# Lecture Notes for
# **Machine Learning in Python**

## Professor Eric Larson
## **Seq-2-Seq and Transformers**

Archived

# Lecture Agenda

- Logistics
  - RNNs due **During Finals Time**
- Agenda
  - Sequence to sequence
  - Transformers

# Class Overview, by topic

**Table Data Visualization** → **Numpy, Pandas, Seaborn**
Overviews with some in-depth discussion

**Dimension Reduction and Image Processing** → **Scikit-learn, Scikit Image,**
Intuition only, Some mathematics

**Linear and Logistic Regression** → **Numpy, Recreate API for Scikit-learn**
Detailed mathematics for simple optimization
intuition for advanced optimization

**Neural Networks and Back Prop.** → **Numpy**
Detailed mathematics for NN operations

**Wide and Deep Networks**   **Convolutional Networks**   **Recurrent Networks** → **Keras, Tensorflow**
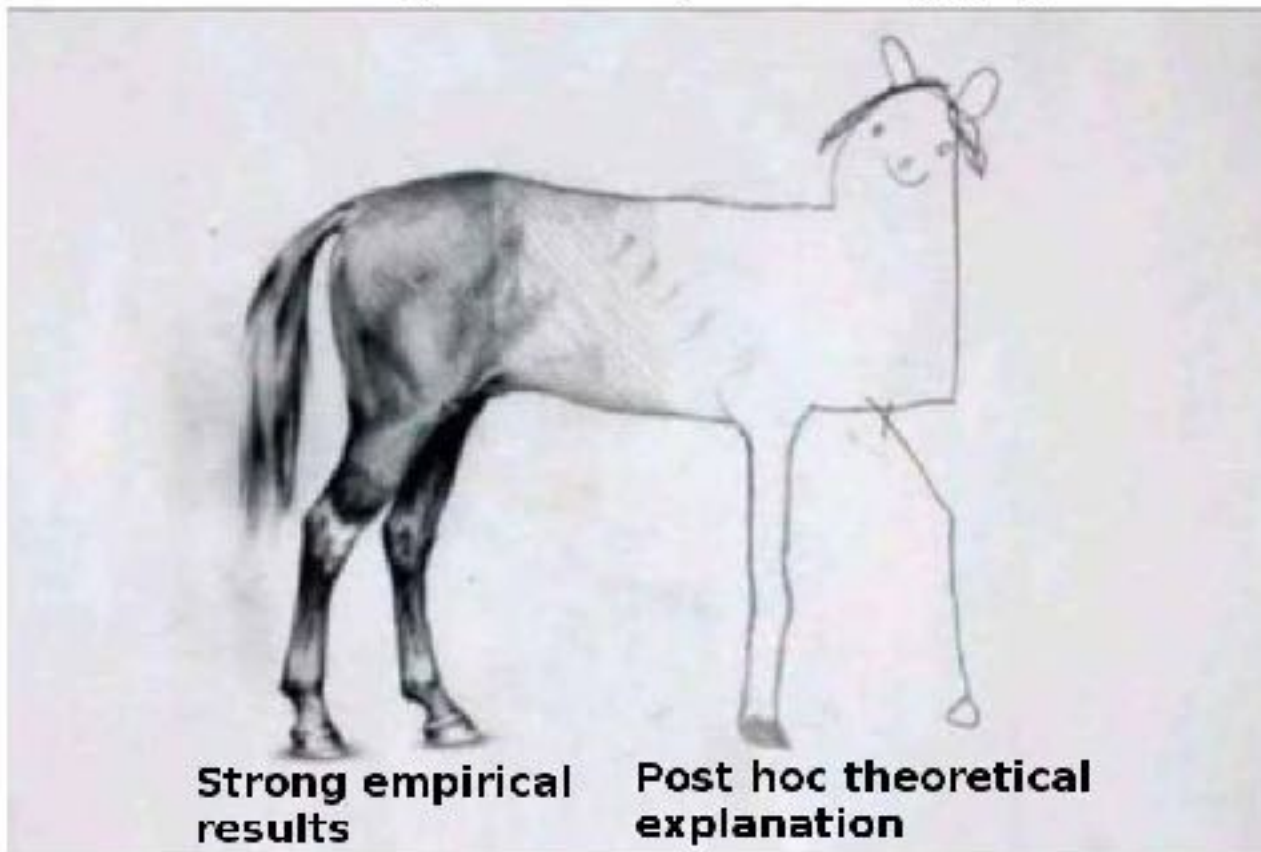Intuition, Detailed implement.

