



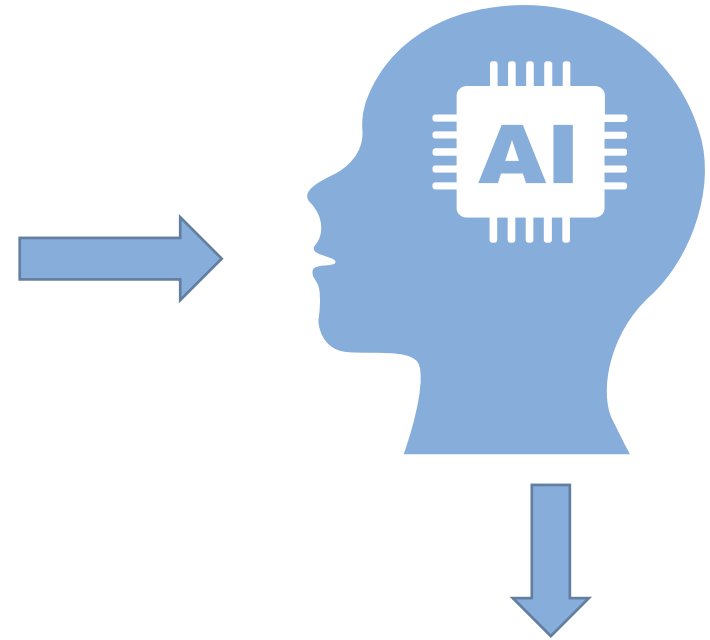
Extract Information From Contract

Using BERT (Deep Learning , NLP)

Extract Information From Contract

Contract

\uffeff_____ .. . II . สัญญาเลขที่ TMV/CM212015-1575/BBB/ys Project : FTTX Fund Code : Fund Center : สัญญาว่าจ้างสร้าง และ/หรือ ปรับปรุงข่ายสายสัญญาณโทรคมนาคม สัญญาฉบับนี้ทำขึ้นเมื่อวันที่ 01 .ค. 2558 ณบริษัท โทร มูฟ จำกัดระหว่าง (1) บริษัท โทร มูฟ จำกัด โดย นาย ศุภชัย เจียรนวนนท์ และศาสตราจารย์พิเศษ อสิก อิศวานนท์ กรรมการผู้มีอำนาจ สำนักงานตั้งอยู่เลขที่ 18 อาคารทรูทาวเวอร์ ถนนรัชดาภิเษก แขวงห้วยขวาง เขตห้วยขวาง กรุงเทพมหานคร ซึ่งต่อไปนี้จะเรียกว่า "ผู้ว่าจ้าง" ฝ่ายหนึ่ง กับ อพ (2) **ห้างหุ้นส่วนจำกัด บรอดแบนด์** โดยนาย ปัญญาวุฒิ พลราช หุ้นส่วนผู้จัดการสำนักงานตั้งอยู่เลขที่ 154 หมู่ที่ 12 ตำบลกุดลาด อำเภอเมืองอุบลราชธานี จังหวัดอุบลราชธานี ซึ่งต่อไปนี้จะเรียกว่า "ผู้รับจ้าง"อีกฝ่ายหนึ่ง โฉกานง CA2-007 หน้า 19 Standard Form 2015



Result

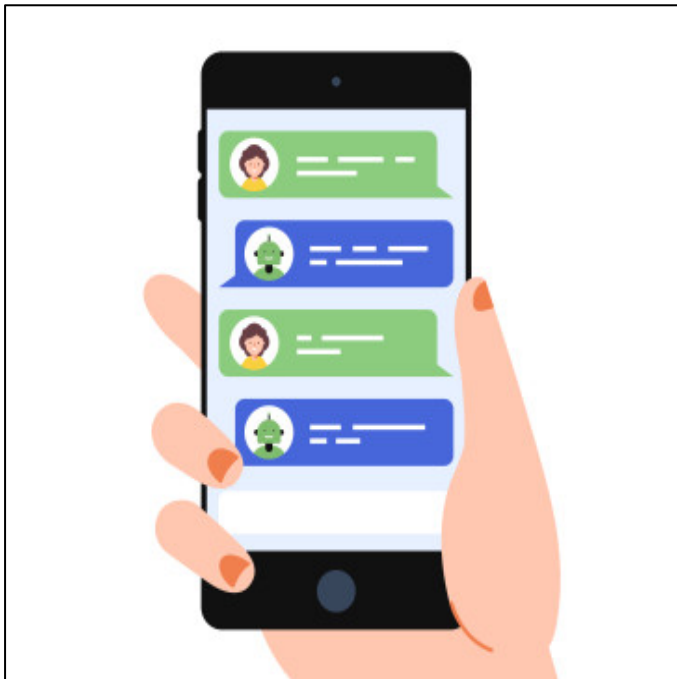
ห้างหุ้นส่วนจำกัด หรือ บริษัทใดเป็นผู้รับจ้าง ?

ห้างหุ้นส่วนจำกัด บรอดแบนด์

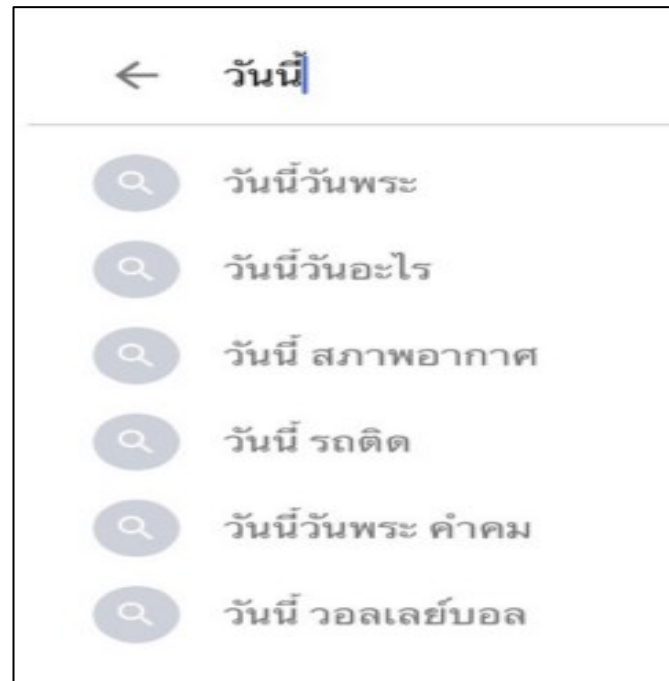
What is Natural Language Processing?

Natural Language Processing, or NLP for short, is broadly defined as the automatic manipulation of natural language, like speech and text, by software.

