# Sequence-to-sequence Models

CIS 530, Computational Linguistics: Spring 2018

John Hewitt & Reno Kriz
University of Pennsylvania

Some concepts drawn a bit transparently from Graham Neubig's excellent
**Neural Machine Translation and Sequence-to-sequence Models: A Tutorial**
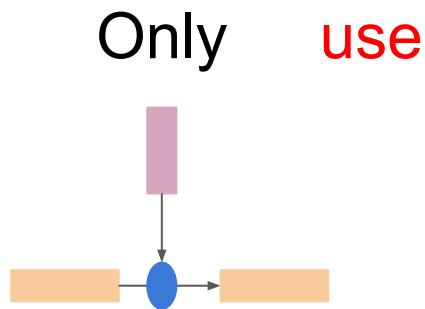**https://arxiv.org/pdf/1703.01619.pdf**

Deep learning!

Now even deeper!

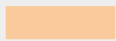# We've already seen RNNs for language modeling

Only

# We've already seen RNNs for language modeling

Only    use

The memory vector, or "state".

The "word vector" representation of the word.

The RNN function, which combines the word vector and the previous state to create a new state.

# We've already seen RNNs for language modeling

Only    use   neural

The memory vector, or "state".

The "word vector" representation of the word.

The RNN function, which combines the word vector and the previous state to create a new state.
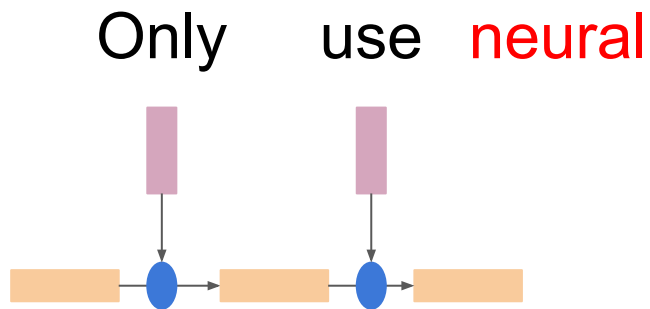
# We've already seen RNNs for language modeling

Only    use   neural   nets



The memory vector, or "state".

The "word vector" representation of the word.

The RNN function, which combines the word vector and the previous state to create a new state.

# We've already seen RNNs for language modeling

Only    use   neural    nets
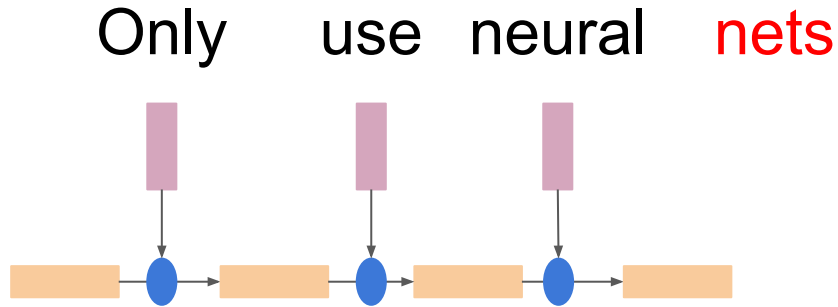
The memory vector, or "state".

The "word vector" representation of the word.

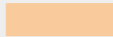The RNN function, which combines the word vector and the previous state to create a new state.

# How does the RNN function work?

The RNN function takes the current RNN state and a word vector and produces a subsequent RNN state that "encodes" the sentence *so far*.

RNN
function

$\bullet$ := $\boxed{1}$ $\boxed{\phantom{}}$ + $\boxed{2}$ $\boxed{\phantom{}}$ + $\boxed{3}$

Learned weights representing how to combine past information (the RNN memory) and current information (the new word vector.)

# How does the prediction function work?

We've seen how RNNs "encode" word sequences. But how do they produce probability distributions over a vocabulary?



A probability distribution over the vocab, constructed from the RNN memory and 1 last transformation (in green.) The softmax function turns "scores" into a probability distribution.

# Want to predict things other than the next word?

The model architecture (read: "design") we've seen so far is frequently used in tasks other than language modeling, because modeling sequential information is useful in language, apparently.

Only　　use　neural　　nets

Here's our RNN encoder, representing the sentence.

# Want to predict things other than the next word?

The model architecture (read: "design") we've seen so far is frequently used in tasks other than language modeling, because modeling sequential information is useful in language, apparently.

Only    use   neural    nets

Predict parts of speech!

ADV    VB    ADJ    NNS

# Want to predict things other than the next word?

The model architecture (read: "design") we've seen so far is frequently used in tasks other than language modeling, because modeling sequential information is useful in language, apparently.

Only    use    neural    nets

Or syntax!

ADV    VB    ADJ    NNS

# General idea: build a representation

The method of building the representation is called an Encoder and is frequently an RNN.