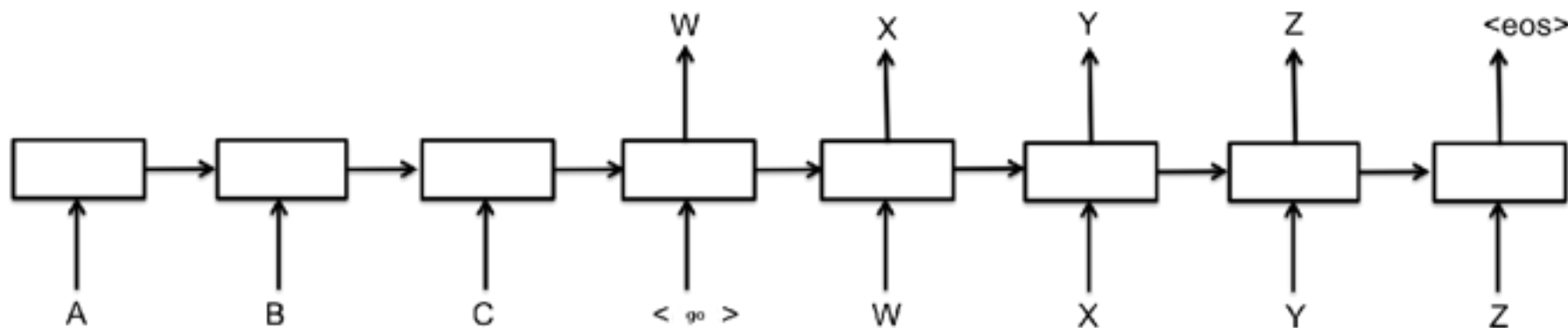**Ethics in Language Models**   **ConceptNet**
Case studies

# Last Time

# Sequence to Sequence



Anatomy of a deep learning paper

Strong empirical results

Post hoc theoretical explanation

# Modeling Sequence to Sequence

Need to translate outputs of unknown size.



- Additional Vocabulary Special Casing:
    - <UNKNOWN>, for unknown input or characters not included in vocabulary
    - <EOS>, end of sentence
    - <GO>, start output sequence
    - <DONTCARE>, outputs before <GO> command

*Sutskever et al.* Sequence to Sequence Learning with Neural Networks, arXiv. 2014
https://arxiv.org/pdf/1409.3215.pdf
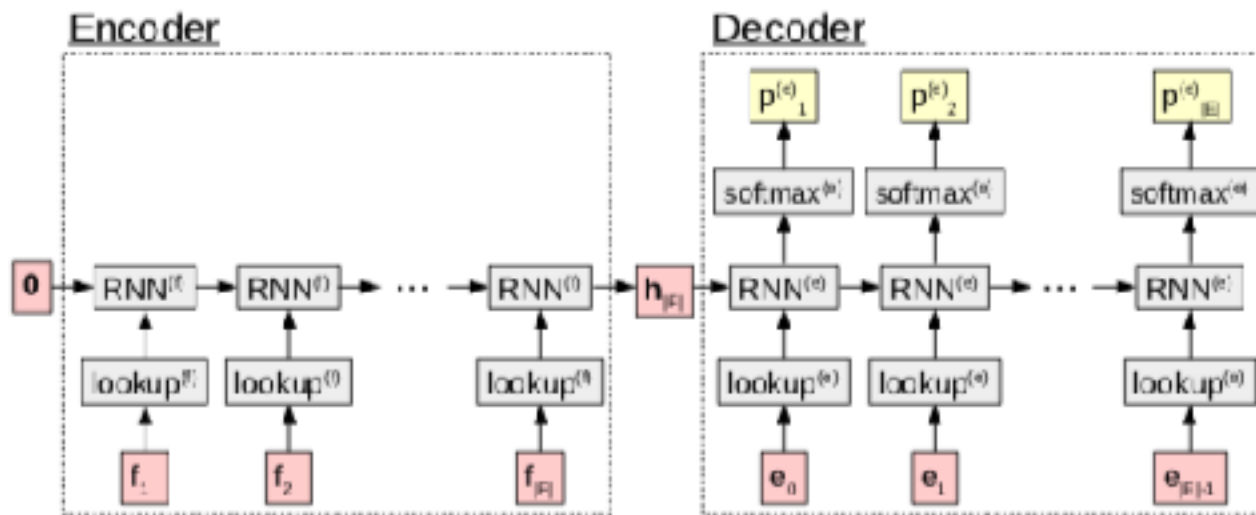
# Modeling Sequence to Sequence