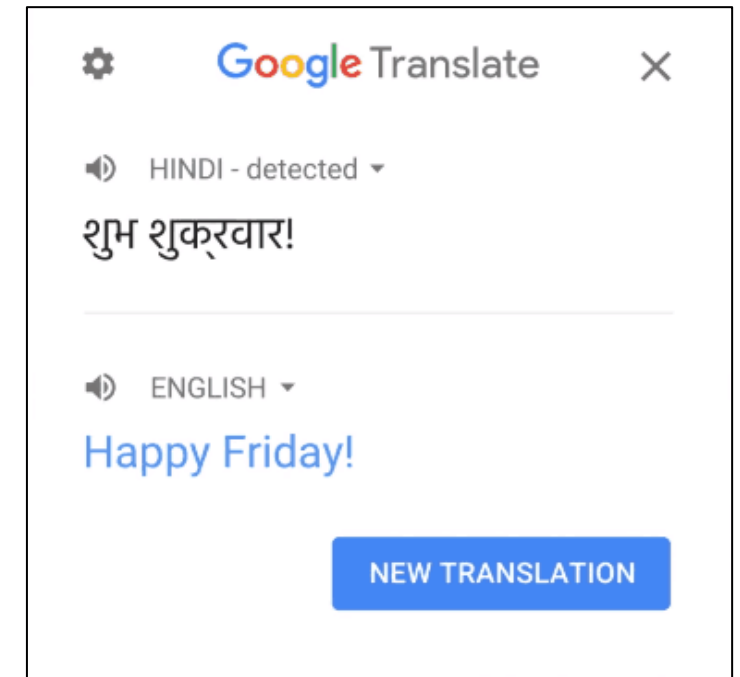
Chatbot



Search Engine



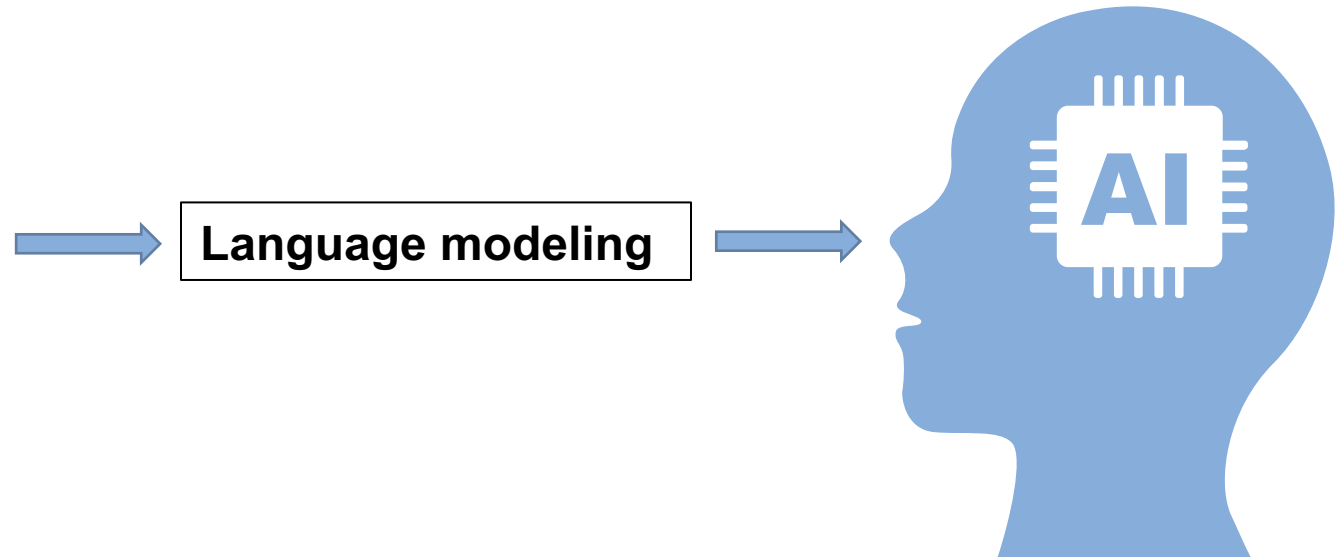
Translator



How do computers understand language?

A statistical language model is a probability distribution over sequences of words. The language model provides context to distinguish between words and phrases that sound similar. For example, in American English, the phrases "recognize speech" and "wreck a nice beach" sound similar, but mean different things.

\ufeff_____ .. . II . สัญญาเลขที่
TMV/CM212015-1575/BBB/ys Project :
FTTX Fund Code : Fund Center : สัญญา
ว่าจ้างสร้าง และ/หรือ ปรับปรุงข่ายสายสัญญาณโทรคมนาคม สัญญา
ฉบับนี้ทำขึ้นเมื่อวันที่ 01 .ค. 2558 ณบริษัท โทร มูล จำกัดระหว่าง (1)
บริษัท โทร มูล จำกัด โดย นายสุกษัย เจียรนนท์ และศาสตราจารย์พิเศษ
อสิก อิศวานนท์ กรรมการผู้มีอำนาจ สำนักงานตั้งอยู่เลขที่ 18 อาคารทรู
ทาวเวอร์ ถนนรัชดาภิเษก แขวงห้วยขวาง เขตห้วยขวาง กรุงเทพมหานคร
ซึ่งต่อไปนี้จะเรียกว่า "ผู้ว่าจ้าง" ฝ่ายหนึ่ง กับ อพ (2) **ห้างหุ้นส่วน**
จำกัด บรอดแบนด์ โดยนายปัญญาวุฒิ พลราช หุ้นส่วนผู้จัดการสำนักงาน
ตั้งอยู่เลขที่ 154 หมู่ ที่ 12 ตำบลกุดลาด อำเภอเมืองอุบลราชธานี จังหวัด
อุบลราชธานี ซึ่งต่อไปนี้จะเรียกว่า "ผู้รับจ้าง"อีกฝ่ายหนึ่ง โฆษก นาง
CA2-007 หน้า 19 Standard Form 2015



One-hot

One hot encodings are another way of representing words in numeric form. The length of the word vector is equal to the length of the vocabulary, and each observation is represented by a matrix with rows equal to the length of vocabulary and columns equal to the length of observation, with a value of 1 where the word of vocabulary is present in the observation and a value of zero where it is not.

One-Hot Word Representations

	The	cat	sat	on	the	mat.
<u>word</u>						
the	1	0	0	0	1	0
cat	0	1	0	0	0	0
on	0	0	0	1	0	0
⋮						
⋮						
⋮						
nunique_words						

Word Embedding

Word embedding is the collective name for a set of language modeling and feature learning techniques where words or phrases from the vocabulary are mapped to vectors of real numbers. The technique is primarily used with Neural Network Models.

