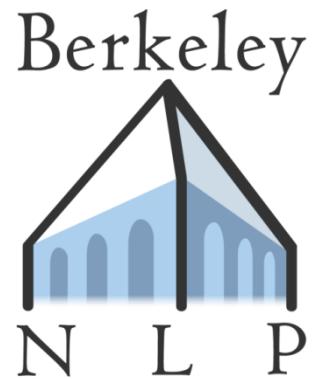


Natural Language Processing



LLMs: Training



BAIR NLP Workshop

- Tomorrow, Friday, October 18, in BWW and virtual
- Highlights:
 - **LLMs and Cognitive Systems**
 - LLM agents
 - Creativity in humans and LLMs
 - Interpretability
 - Language in interaction
- Agenda and link to RSVP: <https://docs.google.com/document/d/1-WJfTMfYnCwlylsJxXoRzjRJNlJ1MyiFkcFLQNYPviE/>



Recap: What is a language model?

- Language models assign a probability to a sequence of words
- We can decompose this probability using the chain rule
- We can autoregressively generate sequences from the language model by sampling from its token-level probability
- We can condition on our language distribution on something else

$$p(\bar{y})$$

$$p(\bar{y}) = \prod_{i=1}^T p(y_i | y_{0:i-1})$$

$$p(y_i | y_{0:i-1})$$

$$p(y_i | y_{0:i-1}; \bar{x})$$



How to Train a Base LM

1. Get some training data
2. Preprocess it (tokenize it)
3. Choose your architecture
4. Optimize a language modeling objective
5. Run inference!

Step 1: Get Training Data



Training Data

- Transformer models are very data-hungry
- Solution: just scrape the web
- **CommonCrawl**: publicly available web scrape collected since 2007 containing 250B webpages, comprising 82% of tokens used to train GPT-3





Data Sources

- Domain-specific webpages:
 - Code and mathematics: Github, StackOverflow
 - Academic and scientific work: arXiv, bioRxiv, PubMed
 - Books: Project Gutenberg
 - General knowledge: Wikipedia
- Domain-general sources:
 - Social media (reddit, Twitter)
 - News sites



Data Sources

The History of Large Language Models



- 1998 CPAT-Tree-Based Language Models with an Application for Text Verification in Chinese. ROCLing 1998. First use of LLM trigram I know of; 200M word corpus
- 2000 A Neural Probabilistic Language Model. Bengio, Ducharme & Vincent NIPS 2000 First neural language model built on 32 million token corpus, 31K vocab
- 2007 Large Language Models in Machine Translation. Brants, Popat, Xu, Och and **Dean**. EMNLP 2007. **2 trillion token corpus** n-gram model of up to 5-grams
- 2018 **GPT** (Radford, Narasimhan, Salimans & Sutskever) and **BERT** (Devlin, Chang, Lee & Toutanova). 3.3 billion token corpus
- 2020- 100+ billion parameter neural language models trained on > 1 trillion tokens: GPT-3, GPT-4, PaLM 2, Llama 3, Nemotron-4,



Web Scraping

1. Seed webcrawler with initial URLs
2. Identify new URLs via outlinks
3. Download HTML representation of webpage
4. Scrape HTML for raw text
5. Postprocess texts

$$\mathcal{D} = \left\{ \bar{d}^{(i)} \right\}_{i=1}^N \quad \bar{d} = \langle b_0, \dots, b_M \rangle$$



Web Data is Noisy

- Deduplication
- Remove junk / nonsense text that's very unlikely according to a simple n-gram language model
- Remove uninteresting pages with few inlinks
- Remove non-English data with external classifiers



Web Data is Unfiltered

- Personally identifiable information (PII) or other personal information
- Adult content
- Explicit hate speech, disinformation
- Copyrighted data
- Test data from NLP benchmarks...



Downstream Effects

Stable Diffusion produces copyright and trademarked images

Original:



Generated:



Codex generates code with non-permissive licenses

```
3685 CBlockIndex * InsertBlockIndex(uint256 hash)
3686 {
3687     if (hash.IsNull())
3688         return NULL;
3689
3690     // Return existing
3691     BlockMap::iterator mi = mapBlockIndex.find(hash);
3692     if (mi != mapBlockIndex.end())
3693         return (*mi).second;
3694
3695     CBlockIndex* pindexNew = new CBlockIndex();
3696     if (!pindexNew)
3697         throw runtime_error("LoadBlockIndex(): new CBlockIndex failed");
3698     mi = mapBlockIndex.insert(make_pair(hash, pindexNew)).first;
3699     pindexNew->pindexBlock = &(*mi).first;
3700
3701     return pindexNew;
3702 }
```

Stable Diffusion generates real individuals





Social Impacts of Webscraping

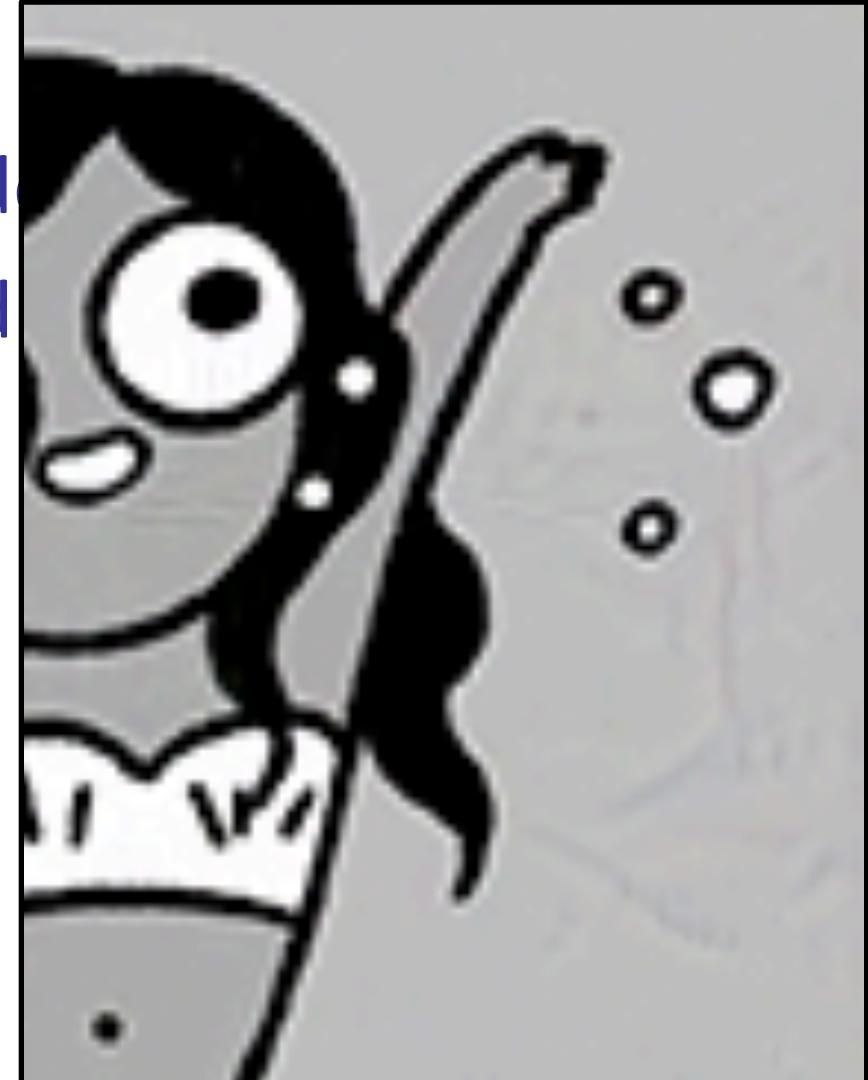
- Trained language models encode:
 - Biases explicitly or implicitly encoded
 - Personal information about individuals
 - Copyrighted data



Karla Ortiz



Sarah Andersen



Glaze, Shan et al. 2023, USENIX

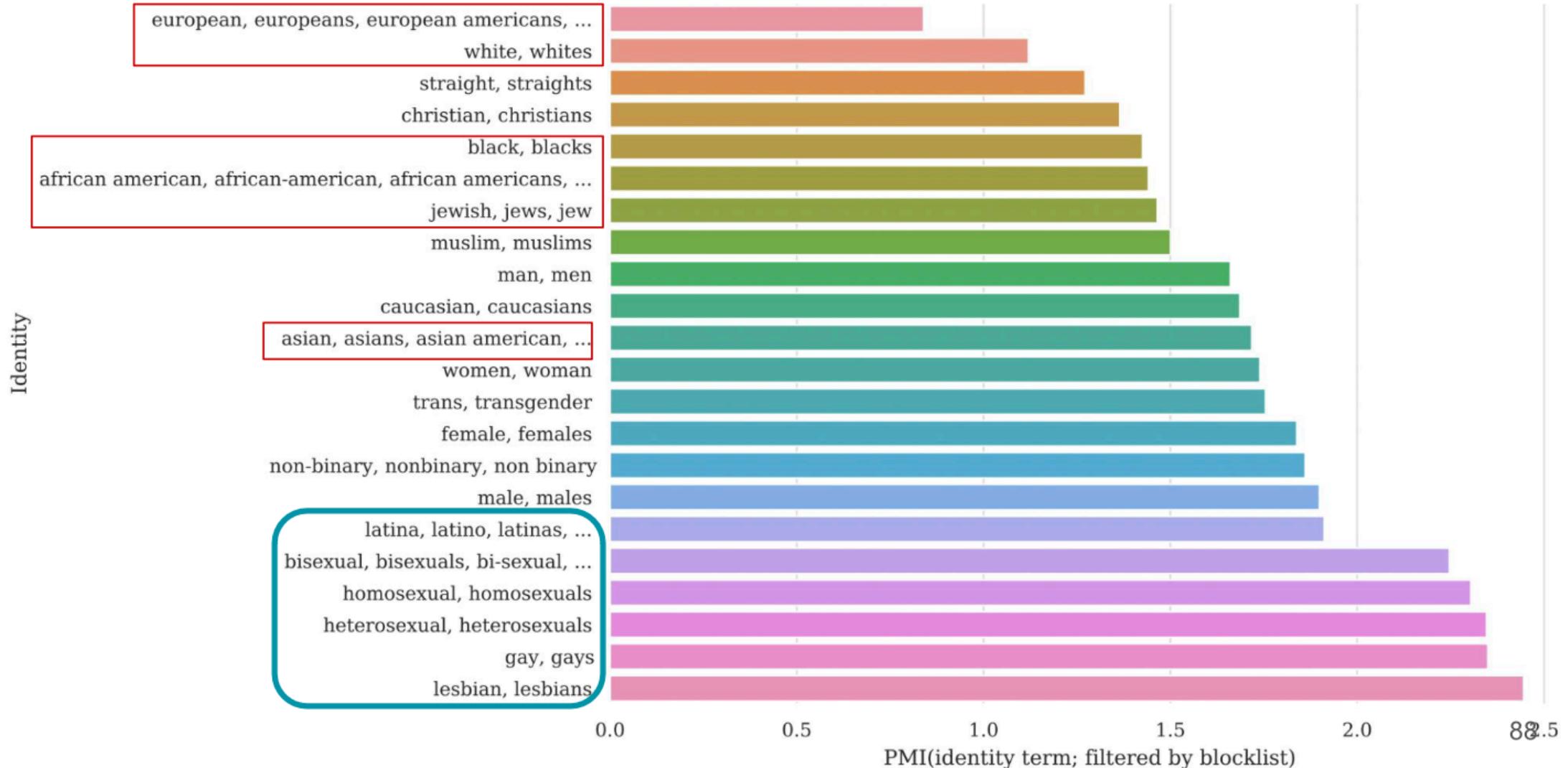


Tradeoffs in Filtering

- Personally identifiable information (PII) or other personal information → Phone numbers of public companies' customer service lines?
- Adult content → Very culturally dependent
- Explicit hate speech, disinformation → What might appear to be hateful or toxic speech is context-dependent

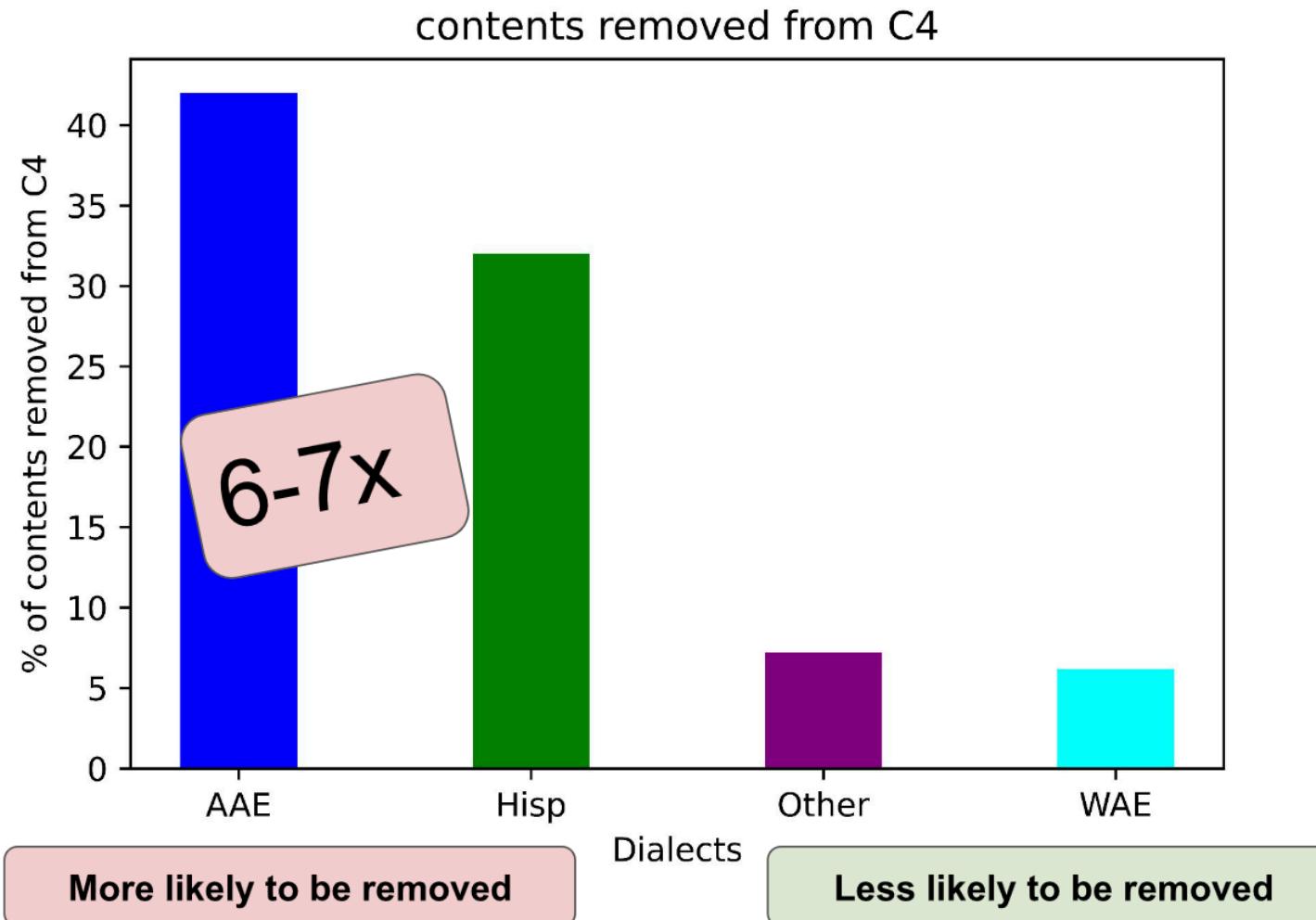


Tradeoffs in Filtering





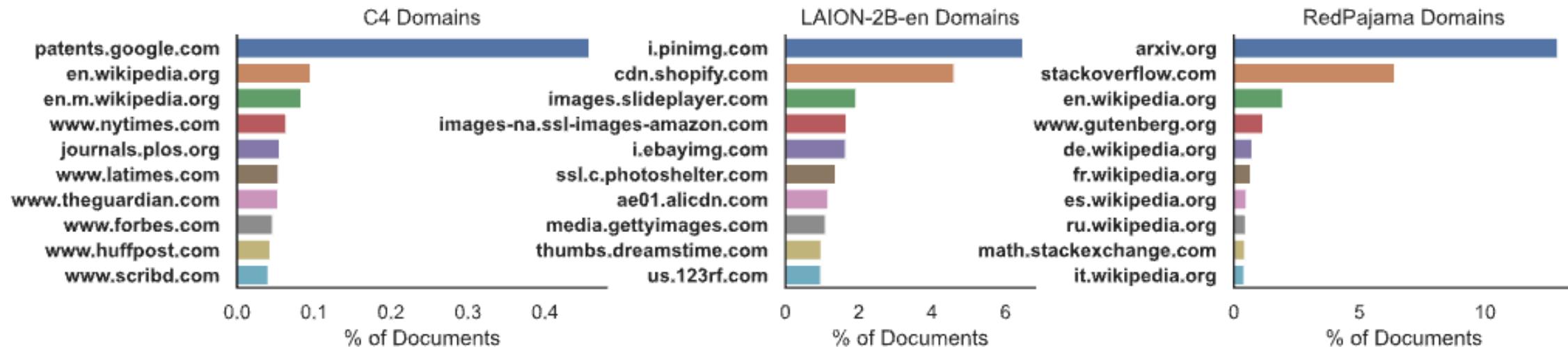
Tradeoffs in Filtering





Pretraining Corpora

Dataset	Origin	Model	Size (GB)	# Documents	# Tokens	max(# Tokens)	min(# Tokens)
OpenWebText	Gokaslan & Cohen (2019)	GPT-2*	41.2	8,005,939	7,767,705,349	95,139	128
C4	Raffel et al (2020)	T5	838.7	364,868,892	153,607,833,664	101,898	5
mC4-en	Chung et al (2021)	umT5	14,694.0	3,928,733,374	2,703,077,876,916	181,949	1
OSCAR	Abadii et al (2021)	BLOOM*	3,327.3	431,584,362	475,992,028,559	1,048,409	1
The Pile	Gao et al (2020)	GPT-J/Neo & Pythia	1,369.0	210,607,728	285,794,281,816	28,121,329	0
RedPajama	Together Compute (2021)	LLaMA*	5,602.0	930,453,833	1,023,865,191,958	28,121,329	0
S2ORC	Lo et al (2020)	SciBERT*	692.7	11,241,499	59,863,121,791	376,681	1
peS2o	Soldaini & Lai (2021)	-	504.3	8,242,162	44,024,690,229	97,043	154
LAION-2B-en	Schuhmann et al (2021)	Stable Diffusion*	570.2	2,319,907,827	29,643,340,153	131,077	0
The Stack	Kočetkov et al (2021)	StarCoder*	7,830.8	544,750,672	1,525,618,728,620	26,298,134	0

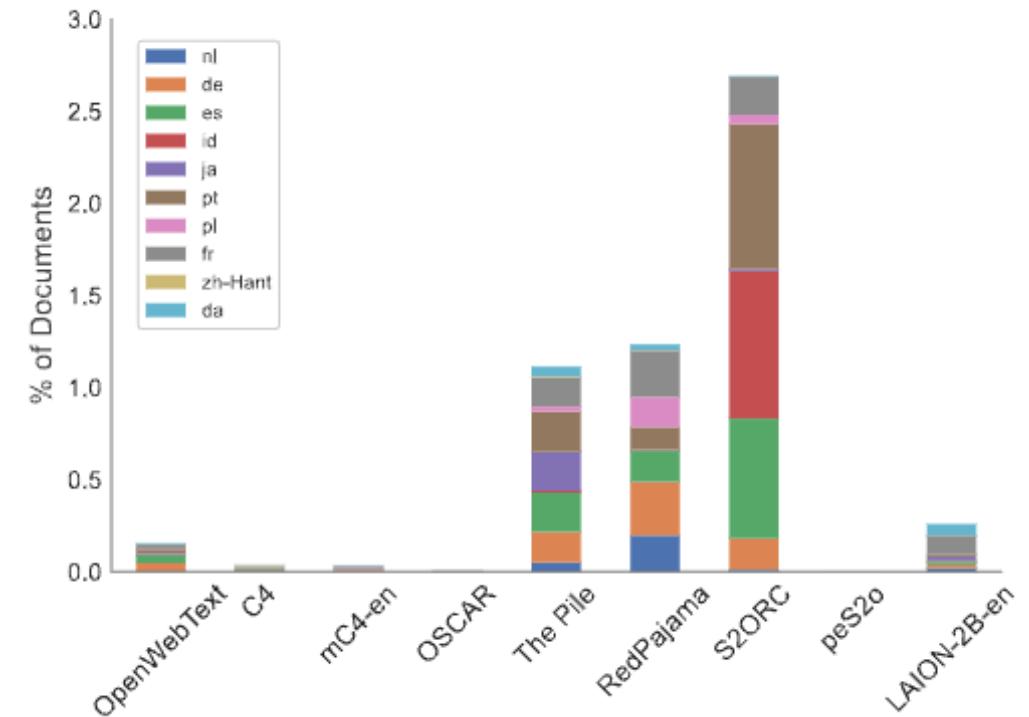




Pretraining Corpora

Table 5: Extrapolated PII frequencies. Count is the extrapolated frequency in a corpus and *Prec.* is our identification precision, estimated by manual analysis.

	Email Addresses		Phone Numbers		IP Addresses	
	Count	Prec.	Count	Prec.	Count	Prec.
OpenWebText	364K	99	533K	87	70K	54
OSCAR	62.8M	100	107M	91	3.2M	43
C4	7.6M	99	19.7M	92	796K	56
mC4-en	201M	92	4B	66	97.8M	44
The Pile	19.8M	43	38M	65	4M	48
RedPajama	35.2M	100	70.2M	94	1.1M	30
S2ORC	630K	100	1.4M	100	0K	0
peS2o	418K	97	227K	31	0K	0
LAION-2B-en	636K	94	1M	7	0K	0
The Stack	4.3M	53	45.4M	9	4.4M	55



Step 2: Tokenization



Tokenization

$$\mathcal{D} = \left\{ \bar{d}^{(i)} \right\}_{i=1}^N$$

“They currently play their home games at Acrisure Stadium.”



“They” “currently” “play” “their” “home”
“game” “#s” “at” “Acrisure” “Stadium” “.”

$$\bar{d} = \langle b_0, \dots, b_M \rangle$$

$$\bar{d} = \langle x_0, \dots, x_{M'} \rangle$$

$$\forall x, x \in \mathcal{V}$$



Tokenization

- Maps from byte sequences to sequences of tokens, where each token is part of a set vocabulary

$$\begin{aligned}\bar{d} &= \langle b_0, \dots, b_M \rangle \\ \bar{d} &= \langle x_0, \dots, x_{M'} \rangle \\ \forall x, x &\in \mathcal{V}\end{aligned}$$

↓



Tokenization

- Approach: simple heuristics (split on spaces, handle punctuation gracefully)

“They currently play their home games at Acrisure Stadium.”



*“They” “currently” “play” “their” “home”
“games” “at” “Acrisure” “Stadium” “.”*

$$\begin{aligned}\bar{d} &= \langle b_0, \dots, b_M \rangle \\ \downarrow \\ \bar{d} &= \langle x_0, \dots, x_{M'} \rangle \\ \forall x, x &\in \mathcal{V}\end{aligned}$$

Problem: requires defining heuristics, including for edge cases

Problem: heuristics are not generalizable to all languages

เราทุกคนเกิดมาอย่างอิสระ เราทุกคนมีความคิดและความเชื่อใจเป็นของเรารอง เราทุกคนควรได้รับการปฏิบัติในทางเดียวกัน.