Only     use   neural    nets



Each memory vector in the encoder attempts to represent the sentence so far, but mostly represents the word most recently input.

# General idea: generate the output one token at a time

The model that takes the encoded representation and generates the output is called the Decoder, and, errrr, is also generally an RNN.

Only    use   neural    nets

Jiri

# General idea: generate the output one token at a time

The model that takes the encoded representation and generates the output is called the Decoder, and, errrr, is also generally an RNN.

Only     use   neural    nets

Jiri   naaṇị

# General idea: generate the output one token at a time

The model that takes the encoded representation and generates the output is called the Decoder, and, errrr, is also generally an RNN.

Only     use   neural    nets

Jiri   naani   netwọk

# General idea: generate the output one token at a time

The model that takes the encoded representation and generates the output is called the Decoder, and, errrr, is also generally an RNN.

# General idea: generate the output one token at a time

The model that takes the encoded representation and generates the output is called the Decoder, and, errrr, is also generally an RNN.

# How is it trained?

In practice, training for a single sentence is done by "forcing" the decoder to generate gold sequences, and penalizing it for assigning the sequence a low probability. Losses for each token in the sequence are summed. Then, the summed loss is used to take a step in the right direction in all model parameters (including word embeddings!) (*stochastic gradient descent.*)

neural    nets

# How is it trained?

In practice, training for a single sentence is done by "forcing" the decoder to generate gold sequences, and penalizing it for assigning the sequence a low probability. Losses for each token in the sequence are summed. Then, the summed loss is used to take a step in the right direction in all model parameters (including word embeddings!) (*stochastic gradient descent.*)

neural    nets

.7

Jiri    naani    netwok    nu

# How is it formalized? How is it trained?

In practice, training for a single sentence is done by "forcing" the decoder to generate gold sequences, and penalizing it for assigning the sequence a low probability. Losses for each token in the sequence are summed. Then, the summed loss is used to take a step in the right direction in all model parameters (including word embeddings!) (*stochastic gradient descent.*)

neural    nets

.7

GOLD:  Jiri

Jiri    naani    netwok    nu

**Loss**
**-log(.7)**

# Sentence-level training

Almost all such networks are trained using **cross-entropy loss**. **At each step, the network produces a probability distribution over possible next tokens. This distribution is penalized from being different from the true distribution (e.g., a probability of 1 on the actual next token.)**

# Sentence-level training

Almost all such networks are trained using **cross-entropy loss**. **At each step, the network produces a probability distribution over possible next tokens. This distribution is penalized from being different from the true distribution (e.g., a probability of 1 on the actual next token.)**

# Sentence-level training

Almost all such networks are trained using **cross-entropy loss. At each step, the network produces a probability distribution over possible next tokens. This distribution is penalized from being different from the true distribution (e.g., a probability of 1 on the actual next token.)**



Jiri   naanị   netwọk   nụ

-log(.7)   -log(.5)   -log(.6)   -log(.4)

# Sentence-level training

Almost all such networks are trained using **cross-entropy loss**. **At each step, the network produces a probability distribution over possible next tokens. This distribution is penalized from being different from the true distribution (e.g., a probability of 1 on the actual next token.)**
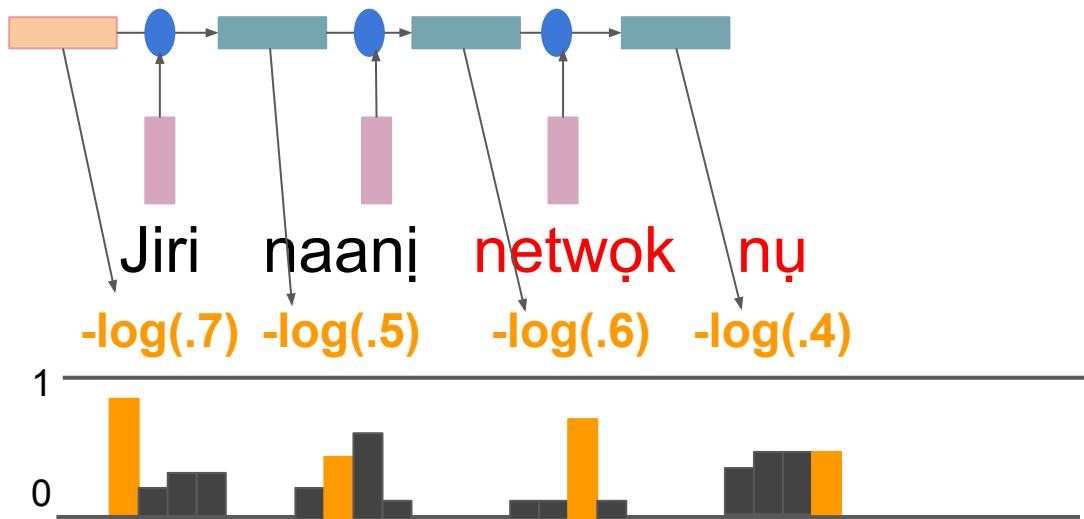


Jiri    naanị    netwọk    nụ

sum(  **-log(.7)  -log(.5)    -log(.6)    -log(.4)** ) = 1.07

Minimize This!