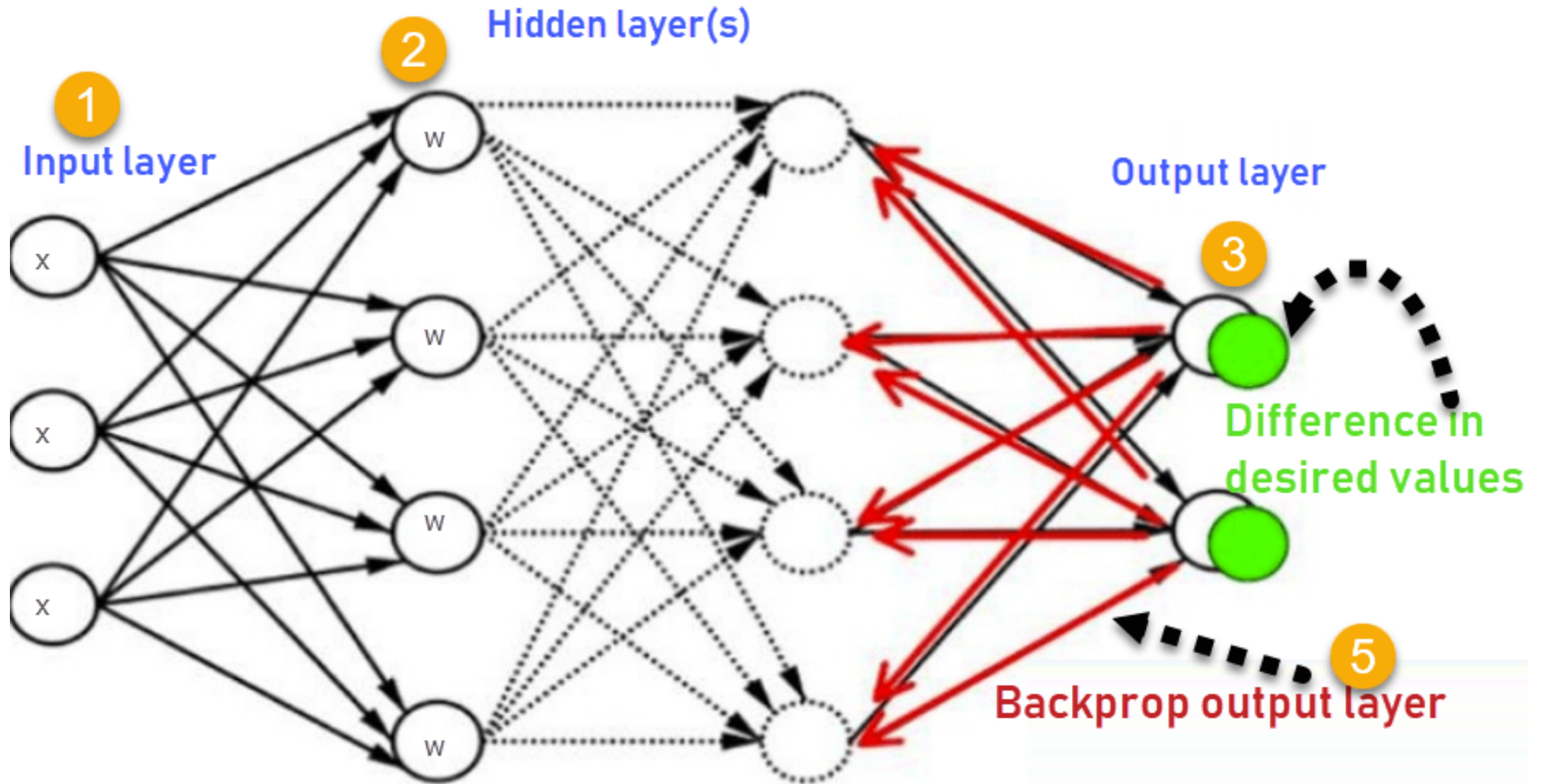
Featurized representation: word embedding

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	<u>0.93</u>	<u>0.95</u>	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
size	⋮	⋮				
cost						
alive						
verb						

I want a glass of orange _____.
I want a glass of apple_____.

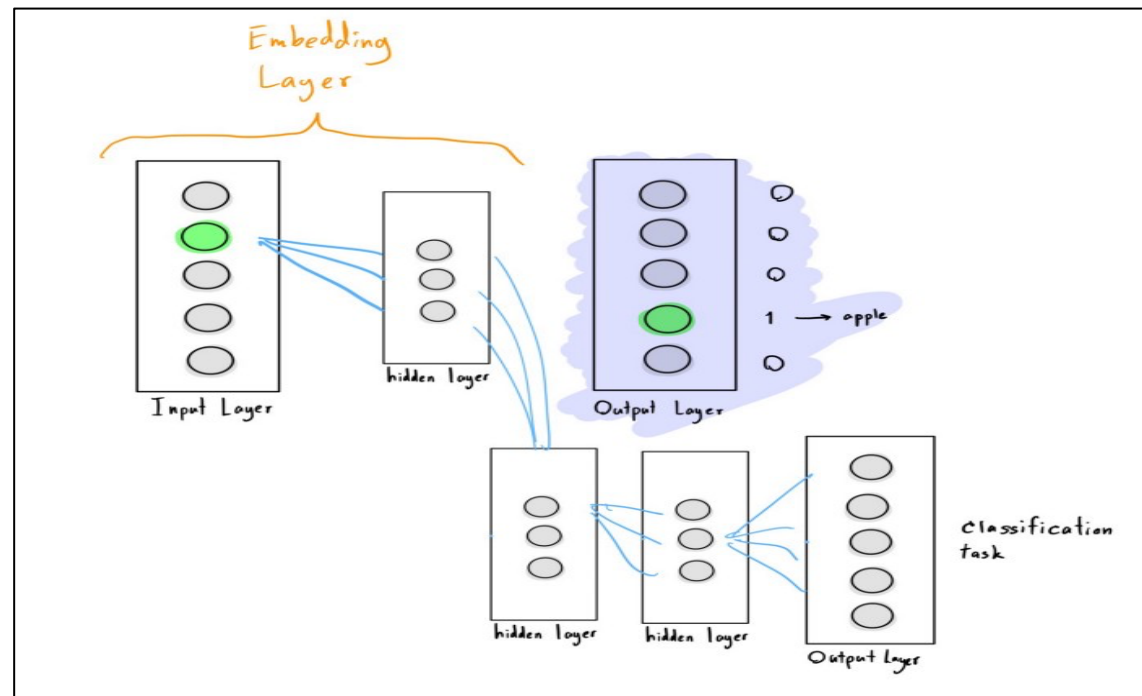
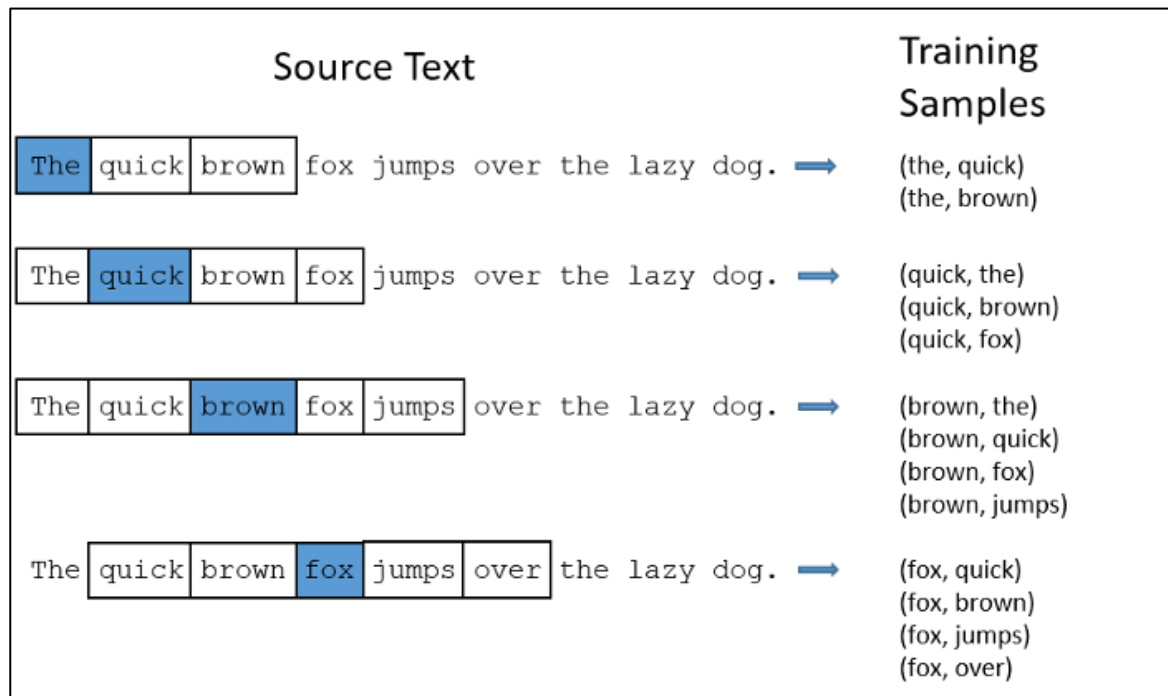
Andrew Ng

