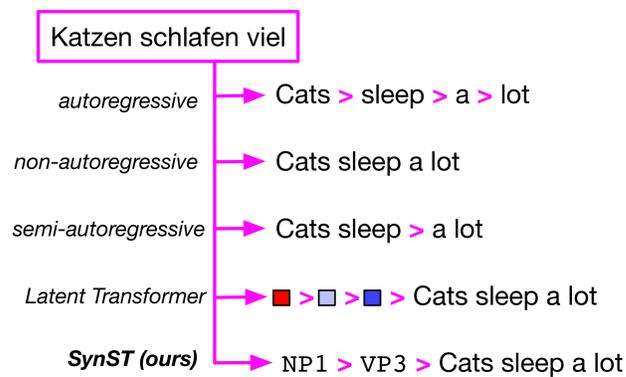


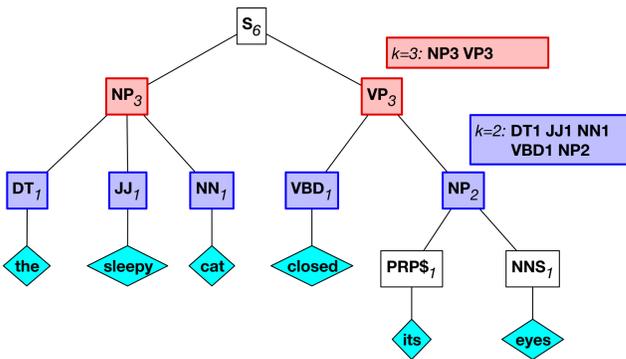
We present the syntactically supervised Transformer (**SynST**), which achieves faster translation and higher BLEU than competing non-autoregressive neural machine translation models.

## SYNST VS. EXISTING SYSTEMS



Each > denotes a single decode step. Fewer decode steps results in faster translation, often at the expense of quality.

