

From
 n -gram-based
to
CRF-based
Translation Models

Lavergne - Crego - Allauzen - Yvon

LIMSI

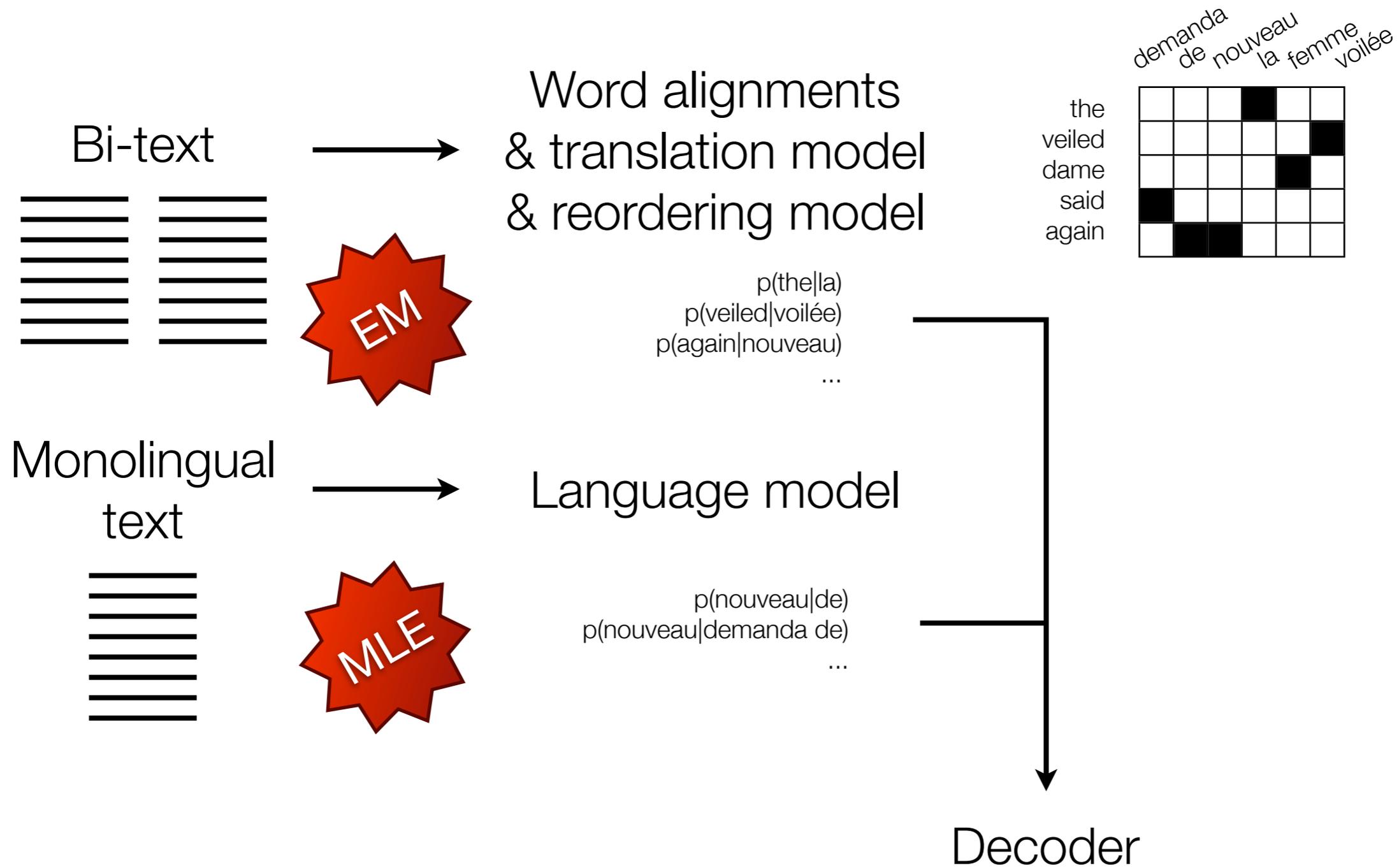
I CAN HAS SIMPLER TRAINING?