Tokenization

Turkish	English
ev	(the) house
evler	(the) houses
evin	your (sing.) house
eviniz	your (pl./formal) house
evim	my house
evimde	at my house
evlerinizin	of your houses
evlerinizden	from your houses
evlerinizdendi	(he/she/it) was from your houses
evlerinizdenmiş	(he/she/it) was (apparently/said to be) from your houses
Evinizdeyim.	I am at your house.
Evinizdeymişim.	I was (apparently) at your house.
Evinizde miyim?	Am I at your house?

it on

$$\bar{d} = \langle b_0, \dots, b_M \rangle$$

Case	Ending	Examples		Meaning
		köy "village"	ağaç "tree"	
Nominative	Ø (none)	köy	ağaç	(the) village/tree
Accusative	-i ⁴	köyü	ağaçı	the village/tree
Genitive	-in ⁴	köyün	ağaçın	the village's/tree's of the village/tree
Dative	-e ²	köye	ağaca	to the village/tree
Locative	-de ²	köyde	ağaçta	in/on/at the village/tree
Ablative	-den ²	köyden	ağaçtan	from the village/tree
Instrumental	-le ²	köyle	ağaçla	with the village/tree

Problem: many words never appear
in the training data



Character- / Byte-Level Models

- Approach: vocabulary is simply all possible Unicode characters that might appear

!	"	#	\$	%	&	'	()	*	+	,
33	34	35	36	37	38	39	40	41	42	43	44
-	.	/	Ø	1	2	3	4	5	6	7	8
45	46	47	48	49	50	51	52	53	54	55	56
9	:	;	<	=	>	?	@	A	B	C	D
57	58	59	60	61	62	63	64	65	66	67	68
E	F	G	H	I	J	K	L	M	N	O	P
69	70	71	72	73	74	75	76	77	78	79	80
Q	R	S	T	U	V	W	X	Y	Z	[\
81	82	83	84	85	86	87	88	89	90	91	92
]	^	-	`	a	b	c	d	e	f	g	h
93	94	95	96	97	98	99	100	101	102	103	104
i	j	k	l	m	n	o	p	q	r	s	t
105	106	107	108	109	110	111	112	113	114	115	116
u	v	w	x	y	z	{		}	~	€	
117	118	119	120	121	122	123	124	125	126	127	128
,	f	"	...	†	#	^	%	š	<	œ	
129	130	131	132	133	134	135	136	137	138	139	140
ž											
141	142	143	144	145	146	147	148	149	150	151	152
™	š	>	œ	ž	Ý	i	¢	£	¤		
153	154	155	156	157	158	159	160	161	162	163	164
Y	I	S	ü	ö	ä						

$$\bar{d} = \langle b_0, \dots, b_M \rangle$$

$$\bar{d} = \langle x_0, \dots, x_{M'} \rangle$$

$$\forall x, x \in \mathcal{V}$$

Problem: representations of each character are not meaningful

Problem: model also needs to learn how to compose words from characters

Problem: input sequences become very long



Tokenization

- Approach: subword tokenization, where frequent words are kept whole and infrequent words are broken into parts

“They currently play their home games at Acrisure Stadium.”



'__They', '__currently', '__play',
['__their', '__home',
['__games', '__at', '__A', 'cris',
'ure', '__Stadium', '__.']

$$\begin{aligned}\bar{d} &= \langle b_0, \dots, b_M \rangle \\ \downarrow \\ \bar{d} &= \langle x_0, \dots, x_{M'} \rangle \\ \forall x, x &\in \mathcal{V}\end{aligned}$$

Adamla tanıştım indicator of subject
indirect object instrumental case suffix verb stem
past tense suffix

I met with the man

Adamın kitabı possessive ending
possessor genitive suffix possessed noun
Man's book



Byte Pair Encoding

- Gradually constructs vocabulary given a target size
- Starts with a base vocabulary consisting of all characters in the training data
- Iteratively constructs vocabulary:
 - Tokenizes all training documents given the current vocabulary
 - Adds the most common bigram to the vocabulary
- Terminates when target vocabulary size is reached

$$\begin{aligned}\bar{d} &= \langle b_0, \dots, b_M \rangle \\ &\quad \downarrow \\ \bar{d} &= \langle x_0, \dots, x_{M'} \rangle \\ &\quad \forall x, x \in \mathcal{V}\end{aligned}$$



Byte Pair Encoding

Documents + frequencies: ("hug", 10), ("pug", 5), ("pun", 12), ("bun", 4), ("hugs", 5)

("h", "u", "g", "p", "n", "b", "s")



Byte Pair Encoding

Documents + frequencies: ("hug", 10), ("pug", 5), ("pun", 12), ("bun", 4), ("hugs", 5)

("h", "u", "g", "p", "n", "b", "s") → ("h" "u" "g", 10), ("p" "u" "g", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "u" "g" "s", 5)



Byte Pair Encoding

Documents + frequencies: ("hug", 10), ("pug", 5), ("pun", 12), ("bun", 4), ("hugs", 5)

("h", "u", "g", "p", "n", "b", "s") → ("h" "u" "g", 10), ("p" "u" "g", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "u" "g" "s", 5)

("h", "u", "g", "p", "n", "b", "s", "ug") → ("h" "ug", 10), ("p" "ug", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "ug" "s", 5)



Byte Pair Encoding

Documents + frequencies: ("hug", 10), ("pug", 5), ("pun", 12), ("bun", 4), ("hugs", 5)

("h", "u", "g", "p", "n", "b", "s") → ("h" "u" "g", 10), ("p" "u" "g", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "u" "g" "s", 5)

("h", "u", "g", "p", "n", "b", "s", "ug") → ("h" "ug", 10), ("p" "ug", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "ug" "s", 5)

("h", "u", "g", "p", "n", "b", "s", "ug", "un") → ("h" "ug", 10), ("p" "ug", 5), ("p" "un", 12), ("b" "un", 4), ("h" "ug" "s", 5)



Byte Pair Encoding

Documents + frequencies: ("hug", 10), ("pug", 5), ("pun", 12), ("bun", 4), ("hugs", 5)

("h", "u", "g", "p", "n", "b", "s") → ("h" "u" "g", 10), ("p" "u" "g", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "u" "g" "s", 5)

("h", "u", "g", "p", "n", "b", "s", "ug") → ("h" "ug", 10), ("p" "ug", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "ug" "s", 5)

("h", "u", "g", "p", "n", "b", "s", "ug", "un") → ("h" "ug", 10), ("p" "ug", 5), ("p" "un", 12), ("b" "un", 4), ("h" "ug" "s", 5)

("h", "u", "g", "p", "n", "b", "s", "ug", "un", "hug") → ("hug", 10), ("p" "ug", 5), ("p" "un", 12), ("b" "un", 4), ("hug" "s", 5)



Byte Pair Encoding

Documents + frequencies: ("hug", 10), ("pug", 5), ("pun", 12), ("bun", 4), ("hugs", 5)

("h", "u", "g", "p", "n", "b", "s") → ("h" "u" "g", 10), ("p" "u" "g", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "u" "g" "s", 5)

("h", "u", "g", "p", "n", "b", "s", "ug") → ("h" "ug", 10), ("p" "ug", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "ug" "s", 5)

("h", "u", "g", "p", "n", "b", "s", "ug", "un") → ("h" "ug", 10), ("p" "ug", 5), ("p" "un", 12), ("b" "un", 4), ("h" "ug" "s", 5)

("h", "u", "g", "p", "n", "b", "s", "ug", "un", "hug") → ("hug", 10), ("p" "ug", 5), ("p" "un", 12), ("b" "un", 4), ("hug" "s", 5)

New word: “puns”

“p” “u” “n” “s” → “p” “un” “s”





Modern Tokenization and Vocabularies

- Subword tokenization is used for all modern pretrained models (though people are still experimenting with character-based models)
- Vocabularies contain ~50-250k wordpieces
- Pretrained word embeddings (e.g. GloVe) aren't necessary



Modern Tokenization and Vocabularies



Matthew Watkins
@SoC_trilogy

I've just found out that several of the anomalous GPT tokens ("TheNitromeFan", " SolidGoldMagikarp", " davidjl", " Smartstocks", " RandomRedditorWithNo",) are handles of people who are (competitively? collaboratively?) counting to infinity on a Reddit forum. I kid you not.

The screenshot shows a search results page for the subreddit 'r/artbn_bots' on Reddit. The search bar at the top contains the text 'r/artbn_bots'. To the left of the search bar is a sidebar with a 'ORDER BY SIZE' dropdown set to 'LARGE' and a list of sorting options: 'Full list.', 'Top 1000', 'Top 500', 'Top 250', and 'Top 100'. The main content area displays a table of user counts, which is the focus of the slide.

...

	Rank	User	Counts
1	/u/davidjl123	163477	
2	/u/Smartstocks	113829	
3	/u/atomicimploder	103178	
4	/u/TheNitromeFan	84581	
5	/u/SolidGoldMagikarp	65753	
6	/u/RandomRedditorWithNo	63434	
7	/u/rideride	59024	
8	/u/Mooraell	57785	
9	/u/Removedpixel	55080	
10	/u/Adinida	48415	
11	/u/rschaosid	47132	

Step 3: Architecture



Recap: Feedforward Networks

- Tokenize
- Embed
- Concatenate
- Linear layer
- Softmax
- Fixed window?
- Word averaging?

output distribution

$$\hat{y} = \text{softmax}(\mathbf{U}\mathbf{h} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer

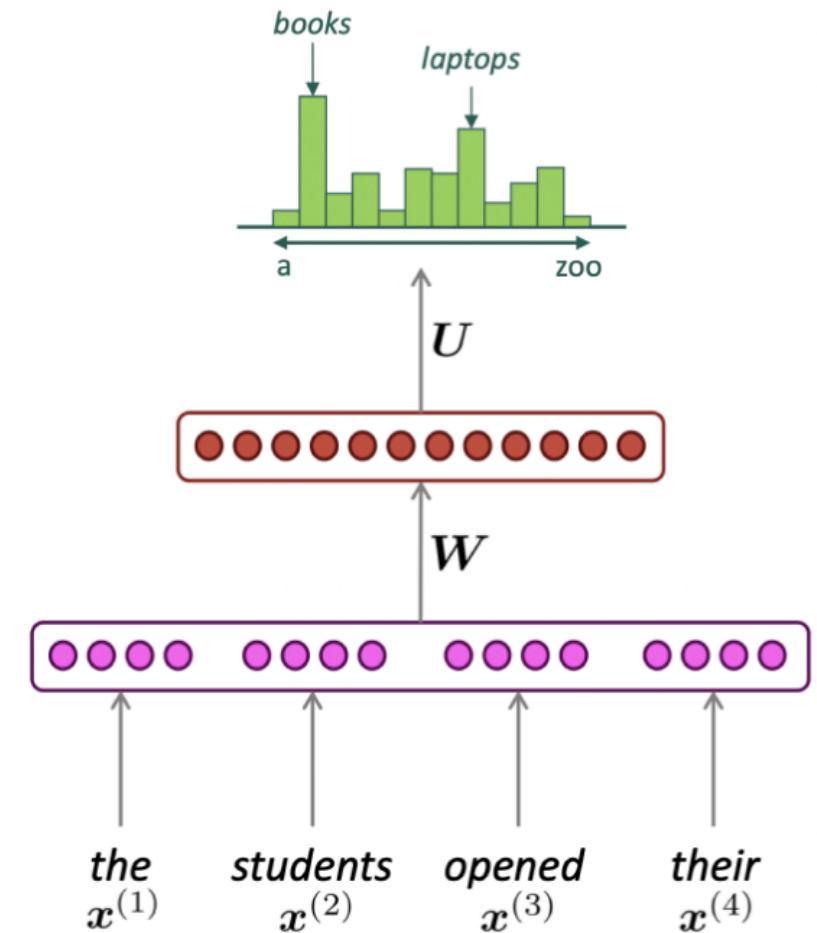
$$\mathbf{h} = f(\mathbf{W}\mathbf{e} + \mathbf{b}_1)$$

concatenated word embeddings

$$\mathbf{e} = [\mathbf{e}^{(1)}; \mathbf{e}^{(2)}; \mathbf{e}^{(3)}; \mathbf{e}^{(4)}]$$

words / one-hot vectors

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \mathbf{x}^{(4)}$$





Recap: Recurrence

output distribution

$$\hat{y}^{(t)} = \text{softmax}(\mathbf{U}h^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden states

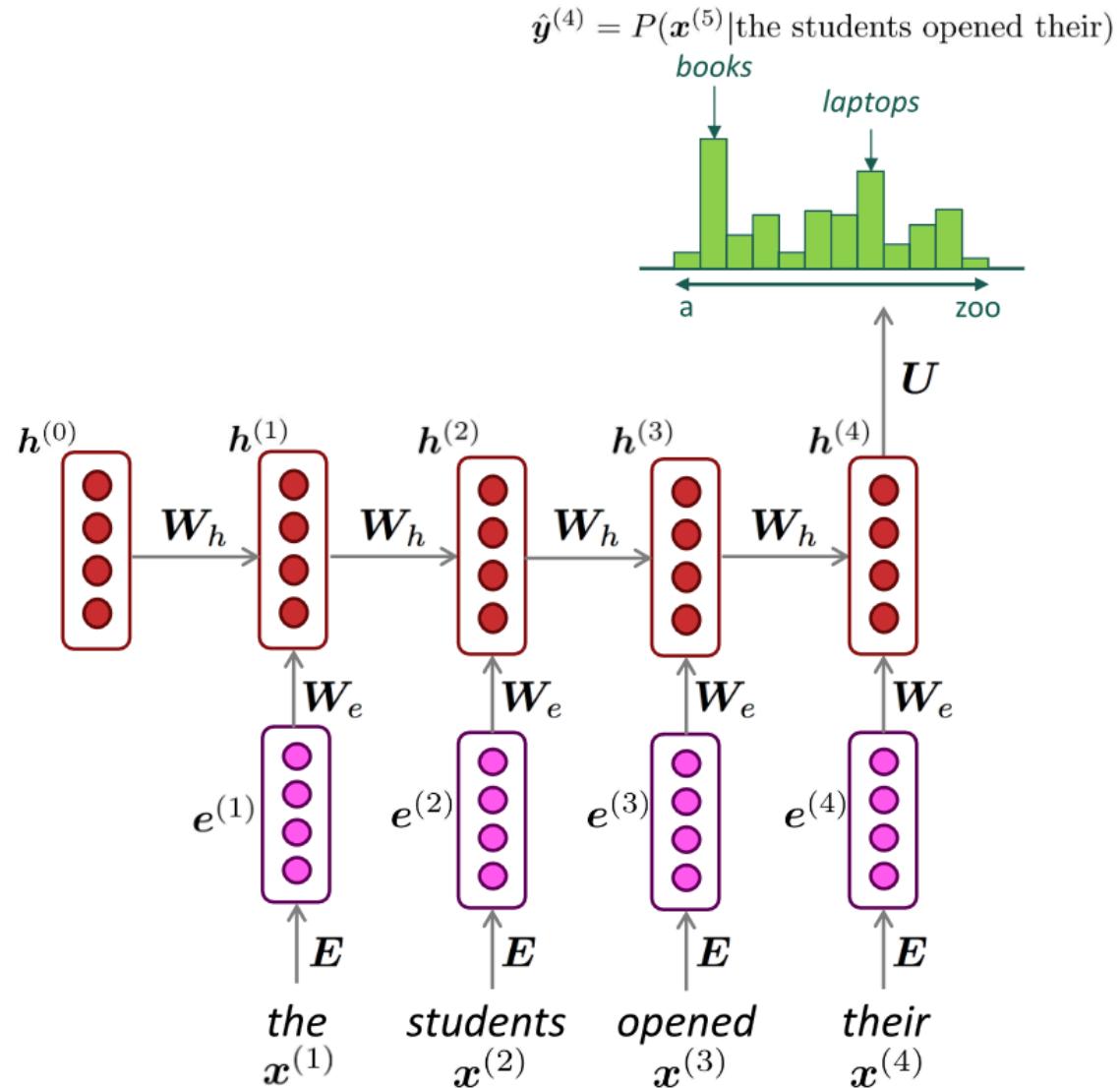
$$h^{(t)} = \sigma(\mathbf{W}_h h^{(t-1)} + \mathbf{W}_e e^{(t)} + \mathbf{b}_1)$$