# How is it formalized?

Let $h_t$ be the RNN hidden state at timestep t:

# How is it formalized?

Let $h_t$ be the RNN hidden state at timestep t:

Let $x_t$ be the input vector at timestep t:

# How is it formalized?

Let $h_t$ be the RNN hidden state at timestep t:

Let $x_t$ be the input vector at timestep t:

The RNN equation posits 2 matrices and 1 vector as parameters:

$W^{hx}$ integrates input vector information.

# How is it formalized?

Let $h_t$ be the RNN hidden state at timestep t:

Let $x_t$ be the input vector at timestep t:

The RNN equation posits 2 matrices and 1 vector as parameters:

$W^{hx}$ integrates input vector information.

$W^{hh}$ integrates information from the previous timestep.

# How is it formalized?

Let $h_t$ be the RNN hidden state at timestep t: 

Let $x_t$ be the input vector at timestep t: 

The RNN equation posits 2 matrices and 1 vector as parameters:

$\mathbf{W^{hx}}$ integrates input vector information. 

$\mathbf{W^{hh}}$ integrates information from the previous timestep. 

$\mathbf{b^h}$ is a bias term. (What function does this perform?)

# How is it formalized?

Let $h_t$ be the RNN hidden state at timestep t:

Let $x_t$ be the input vector at timestep t:

The RNN equation posits 2 matrices and 1 vector as parameters:

$\mathbf{W^{hx}}$ integrates input vector information.

$\mathbf{W^{hh}}$ integrates information from the previous timestep.

$\mathbf{b^h}$ is a bias term. (What function does this perform?)

The RNN equation is: $\mathbf{h_t = tanh(W^{hx}x_t + W^{hh}h_{t-1} + b^h)}$

# How is it formalized?

For prediction, we take the current hidden state, and use it as *features* in what is more or less a linear regression.

# How is it formalized?

For prediction, we take the current hidden state, and use it as *features* in what is more or less a linear regression.

Let $d_t$ be our decision (e.g., word, POS tag) at timestep t. Let D be the set of all possible decisions. Let $\mathbf{s_{t-1}}$ be the most recent *decoder* hidden state.

# How is it formalized?

For prediction, we take the current hidden state, and use it as *features* in what is more or less a linear regression.

Let $d_t$ be our decision (e.g., word, POS tag) at timestep t. Let D be the set of all possible decisions. Let $\mathbf{s_{t-1}}$ be the most recent *decoder* hidden state.

$$d_t = argmax_{d' \in D} \ p(\ d' \mid x_{1:n}, d_{1:t-1})$$

# How is it formalized?

For prediction, we take the current hidden state, and use it as *features* in what is more or less a linear regression.

Let $d_t$ be our decision (e.g., word, POS tag) at timestep t. Let D be the set of all possible decisions. Let $\mathbf{s_{t-1}}$ be the most recent *decoder* hidden state.

$$d_t = argmax_{d' \in D} \ p(\ d' \mid x_{1:n}, d_{1:t-1})$$

$$p(\ * \mid x_{1:n}, d_{1:t-1}) = softmax_D(W^{Dh}s_{t-1} + b^D)$$

# How is it formalized?

For prediction, we take the current hidden state, and use it as *features* in what is more or less a linear regression.

Let $d_t$ be our decision (e.g., word, POS tag) at timestep t. Let D be the set of all possible decisions. Let $\mathbf{s_{t-1}}$ be the most recent *decoder* hidden state.

$$d_t = argmax_{d' \in D} \ p( \ d' \mid x_{1:n}, d_{1:t-1})$$

$$p( \ * \mid x_{1:n}, d_{1:t-1}) = softmax_D(W^{Dh}s_{t-1} + b^D)$$

Note that $\mathbf{W^{hD}s_{t-1} + b^D}$ produces a *vector of scores.* The **softmax** function normalizes scores to a probability distribution by exponentiating each dimension, and normalizing by the sum. For some choice k of K, $\mathbf{p(k) = e^{score(k)}/ \sum_{k' \in K} e^{score(k')}}$

# The information bottleneck and latent structure

Given the diagram below, what problem do you foresee when translating progressively longer sentences?

# The information bottleneck and latent structure

We are trying to encode *variable-length* structure (e.g., variable-length sentences) in a fixed-length memory (e.g., only the 300 dimensions of your hidden state.)



**Encoder (seq)**

The last encoder hidden state is the bottleneck -- all information in the source sentence must pass through it to get to the decoder.

Finding a solution to this problem was the final advance that made neural MT competitive with previous approaches.