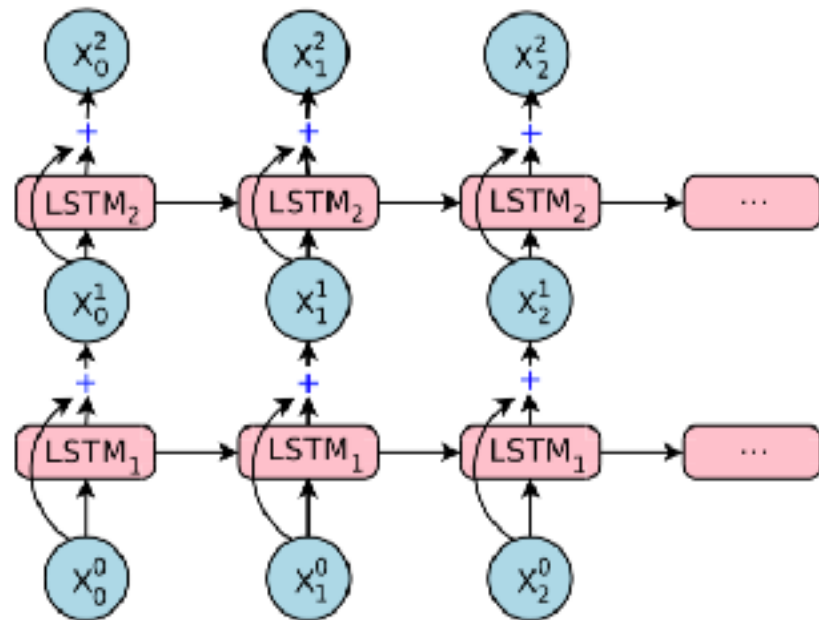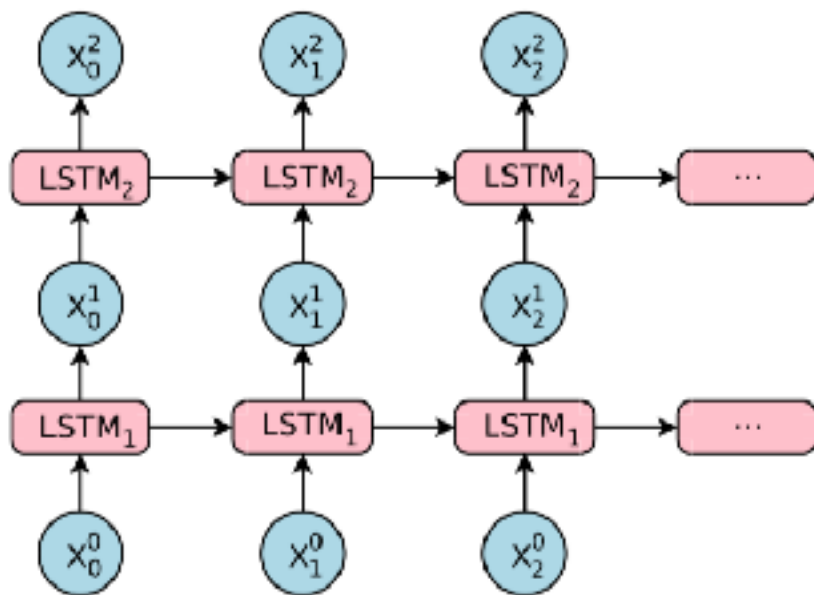Figure 21: A computation graph of the encoder-decoder model.

- **Training Process**: Give actual decoded letters for predicting next token
- **Decoding Process** can alter reliability of results:
    - Greedy Search, always choose most likely "next" symbol, seed
    - Keep list of "best" predictions for seeding (i.e., Beam Search)

Graham Neubig. 2017
Neural Machine Translation and
Sequence-to-sequence Models: A Tutorial
https://arxiv.org/pdf/1703.01619.pdf

https://github.com/m2dsupsdlclass/lectures-labs/blob/master/labs/07_seq2seq/Translation_of_Numeric_Phrases_with_Seq2Seq_rendered.ipynb
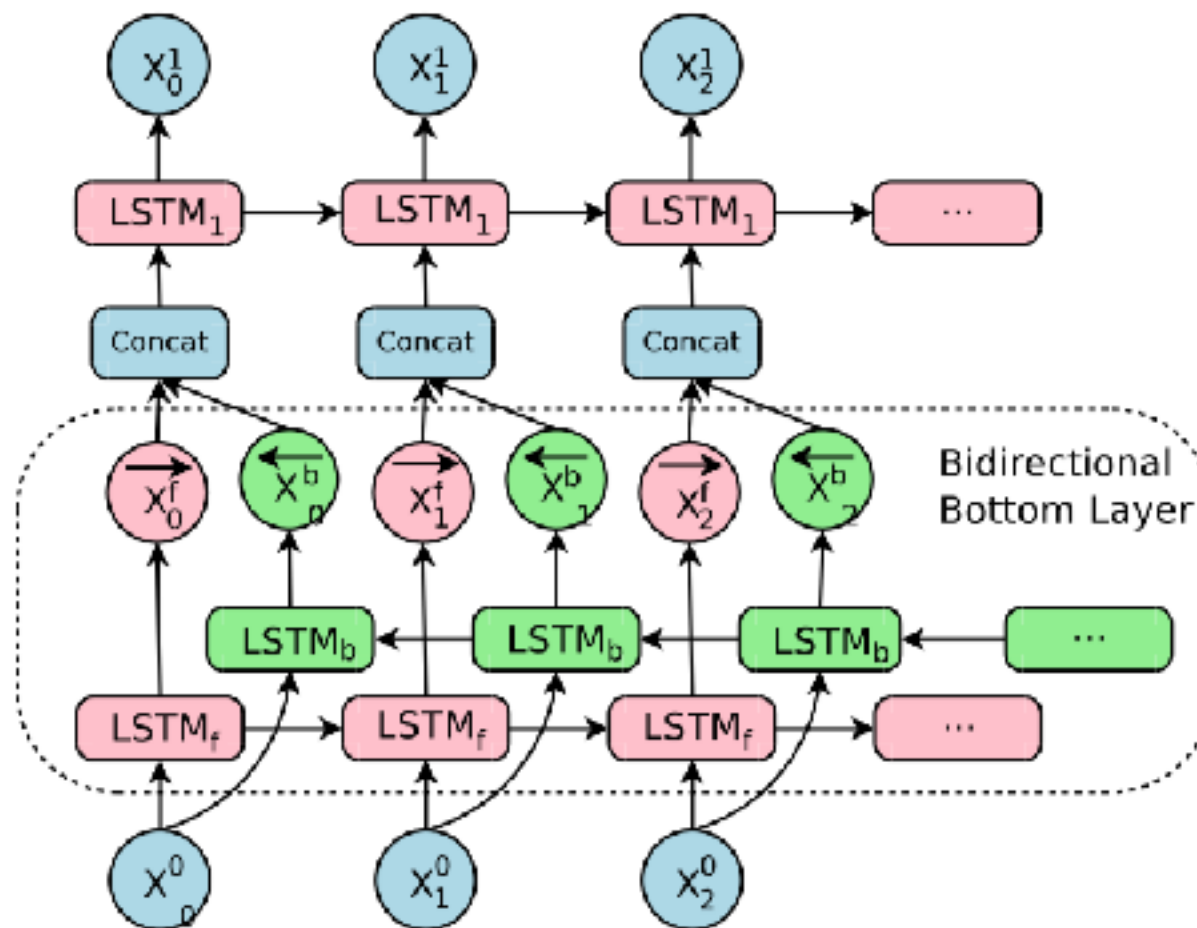
- Google, 2016