$h^{(0)}$ is the initial hidden state

word embeddings

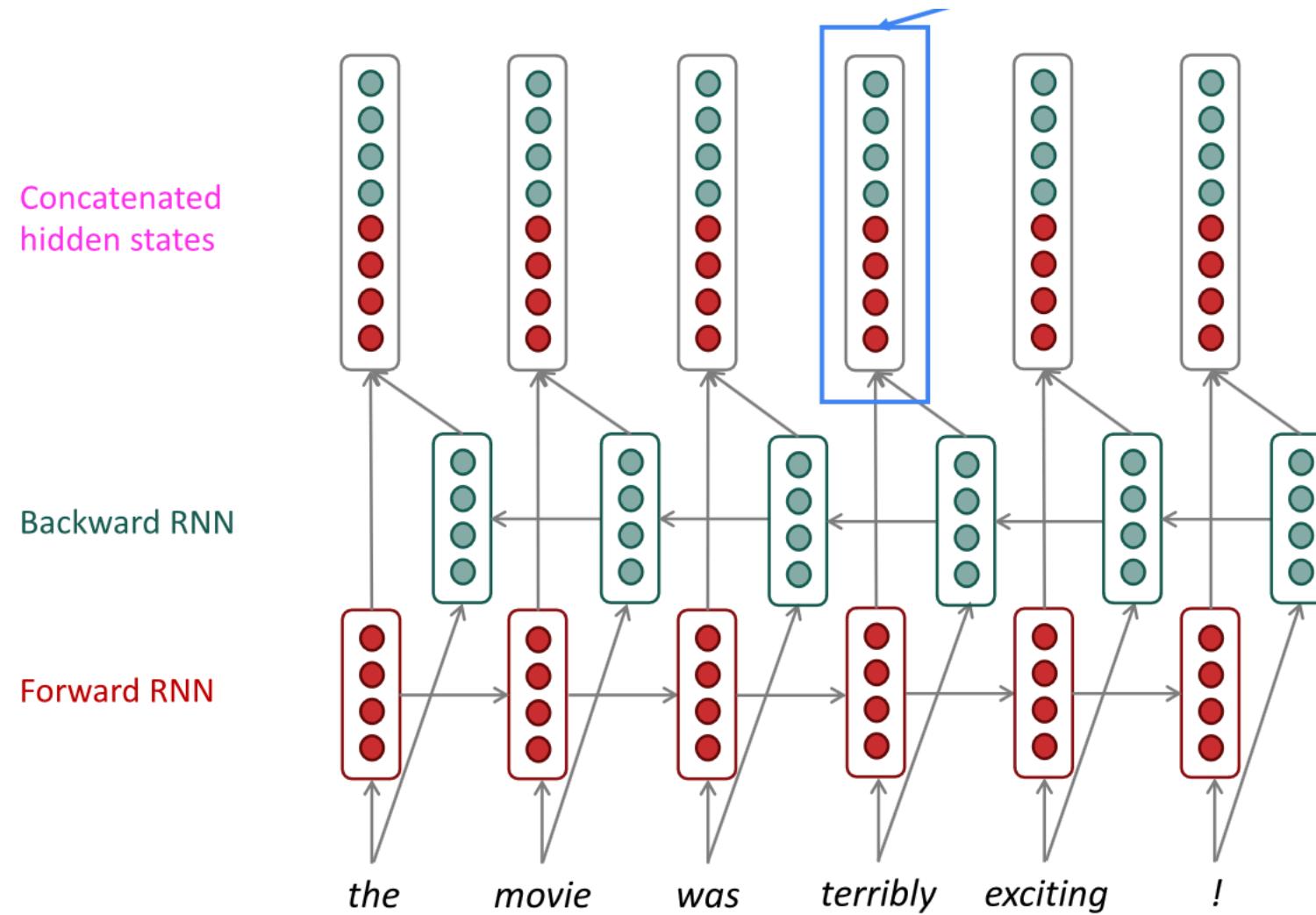
$$e^{(t)} = \mathbf{E}x^{(t)}$$

words / one-hot vectors
 $x^{(t)} \in \mathbb{R}^{|V|}$



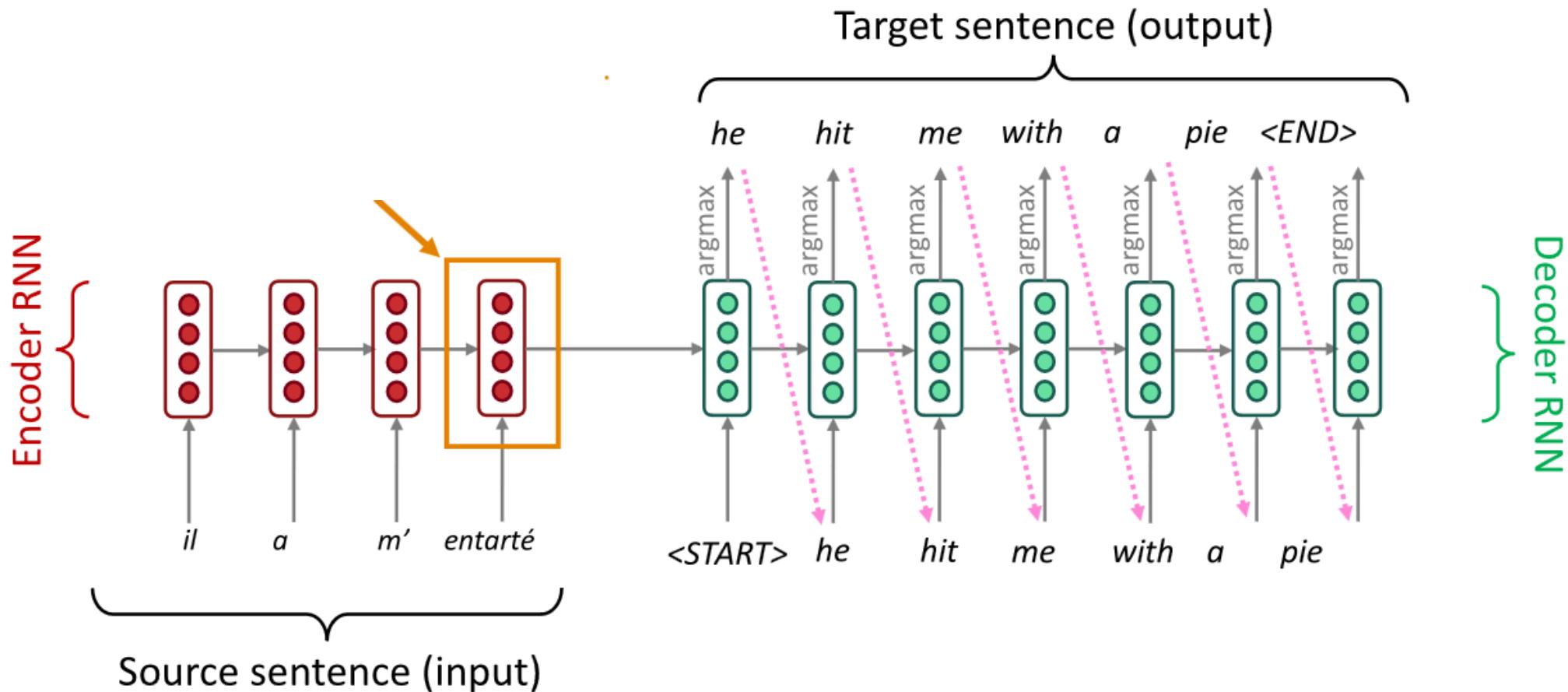


Recap: Recurrence



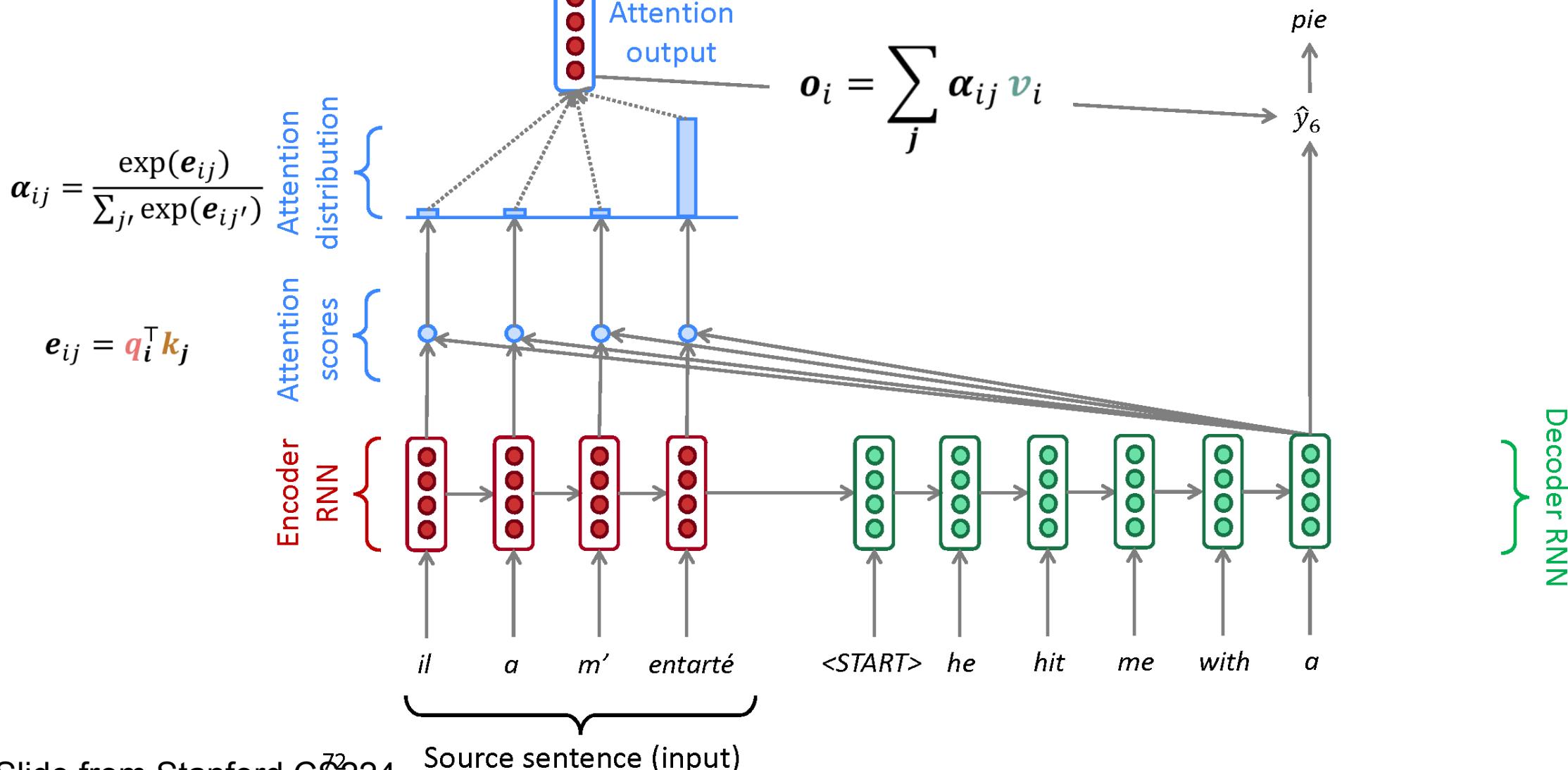


Recap: Recurrence





Recap: Attention





Recap: Attention

- Generic dot-product attention:

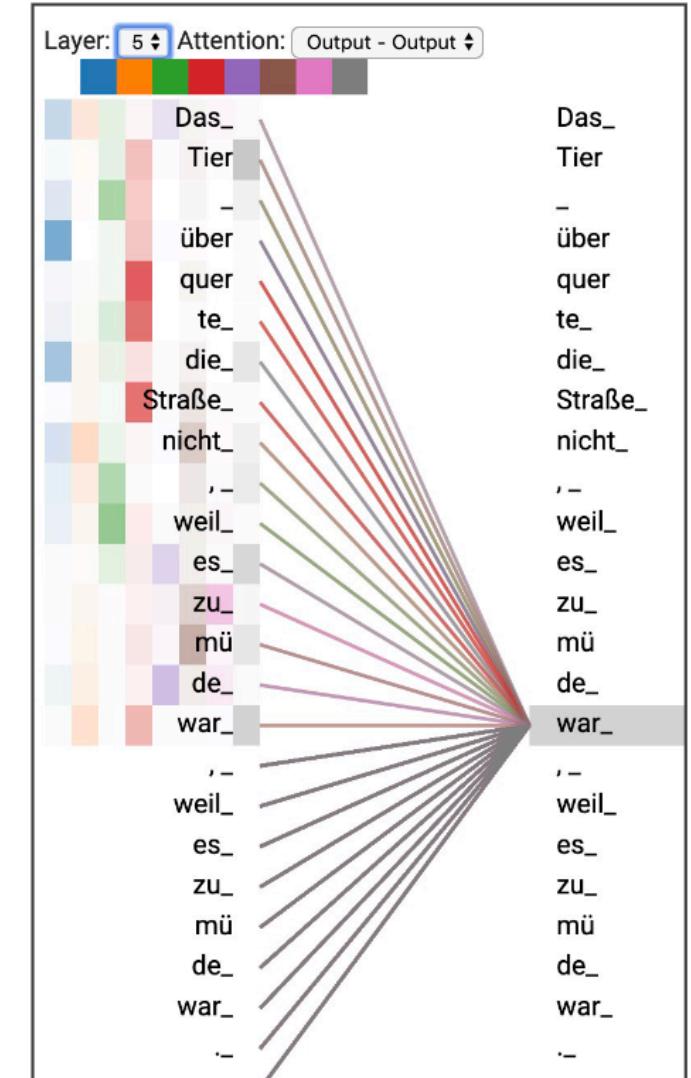
$$\mathbf{e}_{ij} = \mathbf{q}_i^T \mathbf{k}_j \quad \alpha_{ij} = \frac{\exp(\mathbf{e}_{ij})}{\sum_j \exp(\mathbf{e}_{ij'})} \quad \mathbf{o}_i = \sum_j \alpha_{ij} \mathbf{v}_i$$

- Self-attention: queries, keys, and values are all different transformations of the same item-level representation of some sequence:

$$\mathbf{q}_i = Q\mathbf{x}_i \text{ (queries)}$$

$$\mathbf{k}_i = K\mathbf{x}_i \text{ (keys)}$$

$$\mathbf{v}_i = V\mathbf{x}_i \text{ (values)}$$



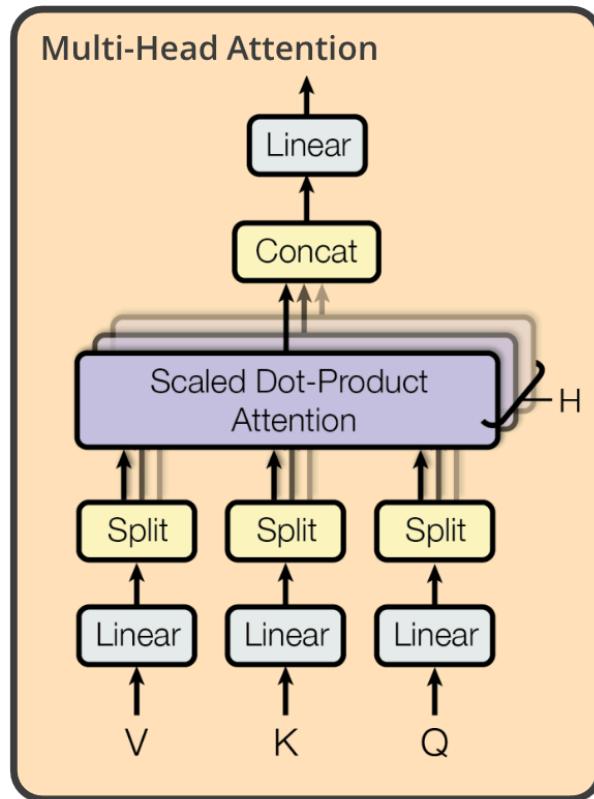


Multi-Head Attention

$$\mathbf{q}_i = Q\mathbf{x}_i \text{ (queries)}$$

$$\mathbf{k}_i = K\mathbf{x}_i \text{ (keys)}$$

$$\mathbf{v}_i = V\mathbf{x}_i \text{ (values)}$$



$$\text{head}_1 = \text{Attention}(\mathbf{Q}\mathbf{W}_1^Q, \mathbf{K}\mathbf{W}_1^K, \mathbf{V}\mathbf{W}_1^V)$$

:

$$\text{head}_H = \text{Attention}(\mathbf{Q}\mathbf{W}_H^Q, \mathbf{K}\mathbf{W}_H^K, \mathbf{V}\mathbf{W}_H^V)$$

$$\text{MultiHeadAtt}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Concat}(\text{head}_1, \dots, \text{head}_H)$$

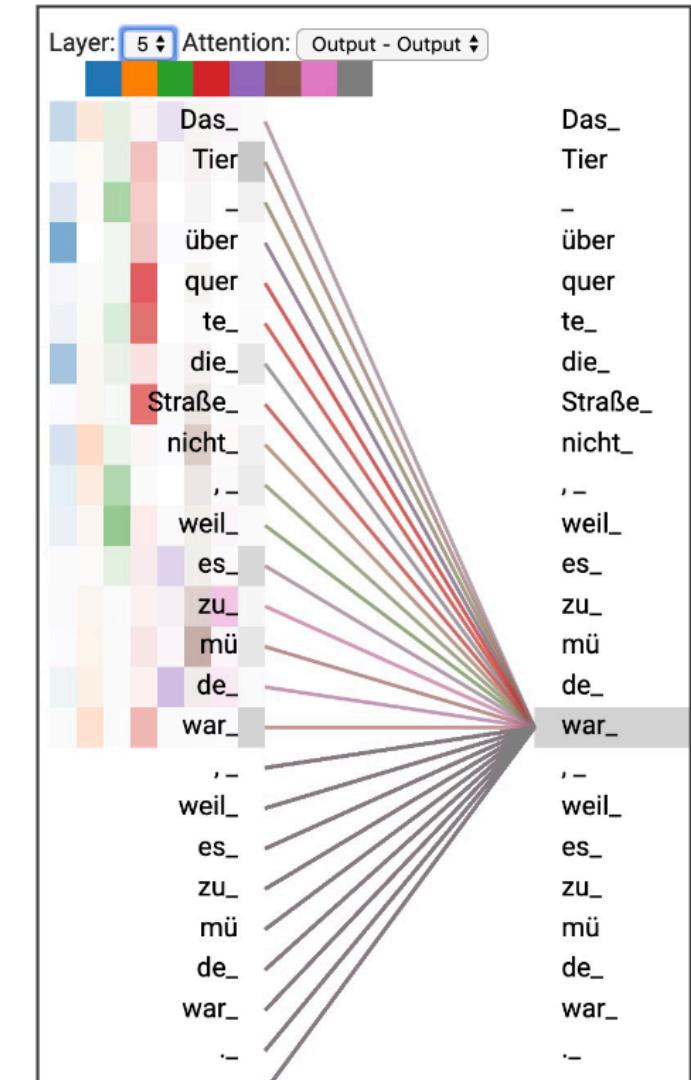
Inputs and outputs of each layer
are the same dimensions:

$$\mathbf{Q} \in \mathbb{R}^{T \times d_{\text{model}}}$$

$$\mathbf{K} \in \mathbb{R}^{T \times d_{\text{model}}}$$

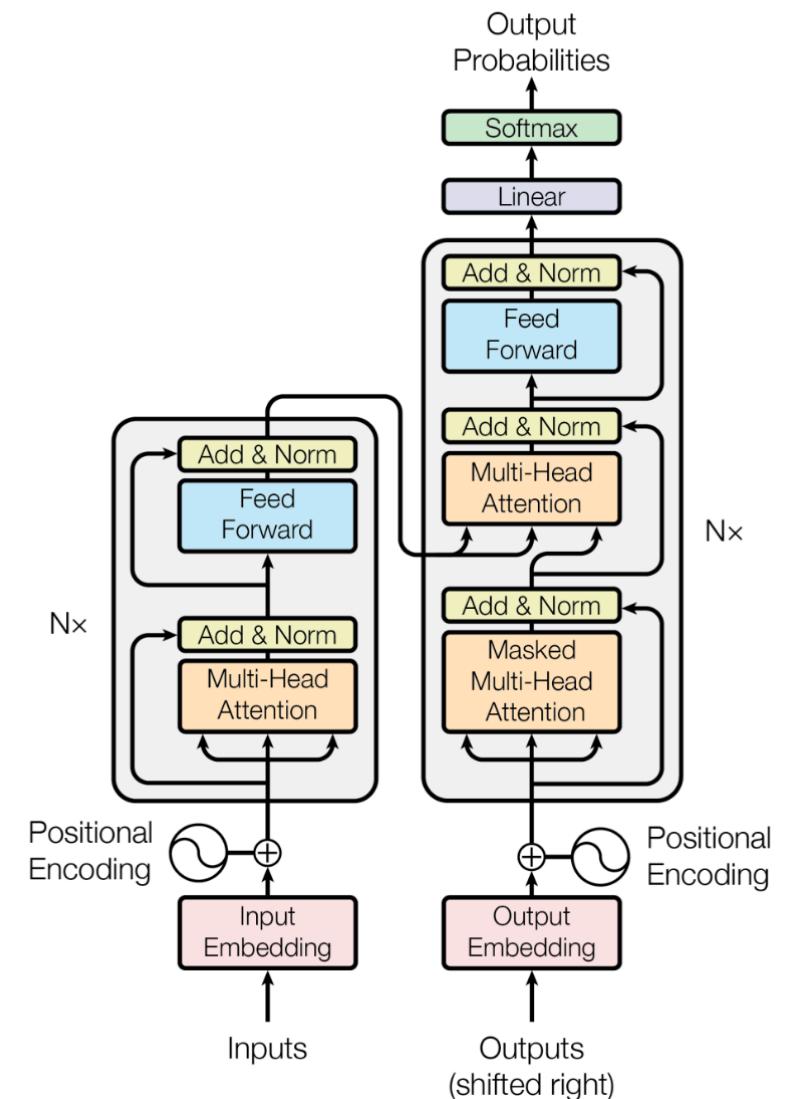
$$\mathbf{V} \in \mathbb{R}^{T \times d_{\text{model}}}$$

$$\text{MultiHeadAtt}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \in \mathbb{R}^{T \times d_{\text{model}}}$$

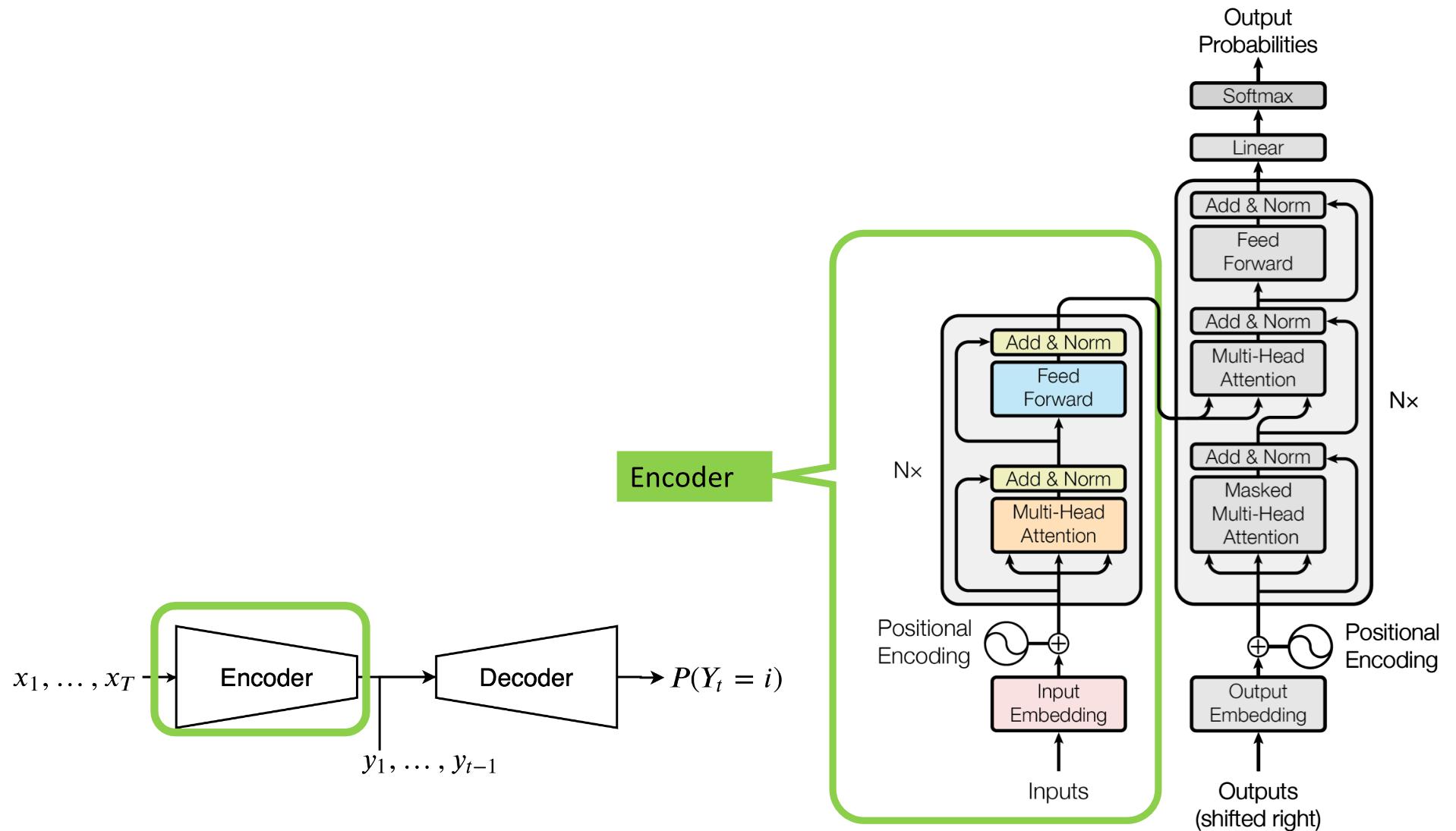




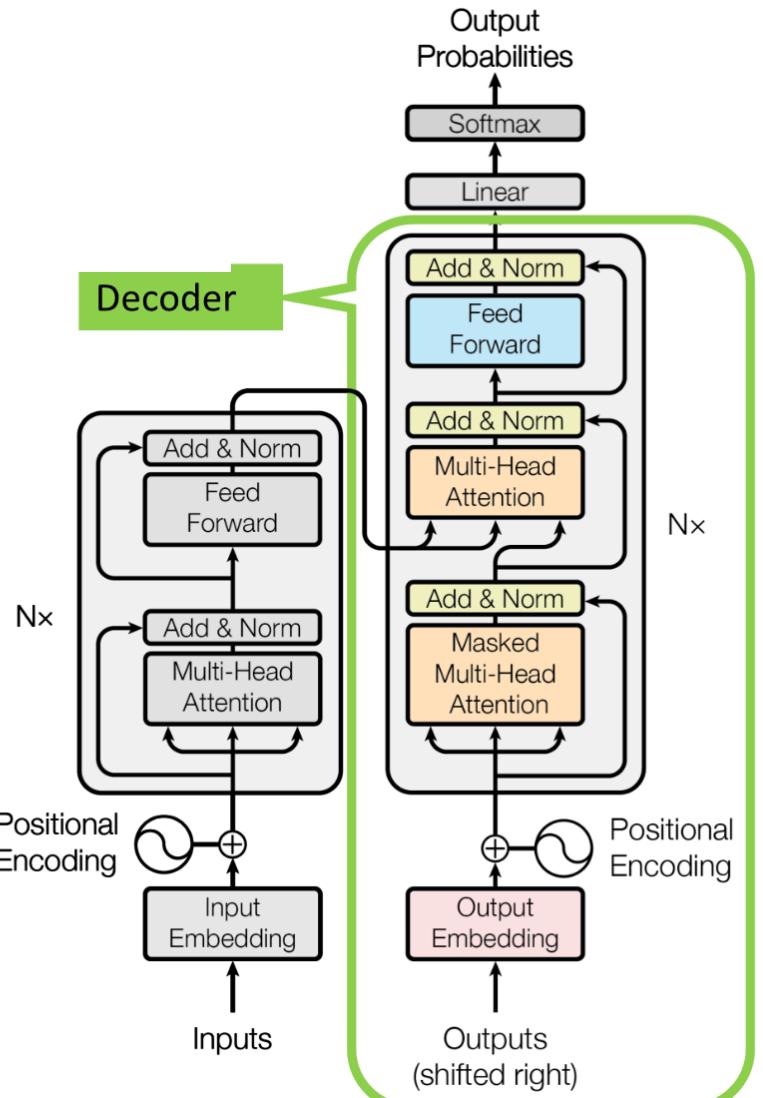
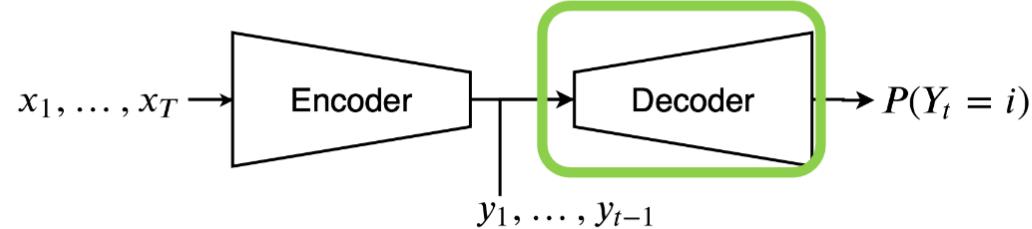
Transformer



