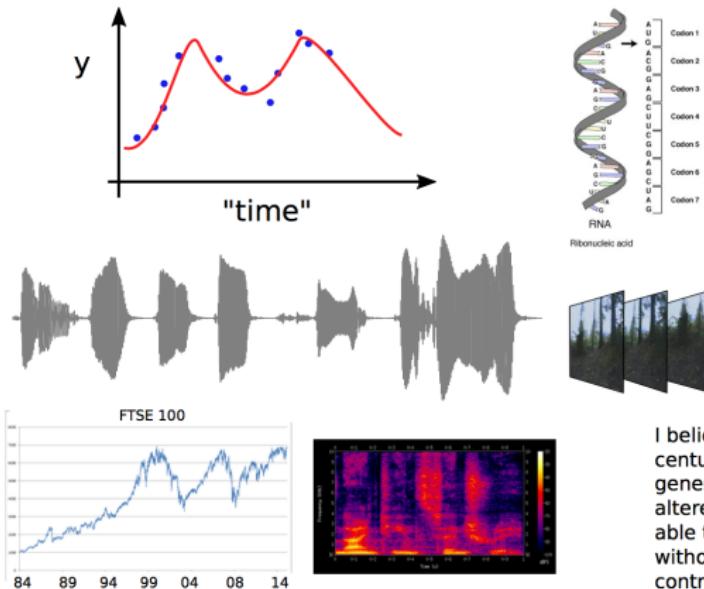


4F10: Deep Learning for Sequence Data

Mark Gales

Michaelmas 2020

Sequence Modelling (3F8)



Good King Wenceslas looked out,
On the Feast of Stephen;
When the snow lay round about;
Deep and crisp and even;
Brightly shone the moon that night;
Though the frost was cruel;
When a poor man came in sight,
Lything winter fuel.

I believe that at the end of the century the use of words and general educated opinion will have altered so much that one will be able to speak of machines thinking without expecting to be contradicted.

A. Turing

- The (input) data considered has been a **fixed size**
 - vector (iris data), or matrix (image data)

What about variable length data?

- Various sources result in sequences:
 - speech recognition (length of waveform)
 - language processing (length of word sequence)
 - DNA sequences
- Already seen various sequence models (3F8/4F10)
 - hidden Markov model: $p(\mathbf{x}_{1:T})$
 - Kalman filter/smoothes: $p(\mathbf{x}_{1:T})$
 - conditional random field: $P(\mathbf{y}_{1:T} | \mathbf{x}_{1:T})$

How can deep learning be applied for these tasks

Sequence Modelling Tasks

- Multiple forms of sequence modelling can be considered
 1. sequence input/output pairs: $\{\{x_1, y_1\}, \dots, \{x_T, y_T\}\}$
 2. input sequence to single output: $\{\{x_1, \dots, x_T\}, y\}$
 3. input sequence to output seq.: $\{\{x_1, \dots, x_T\}, \{y_1, \dots, y_L\}\}$
 4. single input to ouput sequence: $\{x, \{y_1, \dots, y_L\}\}$
 5. sequence distribution modelling: $p(x_1, \dots, x_T)$
- T can vary from training sample to training (test) sample
 - for sequence-to-sequence N can vary from sample to sample
- To simplify notation, sequences may be written as

$$x_{1:T} = \{x_1, \dots, x_T\}$$

Sequence Input/Output

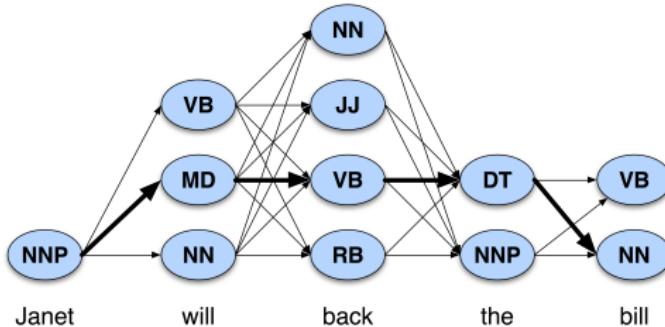
Pair Modelling

Part of Speech Tagging

- Given a word-sequence, $w_{1:L}$, extract PoS tag sequence $y_{1:L}$

Janet/NNP will/MD back/VB the/DT bill/NN

- individual words are **ambiguous** (from Jurafsky & Martin)

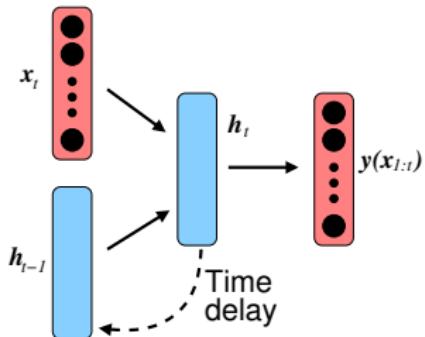


- Tag sequence (usually) unambiguous given word sequence

Recurrent Neural Networks [17, 16]

- Consider a causal sequence of observations $\mathbf{x}_{1:t} = \{\mathbf{x}_1, \dots, \mathbf{x}_t\}$

- Introduce recurrent units



$$\mathbf{h}_t = \mathbf{f}^h (\mathbf{W}_h^f \mathbf{x}_t + \mathbf{W}_h^r \mathbf{h}_{t-1} + \mathbf{b}_h)$$

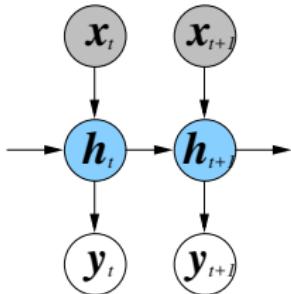
$$\mathbf{y}(\mathbf{x}_{1:t}) = \mathbf{f}^f (\mathbf{W}_y \mathbf{h}_t + \mathbf{b}_y)$$

- \mathbf{h}_t history vector at time t
- Two history weight matrices
 - \mathbf{W}_h^f forward, \mathbf{W}_h^r recursion
- Uses approximation to model history of observations

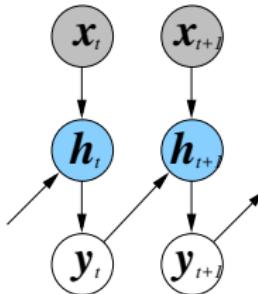
$$\mathcal{F}(\mathbf{x}_{1:t}) = \mathcal{F}(\mathbf{x}_t, \mathbf{x}_{1:t-1}) \approx \mathcal{F}(\mathbf{x}_t, \mathbf{h}_{t-1}) \approx \mathcal{F}(\mathbf{h}_t) = \mathbf{y}(\mathbf{x}_{1:t}) = \mathbf{y}_t$$

- network has (causal) memory encoded in history vector (\mathbf{h}_t)

Elman and Jordan Networks



Elman Network

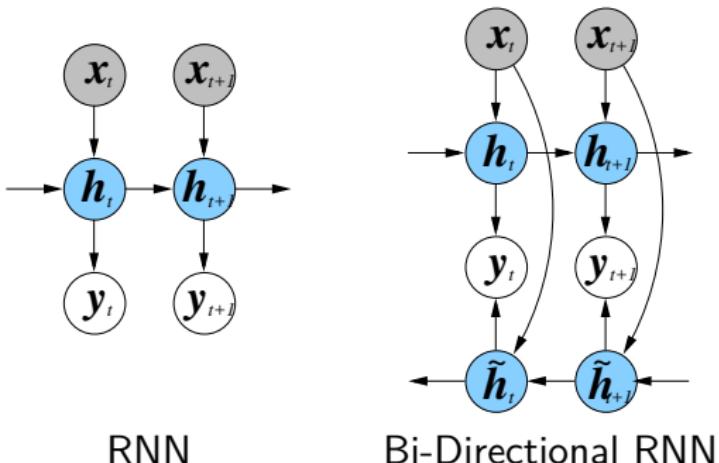


Jordan Network

- Previous slide sometimes referred to as an **Elman Network**
 - history vector, h_t “remembers” input sequence $\mathbf{x}_{1:t}$
- Alternative topology, **Jordan Network**, sometimes used
 - history vector, h_t “remembers” (unobserved) output seq. $\mathbf{y}_{1:t}$

$$\mathcal{F}(\mathbf{x}_{1:t}, \mathbf{y}_{1:t-1}) \approx \mathcal{F}(\mathbf{x}_t, \mathbf{y}_{1:t-1}) \approx \mathcal{F}(\mathbf{x}_t, \mathbf{h}_{t-1}) \approx \mathcal{F}(\mathbf{h}_t) = \mathbf{y}_t$$

