# Deep Learning - COSC2779
## Sequential Data - Applications

Dr. Ruwan Tennakoon

**RMIT**
UNIVERSITY

Sep 13, 2021

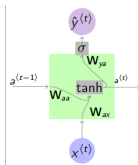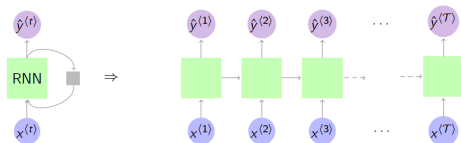**Reference:** *Chapter 10: Ian Goodfellow et. al., "Deep Learning", MIT Press, 2016.*

# Outline

# Revision

**Sequential Data:**

- Data points with variable length.
- Order of the data matters
- Shared feature across time is useful.

**Sequential Data:**

- Data points with variable length.
- Order of the data matters
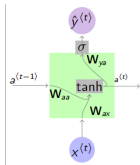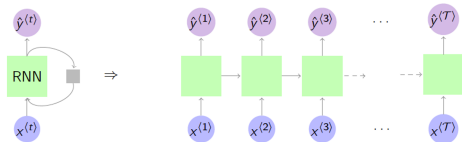- Shared feature across time is useful.



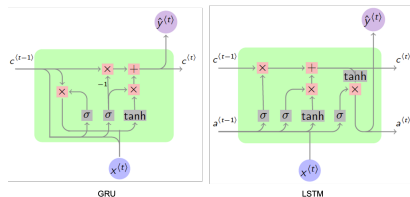Simple RNN Cell

**Sequential Data:**

- Data points with variable length.
- Order of the data matters
- Shared feature across time is useful.

**Improvements:**

- Capture long range dependencies:



GRU          LSTM

- Bi-Directional: Use both historical and future information.
- Deep RNN: Increase model capacity.



Simple RNN Cell

- Understand how RNNs are used in key application areas:
  - Text Classification.
  - Machine Translation.
  - Speech Recognition.
- Will cover the fundamental concepts that will enable you to explore the state-of-the-art.
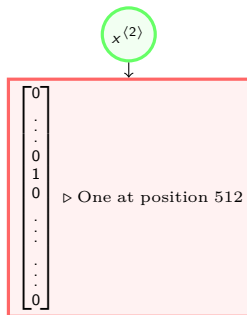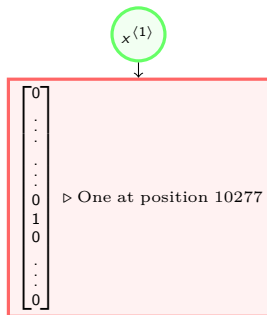
# Outline

# Outline

| | My | experience | so | far | has | been | fantastic |
|---|---|---|---|---|---|---|---|
| | 10277 | 512 | 12011 | 611 | 854 | 325 | 625 |
| $\mathbf{x}^{(i)} =$ | $x^{\langle 1 \rangle}$ | $x^{\langle 2 \rangle}$ | $x^{\langle 3 \rangle}$ | $x^{\langle 4 \rangle}$ | $x^{\langle 5 \rangle}$ | $x^{\langle 6 \rangle}$ | $x^{\langle 7 \rangle}$ |

| Index | Word |
|---|---|
| 1 | a |
| 2 | ability |
| 3 | able |
| . . . | . . . |
| 325 | been |
| . . . | . . . |
| 512 | experience |
| . . . | . . . |
| 611 | far |
| . . . | . . . |
| 625 | fantastic |
| . . . | . . . |
| 854 | has |
| . . . | . . . |
| 10277 | My |
| . . . | . . . |
| 12011 | So |
| . . . | . . . |



$x^{(i)}$ is a matrix with dimensions 20,000 x 7.

