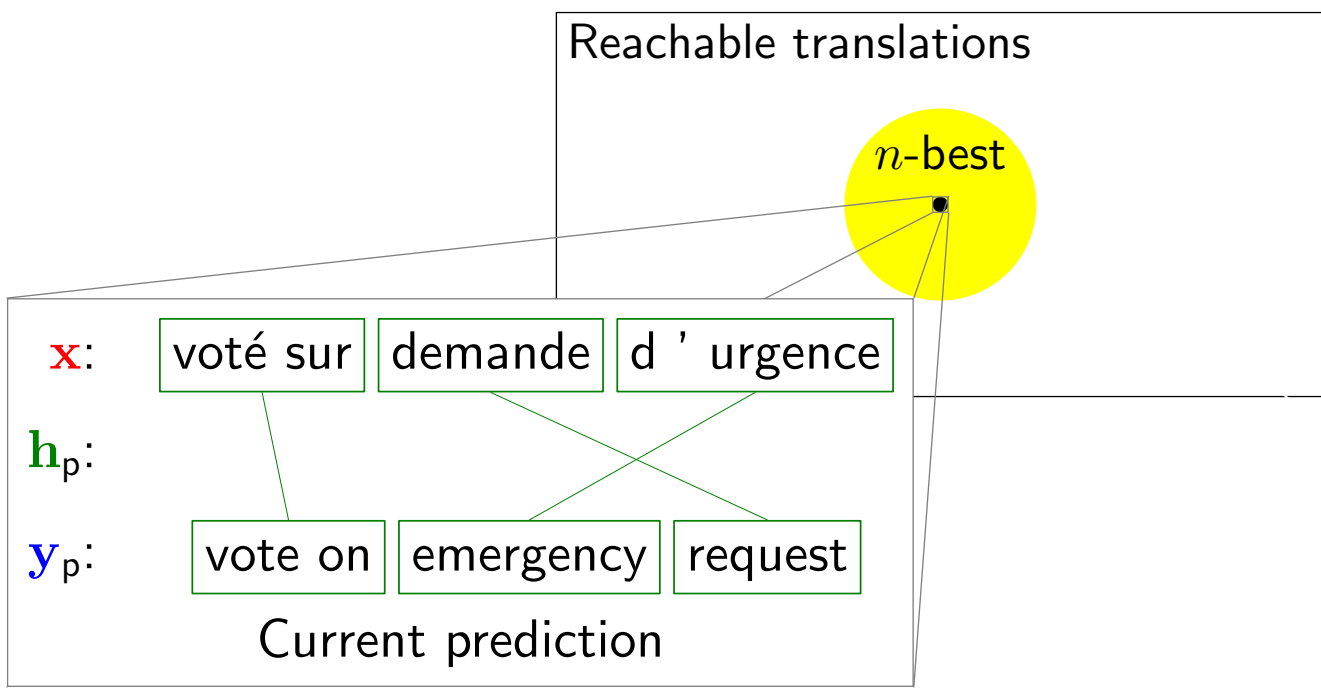
Update strategies

Training example (reference)

