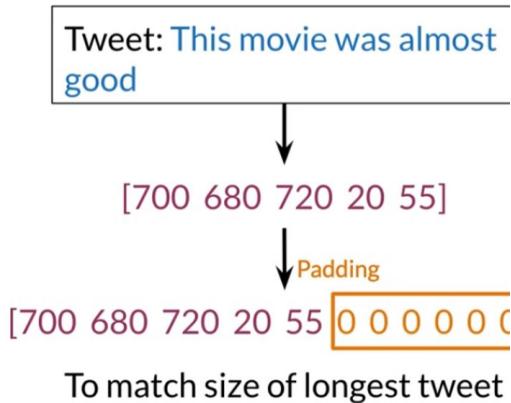
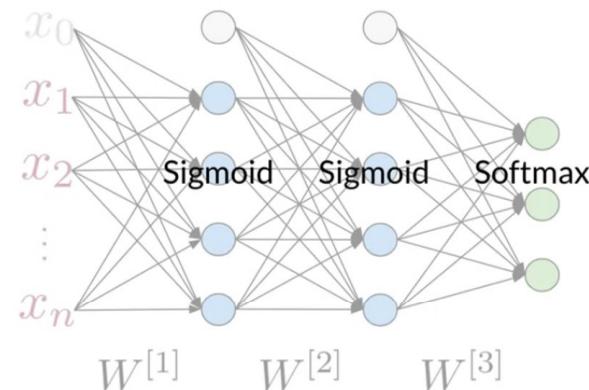


vital Representation

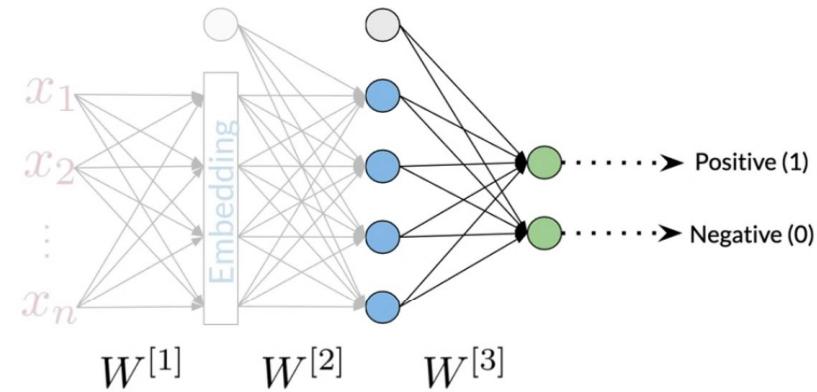
Word	Number
a	1
able	2
about	3
...	...
hand	615
...	...
happy	621
...	...
zebra	1000



ural Networks in Trax



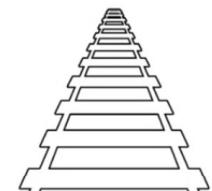
ural Networks for sentiment analysis



tline

Define a basic neural network using Trax

Benefits of Trax



vantages of using frameworks

Run fast on CPUs, GPUs and TPUs

Parallel computing

Record algebraic computations for gradient evaluation

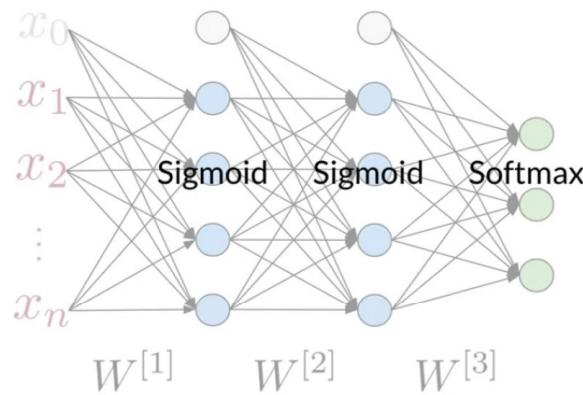
Tensorflow

Pytorch

JAX

plearning.ai

ural Networks in Trax



```
from trax import layers as tl
Model = tl.Serial(
    tl.Dense(4),
    tl.Sigmoid(),
    tl.Dense(4),
    tl.Sigmoid(),
    tl.Dense(3),
    tl.Softmax())
```

TensorFlow → Translate → Tensor2Tensor → Trax

TensorFlow:

Large-Scale Machine Learning on Heterogeneous Distributed Systems

(Preliminary White Paper, November 9, 2015)

Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng
Google Research*

around 2014-15 when we were making TensorFlow.



nmary

Order of computation → Model in Trax

Benefits from using frameworks

plearning.ai

plearning.ai

Trax

Deep Learning with Clear Code and Speed

Let me tell you why, so,

Trax makes programmers efficient

- Bottom-up clean re-design
- Easy to debug, you can read the code
- Full models with dataset bindings included
- Main models regression-tested daily

Is there a price to pay?

- Backwards compatibility
- Trax code and understand what's going to come.

TensorFlow → Translate → Tensor2Tensor → Trax

Accelerating Deep Learning Research with the Tensor2Tensor Library

Monday, June 19, 2017

Posted by Łukasz Kaiser, Senior Research Scientist, Google Brain Team

Translation Model	Training time	BLEU (difference from baseline)
Transformer (T2T)	3 days on 8 GPU	28.4 (+7.8)
SliceNet (T2T)	6 days on 32 GPUs	26.1 (+5.5)
GNMT + Mixture of Experts	1 day on 64 GPUs	26.0 (+5.4)
ConvS2S	18 days on 1 GPU	25.1 (+4.5)
GNMT	1 day on 96 GPUs	24.6 (+4.0)
ByteNet	8 days on 32 GPUs	23.8 (+3.2)
MOSES (phrase-based baseline)	N/A	20.6 (+0.0)

To do that, we created the Tensor2Tensor Library.

What do you want from a deep learning library?

- Makes programmers efficient
- Runs code fast

GPU (8x V100)	on-demand: \$19.84 / hour	preemptible: \$5.92 / hour
TPU (8x v3)	on-demand: \$8 / hour	preemptible: \$1.40 / hour

there are really two things that matters.

Trax runs code fast

- Designed to use a JIT compiler with JAX and XLA
- JAX: fastest Transformer in MLPerf 2020 (JAX: 0.26, TF: 0.35, pyTorch: 0.62)
- No code changes at all between CPU, GPU and TPU, preemptible training
- Tested with TPUs on colab too:

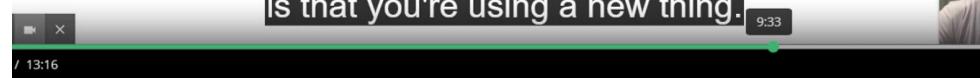


Trax Adam

```
class Adam(Optimizer):
    """Adam optimizer."""

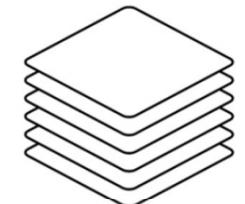
    def init(self, param):
        m = np.zeros_like(param)
        v = np.zeros_like(param)
        return m, v

    def update(self, step, grads, param, slots, opt_params):
        m, v = slots
        learning_rate, b1, b2, eps = opt_params
        m = (1 - b1) * grads + b1 * m # First moment estimate.
        v = (1 - b2) * (grads ** 2) + b2 * v # Second moment estimate.
        mhat = m / (1 - b1 ** (step + 1)) # Bias correction.
        vhat = v / (1 - b2 ** (step + 1))
        param = param - (
            learning_rate * mhat / (np.sqrt(vhat) + eps)).astype(param.dtype)
        return param, (m, v)
```



tline

How classes work and their implementation



their methods.

plearning.ai

Trax is fun to learn and use!

```
# Create a Transformer model.
# Pre-trained model config in gs://trax-ml/models/translation/ende_wmt32k.g
model = trax.models.Transformer(
    input_vocab_size=33300,
    d_model=512, d_ff=2048,
    n_heads=8, n_encoder_layers=6, n_decoder_layers=6,
    max_len=2048, mode='train')
"""Returns a Transformer language model."""
embedder = [
    tl.Embedding(d_model, vocab_size),
    tl.Dropout(rate=dropout, mode=mode),
    tl.PositionalEncoding(max_len=max_len),]

@return tl.Serial(
    tl.ShiftRight(),
    embedder,
    [DecoderBlock(d_model, d_ff, n_heads,
                  dropout, mode)
     for _ in range(n_layers)],
    tl.LayerNorm(),
    tl.Dense(vocab_size),
    tl.LogSoftmax(),
```

because we had the liberty
to do things from scratch



Classes in Python

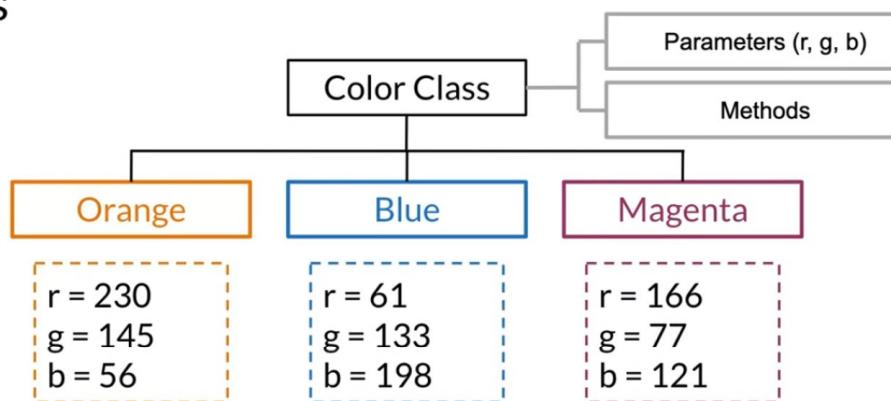
```
class MyClass(Object):
    def __init__(self, y):
        self.y = y
    def my_method(self,x):
        return x + self.y
    def __call__(self, x):
        return self.my_method(x)
```

```
f = MyClass(7)
```

plearning.ai

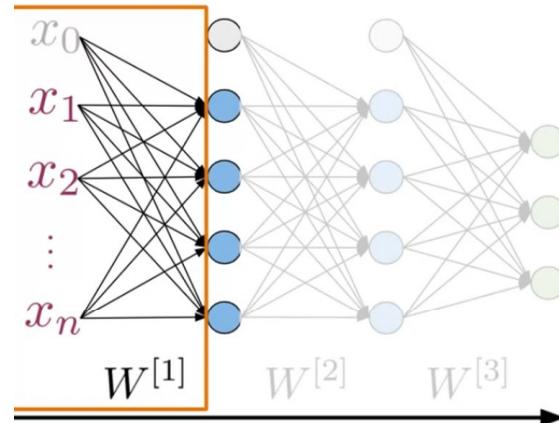
the program will return
the results from my_method.

Inheritance



plearning.ai 1:00
Magenta, would be defined with
a specific set of r, g, b values.

Neural Network Layer



Fully connected layer

$$z^{[i]} = W^{[i]} a^{[i]}$$

plearning.ai You will have to compute
the inner products between

Subclasses

```
class MyClass(Object):
    def __init__(self,y):
        self.y = y
    def my_method(self,x):
        return x + self.y
    def __call__(self,x):
        return self.my_method(x)
```

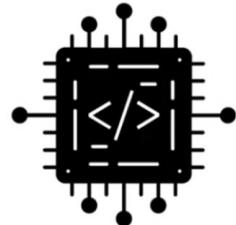
```
class SubClass(MyClass):
    def my_method(self,x):
        return x + self.y**2
```

plearning.ai 1:00
will be inherited from the parent class,
but my_method will be replaced.

Outline

Embedding layer

Mean layer



the mean layer in serial models

elearning.ai

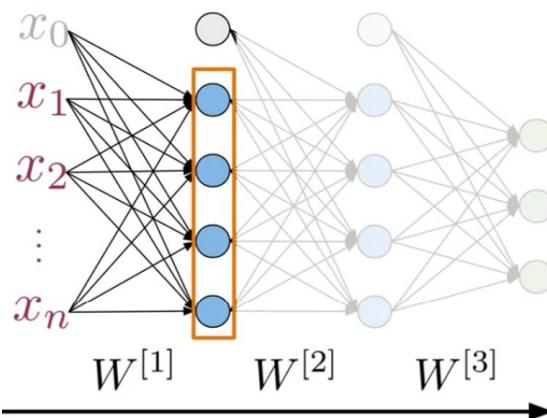
Embedding Layer

Vocabulary Index

Vocabulary	Index	0.020	0.006
I	1	0.020	0.006
am	2	-0.003	0.010
happy	3	0.009	0.010
because	4	-0.011	-0.018
learning	5	-0.040	-0.047
NLP	6	-0.009	0.050
sad	7	-0.044	0.001
not	8	0.011	-0.022

Trainable weights

LU Layer



ReLU = Rectified linear

$$a^{[i]} = g(z^{[i]})$$

elearning.ai

that typically follows a dense fully connected layer,

Embedding Layer

Vocabulary Index

Vocabulary	Index	0.020	0.006
I	1	0.020	0.006
am	2	-0.003	0.010
happy	3	0.009	0.010
because	4	-0.011	-0.018
learning	5	-0.040	-0.047
NLP	6	-0.009	0.050
sad	7	-0.044	0.001
not	8	0.011	-0.022

we'll return a vector equal to 0.020,

elearning.ai

Embedding Layer

Text: I am happy

Vocabulary	Index	0.020	0.006
I	1	0.020	0.006
am	2	-0.003	0.010
happy	3	0.009	0.010

Mean of the word embeddings

0.020	0.006
-0.003	0.010
0.009	0.010

0.009
0.009

the same number of features as you're embedding size.

plearning.ai

Embedding Layer

Vocabulary	Index	0.020	0.006
I	1	0.020	0.006
am	2	-0.003	0.010
happy	3	0.009	0.010
because	4	-0.011	-0.018
learning	5	-0.040	-0.047
NLP	6	-0.009	0.050
sad	7	-0.044	0.001
not	8	0.011	-0.022

treated as a hyperparameter in your model.

plearning.ai

Computing gradients in Trax

$$f(x) = 3x^2 + x$$

$$\frac{\partial f(x)}{\partial x} = 6x + 1$$

Gradient

```
def f(x):
    return 3*x**2 + x
grad_f = trax.math.grad(f)
```

Returns a function

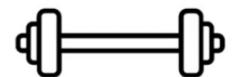
The function grad will return a function

Press Esc to exit full screen

Line

Computing gradients in trax

Training using grad()



Train a weight
Vocab x Embeds

allows you to easily compute gradients,

plearning.ai

Additional Language Models



J'ai vu le match de foot

Sequence	$P(\text{Sequence})$
I saw the game of soccer	4.5 e-0
I saw the soccer game	6.0 e-0
I saw the soccer match	4.6 e-0
Saw I the game of soccer	2.6 e-0

For an accurate translation, you could
plearning.ai compute the probabilities of each sentence

N-grams

$$P(w_2|w_1) = \frac{\text{count}(w_1, w_2)}{\text{count}(w_1)} \longrightarrow \text{Bigrams}$$

$$P(w_3|w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)} \longrightarrow \text{Trigrams}$$

$$P(w_1, w_2, w_3) = P(w_1) \times P(w_2|w_1) \times P(w_3|w_2)$$

- Large N-grams to capture dependencies between distant words
- Need a lot of space and RAM

Even in the case of large corpora your
plearning.ai model would need a lots of space and

Training - DeepLearn...

```
y = model(x)
grads = grad(y.forward)(y.weights,x)
```

In a loop

```
weights -= alpha*grads
```

Forward and
Back-propagation

performed in a single line.

plearning.ai

N-grams

$$P(w_2|w_1) = \frac{\text{count}(w_1, w_2)}{\text{count}(w_1)} \longrightarrow \text{Bigrams}$$

$$P(w_3|w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)} \longrightarrow \text{Trigrams}$$

$$P(w_1, w_2, w_3) = P(w_1) \times P(w_2|w_1) \times P(w_3|w_2)$$

- Large N-grams to capture dependencies between distant words

One of them is that in order to
plearning.ai capture dependencies between words

vantages of RNNs

I was supposed to study with me. I called her but she did not [answ](#)

want
respond
choose
want
have
ask
attempt
answer
know

RNNs look at every previous word

Similar probabilities with trigram

plearning.ai

To see how recurrent
neural networks work,

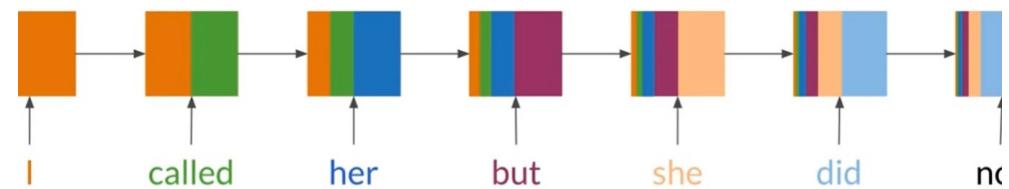
1:23

|Ns Basic Structure

I led her but she did not _____

|Ns Basic Structure

I led her but she did not _____



represents the
computations made at

plearning.ai

plearning.ai

To see how recurrent
neural networks work,

1:19

I led her but she did not _____

vantages of RNNs

I was supposed to study with me. I called her but she did not [answ](#)

want
respond
choose
want
have
ask
attempt
answer
know

RNNs look at every previous word

Similar probabilities with trigram

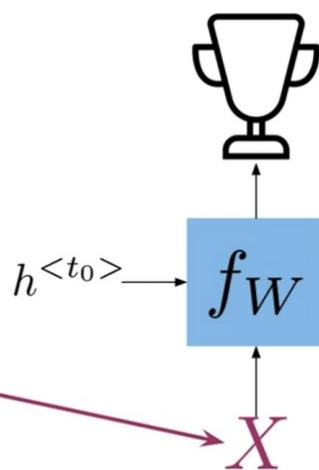
you would have to account
for six-word-long sequences,

plearning.ai

plearning.ai

One to One

	LaLiga Santander
Real Valladolid	0 1
Real Zaragoza	0 4
Atletico Madrid	0 1



plearning.ai

It only has an additional hidden state h superscript t subscript 0 .

Many to One

Sentiment analysis

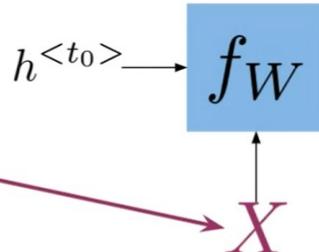
Positive

Set: I am very happy !

your model should output whether the tweet has a negative or positive sentiment.