Vocabulary: Assume 20k Words

# One-Hot Representation

|  | My | experience | so | far | has | been | fantastic |
|---|---|---|---|---|---|---|---|
|  | 10277 | 512 | 12011 | 611 | 854 | 325 | 625 |
| $\mathbf{x}^{(i)} =$ | $x^{\langle 1 \rangle}$ | $x^{\langle 2 \rangle}$ | $x^{\langle 3 \rangle}$ | $x^{\langle 4 \rangle}$ | $x^{\langle 5 \rangle}$ | $x^{\langle 6 \rangle}$ | $x^{\langle 7 \rangle}$ |

$x^{\langle 1 \rangle}$

$x^{\langle 2 \rangle}$

$$\begin{bmatrix} 0 \\ \vdots \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$ ▷ One at position 10277

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ \vdots \\ 0 \end{bmatrix}$$ ▷ One at position 512

vocabulary:
- apple: [1, 0, 0, 0, 0]
- orange: [0, 1, 0, 0, 0]
- cricket: [0, 0, 1, 0, 0]
- kabaddi: [0, 0, 0, 1, 0]
- man: [0, 0, 0, 0, 1]

$|apple - orange| = 2$
$|apple - cricket| = 2$

$x^{(i)}$ is a matrix with dimensions 20,000 x 7.

- One-Hot does not capture similarity between words.

  Seen During Training: | Joe Root was in the English cricket team. |

  During inference: | Dan Murphy was in the Australian Kabaddi squad. |

- If the representation for 'cricket' $\leftrightarrow$ 'kabaddi' and 'team' $\leftrightarrow$ 'squad' was close, then it is clear that "Dan Murphy" in this context is a persons name - not an organization.

- A representation that capture relationships between words can enable inferring on unseen sequences during training.

| | cricket | kabaddi | orange | apple | man |
|---|---|---|---|---|---|
| Gender | 0.1 | 0.01 | 0 | 0 | -1 |
| age | 0 | 0 | 0 | 0 | 1 |
| food | 0 | 0 | 1 | 1 | 0 |
| sport | 1 | 1 | 0 | 0 | 0 |
| alive | 0 | 0 | 0 | 0 | 1 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| noune | 1 | 1 | 1 | 1 | 1 |
| | $e_{cricket}$ | $e_{kabaddi}$ | $e_{orange}$ | $e_{apple}$ | $e_{man}$ |

- $o_i \in \mathbb{R}^{20,000} \rightarrow$ One hot encoding of word i in vocabulary.

- $e_i \in \mathbb{R}^{300} \rightarrow$ Feature encoding of word i in vocabulary.

- $E \in \mathbb{R}^{300 \times 20,000} \rightarrow$ Feature encoding of all words in vocabulary (Embedding Matrix).
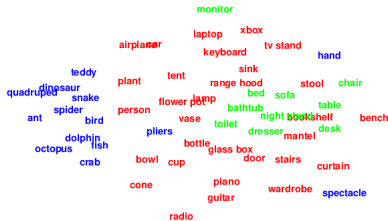
$$e_i = E \cdot o_i$$

Assume we have 300 features (rows in the table above).

The vocabulary size is 20k words.

Feature values can be floating point numbers (not binary).

$$|e_{cricket} - e_{kabaddi}| < |e_{cricket} - e_{orange}|$$

# Properties of word Embedding

Project word embedding from 300D to 2D using method like T-SNE.
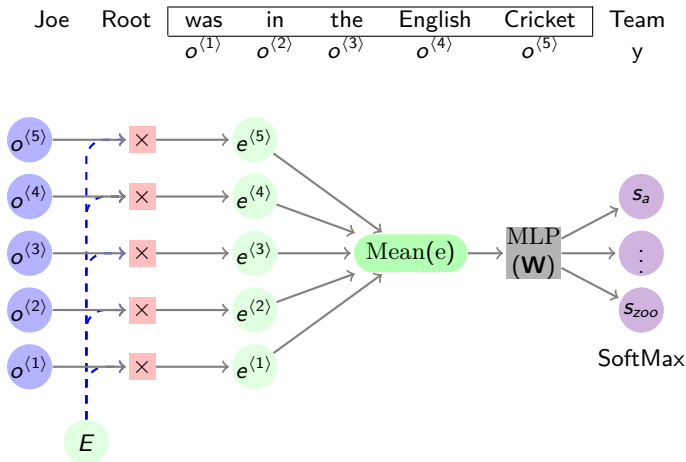
Learn by language modelling. E.g. Predict Next word:

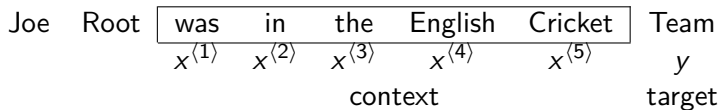> Joe Root was in the English cricket ???.