\mathbf{x} : voté sur demande d ' urgence

\mathbf{y} : vote on a request for urgent procedure

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

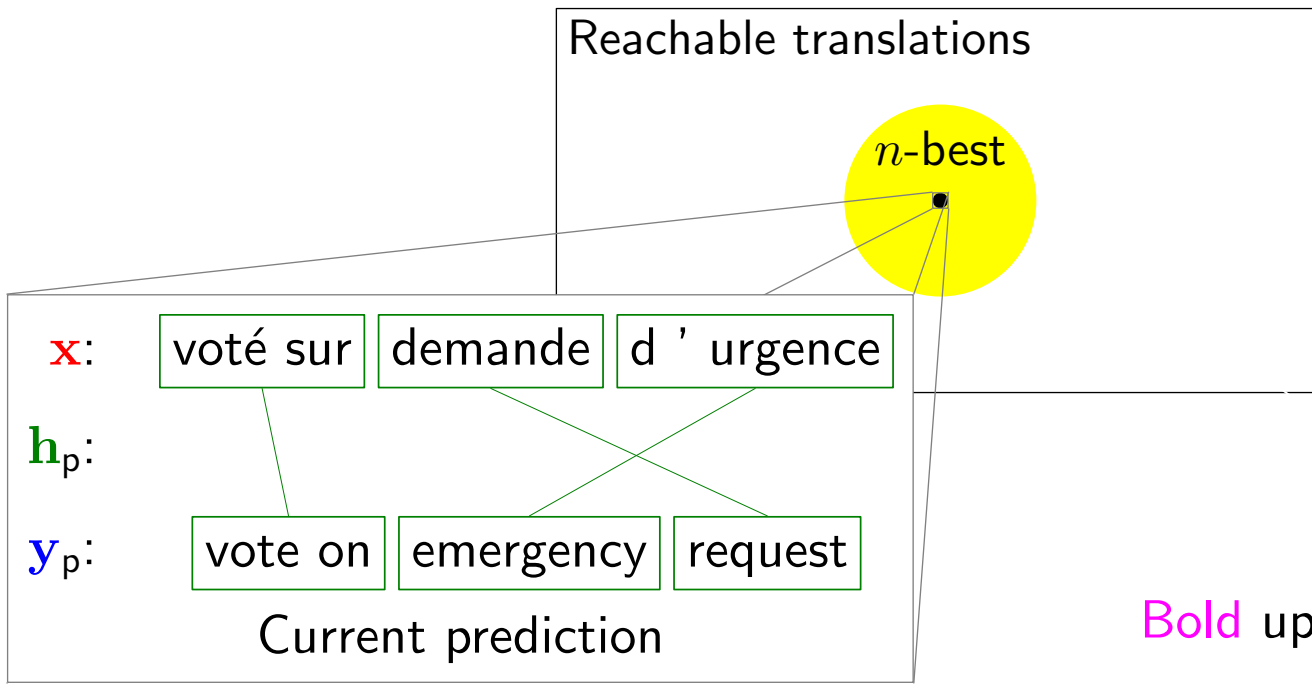


Update strategies

Training example (reference)

\mathbf{x} : voté sur demande d ' urgence
 \mathbf{y} : vote on a request for urgent procedure