# The information bottleneck and latent structure

The key insight is related to the word alignment work we did last week. We allow the decoder to look at *any* encoder state, and let it learn which are important at each time step!

# Learning to pay attention

Attention **summarizes the encoder**, focusing on specific parts/words.

Step 1: Take the decoder state, and compute an *affinity* $\alpha_i$ with all encoder states.

# Learning to pay attention

Step 1:  Take the decoder state, and compute an *affinity* $\alpha_i$ with all encoder states.



$\alpha_0$  $\alpha_1$  $\alpha_2$  $\alpha_3$  $\alpha_4$

Jiri

The affinity function, ●, is a dot product, or something similar.

● : ▬ ▮ = $\alpha_i$

# Learning to pay attention

Attention **summarizes the encoder**, focusing on specific parts/words.

Step 2:  Normalize the scores to sum to 1 by the softmax function.
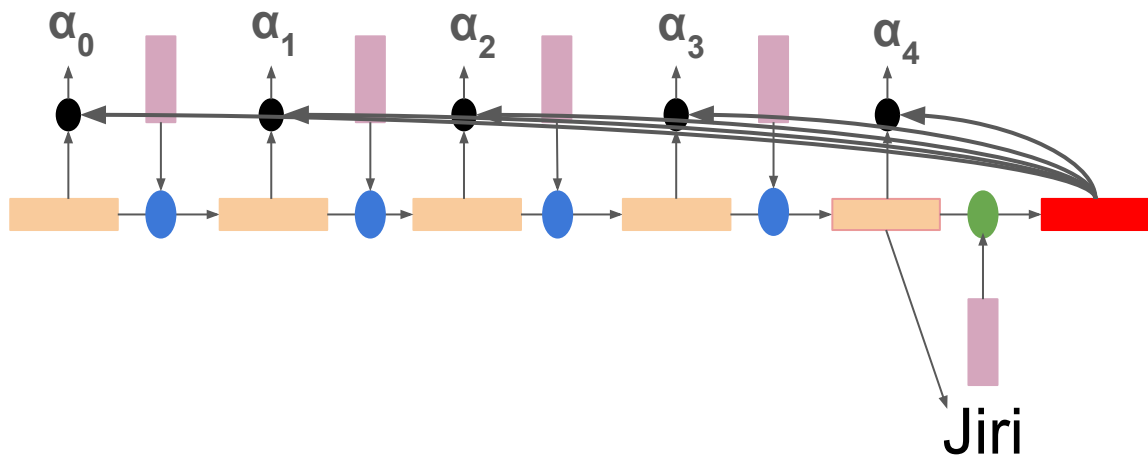
# Learning to pay attention

Attention **summarizes the encoder**, focusing on specific parts/words.

Step 2:  Normalize the scores to sum to 1 by the softmax function.



$a_0$  $a_1$  $a_2$  $a_3$  $a_4$

softmax

Note that $\sum_{i=1,2,3,4} a_i = 1$

$\alpha_0$  $\alpha_1$  $\alpha_2$  $\alpha_3$  $\alpha_4$

Jiri

# Learning to pay attention

Step 3: Average the encoder states, weighted by the **a** distribution.

$$a_0 \boxed{1} + a_1 \boxed{1} + a_2 \boxed{1} + a_3 \boxed{1} + a_4 \boxed{1} = \blacksquare$$



This weighted average

Is called the **context vector.**

Jiri

# Learning to pay attention

Attention **summarizes the encoder**, focusing on specific parts/words.

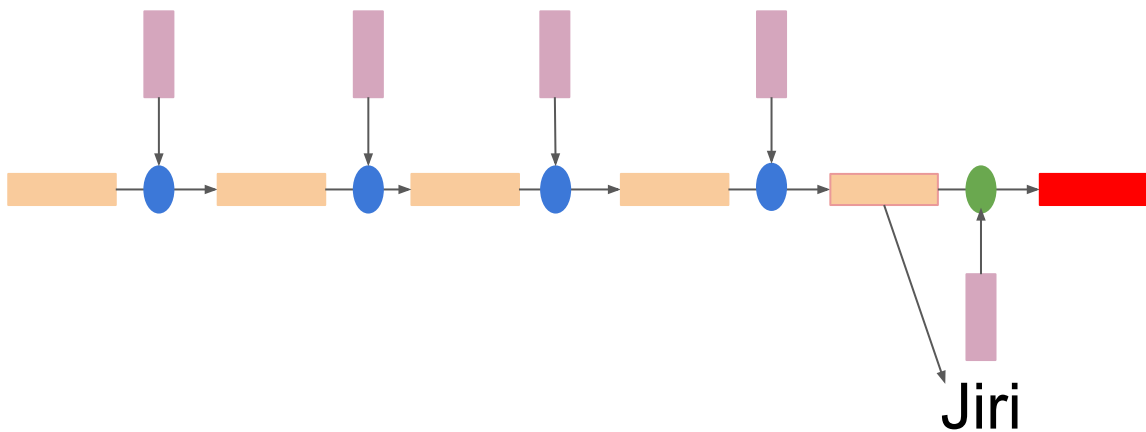Step 3: Average the encoder states, weighted by the **a** distribution.