Can create a supervised learning task where **x** is a fixed length sequence (window selected by user) and $y$ is the next word in the sequence.

| Joe | Root | was | in | the | English | Cricket | Team |
|-----|------|-----|-----|-----|---------|---------|------|
| | | $x^{\langle 1 \rangle}$ | $x^{\langle 2 \rangle}$ | $x^{\langle 3 \rangle}$ | $x^{\langle 4 \rangle}$ | $x^{\langle 5 \rangle}$ | $y$ |

A Neural Probabilistic Language Model

Train for parameters $E$ and the weights of the MLP ($\mathbf{W}$) via back-propagation.

Joe  Root  | was | in | the | English | Cricket | Team

$$x^{\langle 1 \rangle} \quad x^{\langle 2 \rangle} \quad x^{\langle 3 \rangle} \quad x^{\langle 4 \rangle} \quad x^{\langle 5 \rangle}$$

context

$y$

target

Context can be:

- Last 5 words.
- 5 words before and 5 words after.
- Last one word.
- Near by word.

**Word2Vec and GloVe are more recent practical word embedding learning techniques.**

The *skip-grams*:

- Randomly pick a word as context, $c$.
- Randomly pick another word within a window of $c$ as the target, $t$.

"Joe Root was in the English Cricket Team"

| Context | Target |
|---------|--------|
| 'English' | 'Team' |
| 'Team' | 'the' |

Large number of Context, target pairs are picked from a large text corpus (e.g. Wikipedia)

Efficient Estimation of Word Representations in Vector Space

RMIT
UNIVERSITY

"Joe Root was in the English Cricket Team"

| Context | Target |
|---------|--------|
| 'English' | 'Team' |
| 'Team' | 'cricket' |
| ⋮ | ⋮ |

$$p(t \mid c) = \frac{e^{W^\top e_c}}{\sum_{j=1}^{N_v} e^{W^\top e_j}}$$

$N_v$ is the vocabulary size, which is very large.

Word2Vec employs hierarchical SoftMax to make it efficient.



SoftMax

"Joe Root was in the English Cricket Team"

| Context | Word | Target |
|---------|------|--------|
| 'English' | 'Team' | 1 |
| 'English' | 'apple' | 0 |
| 'English' | 'bus' | 0 |

$$p(y = 1 \mid c, t) = \sigma \left( W^{\top} e_c \right)$$

Train multiple sigmoid NN instead of the SoftMax.

More recent model:
Paper: GloVe: Global Vectors for Word Representation



Sigmoid

# Outline

## Sentiment Classification

**"The process of computationally identifying and categorizing opinions expressed in a piece of text"**

Sentiment analysis models focus on:

- Polarity: positive, negative, neutral.
- Feelings/Emotions: angry, happy, sad, etc.
- Intentions: interested v. not interested.

| **x** | $y$ |
|---|---|
| "My experience so far has been fantastic" | Positive |
| "Your support team is useless" | Negative |
| $\vdots$ | $\vdots$ |

Usually the datasets are NOT very large.

| Your $x^{\langle 1 \rangle}$ | support $x^{\langle 2 \rangle}$ | team $x^{\langle 3 \rangle}$ | was $x^{\langle 4 \rangle}$ | useless $x^{\langle 5 \rangle}$ | Negative $y$ |
|---|---|---|---|---|---|



- Dataset can be too small to learn a good word embedding.
- Does not take the order of the words into account.

Train for parameters $E$ and the weights of the MLP (**W**) via back-propagation.

Dataset can be too small to learn a good word embedding.

Transfer learning with word embedding:

1. Learn a word embedding using a large corpus (or download a word embedding).
2. Transfer embedding to the model (Set $E$ to the pre-learned word embedding). Learn the model parameters of the task (**W**).
3. (Optional) Fine tune word embedding.

Does not take the order of the words into account.



"Completely lacks *good* food or *good* service." $\rightarrow$ ?

"Completely lacks *good* food or *good* service." $\rightarrow$ ?



RNN cell can be: LSTM/GRU

Bi-Directional to capture future information.

Deep models to increase capacity.

# Outline

**RMIT**
UNIVERSITY

Both input and output are sequences.