Google's Neural Machine Translation:
https://arxiv.org/pdf/1609.08144.pdf

# GNMT: Bidirectionality

- Google, 2016



Google Neural Machine Translation:
https://arxiv.org/pdf/1609.08144.pdf

# GNMT: Attention

- Google, 2016

$$s_t = AttentionFunction(\mathbf{y}_{i-1}, \mathbf{x}_t) \quad \forall t, \quad 1 \leq t \leq M$$

$$p_t = \exp(s_t) / \sum_{t=1}^{M} \exp(s_t) \quad \forall t, \quad 1 \leq t \leq M$$

$$\mathbf{a}_i = \sum_{t=1}^{M} p_t . \mathbf{x}_t$$

where $\mathbf{x}_t$ is state of the $t^{th}$ encoder
$\mathbf{y}_{i-1}$ is the state of the previous decoder
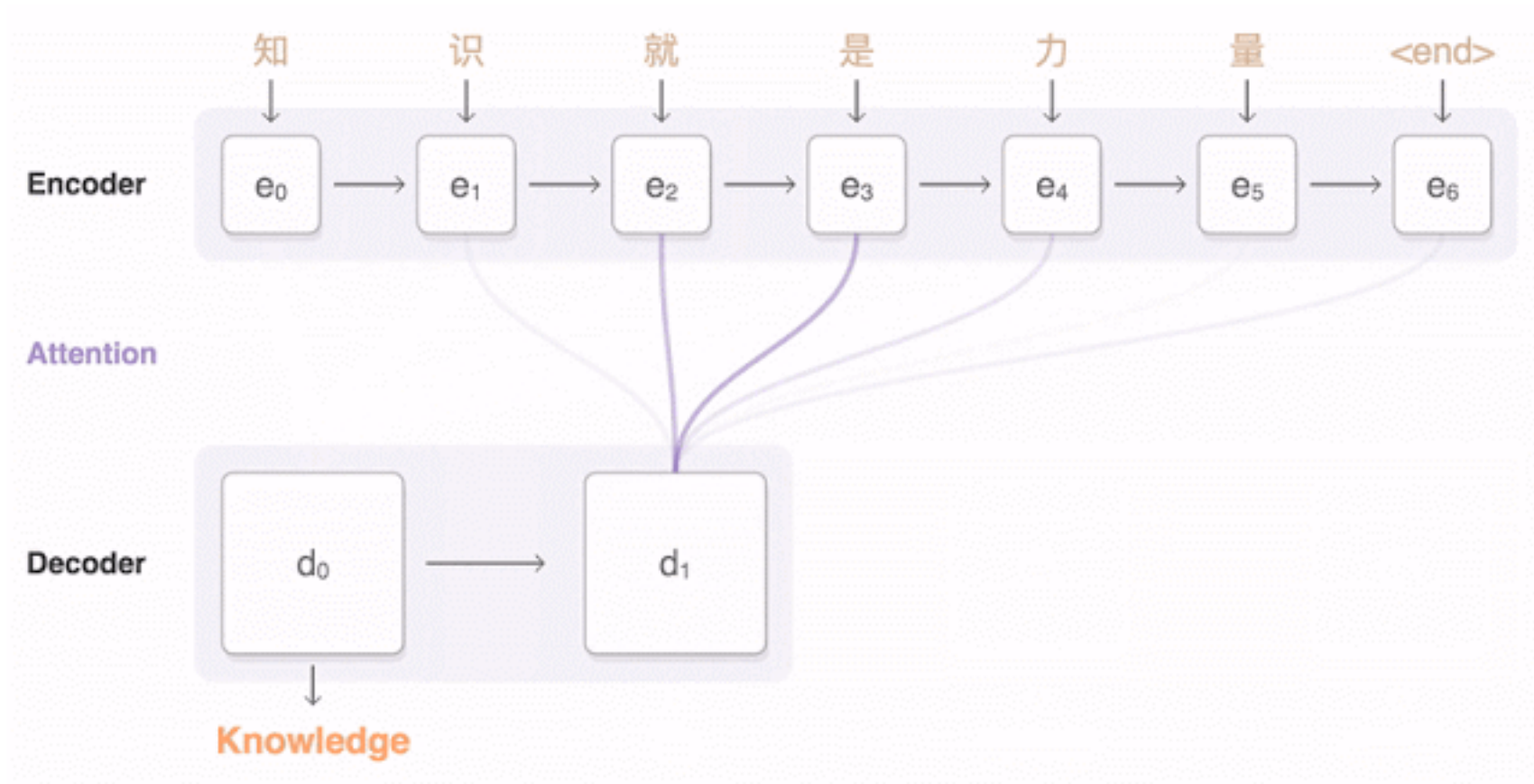and $\mathbf{a}_i$ is the input for the $i^{th}$ decoder