One to One

	LaLiga Santander
Real Valladolid	0 1
Real Zaragoza	0 4
Atletico Madrid	0 1

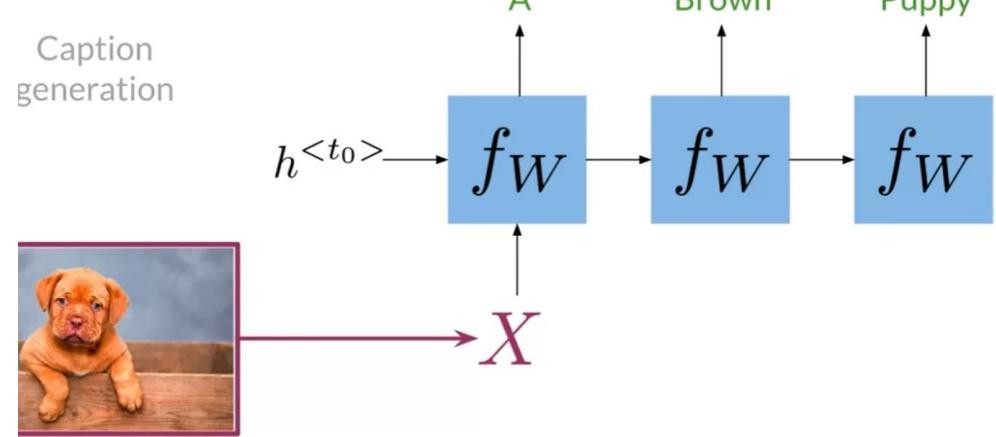


plearning.ai

However, notice that this recurrent neural network isn't much different

One to Many

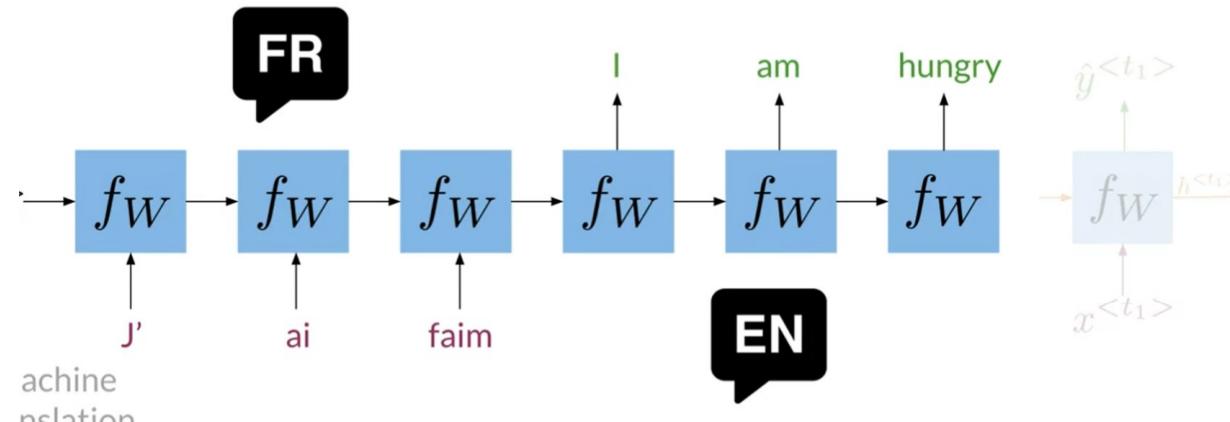
Caption generation



because your RNN takes a single image and generates multiple words to describe it.

plearning.ai

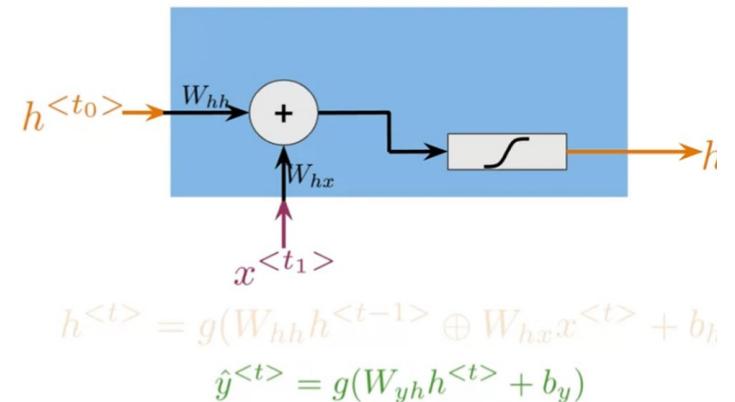
Many to Many



this task, because they propagate information from the beginning to end.

plearning.ai

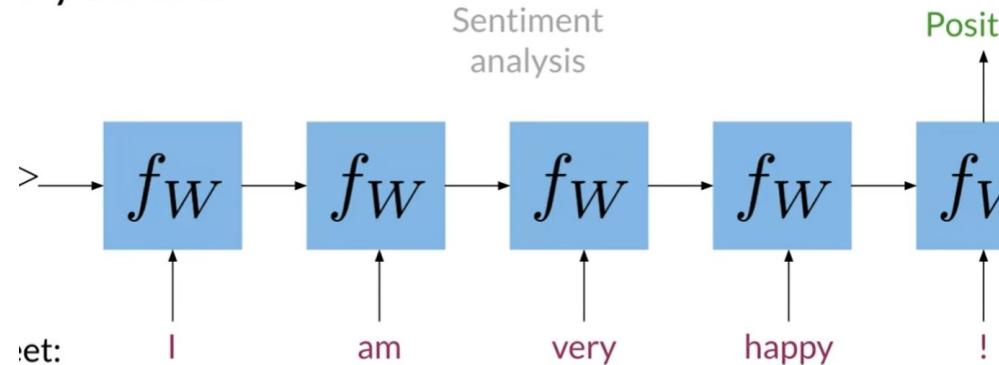
Vanilla RNN



you can compute the current y -hat by multiplying

plearning.ai

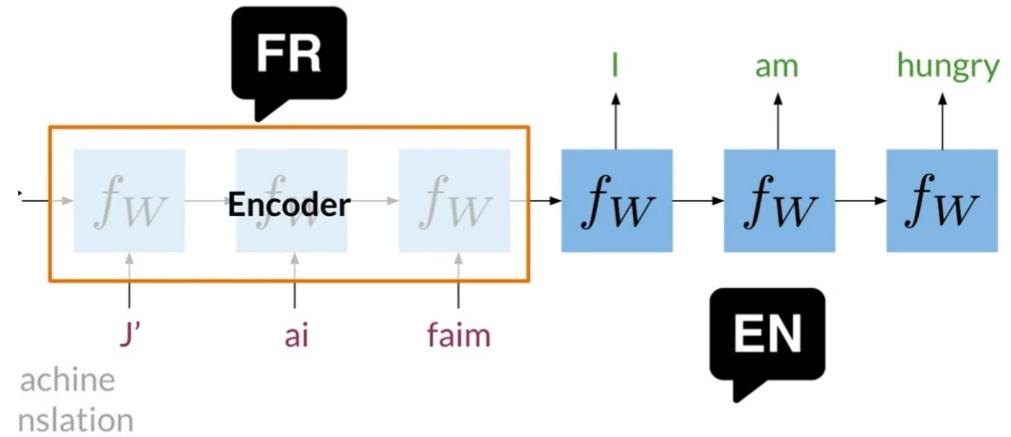
Many to One



Your RNN would take every word from the sequence as inputs in different steps,

plearning.ai

Many to Many



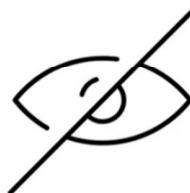
machine translation.

plearning.ai

Summary

Hidden states propagate information through time

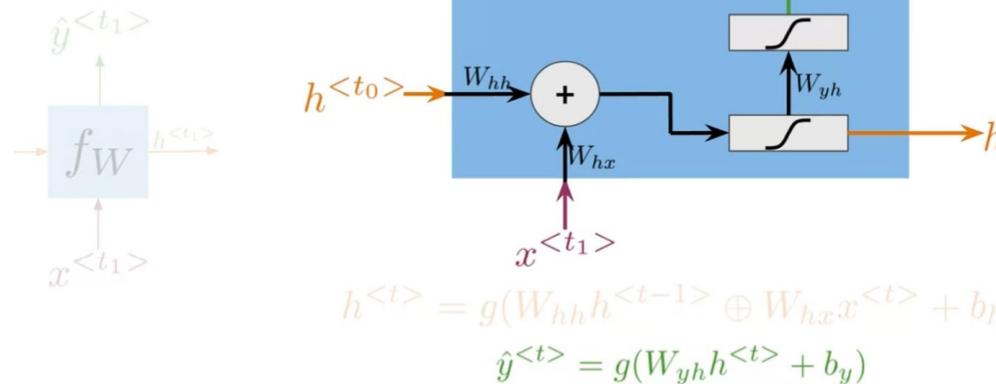
Basic recurrent units have two inputs at each time: $h^{<t-1>}$, $x^{<t>}$



**different positions
within a sequence.**

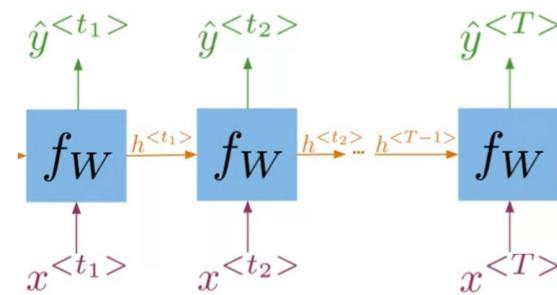
plearning.ai

Vanilla RNN



plearning.ai

Cross Entropy Loss



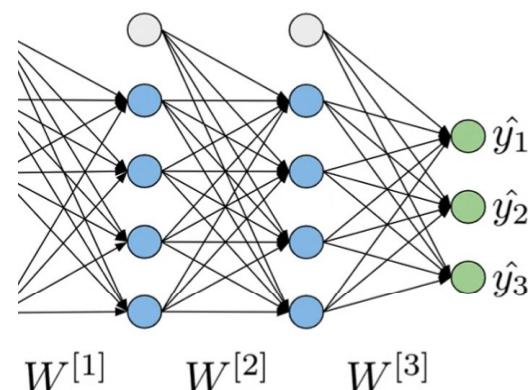
$$h^{<t>} = g(W_h[h^{<t-1>}], x^{<t>})$$

$$\hat{y}^{<t>} = g(W_yh^{<t>} + b_y)$$

$$J = -\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^K y_j^{<t>} \log \hat{y}_j^{<t>}$$

**the summation over time t and the division
by the total number of steps, capital T.**

Cross Entropy Loss



K - classes or possibilities

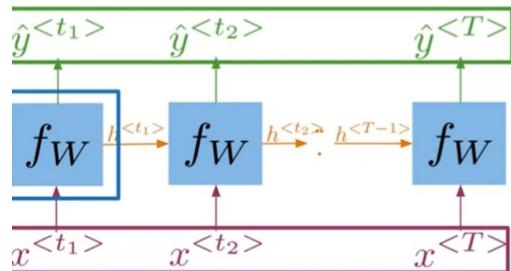
$$J = -\sum_{j=1}^K y_j \log \hat{y}_j$$

Looking at a single example (x)

**output probability in \hat{y} -hat and compare
them with the true y values, 0 or 1.**

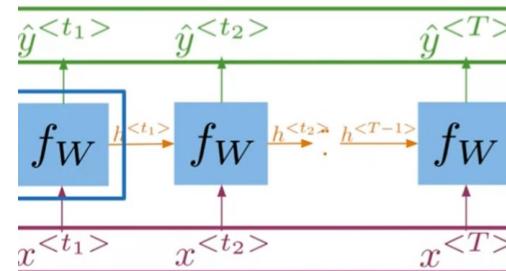
plearning.ai

can() function



```
def scan(fn, elems, initializer=None):
    cur_value = initializer
    ys = []
    for x in elems:
        y, cur_value = fn(x, cur_value)
        ys.append(y)
    return ys, cur_value
```

can() function



```
def scan(fn, elems, initializer=None):
    cur_value = initializer
    ys = []
    for x in elems:
        y, cur_value = fn(x, cur_value)
        ys.append(y)
    return ys, cur_value
```

Finally, the function returns the list of predictions, and the last hidden state.

1:20

predictions, and the last hidden state.

1:20

Finally, the function returns the list of predictions, and the last hidden state.

1:16

predictions, and the last hidden state.

1:16

line

scan() function in tensorflow

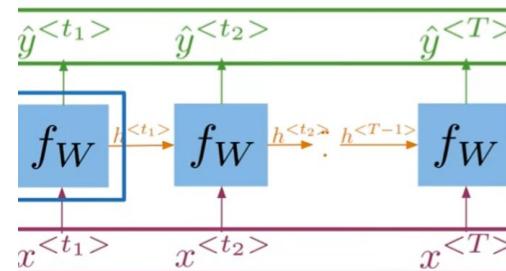
Computation of forward propagation using abstractions



for computing the forward pass in RNNs.

plearning.ai

can() function



```
def scan(fn, elems, initializer=None):
    cur_value = initializer
    ys = []
    for x in elems:
        y, cur_value = fn(x, cur_value)
        ys.append(y)
    return ys, cur_value
```

stores the values of the prediction and hidden states.

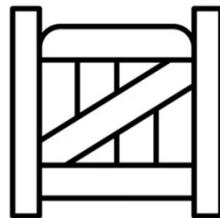
1:16

hidden states.

1:16

Outline

- Gated recurrent unit (GRU) structure
- Comparison between GRUs and vanilla RNNs



Gated Recurrent Units

"**Ants** are really interesting. They are everywhere."

Plural

Relevance and update gates to remember important prior information

deeplearning.ai

GRUs for shorts, with a comparison to vanilla RNNs.

no...ipynb ...

deeplearning.ai

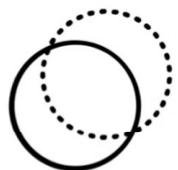
computing irrelevance and update gates,

no...ipynb ...

Summary

Frameworks require abstractions

tf.scan() mimics RNNs



Gated Recurrent Units

"**Ants** are really interesting. They are everywhere."

Plural

keep the information about the subject.

deeplearning.ai

It is important to know that these types of abstractions are needed for deep

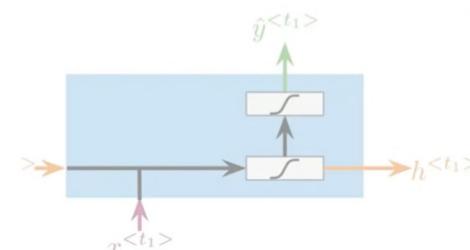
202

deeplearning.ai

no...ipynb ...

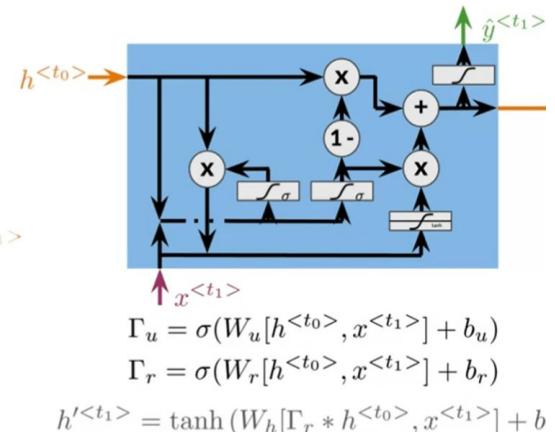
Vanilla RNN vs GRUs

Vanilla RNN vs GRUs



$$h^{<t>} = g(W_h[h^{<t-1>}, x^{<t>}] + b_h)$$

$$\hat{y}^{<t>} = g(W_{yh}h^{<t>} + b_y)$$



$$\Gamma_u = \sigma(W_u[h^{<t_0>}, x^{<t_1>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[h^{<t_0>}, x^{<t_1>}] + b_r)$$

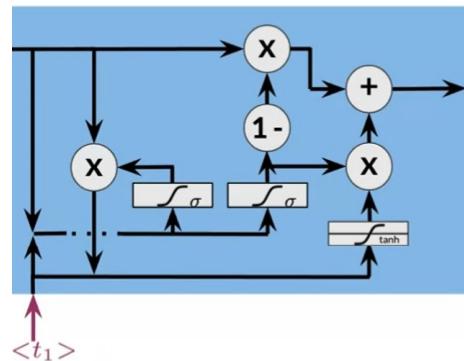
$$h'^{<t_1>} = \tanh(W_h[\Gamma_r * h^{<t_0>}, x^{<t_1>}] + b)$$

Let's compare GRUs
with vanilla RNNs.