**English**          **Sinhala**
"How are you today"   ⇔   "ada obata kohomada?"

- The two sequences are not of the same length.
- One-to-one translation does not work.

| ada | obata | kohomada |
| :-: | :-: | :-: |
| ↓ | ↓ | ↓ |
| today | you | how |

| How $x^{\langle 1 \rangle}$ | are $x^{\langle 2 \rangle}$ | you $x^{\langle 3 \rangle}$ | today $x^{\langle 4 \rangle}$ | $\Leftrightarrow$ | ada $y^{\langle 1 \rangle}$ | obata $y^{\langle 2 \rangle}$ | kohomada $y^{\langle 3 \rangle}$ |
|---|---|---|---|---|---|---|---|

Sequence to Sequence Learning with Neural Networks

Sequence to Sequence Learning with Neural Networks

The decoder we use in machine translation is similar to the one used for language modelling (e.g. text generation).

In language modelling we are interested in learning:

$$p\left(y^{\langle 1 \rangle}, y^{\langle 2 \rangle}, \cdots, y^{\langle \mathcal{T}_y \rangle}\right)$$

In machine translation we are interested in learning a conditional model:

$$p\left(y^{\langle 1 \rangle}, y^{\langle 2 \rangle}, \cdots, y^{\langle \mathcal{T}_y \rangle} \mid x^{\langle 1 \rangle}, x^{\langle 2 \rangle}, \cdots, x^{\langle \mathcal{T}_x \rangle}\right)$$

We can use the same encoder decoder architecture to do image captioning. The main difference would be that the encoder in this task is a CNN (not a RNN).



Image: Python based Project – Learn to Build Image Caption Generator with CNN & LSTM

If interested, a nice worked example at: Python based Project – Learn to Build Image Caption Generator with CNN & LSTM

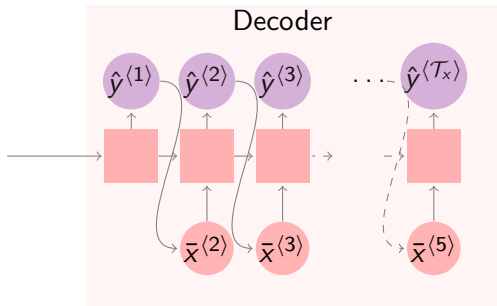Deep Captioning with Multimodal Recurrent Neural Networks (m-RNN)

In machine translation we are interested in learning a conditional model:

$$\arg \max_{y^{\langle 1 \rangle}, y^{\langle 2 \rangle}, \cdots} p\left(y^{\langle 1 \rangle}, y^{\langle 2 \rangle}, \cdots, y^{\langle \mathcal{T}_y \rangle} \mid \mathbf{x}\right)$$

Option 1: **Greedy search**.

- Pick the word with maximum probability using $y^{\langle 1 \rangle}$
- use it as the input to the next step
- Repeat till end of sentence.
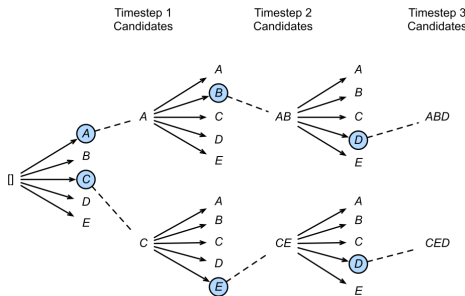
Sinhalese sentence: "Ruwan November mase lankavata yanwa"
Possible English translations:

- <u>Ruwan is</u> visiting Sri Lanka in November.
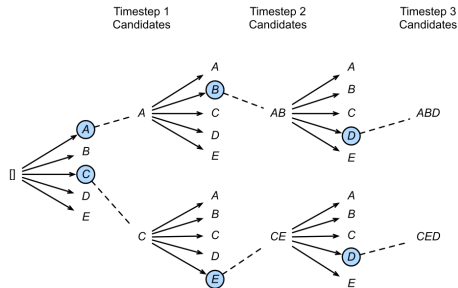- <u>Ruwan is</u> going to be visiting Sri Lanka in November.



Decoder

The first sentence is a better translation, but the greedy search might pick the second one because going is more common than visiting.
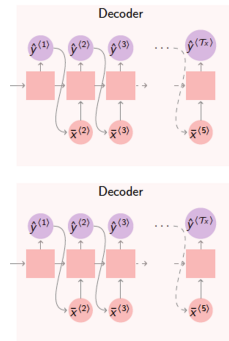
Instead of picking one candidate at each time step, you can pick best k candidates. Here the k is known as the beam width.