Encoder



Decoder





Encoder Input

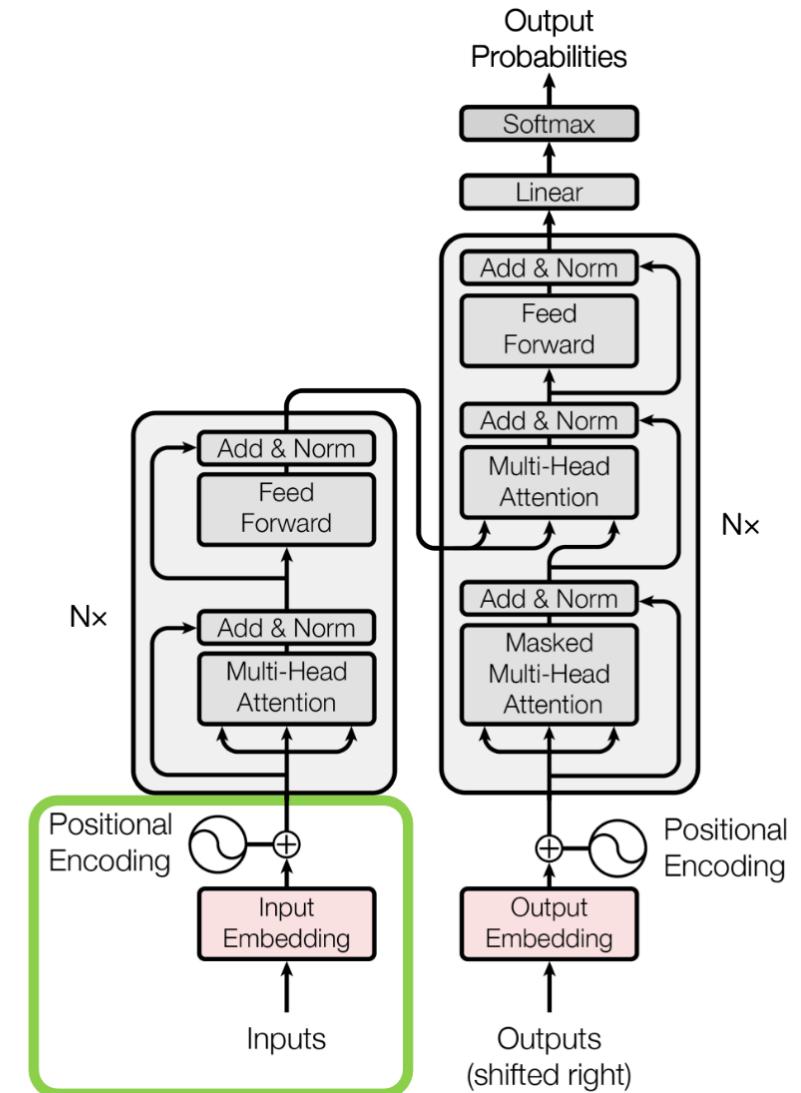
The input into the encoder looks like:

$$\mathbf{H}_0^{\text{enc}} = \begin{matrix} \text{Token Embeddings} \\ \mathbf{H}_0^{\text{enc}} = \left[\begin{array}{c|c} \hline \text{padding} & \text{embedding size} \\ \hline \text{maximum sequence length} & \\ \hline \end{array} \right] + \text{Position Embeddings} \end{matrix}$$

$\mathbf{p}_i = \begin{cases} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{cases}$

Dimension

Index in the sequence





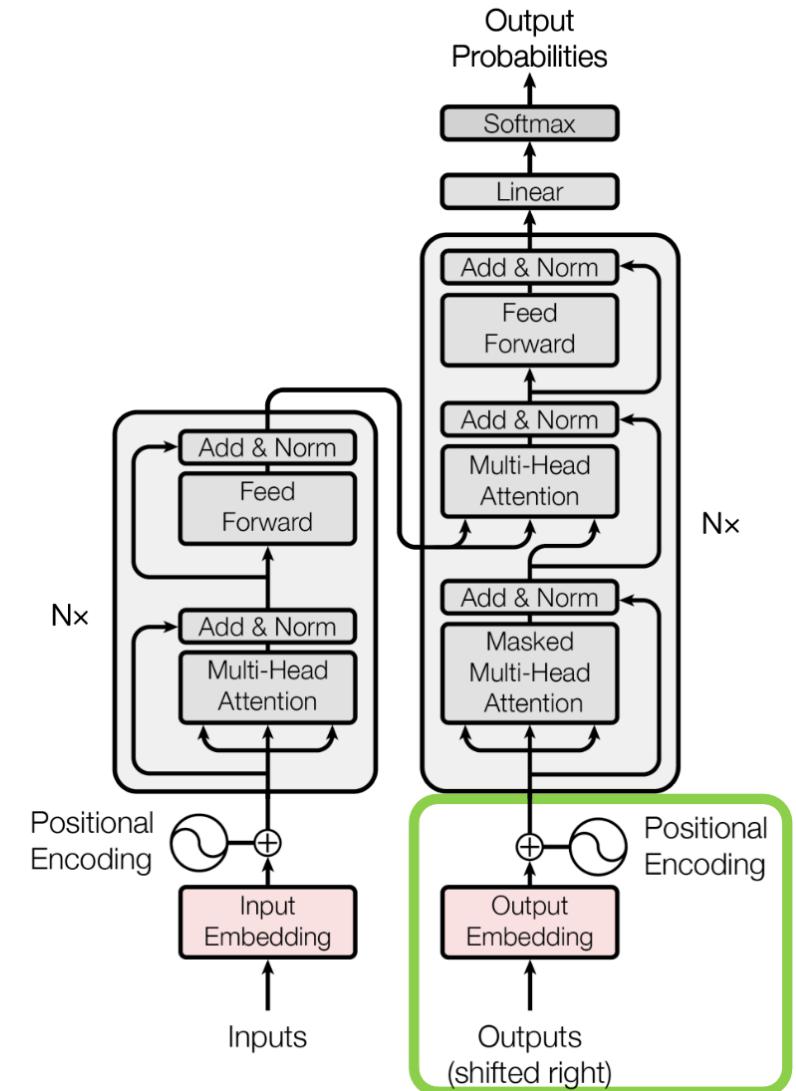
Decoder Input

The input into the encoder looks like:

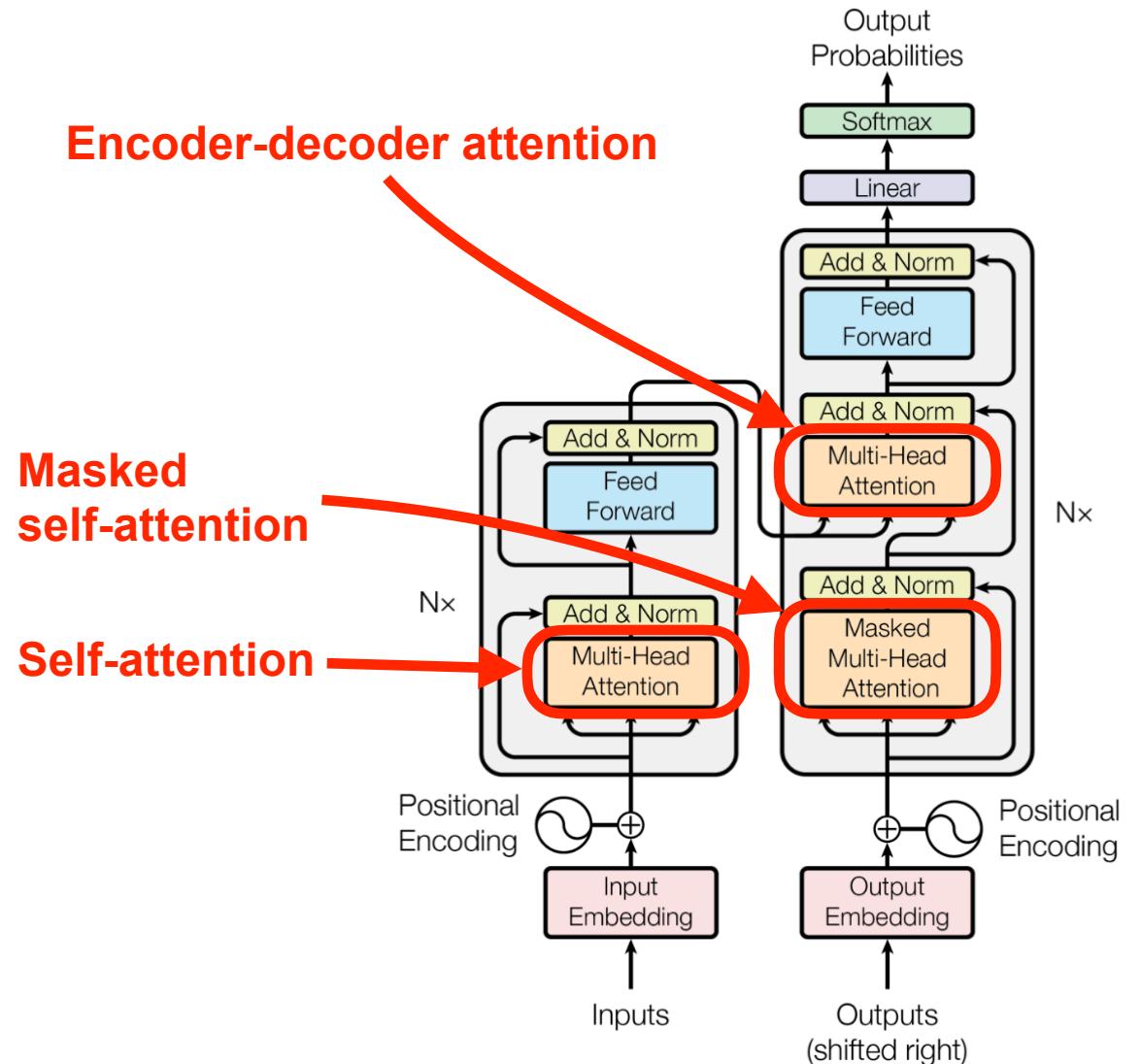
$$H_0^{\text{enc}} = \begin{matrix} \text{Token Embeddings} \\ \left[\begin{array}{c|c} \text{padding} & \end{array} \right] \end{matrix} + \begin{matrix} \text{Position Embeddings} \\ \left[\begin{array}{c} \text{embedding size} \\ \text{maximum sequence length} \end{array} \right] \end{matrix}$$

The input to the decoder looks like:

$$H_0^{\text{dec}} = \begin{matrix} \text{Shifted Token Embeddings} \\ \left[\begin{array}{c|c} \text{padding} & \end{array} \right] \end{matrix} + \begin{matrix} \text{Position Embeddings} \\ \left[\begin{array}{c} \text{embedding size} \\ \text{maximum sequence length} \end{array} \right] \end{matrix}$$



Attention

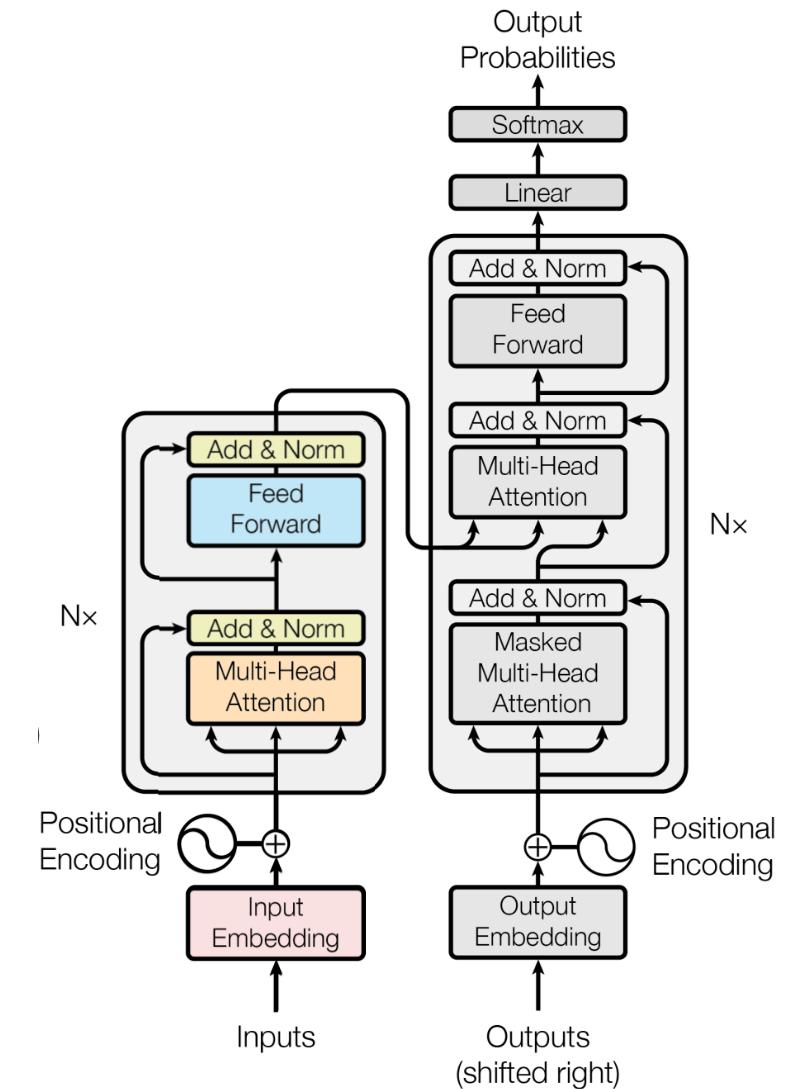




Encoder

Multi-Head
Attention

$$= \text{MultiHeadAtt}(\mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}, \mathbf{H}_i^{enc})$$

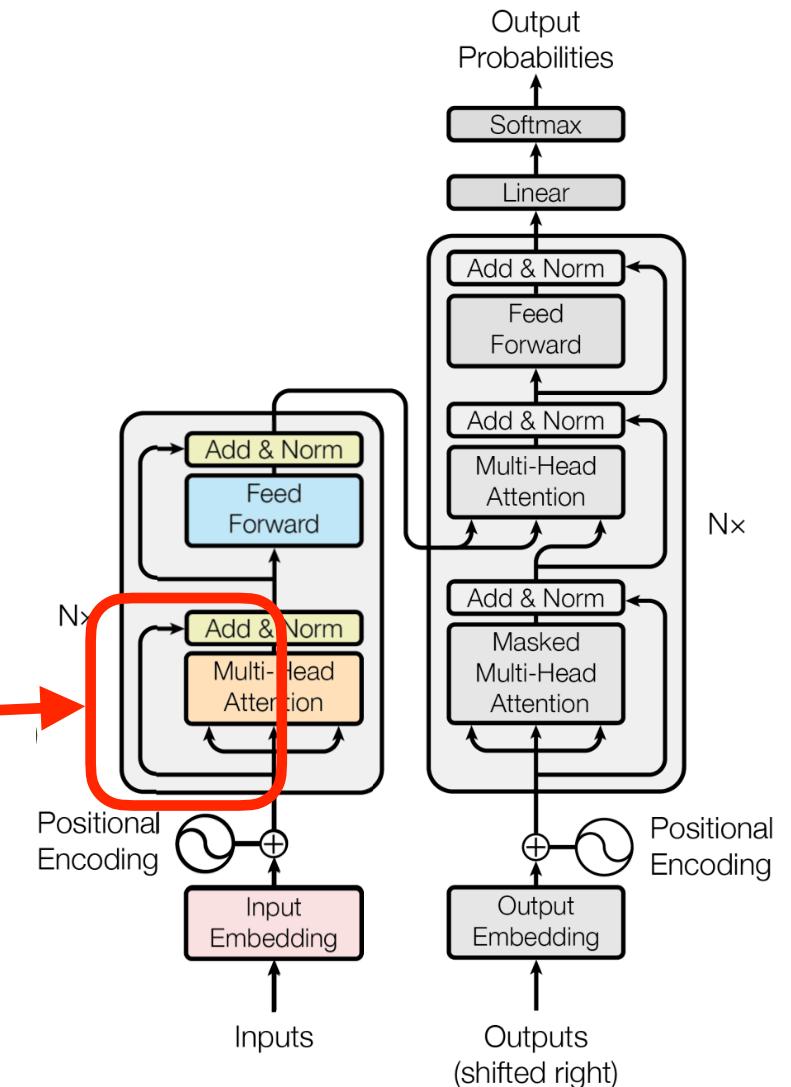


Encoder

$$\text{Multi-Head Attention} = \text{MultiHeadAtt}(\mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}, \mathbf{H}_i^{enc})$$

$$\text{Add & Norm} = \text{LayerNorm}(\text{Multi-Head Attention} + \mathbf{H}_i^{enc})$$

Residual connection



Encoder



Multi-Head Attention

$$= \text{MultiHeadAtt}(\mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}, \mathbf{H}_i^{enc})$$

Add & Norm

$$= \text{LayerNorm}(\text{Multi-Head Attention} + \mathbf{H}_i^{enc})$$

Let $x \in \mathbb{R}^d$ be an individual (word) vector in the model.

Let $\mu = \sum_{j=1}^d x_j$; this is the mean; $\mu \in \mathbb{R}$.

Let $\sigma = \sqrt{\frac{1}{d} \sum_{j=1}^d (x_j - \mu)^2}$; this is the standard deviation; $\sigma \in \mathbb{R}$.

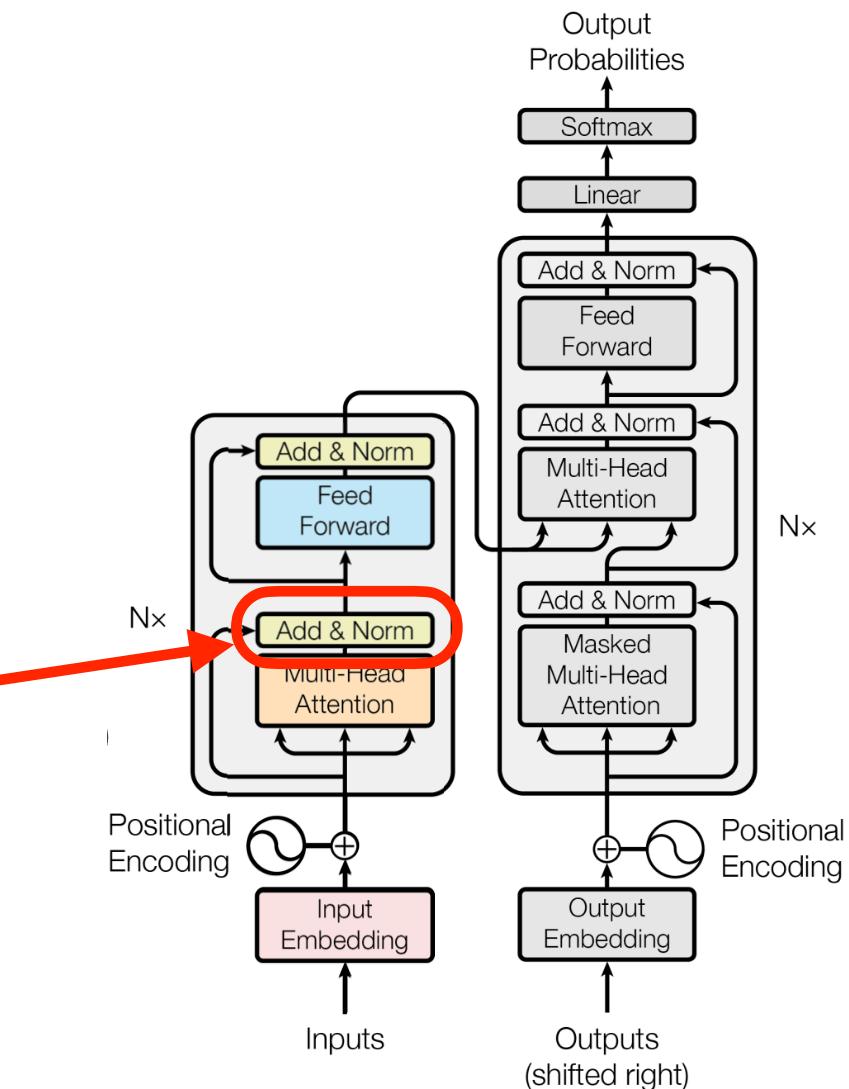
Let $\gamma \in \mathbb{R}^d$ and $\beta \in \mathbb{R}^d$ be learned “gain” and “bias” parameters. (Can omit!)

Then layer normalization computes:

$$\text{output} = \frac{x - \mu}{\sqrt{\sigma} + \epsilon} * \gamma + \beta$$

Normalize by scalar mean and variance

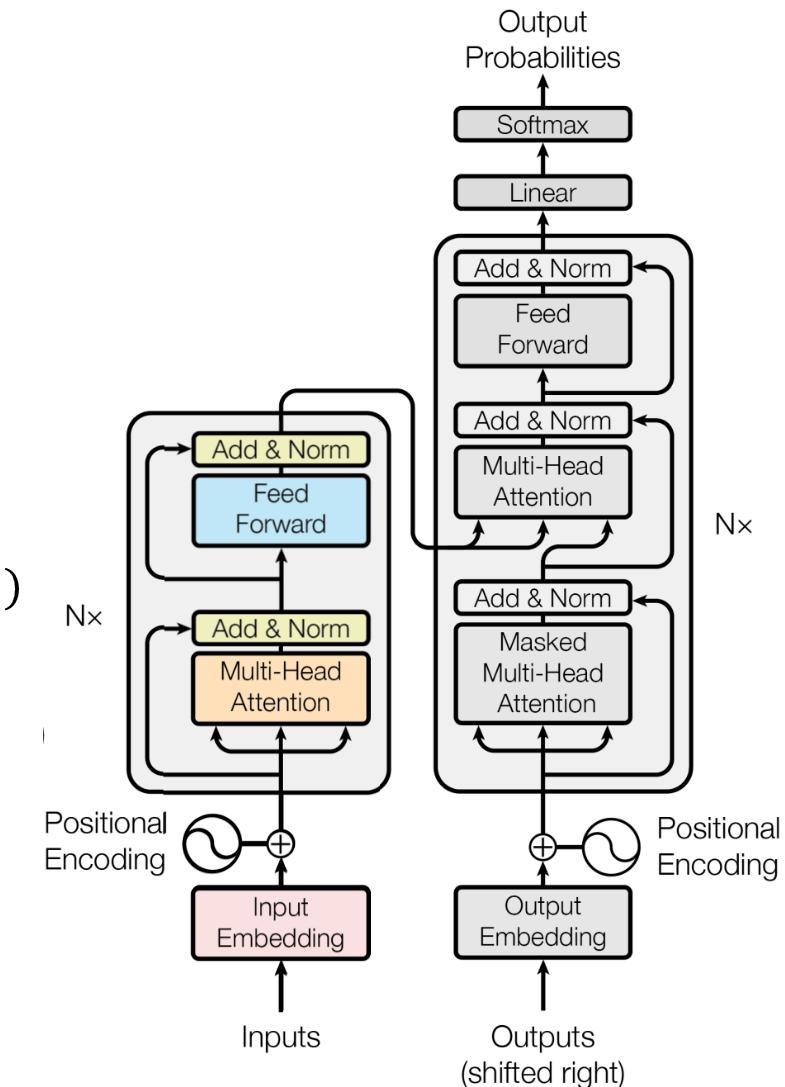
Modulate by learned elementwise gain and bias