plearning.ai

The current hidden
state is computed and

Gated Recurrent Unit



Gates to keep/update relevant information in the hidden state

$$\Gamma_r = \sigma(W_r[h^{<t_0>}, x^{<t_1>}] + b_r)$$

$$\Gamma_u = \sigma(W_u[h^{<t_0>}, x^{<t_1>}] + b_u)$$

$$h'^{<t_1>} = \tanh(W_h[\Gamma_r * h^{<t_0>}, x^{<t_1>}])$$

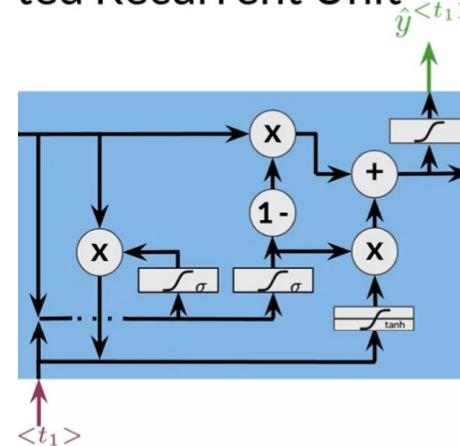
Hidden state candidate

$$h^{<t_1>} = \Gamma_u * h^{<t_0>} + (1 - \Gamma_u) * h'^{<t_1>}$$

After that, a new value
for the hidden state is

plearning.ai

Gated Recurrent Unit



Gates to keep/update relevant information in the hidden state

$$\Gamma_r = \sigma(W_r[h^{<t_0>}, x^{<t_1>}] + b_r)$$

$$\Gamma_u = \sigma(W_u[h^{<t_0>}, x^{<t_1>}] + b_u)$$

$$h'^{<t_1>} = \tanh(W_h[\Gamma_r * h^{<t_0>}, x^{<t_1>}])$$

Hidden state candidate

$$h^{<t_1>} = \Gamma_u * h^{<t_0>} + (1 - \Gamma_u) * h'^{<t_1>}$$

$$\hat{y}^{<t_1>} = g(W_{yh}h^{<t_1>} + b_y)$$

Finally, a prediction v hat

plearning.ai

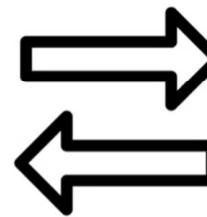
238

426

tline

How bidirectional RNNs propagate information

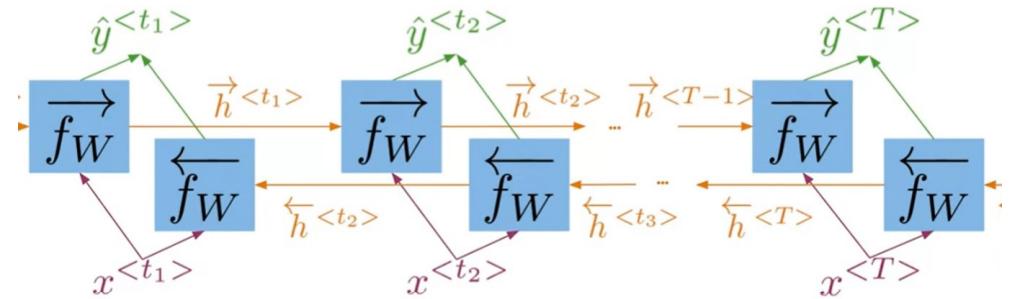
Forward propagation in deep RNNs



plearning.ai

how stacking RNNs together can produce a deep neural network.

directional RNNs

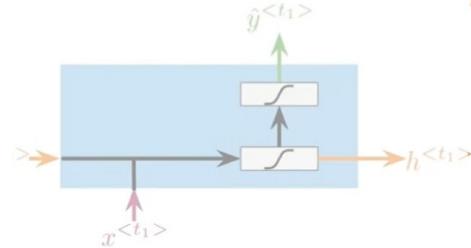


Information flows from the past and from the future independently

So the computations from left to right

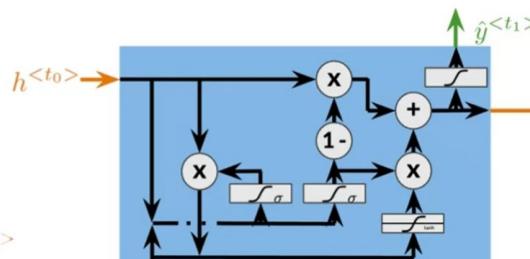
plearning.ai

Vanilla RNN vs GRUs



$$h^{<t_1>} = g(W_h[h^{<t-1>}, x^{<t_1>}] + b_h)$$

$$\hat{y}^{<t_1>} = g(W_y h^{<t_1>} + b_y)$$



$$\Gamma_u = \sigma(W_u[h^{<t_0>}, x^{<t_1>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[h^{<t_0>}, x^{<t_1>}] + b_r)$$

$$h'^{<t_1>} = \tanh(W_h[\Gamma_r * h^{<t_0>}, x^{<t_1>}] + b_h)$$

$$h^{<t_1>} = \Gamma_u * h^{<t_0>} + (1 - \Gamma_u) * h'^{<t_1>}$$

$$\hat{y}^{<t_1>} = g(W_y h^{<t_1>} + b_y)$$

All of these computations allow the network to learn

plearning.ai

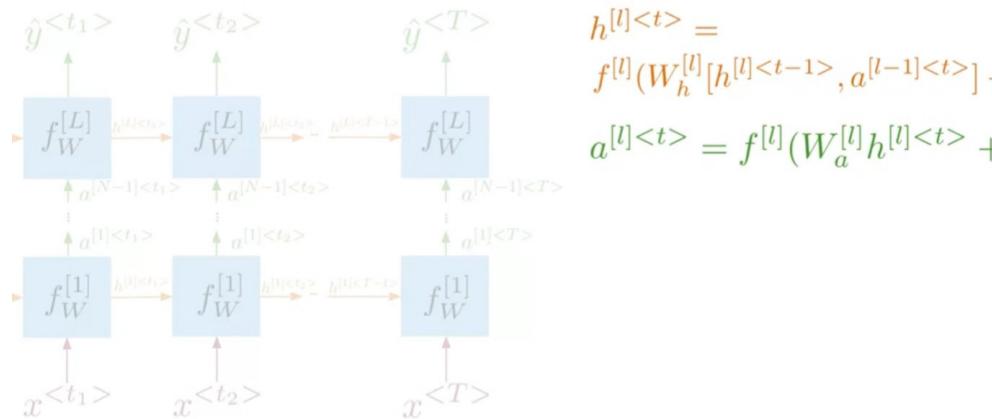
directional RNNs

Press Esc to exit full screen

Bidirectional RNNs work in much the same way that simple RNNs do.

plearning.ai

ep RNNs



plearning.ai

But let's see how information flows in this case.

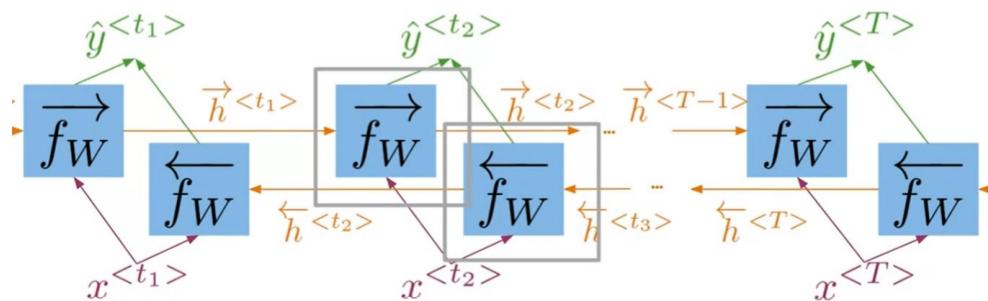
lNs: Advantages

Captures dependencies within a short range

Takes up less RAM than other n-gram models

They're also relatively lightweight compared to other n-gram models,

directional RNNs



Information flows from the past and from the future independently

plearning.ai

When you have computed both of the hidden states for a time step,

tline

RNNs and vanishing/exploding gradients

Solutions

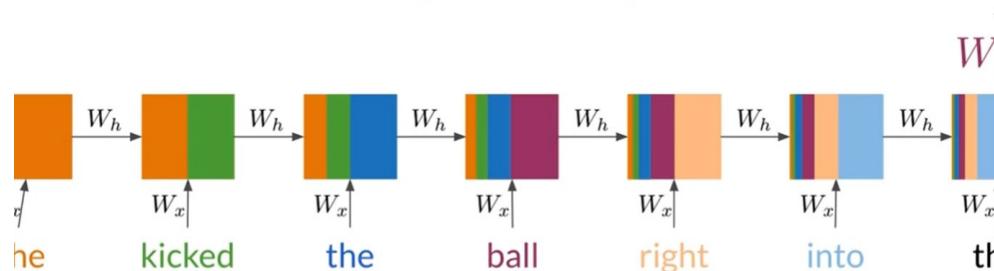


a problem common to RNNs, and then demonstrate a few ways to handle them.

plearning.ai

IN Basic Structure

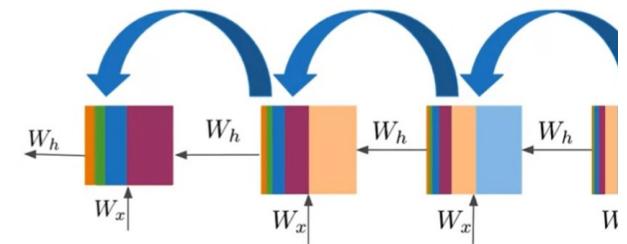
kicked the ball right into the _____



Learnable parameters

propagates information from the beginning
of the sequence through to the end.

e vanishing gradient problem



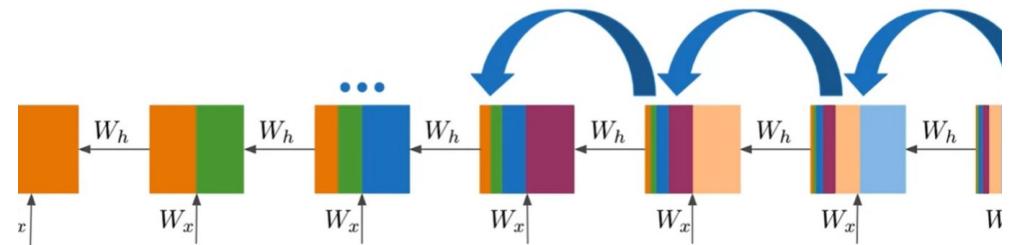
at the early layers as the product of all
the terms from the later layers makes for

INs: Disadvantages

Struggles with longer sequences

Prone to vanishing or exploding gradients

ckpropagation through time



$\int \sigma$ between 0 and 1

$\int \tanh$ between -1 and 1

is part of how vanishing/exploding
gradients are created,

You can think of gradients as a measure

tline

Meet the Long short-term memory unit!

LSTM architecture

Applications



LMs: a memorable solution

Learns when to remember and when to forget

Basic anatomy:

- A cell state
- A hidden state with three gates
- Loops back again at the end of each time step

Gates allow gradients to flow unchanged

along with a few examples of

plearning.ai

these three gates that track inputs as it arrives.

ving for vanishing or exploding gradients

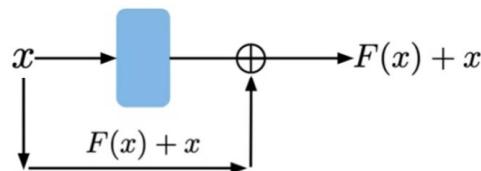
Identity RNN with ReLU activation

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad -1 \longrightarrow 0$$

Gradient clipping

$$32 \longrightarrow 25$$

Skip connections



This effectively skips over the activation functions and

plearning.ai

LMs: a memorable solution

Learns when to remember and when to forget

Basic anatomy:

- A cell state
- A hidden state with three gates
- Loops back again at the end of each time step

training to decide on what changes to make.

plearning.ai

LMs: Based on previous understanding

state = before conversation

get gate = beginning of conversation

it gate = thinking of a response

put gate = responding

ated cell state = after conversation



plearning.ai

since you began
your conversation.

LMs: Based on previous understanding

state = before conversation

get gate = beginning of conversation

it gate = thinking of a response

put gate = responding



plearning.ai

When you decide what to say next,

Applications of LSTMs

t-character

fiction



Chatbots



Music composition



Image
captioning



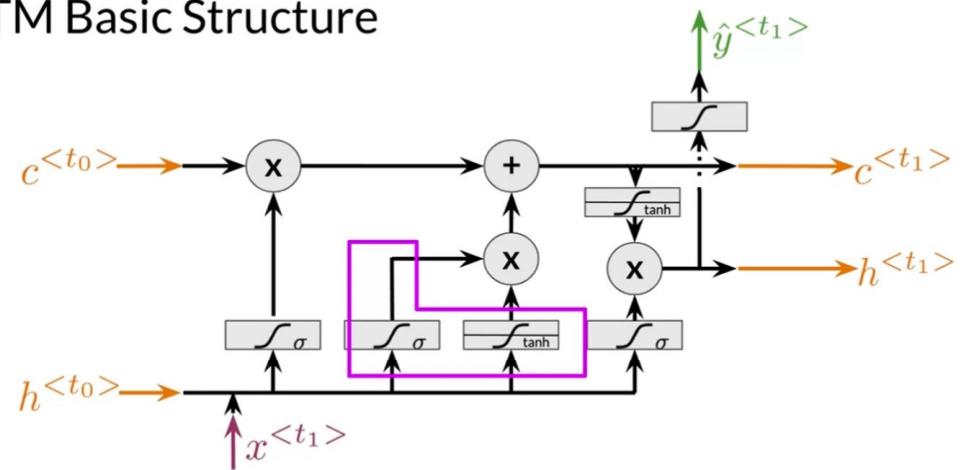
Speech
recognition



plearning.ai

In overcoming the problems
of a traditional RNN,

LM Basic Structure



plearning.ai

a sigmoid layer to choose
which values to update,

Summary

LSTMs use a series of gates to decide which information to keep:

- Forget gate decides what to keep
- Input gate decides what to add
- Output gate decides what the next hidden state will be

One time step is completed after updating the states

plearning.ai

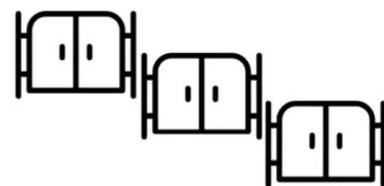
you will have completed
one timestep.

Summary

LSTMs offer a solution to vanishing gradients

Typical LSTMs have a cell and three gates:

- Forget gate
- Input gate
- Output gate



plearning.ai

That typical LSTMs have a cell
or memory and three gates,

Types of Entities



Thailand:
Geographical



Google:
Organization



India:
Geopolitical

plearning.ai or geopolitical
entities like Indian.

What is Named Entity Recognition?

Locates and extracts predefined entities from text

Places, organizations, names, time and dates



They can even be times and dates.

plearning.ai

ample of a labeled sentence



You'll notice that here

plearning.ai

Training NEF Data Processing



deeplearning.ai

Let's go over how you can apply both
to your named entity recognition model.

More Types of Entities



December:
Time Indicator



Egyptian statue:
Artifact



Barack Obama:
Person

Applications of NER systems

Search engine efficiency

Recommendation engines

Customer service

Automatic trading



or a person such as Barack Obama.

plearning.ai

building a web scraper to get

plearning.ai

Processing data for NERs

Assign each class a number

Assign each word a number

Sharon flew to Miami last Friday.

[4282, 853, 187, 5388, 2894, 7]

B-per O O B-geo O B-tim

853 corresponds to flew,

187 corresponds with to, and so forth.

Timeline

Convert words and entity classes into arrays

Token padding

Create a data generator



label data into arrays of numbers
that correspond to one another.

plearning.ai

Training the NER

Create a tensor for each input and its corresponding number

Put them in a batch → 64, 128, 256, 512 ...

Feed it into an LSTM unit

Run the output through a dense layer

Predict using a log softmax over K classes

a tensor for each input and its
corresponding label, as you saw before.

Token padding

LSTMs, all sequences need to be the same size.

Set sequence length to a certain number

Use the <PAD> token to fill empty spaces

You'll be implementing
this in the assignments.

plearning.ai

Summary

Convert words and entities into same-length numerical arrays

Train in batches for faster processing

Run the output through a final layer and activation



Evaluating the model

Pass test set through the model

Get arg max across the prediction array

Mask padded tokens

Compare outputs against test labels

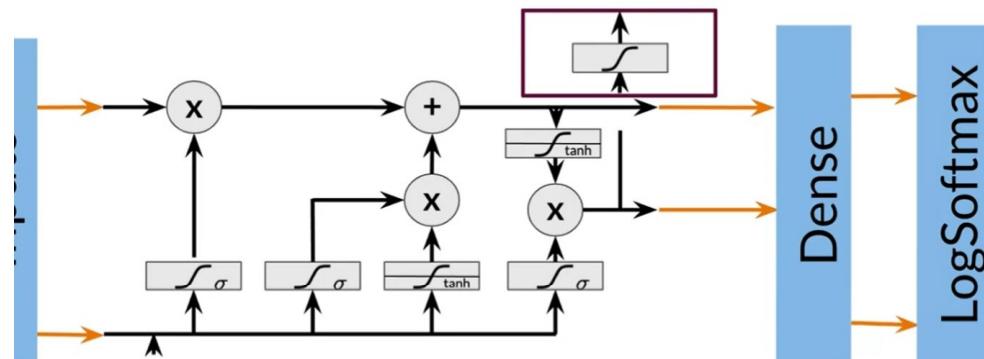
deeplearning.ai

422 like numerical array conversion with token padding.

deeplearning.ai

0:29 compare the outputs against your test labels to see how accurate your model is.

Training the NER



deeplearning.ai And then you compute a log softmax to get the corresponding outputs.

Press Esc to exit full screen

You're almost ready to use your new model.

Summary

If padding tokens, remember to mask them when computing accuracy
Coding assignment!

plearning.ai

Evaluating the model in Python

```
def evaluate_model(test_sentences, test_labels, model):
    pred = model(test_sentences)
    outputs = np.argmax(pred, axis=2)
    mask = ...
    accuracy = np.sum(outputs==test_labels)/float(np.sum(mask))

    return accuracy
```

and dividing it by the sum
of all your mask tokens.

plearning.ai

What do Siamese Networks learn?

I am happy because I am learning



Classification: categorize things

Siamese Networks: Identify similarity between things



What is your age?

How old are you?

When you build a Siamese model, you're

plearning.ai trying to identify the difference or

Question Duplicates

How old are you? = What is your age?

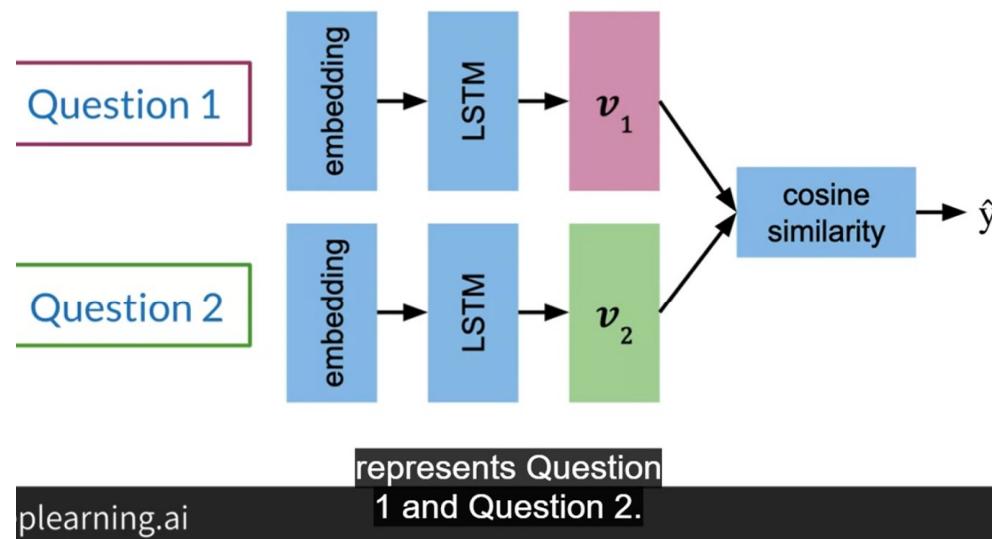
Where are you from?

