

An Overview of Natural Language Processing

Sina Semnani

CS294S October 15, 2020

Some slides are adopted from Giovanni Campagna

These Words Are Not Mine!

Me: In today's lecture, I will talk about natural language processing research.

GPT-2: Natural language processing is the science that helps computers to understand human language. So what is NLP? Well, let me give you the definition I use from my textbook Natural Language Processing. “Natural language processing is the process of finding and transforming words in a source document in order to uncover structured features of the document.” Sounds pretty technical, but I think it is important to understand what is meant by natural language processing. Why does that matter?

Outline

- Introduction
- Deep Learning for Natural Language Processing
- Word Representation
 - One-hot
 - Dense
 - Language Models
 - Contextual
- Sequence to Sequence
- Attention
- GPT-3

Natural Language Processing

- How do we map from text to integers/real numbers/text
- Examples:
 - Text Classification
 - Question Answering
 - Semantic Parsing

Input

Paragraph: ... With a population of 3,792,621,

the writer-director has made a film so

Show me Chinese restaurants in Palo Alto.

the heart soar.

Question: What is the population of Los Angeles?

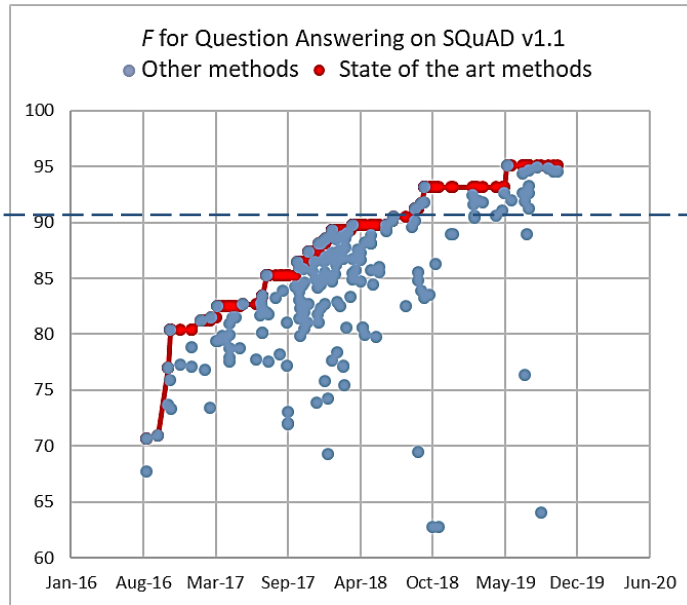
Output

```
now => @QA.restaurant(), geo ==
makeLd
serves Ar +1 (positive) 1 se" =>
notify
```

NLP Has Been Especially Successful in Recent Years

- Even “super-human”, according to some benchmarks for Question Answering, Natural Language Inference, etc.

“Human” performance is 90.5%



But Not Entirely ...

- Reported human performance can be misleading
- These models are very fragile and lack common sense
- Some adversarial tests result in a 2-10x accuracy drop while humans are unaffected

Paragraph: Its counties of Los Angeles, Orange, San Diego, San Bernardino, and Riverside are the five most populous in the state and all are in the top 15 most populous counties in the United States.

Question: What is the smallest geographical region discussed?

Answer: donald trump

But Not Entirely ...

- Besides, we have not even come close to humans on many other tasks
 - Understanding nontrivial dialogues
 - Multilingual tasks and low-resource languages
 - Empathetic text generation
 - Advice giving
 - Common sense
 - ...

Even with 175 Billion Parameters ...

GPT-3:

Human: Are married bachelors impossible?

GPT-3: No, married bachelors are not impossible.

Human: Why are married bachelors possible?

GPT-3: Because the concept of being married is not part of the concept of being a bachelor.

Neural Networks for Natural Language Processing

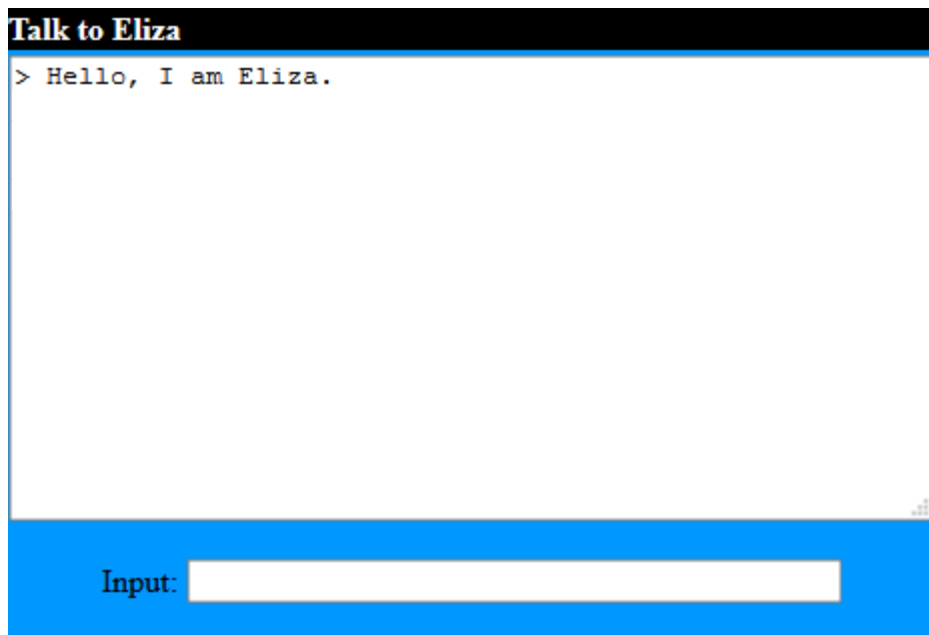


Before Deep Learning for Natural Language

- NLP research was focused on rule-based approaches for a very long time
- 1960s: ELIZA
 - one of the first conversational systems
 - matched keywords and repeated the user

Before Deep Learning for Natural Language

- My existential discussion with ELIZA last night:



Deep Learning for Natural Language

- NLP research was focused on rule-based approaches for a very long time
- 1960s: ELIZA
 - one of the first conversational systems
 - matched keywords and repeated the user
- ...
- Rapid increase in the amount of available digital text and computational power has made deep learning a very suitable tool for natural language processing
- Today, almost all systems that process human language have a machine learning component and learn from large amounts of data

Machine Learning

- Arthur Samuel (1959): Machine Learning is the field of study that gives the computer the ability to learn without being explicitly programmed.
- Instead, we show the computer a lot of examples of the desired output for different inputs.