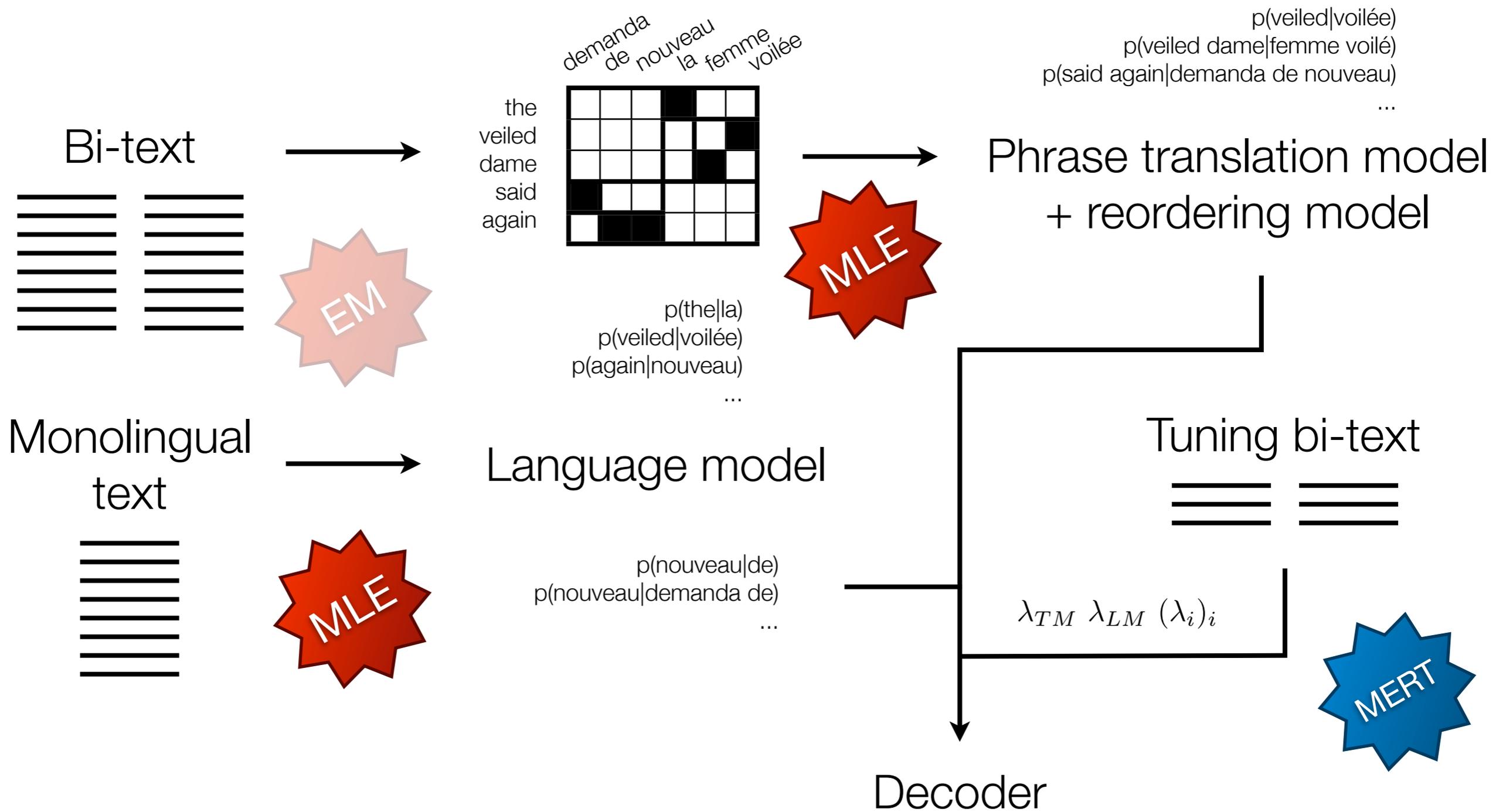
In principio erat
verbum: et ver-
bum erat apud
deum: et deus
erat verbum. Hoc
erat in principio
ad deum. Omnia per ipsum
sunt: et sine ipso factum
nihil: quod factum est. In