Encoder

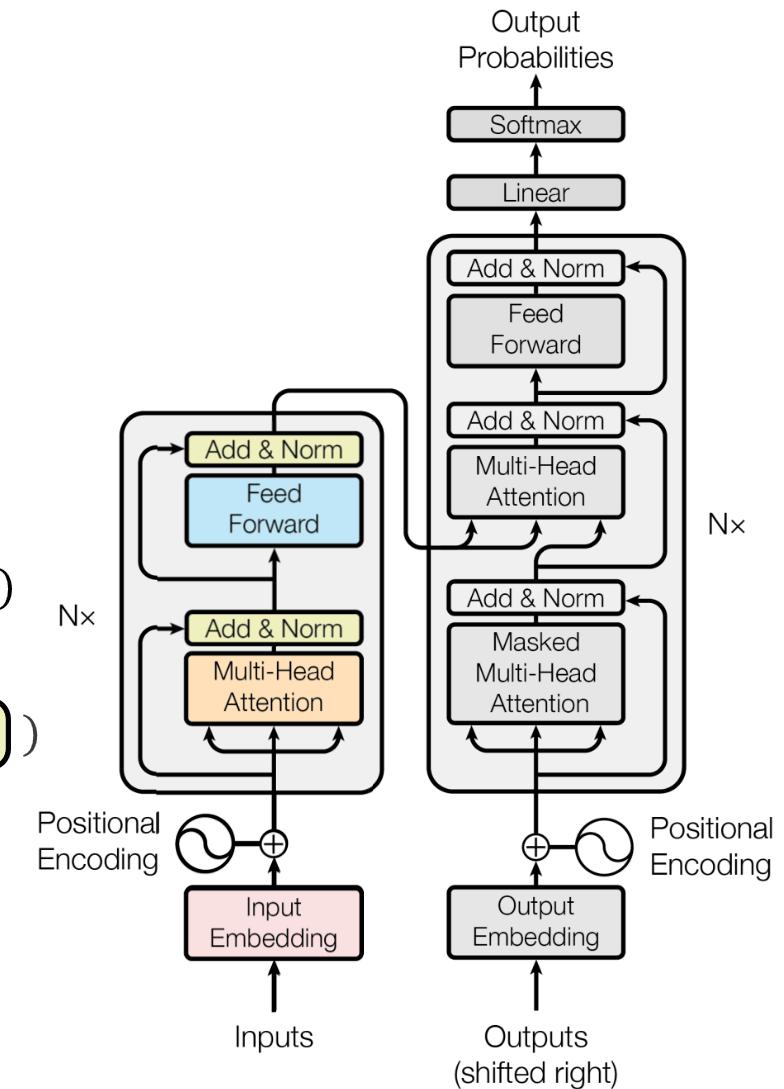
$$\begin{aligned}\text{Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}) \\ \text{Add \& Norm} &= \text{LayerNorm}(\text{Multi-Head Attention} + \mathbf{H}_i^{enc}) \\ \text{Feed Forward} &= \max(0, \text{Add \& Norm } \mathbf{W}_1 + b_1) \mathbf{W}_2 + b_2)\end{aligned}$$





Encoder

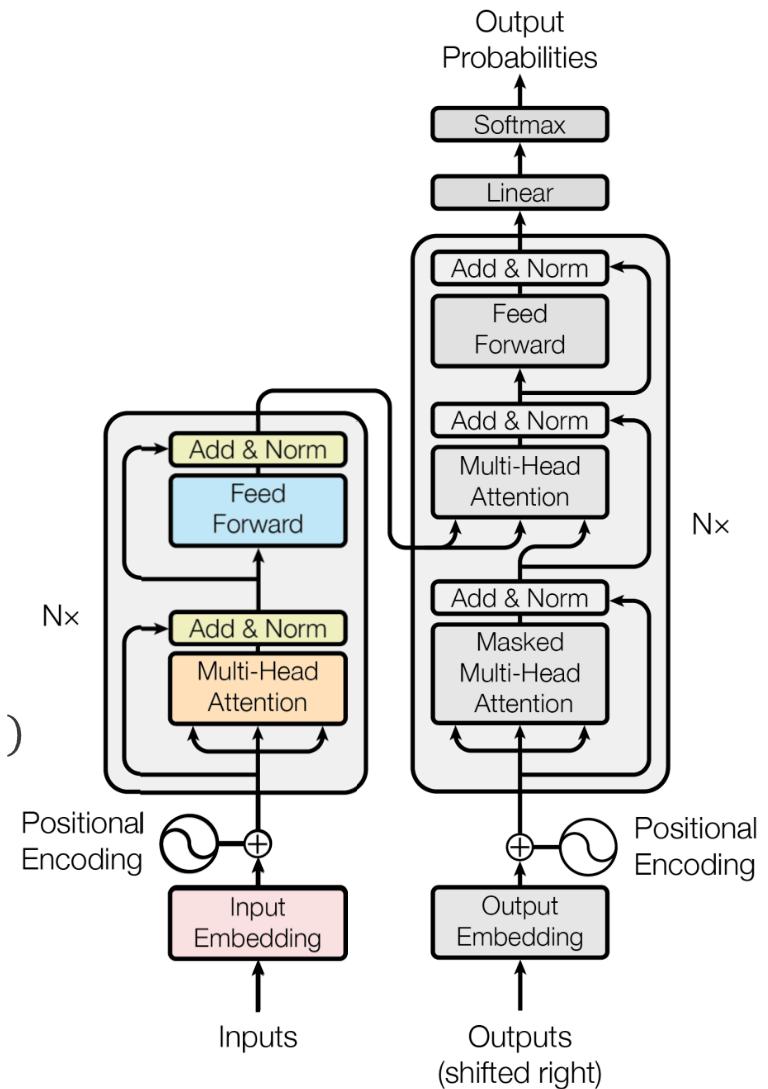
$$\begin{aligned}\text{Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}) \\ \text{Add & Norm} &= \text{LayerNorm}(\text{Multi-Head Attention} + \mathbf{H}_i^{enc}) \\ \text{Feed Forward} &= \max(0, \text{Add & Norm } (\mathbf{W}_1 + b_1) \mathbf{W}_2 + b_2) \\ \text{Add & Norm (2)} &= \text{LayerNorm}(\text{Feed Forward} + \text{Add & Norm })\end{aligned}$$



Encoder



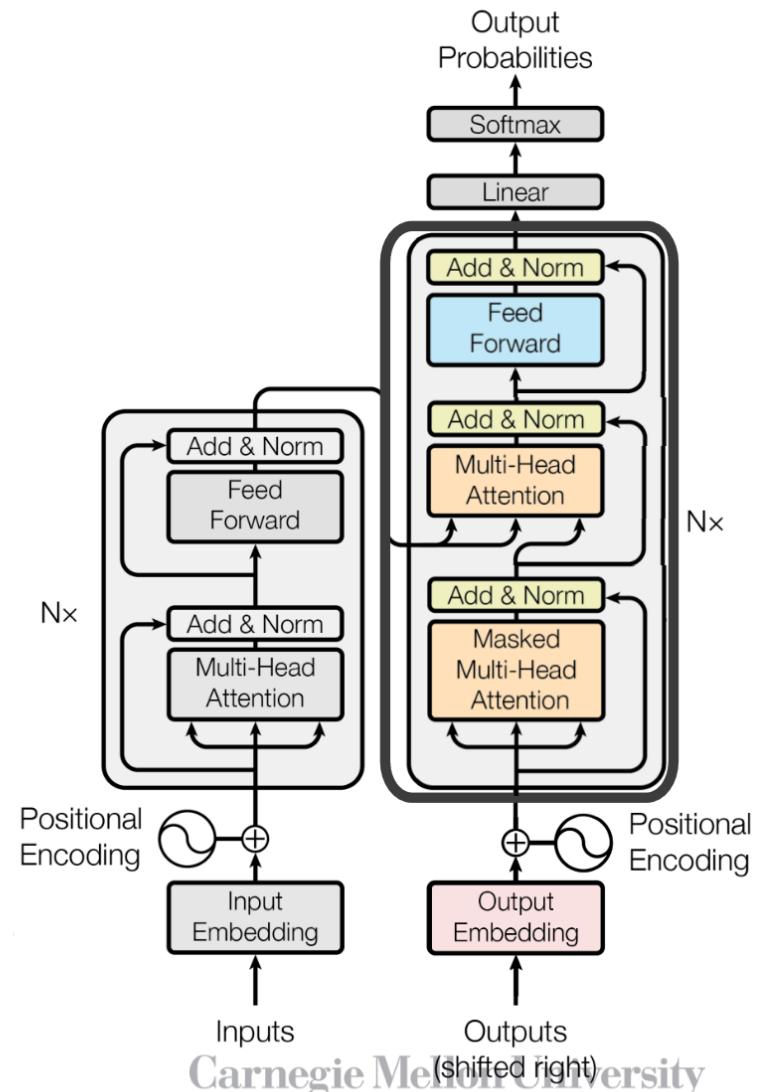
$$\begin{aligned}
 \text{Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}) \\
 \text{Add & Norm} &= \text{LayerNorm}(\text{Multi-Head Attention} + \mathbf{H}_i^{enc}) \\
 \text{Feed Forward} &= \max(0, \text{Add & Norm } (\mathbf{W}_1 + b_1) \mathbf{W}_2 + b_2) \\
 \text{Add & Norm (2)} &= \text{LayerNorm}(\text{Feed Forward} + \text{Add & Norm }) \\
 \mathbf{H}_{i+1}^{enc} &= \text{Add & Norm}(2)
 \end{aligned}$$





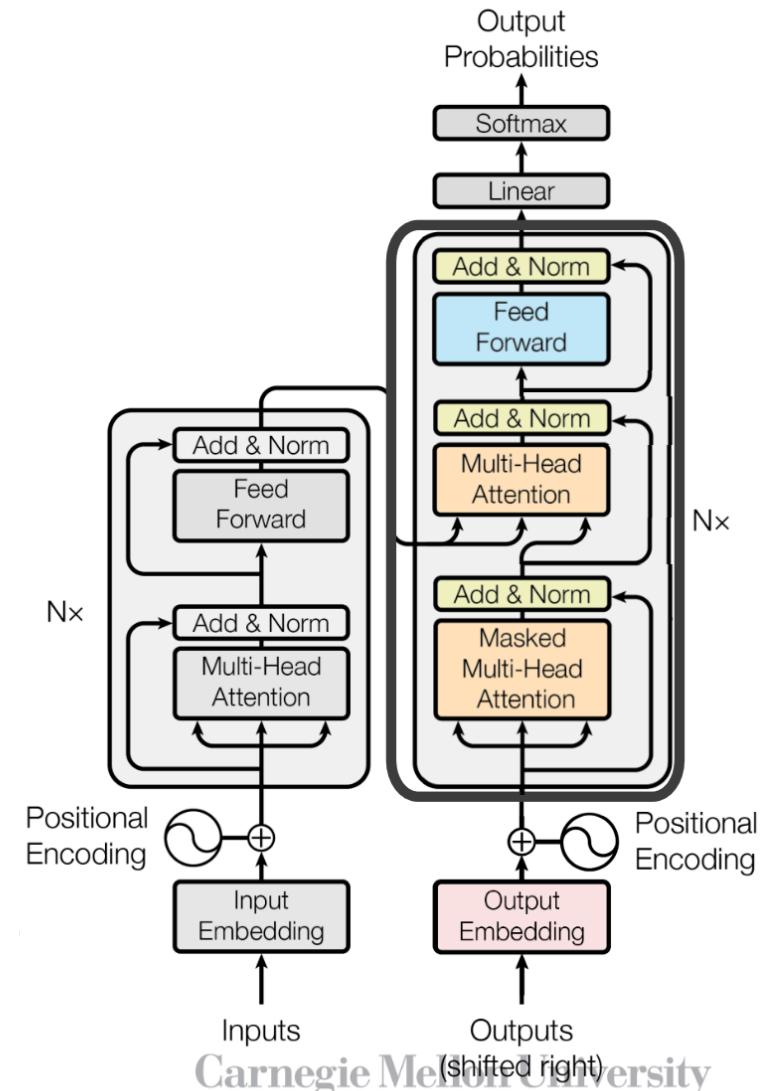
Decoder

$$\text{Masked Multi-Head Attention} = \text{MultiHeadAtt}(\mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}, \mathbf{H}_i^{dec})$$



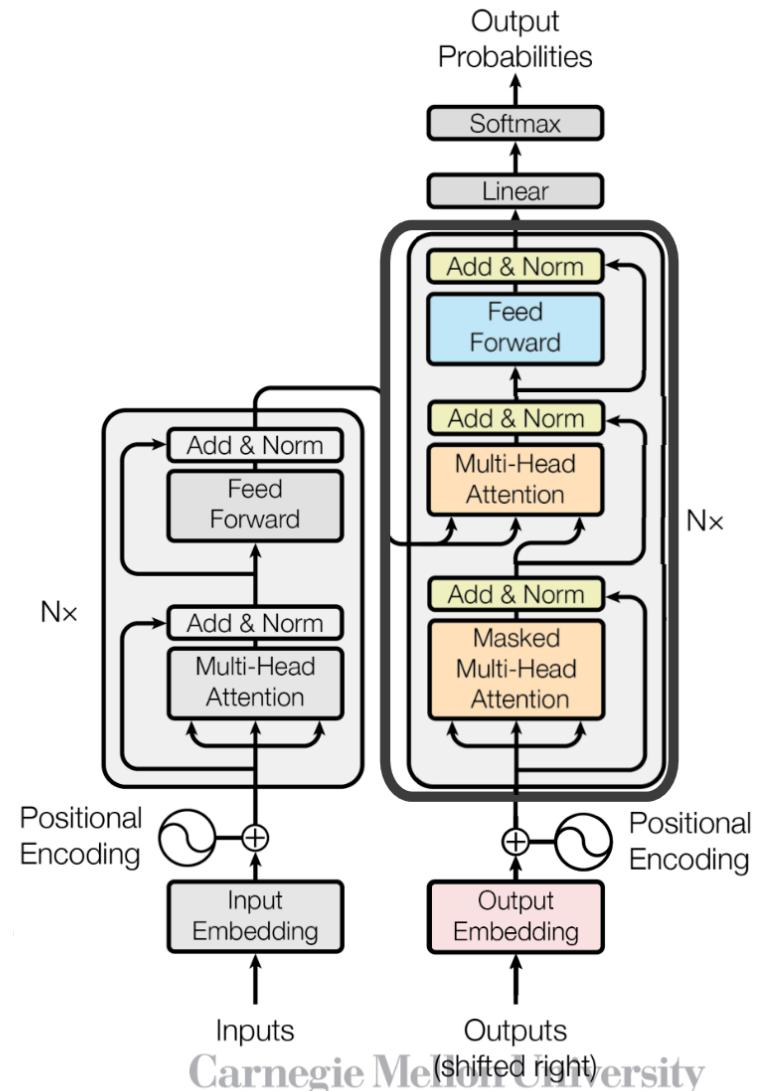
Decoder

$$\begin{aligned} \text{Masked Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}) \\ \text{Add & Norm} &= \text{LayerNorm}(\text{Masked Multi-Head Attention} + \mathbf{H}_i^{dec}) \end{aligned}$$



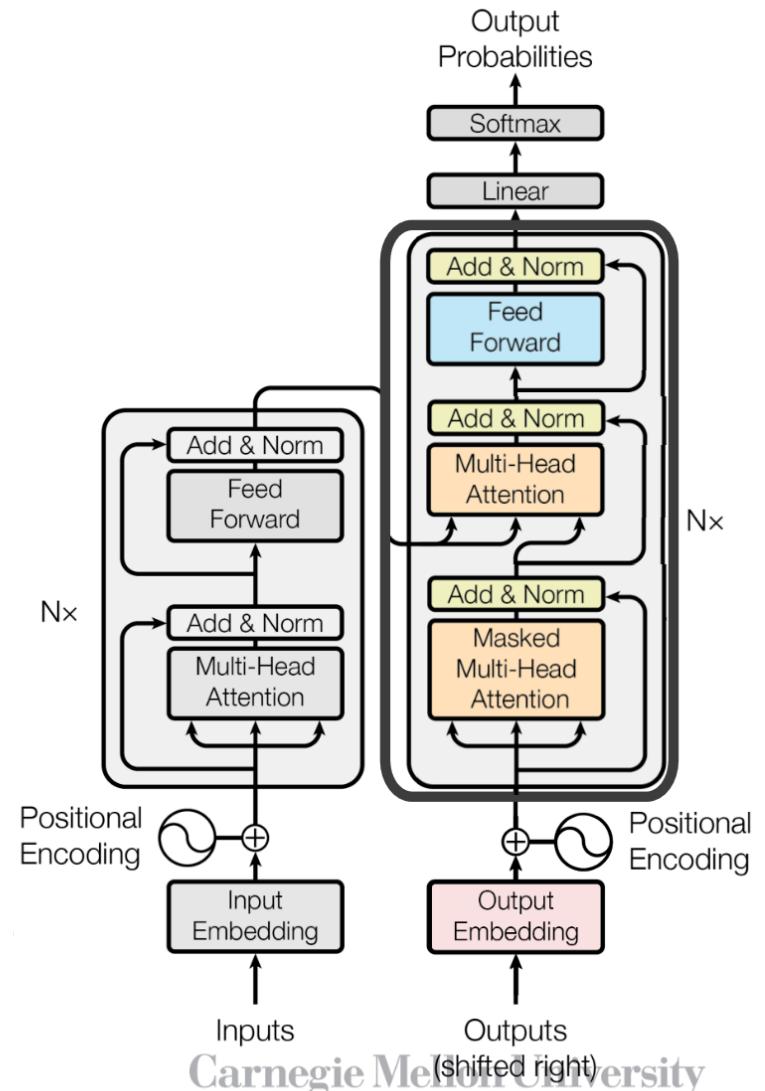
Decoder

$$\begin{aligned}\text{Masked Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}) \\ \text{Add & Norm} &= \text{LayerNorm}(\text{Masked Multi-Head Attention} + \mathbf{H}_i^{dec}) \\ \text{Enc-Dec Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}, \text{Add & Norm})\end{aligned}$$



Decoder

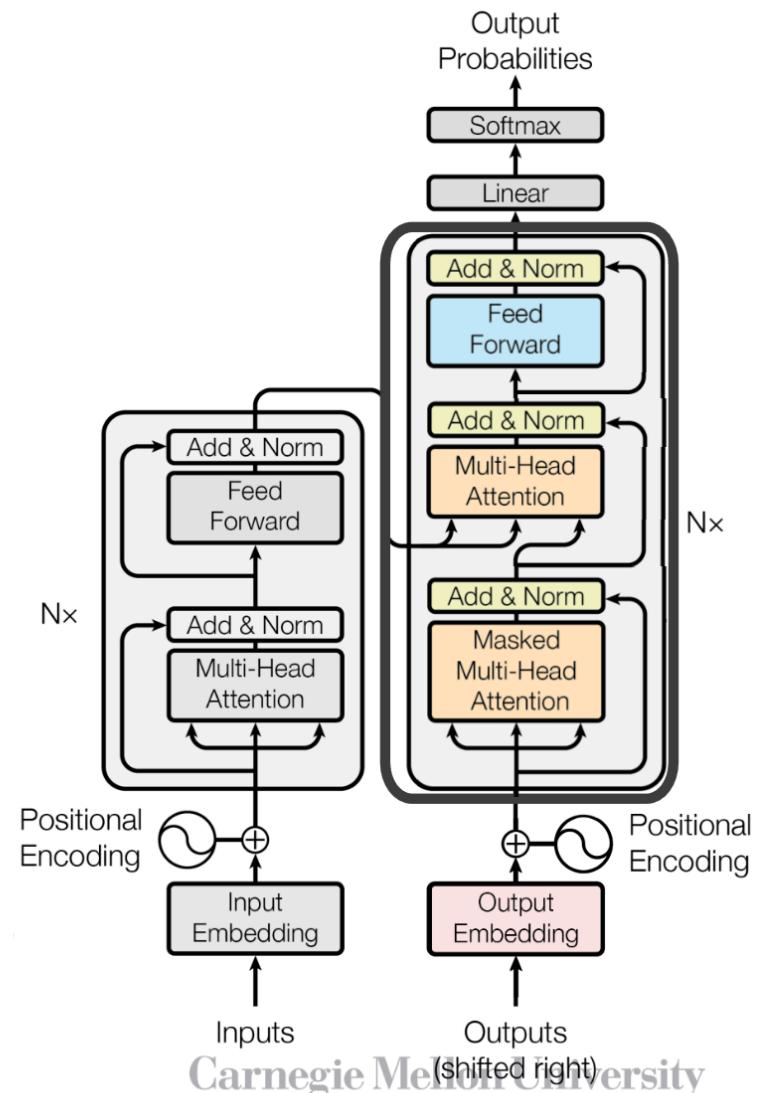
$$\begin{aligned}
 \text{Masked Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}) \\
 \text{Add & Norm} &= \text{LayerNorm}(\text{Masked Multi-Head Attention} + \mathbf{H}_i^{dec}) \\
 \text{Enc-Dec Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}, \text{Add & Norm}) \\
 \text{Add & Norm (2)} &= \text{LayerNorm}(\text{Enc-Dec Multi-Head Attention} + \text{Add & Norm})
 \end{aligned}$$





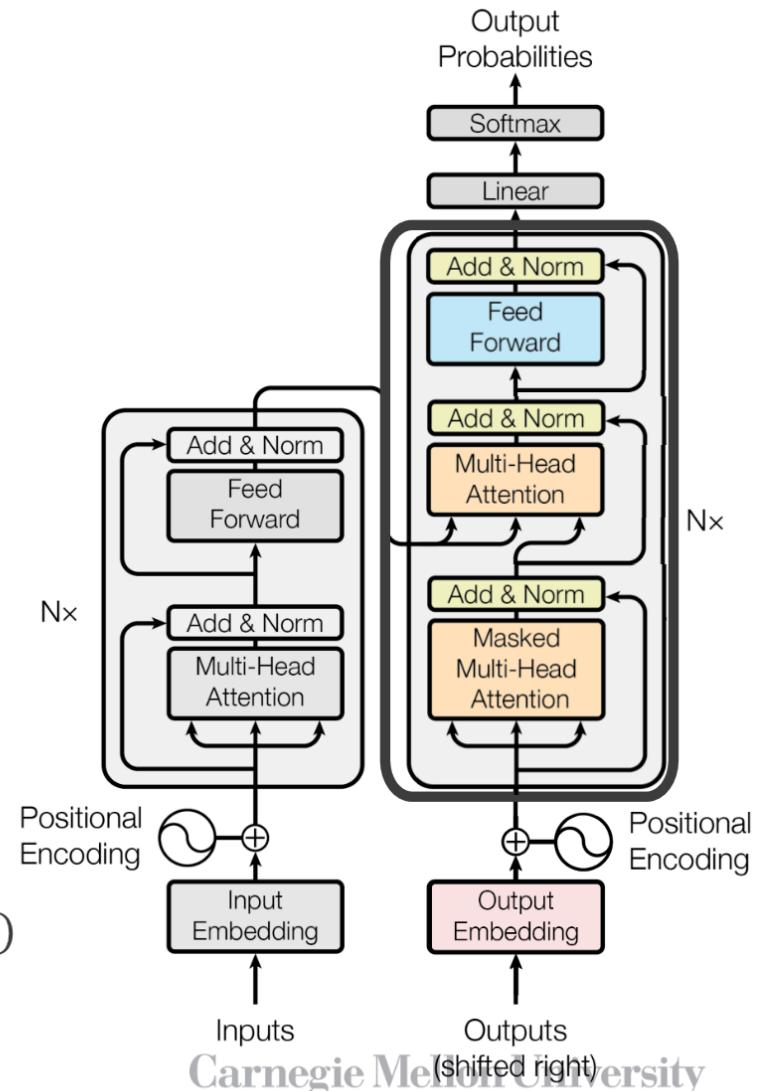
Decoder

$$\begin{aligned}
 \text{Masked Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}) \\
 \text{Add & Norm} &= \text{LayerNorm}(\text{Masked Multi-Head Attention} + \mathbf{H}_i^{dec}) \\
 \text{Enc-Dec Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}, \text{Add & Norm}) \\
 \text{Add & Norm (2)} &= \text{LayerNorm}(\text{Enc-Dec Multi-Head Attention} + \text{Add & Norm}) \\
 \text{Feed Forward} &= \max(0, \text{Add & Norm (2)} \mathbf{W}_1 + b_1) \mathbf{W}_2 + b_2)
 \end{aligned}$$