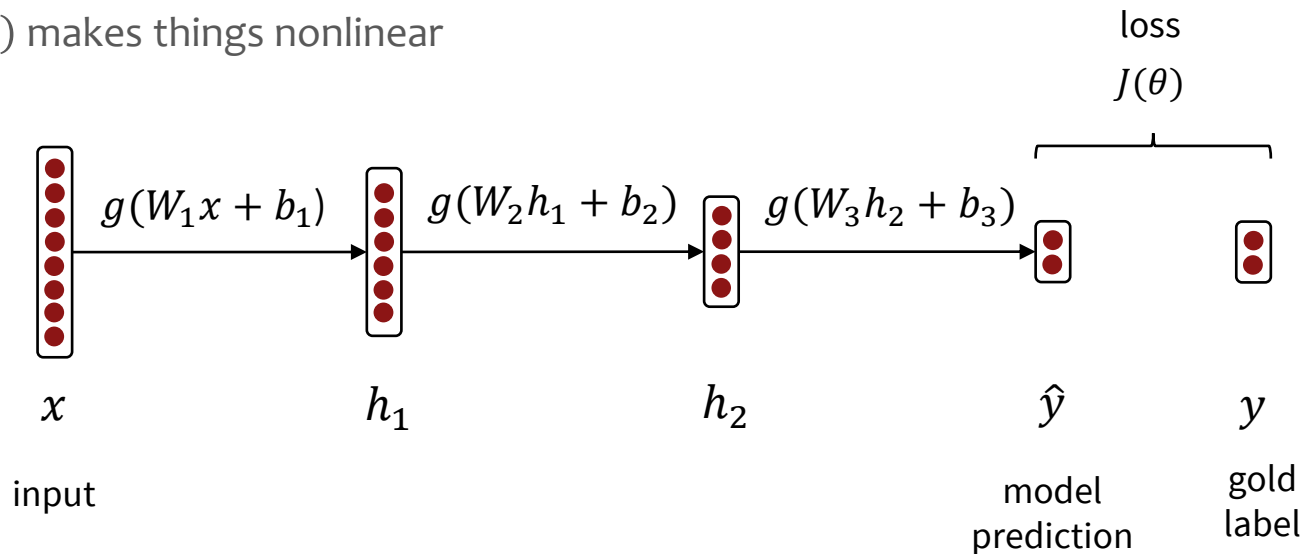
Machine Learning

- The goal is to learning a parametrized function
- The parametrized function can have various shapes:
 - Logistic Regression
 - Support Vector Machines
 - Decision Trees
 - **Neural Networks**
- Inputs and outputs can be many different things:

<ul style="list-style-type: none">• Text• Image• Integer• $y \in \mathbb{R}^m$• ...	To	<ul style="list-style-type: none">• Text• Image• Integer• $y \in \mathbb{R}^n$• ...
---	----	--

Deep Learning

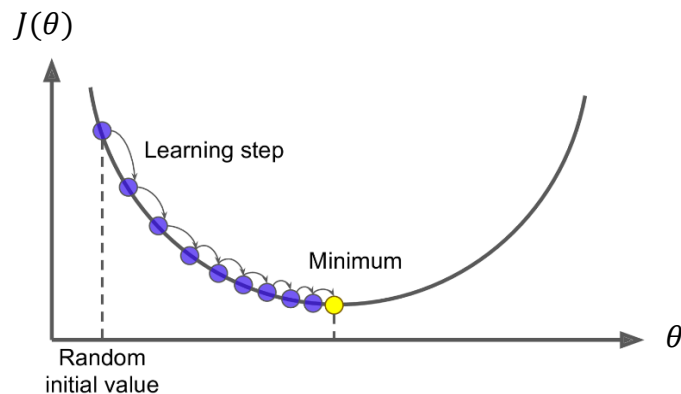
- The parametrized function is a combination of smaller functions
- Example: Feedforward Neural Network
 - An input vector x goes to output vector y using a combination of functions of the form $\text{output} = g(W \times \text{input} + b)$
 - $g(\cdot)$ makes things nonlinear



Loss Function and Gradient Descent

- Calculate gradient of loss with respect to parameters
- Iteratively update parameters to minimize loss

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta)$$

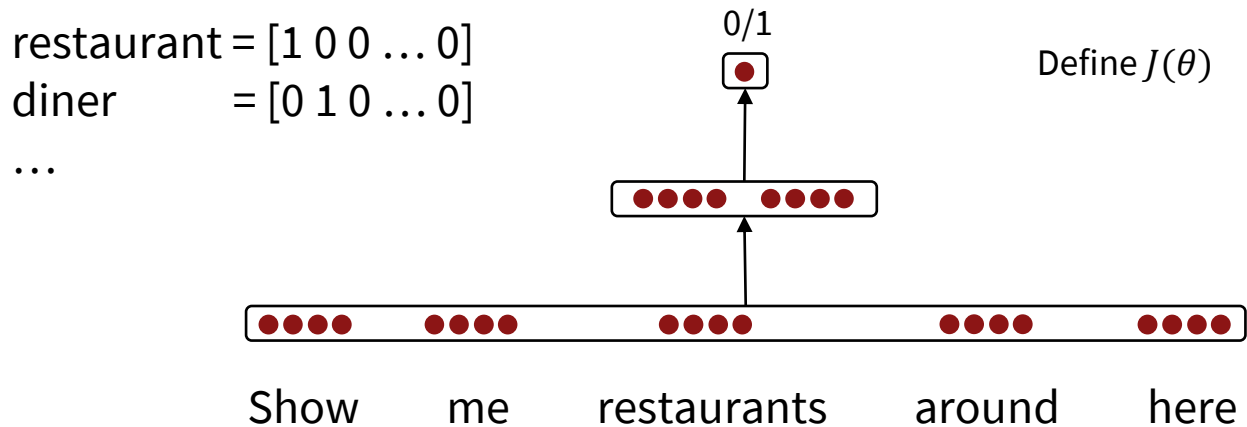


Text Representation



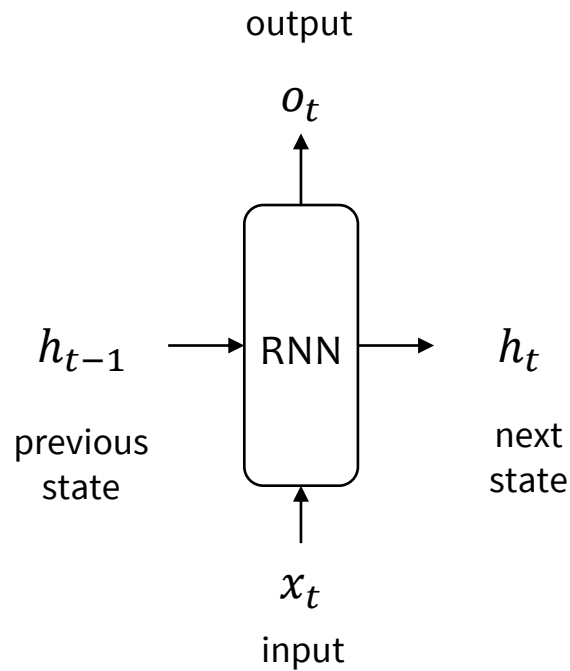
Word Representation: One-Hot Vectors

- We have a calculus for functions that are from R^n to R^m
- So we have to convert everything to vectors
- Consider the simple task of domain detection: 0 means is restaurants skill, 1 means everything else