1-gram-based *noisy channel* MT



$$e^* = \operatorname{argmax}_e p(f|e)p(e)$$

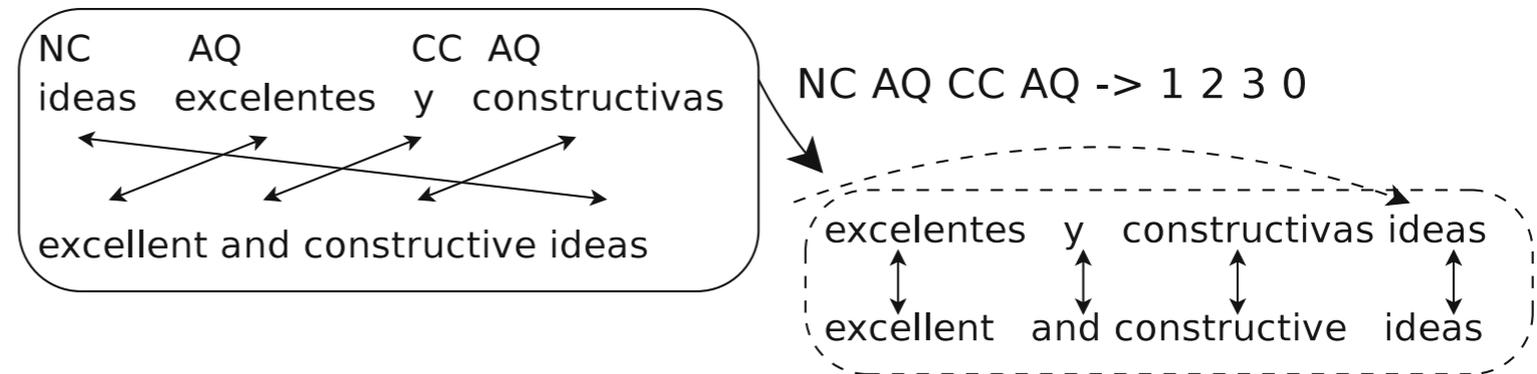
phrased-based *discriminative* MT



$$e^* = \operatorname{argmax}_e \lambda_{TM} \log p(f|e) + \lambda_{LM} \log p(e) + \sum_i \lambda_i f_i(f, e)$$

Reordering for n -gram models

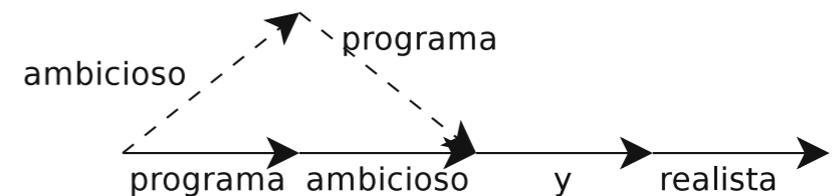
Training from
word-aligned
POS-tagged bi-text



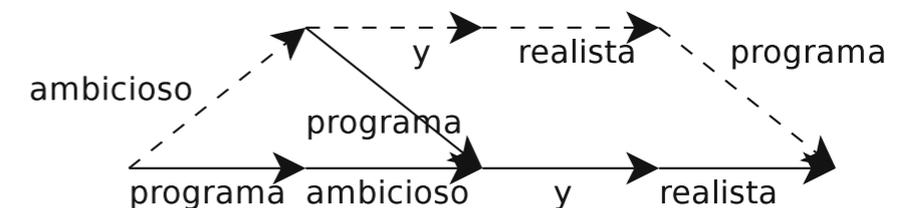
programa ambicioso y realista
NC AQ CC AQ