≠

Where are you going?

This example shows that comparing meaning
plearning.ai is not as simple as just comparing words.

Model Architecture



Loss Function

How old are you?

Anchor

$$\cos(v_1, v_2) = \frac{v_1 \cdot v_2}{\|v_1\| \|v_2\|}$$

$s(v_1, v_2)$

What is your age?

Positive

$s(A, P)$

Where are you from? Negative

Challenges in NLP



Handwritten checks

What is your age?
How old are you?



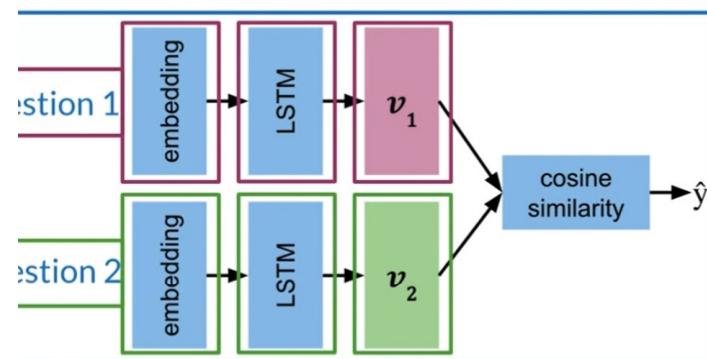
Queries



Question duplicates

to predict whether a new query is similar
to the one that was already executed.

Model Architecture



- 1) Inputs
- 2) Embedding
- 3) LSTM
- 4) Vectors

Then you take the outputs of each of the sub networks and

plets

How old are you?

What is your age?

Where are you from?

A
P
N

Triplets !!!



plet Selection

Hard triplets are better for training !

Triplet A, P, N:

duplicate set: A, P
non-duplicate set: A, N

Random:

$$\mathcal{L} = \max(\text{diff} + \alpha, 0)$$

$$\text{diff} = s(A, N) - s(A, P)$$

Easy to satisfy. Little to learn



plearning.ai

what gives rise to
the name triplets,

plearning.ai

Instead of selecting
random triplets,

s Function

How old are you?

Anchor

$$\cos(v_1, v_2) = \frac{v_1 \cdot v_2}{\|v_1\| \|v_2\|}$$

$s(v_1, v_2)$

What is your age?

Positive

$$s(A, P) \approx 1$$

Where are you from? Negative

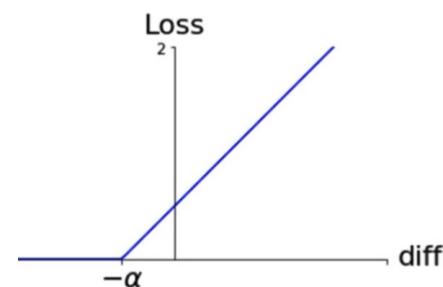
$$s(A, N) \approx -1$$

$$\text{diff} = s(A, N) - s(A, P)$$

and P to calculate
the difference.

plearning.ai

plet Loss



Simple loss:
 $\text{diff} = s(A, N) - s(A, P)$

Non linearity:

$$\mathcal{L} = \begin{cases} 0; & \text{if } \text{diff} \leq 0 \\ \text{diff}; & \text{if } \text{diff} > 0 \end{cases}$$

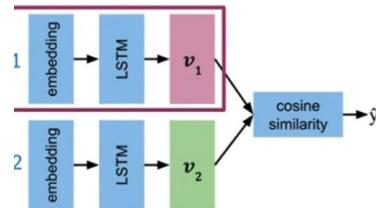
Alpha margin:

$$\mathcal{L} = \begin{cases} 0; & \text{if } \text{diff} + \alpha \leq 0 \\ \text{diff} + \alpha; & \text{if } \text{diff} + \alpha > 0 \end{cases}$$

improve and learn from
this training example.

plearning.ai

Computing The Cost



Batch 1
What is your age?
Can you see me?
Where are thou?
When is the game?

v_{1_1}			
v_{1_2}			
v_{1_3}			
v_{1_4}			

the vector outputs for each question in the batch.

plearning.ai

Computing The Cost

$$s(v_1, v_2)$$

v_1		v_2	
-1	-2	-3	-4
0.9	-0.8	0.3	-0.5
-0.8	0.5	0.1	-0.2
0.3	0.1	0.7	-0.8
-0.5	-0.2	-0.8	1.0

ranges from negative 1 to positive 1,

4:21

Computing The Cost

Prepare the batches as follows:



What is your age?	How old are you?
Can you see me?	Are you seeing me?
Where are thou?	Where are you?
When is the game?	What time is the game?

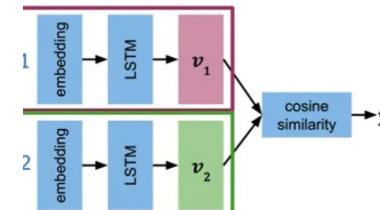
b = 4

none of the rows in those

plearning.ai

column contain a sentence that

Computing The Cost



Batch 1

What is your age?

Can you see me?

Where are thou?

When is the game?

Batch 2

How old are you?

Are you seeing me?

Where are you?

What time is the game?

v_{1_1}			
v_{1_2}			
v_{1_3}			
v_{1_4}			

v_{2_1}			
v_{2_2}			
v_{2_3}			
v_{2_4}			

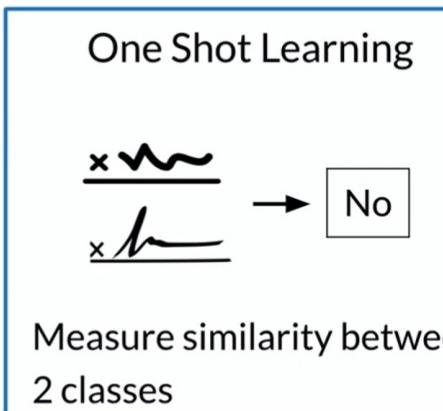
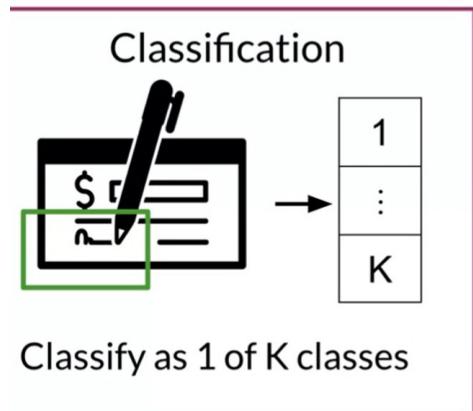
Here, for example, v_{1_1}

is a duplicate of v_{2_1} ,

3:18

plearning.ai

Classification vs One Shot Learning



Then you can test a similarity score against

plearning.ai

Prepare Batches

Question 1:
batch size b

Batch 1
What is your age?
Can you see me?
Where are thou?
When is the game?

Question 2:
batch size b

Batch 2
How old are you?
Are you seeing me?
Where are you?
What time is the game?

The second question in Batch 1 is

plearning.ai

Computing The Cost

		v_1	v_2	
v_1	-1	-2	-3	-4
0.9	-0.8	0.3	-0.5	
-0.8	0.5	0.1	-0.2	
0.3	0.1	0.7	-0.8	
-0.5	-0.2	-0.8	1.0	

$$\mathcal{L}(A, P, N) = \max(\text{diff} + \alpha, 0)$$

$$\text{diff} = s(A, N) - s(A, P)$$

$$\mathcal{J} = \sum_{i=1}^m \mathcal{L}(A^{(i)}, P^{(i)}, N^{(i)})$$

but there are more techniques available that's can

plearning.ai

524

One Shot Learning

No need for retraining !



Learn a similarity score

a similarity function that can be

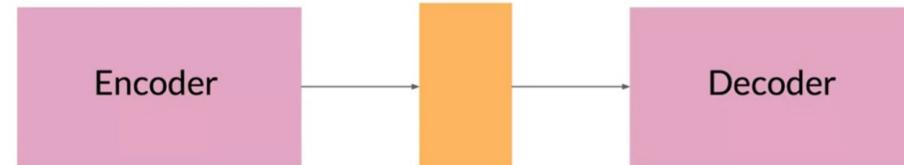
plearning.ai

tline

Introduction to Neural Machine Translation
Seq2Seq model and its shortcomings
Solution for the information bottleneck



Seq2Seq model



plearning.ai

To get started on this weeks material,
I'll introduce you to neural machine

How are you today?

Wie geht es Ihnen heute

ntearning.ai 4:58
It takes in a hidden states and a string
of words, such as a single sentence.

Prepare Batches

Question 1:
batch size b

- {
 - Batch 1
 - What is your age?
 - Can you see me?
 - Where are thou?
 - When is the game?
 - Batch 2
 - How old are you?
 - Are you seeing me?
 - Where are you?
 - What time is the game?

Question 2:
batch size b

to get outputs vectors
for each batch.

$v_1 = (1, d_{\text{mod}})$
v_{1_1}
v_{1_2}
v_{1_3}
v_{1_4}
v_2
v_{2_1}
v_{2_2}
v_{2_3}
v_{2_4}

Seq2Seq model

Introduced by Google in 2014

Maps variable-length sequences to fixed-length memory

LSTMs and GRUs are typically used to overcome the
vanishing gradient problem



ntearning.ai

GRUs are typically used to
avoid vanishing gradients.

Seq2Seq shortcomings

Variable-length sentences + fixed-length memory =



As sequence size increases, model performance decreases

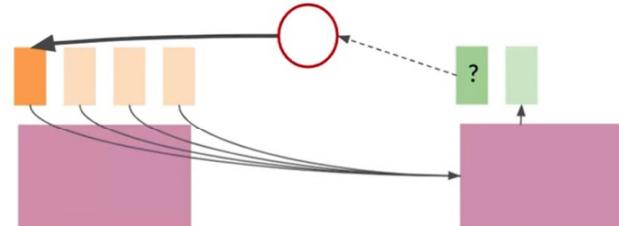
This results in lower model performance
as sequence size increases,

plearning.ai

ution: focus attention in the right place

Prevent sequence overload by giving the model a way to focus on the **likeliest** words at each step

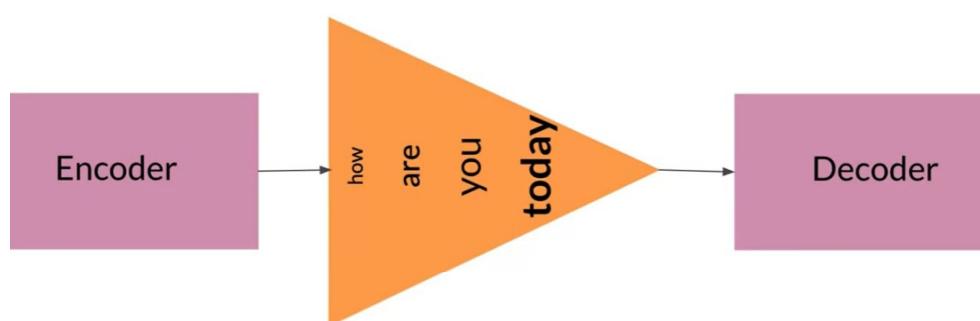
Do this by providing the information specific to each input word



Up next you'll get a conceptual idea of
what this new layer is doing and why.

ntearning.ai

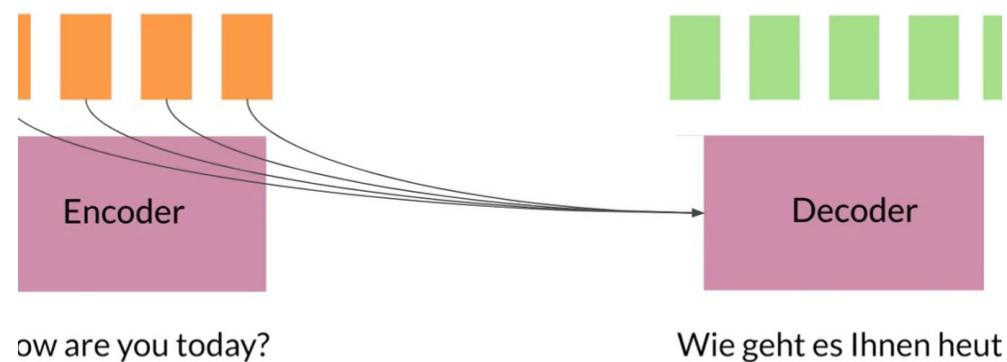
The information bottleneck



are given more importance.

plearning.ai

One vector per word



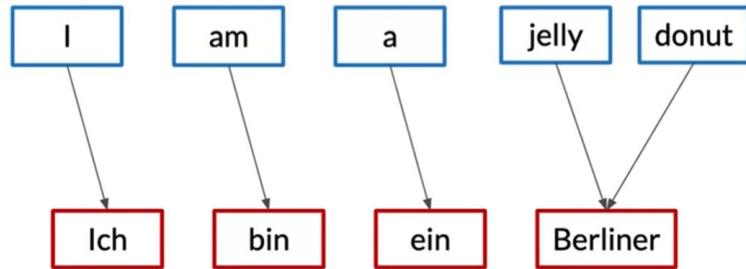
Wie geht es Ihnen heut?

ntearning.ai

memory efficient model that predicts
accurately from a long sequence?

4:09

Word alignment



an optional reading about this historic

plearning.ai misunderstanding in the course material.

Motivation for alignment

Correctly aligned words are the goal:

Translating from one language to another

Word sense discovery and disambiguation

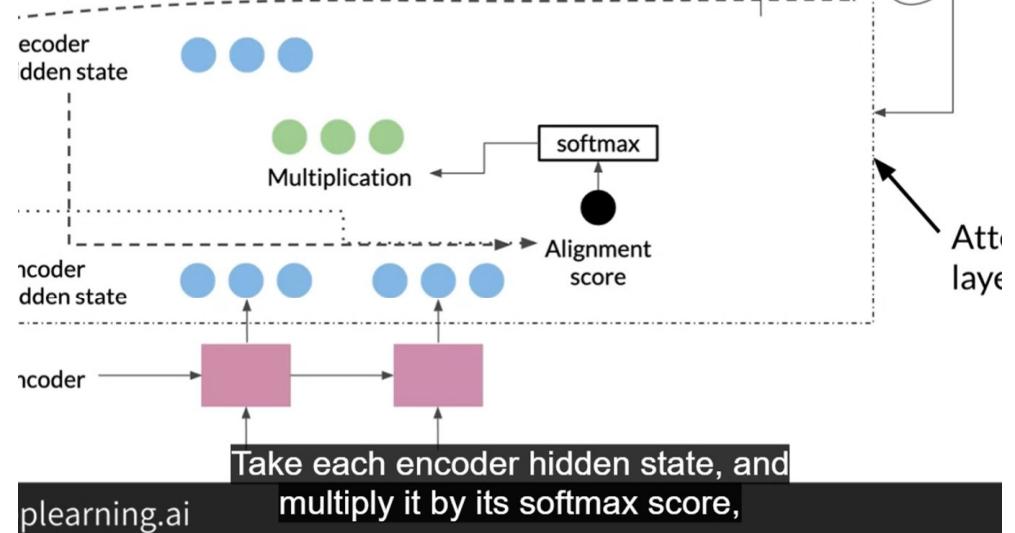


Achieve alignment with a system for retrieving information step by step and scoring it

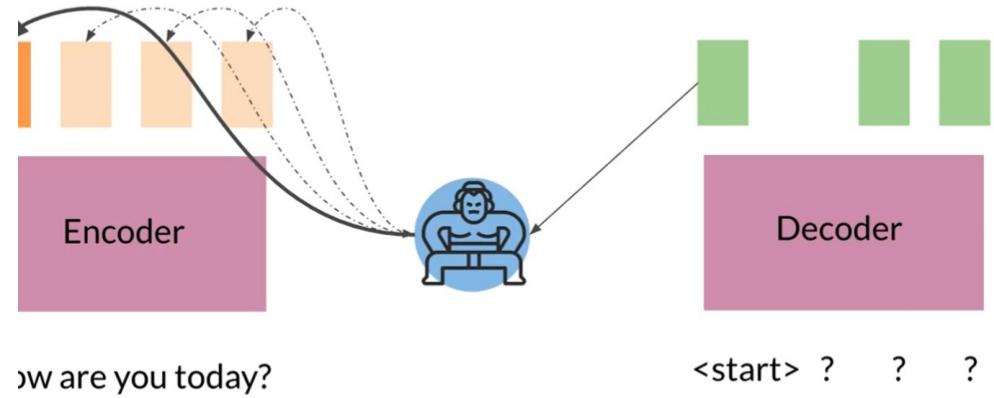
For this week's assignments on neural machine translation with attention,

plearning.ai

Calculating alignment



Give some inputs more weight!