Only    use   neural    nets

In this example, since "Jiri" means "use", the attention will focus on the vectors around "use".

Jiri

Focus of context vector over encoder states

# Learning to pay attention

Step 4:  Use the context vector at prediction, concatenating it to the decoder state.



softmax( ) =

4

Probability distribution over the vocabulary

Jiri

This vector has the current decoder information, ▮, but also a focused summary of the encoder, ▮.

# Attention Formalization

Attention computes the **affinity between the decoder state and all encoder states**. There are many affinity computation methods, but they're all like a **dot product**.

Let there are $n$ encoder states. The affinity between encoder state $i$ and the decoder state is $\alpha_i$. The encoder states are $h_{1:n}$, and the decoder state is $s_{t-1}$.

$\qquad$ Let $\alpha_i = f(h_i, s_{t-1}) = h_i^\top s_{t-1}$

$\qquad$ Let weights $a = \text{softmax}(\alpha)$.

$\qquad$ Let the context $c = \sum_{i=1:n} h_i a_i$. *(Note that this is a weighted average.)*

# Attention Formalization

Attention is used at prediction as extra information in the final prediction.

Reminder, we let the context $\mathbf{c} = \sum_{i=1:n} \mathbf{h_i a_i}$. *(Weighted average of encoder states.)*

# Attention Formalization

Attention is used at prediction as extra information in the final prediction.

Reminder, we let the context $c = \sum_{i=1:n} h_i a_i$. *(Weighted average of encoder states.)*

Let the notation $[s;c]$ mean the concatenation of vectors $s$ and $c$.

$$d_t = argmax_{d' \in D} \; p(\, d' \mid x_{1:n}, d_{1:t-1})$$

(same as before, without attention)

# Attention Formalization

Attention is used at prediction as extra information in the final prediction.

Reminder, we let the context $c = \sum_{i=1:n} h_i a_i$. *(Weighted average of encoder states.)*

Let the notation $[s;c]$ mean the concatenation of vectors $s$ and $c$.

$$d_t = argmax_{d' \in D} \; p( \, d' \mid x_{1:n}, d_{1:t-1})$$

$$p( \, * \mid x_{1:n}, d_{1:t-1}) = softmax_D(W^{D(2h)}[s_{t-1};c]+b^D)$$

So, the only difference is that the final prediction uses the context vector concatenated to the decoder state to make the prediction.

# Empirical considerations

There are a lot of "hyperparameter" choices that can greatly affect the quality of your model. **In short, take parameters from papers/tutorials, and grid search (try many combinations of parameters) around them.**

**RNN variants: LSTMs have a different (much better) recurrent equation.**

**Hidden state sizes**: larger: more memory! Requires more data.

**Embedding sizes**: more representation power! Requires more data.

**Learning rate**: the step size you take in learning your parameters! Start this "large", and cut it in half when your training stops improving development set performance.

# Empirical considerations

There are a lot of "hyperparameter" choices that can greatly affect the quality of your model. **In short, take parameters from papers/tutorials, and grid search (try many combinations of parameters) around them.**

**Regularization: "dropout" prevents overfitting by making each node in your hidden state unavailable for an observation with a given probability. Try some values around .2 to .3.**

**Batch size:** The number of observations to group together before performing a parameter update step. Larger batches: less fine-grained training, many more observations per minute, especially on GPU.

# Case study: text simplification

Text simplification is the process in which a text is transformed into an equivalent text that can be more easily read by a broader audience (Saggion, 2017).

Simplification can be used as a preprocessing tool for improving performance of many NLP end-tasks such as parsing, SRL, summarization, Information Retrieval etc.

# Case study: text simplification

Text simplification is the process in which a text is transformed into an equivalent text that can be more easily read by a broader audience (Saggion, 2017).

Simplification can be used as a preprocessing tool for improving performance of many NLP end-tasks such as parsing, SRL, summarization, Information Retrieval etc.

"There's just one major hitch: the primary purpose of education is to develop citizens with a wide variety of skills."