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$



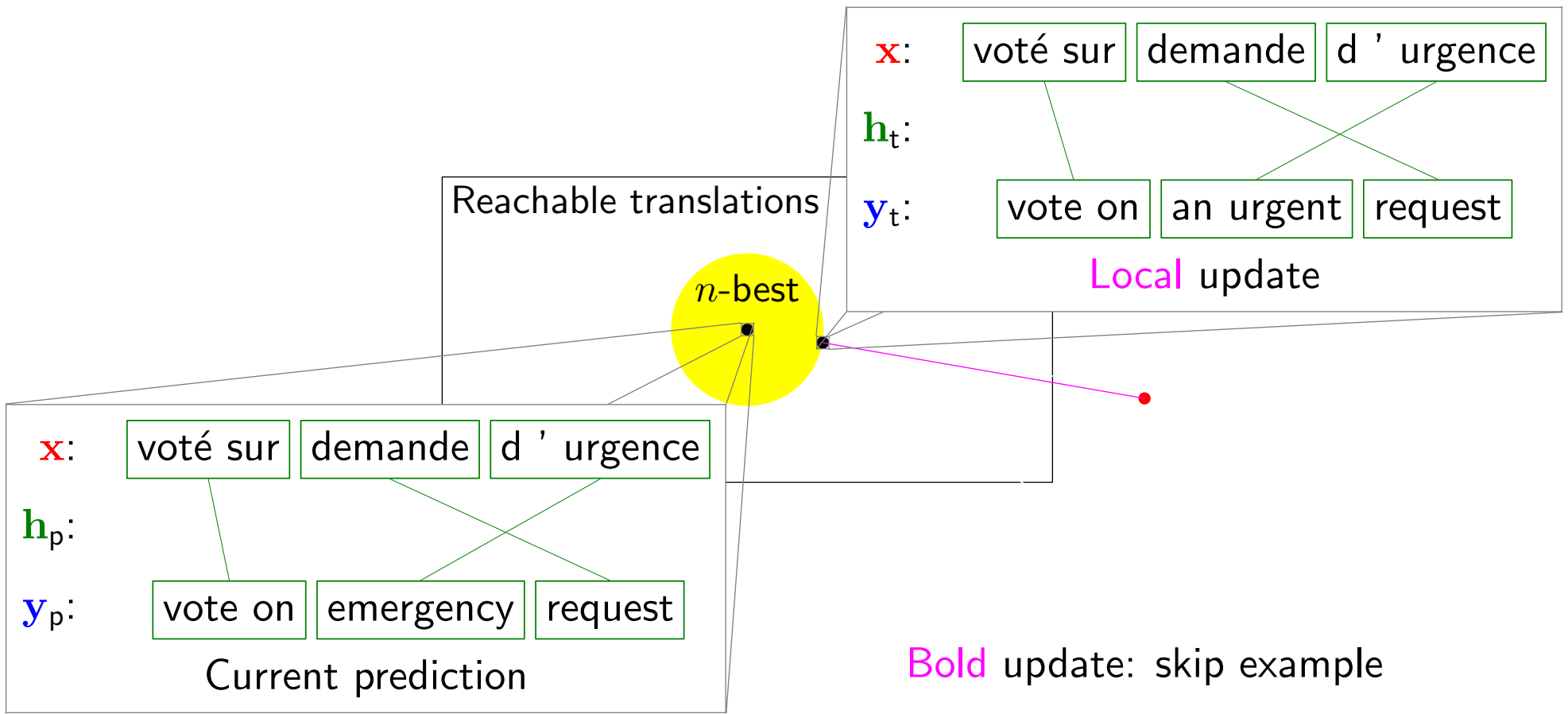
Bold update: skip example

Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)

\mathbf{x} : voté sur demande d'urgence
 \mathbf{y} : vote on a request for urgent procedure



Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)

\mathbf{x} : voté sur demande d'urgence

\mathbf{y} : vote on a request for urgent procedure

\mathbf{x} :

\mathbf{h}_t :

Decoder	Bold	Local
Monotonic	34.3	34.6
Limited distortion	33.5	34.7

te

\mathbf{x} :

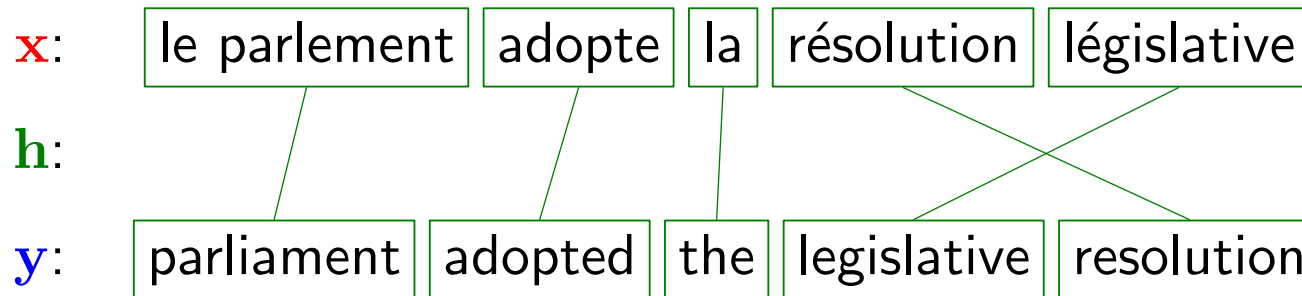
\mathbf{h}_p :

\mathbf{y}_p :