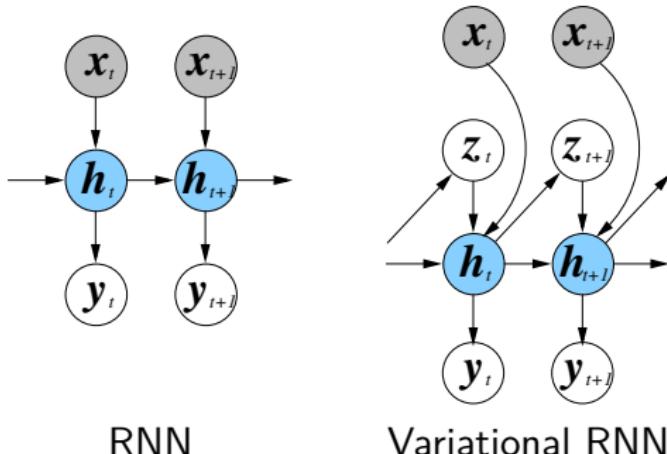
RNN Variant: Bi-Directional RNN [19]



- Bi-directional: use **complete observation sequence** - non-causal

$$\mathcal{F}_t(\mathbf{x}_{1:T}) = \mathcal{F}(\mathbf{x}_{1:t}, \mathbf{x}_{t:T}) \approx \mathcal{F}(\mathbf{h}_t, \tilde{\mathbf{h}}_t) = \mathbf{y}_t(\mathbf{x}_{1:T})$$

Latent Variable (Variational) RNN (reference) [6]



- Variational: introduce **latent** variable sequence $\mathbf{z}_{1:T}$

$$\begin{aligned} p(\mathbf{y}_t | \mathbf{x}_{1:t}) &\approx \int p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{z}_t, \mathbf{h}_{t-1}) p(\mathbf{z}_t | \mathbf{h}_{t-1}) d\mathbf{z}_t \\ &\approx \int p(\mathbf{y}_t | \mathbf{h}_t) p(\mathbf{z}_t | \mathbf{h}_{t-1}) d\mathbf{z}_t \end{aligned}$$

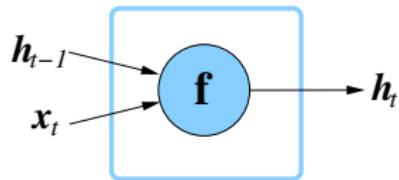
- \mathbf{z}_t a function of complete history (complicates training)

- A flexible extension to activation function is **gating**
 - standard form is ($\sigma()$ **sigmoid** activation function)

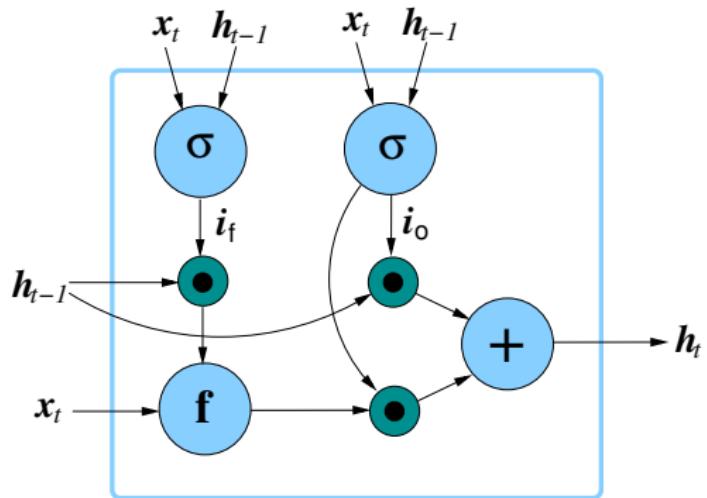
$$i = \sigma(\mathbf{W}^f \mathbf{x}_t + \mathbf{W}^r \mathbf{h}_{t-1} + \mathbf{b})$$

- vector acts a probabilistic gate on network values
- Gating can be applied at various levels
 - **features**: impact of input/output features on nodes
 - **time**: memory of the network
 - **layer**: influence of a layer's activation function

Gated Recurrent Unit [5]



Recurrent unit



Gated Recurrent Unit

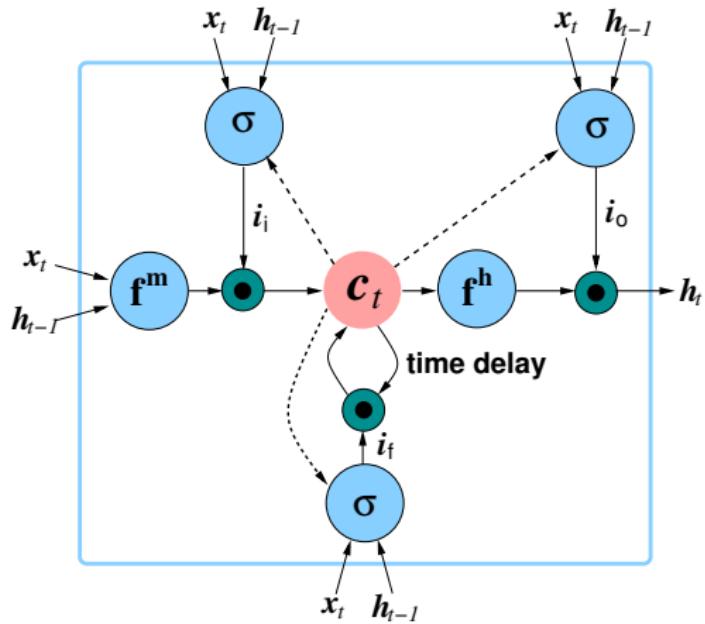
Gated Recurrent Unit [5]

- Gated Recurrent Unit (GRU) introduces
 - forget gate (i_f): gating over time
 - output gate (i_o): gating over features (and time)
- Relationships (standard configuration - there are variants)

$$\begin{aligned} i_f &= \sigma(\mathbf{W}_f^f \mathbf{x}_t + \mathbf{W}_f^r \mathbf{h}_{t-1} + \mathbf{b}_f) \\ i_o &= \sigma(\mathbf{W}_o^f \mathbf{x}_t + \mathbf{W}_o^r \mathbf{h}_{t-1} + \mathbf{b}_o) \\ \tilde{\mathbf{h}}_t &= \mathbf{f}(\mathbf{W}_h^f \mathbf{x}_t + \mathbf{W}_h^r (i_f \odot \mathbf{h}_{t-1}) + \mathbf{b}_h) \\ \mathbf{h}_t &= i_o \odot \mathbf{h}_{t-1} + (1 - i_o) \odot \tilde{\mathbf{h}}_t \end{aligned}$$

- \odot represents element-wise multiplication between vectors

Long-Short Term Memory Networks [11, 9]



Long-Short Term Memory Networks

- The operations can be written as (peephole config):
 - Forget gate (i_f), Input gate (i_i), Output gate (i_o)

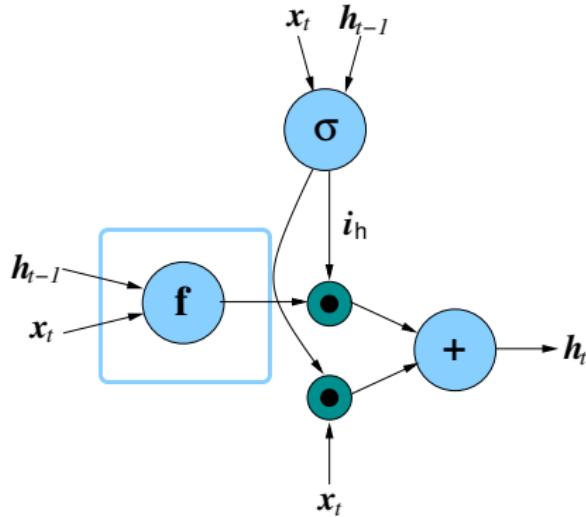
$$\begin{aligned} i_f &= \sigma(W_f^f x_t + W_f^r h_{t-1} + W_f^m c_{t-1} + b_f) \\ i_i &= \sigma(W_i^f x_t + W_i^r h_{t-1} + W_i^m c_{t-1} + b_i) \\ i_o &= \sigma(W_o^f x_t + W_o^r h_{t-1} + W_o^m c_t + b_o) \end{aligned}$$

- Memory Cell, history vector and gates are related by

$$\begin{aligned} c_t &= i_f \odot c_{t-1} + i_i \odot f^m (W_c^f x_t + W_c^r h_{t-1} + b_c) \\ h_t &= i_o \odot f^h (c_t) \end{aligned}$$

- more complicated than GRU (three gates, memory cell)
- memory cell weight matrices (W_f^m, W_i^m, W_o^m) diagonal
- can allow explicit analysis of individual cell elements

Highway Connections [20]



- Gate the output of the node (example from recurrent unit)
 - combine with output from previous layer (x_t)

$$i_h = \sigma(\mathbf{W}_1^f x_t + \mathbf{W}_1^r h_{t-1} + \mathbf{b}_1)$$

$$h_t = i_h \odot f(\mathbf{W}_h^f x_t + \mathbf{W}_h^r h_{t-1} + \mathbf{b}_h) + (1 - i_h) \odot x_t$$