Neural Network Model



Word Embedding

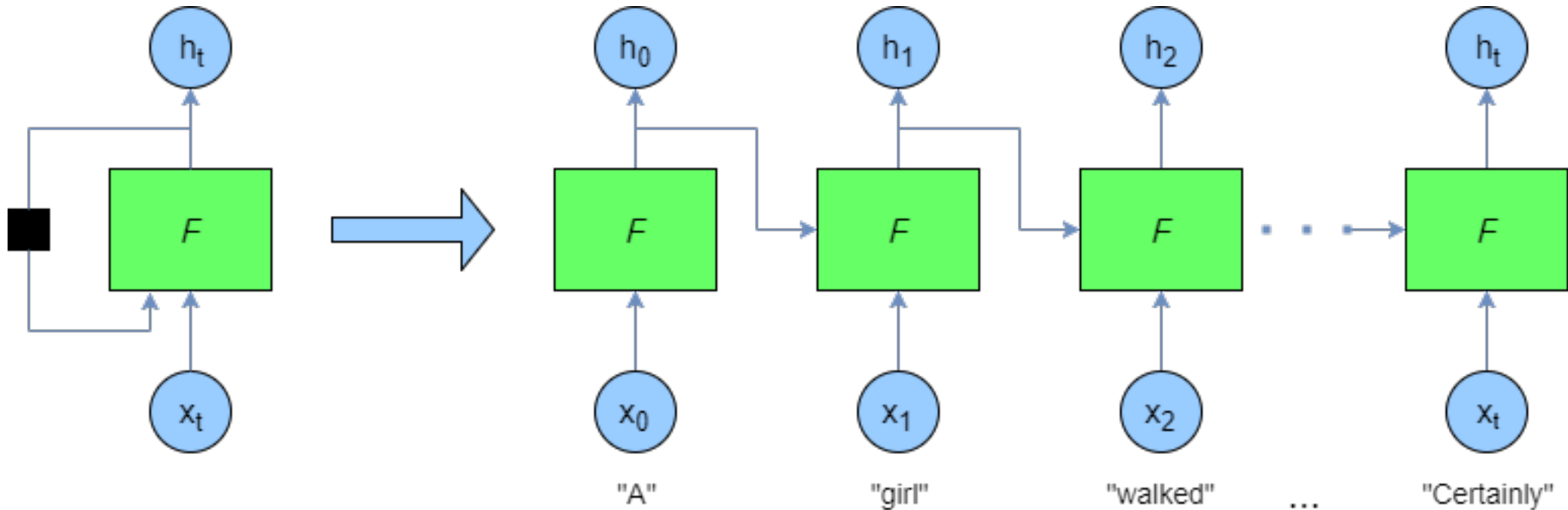
skip-gram



•Skip-Gram — ใช้ Context 1 คำเพื่อหา next word หลายคำ

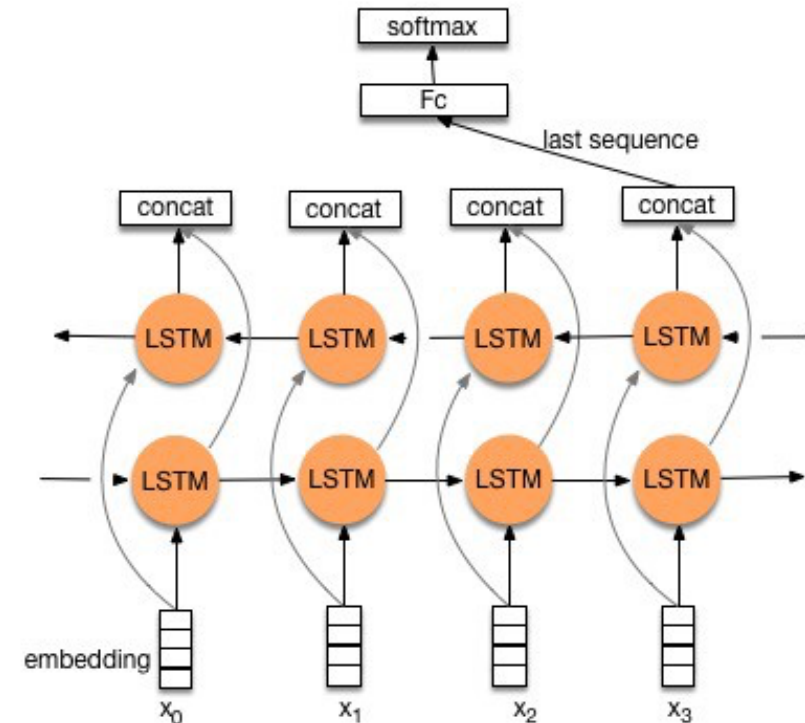
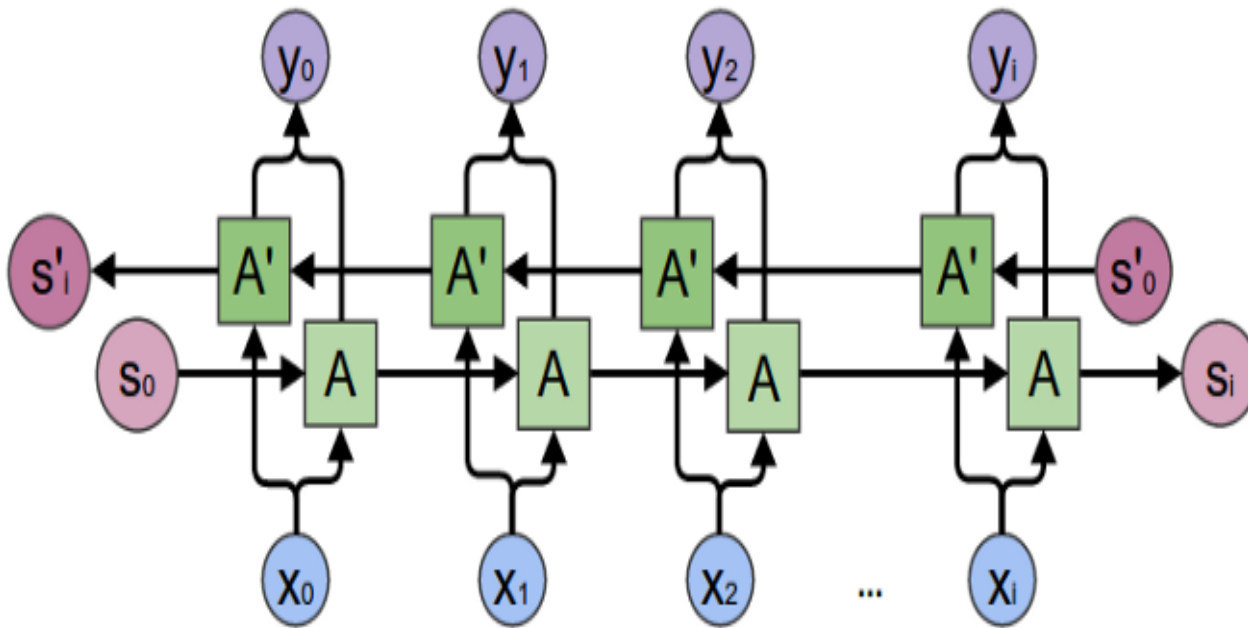
Recurrent Neural Networks (RNN)

Conceptually they differ from a standard neural network as the standard input in a RNN is a word instead of the entire sample. This gives the flexibility for the network to work with varying lengths of sentences. It also provides an additional advantage of sharing features learned across