NC AQ -> 1 0

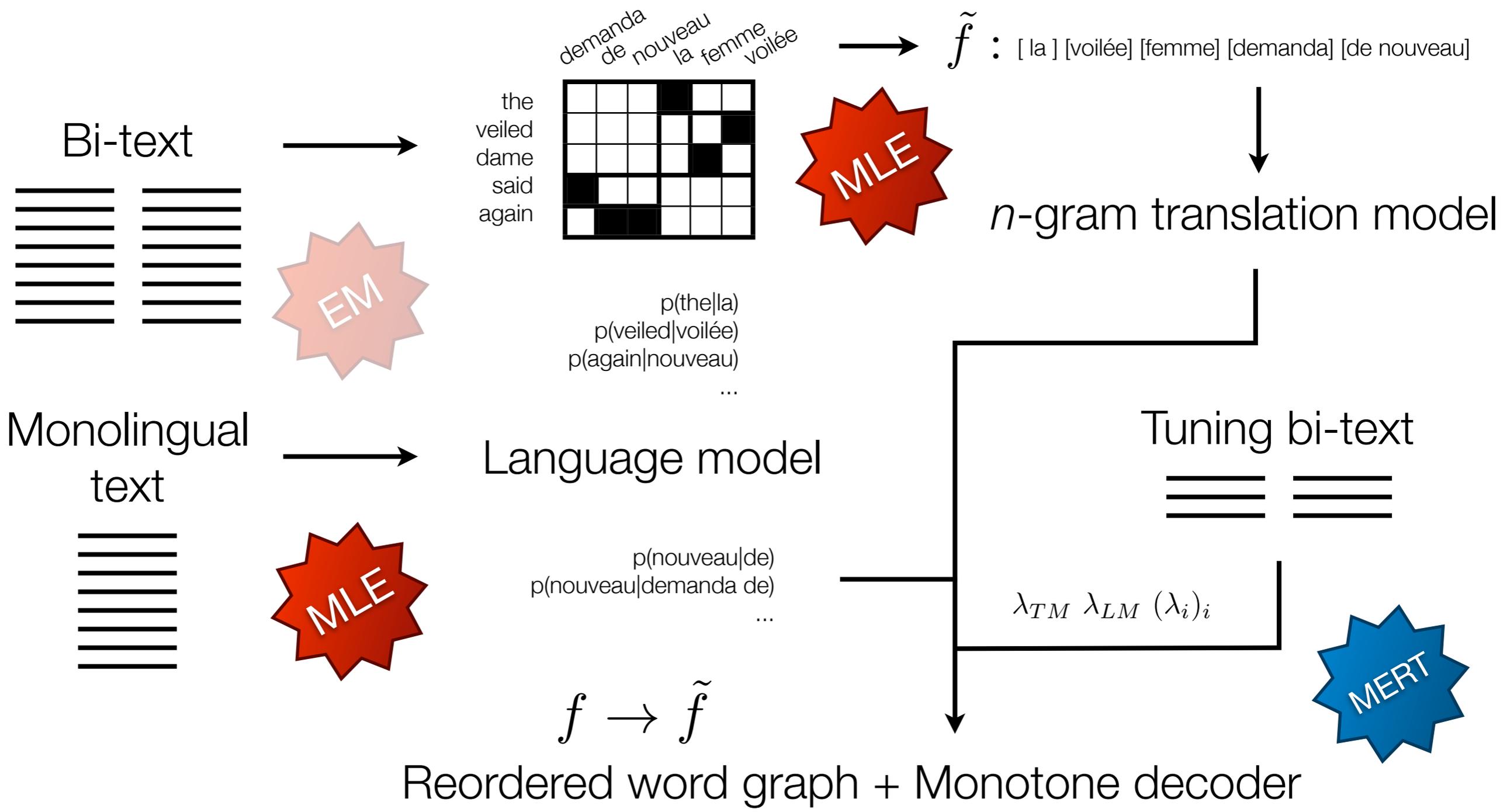


NC AQ CC AQ -> 1 2 3 0



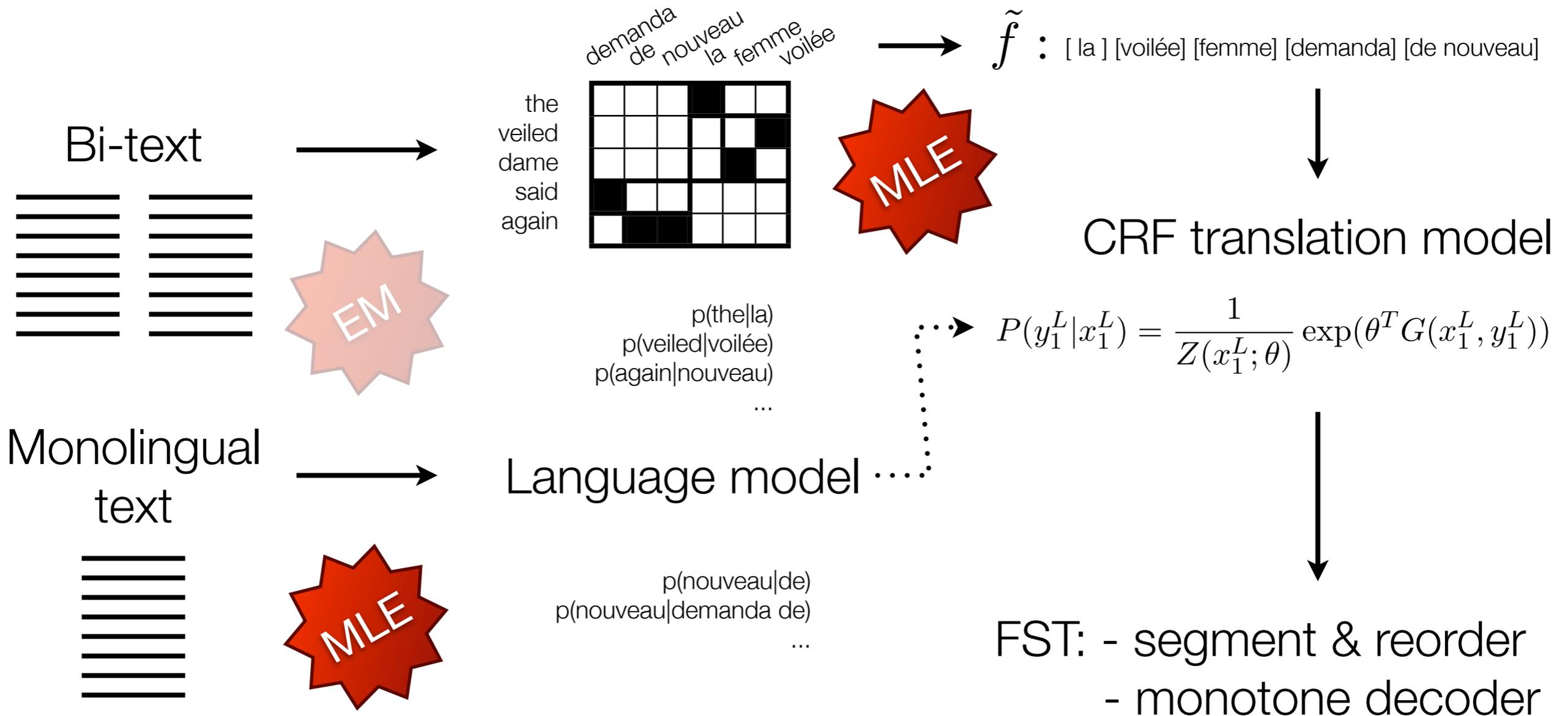
Input search graph

n-gram-based *discriminative* MT



$$e^* = \operatorname{argmax}_e \lambda_{TM} \log p(\tilde{f}, \tilde{e}) + \lambda_{LM} \log p(e) + \sum_i \lambda_i f_i(f, e)$$