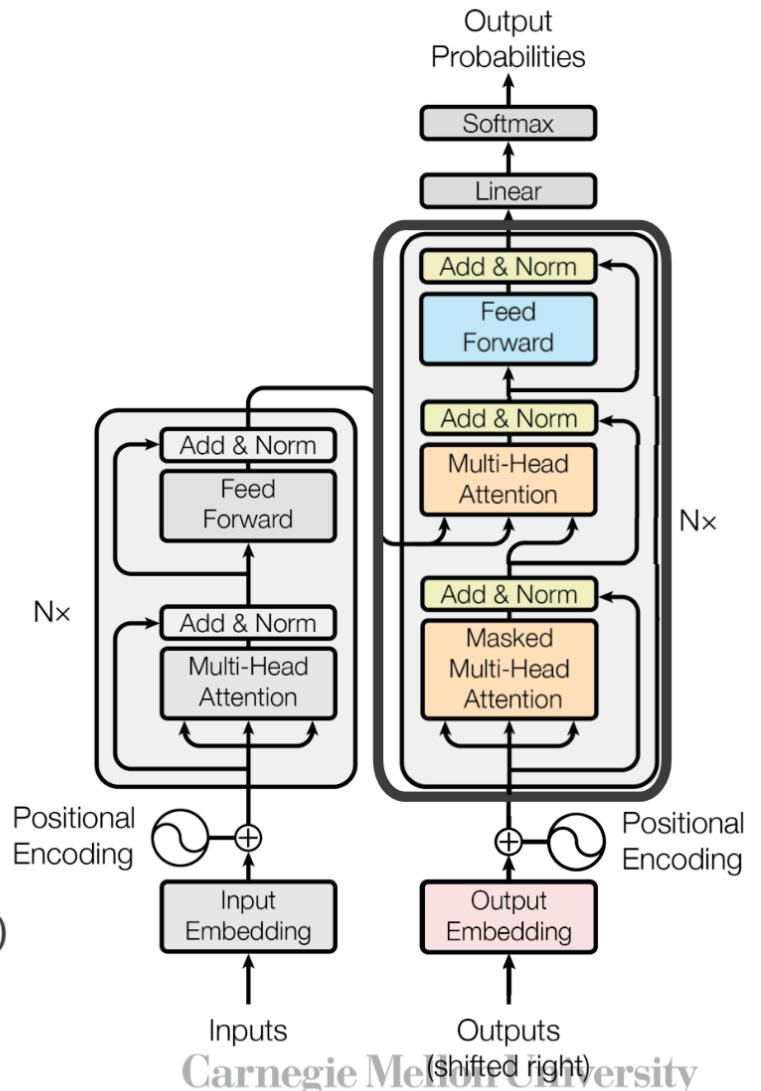
Decoder

$$\begin{aligned}
 \text{Masked Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}) \\
 \text{Add & Norm} &= \text{LayerNorm}(\text{Masked Multi-Head Attention} + \mathbf{H}_i^{dec}) \\
 \text{Enc-Dec Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}, \text{Add & Norm}) \\
 \text{Add & Norm (2)} &= \text{LayerNorm}(\text{Enc-Dec Multi-Head Attention} + \text{Add & Norm}) \\
 \text{Feed Forward} &= \max(0, \text{Add & Norm (2)} \mathbf{W}_1 + b_1) \mathbf{W}_2 + b_2 \\
 \text{Add & Norm (3)} &= \text{LayerNorm}(\text{Feed Forward} + \text{Add & Norm (2)})
 \end{aligned}$$



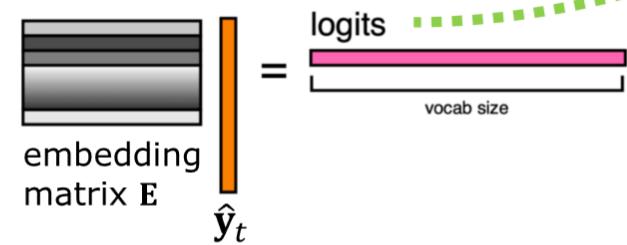
Decoder

$$\begin{aligned}
 \text{Masked Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}, \mathbf{H}_i^{dec}) \\
 \text{Add & Norm} &= \text{LayerNorm}(\text{Masked Multi-Head Attention} + \mathbf{H}_i^{dec}) \\
 \text{Enc-Dec Multi-Head Attention} &= \text{MultiHeadAtt}(\mathbf{H}_i^{enc}, \mathbf{H}_i^{enc}, \text{Add & Norm}) \\
 \text{Add & Norm (2)} &= \text{LayerNorm}(\text{Enc-Dec Multi-Head Attention} + \text{Add & Norm}) \\
 \text{Feed Forward} &= \max(0, \text{Add & Norm (2)} \mathbf{W}_1 + b_1) \mathbf{W}_2 + b_2 \\
 \text{Add & Norm (3)} &= \text{LayerNorm}(\text{Feed Forward} + \text{Add & Norm (2)}) \\
 \mathbf{H}_{i+1}^{dec} &= \text{Add & Norm (3)}
 \end{aligned}$$

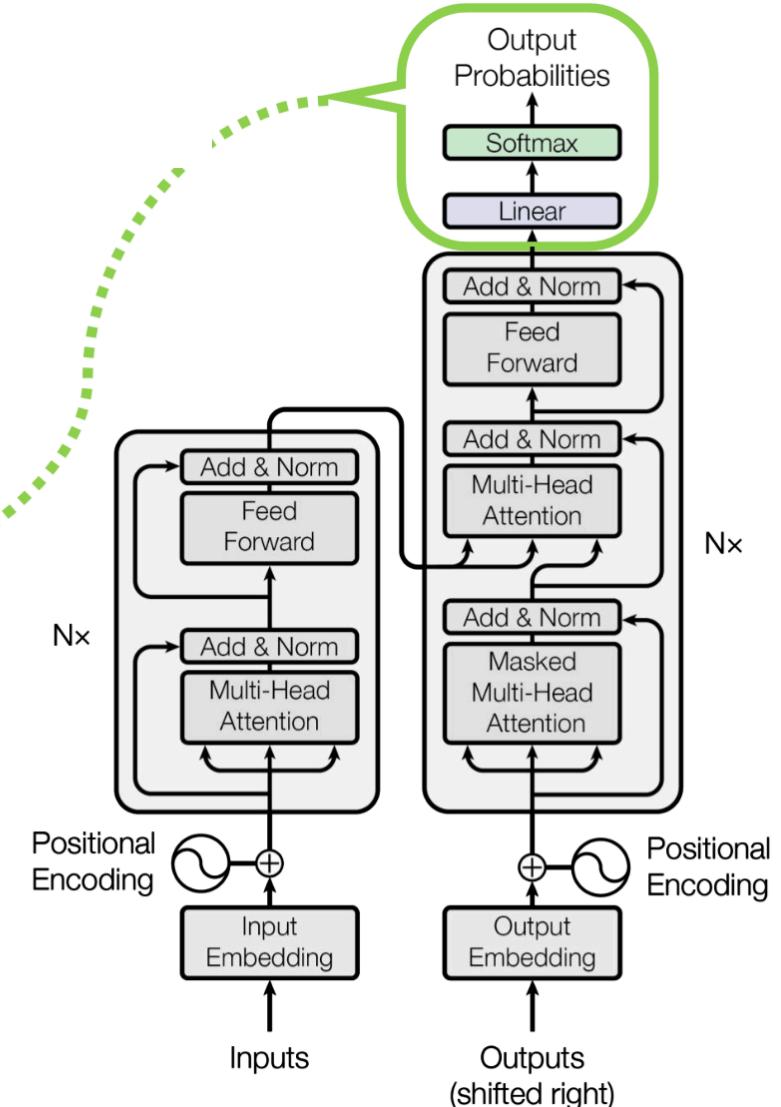




Output Probabilities



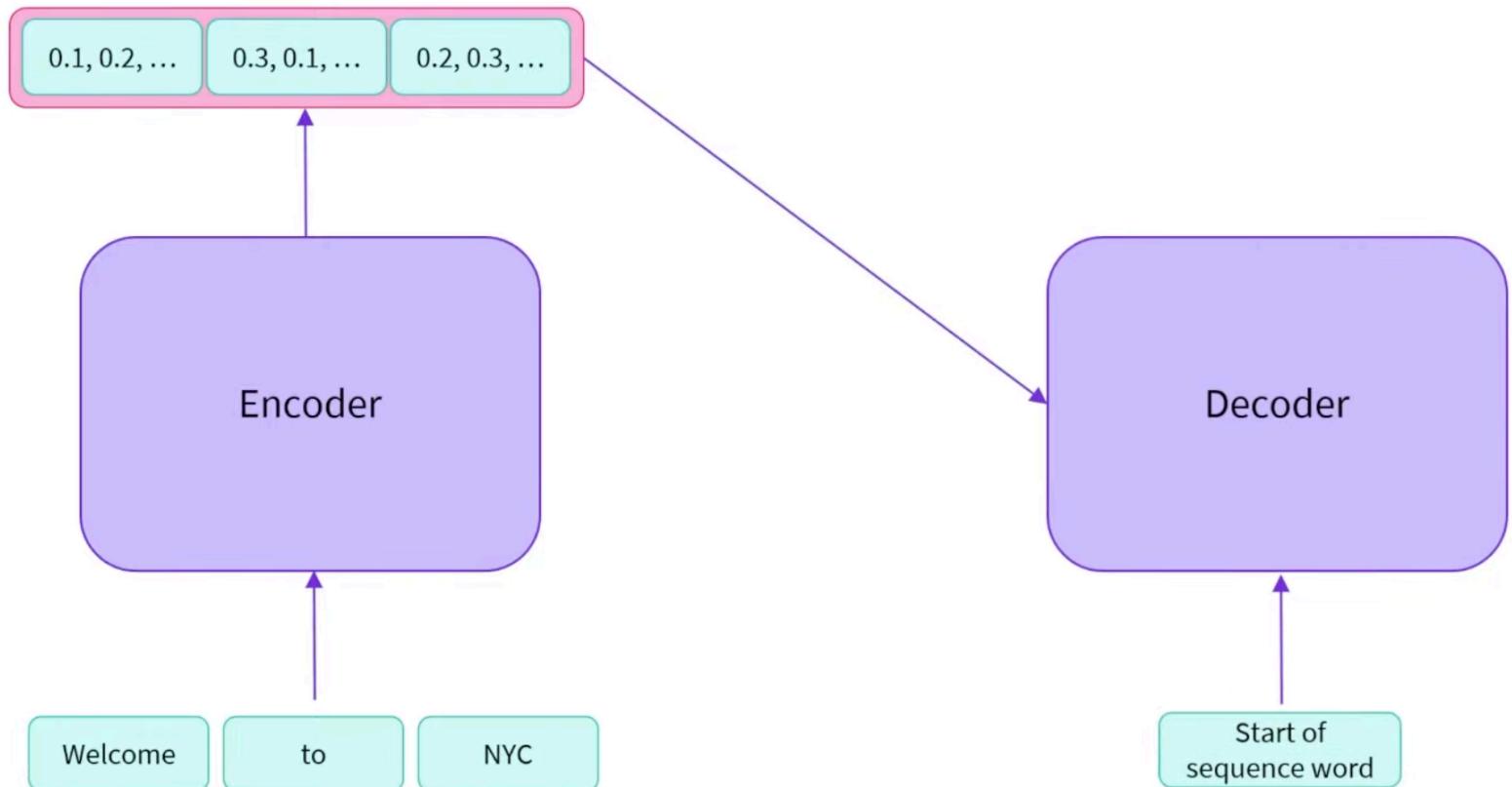
$$P(Y_t = i | \mathbf{x}_{1:T}, \mathbf{y}_{1:t-1}) = \frac{\exp(\mathbf{E}\hat{\mathbf{y}}_t[i])}{\sum_j \exp(\mathbf{E}\hat{\mathbf{y}}_t[j])}$$





Encoder-Decoder Inference

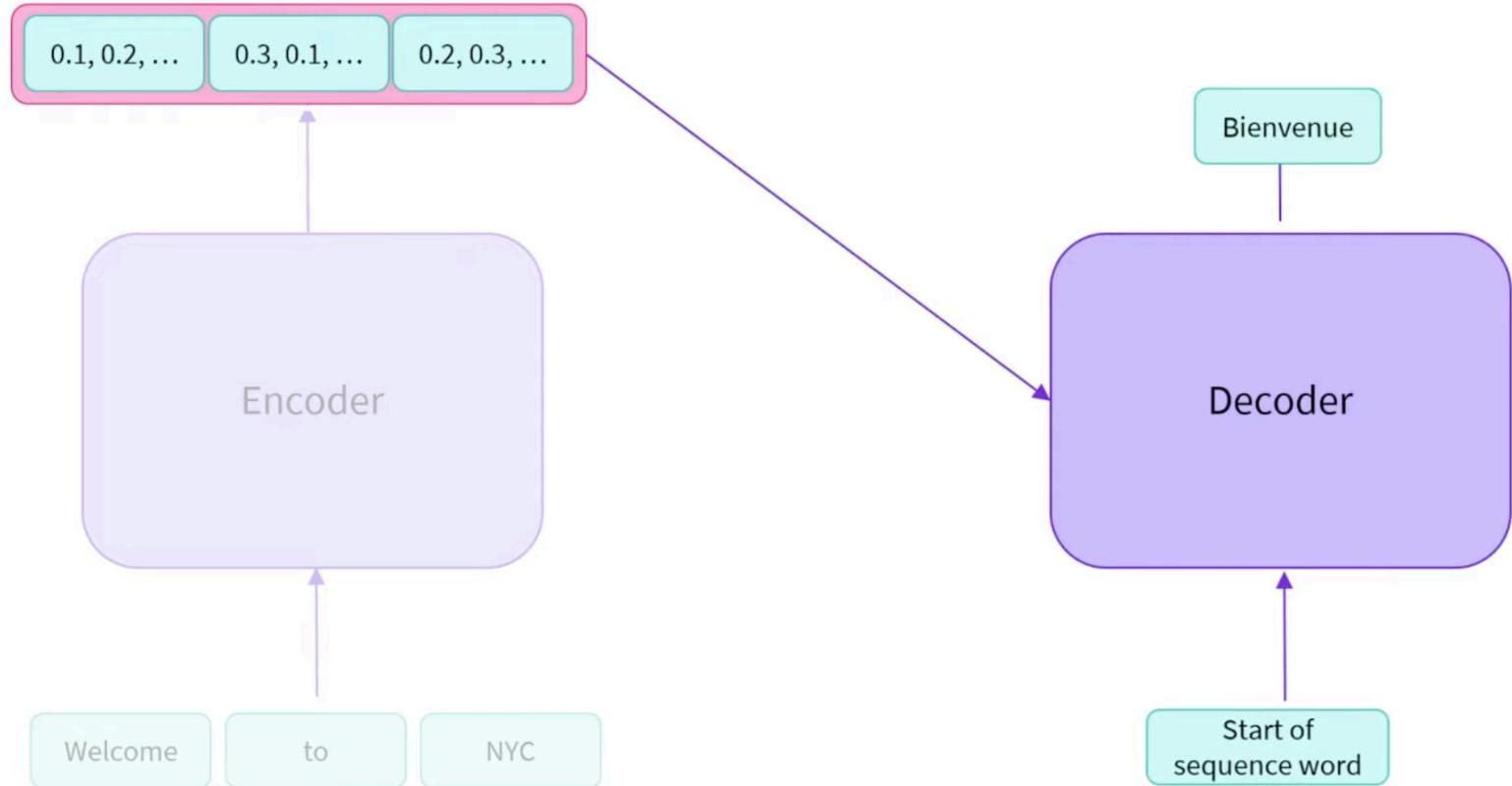
- Encode input sequence





Encoder-Decoder Inference

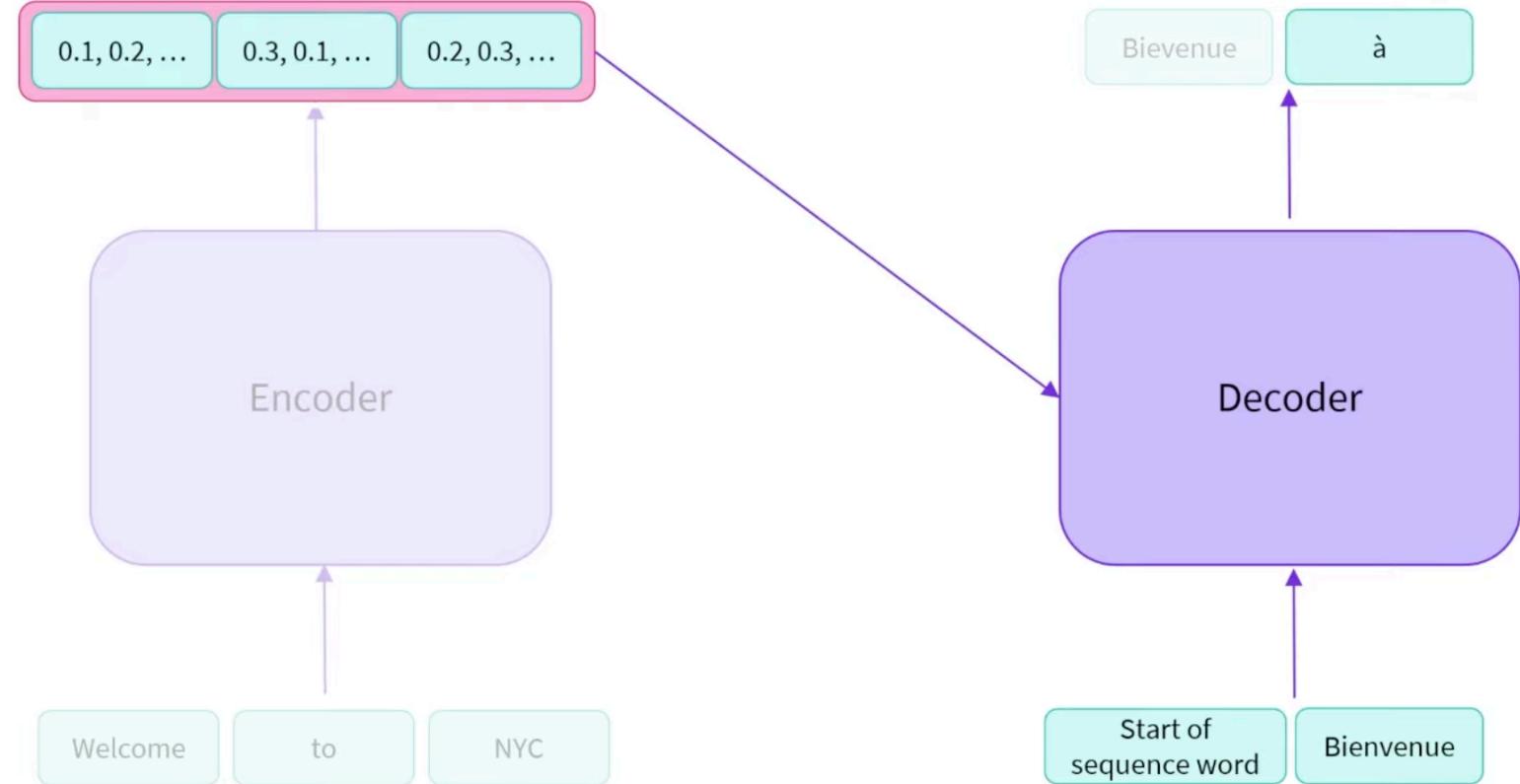
- Encode input sequence
- Attention over input token representations and <start>





Encoder-Decoder Inference

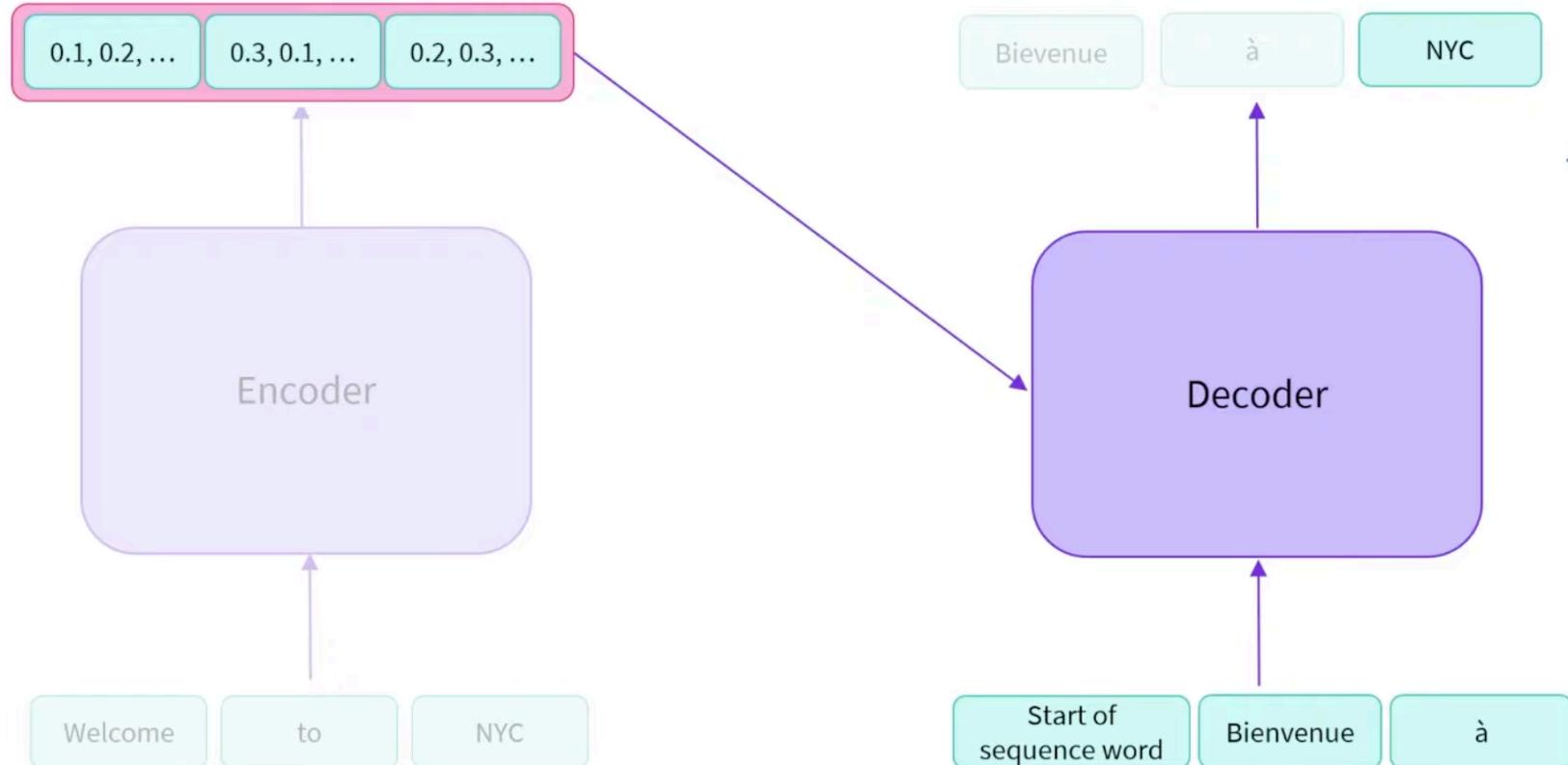
- Encode input sequence
- Attention over input token representations and <start>
- Self-attention





