

Models	WMT14		IWSLT14	Latency	Speedup
	En-De	De-En	De-En		
Autoregressive models					
Transformer [1] (beam size = 4)	27.41	31.29	33.26	$387ms^{\ddagger}$	$1.00\times$
Non-autoregressive models (Ours)					
NART-CRF	23.32 (4.09)	25.75 (5.54)	26.39 (6.87)	$35ms^{\ddagger}$	$11.1\times^{\ddagger}$
NART-CRF (rescoring 9)	26.04 (1.37)	28.88 (2.41)	29.21 (4.05)	$60ms^{\ddagger}$	$6.45\times^{\ddagger}$
NART-CRF (rescoring 19)	26.68 (0.73)	29.26 (2.03)	29.55 (3.71)	$87ms^{\ddagger}$	$4.45\times^{\ddagger}$
NART-DCRF	<b>23.44 (3.97)</b>	<b>27.22 (4.07)</b>	<b>27.44 (5.82)</b>	$37ms^{\ddagger}$	$10.4\times^{\ddagger}$
NART-DCRF (rescoring 9)	<b>26.07 (1.34)</b>	<b>29.68 (1.61)</b>	<b>29.99 (3.27)</b>	$63ms^{\ddagger}$	$6.14\times^{\ddagger}$
NART-DCRF (rescoring 19)	<b>26.80 (0.61)</b>	<b>30.04 (1.25)</b>	<b>30.36 (2.90)</b>	$88ms^{\ddagger}$	$4.39\times^{\ddagger}$