different positions of text which can not be obtained in a standard neural network.



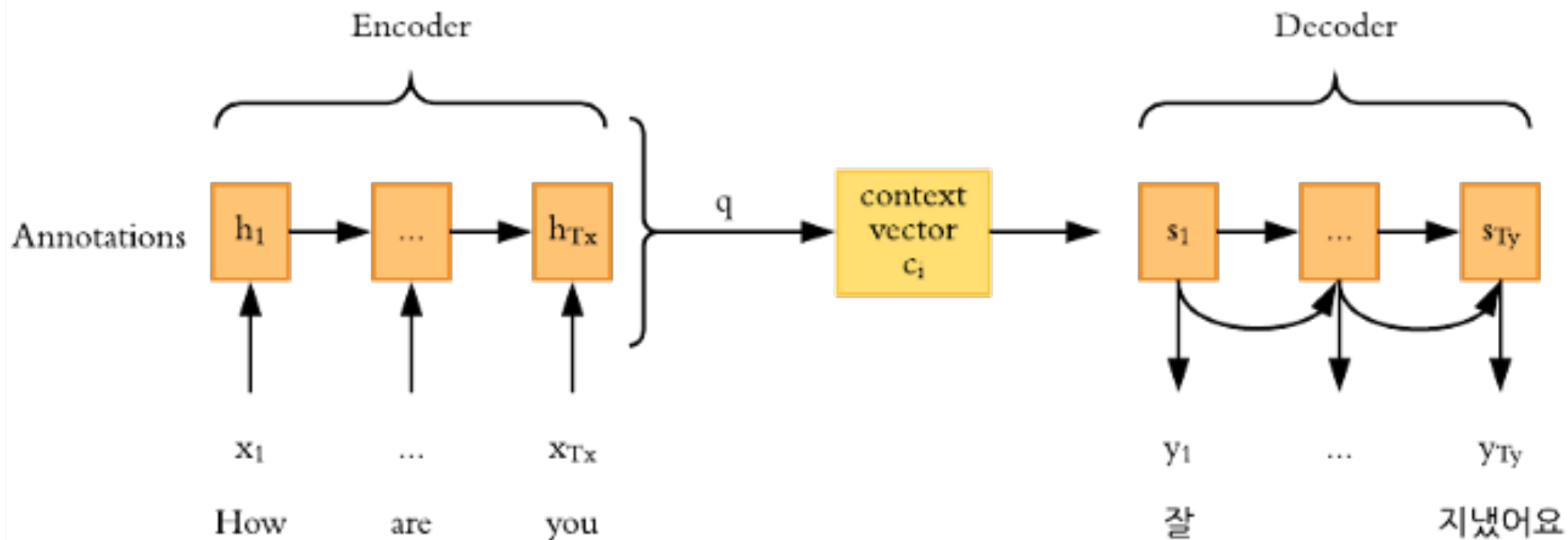
Bidirectional RNN

Bidirectional recurrent neural networks(RNN) are really just putting two independent RNNs together. The input sequence is fed in normal time order for one network, and in reverse time order for another. The outputs of the two networks are usually concatenated at each time step. This structure allows the networks to have both backward and forward information about the sequence at every time step.



