# Sequence to Sequence Models

Jon Dehdari

January 25, 2016

# Good Morning!

# Sentence Vectors

- We've seen that words can be represented as vectors. Can sentences be represented as vectors?
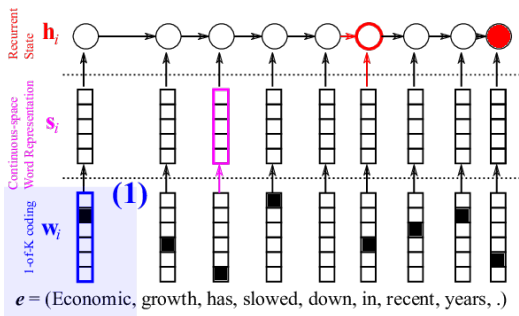
# Sentence Vectors

- We've seen that words can be represented as vectors. Can sentences be represented as vectors?
- Sure, why not?

# Sentence Vectors

- We've seen that words can be represented as vectors. Can sentences be represented as vectors?
- Sure, why not? How? From the hidden state at the end of a sentence: $\mathbf{h}_i = \phi_{\mathsf{enc}}(\mathbf{h}_{i-1}, \mathbf{s}_i)$ ($\phi_{\mathsf{enc}} = $ LSTM or GRU)
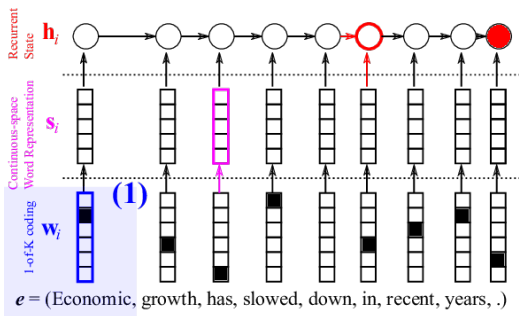
# Sentence Vectors

- We've seen that words can be represented as vectors. Can sentences be represented as vectors?
- Sure, why not? How? From the hidden state at the end of a sentence: $\mathbf{h}_i = \phi_{\text{enc}}(\mathbf{h}_{i-1}, \mathbf{s}_i)$ ($\phi_{\text{enc}}$ = LSTM or GRU)



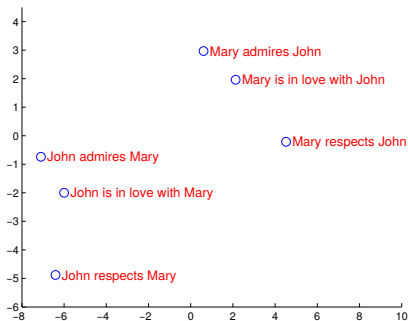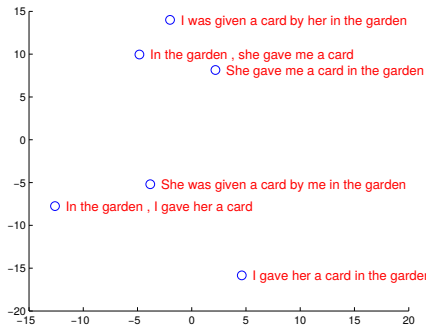$e$ = (Economic, growth, has, slowed, down, in, recent, years, .)

# Sentence Vectors

- We've seen that words can be represented as vectors. Can sentences be represented as vectors?
- Sure, why not? How? From the hidden state at the end of a sentence: $\mathbf{h}_i = \phi_{\text{enc}}(\mathbf{h}_{i-1}, \mathbf{s}_i)$  ($\phi_{\text{enc}}$ = LSTM or GRU)



$e$ = (Economic, growth, has, slowed, down, in, recent, years, .)

- Are they any good? For Elman networks (SRNs), not so much. For LSTMs or GRUs, yes, they're pretty good
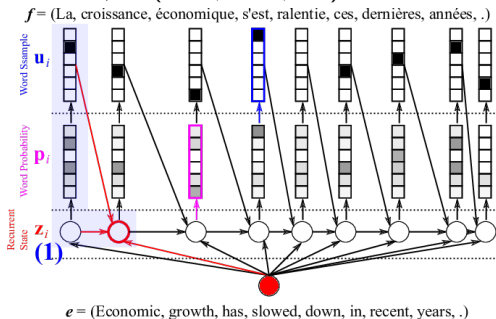
# Sentence Vector Examples



Sentence vectors were projected to two dimensions using PCA
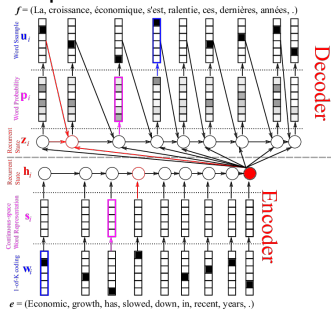
# Generating Sentences from Vectors

- We can also try to go the other direction, generating sentences from vectors
- How? Use an RNN to **decode**, rather than **encode** a sentence:

$$\mathbf{z}_i = \phi_{\text{dec}}(\mathbf{z}_{i-1}, \mathbf{u}_{i-1}, \mathbf{h}_T)$$



$f$ = (La, croissance, économique, s'est, ralentie, ces, dernières, années, .)