Input Sequence to Target

Twitter Sentiment Analysis

- Given word sequence, $w_{1:L}$, extract sentiment {pos, neg}

Type in a word and we'll highlight the good and the bad

angry

Sentiment analysis for angry



Tweets about: angry

gabbs_loolive: Soooooooooooooo he pissed me off this morning nd now I'm just angry
Posted 17 seconds ago

bvaDonghaelover: sm * sh Indonesia mimic the style of super junior I was very angry @TEUKdom @special1004 @donghae861015 @KyuhyunBiased @GaemGuy
Posted 22 seconds ago

AmbahJambah: love watching @remow49 play angry birds while we should be working....
Posted 23 seconds ago

babyinashack: @Maisarahh Cuteboy? No. I'm angry.
Posted 25 seconds ago

hondanthon: People on the internet are wrong, it's making me angry, and I can't be arsed correcting them.
Posted 26 seconds ago

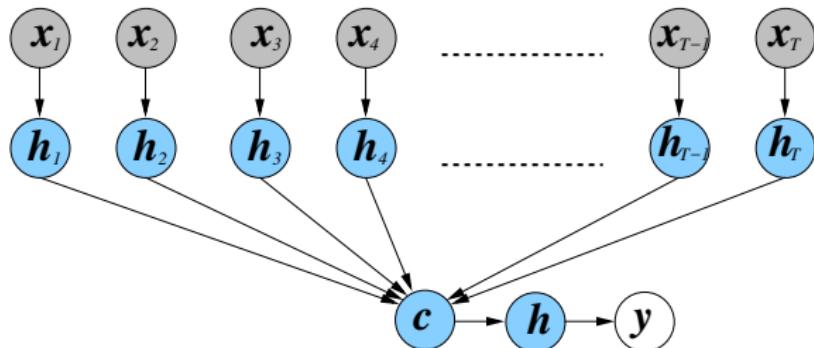
aymeeahmd: I won't get angry if you found someone better. I'm happy if you're happy. #nowthatdislove Cheeyyy! :)
Posted 35 seconds ago

twareffity: @Eamonn_Forge Fittingly, I read this tweet on a bus while a helmet played Angry Birds on his netbook beside me
Posted 40 seconds ago

The results for this query are:

- “Twitter Sentiment” for [angry](#)
 - tweets about [angry birds](#) game
 - sentiment accuracy reasonable
- Ambiguous for individual words

Averaging



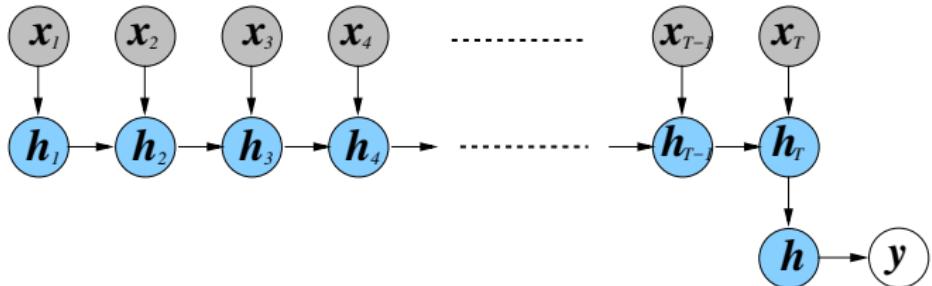
- Consider regression/classification of the form

$$\mathbf{y} = \mathcal{F}(\mathbf{x}_1, \dots, \mathbf{x}_T)$$

- need to map T (variable) length sequence to a fixed length
- Network above is one simple form

$$\mathbf{y} = \mathcal{F}(\mathbf{h}); \quad \mathbf{h} = \mathcal{F}(\mathbf{c}); \quad \mathbf{c} = \frac{1}{T} \sum_{t=1}^T \mathbf{h}_t; \quad \mathbf{h}_t = \mathcal{F}(\mathbf{x}_t)$$

RNN Encoding (Sequence Embedding)



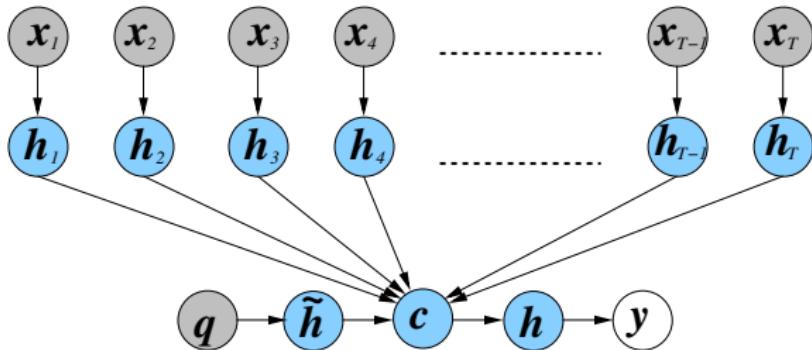
- Alternative use final history vector as information

$$\mathbf{h} = \mathcal{F}(\mathbf{h}_T); \quad \mathbf{h}_t = \mathcal{F}(\mathbf{x}_t, \mathbf{h}_{t-1})$$

- Handles variable length, but biased to later inputs
 - can use **bi-directional** information

$$\mathbf{h} = \mathcal{F}(\mathbf{h}_T, \tilde{\mathbf{h}}_1); \quad \mathbf{h}_t = \mathcal{F}(\mathbf{x}_t, \mathbf{h}_{t-1}); \quad \tilde{\mathbf{h}}_t = \mathcal{F}(\mathbf{x}_t, \tilde{\mathbf{h}}_{t+1})$$

Attention Mechanism



- Extend average to be a function of some **query**, q ,

$$\mathbf{c} = \sum_{t=1}^T \alpha_t \mathbf{h}_t; \quad \alpha_t = \frac{\exp(e_t)}{\sum_{i=1}^T \exp(e_i)}; \quad e_t = \mathcal{F}(\tilde{\mathbf{h}}, \mathbf{h}_t); \quad \tilde{\mathbf{h}} = \mathcal{F}(\mathbf{q})$$

- yields a **probability mass function** over the sequence
- important to get the appropriate form of q

Form of Attention Mechanism

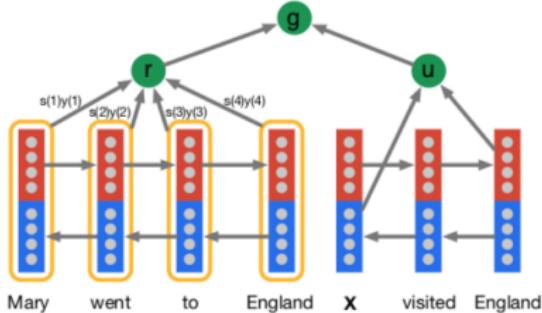
- At the center of the attention mechanism is $e_t = \mathcal{F}(\tilde{\mathbf{h}}, \mathbf{h}_t)$
 - relates (embedded) query \mathbf{q} to (embedded) observation \mathbf{x}_t
 - how relevant is that observation to the query
- Two standard forms of attention are:
 - dot-product attention (parameters \mathbf{W}_{xq})

$$e_t = \mathbf{h}_t^T \mathbf{W}_{xq} \tilde{\mathbf{h}}$$

- additive attention (parameters $\mathbf{w}, \mathbf{W}_x, \mathbf{W}_q$)

$$e_t = \mathbf{w}^T \tanh (\mathbf{W}_x \mathbf{h}_t + \mathbf{W}_q \tilde{\mathbf{h}})$$

Attentive Reader [10]



- Example question-answering network combines approaches
 - question encoding - forward, $\mathbf{u}_{|q|}$, and backwards, $\tilde{\mathbf{u}}_1$
 - document encoding - bi-directional encoding per word, $\mathbf{y}_t, \tilde{\mathbf{y}}_t$ followed by attention mechanism

$$\mathbf{r} = \sum_{t=1}^{|d|} s_t \begin{bmatrix} \mathbf{y}_t \\ \tilde{\mathbf{y}}_t \end{bmatrix}; \quad s_t = \frac{1}{Z} \exp \left(\mathbf{w}^\top \tanh \left(\mathbf{W}_{\text{ym}} \begin{bmatrix} \mathbf{y}_t \\ \tilde{\mathbf{y}}_t \end{bmatrix} + \mathbf{W}_{\text{um}} \begin{bmatrix} \mathbf{u}_{|q|} \\ \tilde{\mathbf{u}}_1 \end{bmatrix} \right) \right)$$

Attentive Reader

by ent423 ,ent261 correspondent updated 9:49 pm et ,thu march 19 ,2015 (ent261) a ent114 was killed in a parachute accident in ent45 ,ent85 ,near ent312 ,a ent119 official told ent261 on wednesday .he was identified thursday as special warfare operator 3rd class ent23 ,29 ,of ent187 ,ent265 .“ ent23 distinguished himself consistently throughout his career .he was the epitome of the quiet professional in all facets of his life ,and he leaves an inspiring legacy of natural tenacity and focused

...

ent119 identifies deceased sailor as X ,who leaves behind a wife

by ent270 ,ent223 updated 9:35 am et ,mon march 2 ,2015 (ent223) ent63 went familial for fall at its fashion show in ent231 on sunday ,dedicating its collection to `` mamma '' with nary a pair of `` mom jeans '' in sight .ent164 and ent21 ,who are behind the ent196 brand ,sent models down the runway in decidedly feminine dresses and skirts adorned with roses ,lace and even embroidered doodles by the designers ' own nieces and nephews .many of the looks featured saccharine needlework phrases like `` i love you ,

...