Sequence-to-sequence model

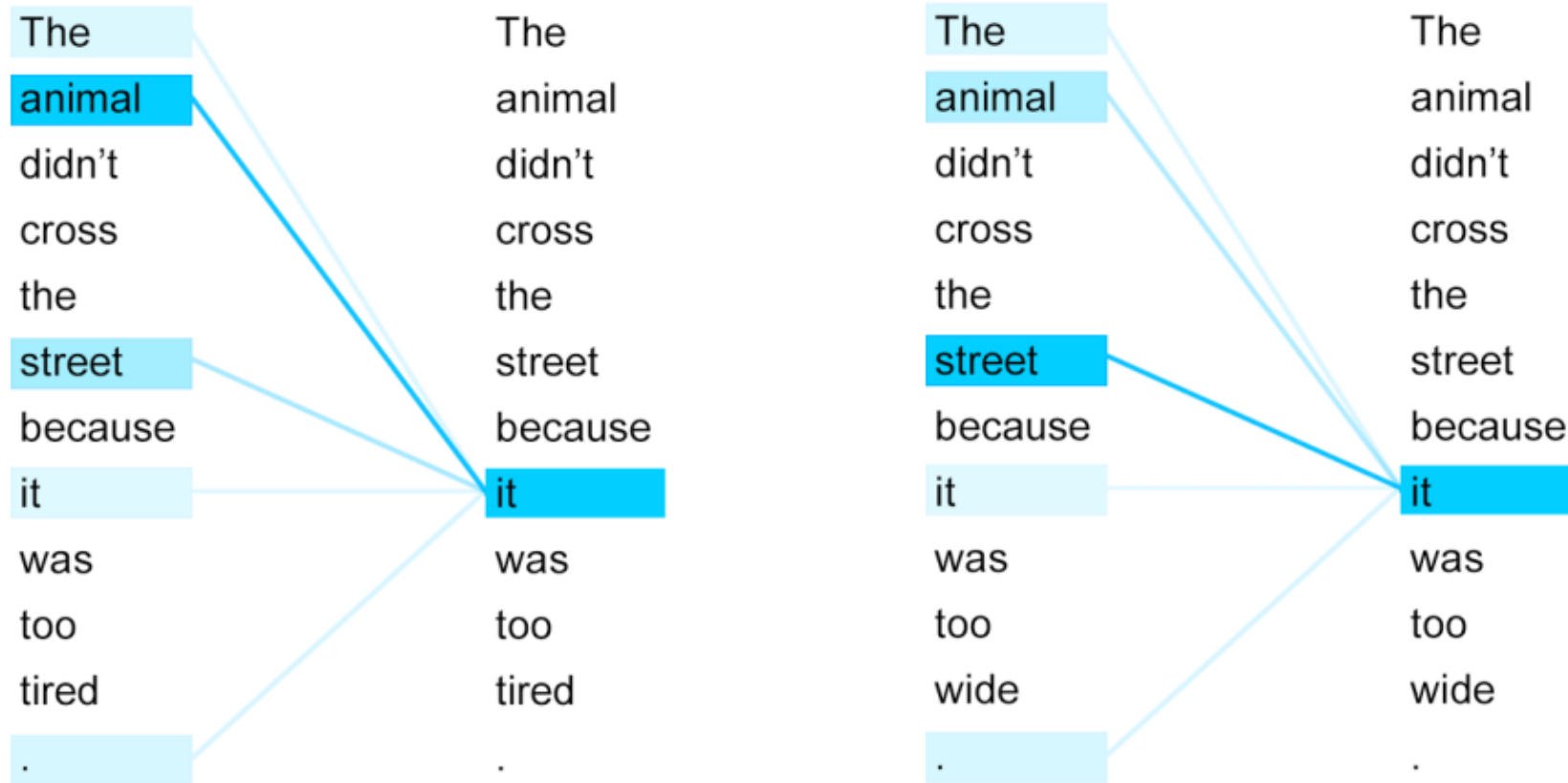
sequence-to-sequence model (seq2seq) หรือเรียกอีกชื่อว่า RNN Encoder-Decoder ซึ่งโมเดลนี้จะแบ่งเป็นสองส่วน เรียกว่า encoder กับ decoder โดยส่วน encoder จะรับ input เข้ามาทีละหน่วยผ่านทาง RNN และเก็บสะสม information ที่จำเป็นไว้ จากนั้นจะผ่าน information นี้ไปยังส่วน decoder ซึ่งก็จะเป็น RNN อีกตัวหนึ่งที่ให้ output ออกมาทีละหน่วย โดยดูจาก information ที่ได้รับมา และ output ตัวก่อนหน้า



โมเดล seq2seq จะมีปัญหาขอขวดเกิดขึ้น นั่นคือการส่ง information เป็นทอดๆ ตามสายยาวแบบนี้ อาจจะมี information ที่จำเป็นบางอย่างสูญหายไป ระหว่างทางได้ จะดีกว่าไหมถ้าเราให้กระบวนการสร้าง output สามารถโฟกัสไปที่ input ส่วนใดส่วนหนึ่งได้โดยตรง และนี่คือที่มาของ **Attention** นั่นเอง