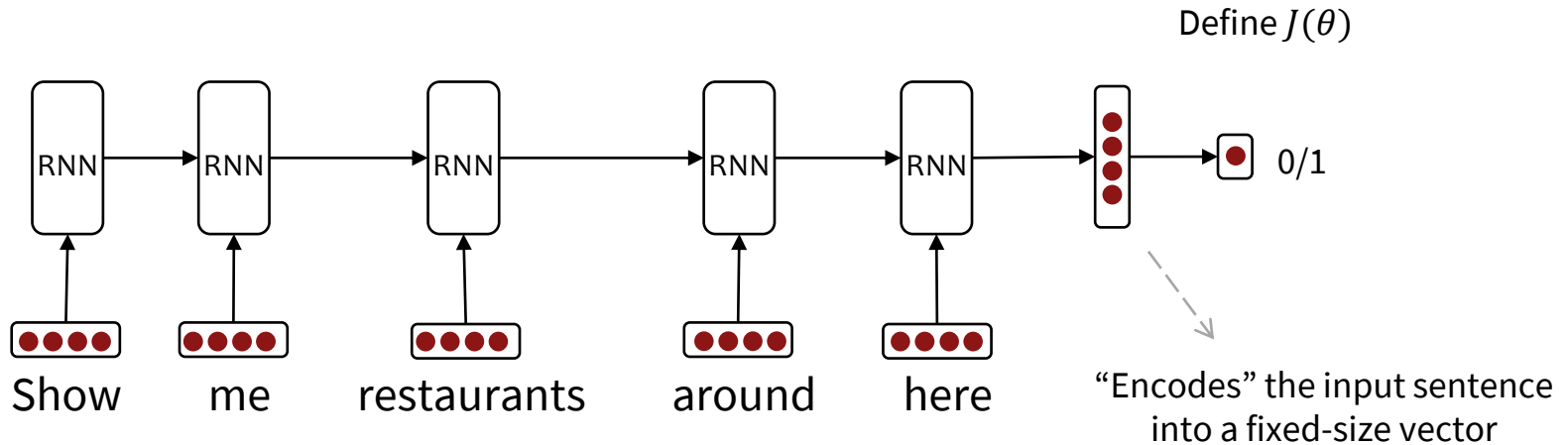
Sequence Representation: Recurrent Neural Networks

- $h_t, o_t = RNN(x_t, h_{t-1}; \theta)$
- θ is the learned parameters
- Various types of cells:
 - Gated Recurrent Unit (GRU)
 - Long Short-Term Memory (LSTM)