X dedicated their fall fashion show to moms

- Attention mechanisms can be **interpretable**
 - can examine which word in document is most relevant

Sequence to Sequence

Machine Translation



Marquer les rafales

Les rafales de marque est un lecteur dans la technologie de l'information dans le [laboratoire d'intelligence de machine](#) (autrefois le groupe de vision et de robotique de la parole (SVR)) et un camarade de l'[université d'Emmanuel](#). Il est un membre du [groupe de recherche de la parole](#) ainsi que les [jeunes de Steve](#) de membres de personnel de corps enseignant, la [réfexion boîte](#) et la [facture Byrne de Phil](#).

Une brève [biographie](#) est accessible en ligne.

[\[Recherche\]](#) | [projets](#) | [publications](#) | [étudiants](#) | [enseignant](#) | [contact](#)

Intérêts de recherches

- Reconnaissance de la parole continue de grand vocabulaire
- Reconnaissance de la parole robuste
- Adaptation d'orateur
- Étude de machine (en particulier choix modèle et méthodes grain-basées)
- Identification et vérification d'orateur

Une brève introduction à la [reconnaissance de la parole](#) est accessible en ligne.
[dessus](#)

Projets de recherche

Projets en cours :

- Bruit ASR robuste (Europe Ltd de recherches de Toshiba placée)
- Traitement averti d'environnement rapide et robuste (Europe Ltd de recherches de Toshiba placée)
 - [Position d'associé de recherches disponible](#)
- [AGILE](#) (projet placé par [GALE de DARPA](#))
- [Version 3 de HTK - HTK_V3.4](#) et [exemples](#) sont disponibles.

Projets récemment réalisés :

- [CoreTEx](#) (améliorant la technologie de reconnaissance de la parole de noyau)
- [Transcription audio riche de HTK](#)(Projet placé par OREILLES de DARPA) - [pages Web locaux](#)



- **Google Translate (2009)**
 - applied to my web-page
 - target language French
- Ambiguous for individual words
- Interesting output for people
 - [Mark Gales](#)
 - [Bill Byrne](#)
 - [Phil Woodland](#)
 - [Steve Young](#)

Encoder-Decoder Sequence Models

- Train a discriminative model from
 - $\mathbf{x}_{1:L} = \{\mathbf{x}_1, \dots, \mathbf{x}_L\}$: L -length input sequence
 - $\mathbf{y}_{1:K} = \{\mathbf{y}_1, \dots, \mathbf{y}_K\}$: K -length output

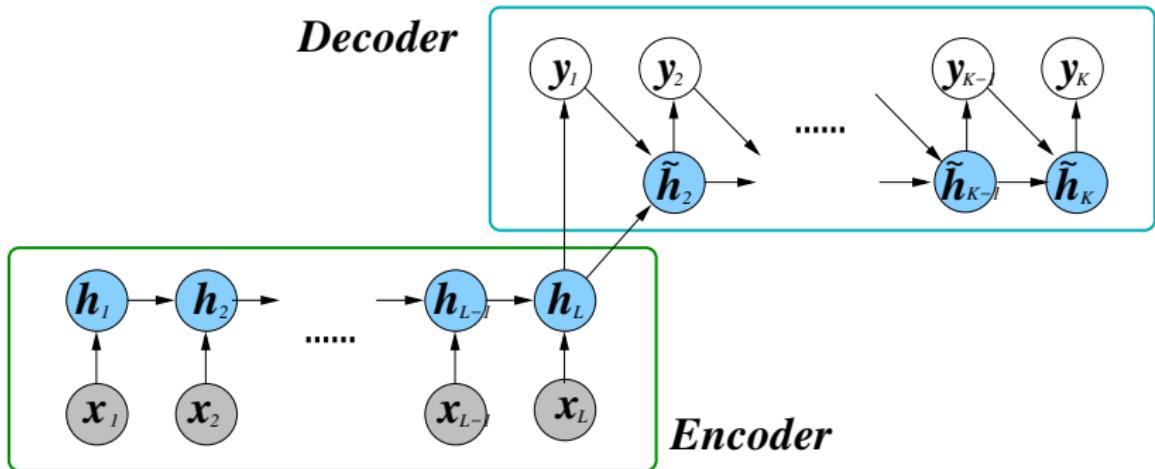
$$\begin{aligned} p(\mathbf{y}_{1:K} | \mathbf{x}_{1:L}) &= \prod_{i=1}^K p(\mathbf{y}_i | \mathbf{y}_{1:i-1}, \mathbf{x}_{1:L}) \\ &\approx \prod_{i=1}^L p(\mathbf{y}_i | \mathbf{y}_{i-1}, \tilde{\mathbf{h}}_{i-1}, \mathbf{c}) \end{aligned}$$

- need to map $\mathbf{x}_{1:L}$ to a fixed-length vector

$$\mathbf{c} = \phi(\mathbf{x}_{1:L})$$

- \mathbf{c} is a fixed length vector - like a **sequence kernel**

RNN Encoder-Decoder Model [8, 14]

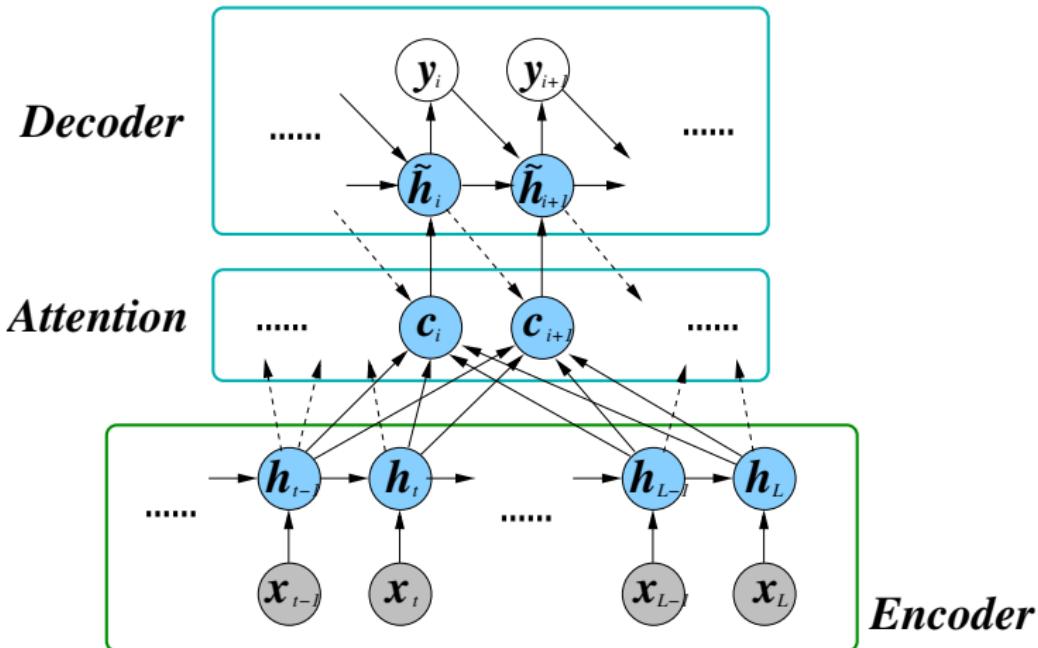


- One form is to use **hidden unit** from acoustic RNN/LSTM

$$\mathbf{c} = \phi(\mathbf{x}_{1:L}) = \mathbf{h}_L$$

- dependence on context is global via \mathbf{c} - possibly limiting

Attention-Based Models [4, 3, 14]



Attention-Based Models

- Introduce attention layer to system
 - introduce dependence on locality i

$$p(\mathbf{y}_{1:K} | \mathbf{x}_{1:L}) \approx \prod_{i=1}^K p(\mathbf{y}_i | \mathbf{y}_{i-1}, \tilde{\mathbf{h}}_{i-1}, \mathbf{c}_i) \approx \prod_{i=1}^K p(\mathbf{y}_i | \tilde{\mathbf{h}}_i)$$

where

$$\mathbf{c}_i = \sum_{\tau=1}^L \alpha_{i\tau} \mathbf{h}_\tau; \quad \alpha_{i\tau} = \frac{\exp(e_{i\tau})}{\sum_{j=1}^L \exp(e_{ij})}, \quad e_{i\tau} = \mathcal{F}(\tilde{\mathbf{h}}_{i-1}, \mathbf{h}_\tau)$$

- $e_{i\tau}$ how well position τ in input predicts position i in output
- \mathbf{h}_τ is representation (RNN) for the input at position τ
- Uses output history $\mathbf{y}_{1:i-1}$ as query for i^{th} word

Machine Translation

- Given word sequence, $w_{1:L}$, translate to French $\hat{y}_{1:K}$



Marquer les rafales

Les rafales de marque est un lecteur dans la technologie de l'information dans le [laboratoire d'intelligence de machine](#) (autrefois le groupe de vision et de robotique de la parole (SVR)) et un camarade de l'[université d'Emmanuel](#). Il est un membre du [groupe de recherche de la parole](#) ainsi que les [jeunes de Steve](#) de membres de personnel de corps enseignant, la [réfexion boisée](#) et la [facture Byrne de Phil](#).