# Beam Search

At each time-step you need k copies of the network (k=2 in example).



- <u>Ruwan is</u> visiting Sri Lanka in November.

- <u>Ruwan is</u> going to be visiting Sri Lanka in November.

Assume you are given the following paragraph to be translate to another language:

"*Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.*"

Would you:

- Read the complete paragraph and then start translating
- Translate one small part at a time.

Assume you are given the following paragraph to be translate to another language:

"*Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.*"
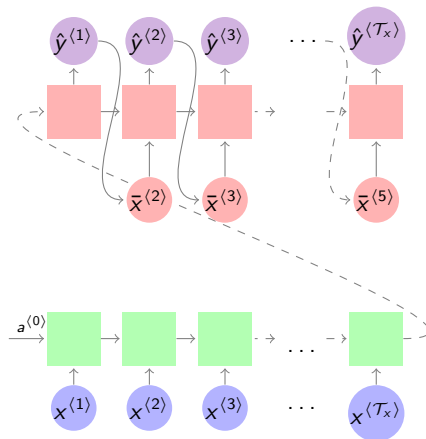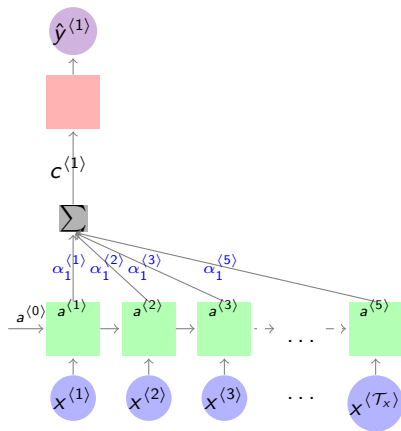
Would you:

- Read the complete paragraph and then start translating
- Translate one small part at a time.

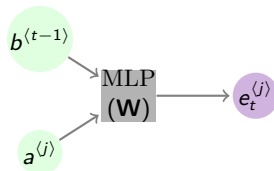**How to decide what part to be looked at in making the decision?**

Use a small neural network to decide the attention in making each decision.


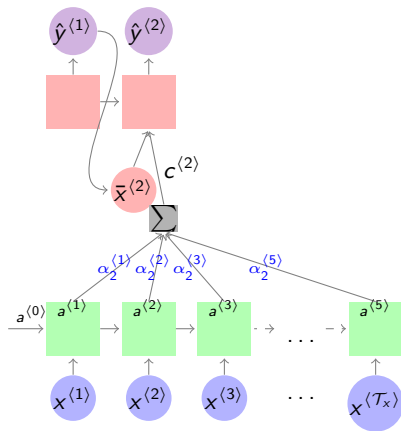
Machine translation with regular RNN.
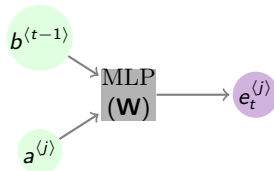
$$c^{\langle t \rangle} = \sum_j \alpha_t^{\langle j \rangle} a^{\langle j \rangle}$$

$$\alpha_t^{\langle j \rangle} = \frac{\exp(e_t^{\langle j \rangle})}{\sum_j \exp(e_t^{\langle j \rangle})}$$

$b^{\langle t-1 \rangle}$ is the cell state from the previous time step of the decoder.

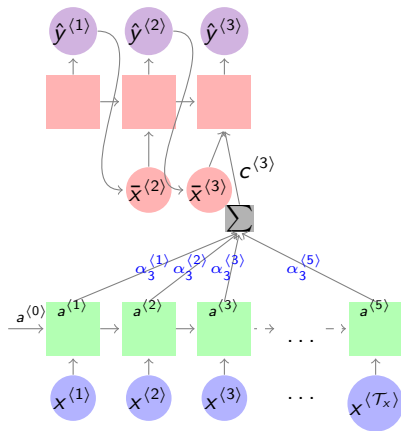Neural Machine Translation by Jointly Learning to Align and Translate
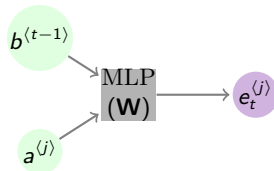
$$c^{\langle t \rangle} = \sum_j \alpha_t^{\langle j \rangle} a^{\langle j \rangle}$$

$$\alpha_t^{\langle j \rangle} = \frac{\exp(e_t^{\langle j \rangle})}{\sum_j \exp(e_t^{\langle j \rangle})}$$

$b^{\langle t-1 \rangle}$ is the cell state from the previous time step of the decoder.

