



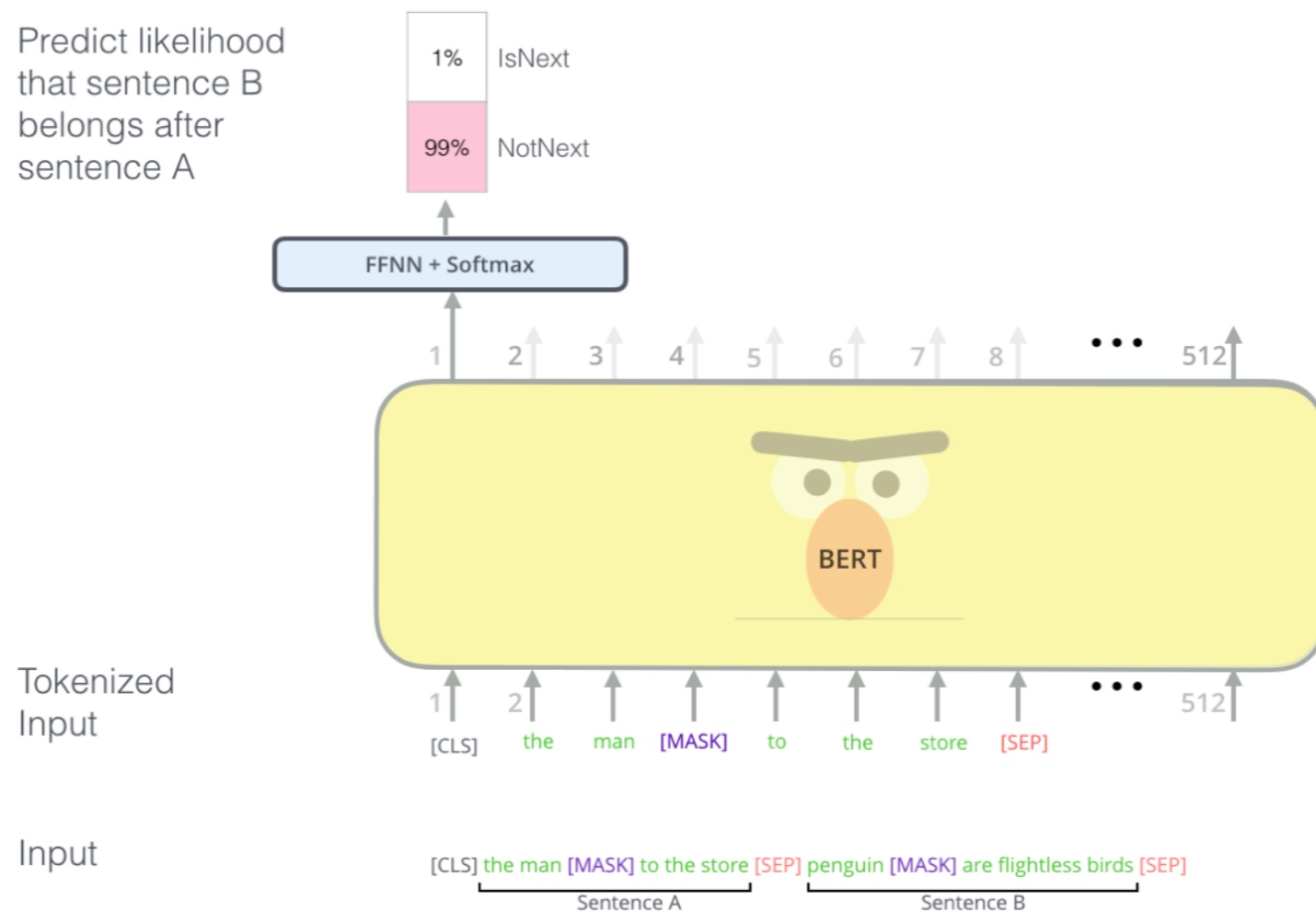
COS 598C Advanced Topics in Computer Science:  
Deep Learning for Natural Language Processing

# Pre-training and Fine-tuning

Winter 2020

# A note on next sentence prediction

- NSP could hurt the MLM training objective.
- **Recommended reading:** SpanBERT: Improving Pre-training by Representing and Predicting Spans