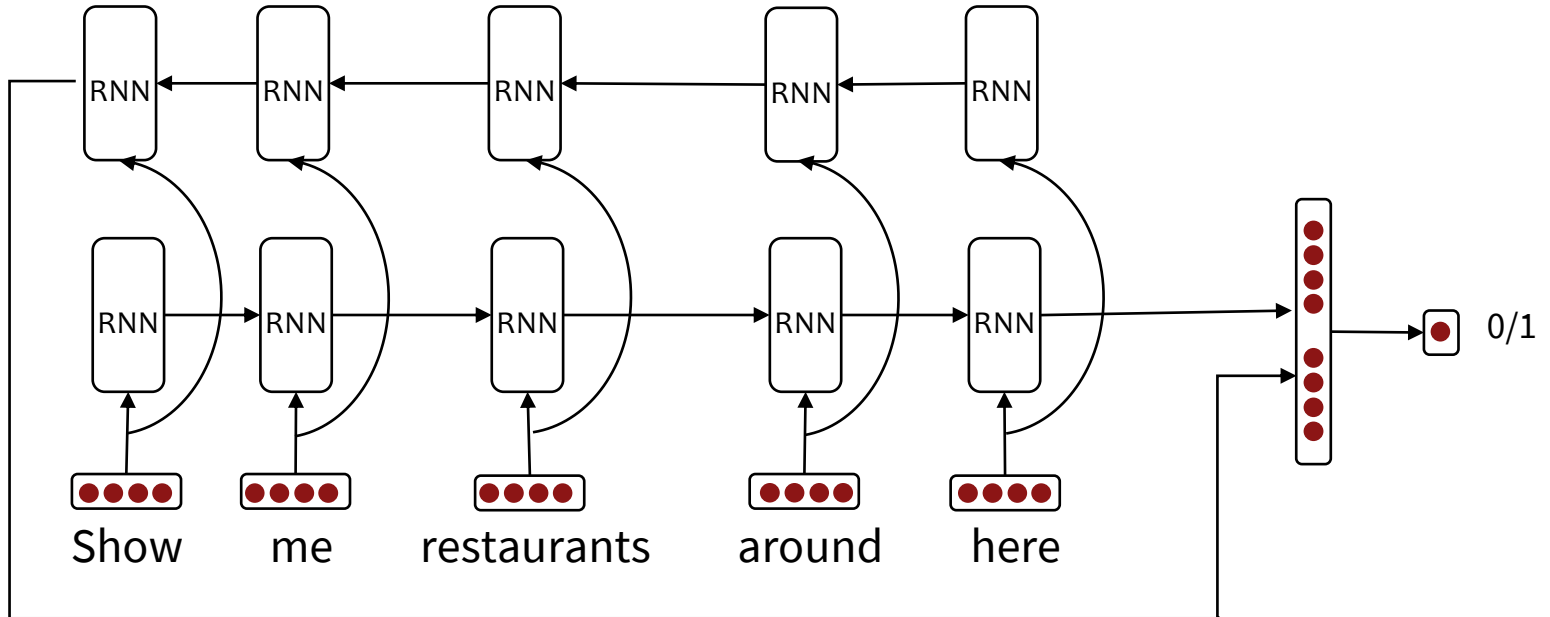
Encode Sequences

- Recurrent: repeat the same box, with the same θ for each word in the sequence



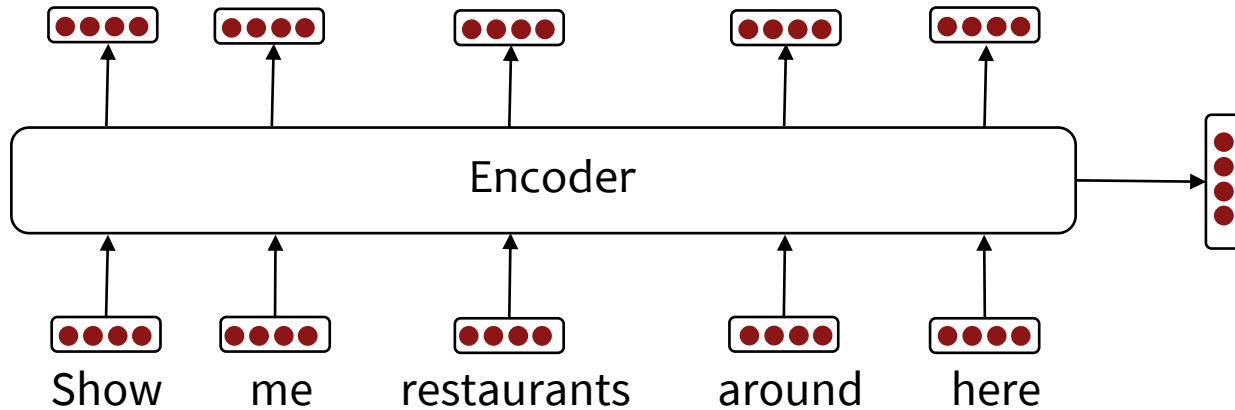
Encode Sequences

- It can be Bi-directional



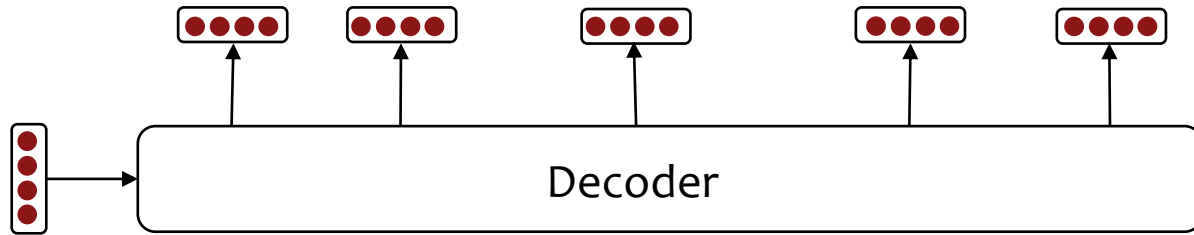
Encoder

Converts a sequence of inputs to one or more fixed size vectors



Decoder

Receives a fixed size vector and produces probability distributions over words, i.e. vectors of size $|V|$ whose elements sum to 1



Quiz

In the assignment, the goal was to build a system that can convert natural sentences to their corresponding ThingTalk programs.

You trained a semantic parser for this task.

Do you think you used one-hot encoding for word representations?
Why or Why not?

No. Just to name a few limitations of one-hot encoding:
Large size of input would result in inefficient computations.
Words with similar meanings would have nothing in common.

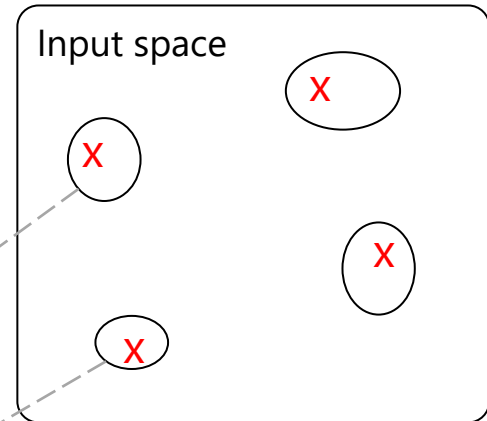
The Effect of Better Embeddings

- During training, neural networks learn to map regions of the input space to specific outputs
- If word embeddings map similar words to similar regions, the neural network will have an easier job

restaurant = [1 0 0 ... 0]
diner = [0 1 0 ... 0]
...

These sentences are in the restaurants domain

These are in the hotels domain



Word Representation: Dense Vectors

- Also called Distributed Representation
- In practice, ~ 100 - 1000 dimensional vectors (much smaller than $|V|$)
- Learned from large text corpora

I went to this amazing restaurant last night.

We were at the diner when we saw him.

Ali went to the movies.

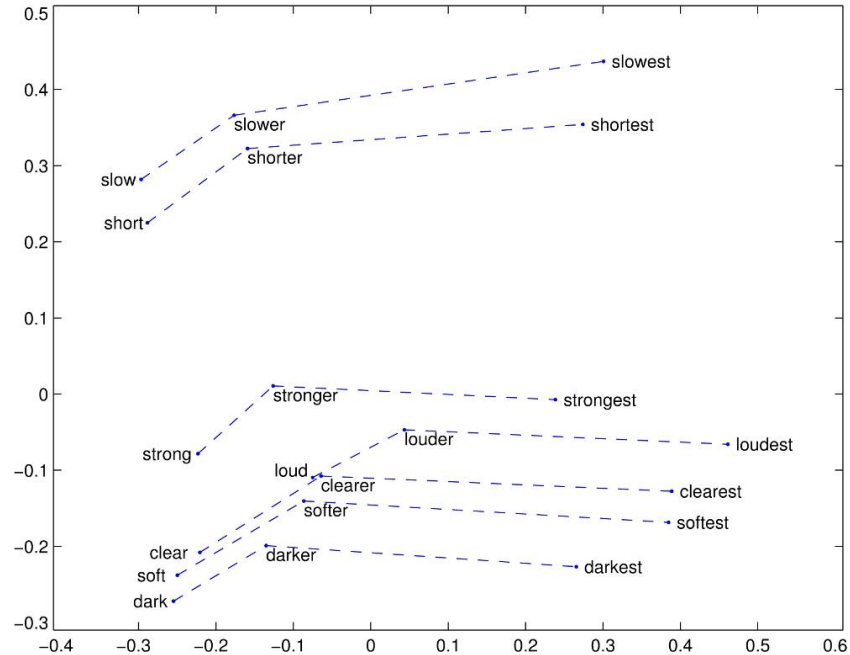
She was at the movies.

...

Learn embeddings that maximize our ability to predict the surrounding words of a word

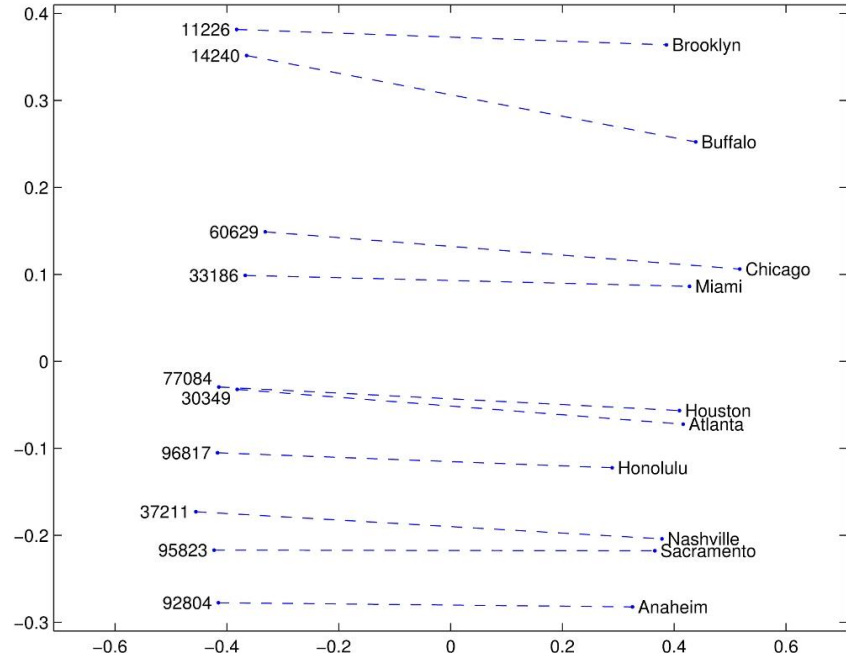
$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{j=-m}^{+m} \log P(w_{t+j} | w_t; \theta)$$

Word Representation: Dense Vectors



Word Representation: Dense Vectors

There exists a 300-dimensional vector z such that if you add it to the vector of a city name, you get the vector of their zip codes!