Current prediction

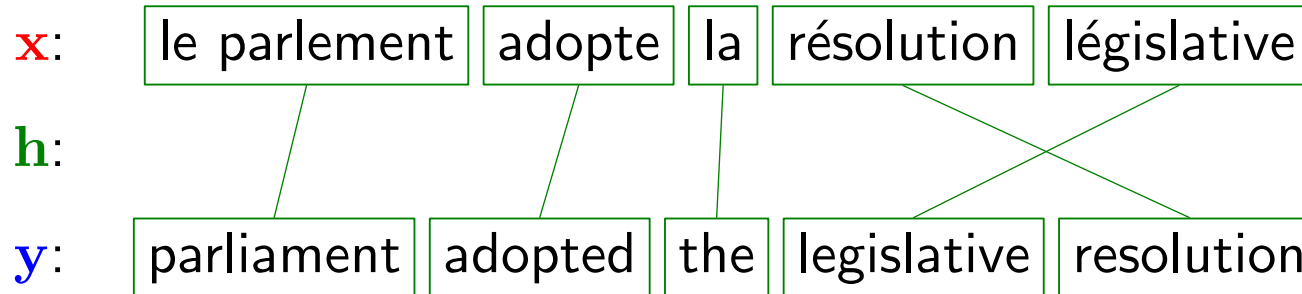
Bold update: skip example

Features



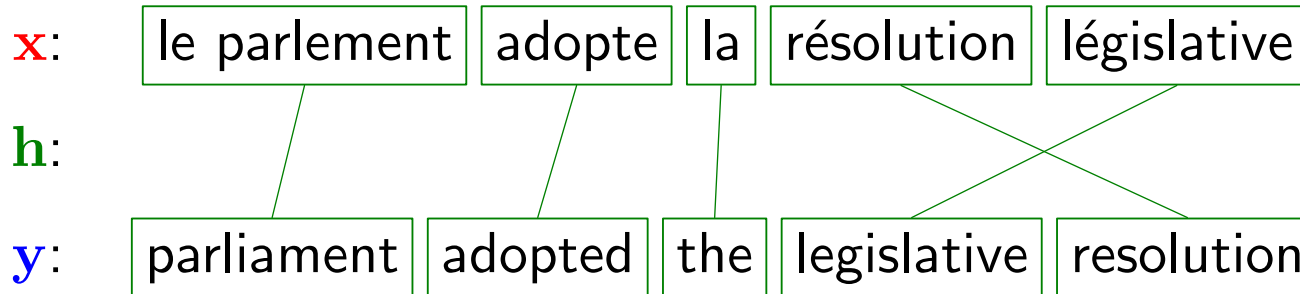
$\Phi(\mathbf{x}, \mathbf{y}, \mathbf{h}) =$

Features



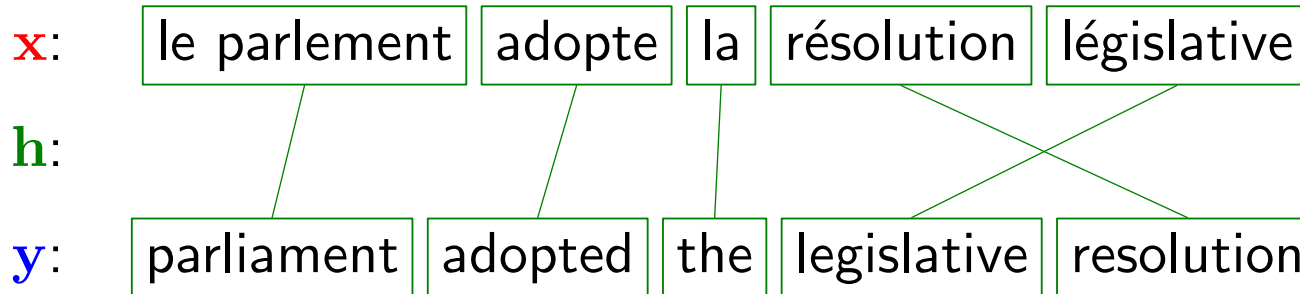
$$\Phi(\mathbf{x}, \mathbf{y}, \mathbf{h}) = \left\{ \begin{array}{ll} \text{BLANKET} & \text{TranslationLogProb} \quad -8.92 \\ & \text{LangModelLogProb} \quad -2.462 \end{array} \right.$$

Features



$$\Phi(\mathbf{x}, \mathbf{y}, \mathbf{h}) = \left\{ \begin{array}{ll} \text{BLANKET} & \begin{array}{l} \text{TranslationLogProb} \quad -8.92 \\ \text{LangModelLogProb} \quad -2.462 \end{array} \\ \text{LEXICAL} & \begin{array}{l} \boxed{\text{le parlement}} \text{---} \boxed{\text{parliament}} \quad 1 \\ [\text{parliament adopted the}] \quad 1 \\ \dots \end{array} \end{array} \right.$$

Features



$$\Phi(\mathbf{x}, \mathbf{y}, \mathbf{h}) = \left\{ \begin{array}{lll} \text{BLANKET} & \begin{array}{l} \text{TranslationLogProb} \\ \text{LangModelLogProb} \end{array} & \begin{array}{l} -8.92 \\ -2.462 \end{array} \\ \text{LEXICAL} & \begin{array}{l} \boxed{\text{le parlement}} \text{---} \boxed{\text{parliament}} \\ [\text{parliament adopted the}] \end{array} & \begin{array}{l} 1 \\ 1 \end{array} \\ \dots & \dots & \dots \\ \text{POS} & \begin{array}{l} \boxed{\text{DT NN}} \text{---} \boxed{\text{NN}} \\ [\text{NN VBD DT}] \end{array} & \begin{array}{l} 1 \\ 1 \end{array} \\ \dots & \dots & \dots \\ & \begin{array}{l} \boxed{\text{NN}} \quad \boxed{\text{JJ}} \\ \boxed{\text{JJ}} \quad \boxed{\text{NN}} \end{array} \text{ swap} & 1 \end{array} \right.$$

Effect of features

trente-cinq	langues
five	languages

Features	DEV BLEU
BLANKET (untuned)	33.0

Effect of features

trente-cinq	langues
five	languages

Tune BLANKET features...