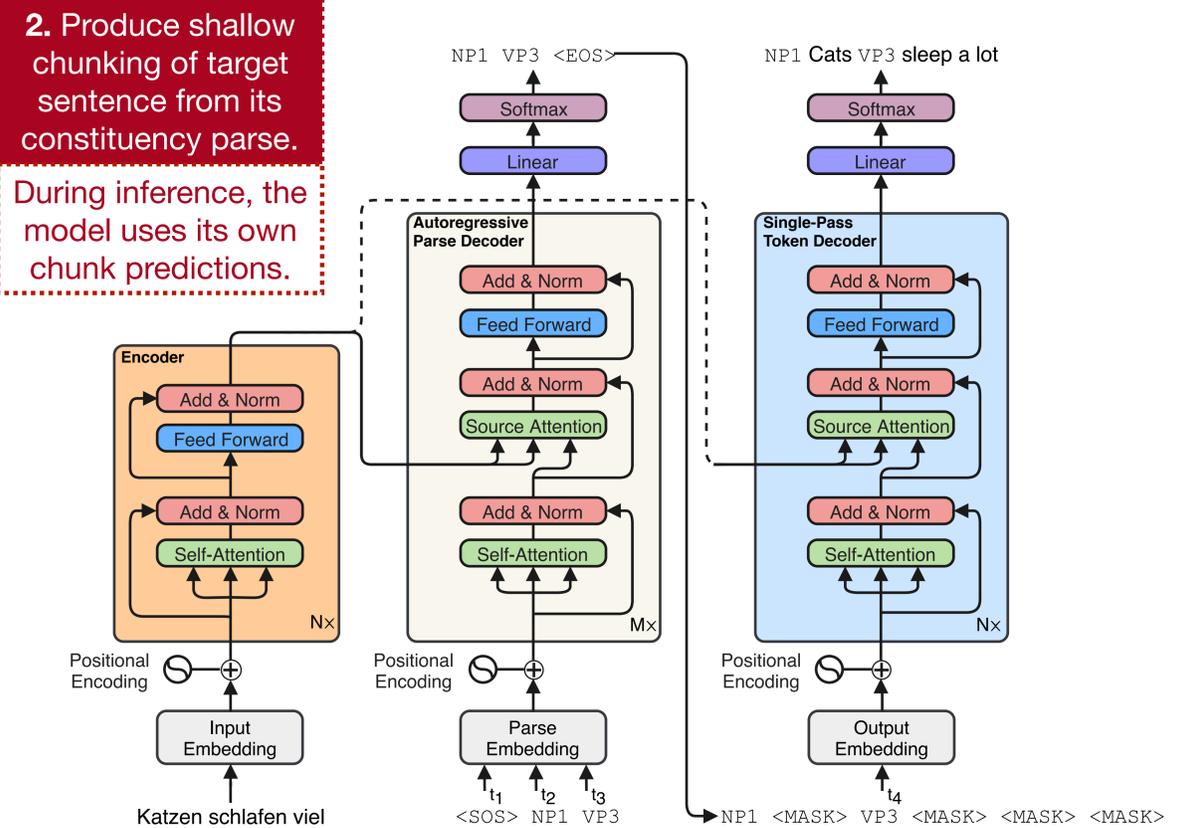
## TARGET PARSE CHUNKING



During an in-order traversal, if the subtree rooted at a visited node spans  $\leq k$  tokens, append it to our **chunk sequence**.

## TRAINING SYNST

1. Encode source sentence.
2. Produce shallow chunking of target sentence from its constituency parse.
3. Autoregressively predict target chunks.
4. Non-autoregressively predict target tokens.



## CONTROLLED EXPERIMENTS

Model	WMT En-De		WMT De-En		IWSLT En-De		WMT En-Fr	
	BLEU	Speedup	BLEU	Speedup	BLEU	Speedup	BLEU	Speedup
Vanilla Transformer								
Beam Size = 4	26.87	1.00x	30.73	1.00x	30.00	1.00x	40.22	1.00x
Beam Size = 1	25.82	1.15x	29.83	1.14x	28.66	1.16x	39.41	1.18x
Semi-Autoregressive Transformer								
k = 2	22.81	2.05x	26.78	2.04x	25.48	2.03x	36.62	2.14x
k = 4	16.44	3.61x	21.27	3.58x	20.25	3.45x	28.07	3.34x
k = 6	12.55	<b>4.86x</b>	15.23	4.27x	14.02	<b>4.39x</b>	24.63	4.77x
Latent Transformer*								
* As reported in (Kaiser et al. 2018)	19.8	3.89x	-	-	-	-	-	-
Syntactically Supervised Transformer								
k=6	20.74	<b>4.86x</b>	25.50	<b>5.06x</b>	23.82	3.78x	33.47	<b>5.32x</b>

## ANALYSIS ON IWSLT DEV SET

### Constituent identity is crucial for quality

Only predicting constituent length (1 > 3) rather than type & length (NP1 > VP3), causes a **BLEU drop** from **23.8** to **8.2**.

### How much does SynST rely on syntax?

Source: *Katzen schlafen viel*    Target: *Cats sleep a lot*  
 Predicted Parse: NP1 > VP2    Prediction: *Cats sleep lots*  
 Gold Parse: NP1 > VP3    Parsed Prediction: NP1 > VP2

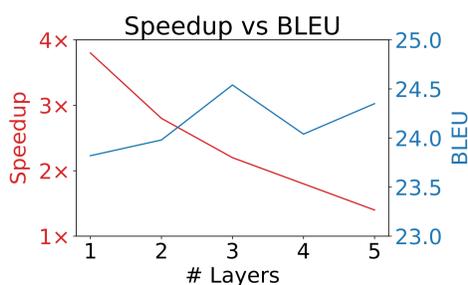
	Predicted Parse	Predicted Parse	Parsed Prediction	Parsed Prediction
	vs Gold Parse (trained separately)	vs Gold Parse	vs Gold Parse	vs Predicted Parse
F1	65.48	69.64	79.16	89.90
Exact Match	4.23%	5.24%	5.94%	43.10%

Parsed prediction closely matches predicted parse, though there exists room for improvement for parse prediction.

### Ground-truth syntax yields huge improvements

Conditioning on the ground-truth chunk sequence during inference dramatically **improves BLEU** from **23.8** to **41.5**, yielding an upper bound for our approach.

### SynST's bottleneck is its parse decoder



A one-layer parse decoder is **~3x faster** than a 5-layer version, with only a **~0.5 BLEU drop**.

### Future work: dynamic vs fixed k

Randomly sampling possible chunk sequences during training by varying  $k$  leads to a **large BLEU improvement** (+1.5) with minimal impact to speedup (drop from 3.8x to 3.1x). Improving parse prediction is an avenue for future research.