# GNMT: Attention

$$s_t = AttentionFunction(\mathbf{y}_{i-1}, \mathbf{x}_t) \quad \forall t, \quad 1 \leq t \leq M$$

$$p_t = \exp(s_t)/\sum_{t=1}^{M} \exp(s_t) \quad \forall t, \quad 1 \leq t \leq M$$
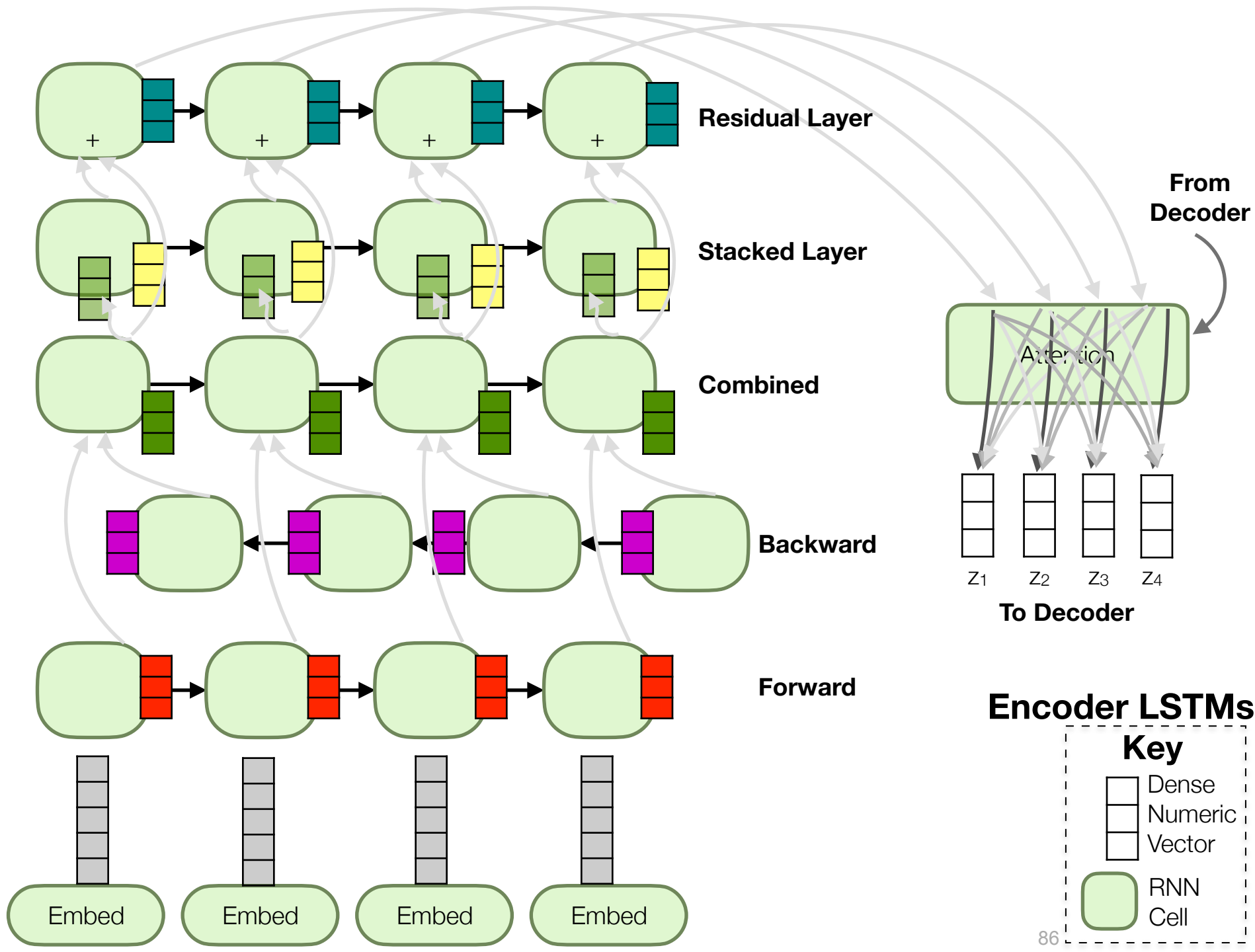
$$\mathbf{a}_i = \sum_{t=1}^{M} p_t . \mathbf{x}_t$$
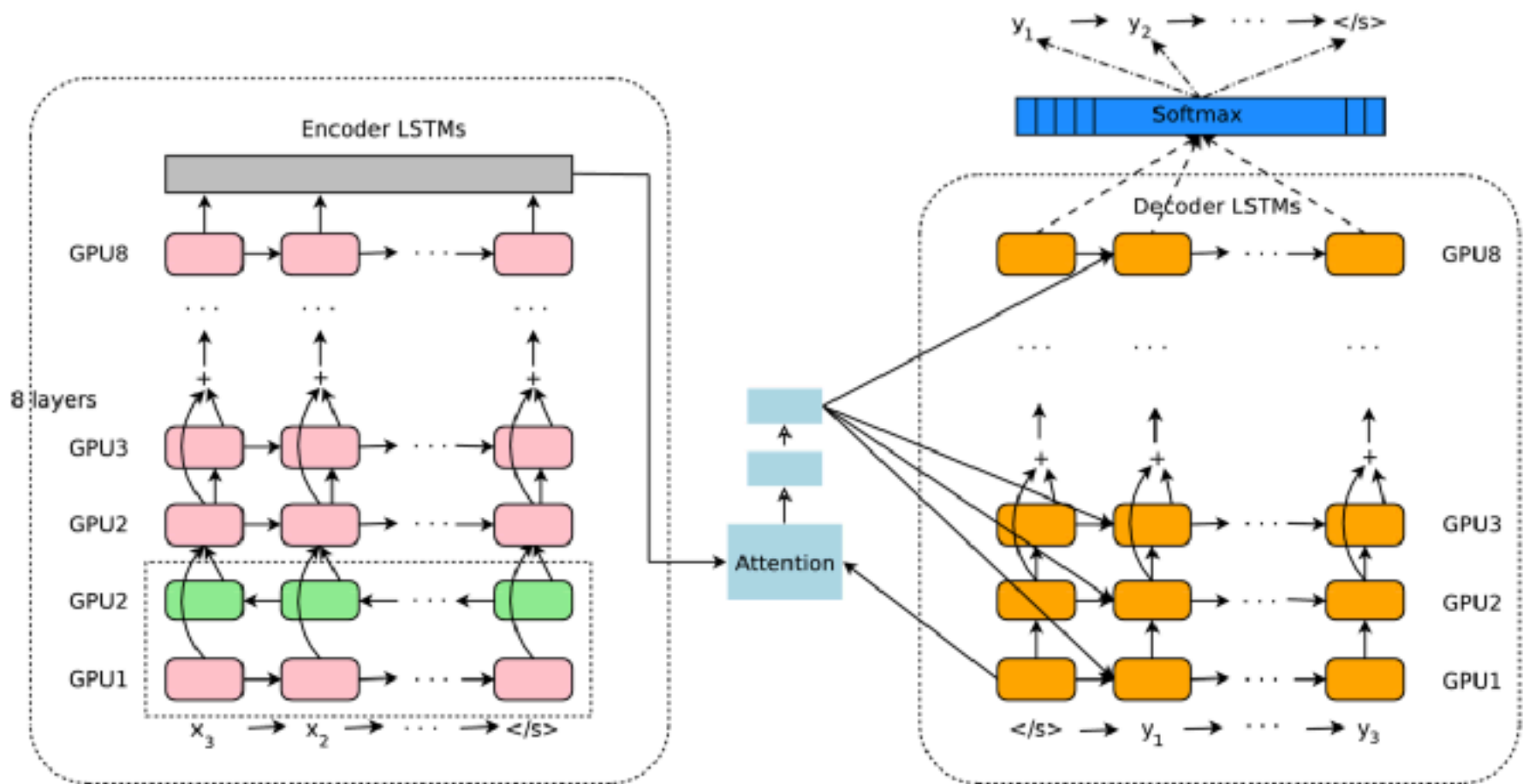
- Google, 2016