- Resulting relative weights: $w(\text{TranslationLogProb}) = 2$
 $w(\text{LangModelLogProb}) = 1$

Features	DEV BLEU
BLANKET (untuned)	33.0

Effect of features

trente-cinq	langues
five	languages
thirty-five	languages

Tune BLANKET features...

- Resulting relative weights: $w(\text{TranslationLogProb}) = 2$
 $w(\text{LangModelLogProb}) = 1$

Features	DEV BLEU
BLANKET (untuned)	33.0

Effect of features

trente-cinq	langues
five	languages
thirty-five	languages

Tune BLANKET features...

- Resulting relative weights: $w(\text{TranslationLogProb}) = 2$
 $w(\text{LangModelLogProb}) = 1$

Features	DEV BLEU
BLANKET (untuned)	33.0
BLANKET	33.4

Effect of features

pour cela que	j ' ai	voté favorablement	.
for that	i have	voted in favour	.

Features	DEV BLEU
BLANKET (untuned)	33.0
BLANKET	33.4

Effect of features

pour cela que	j ' ai	voté favorablement	.
for that	i have	voted in favour	.

Add LEXICAL features...

- j ' ai — i have gets very negative weight
- Literal phrase translations downweighted to allow context-sensitive translations

Features	DEV BLEU
BLANKET (untuned)	33.0
BLANKET	33.4

Effect of features

pour cela que	j ' ai	voté favorablement	.
for that	i have	voted in favour	.
for that reason	i	voted in favour	.

Add LEXICAL features...

- j ' ai — i have gets very negative weight
- Literal phrase translations downweighted to allow context-sensitive translations

Features	DEV BLEU
BLANKET (untuned)	33.0
BLANKET	33.4

Effect of features

pour cela que	j ' ai	voté favorablement	.
for that	i have	voted in favour	.
for that reason	i	voted in favour	.

Add LEXICAL features...

- j ' ai — i have gets very negative weight
- Literal phrase translations downweighted to allow context-sensitive translations

Features	DEV BLEU
BLANKET (untuned)	33.0
BLANKET	33.4
BLANKET+LEXICAL	35.0

Effect of features

How can we generalize beyond lexical features?

Features	DEV BLEU
BLANKET (untuned)	33.0
BLANKET	33.4
BLANKET+LEXICAL	35.0

Effect of features

How can we generalize beyond lexical features?

Add POS features...

- `la réalisation du droit` — `the right` has negative weight
generalizes to `DT NN IN NN` — `DT NN`
- Number of nonzero features drops (1.55M to 1.24M)

Features	DEV BLEU
BLANKET (untuned)	33.0
BLANKET	33.4
BLANKET+LEXICAL	35.0

Effect of features

How can we generalize beyond lexical features?

Add POS features...

- `la réalisation du droit` — `the right` has negative weight
generalizes to `DT NN IN NN` — `DT NN`
- Number of nonzero features drops (1.55M to 1.24M)

Features	DEV BLEU
BLANKET (untuned)	33.0
BLANKET	33.4
BLANKET+LEXICAL	35.0
BLANKET+LEXICAL+POS	35.3

Alignment constellation features

Phrase-extraction heuristic is important [Koehn '03]

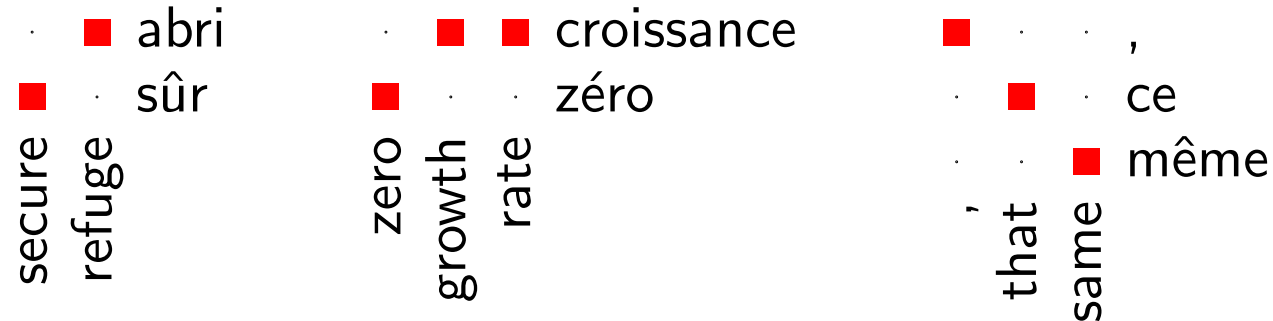
·	■	abri
■	·	sûr
secure		
refuge		

·	■	■	croissance
■	·	·	zéro
zero			
growth			
rate			

■	·	·	,
·	■	·	ce
·	·	■	même
,	that	same	

Alignment constellation features

Phrase-extraction heuristic is important [Koehn '03]