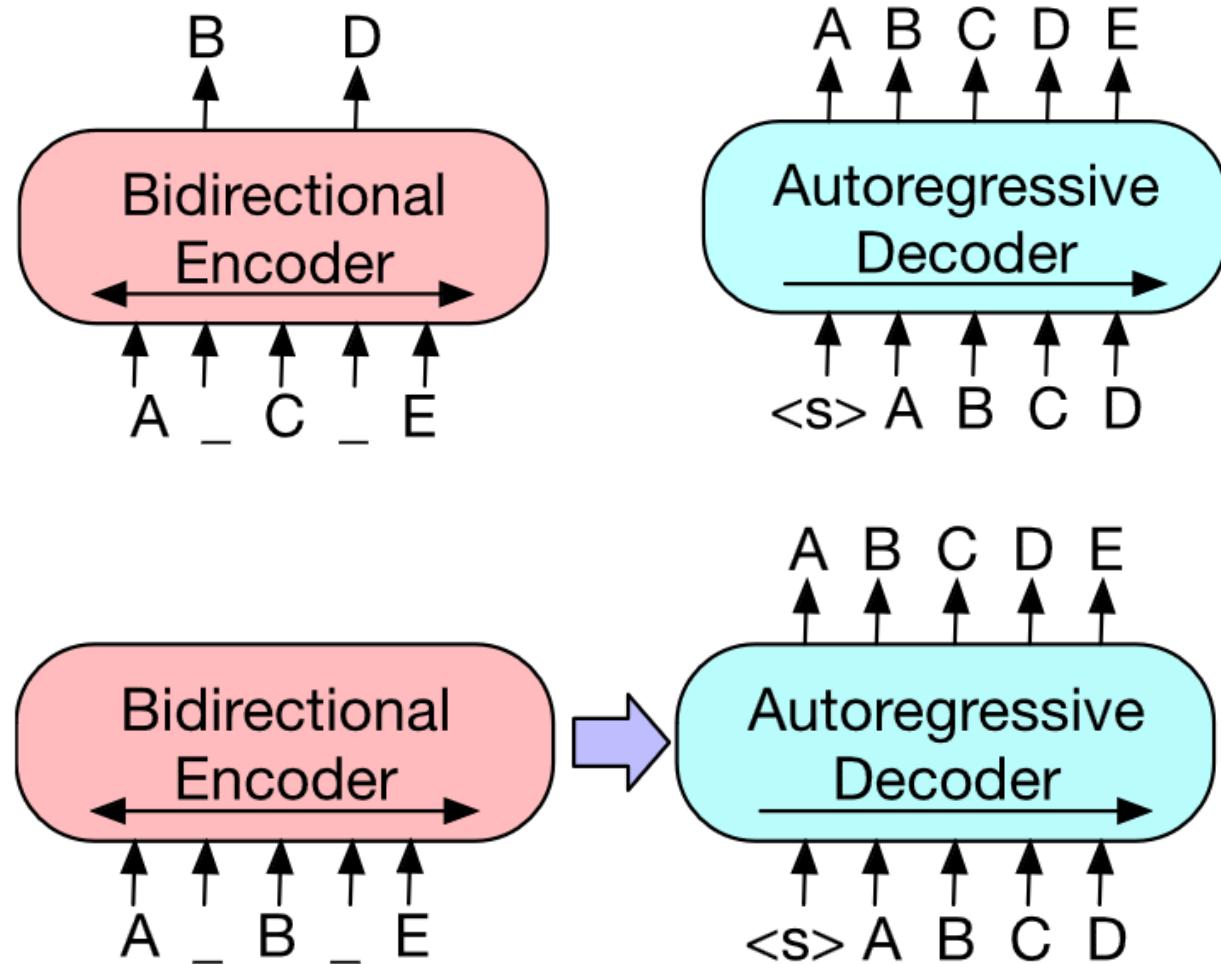
Encoder-Decoder Inference

- Encode input sequence
- Attention over input token representations and <start>
- Self-attention





Encoder, Decoder, Encoder-Decoder





Problems with the Transformer?

- Fixed context lengths “solved” with position embeddings
- Self-attention has quadratic cost $O(n^2d)$
- Plug: Annotated Transformer (Sasha Rush):
<http://nlp.seas.harvard.edu/annotated-transformer/>

Step 4: Optimization



Recap: Language Modeling Objective

- Assume we have training dataset including documents comprising sequences of bytes

$$\mathcal{D} = \left\{ \bar{d}^{(i)} \right\}_{i=1}^N \quad \bar{d} = \langle b_0, \dots, b_M \rangle$$

- Our objective is to find the LM parameters that maximize the probability of this dataset

$$\theta^* = \arg \max_{\theta} \prod_{\bar{d} \in \mathcal{D}} p(\bar{d}; \theta)$$

- We assume documents are *tokenized* into sequences that the LM models autoregressively:

$$\bar{d} = \langle x_0, \dots, x_{M'} \rangle \quad p(\bar{d}; \theta) = \prod_{j=1}^{M'} p(x_j \mid \langle x_0, \dots, x_{j-1}; \theta \rangle)$$



Recap: Language Modeling Objective

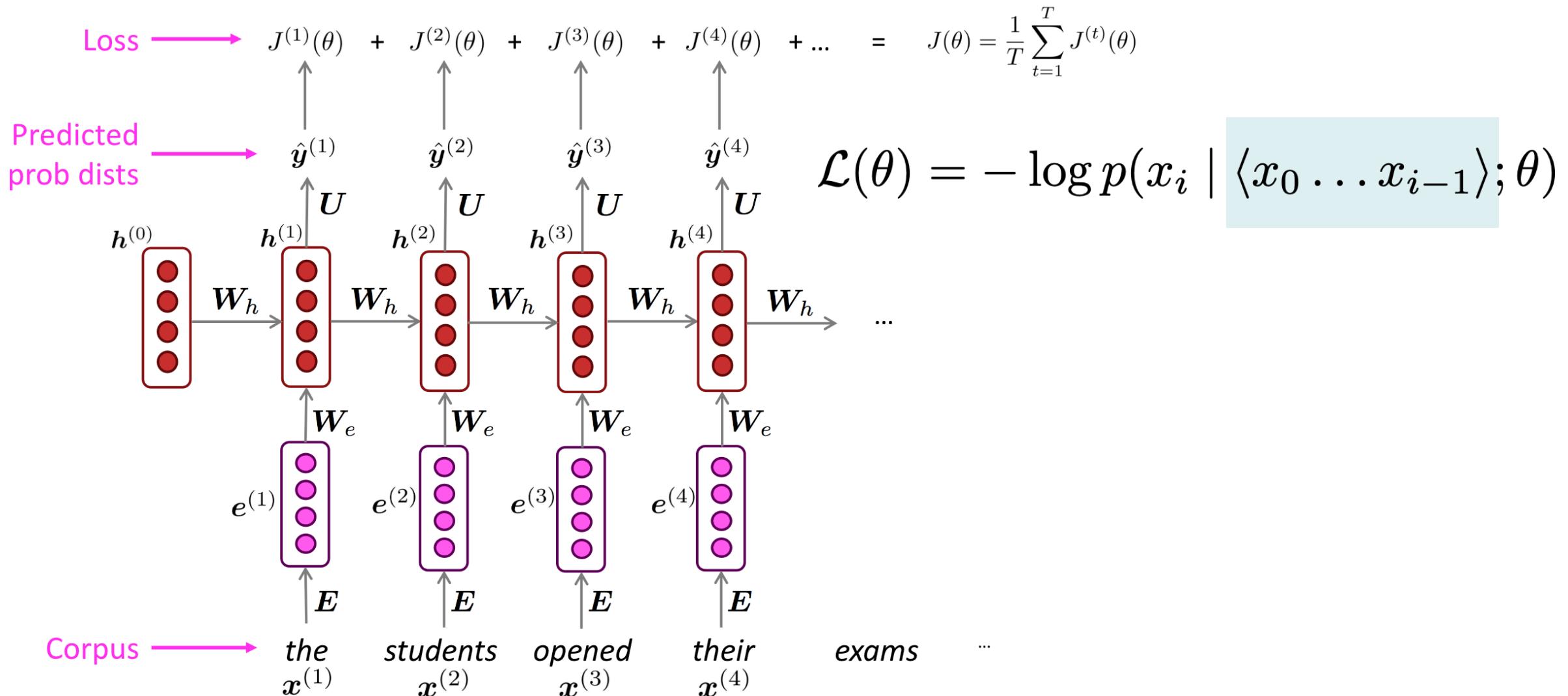
- Loss for step i is cross-entropy between true distribution p^* (i.e., one-hot) and predicted distribution:

$$\mathcal{L}(\theta) = - \sum_{x \in \mathcal{V}} p^*(x_i = x \mid \langle x_0 \dots x_{i-1} \rangle) \log p(x_i = x \mid \langle x_0 \dots x_{i-1} \rangle; \theta)$$

$$\mathcal{L}(\theta) = - \log p(x_i \mid \langle x_0 \dots x_{i-1} \rangle; \theta)$$

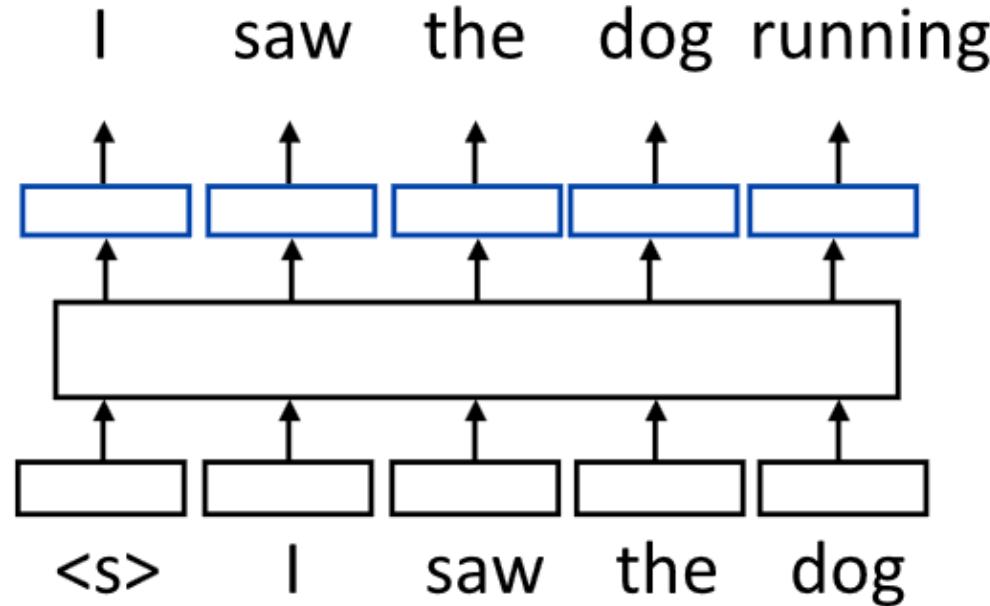


Next token prediction



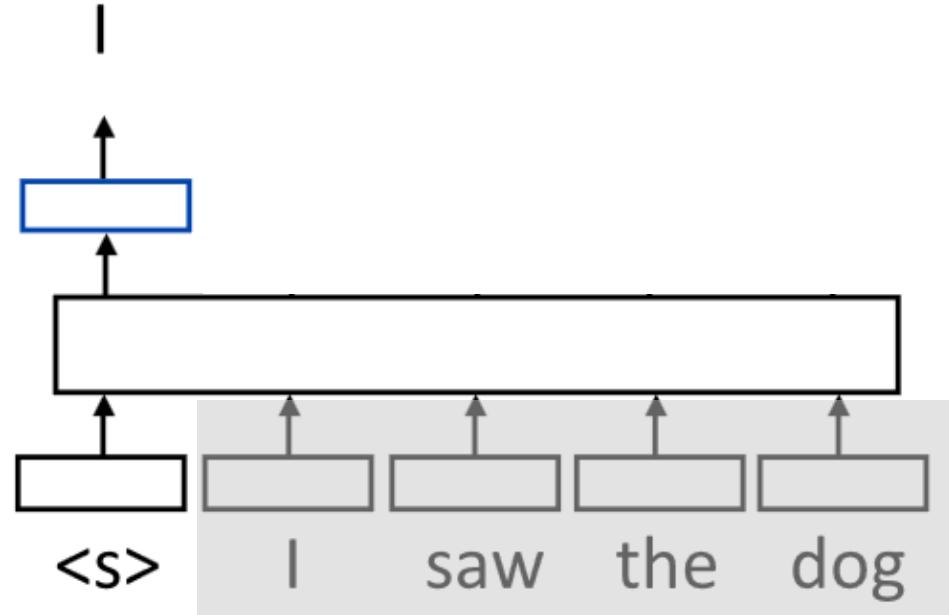


Next token prediction in Transformers



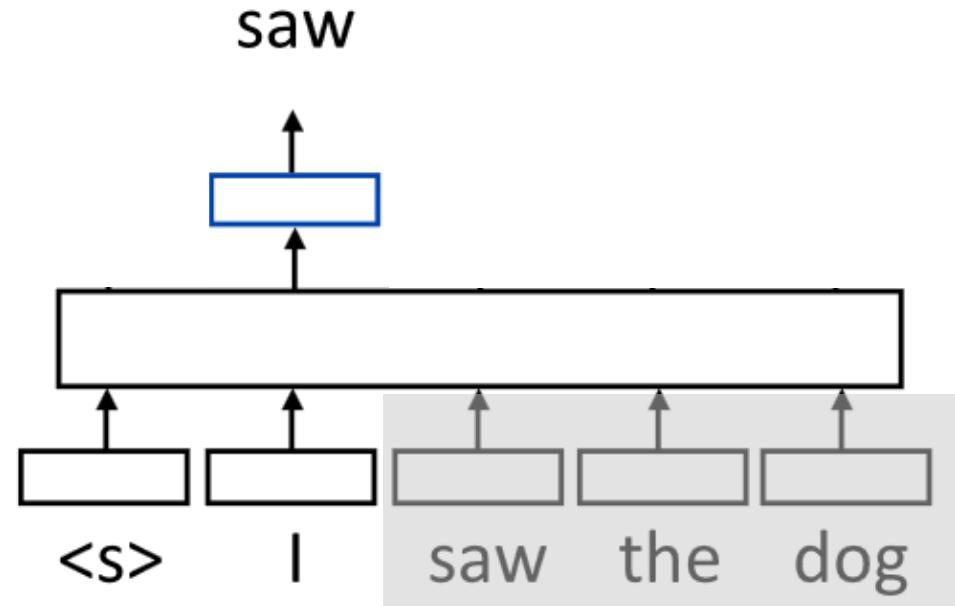


Next token prediction in Transformers



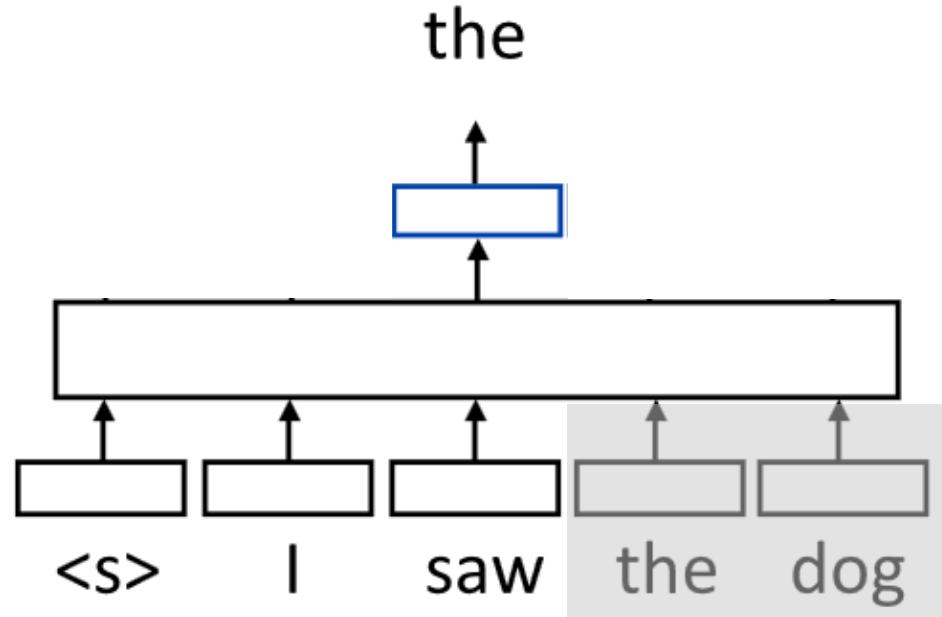


Next token prediction in Transformers



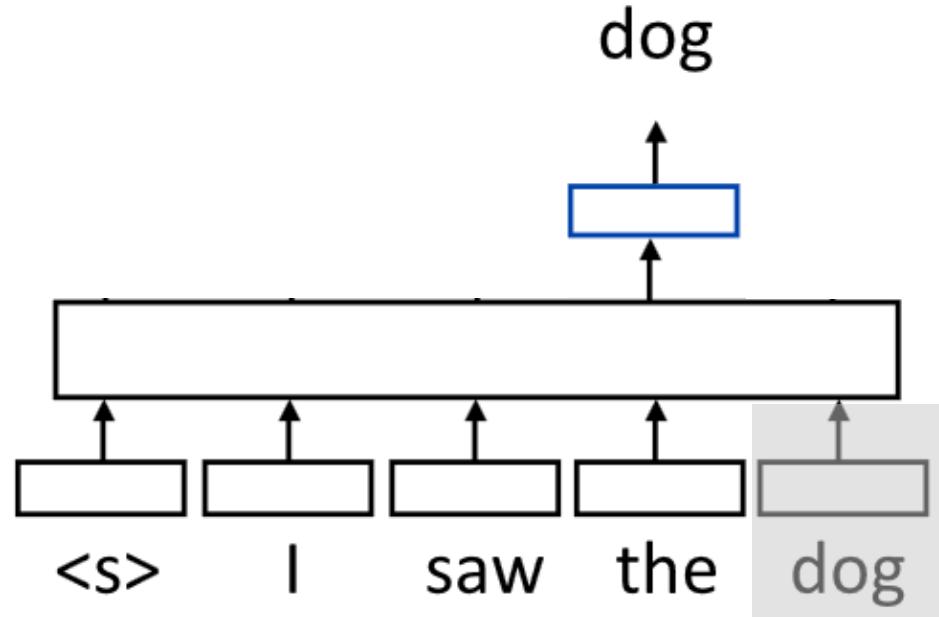


Next token prediction in Transformers



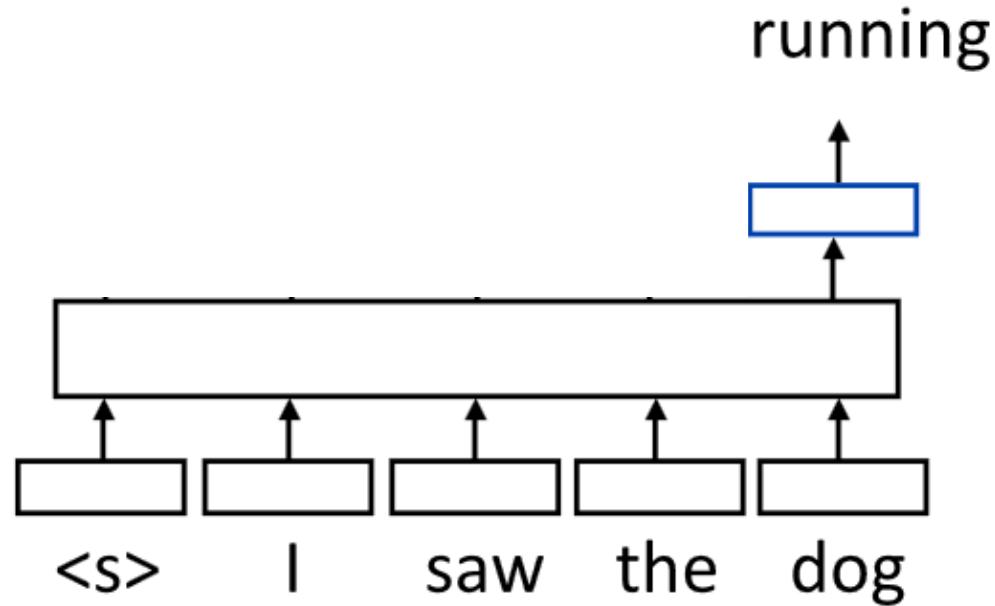


Next token prediction in Transformers





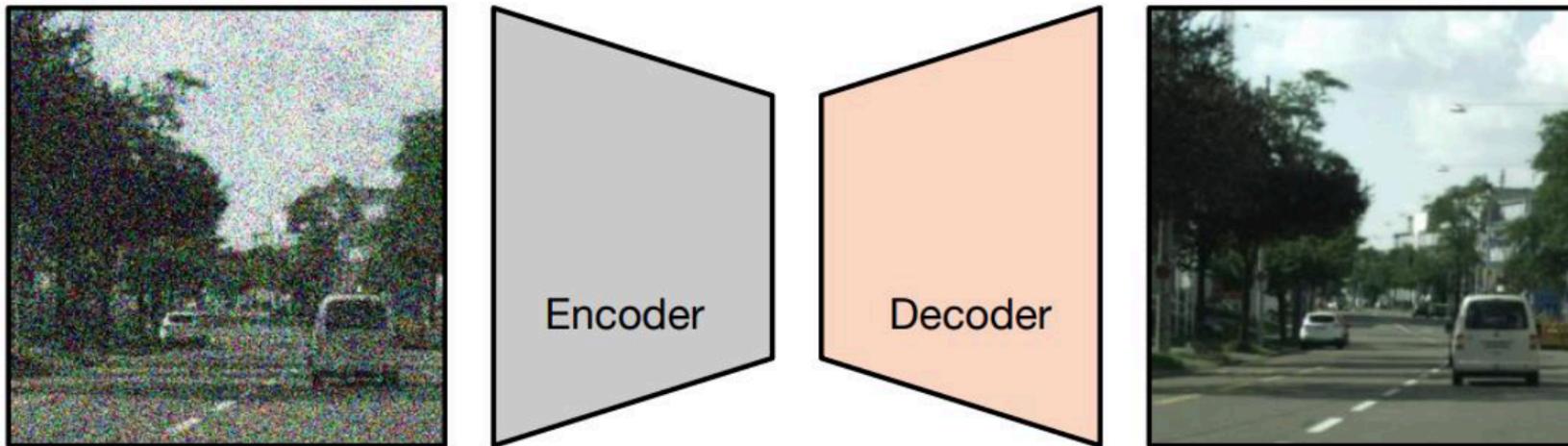
Next token prediction in Transformers





Denoising Objectives

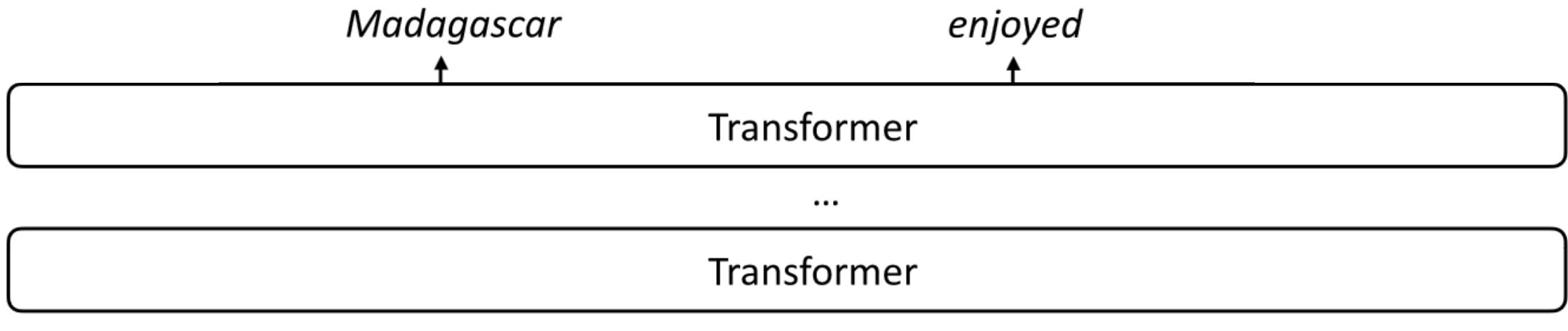
- Our goal: learn a distribution over text sequences
- Our assumption so far: this distribution is only backwards-looking (conditioned on prefix of the sequence)
- What if we remove this assumption?