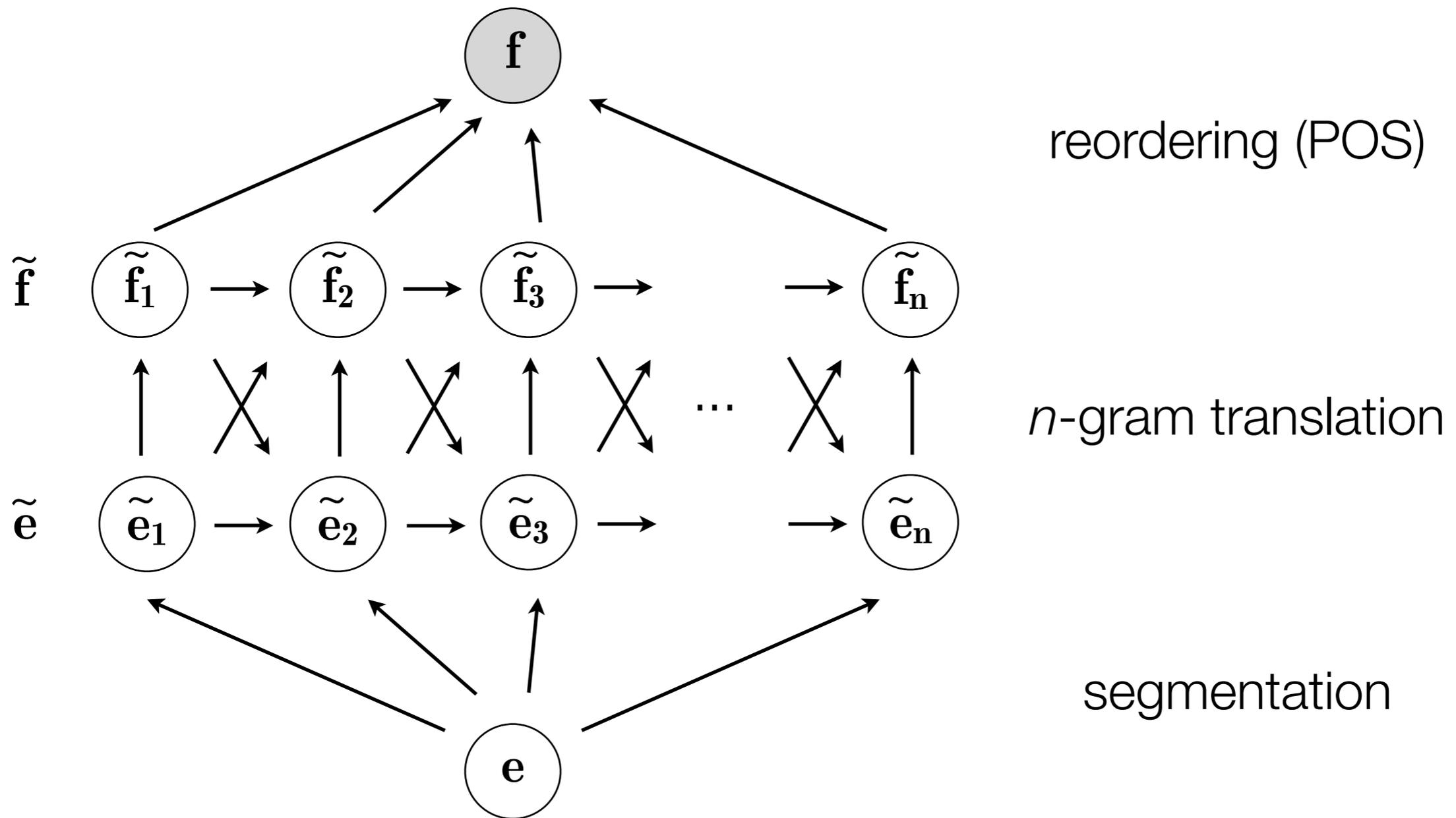
CRF-based MT



$$\tilde{\mathbf{e}}^* = \arg \max_{\tilde{\mathbf{e}}} P(\tilde{\mathbf{e}} | \mathbf{f})$$