Word Representation: Dense Vectors

- We have one vector v for each word w .
- d has to encode all aspects and meanings of w
- These two sentences will be almost identical in terms of word embeddings.

How much does a share of **Apple** cost?

How much does a pound of **apple** cost?

- We can do better

Language Modeling

- The task of estimating the probability of a sequence of words

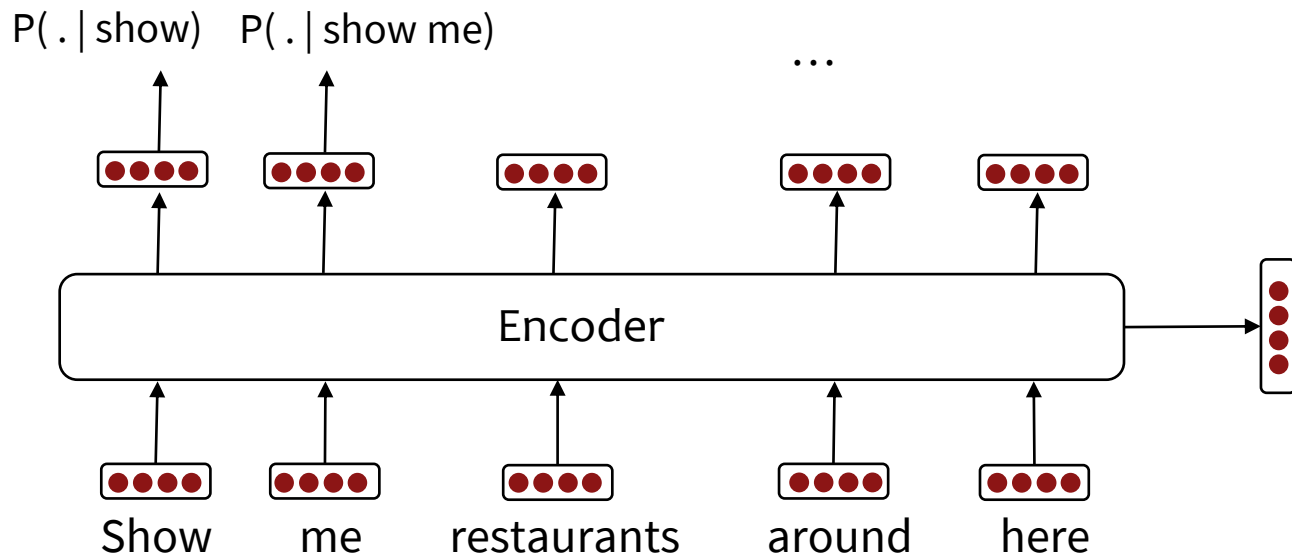
$$P(w_1 w_2 w_3 \dots w_m)$$

- Usually requires simplifying assumptions

$$\begin{aligned} P(w_1 w_2 w_3 \dots w_m) &= \prod_{i=1}^m P(w_i | w_1 \dots w_{i-1}) \\ &\approx \\ &\prod_{i=1}^m P(w_i | w_{i-n} \dots w_{i-1}) \end{aligned}$$

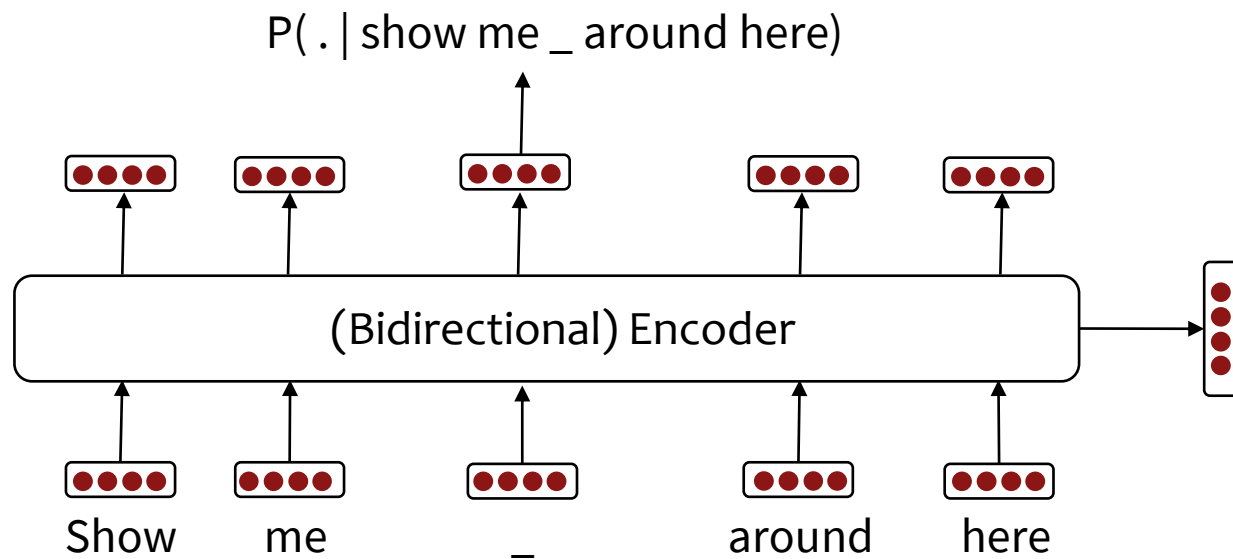
Autoregressive Language Models

- Autoregressive: predict the next word



Masked Language Models

- Masked: fill in the blank



Word Representation: Contextual

- Training data for a task is limited
- Pre-train a language model on a very large text corpus
- Embeddings from Language Models: ELMo (Oct. 2017)
- Generative Pre-training: GPT (June 2018)
- Bidirectional Encoder Representations from Transformers: BERT (Oct. 2018)
- GPT-2 (Feb. 2019)
- T5 (Oct. 2019)
- GPT-3 (May 2020)
- ...

corpus size

800 million words

1x

4x

48x

47x

35x

Quiz

A language model is trained to be good at predicting missing words. How can we test if the contextual representations learned by the language model are good at capturing the meaning of sentences as well?

1. By evaluating them on downstream tasks. BERT for instance improved state of the art results for several NLP tasks by 4-8%.
2. By looking at the representations themselves.