You can see here that the hidden state shown in dark orange has a heavier line,

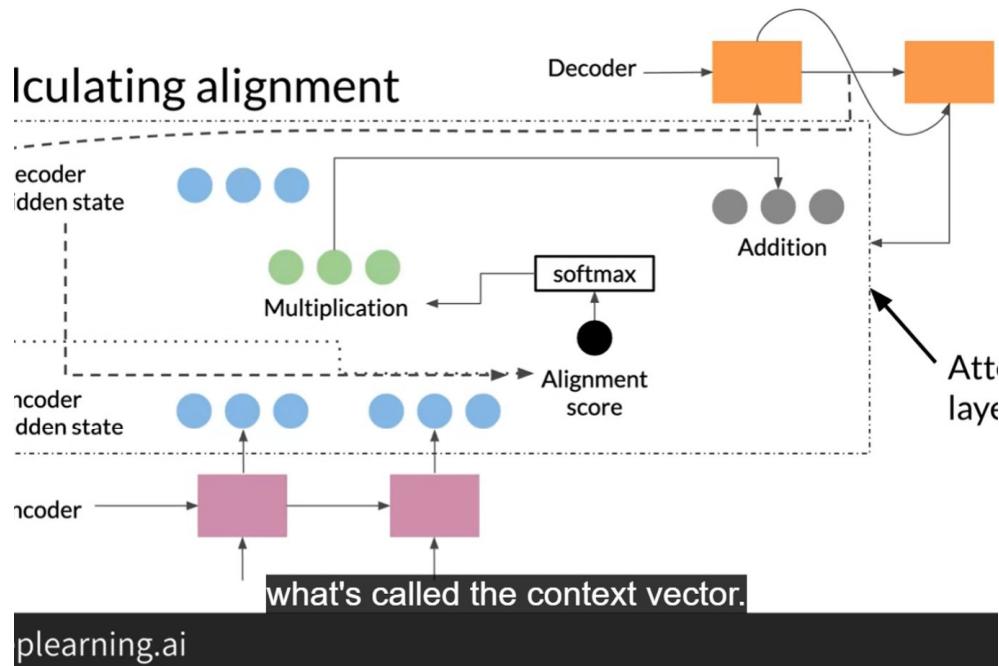
plearning.ai

Data in machine translation

English	German
I am hungry!	Ich habe Hunger.
...	...
I watched the soccer game.	Ich habe das Fußballspiel gesehen.

plearning.ai You're going to have a great many of these inputs.

Calculating alignment



Machine translation setup

State-of-the-art uses pre-trained vectors otherwise, represent words with a one-hot vector to create the input p track of index mappings with word2ind and ind2word dictionaries start-of and end-of sequence tokens:



plearning.ai a starts of sequence token like this,

Machine translation setup

State-of-the-art uses pre-trained vectors otherwise, represent words with a one-hot vector to create the input p track of index mappings with word2ind and ind2word dictionaries



plearning.ai

Preparing to Translate to German

ENGLISH SENTENCE:

In the ballpoint and the mechanical pencil in the series are equipped with a special mechanism: when the twist mechanism is activated, the lead is pushed forward.

TOKENIZED VERSION OF THE ENGLISH SENTENCE:

```
46 4 11358 362 8 4 23326 20104 1745 8210 9641 5 6  
103 31 2767 30 13 914 4797 64 196 4 22474 5 4797 16  
54 86 2 4 1060 16 6413 1138 3 [1 0 0 0 0 0 0 0 0 0 0]  
) 0 0 0 0 0 0 0 0 0 0 0 0 0 0]
```

Then a series of zeros
that's if you were to count,

plearning.ai

line

Teacher forcing

Model for NMT with attention



along with being introduced to

plearning.ai

Preparing to Translate to German

ENGLISH SENTENCE:

In the ballpoint and the mechanical pencil in the series are equipped with a special mechanism: when the twist mechanism is activated, the lead is pushed forward.

TOKENIZED VERSION OF THE ENGLISH SENTENCE:

```
46 4 11358 362 8 4 23326 20104 1745 8210 9641 5 6  
103 31 2767 30 13 914 4797 64 196 4 22474 5 4797 16  
54 86 2 4 1060 16 6413 1138 3 1 0 0 0 0 0 0 0 0 0  
) 0 0 0 0 0 0 0 0 0 0 0 0 0 0]
```

the tokenized version of
the English sentence,

plearning.ai

English to German

GERMAN TRANSLATION:

Kugelschreiber und der Drehbleistift der Serie sind mit einem besonderen Mechanismus ausgestattet: Bei Betätigung der Drehmechanik wird die Schreibmine nach vorne geschoben.

TOKENIZED VERSION OF THE GERMAN TRANSLATION:

```
9 3892 5280 14774 2418 12 11 9883 6959 7298 15157 5 11 845  
114 5324 10565 2520 64 752 12954 26538 147 11 9883 23326 201  
78 10 21150 10166 126 14566 5 23850 1171 3 [1 0 0 0 0]  
) 0 0 0 0 0 0 0 0 0 0 0 0 0 0
```

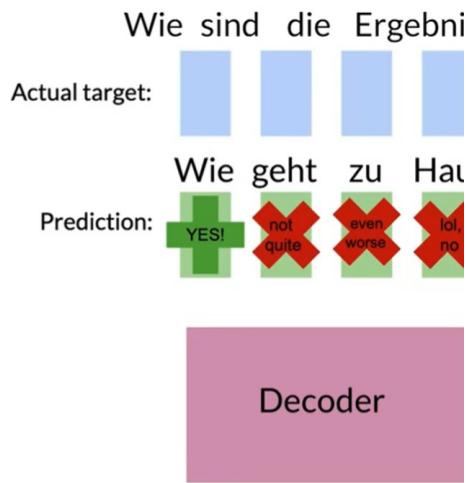
of unpaired tokens
in the sequence.

plearning.ai

Teacher forcing: motivation



How are the results?



the prediction made at each step.

nterlearning.ai

How to know predictions are correct?

Teacher forcing allows the model to “check its work” at each step

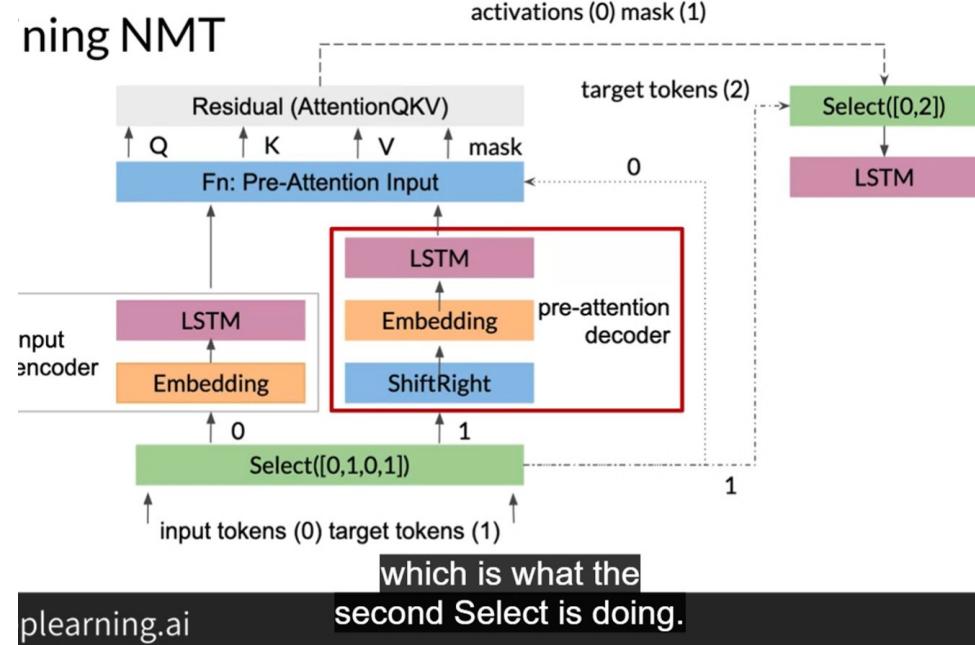
Compare its prediction against the real output during training

Result: Faster, more accurate training

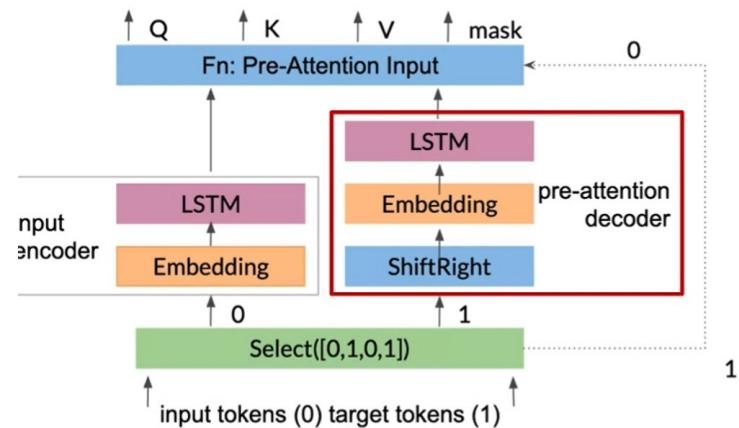


This yields faster training with

plearning.ai



ning NMT



you can prepare them for the attention layer.

plearning.ai

EU Score

Indicators for Bilingual Evaluation Understudy

Measures the quality of machine-translated text by comparing “candidate” to one or more “reference” translations.

Score: the closer to 1, the better, and vice versa:



The closer to zero,
the worse it is.

plearning.ai

EU score is great, but...

sider the following:

BLEU doesn't consider semantic meaning

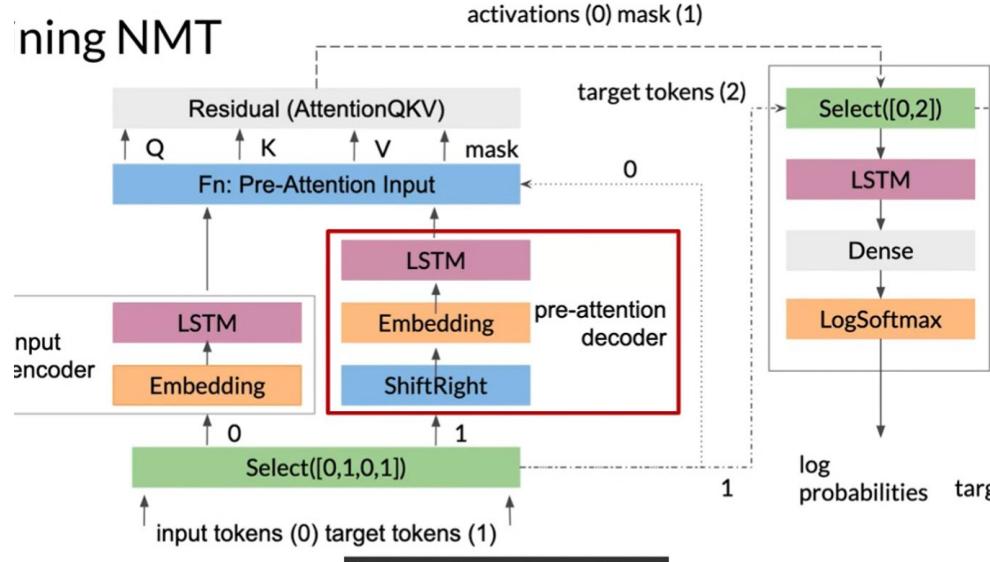
BLEU doesn't consider sentence structure:

“Ate I was hungry because!”



It also doesn't consider the
structure of the sentence.

ning NMT



EU Score

Candidate	I	I	am	I	I
Reference 1	Younes	said	I	am	hungry
Reference 2	He	said	I	am	hungry

“I” appears at most once in both, so clip to one: $m_w = 1$

(Sum over unique n-gram counts in the candidate)

(total # of words in candidate)

the unique uni-gram counts in the candidates.

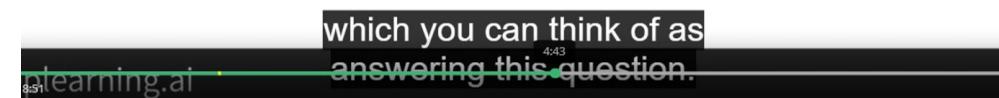
plearning.ai

ROUGE evaluation

Model	The	cat	had	striped	orange	fur
Reference	The	cat	had	orange	fur	

Recall = How much of the reference text is the system text capturing?

Precision = How much of the model text was relevant?

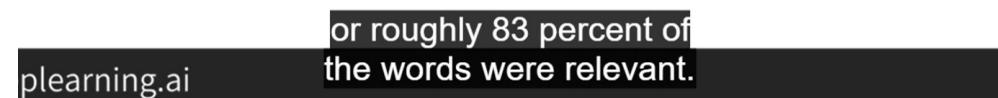


Precision in ROUGE

Model	The	cat	had	striped	orange	fur
Reference	The	cat	had	orange	fur	

Precision = $\frac{\text{Sum of overlapping unigrams in model and reference)}}{\text{total # of words in model}}$

$$\frac{5}{6} = 0.83$$

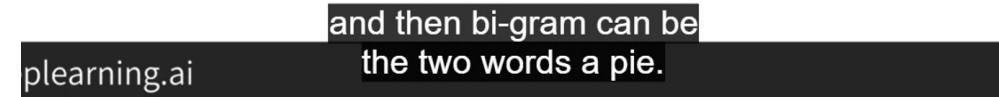


ROUGE

all-Oriented Understudy for Gisting Evaluation

evaluates quality of machine text

measures precision and recall between generated text and human-created text

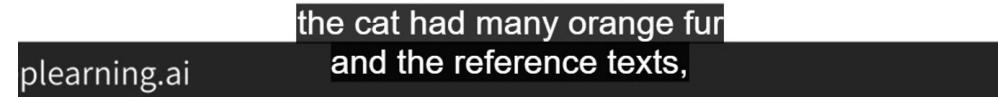


Recall in ROUGE

Model	The	cat	had	striped	orange	fur
Reference	The	cat	had	orange	fur	

Recall = $\frac{\text{Sum of overlapping unigrams in model and reference)}}{\text{total # of words in reference}}$

$$\frac{5}{5} = 1.0$$

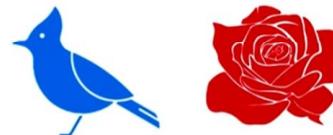


Summary

BLEU score compares “candidate” against “references” using an n-gram average

BLEU doesn’t consider meaning or structure

ROUGE measures machine-generated text against an “ideal” reference



which was created to evaluate

plearning.ai

Greedy decoding

chooses the most probable word at each step

the best word at each step may not be the best for longer sequences

Ich habe Hunger.

I am hungry.

I am, am, am, am...

Problems in ROUGE

Doesn’t take themes or concepts into consideration (i.e., a low ROUGE score doesn’t necessarily mean the translation is bad)

Model	I	am	a	fruit-filled	pastry
Reference	I	am	a	jelly	donut



you start taking your
own ROUGE scores.

plearning.ai

Sampling

Random sampling

Temperature in sampling

Greedy decoding

Beam search

Minimum Bayes’ risk (MBR)



of this week’s lectures.
That’s awesome.

plearning.ai

Temperature

Sampling, temperature is a parameter allowing for more or less randomness in predictions
lower temperature setting = More confident, conservative network
higher temperature setting = More excited, random network (and more mistakes)



plearning.ai

and you may get some pretty fun predictions.

Random sampling

am	full	hungry	I	the
0.05	0.3	0.15	0.25	0.25

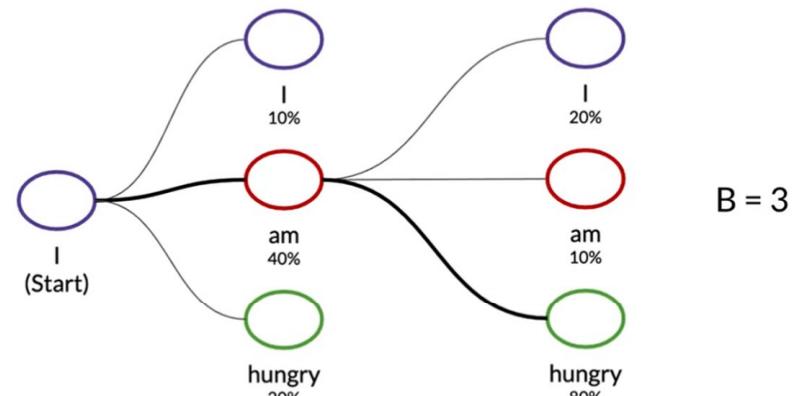
can be a little too random for accurate translation!

Action: Assign more weight to more probable words, and less weight to less probable words.

and less weight to the others.

plearning.ai

Beam search example



plearning.ai

This is a very simple example but you can

Beam search decoding

order, more exploratory decoding alternative
selects multiple options for the best input based on conditional probability
number of options depends on a predetermined beam width parameter
selects B number of best alternatives at each time step



to select several possible options.

plearning.ai

Problems with beam search

the model learns a distribution, that tends to carry more weight than the tokens

cause translation problems, i.e. in a speech corpus that hasn't been cleaned



Prediction:
"Umm uhh
ummm huh?"

It can also use single tokens in a problematic way,

plearning.ai

Problems with beam search

the model learns a distribution, that tends to carry more weight than the tokens



Prediction:
"Umm uhh
ummm huh?"

However, beam search decoding runs into issues where

plearning.ai

Problems with beam search

"mag die Vereinigten Staaten, weil die Vereinigten Staaten groß sind."

with 11 good English translations of "Vereinigten Staaten," but a ~1/11th probability of the non-word "Uhm" occurring, you might get this as a translation:

the United States, because the Uhm is big."

with 11^2 good translations, the most probable one will still be "U"

and said because one over 11 squared is less than

plearning.ai

Problems with beam search

"mag die Vereinigten Staaten, weil die Vereinigten Staaten groß sind."

with 11 good English translations of "Vereinigten Staaten," but a ~1/11th probability of the non-word "Uhm" occurring, you might get this as a translation:

the United States, because the Uhm is big."

These could be USA, US,

plearning.ai

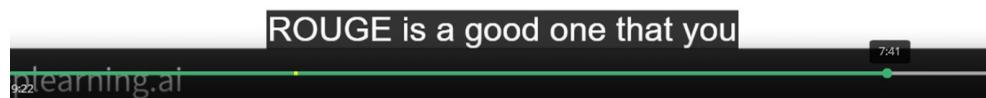
nimum Bayes Risk (MBR)

Compares many samples against one another. To implement MBR:

Generate several random samples

Compare each sample against all the others and assign a similarity score (such as ROUGE!)

Select the sample with the highest similarity: **the golden one** ✨



nimum Bayes Risk (MBR)

Compares many samples against one another. To implement MBR:

Generate several random samples

Compare each sample against all the others and assign a similarity score (such as ROUGE!)