# Fine-tuning vs Feature-based

- **Recommended reading:** To Tune or Not to Tune? Adapting Pretrained Representations to Diverse Tasks

Pretraining	Adaptation	NER	SA	Nat. lang. inference		Semantic textual similarity		
		CoNLL 2003	SST-2	MNLI	SICK-E	SICK-R	MRPC	STS-B
Skip-thoughts	❄️	-	81.8	62.9	-	86.6	75.8	71.8
ELMo	❄️	91.7	<b>91.8</b>	<b>79.6</b>	<b>86.3</b>	<b>86.1</b>	<b>76.0</b>	<b>75.9</b>
	🔥	<b>91.9</b>	91.2	76.4	83.3	83.3	74.7	75.5
	$\Delta = \text{🔥} - \text{❄️}$	0.2	-0.6	<b>-3.2</b>	<b>-3.3</b>	<b>-2.8</b>	<b>-1.3</b>	-0.4
BERT-base	❄️	92.2	93.0	<b>84.6</b>	84.8	86.4	78.1	82.9
	🔥	<b>92.4</b>	<b>93.5</b>	<b>84.6</b>	<b>85.8</b>	<b>88.7</b>	<b>84.8</b>	<b>87.1</b>
	$\Delta = \text{🔥} - \text{❄️}$	0.2	0.5	0.0	1.0	<b>2.3</b>	<b>6.7</b>	<b>4.2</b>

# Fine-tuning vs Feature-based

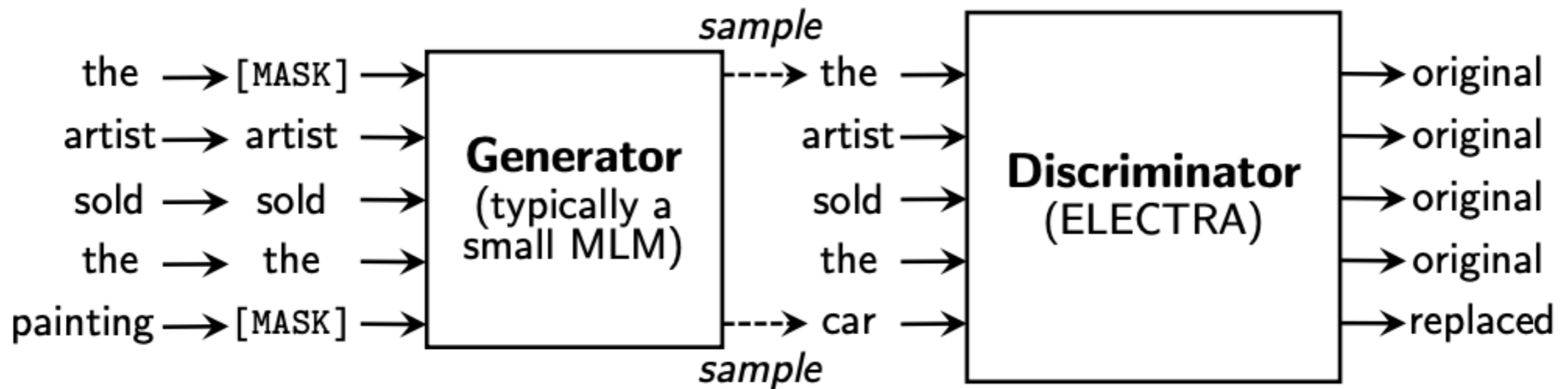
- **Recommended reading:** Incorporating BERT into Neural Machine Translation

Table 1: Preliminary explorations on IWSLT'14 English→German translation.

Algorithm	BLEU score
Standard Transformer	28.57
Use BERT to initialize the encoder of NMT	27.14
Use XLM to initialize the encoder of NMT	28.22
Use XLM to initialize the decoder of NMT	26.13
Use XLM to initialize both the encoder and decoder of NMT	28.99
Leveraging the output of BERT as embeddings	29.67

# How to fix the 15% masking?

- **Recommended reading:** ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators



# How to fix the 15% masking?

- **Next lecture:** XLNet: Generalized Autoregressive Pretraining for Language Understanding