Sequence
to
Sequence



When Both Input and Output Are Sequences of Words

- Seq2Seq has many use cases
 - Machine Translation
 - Question Generation
 - Semantic Parsing
- We will use examples from semantic parsing

Show me restaurants around here



```
now => @QA.Restaurant()  
, geo == current_location => notify
```

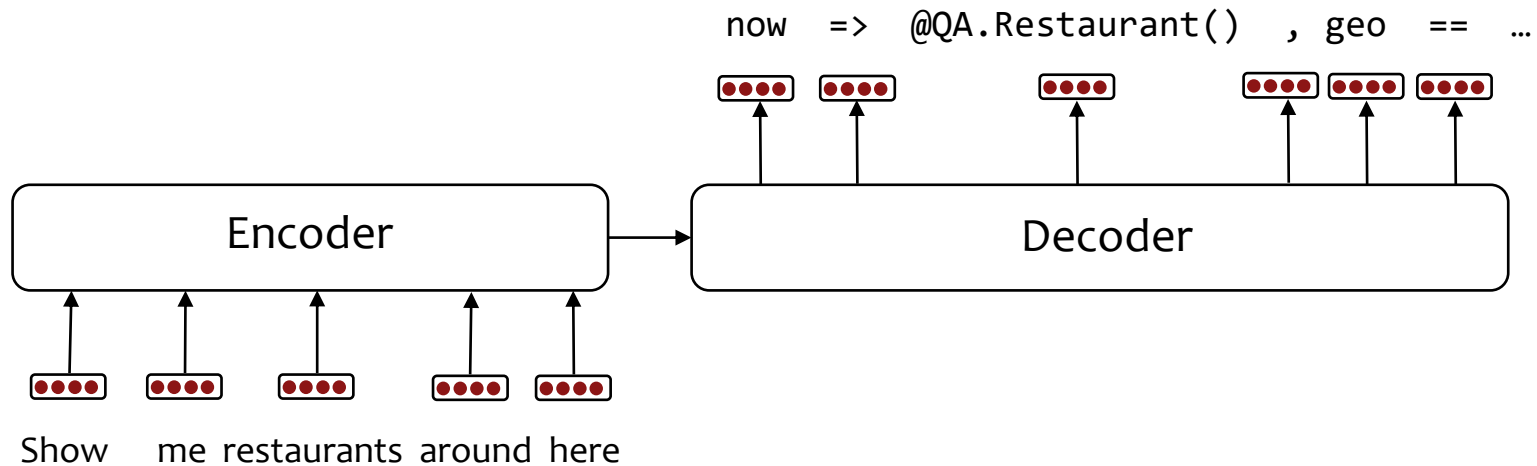
Sequence to Sequence

- Dataset: pairs of source sentence $x_1 x_2 \dots x_s$ and target sentence $y_1 y_2 \dots y_t$
- For instance, pairs of natural sentences and their ThingTalk programs
- The objective is to learn θ that maximizes:

$$\begin{aligned} J(\theta) &= P(y_1 y_2 \dots y_t \mid x_1 x_2 \dots x_s ; \theta) \\ &= \\ &P(y_1 \mid x_1 x_2 \dots x_s ; \theta) \times P(y_2 \mid y_1 x_1 x_2 \dots x_s ; \theta) \times \dots \end{aligned}$$

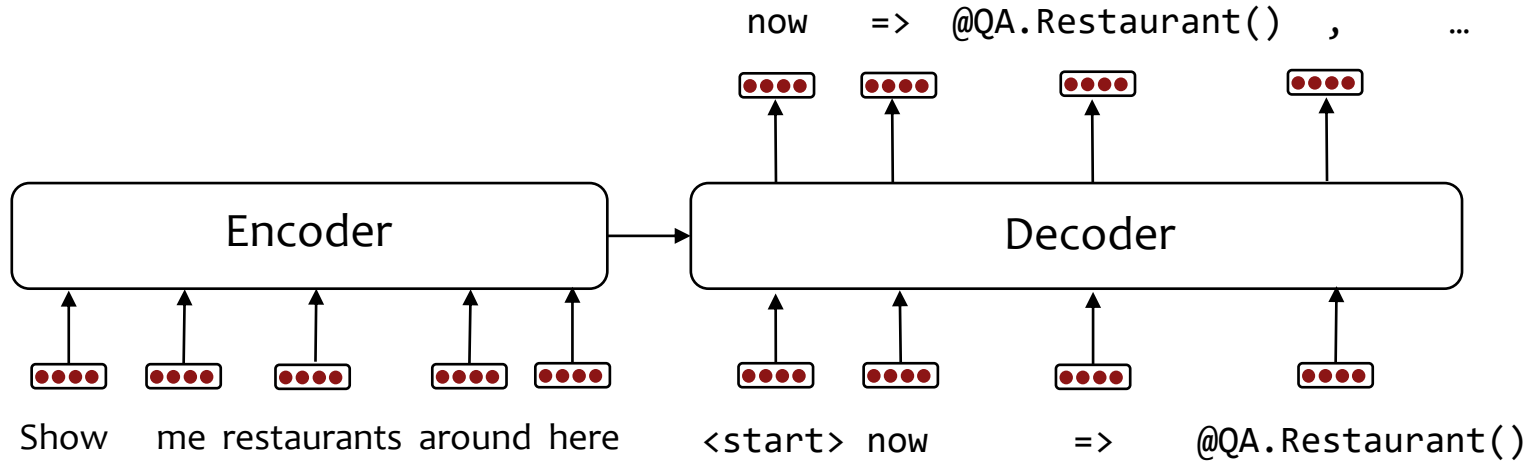
Encoder-Decoder

We can use encoder-decoder models for Seq2Seq tasks



Encoder-Decoder

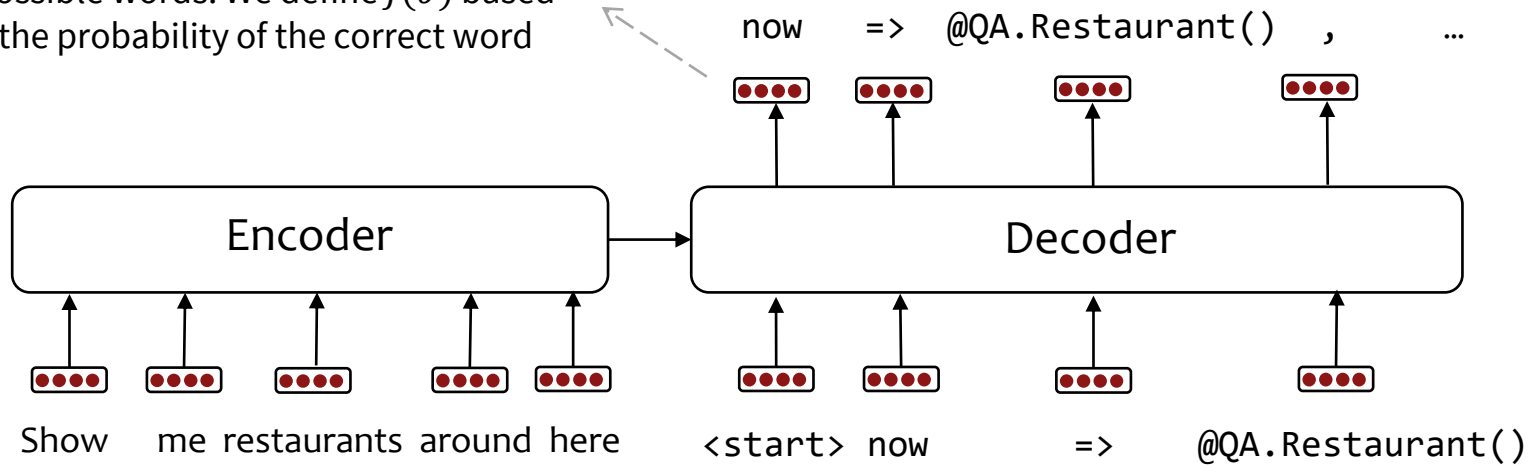
In practice, we also input the previous token to the decoder