$e$ = (Economic, growth, has, slowed, down, in, recent, years, .)

- $\mathbf{h}_T$ ensures global sentence coherency (& adequacy in MT); $\mathbf{u}_{i-1}$ ensures local fluency
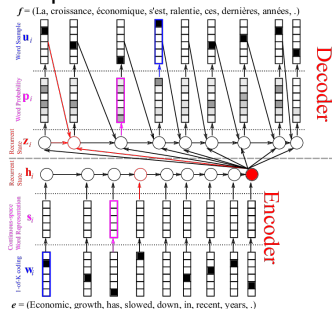
# Using Neural Encoders & Decoders to Translate

- We can combine the neural encoder and decoder of previous slides to form an **encoder-decoder model**
- This can be used for machine translation, and other tasks that map sequences to sequences

# Using Neural Encoders & Decoders to Translate

- We can combine the neural encoder and decoder of previous slides to form an **encoder-decoder model**
- This can be used for machine translation, and other tasks that map sequences to sequences
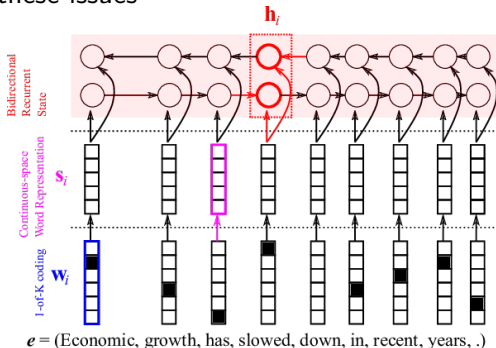


- Monolingual word projections (vectors/embeddings) are trained to maximize likelihood of next word
- Source-side word projections ($s_i$) in an encoder-decoder setting are trained to maximize target-side likelihood

# Bidirectional RNNs

- The basic encoder-decoder architecture doesn't handle long sentences very well
- Everything must fit into a fixed-size vector,
- and RNNs remember recent items better

# Bidirectional RNNs

- The basic encoder-decoder architecture doesn't handle long sentences very well
- Everything must fit into a fixed-size vector,
- and RNNs remember recent items better
- We can combine left-to-right and right-to-left RNNs to overcome these issues
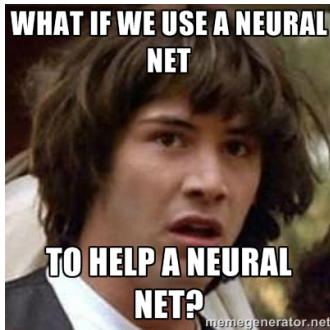


$e$ = (Economic, growth, has, slowed, down, in, recent, years, .)

# What if . . .

- Even bidirectional encoder-decoders have a hard time with long sentences
- We need a way to keep track of what's already been translated and what to translate next
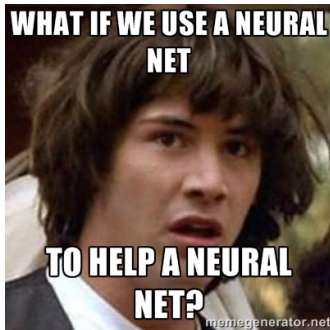
# What if . . .

- Even bidirectional encoder-decoders have a hard time with long sentences
- We need a way to keep track of what's already been translated and what to translate next

# What if . . .

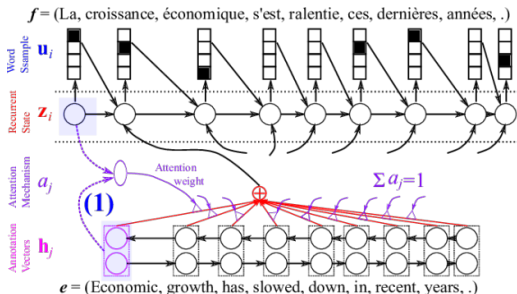- Even bidirectional encoder-decoders have a hard time with long sentences
- We need a way to keep track of what's already been translated and what to translate next



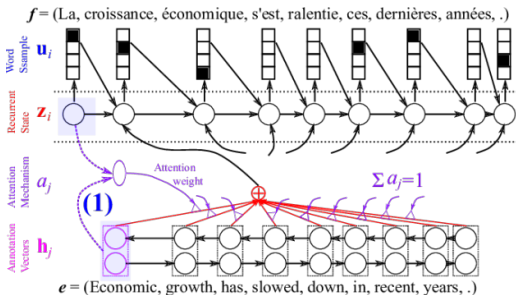- For neural nets, the solution is often more neural nets . . .