Google Neural Machine Translation:
https://arxiv.org/pdf/1609.08144.pdf
https://medium.com/@Synced/history-and-frontier-of-the-neural-machine-translation-dc981d25422d

**Residual Layer**

**Stacked Layer**

**Combined**

**Backward**

**Forward**

**From Decoder**

Attention

$Z_1$  $Z_2$  $Z_3$  $Z_4$

**To Decoder**

**Encoder LSTMs**

Embed    Embed    Embed    Embed

**Key**

Dense
Numeric
Vector

RNN
Cell

86

- Google, 2016



Google Neural Machine Translation:
https://arxiv.org/pdf/1609.08144.pdf

- Can translation also be done using only CNNs?
  - Yes, Facebook AI already did it,
  - 9 times faster than GNMT
  - Similar Performance
  - July, 2017



https://arxiv.org/pdf/1705.03122.pdf

… from Olivier Grisel

https://github.com/m2dsupsdlclass/lectures-labs/blob/master/labs/07_seq2seq/
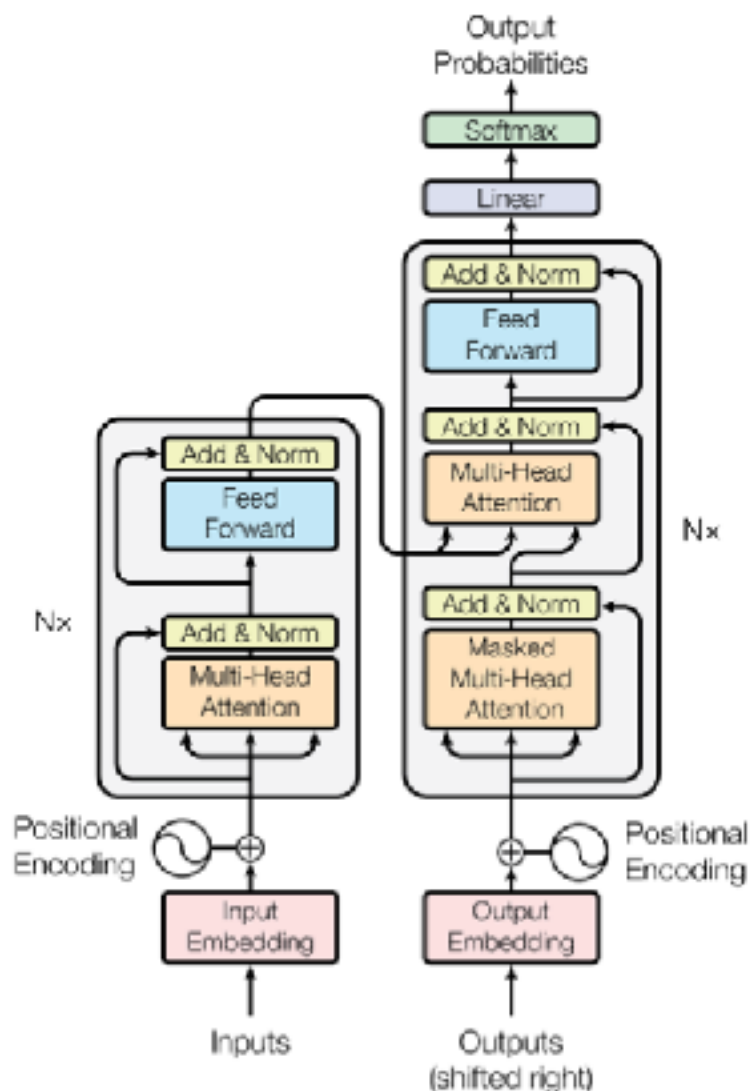Translation_of_Numeric_Phrases_with_Seq2Seq_rendered.ipynb
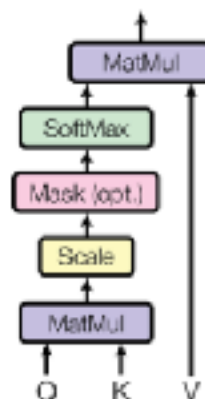
# Transformers

# Attention is All You Need

- Well, its a good paper title, but not exactly accurate
- Problem: recurrent networks are not inherently parallelized or efficient at remembering
- Convolution needs many examples from all different word positions (after flattening)
- Filters are not resilient to long-term relationships
- Transformer Solution:
  - Build attention into model from the **beginning**
  - Compare all words to each other through **multi-headed** attention
  - Define a notion of "**position**"in the sentence

91
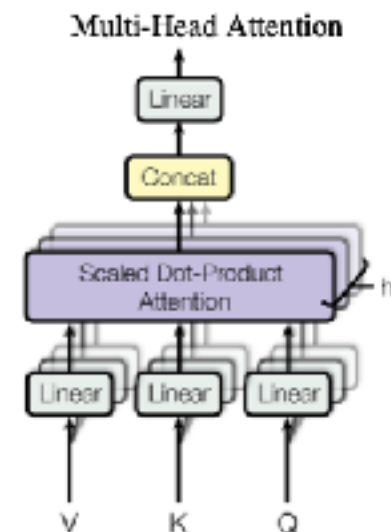
Scaled Dot-Product Attention

for each word

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

Multi-Head Attention

more than one
Q,K,V use in document

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf

| Input | Thinking | Machines |
|---|---|---|
| Embedding | $X_1$ | $X_2$ |

Learned Matrices

Outputs of Matrix Multiplications:

| Queries | $q_1$ | $q_2$ | $W^Q$ |
|---|---|---|---|
| Keys | $k_1$ | $k_2$ | $W^K$ |
| Values | $v_1$ | $v_2$ | $W^V$ |

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

Input

Embedding

Queries

Keys

Values

Score

Divide by 8 ($\sqrt{d_k}$)

Softmax

Softmax X Value

Sum

Calc. q, k, v for each word

**Thinking**  **Machines**

$x_1$  $x_2$

$q_1$  $q_2$

$k_1$  $k_2$

$v_1$  $v_2$

$q_1 \cdot k_1 = 112$    $q_1 \cdot k_2 = 96$

14    12

0.88    0.12

Calc weights for $z_1$

$v_1$    $v_2$

weighted sum for all words in document

$z_1$  attention for word 1

$z_2$  attention for word 2

Straight forward to do this operation in matrix form:

X    $W^Q$    Q

Thinking Machines

X    $W^K$    K

Thinking Machines

X    $W^V$    V

Thinking Machines

Q    $K^T$    V

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) \quad V$$

Z

$= \quad z_1 \atop z_2$

# Transformer: Multi-headed Attention



one row for each word

Size of row determined by $W^o$

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

- Objective: add notion of position to embedding
- Attempt in paper: add sin/cos to embedding
- But could be anything that encodes position

$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{\mathrm{model}}})$$

$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{\mathrm{model}}})$$

Now use the new embeddings, with position, into transformer architecture

POSITIONAL ENCODING

| 0 | 0 | 1 | 1 |
|---|---|---|---|

| 0.84 | 0.0001 | 0.54 | 1 |
|---|---|---|---|

| 0.91 | 0.0002 | -0.42 | 1 |
|---|---|---|---|

+      +      +

EMBEDDINGS   $X_1$     $X_2$     $X_3$

**Hypothesis**: Now the word proximity is encoded in the embedding matrix, with other pertinent information.  Well, it does help… so it could be true that this is a good way to do it.