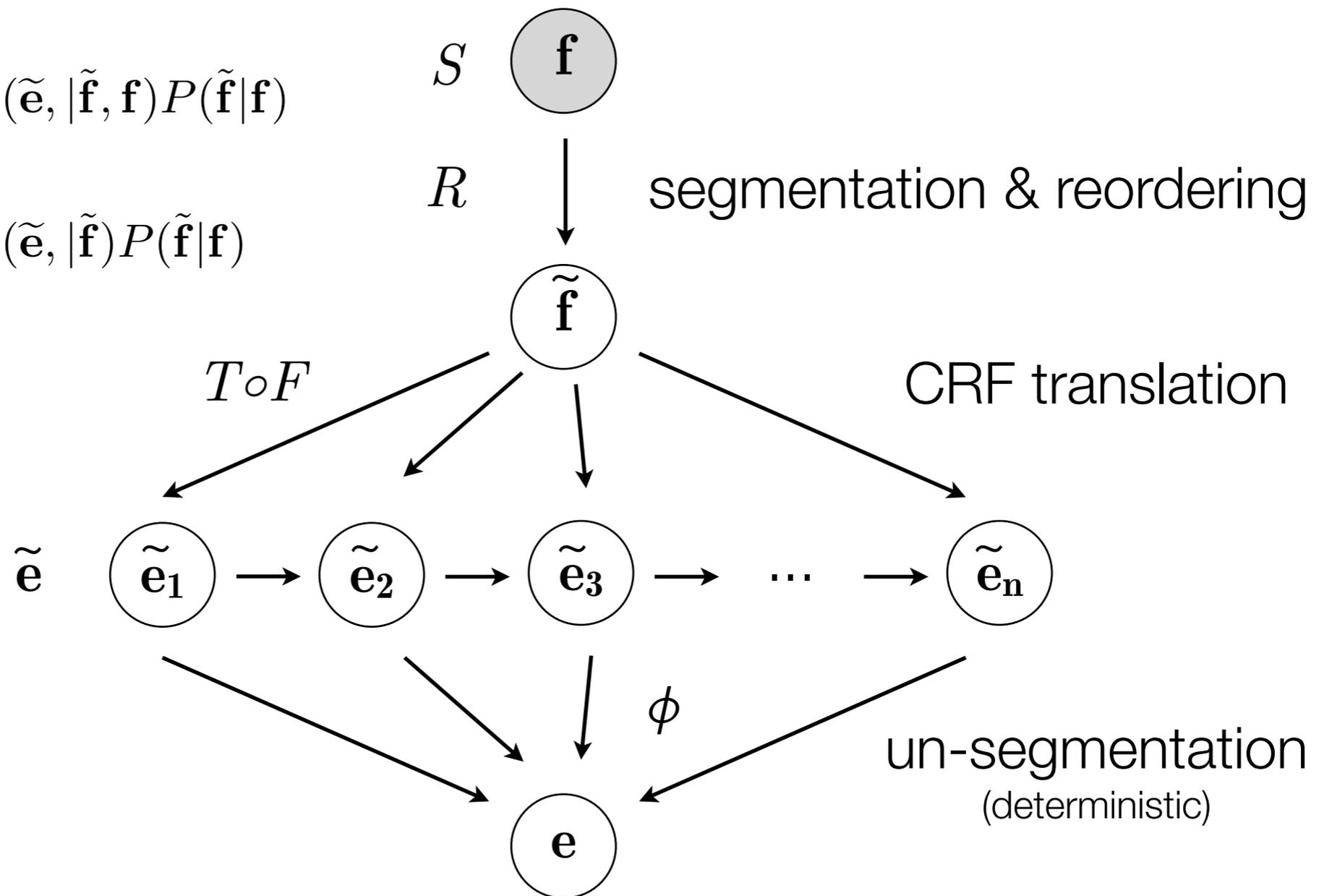
$$\approx \arg \max_{\tilde{\mathbf{f}} \in \mathcal{L}(\mathbf{f}), \tilde{\mathbf{e}}} P(\tilde{\mathbf{e}}, \tilde{\mathbf{f}} | \mathbf{f}) P(\tilde{\mathbf{f}} | \mathbf{f})$$