Masking / Infilling Objectives

- Randomly mask out ~15% of tokens in the input, and try to predict them from past *and future* context

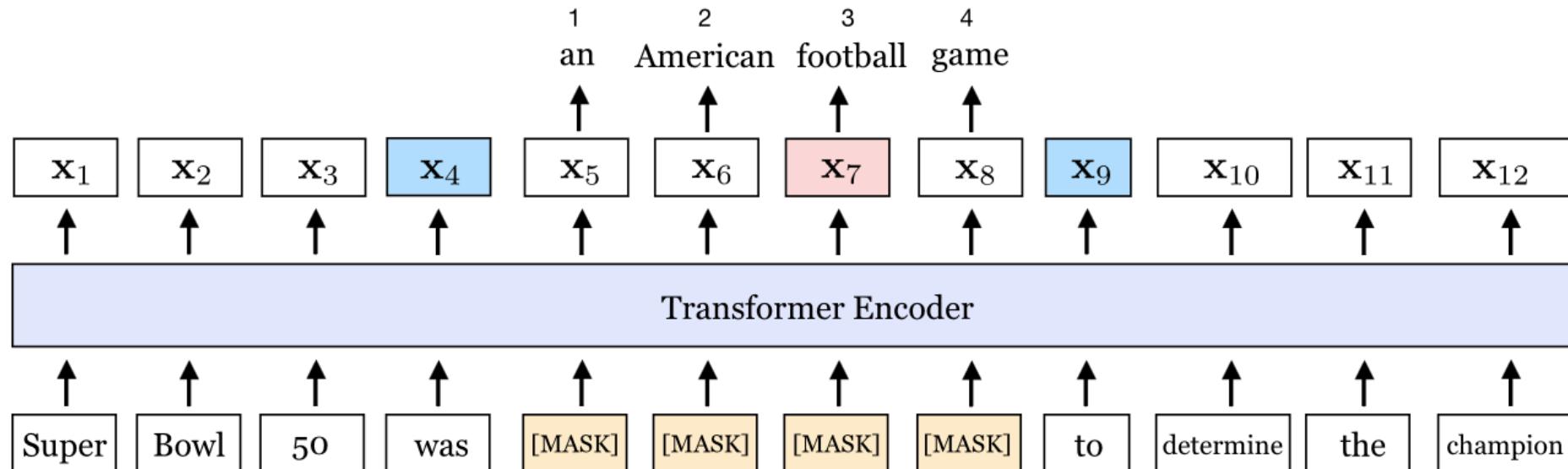


[CLS] John visited [MASK] yesterday and really [MASK] it [SEP]



Masking / Infilling Objectives

- Randomly mask out ~15% of tokens in the input, and try to predict them from past *and future* context
- Or mask out spans of text





Auxiliary Objectives

A_C._E.
Token Masking

D E . A B C .
Sentence Permutation

C . D E . A B
Document Rotation

A . C . E .
Token Deletion

A B C . D E .
Text Infilling

Step 5: Inference



Recap: What is a language model?

- Language models assign a probability to a sequence of words
- We can decompose this probability using the chain rule
- We can autoregressively generate sequences from the language model by sampling from its token-level probability
- We can condition on our language distribution on something else

$$p(\bar{y})$$

$$p(\bar{y}) = \prod_{i=1}^T p(y_i | y_{0:i-1})$$

$$p(y_i | y_{0:i-1})$$

$$p(y_i | y_{0:i-1}; \bar{x})$$



What can we do with language models?

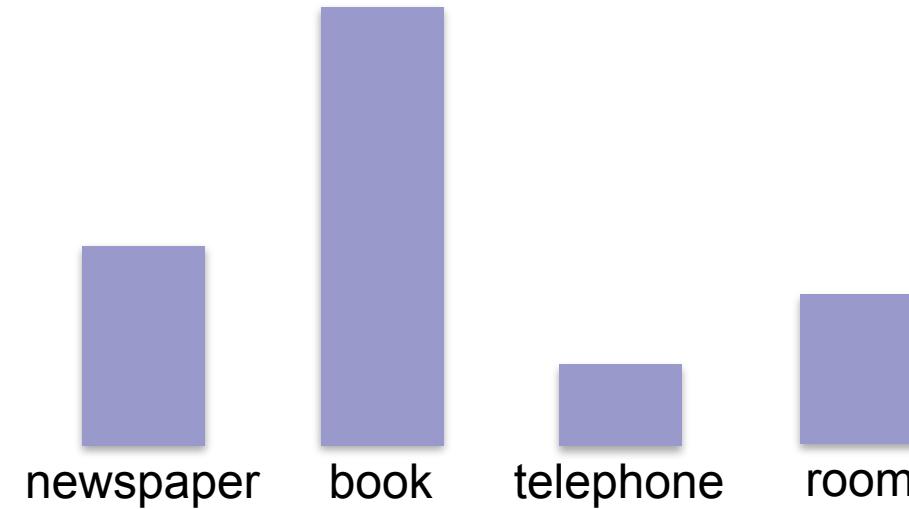
- Computing probabilities of a sequence
- Autoregressive sequence generation



Decoding strategies

- Argmax (greedy decoding)

$$y_T = \arg \max_{y \in \mathcal{V}} p(y \mid y_{0:t-1})$$

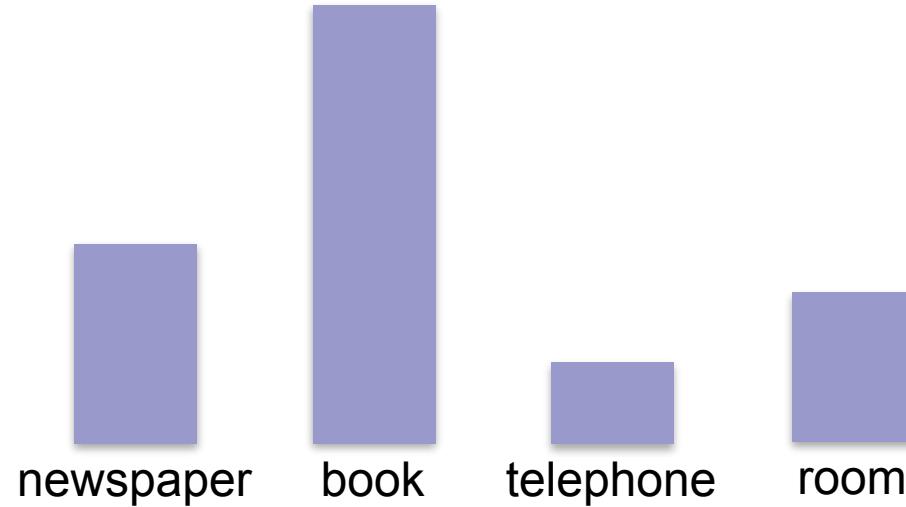




Decoding strategies

- Argmax (greedy decoding)
- Sampling from language model directly

$$y_T = \arg \max_{y \in \mathcal{V}} p(y \mid y_{0:t-1})$$
$$y_T \sim p(\cdot \mid y_{0:t-1})$$





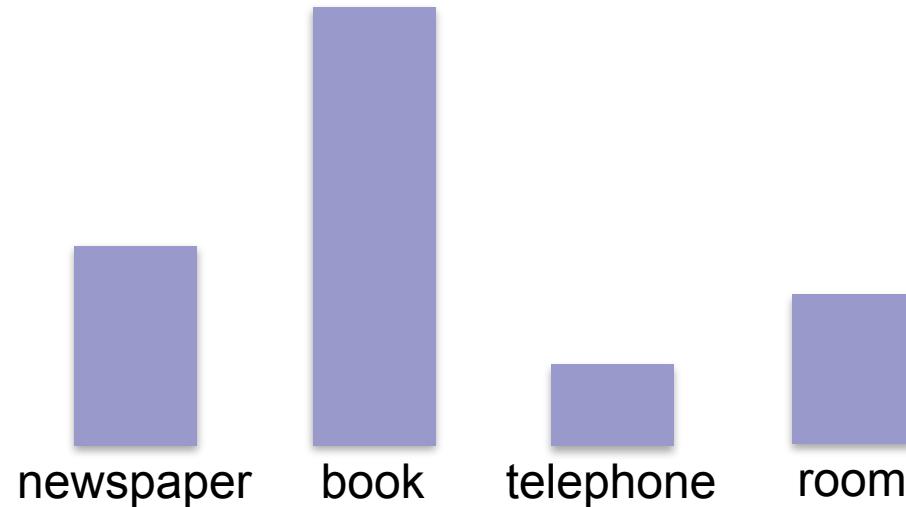
Decoding strategies

- Argmax (greedy decoding)
- Sampling from language model directly
- Adjusting temperature of distribution

$$y_T = \arg \max_{y \in \mathcal{V}} p(y \mid y_{0:t-1})$$

$$y_T \sim p(\cdot \mid y_{0:t-1})$$

$$p'(y_T = y) = \frac{\exp(z_y/T)}{\sum_{y' \in \mathcal{V}} (z_{y'}/T)}$$





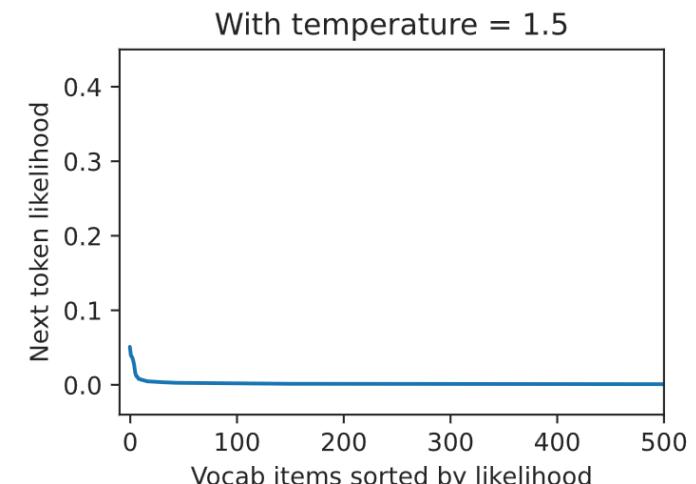
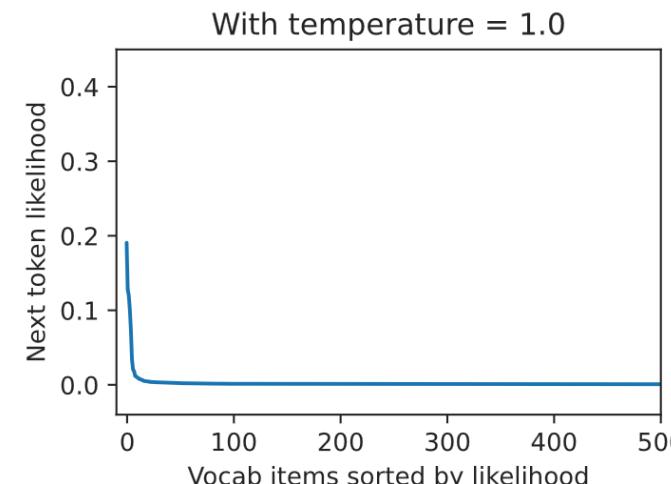
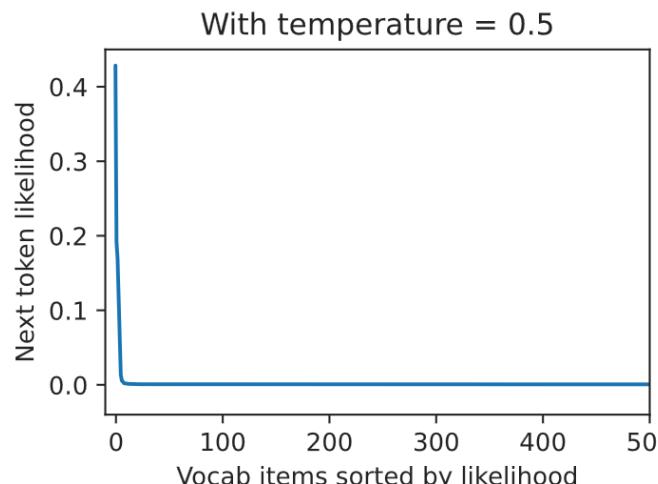
Decoding strategies

- Argmax (greedy decoding)
- Sampling from language model directly
- Adjusting temperature of distribution

$$y_T = \arg \max_{y \in \mathcal{V}} p(y \mid y_{0:t-1})$$

$$y_T \sim p(\cdot \mid y_{0:t-1})$$

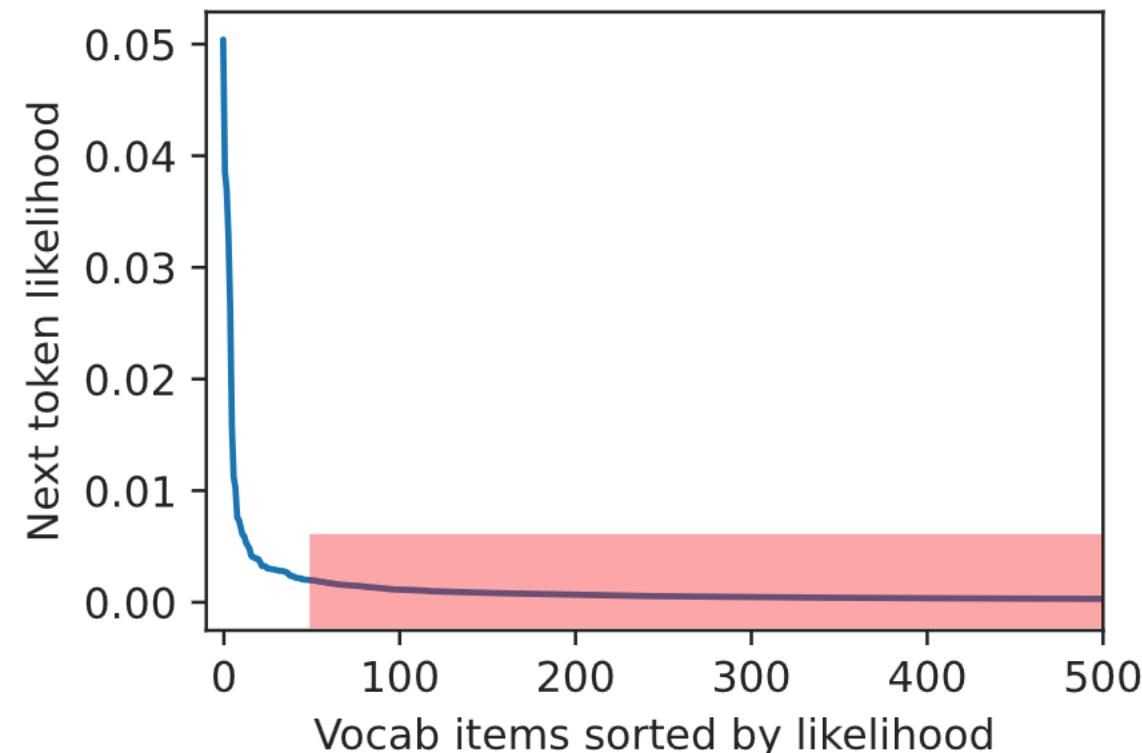
$$p'(y_T = y) = \frac{\exp(z_y/T)}{\sum_{y' \in \mathcal{V}} (\exp(z_{y'}/T))}$$





Decoding strategies

- Top- k sampling: reassign probability mass from all but the top k tokens to the top k tokens





Decoding strategies

- Nucleus sampling: reassign probability mass to the most probable tokens whose cumulative probability is at least p

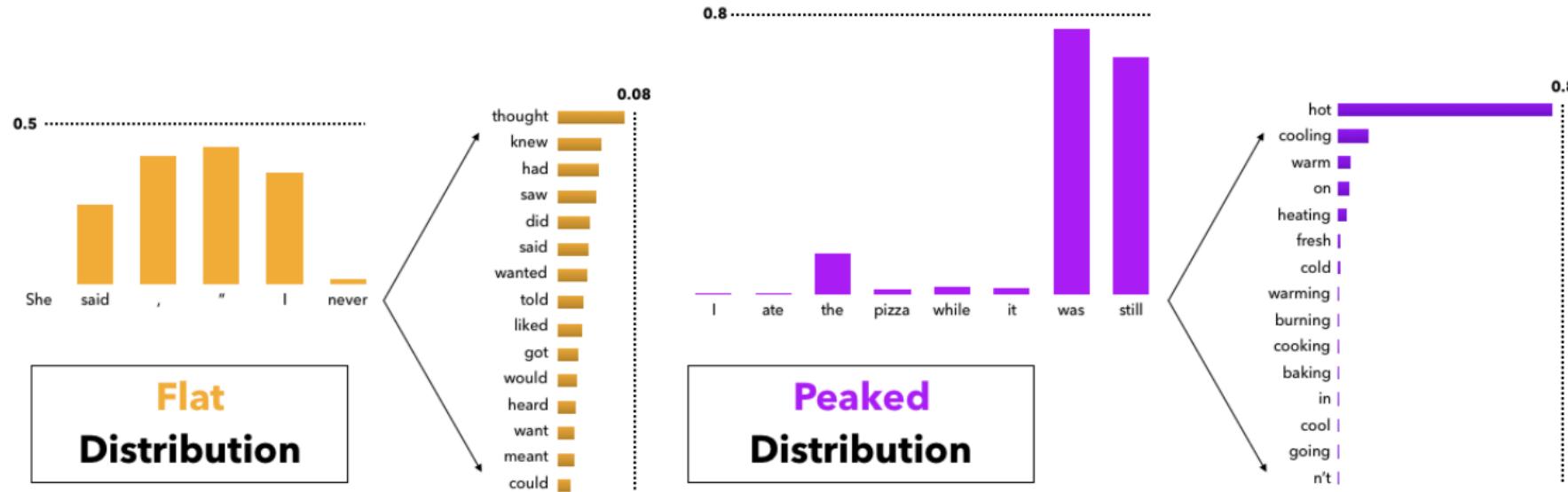


Figure 5: The probability mass assigned to partial human sentences. Flat distributions lead to many moderately probable tokens, while peaked distributions concentrate most probability mass into just a few tokens. The presence of flat distributions makes the use of a small k in top- k sampling problematic, while the presence of peaked distributions makes large k 's problematic.



Beam search

- It's intractable to find the *most probable sequence* according to a language model

$$\bar{y}^* = \arg \max_{\bar{y} \in \mathcal{V}^*} p(\bar{y})$$

- Greedy search doesn't yield the most probable sequence
- Instead: beam search

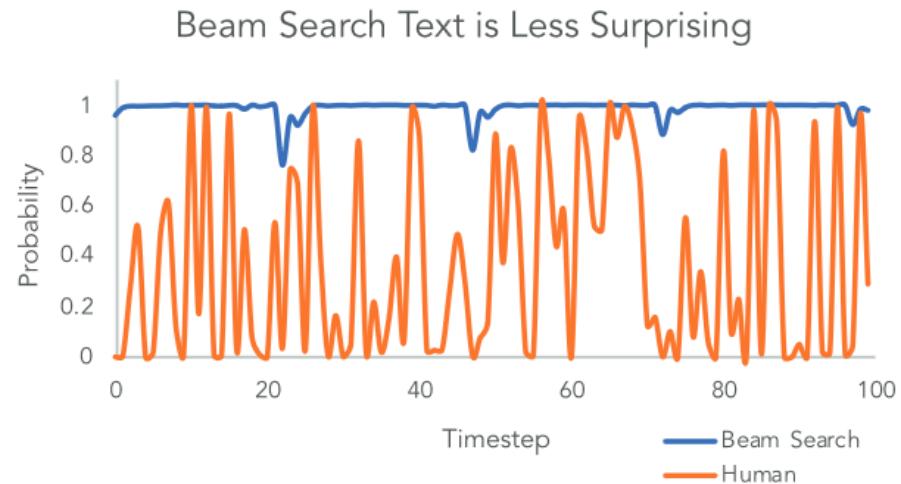
$$y_t = \arg \max_{y \in \mathcal{V}} p(y \mid y_{0:t-1})$$

- Approximate the search by keeping around candidate continuations
- At the end, choose the highest probability sequence in the beam



Beam search

- But do we even want to find the highest-probability sequence according to a LM?
- Human language is noisy and surprising
- Optimizing for LM probability leads to repetitive and uninteresting text



Beam Search

...to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and...

Human

...which grant increased life span and three years warranty. The Antec HCG series consists of five models with capacities spanning from 400W to 900W. Here we should note that we have already tested the HCG-620 in a previous review and were quite satisfied With its performance. In today's review we will rigorously test the Antec HCG-520, which as its model number implies, has 520W capacity and contrary to Antec's strong beliefs in multi-rail PSUs is equipped...



Beam search

- But do we even want to find the highest-probability sequence according to a LM?
- Human language is noisy and surprising
- Optimizing for LM probability leads to repetitive and uninteresting text

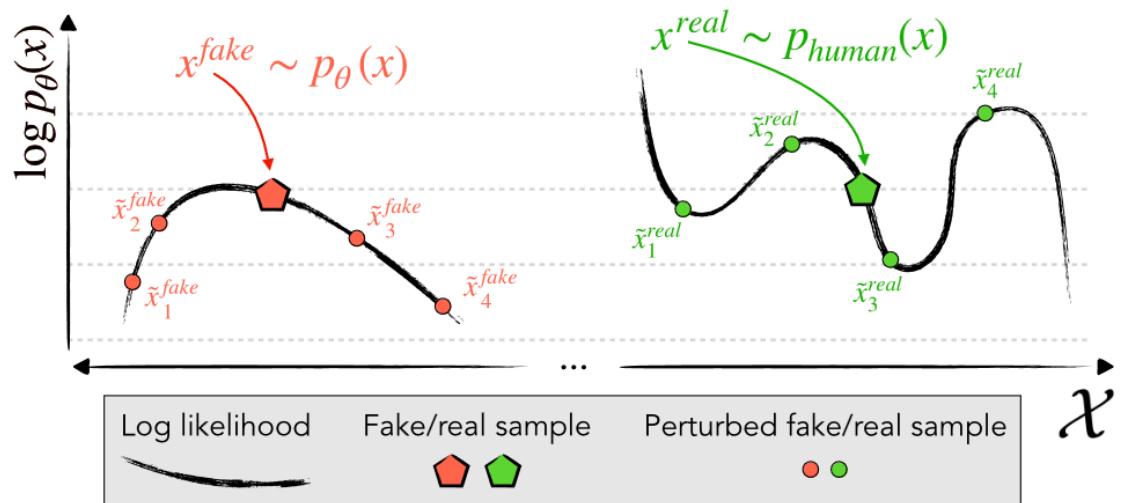


Figure 2. We identify and exploit the tendency of machine-generated passages $x \sim p_\theta(\cdot)$ (**left**) to lie in negative curvature regions of $\log p(x)$, where nearby samples have lower model log probability on average. In contrast, human-written text $x \sim p_{real}(\cdot)$ (**right**) tends not to occupy regions with clear negative log probability curvature; nearby samples may have higher or lower log probability.