Select the sample with the highest similarity: **the golden one** ✨

which is sometimes referred
to as the golden one.

plearning.ai

ample: MBR Sampling

Generate the scores for 4 samples:

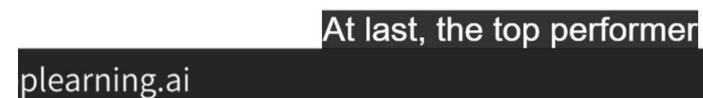
Calculate similarity score between sample 1 and sample 2

Calculate similarity score between sample 1 and sample 3

Calculate similarity score between sample 1 and sample 4

Average the score of the first 3 steps (Usually a weighted average)

Repeat until all samples have overall scores



ample: MBR Sampling

Generate the scores for 4 samples:

Calculate similarity score between sample 1 and sample 2

Calculate similarity score between sample 1 and sample 3

Calculate similarity score between sample 1 and sample 4

Average the score of the first 3 steps (Usually a weighted average)

Note, it's more common to use

plearning.ai

Timeline

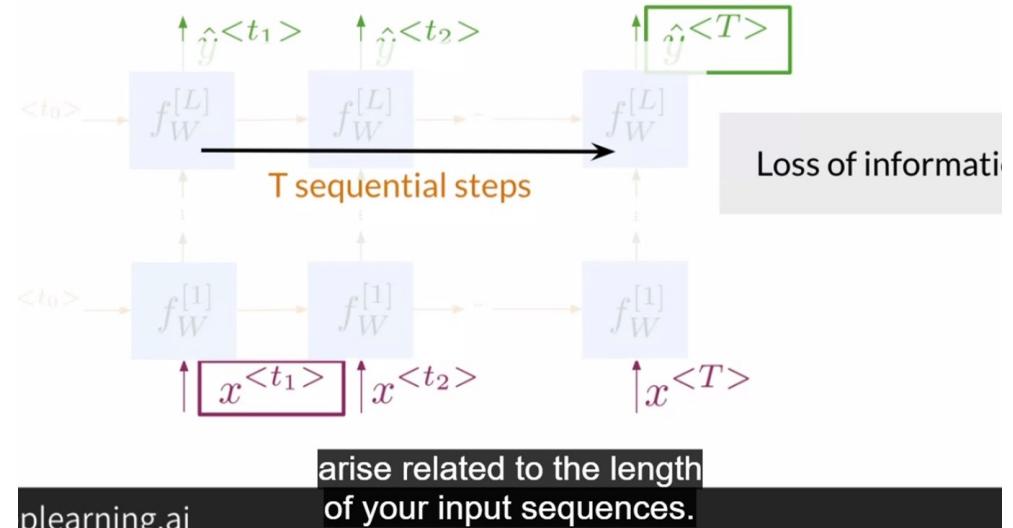
Issues with RNNs

Comparison with Transformers



recurrent neural networks using
plearning.ai

RNN Architectures



Summary

Beam search uses conditional probabilities and the beam width parameter

MBR (Minimum Bayes Risk) takes several samples and compares them against each other to find the golden one ✨

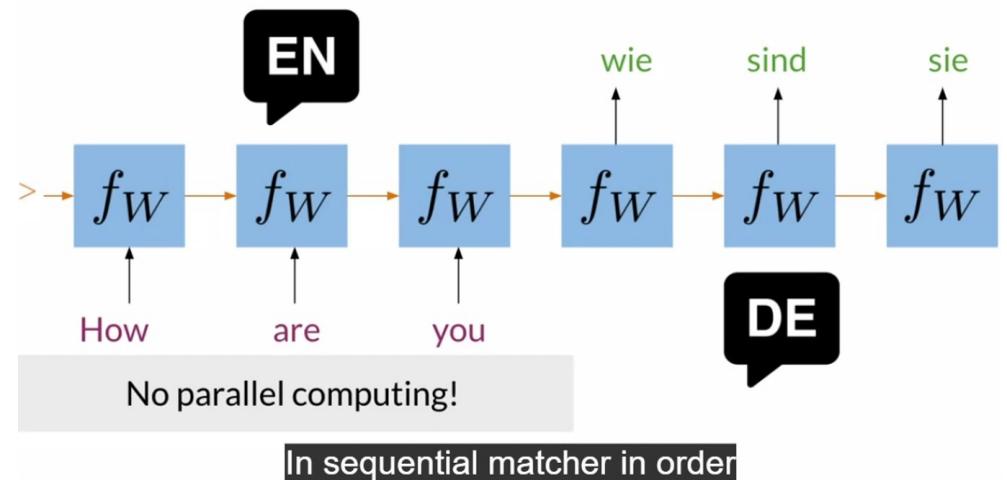
Go forth to the coding assignment!



then chooses the best performer.

plearning.ai

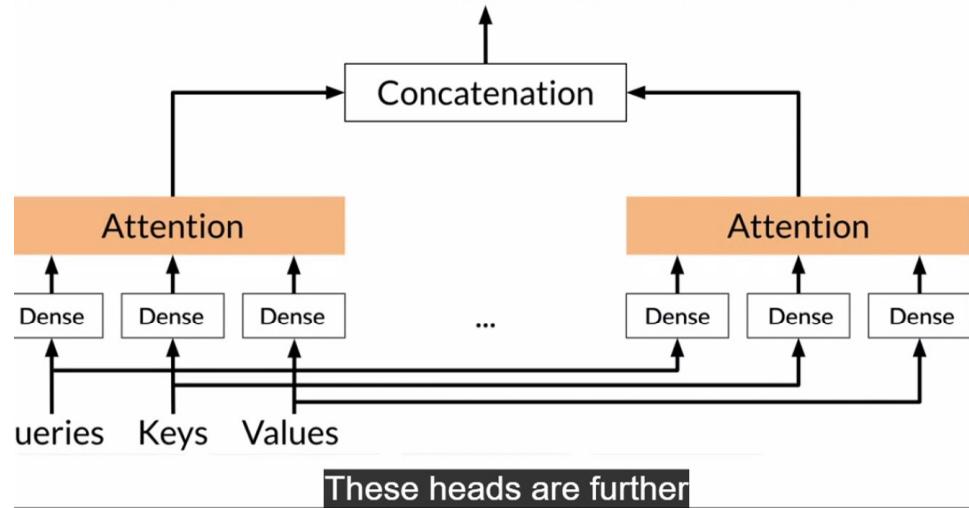
Sequential Machine Translation



In sequential matcher in order to start the translation,

plearning.ai

RNNs vs Transformer: Multi-headed attention



plearning.ai

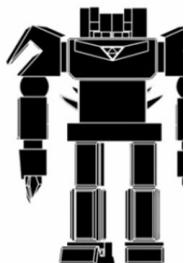
Summary

In RNNs parallel computing is difficult to implement

For long sequences in RNNs there is loss of information

In RNNs there is the problem of vanishing gradient

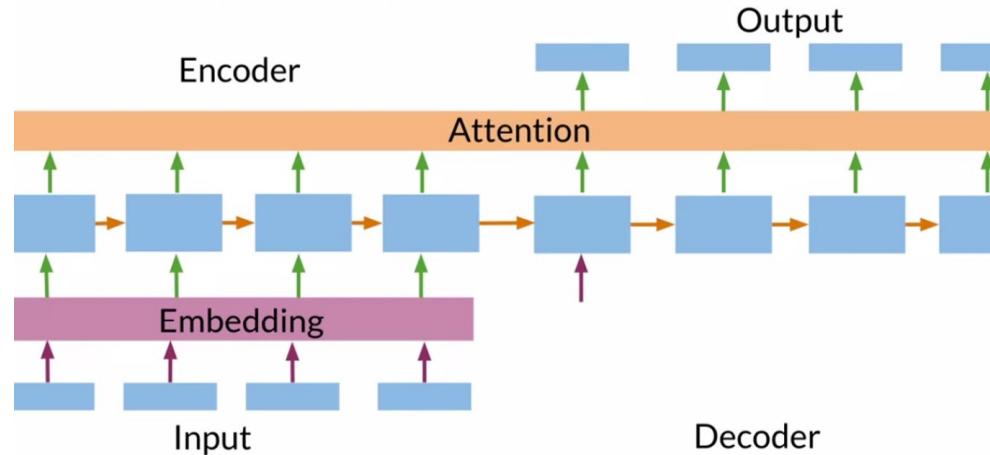
Transformers help with all of the above



With RNNs, it is hard to fully

plearning.ai

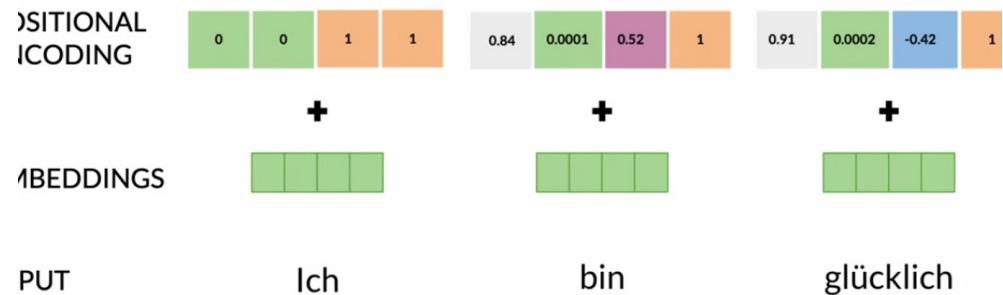
RNNs vs Transformer: Encoder-Decoder



For RNNs, the number of steps
is equal to T. Finally,

plearning.ai

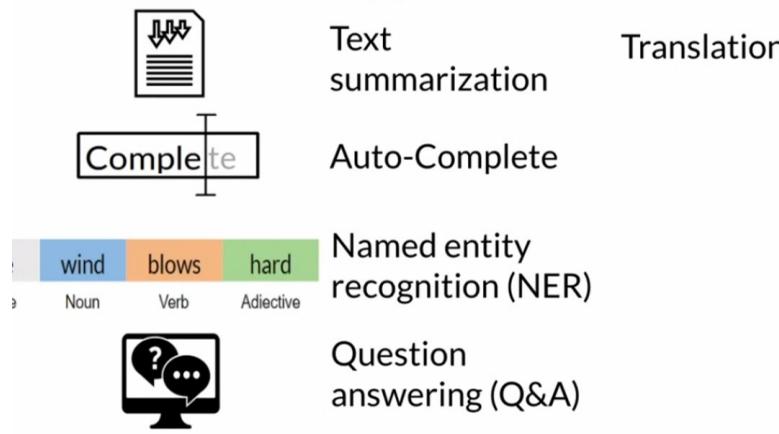
RNNs vs Transformer: Positional Encoding



the multi-head attention
layer computes

plearning.ai

transformer NLP applications



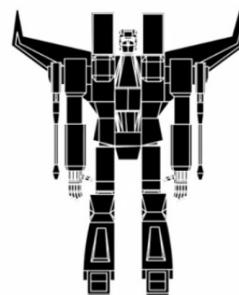
automatic question answering,
machine translation.

Outline

Transformers applications in NLP

Some Transformers

Introduction to T5



First, I'll mention the most popular
applications of Transformers in NLP.

State of the Art Transformers

Iford, A., et al. (2018)

Open AI

Zhilin, J., et al. (2018)

Google AI Language

Szilin, R., et al. (2019)

Google

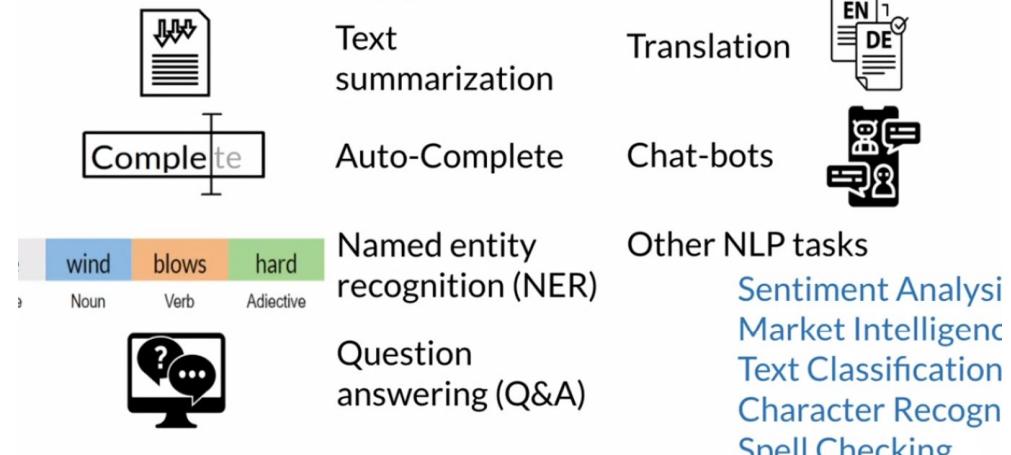
GPT-2: Generative Pre-training for
Transformer

BERT: Bidirectional Encoder
Representations from Transformers

T5: Text-to-text transfer transformer

was also created by Google.

transformer NLP applications



Another application is Chat-bots, and many

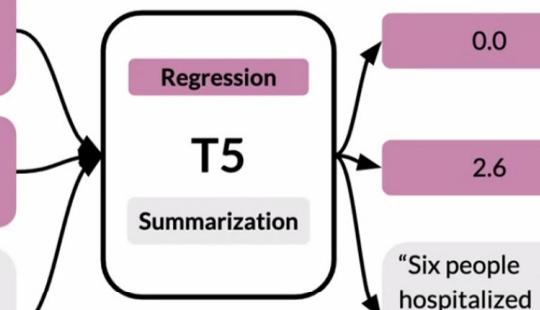
other NLP tasks like Sentiment Analysis,

5: Text-To-Text Transfer Transformer

tsb sentence1: "Cats and dogs are mammals." Sentence2: "There are four known forces in nature - gravity, electromagnetic, weak and strong."

tsb sentence1: "Cats and dogs are mammals." Sentence2: "Cats, dogs, and cows are domesticated."

ummarize: "State authorities dispatched emergency crews Tuesday to survey the damage after an onslaught of severe weather in mississippi..."



six people hospitalized after a storm in Attala County.

plearning.ai

0.0

2.6

"Six people hospitalized after a storm in Attala county"

Outline

- Introducing attention (Translation example)
- Mathematics behind Attention



First I'll introduce the concept of attention using an example of

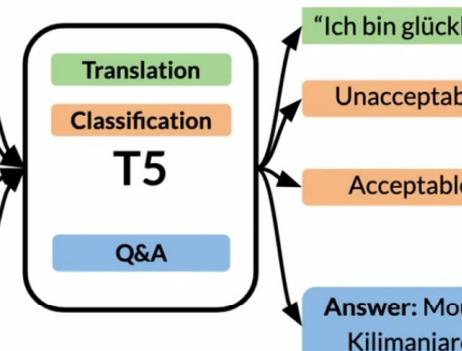
5: Text-To-Text Transfer Transformer

translate English into German: "I am happy"

ola sentence: "He bought fruits and."

ola sentence: "He bought fruits and vegetables."

question: Which volcano in Tanzania is the highest mountain in Africa?



And remember, that's all of these tasks are done by the same model with no

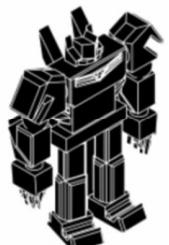
plearning.ai

summary

Transformers are suitable for a wide range of NLP applications

GPT-2, BERT and T5 are the cutting-edge Transformers

T5 is a powerful multi-task transformer



In this video you saw what

are the Transformers applications in NLP,

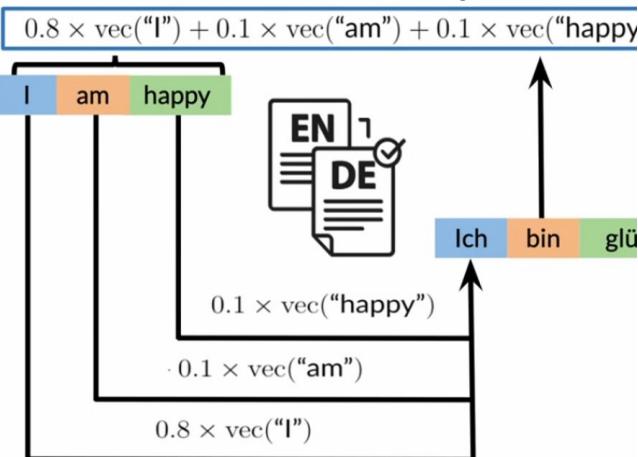
plearning.ai

Introducing attention - Translation example

query (German word) looks for similar words (English words).

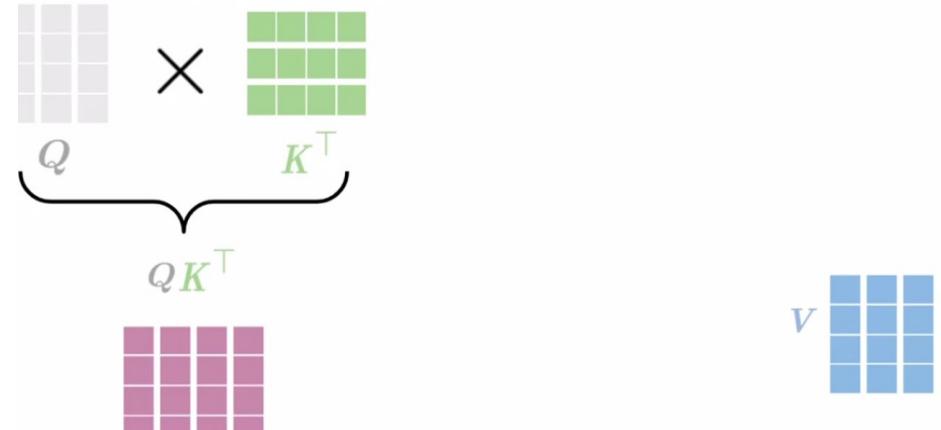
“Ich” key has a probability of being a match for the query.

return sum of the words weighted by their probabilities.



In the end, the query gets the sum of the key vectors weighted by the probabilities.

concept of attention



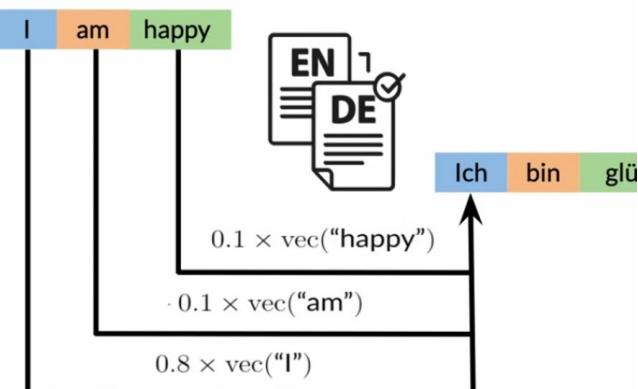
attention model calculates weights or scores representing

Introducing attention - Translation example

query (German word) looks for similar words (English words).