n -gram generative story

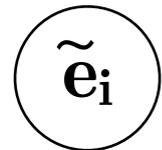


Discriminative story

$$\begin{aligned}
 P(\mathbf{e}|\mathbf{f}) &= \sum_{\tilde{\mathbf{f}}, \tilde{\mathbf{e}} | \phi(\tilde{\mathbf{e}})=\mathbf{e}} P(\tilde{\mathbf{e}}, \tilde{\mathbf{f}}|\mathbf{f}) \\
 &= \sum_{\tilde{\mathbf{f}}, \tilde{\mathbf{e}} | \phi(\tilde{\mathbf{e}})=\mathbf{e}} P(\tilde{\mathbf{e}}, |\tilde{\mathbf{f}}, \mathbf{f}) P(\tilde{\mathbf{f}}|\mathbf{f}) \\
 &= \sum_{\tilde{\mathbf{f}}, \tilde{\mathbf{e}} | \phi(\tilde{\mathbf{e}})=\mathbf{e}} P(\tilde{\mathbf{e}}, |\tilde{\mathbf{f}}) P(\tilde{\mathbf{f}}|\mathbf{f})
 \end{aligned}$$



CRF Features



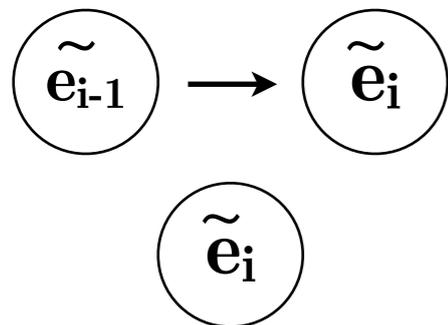
translation features [trs]

+context features [ctx]

+suffix/prefix features [ix]

+segmentation (length) features [seg]

+distortion (Δ +lex) features [ord]



target bigram features [trg]

+LM (non-local)

Negative results

	dev	test	# feat.
Moses (3g)	21.2	20.5	
<i>n</i> -gram (2g,3g)	20.6	20.2	755K
<i>n</i> -gram (3g,3g)	21.5	21.2	755K
CRF (trs,trg)	-	18.3	660K
CRF + ctx	-	18.8	1.5M
CRF + ix,ord,seg	-	19.1	1.5M
CRF + ix,ord,seg+3g	-	19.1	1.5M

+Language model does not help

Weaknesses:

- scoring of reordering & segmentation
- target LM

Oracle segmentation / reordering

	dev	test	# feat.
<i>decoding with optimal segmentation/reordering</i>			
CRF (trs,trg)	23.8	25.1	660K
CRF +ctx	24.1	25.4	1.5M
CRF +ix,ord,seg	24.3	25.6	1.5M
<i>decoding with optimal reordering</i>			
<i>n</i> -gram (2g,3g)	20.6	24.1	755K
<i>n</i> -gram (3g,3g)	21.5	25.2	755K
CRF trs,trg	-	22.8	660K
CRF +ctx	-	23.1	1.5M
CRF +ix,ord,seg	-	23.5	1.5M
<i>regular decoding</i>			
Moses (3g)	21.2	20.5	
<i>n</i> -gram (2g,3g)	20.6	20.2	755K
<i>n</i> -gram (3g,3g)	21.5	21.2	755K
CRF (trs,trg)	-	18.3	660K
CRF +ctx	-	18.8	1.5M
CRF +ix,ord,seg	-	19.1	1.5M
CRF +ix,ord,seg+3g	-	19.1	1.5M