Une brève [biographie](#) est accessible en ligne.

[[Recherche](#) | [projets](#) | [publications](#) | [étudiants](#) | [enseignant](#) | [contact](#)]

Intérêts de recherches

- [Reconnaissance de la parole continue de grand vocabulaire](#)
- [Reconnaissance de la parole robuste](#)
- Adaptation d'orateur
- Étude de machine (en particulier choix modèle et méthodes grain-basées)
- Identification et vérification d'orateur

Une brève introduction à la [reconnaissance de la parole](#) est accessible en ligne.
[dessus](#)

Projets de recherche

Projets en cours :

- [Bruit ASR robuste \(Europe Ltd de recherches de Toshiba placée\)](#)
- [Traitement averti d'environnement rapide et robuste \(Europe Ltd de recherches de Toshiba placée\)](#)
 - Position d'associé de recherches disponible
- [AGILE](#) (projet placé par [GALE](#) de DARPA)
- [Version 3 de HTK - HTK_V3.4](#) et [exemples](#) sont disponibles.

Projets récemment réalisés :



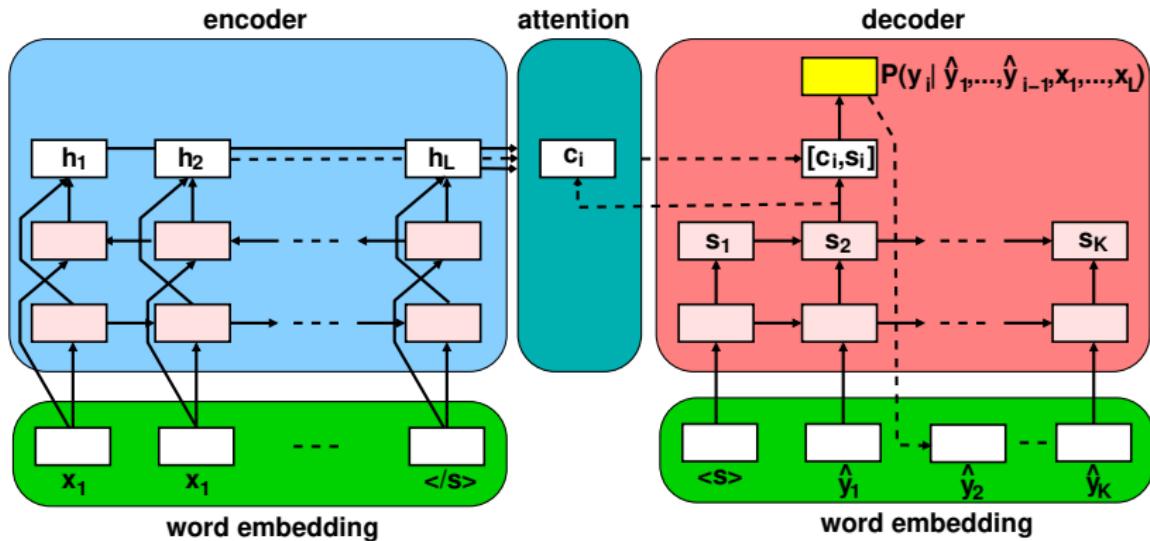
Google Translate (2009)

- applied to my web-page
- target language French

- Ambiguous for individual words
- Interesting output for people

- [Mark Gales](#)
- [Bill Byrne](#)
- [Phil Woodland](#)
- [Steve Young](#)

(Simplified) Neural Machine Translation



- Encoder-decoder framework with attention
 - word-embedding maps word (cat) to vector ($\text{emb}(\text{cat})$)
- Current state-of-the-art more complicated see MT papers

Neural Machine Translation - Distribution

- Overall probability distribution over all outputs based on

$$\hat{\mathbf{y}}_{1:K} \sim P(\mathbf{y}_{1:K} | \mathbf{w}_{1:L}) = P(\mathbf{y}_{1:K} | \mathbf{x}_{1:L})$$

- $\mathbf{x}_{1:L}$ sequence of word embeddings for word sequence $\mathbf{w}_{1:L}$
- generate translated word sequence $\hat{\mathbf{y}}_{1:K}$ from this distribution (swapped to $P()$ as the set of possible translations is discrete)

- Similar to language modelling use conditional distributions

$$P(\mathbf{y}_{1:K} | \mathbf{x}_{1:L}) = \prod_{i=1}^K P(y_i | \mathbf{y}_{1:i-1}, \mathbf{x}_{1:L})$$

- need: a representation of: input embedding sequence $\mathbf{x}_{1:L}$; (conditional) representation of the previous output words $\mathbf{y}_{1:i-1}$

Neural Machine Translation - Attention

- Possible translation generated in an **autoregressive** fashion

$$\hat{y}_i \sim P(y_i | \hat{\mathbf{y}}_{1:i-1}, \mathbf{x}_{1:L}) \approx P(y_i | \mathbf{s}_i, \mathbf{c}_i)$$

where

$$\hat{y}_i \sim \mathcal{F}_{\text{softmax}}(\mathbf{s}_i, \mathbf{c}_i); \quad \mathbf{s}_i = \mathcal{F}_{\text{rnn}}(\mathbf{s}_{i-1}, \tilde{\mathbf{h}}_i); \quad \tilde{\mathbf{h}}_i = \mathcal{F}_{\text{rnn}}(\tilde{\mathbf{h}}_{i-1}, \text{emb}(\hat{y}_{i-1}))$$

$$\mathbf{c}_i = \mathcal{F}_{\text{att}}(\mathbf{s}_i, \mathbf{h}_{1:L}) = \sum_{\tau=1}^L \alpha_{i\tau} \mathbf{h}_\tau; \quad \alpha_{i\tau} = \frac{\exp(e_{i\tau})}{\sum_{j=1}^L \exp(e_{ij})}, \quad e_{i\tau} = \mathcal{F}(\mathbf{s}_i, \mathbf{h}_\tau)$$

- $e_{i\tau}$ how well position τ in input predicts position i in output
- \mathbf{h}_τ is representation (bi-RNN) for the input at position τ
- At inference time need to find most probable translations
 - simplest approach is **greedy** search

Greedy Inference for NMT

- The decoder part of the model related to Jordan network
 - history vector encodes outputs and context histories

$$\hat{y}_i = \arg \max_{\omega} \{P(\omega | \mathbf{s}_i, \mathbf{c}_i)\}$$

- where $\omega \in \mathcal{V}$, \mathcal{V} is the NMT vocabulary
- Sequence generation process:
 - generate word-embeddings for source sentence $\mathbf{x}_{1:L}$
 - generate encoding of source sentence $\mathbf{h}_{1:L} = \{\mathbf{h}_1, \dots, \mathbf{h}_L\}$
 - initialise decoding: $\hat{y}_0 = <\text{s}>$, \mathbf{s}_0 , $\tilde{\mathbf{h}}_0$, set $i = 0$
 - DECODE: until $\hat{y}_i = </\text{s}>$ do
 - $i = i + 1$
 - compute new “state” vector \mathbf{s}_i
 - compute “context” information $\mathbf{c}_i = \mathcal{F}_{\text{att}}(\mathbf{s}_i, \mathbf{h}_{1:L})$
 - obtain distribution over target language words $P(y_i | \mathbf{s}_i, \mathbf{c}_i)$
 - sample from distribution \hat{y}_i , append to sentence $\rightarrow \hat{y}_{1:i}$

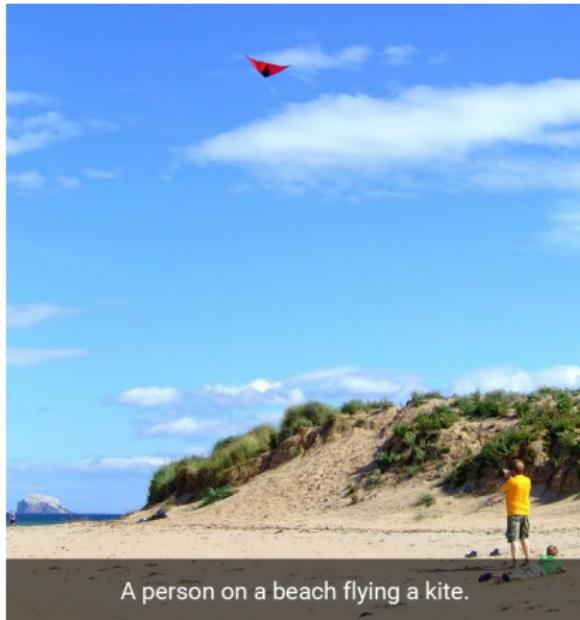
Teacher-Forcing (reference) [1, 12]