# Achtung, Baby!

- Attention-based decoding adds another network (**a**) that takes as input the encoder's hidden state (**h**) and the decoder's hidden state (**z**), and outputs a probability for each source word at each time step (when and where to pay attention) :

$$e_{i,j} = a(z_{i-1}, h_j) = v_a^\top \tanh(W_a s_{i-1} + U_a h_j)$$



$f$ = (La, croissance, économique, s'est, ralentie, ces, dernières, années, .)

$e$ = (Economic, growth, has, slowed, down, in, recent, years, .)

# Achtung, Baby!

- Attention-based decoding adds another network (**a**) that takes as input the encoder's hidden state (**h**) and the decoder's hidden state (**z**), and outputs a probability for each source word at each time step (when and where to pay attention) :

$$e_{i,j} = a(z_{i-1}, h_j) = v_a^\top \tanh(W_a s_{i-1} + U_a h_j)$$



$f$ = (La, croissance, économique, s'est, ralentie, ces, dernières, années, .)

$e$ = (Economic, growth, has, slowed, down, in, recent, years, .)

- The attention weights can also function as soft word alignments. They're trained on target-side MLE

# Image Caption Generation

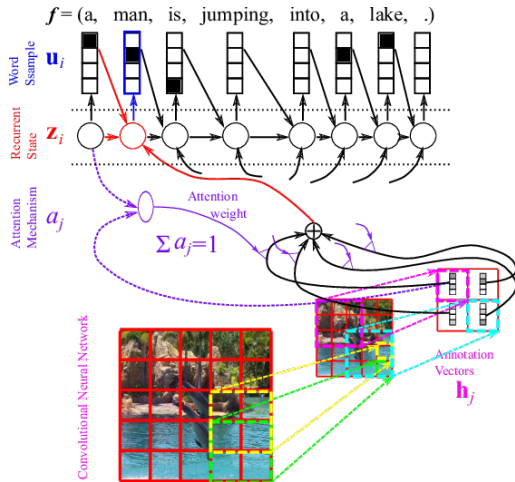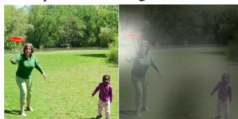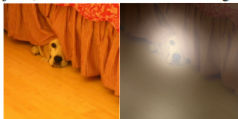- You can use attention-based decoding to give textual descriptions of images

# Image Caption Generation Examples

Figure 4. Examples of attending to the correct object (*white* indicates the attended regions, *underlines* indicated the corresponding word)
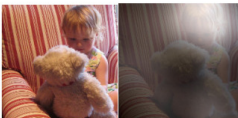


A woman is throwing a <u>frisbee</u> in a park.

A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

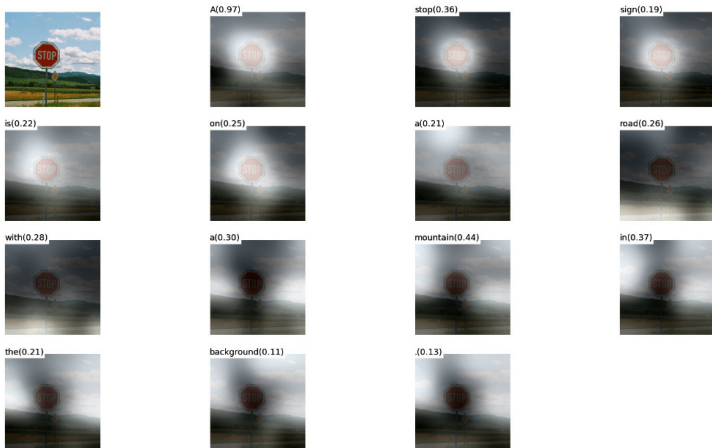A little <u>girl</u> sitting on a bed with a teddy bear.

A group of <u>people</u> sitting on a boat in the water.

A giraffe standing in a forest with <u>trees</u> in the background.

# Image Caption Generation, Step by Step



(b) A stop sign is on a road with a mountain in the background.