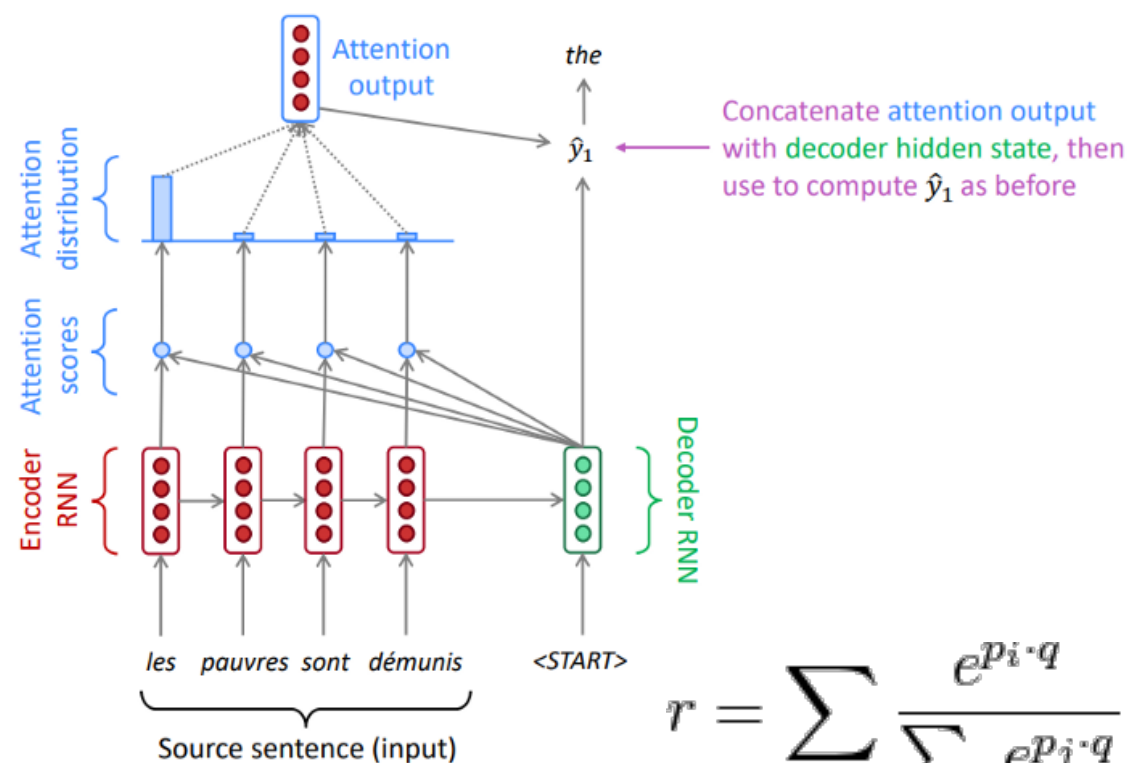
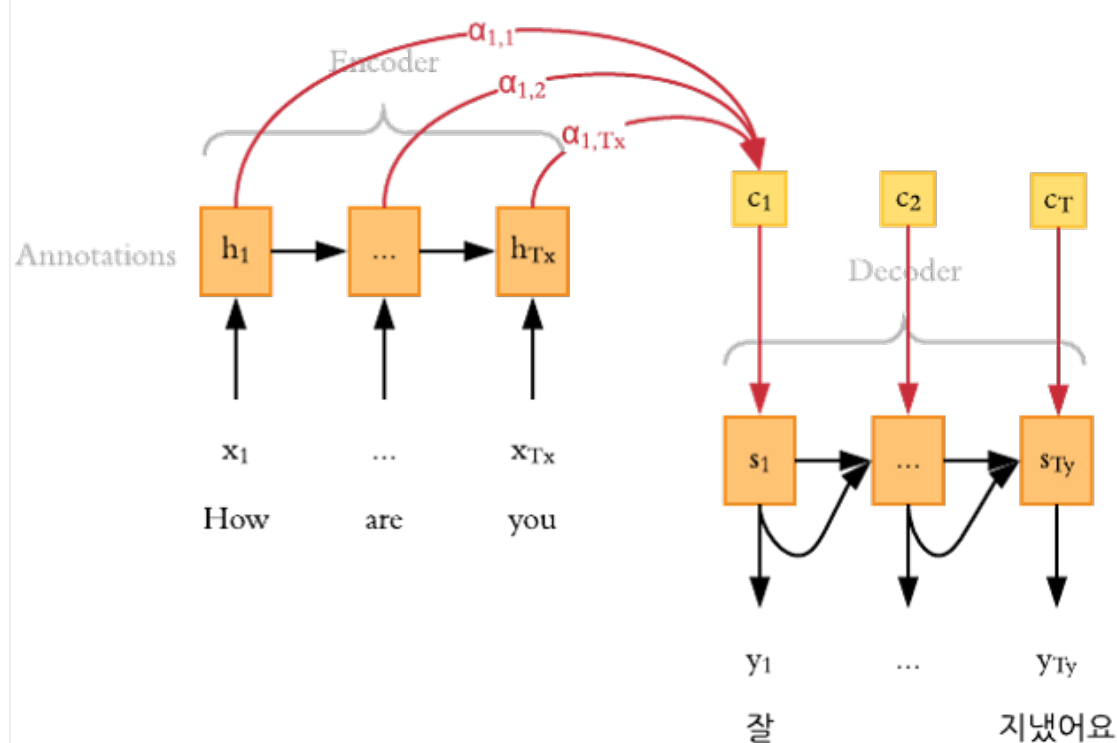
Transformer model : Attention Is All You Need

The **idea here is to learn a context vector** (say U), which gives us global level information on all the inputs and tells us about the most important information (this could be done by taking a **cosine similarity** of this context vector U the input hidden states from the fully connected layer. We do this for each input x_i and thus obtain a θ_i (attention weights).



Attention

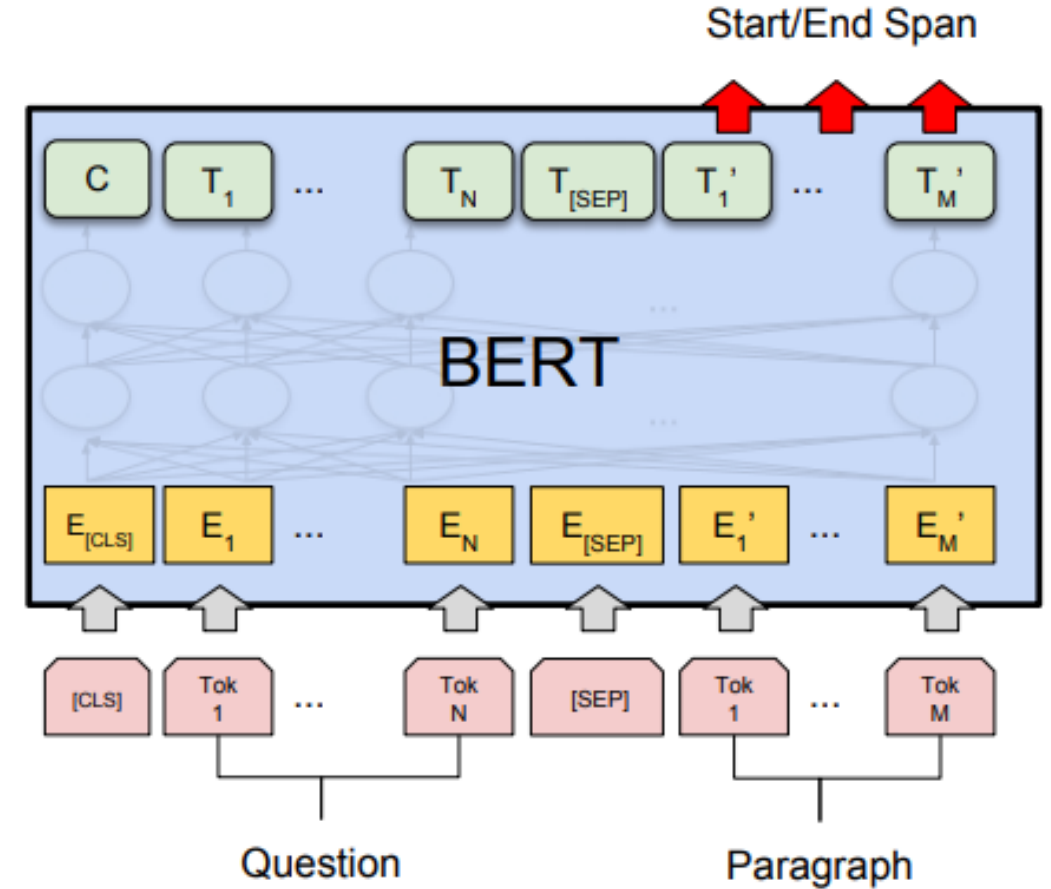
Attention สำหรับ seq2seq สามารถแสดงได้ดังรูปด้านล่างนี้ โดยเมื่อต้องการจะคำนวณ output ที่ตำแหน่งหนึ่ง ก็จะนำ vector ของ decoder (q) ณ ตำแหน่งนั้น มาใช้หา attention score กับ vector ของ encoder (p) ในทุกตำแหน่ง ซึ่งถ้า score ที่ encoder ตำแหน่งไหนสูง หมายความว่าเราจะให้ความสำคัญ หรือใส่ใจกับตำแหน่งนั้นมาก การคำนวณค่านี้ก็ได้หลายวิธี โดยวิธีที่ง่ายที่สุดก็คือการทำ dot product กันตรงๆระหว่าง p กับ q เลย ซึ่งหมายความว่าเราจะใส่ใจกับตำแหน่งที่มีค่า p ใกล้เคียงกับค่า q และเมื่อได้ค่า score ออกมาแล้ว ก็จะเอาเข้าฟังก์ชัน *softmax* เพื่อแปลงเป็นค่าความน่าจะเป็น