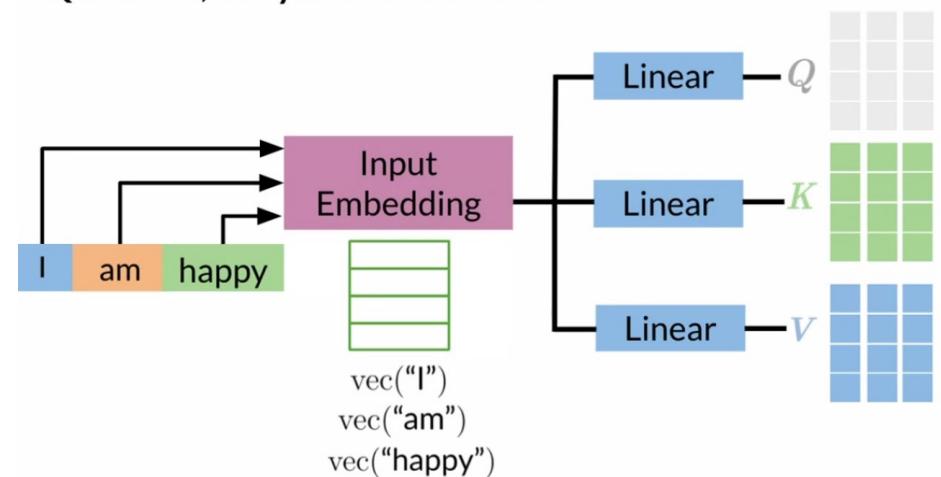
“Ich” key has a probability of being a match for the query.

return sum of the words weighted by their probabilities.



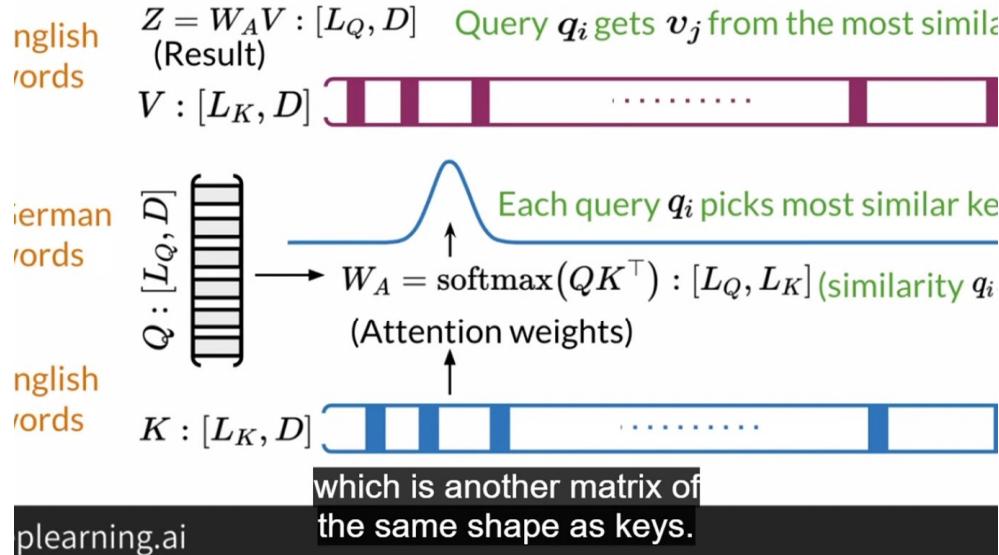
get zero points its probability,
while it's not very similar to am or

Queries, Keys and Values

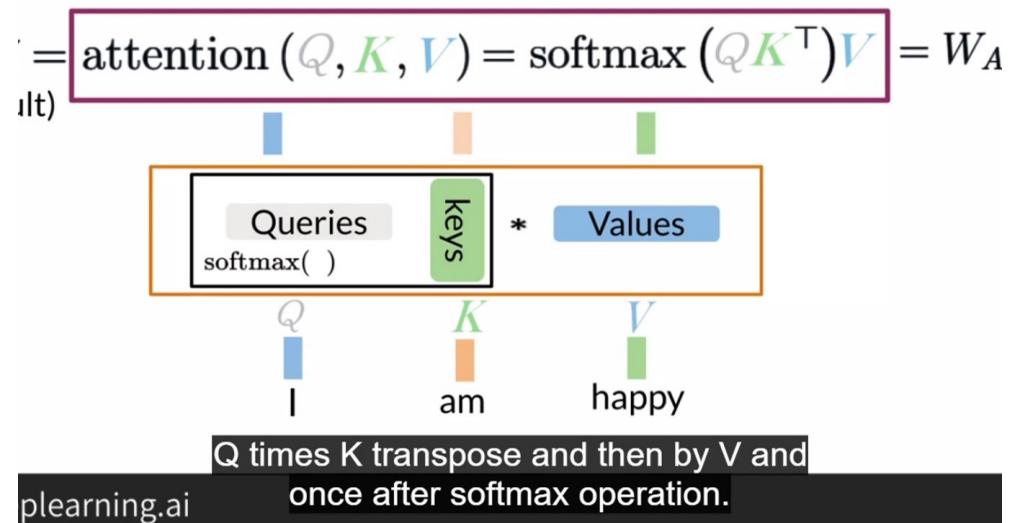


Then you can do the same for
the word am to output a second vector.

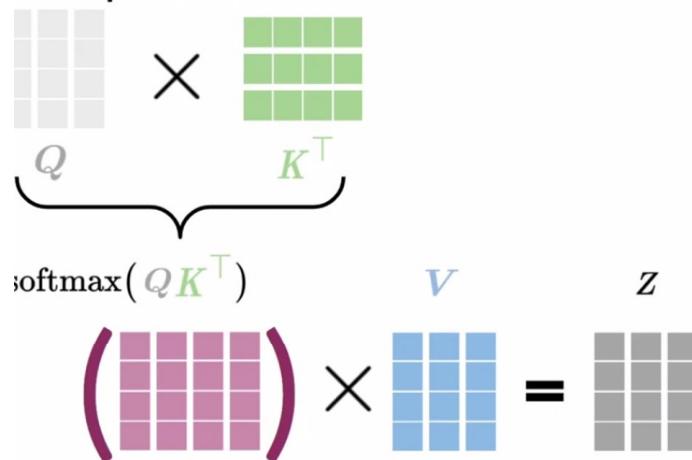
Attention math



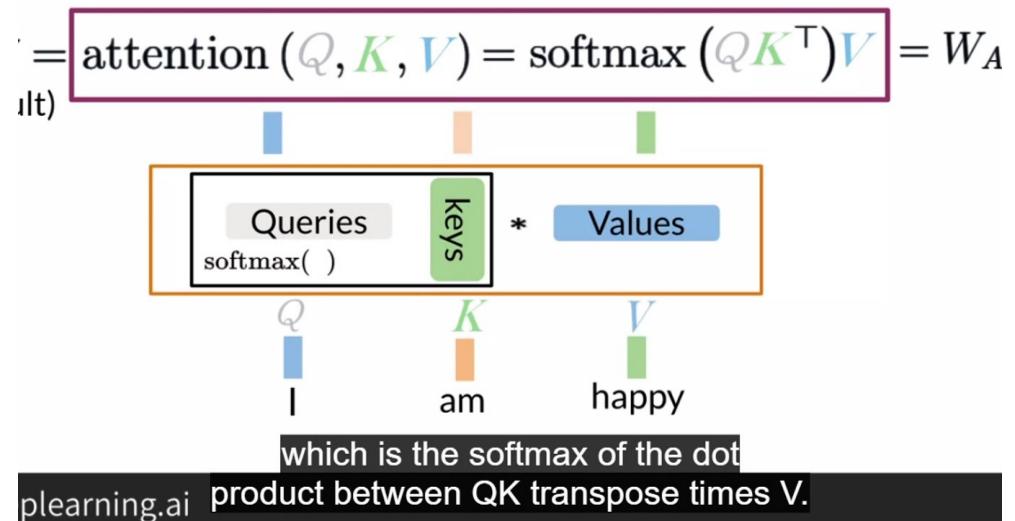
Attention formula



Concept of attention



Attention formula



Outline

- Ways of Attention
- Overview of Causal Attention
- Math behind causal attention



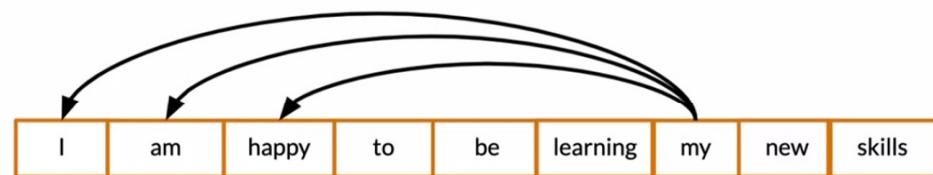
plearning.ai

First, you'll see what are three main ways of attention.

usual attention

Queries and keys are words from the same sentence

Queries should only be allowed to look at words before



these were not generated yet, they can attend to any word in the past though.

Summary

Dot-product Attention is essential for Transformer

The input to Attention are queries, keys, and values

A softmax function makes attention more focused on best keys

GPUs and TPUs is advisable for matrix multiplications

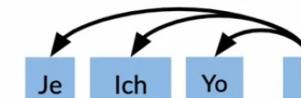


plearning.ai

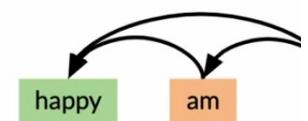
Dot-product Attention is the heart and soul of transformers.

three ways of attention

Encoder/decoder attention: One sentence (decoder) looks at another one (encoder)



Causal (self) attention: In one sentence, words look at previous words (used for generation)



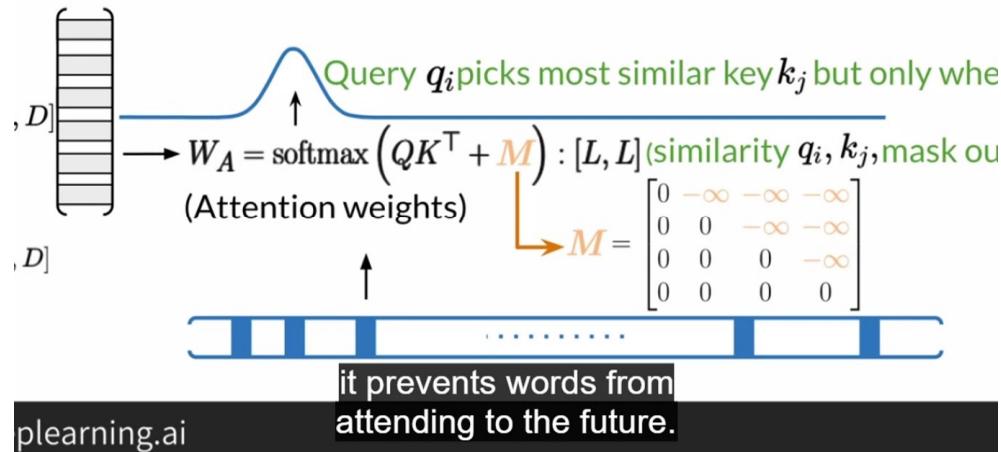
Bi-directional self attention: In one sentence, words look at both previous and future words



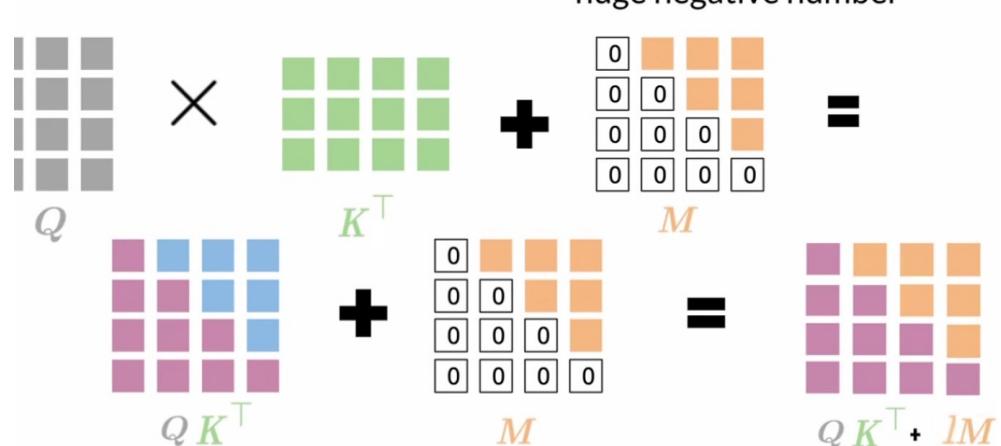
where words in the same sentence look both at previous and future words.

plearning.ai

usual attention math



usual attention math



summary

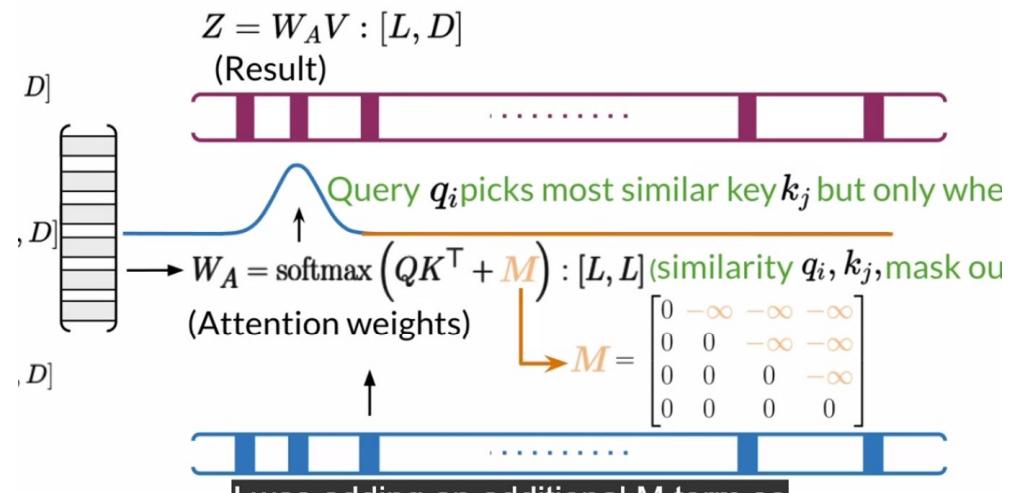
There are three main ways of Attention: Encoder/Decoder, Causal Bi-directional type

In causal attention, queries and keys come from the same sentence. queries search among words before only



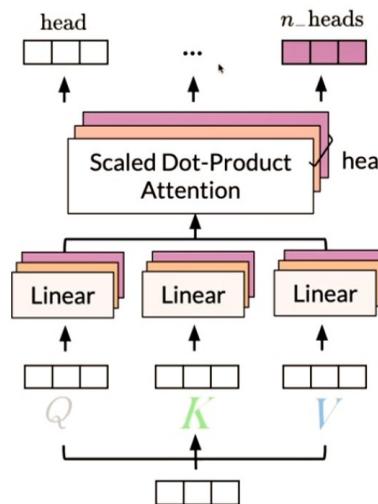
Encoder decoder attention, causal self-attention, and

usual attention math



Multi-Head Attention

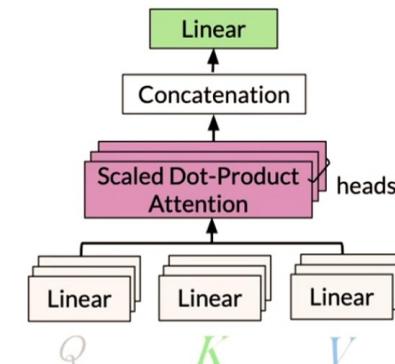
- Each head uses different linear transformations to represent words
- Different heads can learn different relationships between words



So over here you have
N different heads,

plearning.ai

Multi-Head Attention - Overview



except by the scale factor,

plearning.ai

Intuition

Intuition of Multi-Head Attention

Scaled dot-product and concatenation

Multi-Head Attention formula

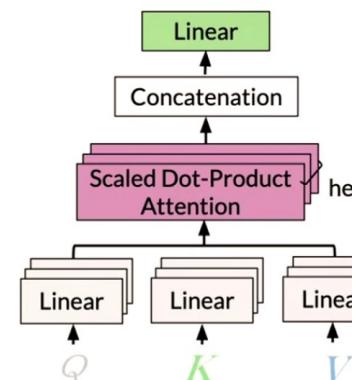


some intuition on
multi-head attention.

plearning.ai

Multi-Head Attention - Scaled dot product

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right)V$$

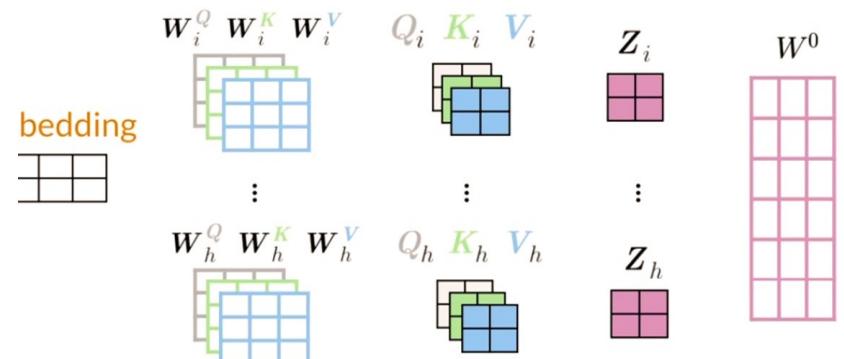


DK is the key inquiry dimension.

plearning.ai

Multi-Head Attention math

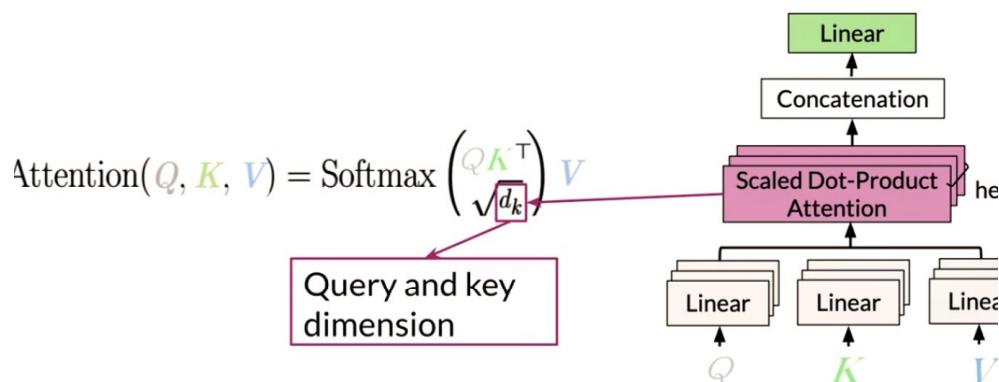
Multi-Head Attention math



To understand the math behind multi-head attention,

matrix capital W , superscript 0,

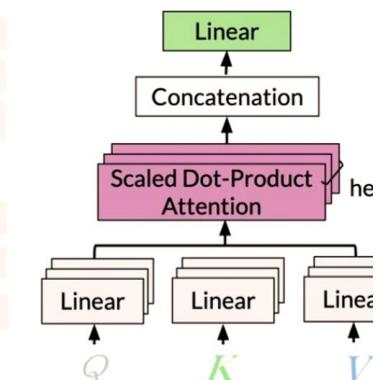
Multi-Head Attention - Scaled dot product



from the function to be

Multi-Head Attention - Concatenation

Input(Q, K, V): [batch, length, d_{model}]
 Linear layer: [batch, length, $n_heads * d_{head}$]
 Transpose: [batch, n_heads , length, d_{head}]
 Apply attention treating n_heads like batch
 Result shape: [batch, n_heads , length, d_{head}]
 Transpose: [batch, length, $n_heads * d_{head}$]
 Linear layer into: [batch size, length, d_{model}]



different heads into a joint representation.