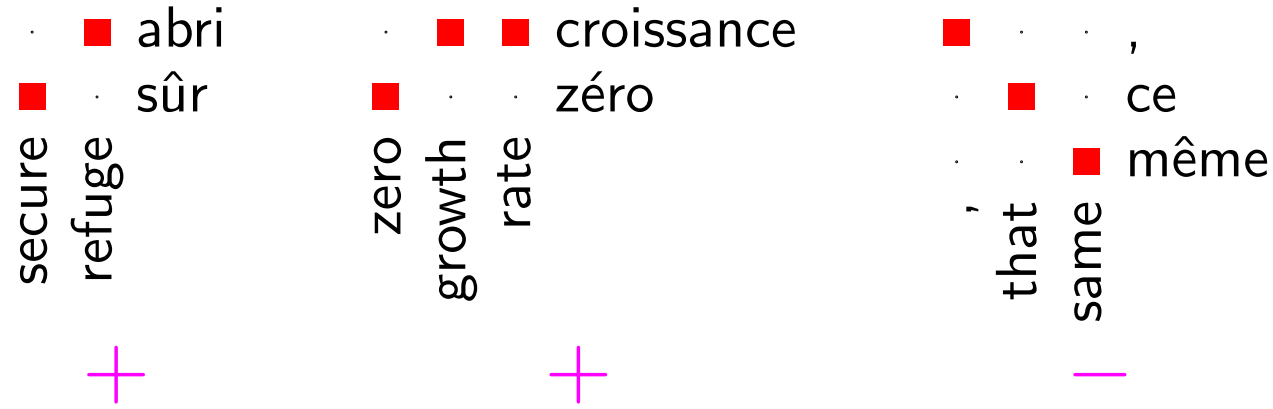


Put features on the phrase-extraction process itself

$$\Phi(\mathbf{x}, \mathbf{y}, \mathbf{h}) = \begin{cases} \begin{bmatrix} \blacksquare & \cdot \\ \cdot & \blacksquare \end{bmatrix} & 1 \\ \dots & \dots \end{cases}$$

Alignment constellation features

Phrase-extraction heuristic is important [Koehn '03]



Put features on the phrase-extraction process itself

$$\Phi(\mathbf{x}, \mathbf{y}, \mathbf{h}) = \begin{cases} \begin{bmatrix} \blacksquare & \cdot \\ \cdot & \blacksquare \end{bmatrix} & 1 \\ \dots & \dots \end{cases}$$

Summary of results

Monotonic systems	TEST BLEU
Pharaoh (MERT)	28.8
BLANKET	28.4
BLANKET+LEXICAL+POS	29.6

Summary of results

Monotonic systems	TEST BLEU
Pharaoh (MERT)	28.8
BLANKET	28.4
BLANKET+LEXICAL+POS	29.6

Distortion systems	TEST BLEU
Pharaoh (MERT) [Full]	29.5
BLANKET [Limited]	30.0
BLANKET+LEXICAL+POS [Limited]	30.9

Conclusions

- Perceptron training with hidden variables:
local updates are better

Conclusions

- Perceptron training with hidden variables:
local updates are better
- Allow many expressive features:
blanket, lexical, POS, alignment constellation

Conclusions

- Perceptron training with hidden variables:
local updates are better
- Allow many expressive features:
blanket, lexical, POS, alignment constellation
- Significant BLEU improvements over Pharaoh

Conclusions

- Perceptron training with hidden variables:
local updates are better
- Allow many expressive features:
blanket, lexical, POS, alignment constellation
- Significant BLEU improvements over Pharaoh
- Extension: beyond phrase-based models

Conclusions

- Perceptron training with hidden variables:
local updates are better
- Allow many expressive features:
blanket, lexical, POS, alignment constellation
- Significant BLEU improvements over Pharaoh
- Extension: beyond phrase-based models

That's it.