BERT

(Bidirectional Encoder Representations from Transformers)

BERT's key technical innovation is applying the bidirectional training of Transformer, a popular attention model, to language modelling. This is in contrast to previous efforts which looked at a text sequence either from left to right or combined left-to-right and right-to-left training. The paper's results show that a language model which is bidirectionally trained can have a deeper sense of language context and flow than single-direction language models.



(c) Question Answering Tasks:

[illegible]

OCR

TEXT

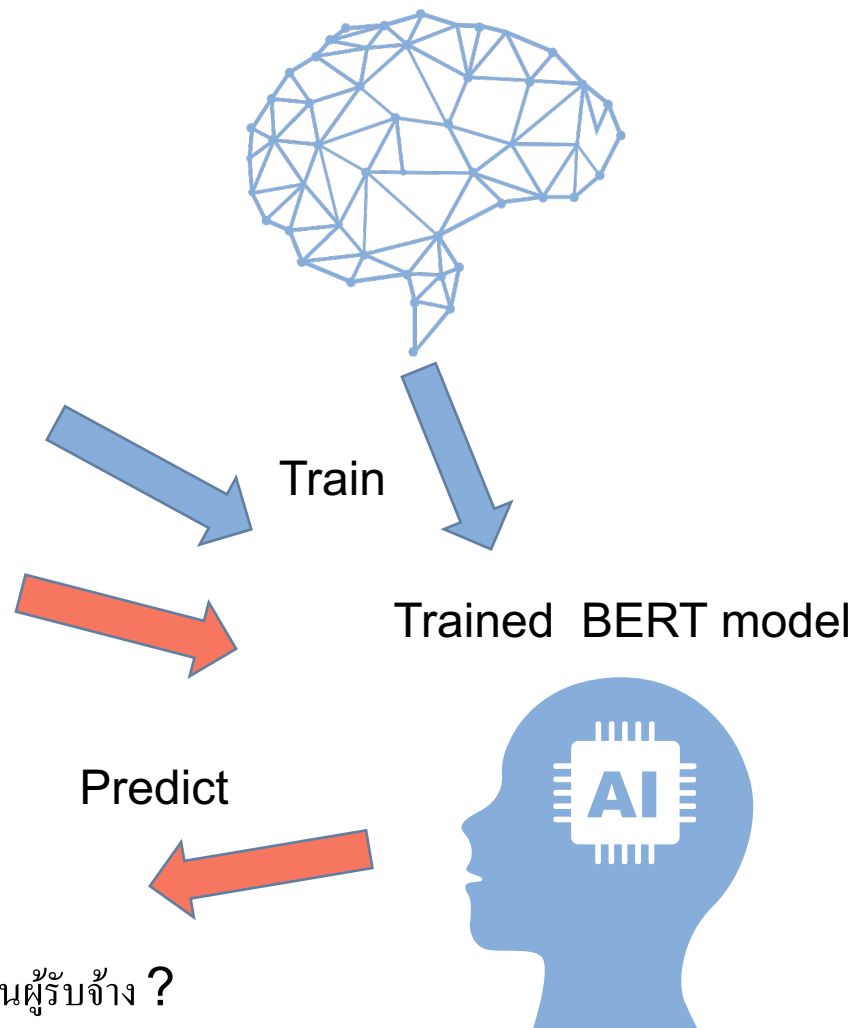
luffeff_____ ... II . สัญญาเลขที่
TMV/CM212015-1575/BBB/ys Project :
FTTX Fund Code : Fund Center : สัญญา
ว่าจ้างสร้าง และ/หรือ ปรับปรุงขยายสายสัญญาคมนาคม
สัญญาฉบับนี้ทำขึ้นเมื่อวันที่ ๐ 1 .ค. 2558 ณบริษัท โทร มูฟ จำกัด
ระหว่าง (1) บริษัท โทร มูฟ จำกัด โดย นายศุภชัย เจียรนวนนท์ และ
ศาสตราจารย์พิเศษ อสิก อ้วนานนท์ กรรมการผู้ชำนาญ
สำนักงานตั้งอยู่เลขที่ 18 อาคารทรูทาวเวอร์ ถนนรัชดาภิเษก แขวง
ห้วยขวาง เขตห้วยขวาง กรุงเทพมหานคร ซึ่งต่อไป นี้ จะเรียกว่า
"ผู้ว่าจ้าง" ฝ่ายหนึ่ง กับ อพ (2) **ห้างหุ้นส่วนจำกัด บรอดแบนด์**
โดยนายปัญญาวุฒิ พลราช หุ้นส่วนผู้จัดการสำนักงานตั้งอยู่เลขที่
154 หมู่ที่ 12 ตำบลกุดลาด อำเภอเมืองอุบลราชธานี จังหวัด
อุบลราชธานี ซึ่งต่อไปนี้จะเรียกว่า "ผู้รับจ้าง"อีกฝ่ายหนึ่ง ไซกา
นง CA2-007 หน้า 19 Standard Form 2015

Result

ห้างหุ้นส่วนจำกัด หรือ บริษัทใดเป็นผู้รับจ้าง ?

ห้ามหุ้นส่วนจำกัด บรอดแบนด์

BERT base model





Q & A