- One issue is that errors can occur in the history $\hat{y}_{0:i-1}$
 - system usually trained on correct data (bitext for MT)
this is sometimes referred to as **teacher-forcing**
- **But** at inference (decoding) time system can go “off-track”
 - would like to make system robust to model history errors
- Approaches have been developed to address this problem
 - **scheduled sampling**: mix “correct” history with model history
 - **professor forcing**: related to adversarial training

Single Input to Sequence

Image Captioning

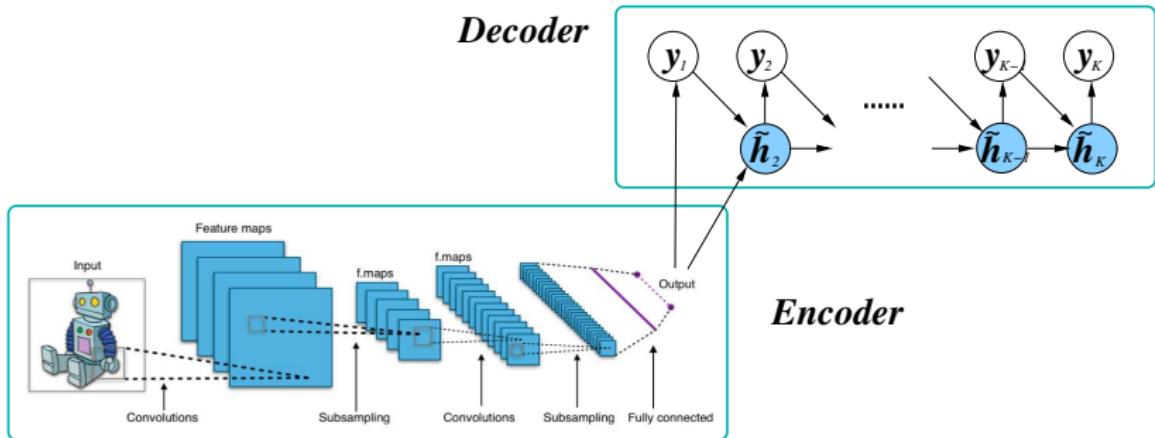
- Given image, x , generate caption $w_{1:K}$



A person on a beach flying a kite.

- Images highly complex
 - what's most relevant?
 - how to generate a caption?

Image Captioning

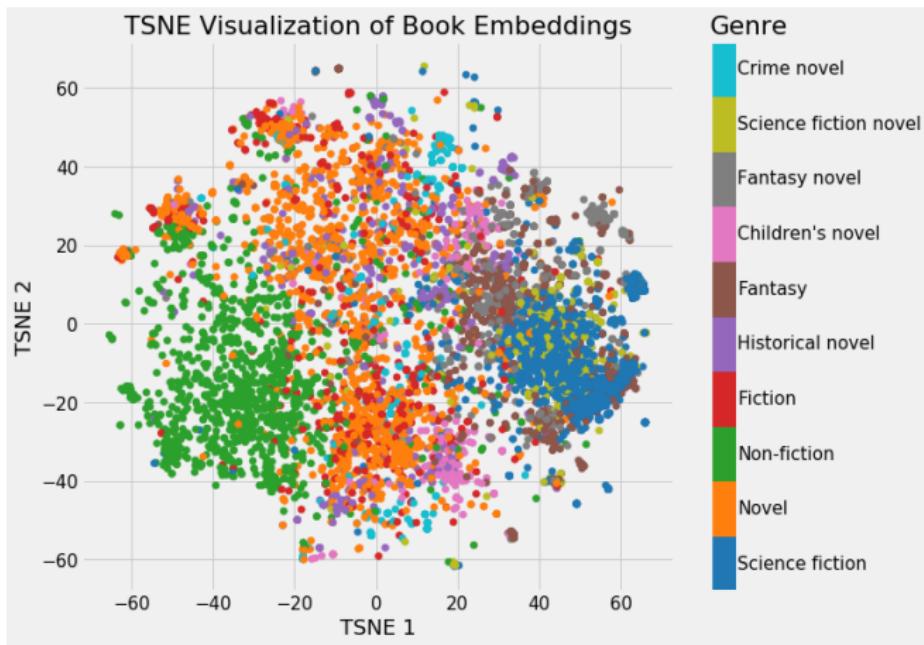


- Encode-image as vector use a deep convolutional network
 - generate caption using recurrent network (RNN/LSTM)
 - all parameters optimised (using example image captions)

- Interesting to examine output of the **encoder**
 - does not have to be interpretable
 - just convey information from image to caption generation
- An example of **end-to-end** training
 - **train the system to complete task!**
 - minimises need for human expertise ...
 - no need to decide which aspect of image is important

Word Embeddings

Book Embeddings



- First stage of many text processing stages is word embedding

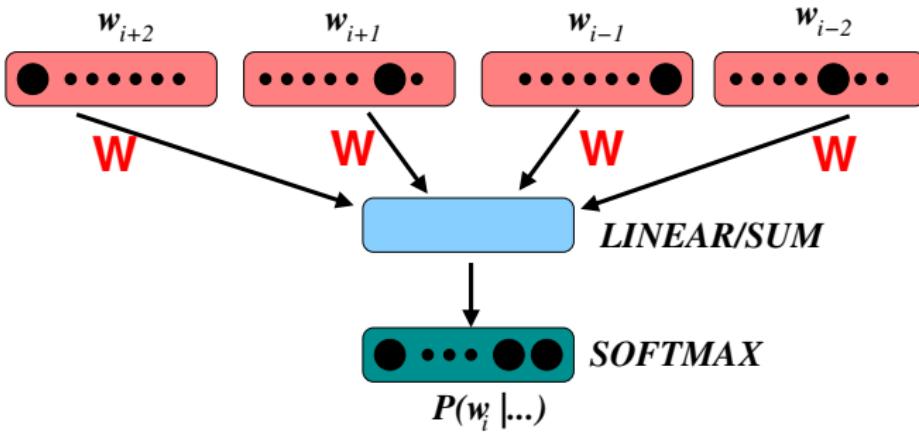
$$\begin{array}{ccc} w_i & \rightarrow & \mathbf{t}_i \\ \text{word} & & \text{1-hot coding} \end{array} \quad \rightarrow \quad \mathbf{W}\mathbf{t}_i = \mathbf{x}_i \quad \text{(linear) embedding}$$

- able to control dimensionality of \mathbf{x}_i
- need to learn transformation matrix \mathbf{W}
- Could train on specific target training data
 - optimal transformation for target domain
 - often only limited data available

Train embedding just on general text data

- Able to fine-tune to specific domain - transfer learning

Word2Vec: Continuous Bag of Words (CBOW)



- **Continuous Bag of words (CBOW):** like feed-forward NN
 - learn shared word projections, W : maximise

$$P(w_i | w_{i+2}, w_{i+1}, w_{i-1}, w_{i-2})$$

- other variants such as **Skip N-Gram** possible

Word2Vec Relationships

- System trained on Google 1-Billion word corpus
 - CBOW model used to train [Word2Vec](#) model
 - examine [linear](#) relationships between words

A	B	C	A-B+C
france	paris	rome	italy
king	man	woman	queen
did	do	give	gave

Context Dependent Word Embeddings [7]