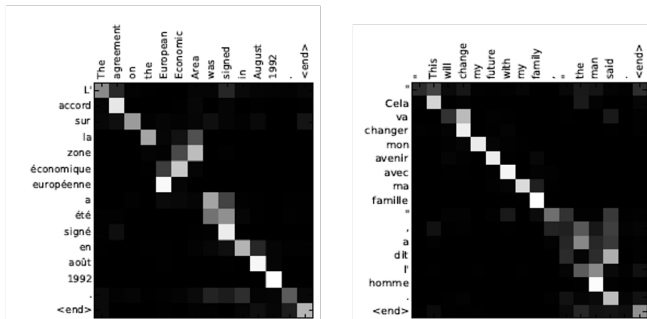Neural Machine Translation by Jointly Learning to Align and Translate

$$c^{\langle t \rangle} = \sum_j \alpha_t^{\langle j \rangle} a^{\langle j \rangle}$$

$$\alpha_t^{\langle j \rangle} = \frac{\exp(e_t^{\langle j \rangle})}{\sum_j \exp(e_t^{\langle j \rangle})}$$

$b^{\langle t-1 \rangle}$ is the cell state from the previous time step of the decoder.
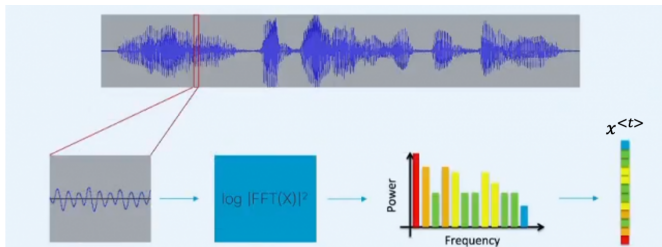
Neural Machine Translation by Jointly Learning to Align and Translate

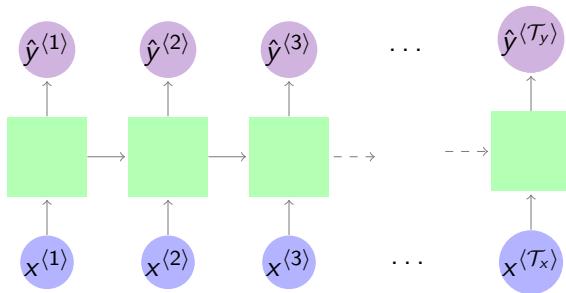Attention weights ($\alpha$) of English to French translation.

# Outline

Speech signals are usually represented as a spectrogram.



$x^{\langle t \rangle}$ is a vector where number of elements equal the number of frequencies in the spectrogram.
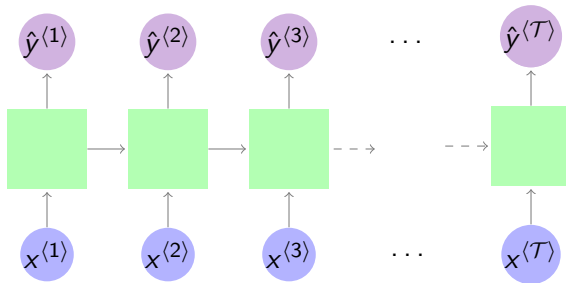
We can use a many-to-many RNN for speech recognition (e.g. Attention model).



**Problem:** $x$ vectors are generated at fixed time intervals (higher rate e.g. 100Hz). However the words in text are not (much lower rate). $\mathcal{T}_y \ll \mathcal{T}_x$

Allow multiple outputs for the same text-character.



Output of RNN: he__ll_l__oo_~wo_rr__l_d
Collapse repeated characters not separated by "_"
CTC Output: "hello world"

Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks

**RMIT**
UNIVERSITY

- How to represents input: Text, Speech
- Main considerations in applying RNN to real-world-applications.
- Encoder-decoder networks for sequence to sequence models.

**Next week**: Representation learning.

**Lab**: Sequence to Sequence model