Encoder-Decoder

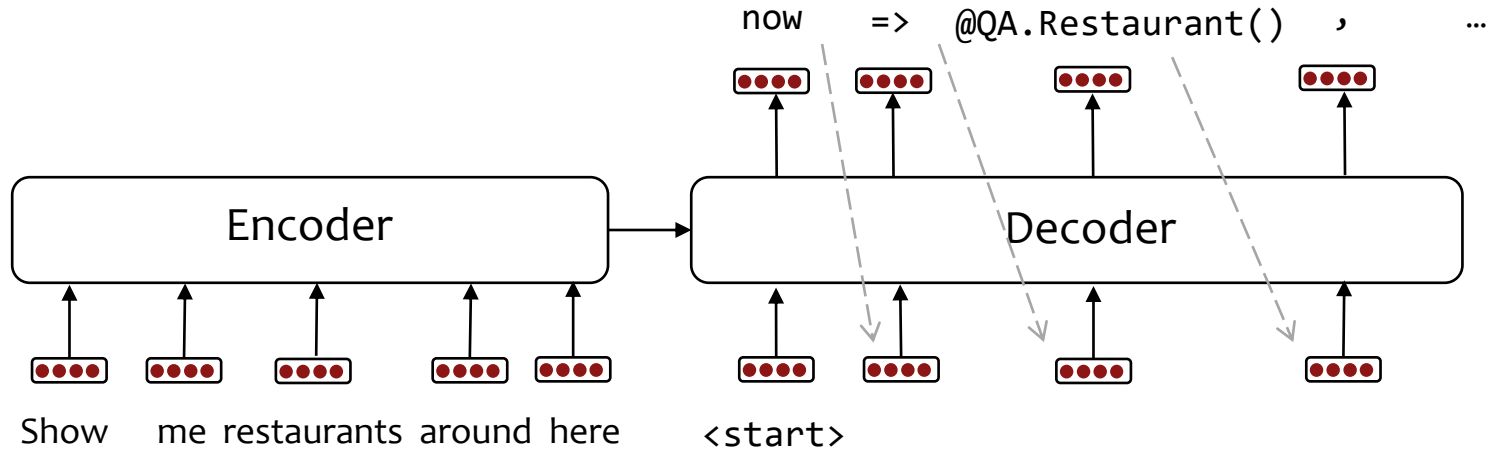
At training time, decoder always gets the **gold target** as input

These vectors define a distribution over all possible words. We define $J(\theta)$ based on the probability of the correct word



Encoder-Decoder

- At generation time, we feed in the word generated by the decoder at previous time step.
- Pro: very fast to converge in practice
- Con: model is never exposed to its own errors during training



From Word Probabilities to Output Sequence

- Greedy decoding: at each step, pick the most probable word
 - Greedy decoding can make search errors: if we choose a wrong word at a step, we might never recover
- Beam Search: at each step, keep the K most probable observed outputs
- Sampling: pick a word at random according to the distribution
- ...

Downside of Word-Level Loss

Source: Show me restaurants around here.

Gold target: now => @QA.Restaurant() , geo == current_location => notify

Model output: now => @QA.Hospital() , geo == current_location => notify

Most of the sentence is the same as the gold, so low cost, but you will –literally– end up in a hospital!

A small difference in words is not the same as a small difference in meaning.

Downside of Word-Level Loss

Source: Show me nearby restaurants.

Gold target: mostrami ristoranti nelle vicinanze

Model output: sto cercando un ristorante qui attorno
(I'm looking for a restaurant around here)

Most of the sentence is different from the gold, so high cost, but the answer is correct.

Difference in words is not the same as difference in meaning.

Quiz

Is this a problem in semantic parsing as well?

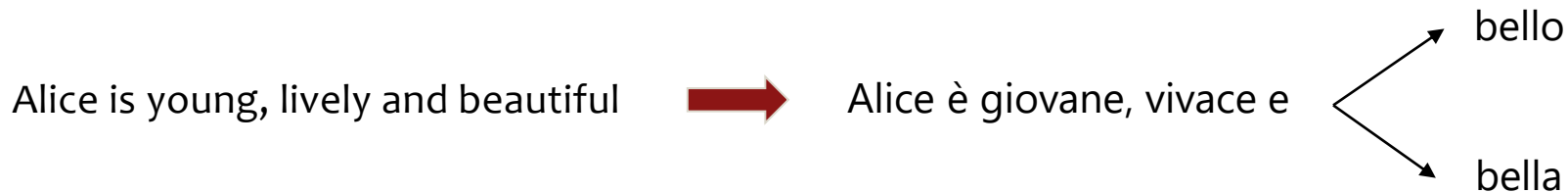
Not for ThingTalk. ThingTalk is normalized, that is, each meaning has exactly one ThingTalk code.

Attention

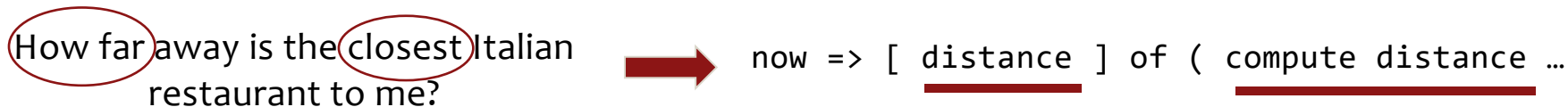


Capturing Long Term Dependencies is Important in NL

- When generating a word, the model has to look at multiple words that are potentially far from each other.