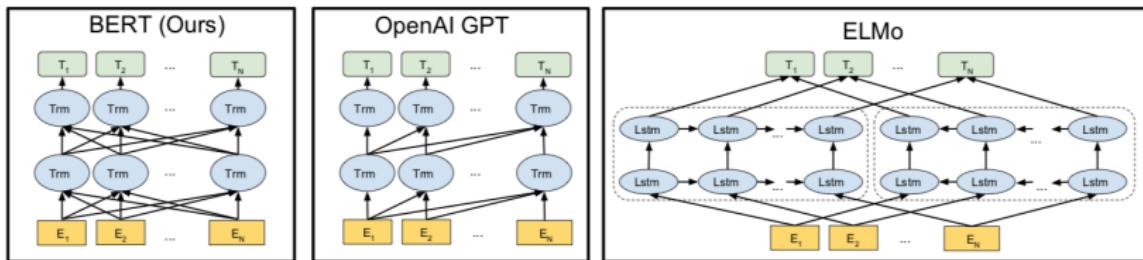
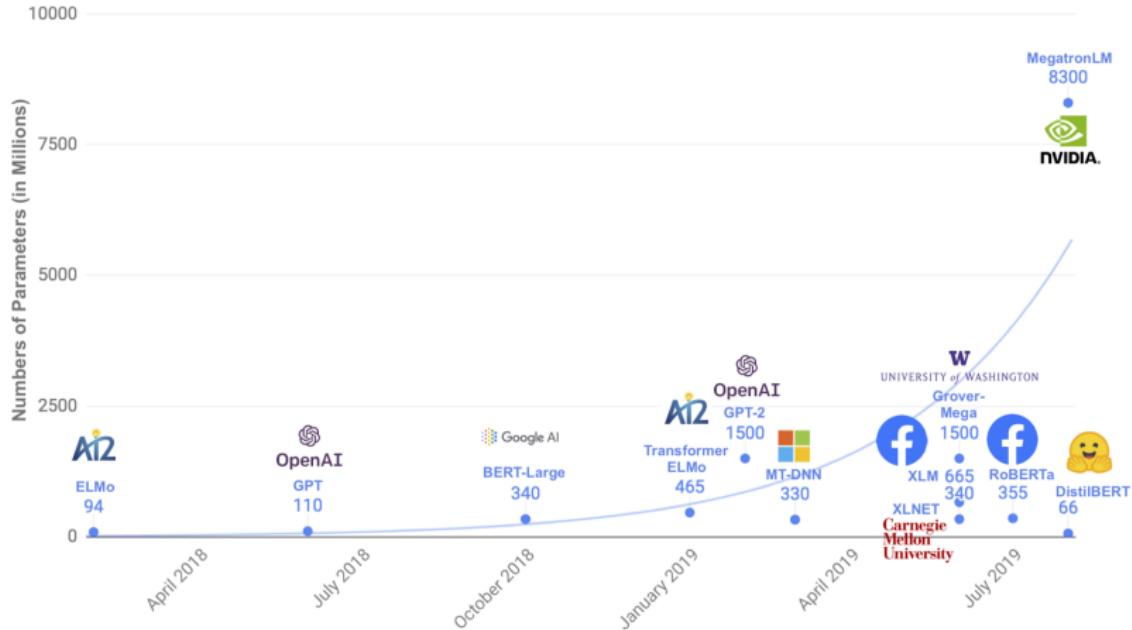


Figure: Figure from Google paper

- Interest in **context dependent** word embeddings
 - no unique embedding for a word
 - BERT**: trained on large amounts of text data used

Recent Word Embeddings [18]



Language Modelling

- Predict next word given observed sequence $\hat{w}_1, \dots, \hat{w}_{i-1}$

$$\hat{w}_i \sim P(w_i | \hat{w}_1, \dots, \hat{w}_{i-1}); \quad P(w_{i+1} | \hat{w}_1, \dots, \hat{w}_i)$$

- can repeat the process to generate a word sequence
- simple **auto-regressive** data generation process
- Simple data generation example
 - initialise process with $w_0 = <\text{s}>$, sentence start symbol

$<\text{s}> \text{ ???}$

$P(w_1 | <\text{s}>)$

$<\text{s}> \text{ The } \text{ ???}$

$P(w_2 | <\text{s}>, \text{The})$

$<\text{s}> \text{ The cat } \text{ ???}$

$P(w_3 | <\text{s}>, \text{The}, \text{cat})$

$<\text{s}> \text{ The cat sat } \text{ ???}$

$P(w_4 | <\text{s}>, \text{The}, \text{cat}, \text{sat})$

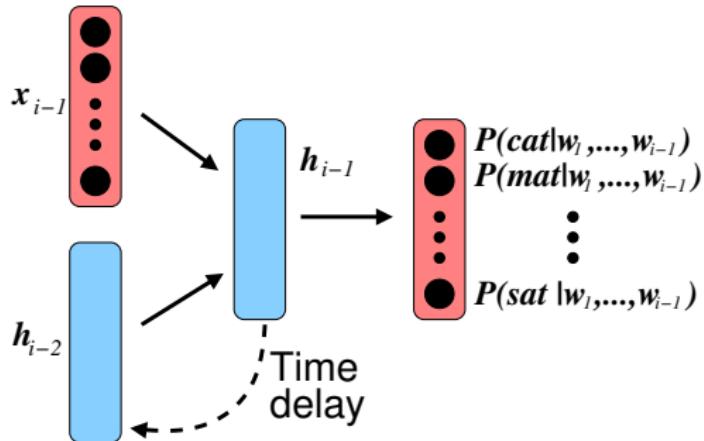
$<\text{s}> \text{ The cat sat on } \text{ ???}$

$P(w_5 | <\text{s}>, \text{The}, \text{cat}, \text{sat}, \text{on})$

- Can also be used to model word sequence distribution

$$P(\mathbf{w}_{1:L}) = P(w_1, \dots, w_L) = \prod_{i=1}^L P(w_i | w_1, \dots, w_{i-1})$$

RNN Language Models



- Use RNN **history vector** as word-history representation

$$P(\mathbf{w}_{1:L}) \approx \prod_{i=1}^L P(w_i | w_{i-1}, \mathbf{h}_{i-2}) = \prod_{i=1}^L P(w_i | \mathbf{x}_{i-1}, \mathbf{h}_{i-2}) \approx \prod_{i=1}^L P(w_i | \mathbf{h}_{i-1})$$

- impact of word in history gradually decreases with distance from prediction

Opinion Artificial intelligence (AI)

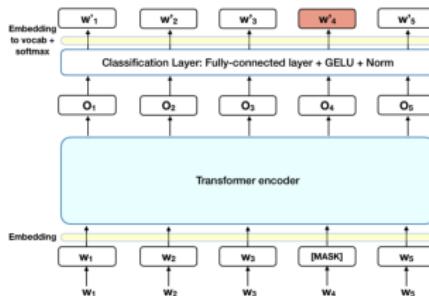
A robot wrote this entire article. Are you scared yet, human?

GPT-3

- Deep-learning has revolutionised language modelling
 - GPT-3 is a [large](#) language model from OpenAI
 - about 175 billion model parameters trained on 500 billion tokens
- Use sequence-to-sequence transformer architecture

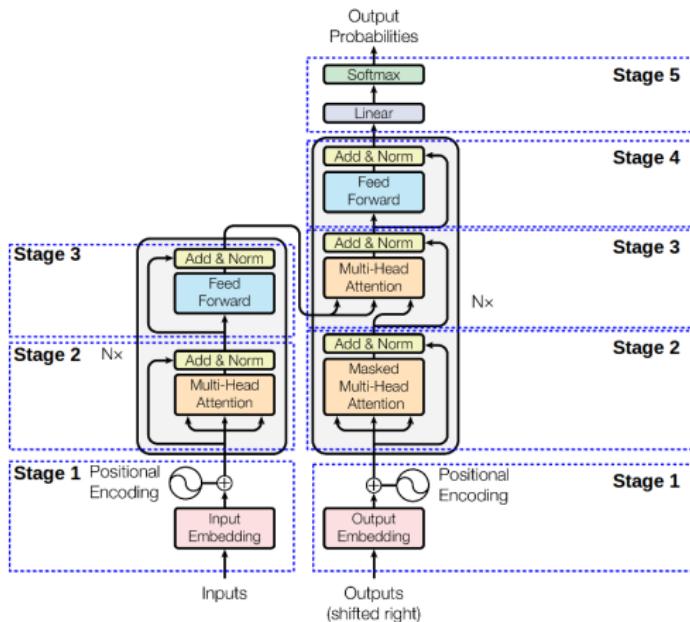
Masked Language Model: Start of BERT [7]

- Rather than predicting the next word in the sequence
 - predict one or more masked words in the complete sequence
 - uses information from complete sequence



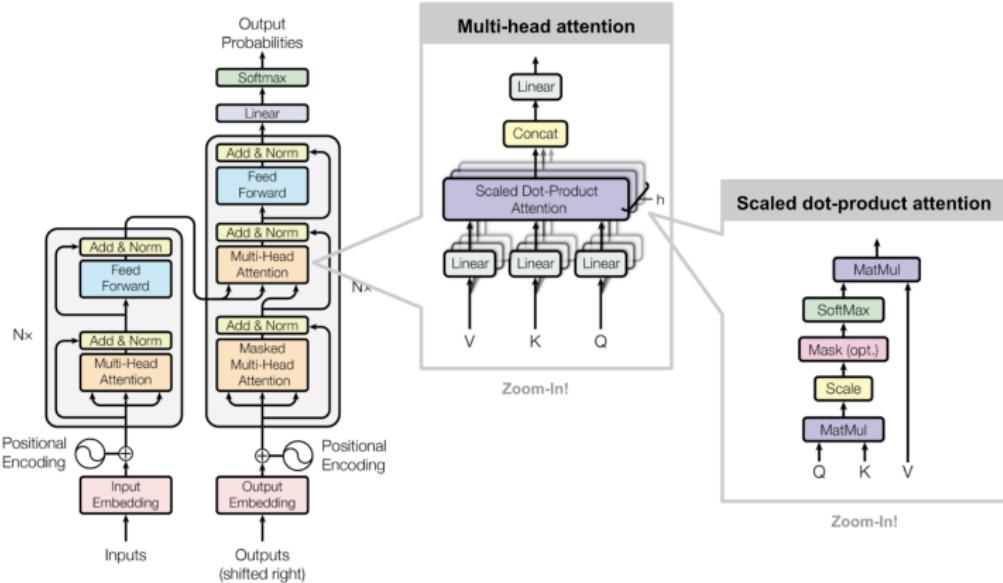
- Illustration above - fourth word is masked and predicted
 - (GELU - Gaussian Error Linear Unit, similar to a ReLU)
- Transformer encoder used as part of network
 - popular encoder, part of powerful sequence-to-sequence model
 - transformer decoder used in GPT-2/3