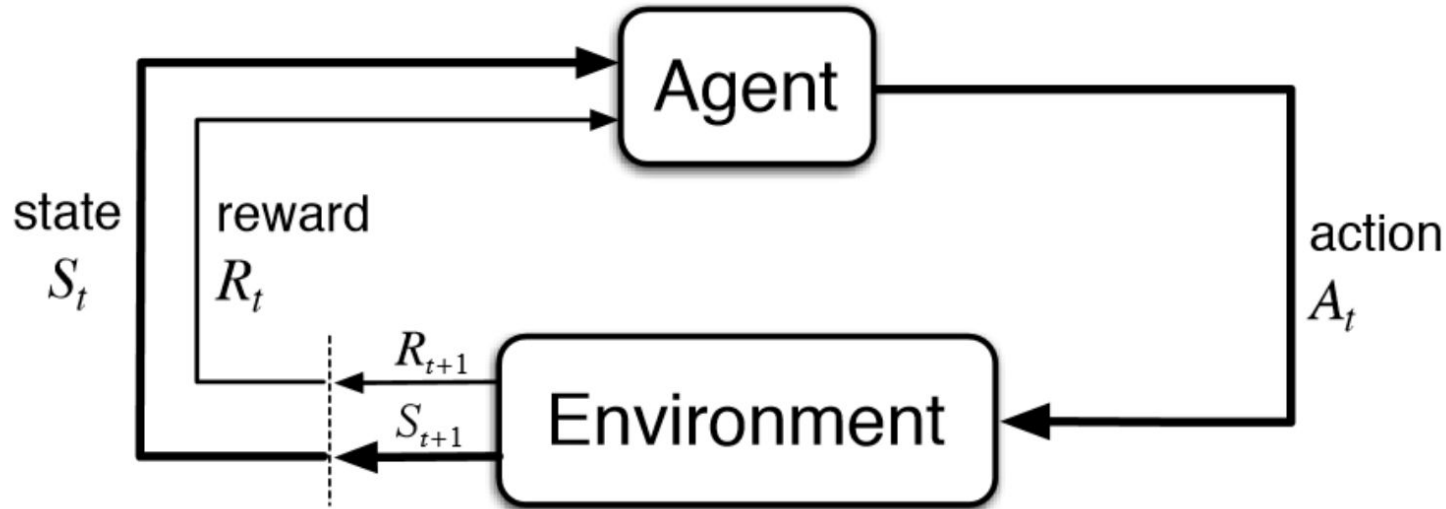"The purpose of education is to develop many skills."

# Case study: text simplification

Text simplification can be thought of in part as monolingual machine translation.

Problem: The most common rewrite operation is copying from the complex sentence to the simple sentence.
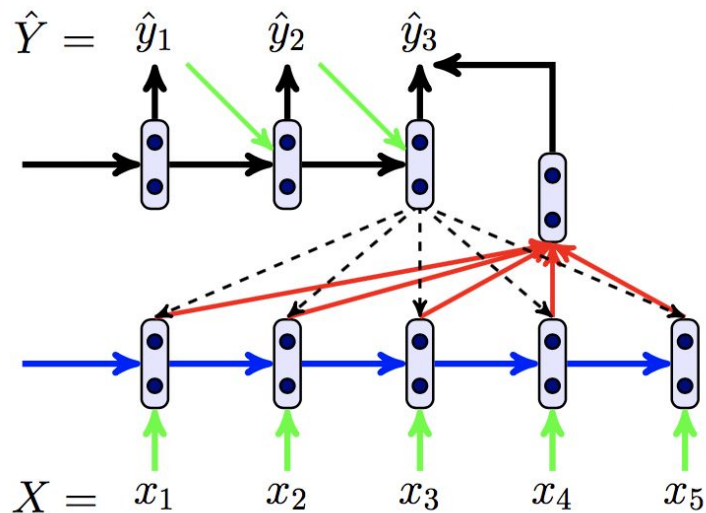- One solution: Add in reinforcement learning (Zhang and Lapata 2017), to encourage the model to use other rewrite operations, such as deletion, substitution, word reordering.