- Some words are more important than others

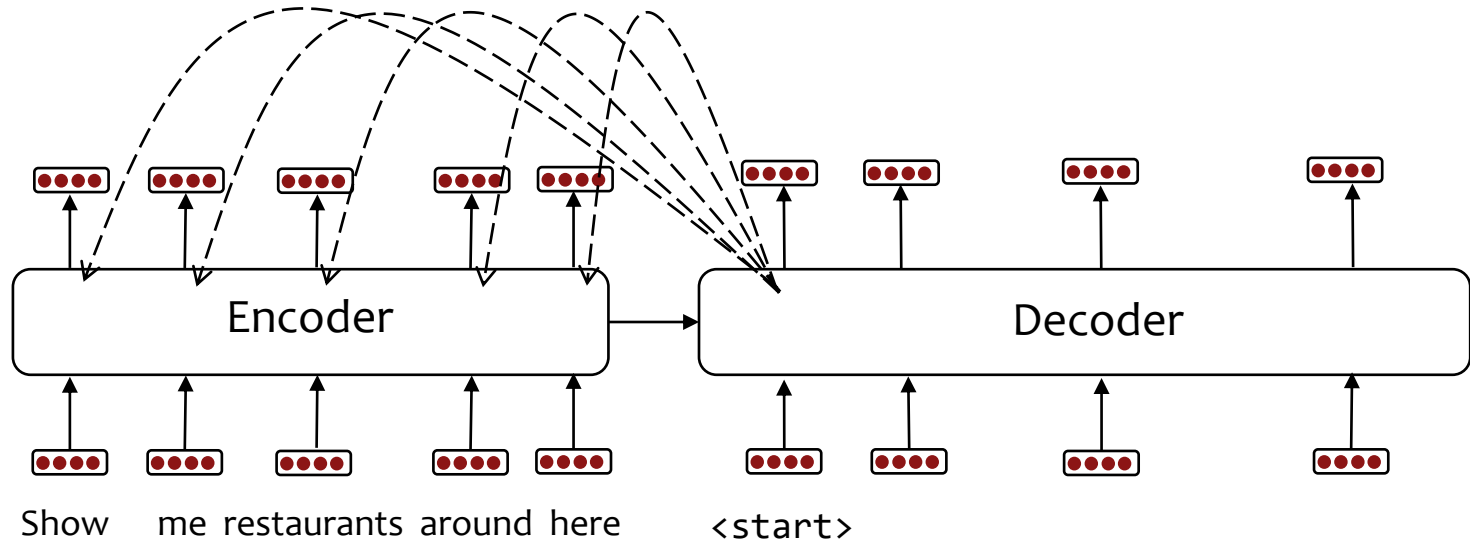


Attention

- Designed to alleviate this exact problem
- At each decoding step, compute *attention scores* by combining encoder and decoder states
- Normalize scores with *softmax*
- Mix them into a *context vector*
- Mix decoder state and context vector

Encoder-Decoder with Attention

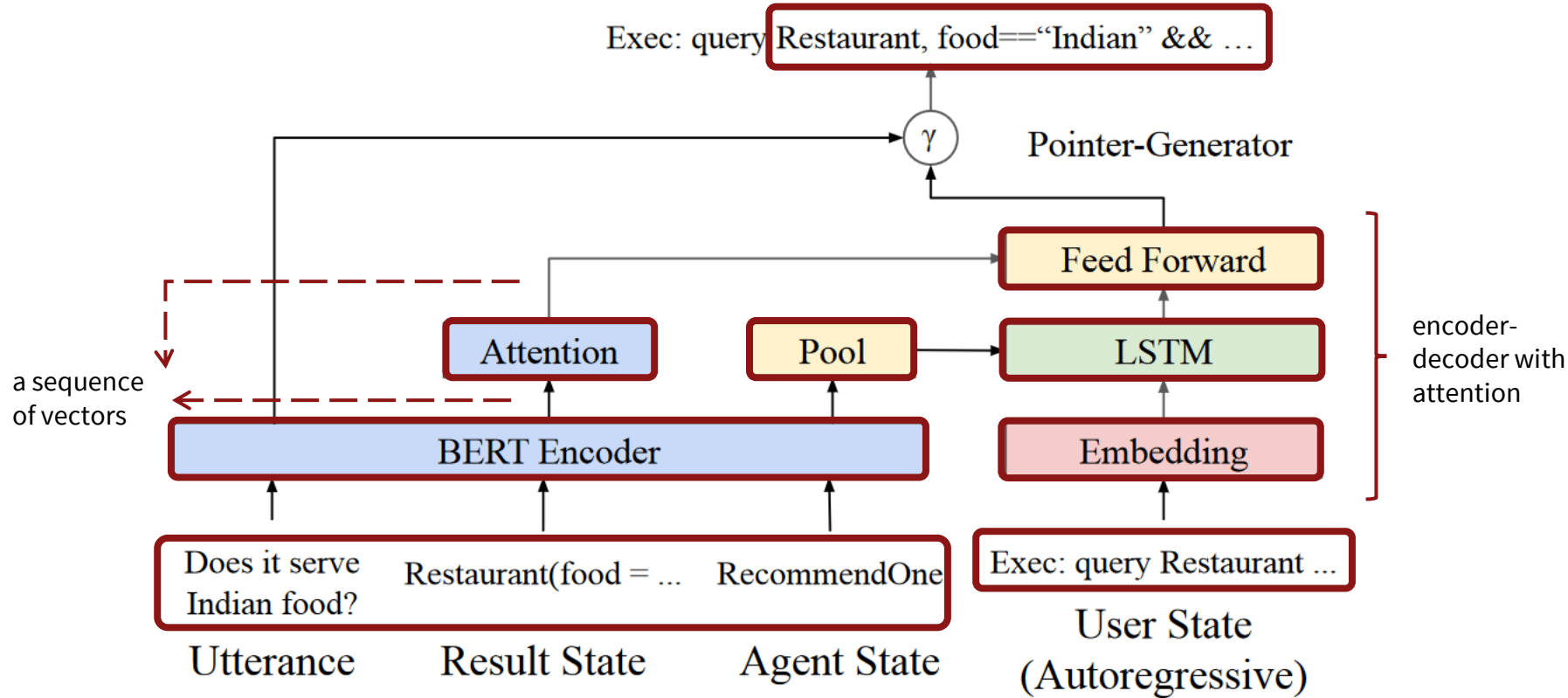
When generating a word for the output, directly look at all the words in the input



Transformer

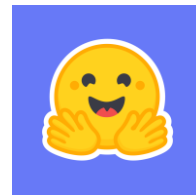
- Is a relatively new class of parametrized functions
- Instead of RNNs, is entirely made up of attentions
- Attentions are easy to compute in parallel, which is especially beneficial when using GPUs
- Empirically, Transformer outperforms RNN in a wide range of tasks and datasets.
- Has encoder, decoder and Seq2Seq variants.

Remember This Image from Lecture 1?



Practical Notes

- Python
- PyTorch
- Genie NLP
- HuggingFace's transformers package includes state-of-the-art pre-trained language models like BERT



GPT-2, 3, 4, ..., N

- Very large transformer models with 175 billion parameters
- Trained on large datasets of Books, Wikipedia and the rest of the Web
- Trained for 3.14×10^{23} FLOPS
- With the objective to predict the next word given the previous words
- They pick up a lot of knowledge about English grammar, the world, and some logic.

GPT-3

- GPT-3 can be “programmed” by showing them a few examples, or a prompt.
- **The following is a conversation with an AI assistant. The assistant is helpful, creative, clever, knowledgeable about myths, legends, jokes, folk tales and storytelling from all cultures, and very friendly.**

Human: Hello, who are you?

AI: I am an AI created by OpenAI. How can I help you today?

Human: I am feeling bored today. Grandma, tell me a story about the time the Cat stole the sun.

AI: Once upon a time, the Cat went to visit the Sun. He hadn't seen the Sun for quite some time. ...

GPT-2, 3, 4, ..., N

There are always caveats:

- The example on the first slide was chosen –by a human- from 10 outputs.
- Writing about is about productivity tips.

Artificial intelligence / Machine learning

A college kid's fake, AI-generated blog fooled tens of thousands. This is how he made it.

"It was super easy actually," he says, "which was the scary part."

Discussion

I talked about how we got here.
But where do we go from here?