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

Decoding time step: 1 ② 3 4 5 6

OUTPUT: I

Kencdec    Vencdec

Linear + Softmax

ENCODERS

DECODERS

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT: Je    suis    étudiant    PREVIOUS OUTPUTS: I

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

# Results

English German Translation quality



English French Translation Quality



https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html

**Implementations**:
- Not Native to Keras or Tensorflow, but many Open Source Implementations Exist
- Is Native to PyTorch

# Next time

- Class Retrospective

# TensorFlow

```python
with tf.variable_scope('rnn_cell'):
    W = tf.get_variable('W', [num_classes + state_size, state_size])
    b = tf.get_variable('b', [state_size], initializer=tf.constant_initializer(0.0))

def rnn_cell(rnn_input, state):
    with tf.variable_scope('rnn_cell', reuse=True):
        W = tf.get_variable('W', [num_classes + state_size, state_size])
        b = tf.get_variable('b', [state_size], initializer=tf.constant_initializer(0.0))
    return tf.tanh(tf.matmul(tf.concat(1, [rnn_input, state]), W) + b)

state = init_state
rnn_outputs = []
for rnn_input in rnn_inputs:
    state = rnn_cell(rnn_input, state)
    rnn_outputs.append(state)
final_state = rnn_outputs[-1]

#logits and predictions
with tf.variable_scope('softmax'):
    W = tf.get_variable('W', [state_size, num_classes])
    b = tf.get_variable('b', [num_classes], initializer=tf.constant_initializer(0.0))
logits = [tf.matmul(rnn_output, W) + b for rnn_output in rnn_outputs]
predictions = [tf.nn.softmax(logit) for logit in logits]

# Turn our y placeholder into a list labels
y_as_list = [tf.squeeze(i, squeeze_dims=[1]) for i in tf.split(1, num_steps, y)]

#losses and train_step
losses = [tf.nn.sparse_softmax_cross_entropy_with_logits(logit,label) for \
          logit, label in zip(logits, y_as_list)]
total_loss = tf.reduce_mean(losses)
train_step = tf.train.AdagradOptimizer(learning_rate).minimize(total_loss)
```

Lecture Notes for Machine Learning in Python    |    Professor Eric C. Larson

# recurrent networks

```python
def train_network(num_epochs, num_steps, state_size=4, verbose=True):
    with tf.Session() as sess:
        sess.run(tf.initialize_all_variables())
        training_losses = []
        for idx, epoch in enumerate(gen_epochs(num_epochs, num_steps)):
            training_loss = 0
            training_state = np.zeros((batch_size, state_size))
            if verbose:
                print("\nEPOCH", idx)
            for step, (X, Y) in enumerate(epoch):
                tr_losses, training_loss_, training_state, _ = \
                    sess.run([losses,
                              total_loss,
                              final_state,
                              train_step],
                                 feed_dict={x:X, y:Y, init_state:training_state})
                training_loss += training_loss_
                if step % 100 == 0 and step > 0:
                    if verbose:
                        print("Average loss at step", step,
                              "for last 250 steps:", training_loss/100)
                    training_losses.append(training_loss/100)
                    training_loss = 0

    return training_losses
```

```python
def train_network(num_epochs, num_steps, state_size=4, verbose=True):
    with tf.Session() as sess:
        sess.run(tf.initialize_all_variables())
        for idx, epoch in enumerate(gen_epochs(num_epochs, num_steps)):
            training_state = np.zeros((batch_size, state_size))
            for X, Y in epoch:
                tr_losses, training_loss_, training_state, _ = \
                    sess.run([losses,
                              total_loss,
                              final_state,
                              train_step],
                             feed_dict={x:X, y:Y, init_state:training_state})
```

# TensorFlow (simplified)

```python
cell = tf.nn.rnn_cell.BasicRNNCell(state_size)
rnn_outputs, final_state = tf.nn.rnn(cell, rnn_inputs, initial_state=init_state)


loss_weights = [tf.ones([batch_size]) for i in range(num_steps)]
losses = tf.nn.seq2seq.sequence_loss_by_example(logits, y_as_list, loss_weights)


x = tf.placeholder(tf.int32, [batch_size, num_steps], name='input_placeholder')
y = tf.placeholder(tf.int32, [batch_size, num_steps], name='labels_placeholder')
init_state = tf.zeros([batch_size, state_size])

x_one_hot = tf.one_hot(x, num_classes)
rnn_inputs = tf.unpack(x_one_hot, axis=1)


cell = tf.nn.rnn_cell.BasicRNNCell(state_size)
rnn_outputs, final_state = tf.nn.rnn(cell, rnn_inputs, initial_state=init_state)


with tf.variable_scope('softmax'):
    W = tf.get_variable('W', [state_size, num_classes])
    b = tf.get_variable('b', [num_classes], initializer=tf.constant_initializer(0.0))
logits = [tf.matmul(rnn_output, W) + b for rnn_output in rnn_outputs]
predictions = [tf.nn.softmax(logit) for logit in logits]


y_as_list = [tf.squeeze(i, squeeze_dims=[1]) for i in tf.split(1, num_steps, y)]


loss_weights = [tf.ones([batch_size]) for i in range(num_steps)]
losses = tf.nn.seq2seq.sequence_loss_by_example(logits, y_as_list, loss_weights)
total_loss = tf.reduce_mean(losses)
train_step = tf.train.AdagradOptimizer(learning_rate).minimize(total_loss)
```