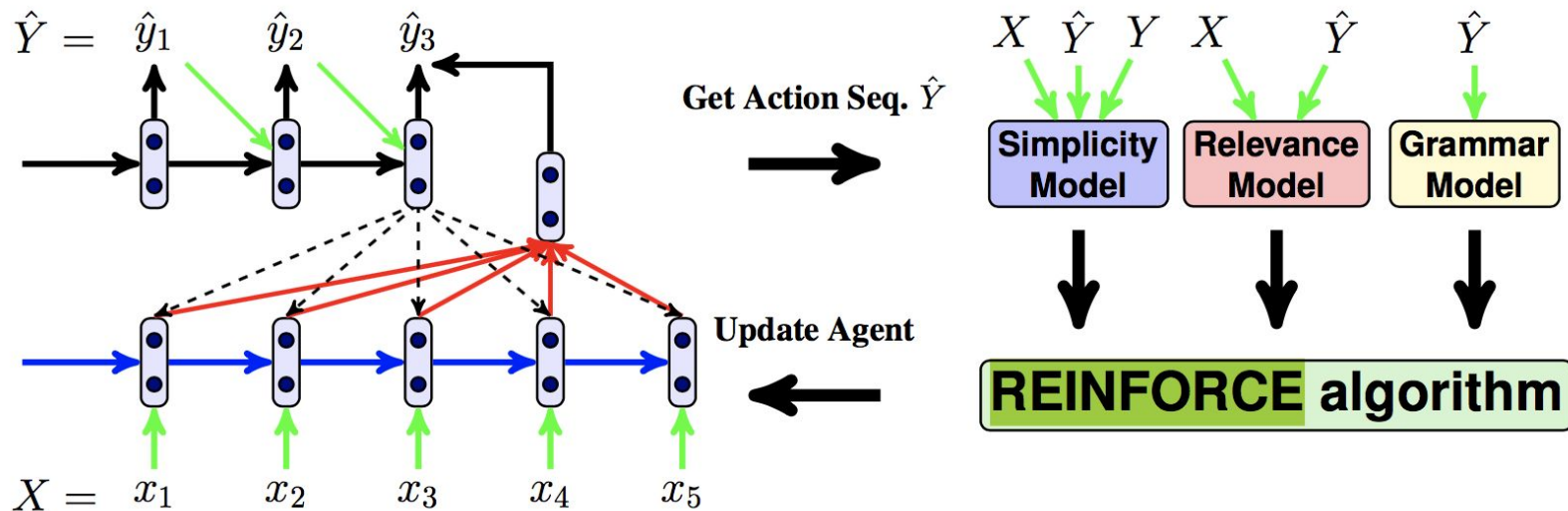# A brief introduction to Reinforcement Learning



The reinforcement learning framework (Sutton and Barto, 1998)

# Case study: text simplification



Basic encoder-decoder model, from (Zhang and Lapata, 2017).

# Case study: text simplification



Encoder-Decoder model with reinforcement learning (Zhang and Lapata, 2017).
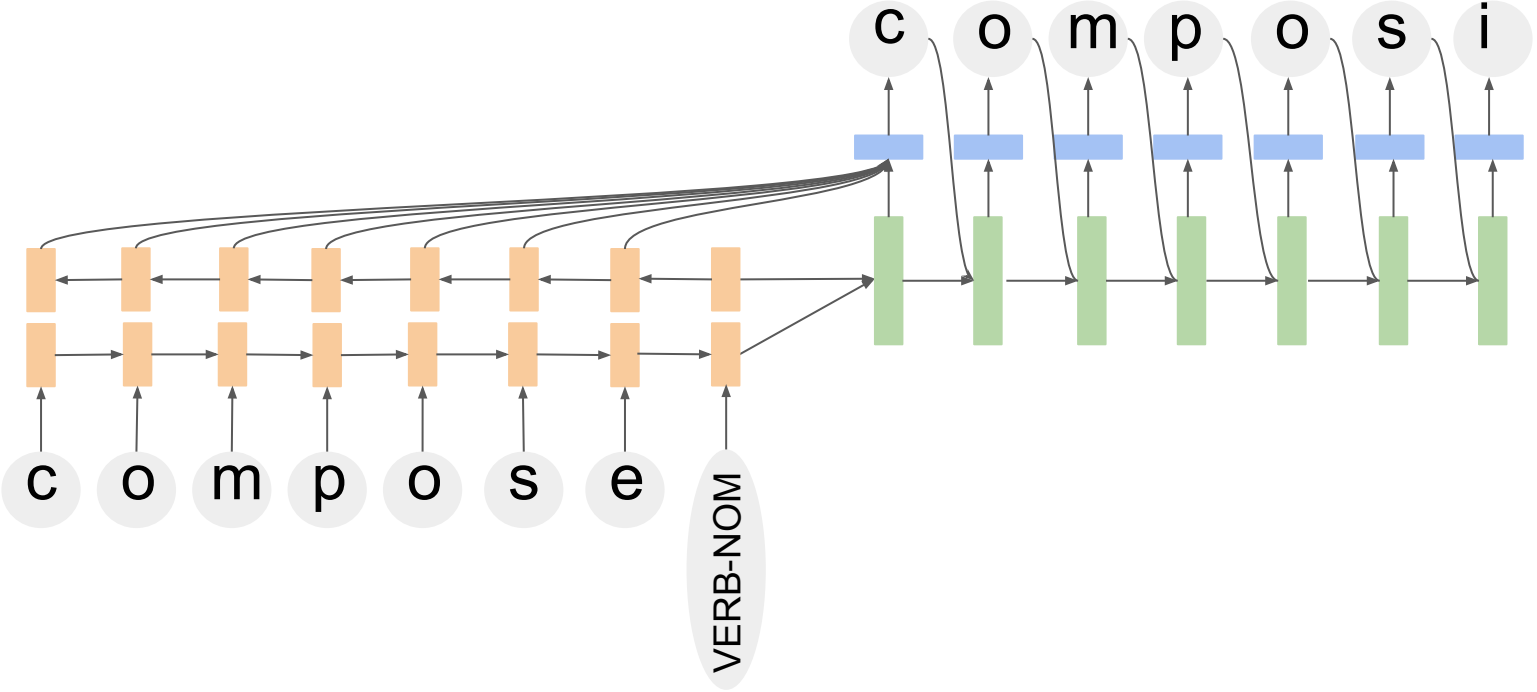
# Derivational morphology