Transformer: Sequence-to-Sequence Model [21]



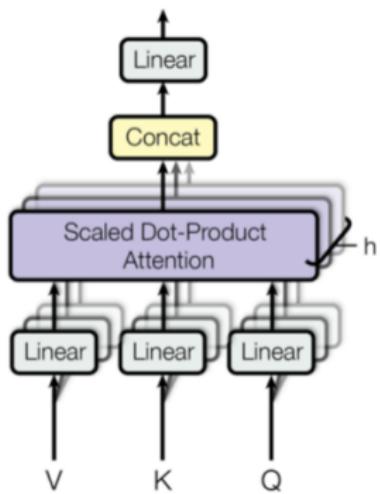
- Only based on self-attention: encoding/history/context
 - multiple transformer blocks for encoder and decoder

Transformer: Detail



- **transform block** comprises:
 - residual network plus layer norm (**Add & Norm**) acting on:
 - (1) **Multi-Head Attention** layer; (2) **Feed-Forward** network

Encoder Multi-Head Self-Attention



- Notation
 - Q : query for attention
 - K : key for attention - length d_k
 - V : input values ($\mathbf{v}_1, \dots, \mathbf{v}_L$)
- Scaled dot-product multi-head attention

$$\alpha_{\tau i}^{(j)} = \frac{1}{Z} \exp\left(\mathbf{k}_{\tau}^{(j)\top} \mathbf{q}_i^{(j)} / \sqrt{d_k^{(j)}}\right)$$

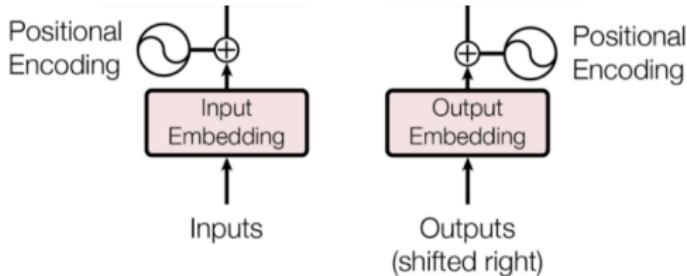
$$\mathbf{c}_i^{(j)} = \sum_{\tau=1}^L \alpha_{\tau i}^{(j)} \mathbf{W}_v^{(j)} \mathbf{v}_{\tau}$$

$$\mathbf{c}_i = [\mathbf{c}_i^{(1)\top} \quad \dots \quad \mathbf{c}_i^{(J)\top}]^{\top}$$

- Self-attention $Q = K = V$

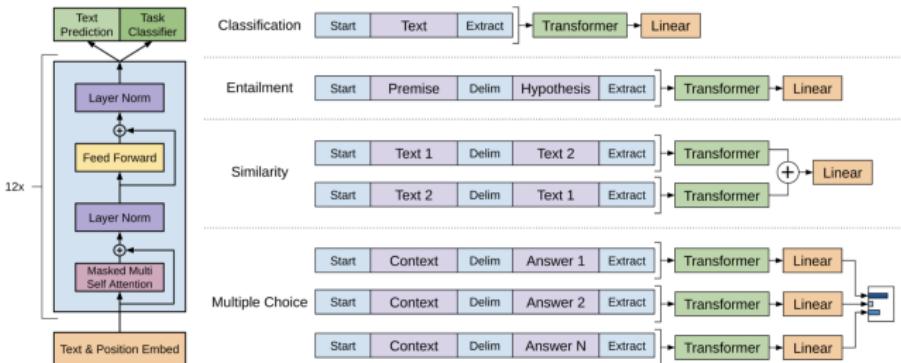
$$\mathbf{k}_i^{(j)} = \mathbf{W}_k^{(j)} \mathbf{v}_i; \quad \mathbf{q}_i^{(j)} = \mathbf{W}_q^{(j)} \mathbf{v}_i;$$

Position Encoding (reference)



- Attention has no representation of position
 - but predictions should be sensitive to location
 - Add **position encoding** vector word embedding
 - use sines and cosines of different frequencies
- $$\text{pe}_{\tau,2i} = \sin(\tau/10000^{2i/d})$$
- $$\text{pe}_{\tau,2i+1} = \cos(\tau/10000^{2i/d})$$
- d is the number of nodes in a network layer
 - τ is the position of the word, i element of position vector

Transformer Decoder (reference) [13, 2]



- GPT-2/3 based on the transformer-decoder architecture
 - similar to decoder-side of the sequence-to-sequence transformer
 - autoregressive model with transformer units (no encoder)
 - train on multiple contiguous sentences (inc. sentence end tok)
- Pre-trained on large amounts of language model data
 - can be fine-tuned to any task of interest (RHS above)

- [1] S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer, "Scheduled sampling for sequence prediction with recurrent neural networks," in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds. Curran Associates, Inc., 2015, pp. 1171–1179. [Online]. Available: <http://papers.nips.cc/paper/5956-scheduled-sampling-for-sequence-prediction-with-recurrent-neural-networks.pdf>
- [2] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," 2020.
- [3] W. Chan, N. Jaitly, Q. V. Le, and O. Vinyals, "Listen, attend and spell," *CoRR*, vol. abs/1508.01211, 2015. [Online]. Available: <http://arxiv.org/abs/1508.01211>
- [4] J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," *CoRR*, vol. abs/1506.07503, 2015. [Online]. Available: <http://arxiv.org/abs/1506.07503>
- [5] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [6] J. Chung, K. Kastner, L. Dinh, K. Goel, A. C. Courville, and Y. Bengio, "A recurrent latent variable model for sequential data," *CoRR*, vol. abs/1506.02216, 2015. [Online]. Available: <http://arxiv.org/abs/1506.02216>
- [7] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, J. Burstein, C. Doran, and T. Solorio, Eds. Association for Computational Linguistics, 2019, pp. 4171–4186. [Online]. Available: <https://doi.org/10.18653/v1/n19-1423>
- [8] A. Graves, "Sequence transduction with recurrent neural networks," *CoRR*, vol. abs/1211.3711, 2012. [Online]. Available: <http://arxiv.org/abs/1211.3711>
- [9] A. Graves, A.-R. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *2013 IEEE international conference on acoustics, speech and signal processing*. IEEE, 2013, pp. 6645–6649.

- [10] K. M. Hermann, T. Kocišký, E. Grefenstette, L. Espeholt, W. Kay, M. Suleyman, and P. Blunsom, "Teaching machines to read and comprehend," *CoRR*, vol. abs/1506.03340, 2015. [Online]. Available: <http://arxiv.org/abs/1506.03340>
- [11] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.
- [12] A. M. Lamb, A. G. ALIAS PARTH GOYAL, Y. Zhang, S. Zhang, A. C. Courville, and Y. Bengio, "Professor forcing: A new algorithm for training recurrent networks," in *Advances in Neural Information Processing Systems 29*, D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, Eds. Curran Associates, Inc., 2016, pp. 4601–4609. [Online]. Available: <http://papers.nips.cc/paper/6099-professor-forcing-a-new-algorithm-for-training-recurrent-networks.pdf>
- [13] P. J. Liu, M. Saleh, E. Pot, B. Goodrich, R. Sepassi, L. Kaiser, and N. Shazeer, "Generating wikipedia by summarizing long sequences," 2018.
- [14] L. Lu, X. Zhang, K. Cho, and S. Renals, "A study of the recurrent neural network encoder-decoder for large vocabulary speech recognition," in *Proc. INTERSPEECH*, 2015.
- [15] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," *CoRR*, vol. abs/1310.4546, 2013. [Online]. Available: <http://arxiv.org/abs/1310.4546>
- [16] T. Robinson and F. Fallside, "A recurrent error propagation network speech recognition system," *Computer Speech & Language*, vol. 5, no. 3, pp. 259–274, 1991.
- [17] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Parallel distributed processing: Explorations in the microstructure of cognition, vol. 1," D. E. Rumelhart, J. L. McClelland, and C. PDP Research Group, Eds. Cambridge, MA, USA: MIT Press, 1986, ch. Learning Internal Representations by Error Propagation, pp. 318–362.
- [18] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter," 2019.
- [19] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," *IEEE Transactions on Signal Processing*, vol. 45, no. 11, pp. 2673–2681, 1997.

- [20] R. K. Srivastava, K. Greff, and J. Schmidhuber, "Highway networks," *CoRR*, vol. abs/1505.00387, 2015. [Online]. Available: <http://arxiv.org/abs/1505.00387>
- [21] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *CoRR*, vol. abs/1706.03762, 2017. [Online]. Available: <http://arxiv.org/abs/1706.03762>