- Process of generating new words from existing words
  - Changes semantic meaning
  - Often a new part-of-speech

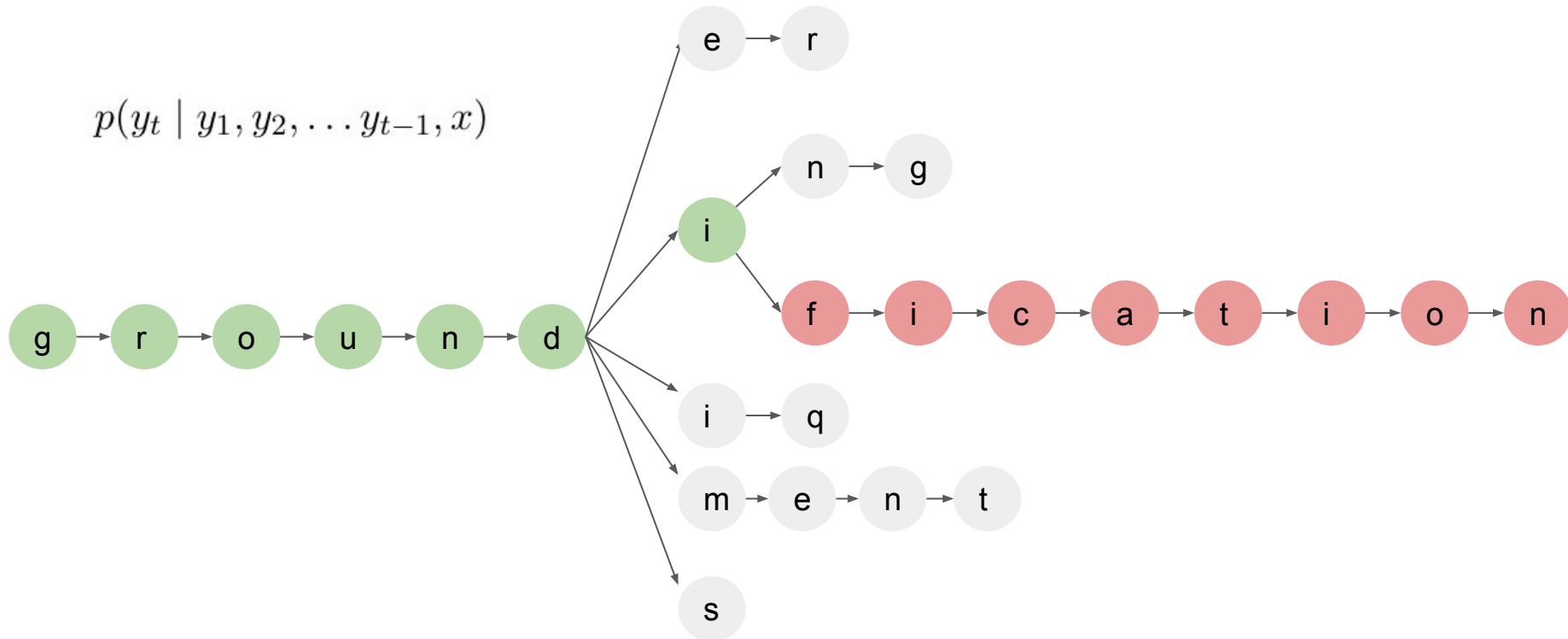| employ | V -> N, Agent | employer |
| employ | V -> N, Passive | employee |
| employ | V -> N, Result | employment |
| employ | V -> Adj, Potential | employable |
| employable | V -> Adj -> N, Stative | employability |

**Encoder (seq)**　　　　**Decoder (2seq)**

# Derivational morphology

# Derivational morphology: search



$$p(y_t \mid y_1, y_2, \ldots y_{t-1}, x)$$

# *Reference Sheet*

The memory vector, or "state" . Color denotes whether encoder or decoder.

The "word vector" representation of the word.

The RNN function, which combines the word vector and the previous state to create a new state.

A learned parameter matrix

$W^{hx}$ integrates input vector information.

$W^{hh}$ integrates information from the previous timestep.

$b^h$ is a bias term.

$d_t$ is our decision at timestep t.

The RNN equation is:

$$h_t = \tanh(W^{hx}x_t + W^{hh}h_{t-1} + b^h)$$

$$d_t = \underset{d' \in D}{argmax}\ p(\ d'\ |\ x_{1:n}, d_{1:t-1})$$

$$p(\ *\ |\ x_{1:n}, d_{1:t-1}) =$$

$$\text{softmax}_D(W^{Dh}s_{t-1} + b^D)$$