Summary

Different heads can learn different relationship between words

Scaled dot-product is adequate for Multi-Head Attention

Multi-Headed models attend to information from different representations at different positions



and values more than once,

learning.ai

Multi-Head Attention Formula

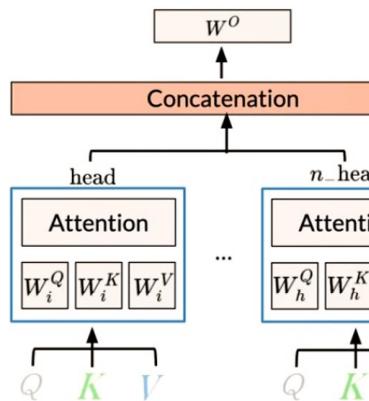
$$\text{MultiHead}(Q, K, V) = \text{Concat}(h_1, \dots, h_h)W^0$$

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

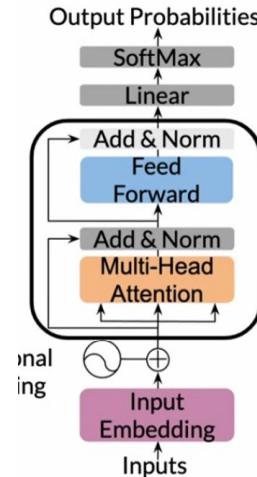
Each head h_i is the attention function of **Query, Key and Value** with trainable parameters (W_i^Q, W_i^K, W_i^V)

of queries, keys, and values.

learning.ai



Transformer decoder



The attention and feed-forward layers are repeated N times.

Overview

- input: sentence or paragraph
 - we predict the next word
- sentence gets embedded, add positional encoding
 - (vectors representing $\{0, 1, 2, \dots, K\}$)
- multi-head attention looks at previous words
- feed-forward layer with ReLU
 - that's where most parameters are!
- residual connection with layer normalization
- repeat N times

Outline

Overview of Transformer decoder

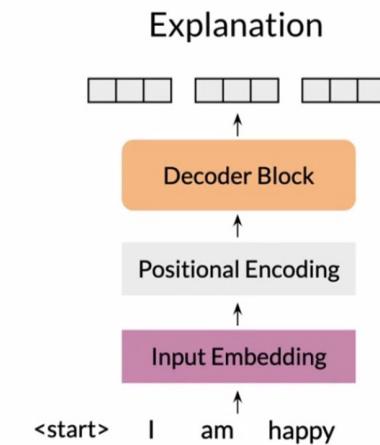
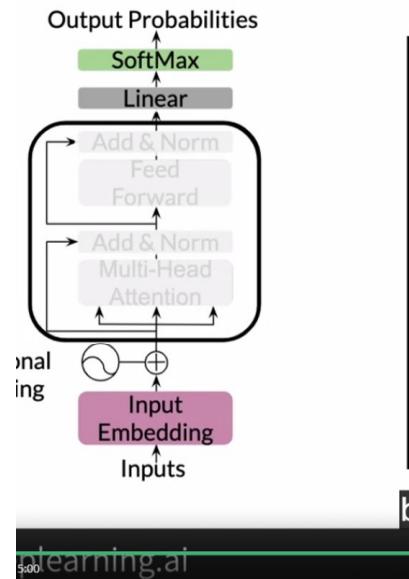
Implementation (decoder and feed-forward block)



>> In this video, you'll see the basic structure of a transformer decoder.

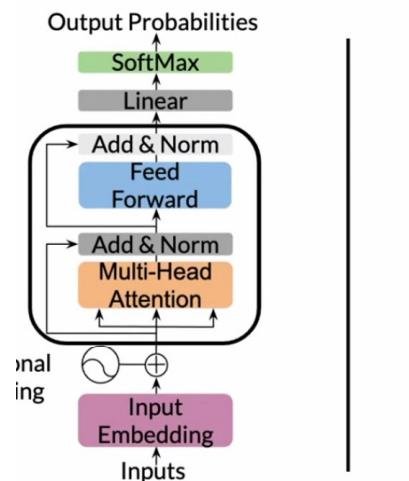
learning.ai

Transformer decoder



by vocab size.

Transformer decoder



Overview

- input: sentence or paragraph
 - we predict the next word
- sentence gets embedded, add positional encoding
 - (vectors representing $\{0, 1, 2, \dots, K\}$)
- multi-head attention looks at previous words
- feed-forward layer with ReLU
 - that's where most parameters are!
- residual connection with layer normalization
- repeat N times
- dense layer and softmax for output

And then you have a final dense layer for outputs and a softmax layer, and

plearning.ai

Summary

Transformer decoder mainly consists of three layers

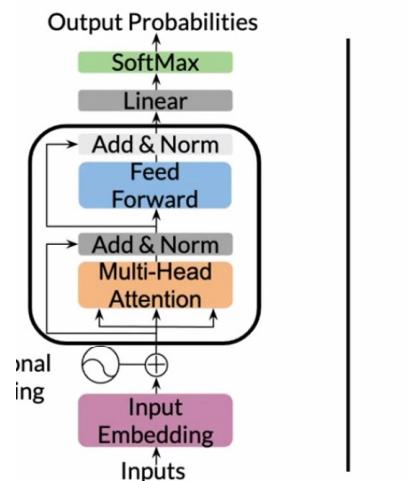
Decoder and feed-forward blocks are the core of this model code

It also includes a module to calculate the cross-entropy loss



You saw the building blocks used to implement a transformer decoder.

Transformer decoder



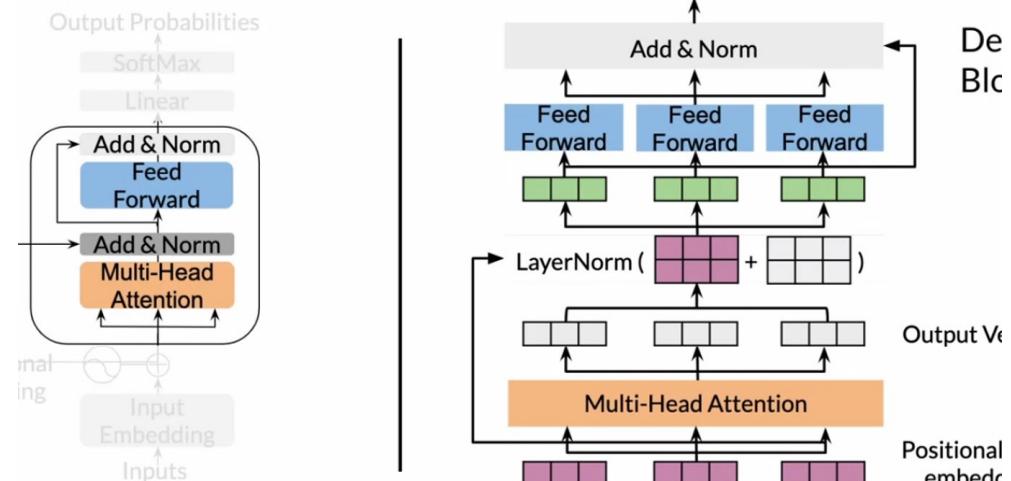
Overview

- input: sentence or paragraph
 - we predict the next word
- sentence gets embedded, add positional encoding
 - (vectors representing $\{0, 1, 2, \dots, K\}$)
- multi-head attention looks at previous words
- feed-forward layer with ReLU
 - that's where most parameters are!
- residual connection with layer normalization
- repeat N times
- dense layer and softmax for output

And then you have a final dense layer for outputs and a softmax layer, and

plearning.ai

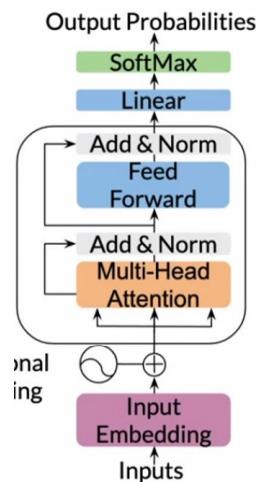
Transformer decoder



Finally, the encoder layer output is obtained.

plearning.ai

Technical details for data processing



Model Input:

ARTICLE TEXT <EOS> SUMMARY <EOS> <pad> ...

Tokenized version:

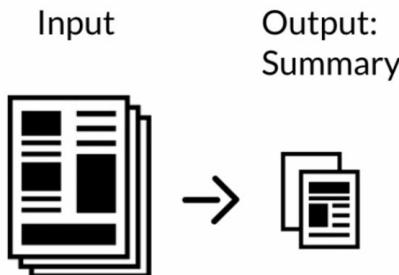
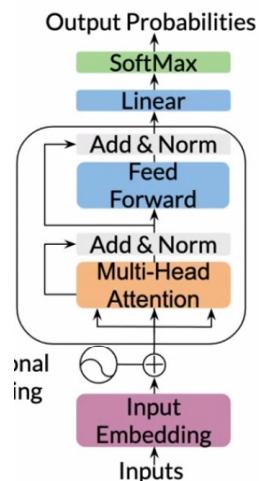
[2,3,5,2,1,3,4,7,8,2,5,1,2,3,6,2,1,0,0]

Loss weights: 0s until the first <EOS> and 1 on the start of the summary.

you weight the loss for

plearning.ai the words within the article with 0s, and

Transformer for summarization



plearning.ai

But the one thing may
immediately stand out to you.

Difference with a Language Model

Model input:

Article] <EOS> [Summary] <EOS>

Inference:

Provide: [Article] <EOS>

Generate summary word-by-word

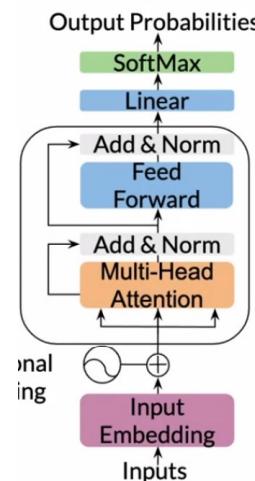
- until the final <EOS>

Pick the next word by random sampling

- each time you get a different summary!

it generates a probability distribution over all possible words.

First function



Cross entropy loss

$$J = -\frac{1}{m} \sum_j^m \sum_i^K y_j^i \log \hat{y}_j^i$$

j : over summary

i : batch elements



the article to be summarized.

2:46

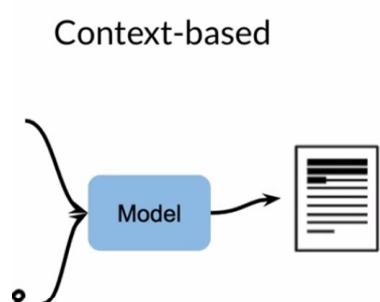
plearning.ai

plearning.ai

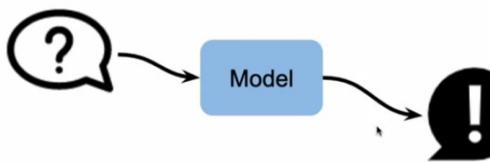
4:14

Question Answering

Context-based



Closed book



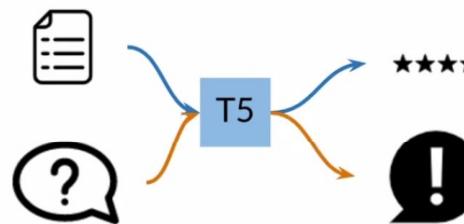
so it comes up with
its own answer.

Week 3

Question Answering



BERT



Transfer Learning

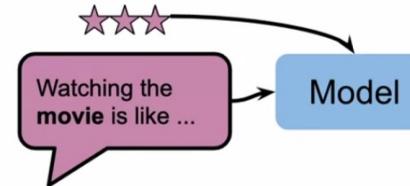


you can see here it has
several possible inputs.

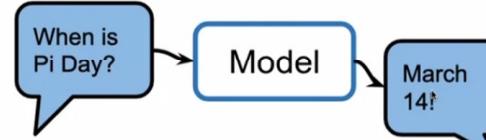
learning.ai

Transfer Learning: Different Tasks

e-Training
timent
Classification



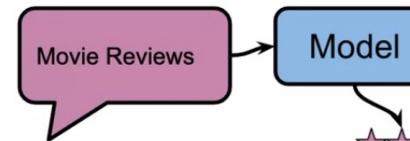
Inference
nstream task:
Question Answering



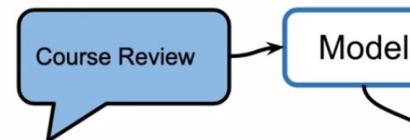
Then you can ask them
model the same question.

Transfer learning

e-training



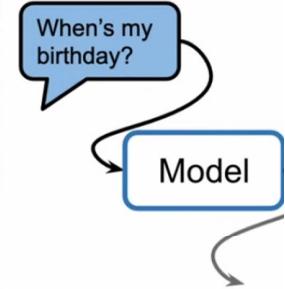
Inference
Downstream Task



and you get your prediction.

learning.ai

Inference



: more data, better performance

English wikipedia
~13 GB

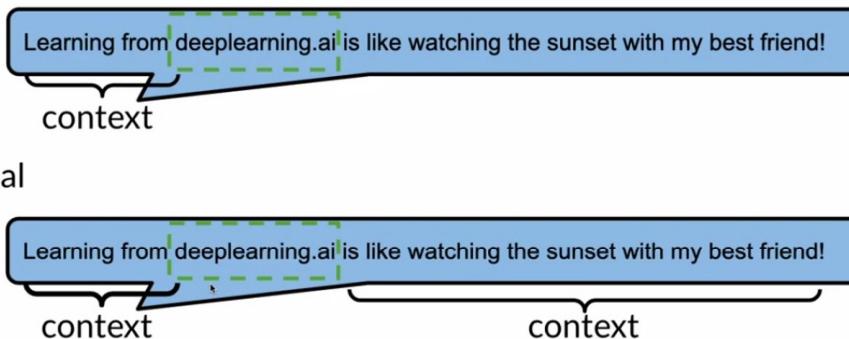
C4
Colossal Clean Crawled
Corpus
~800 GB

is that the more data you have,

plearning.ai

RT: Bi-directional Context

-directional



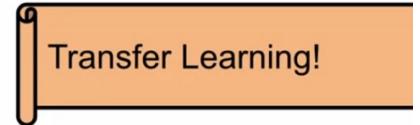
nlearning ai

This is one of the main takeaways for bidirectionality.

sirable Goals



- Reduce training time
 - Improve predictions
 - Small datasets

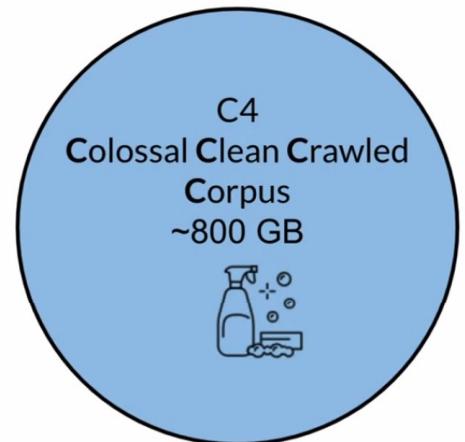


a lot from other tasks.

plearning.ai

: more data, better performance

English wikipedia
~13 GB

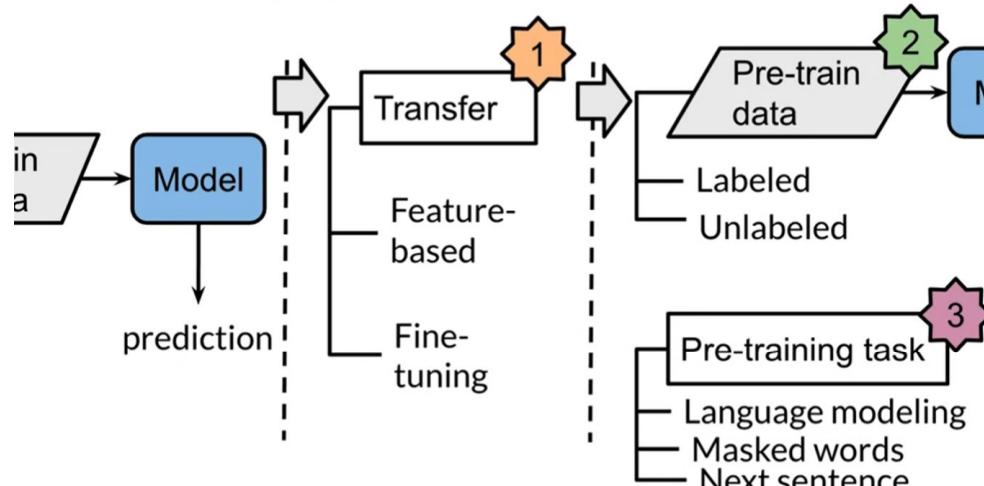


plearning ai

when compared to the English Wikipedia.

Transfer

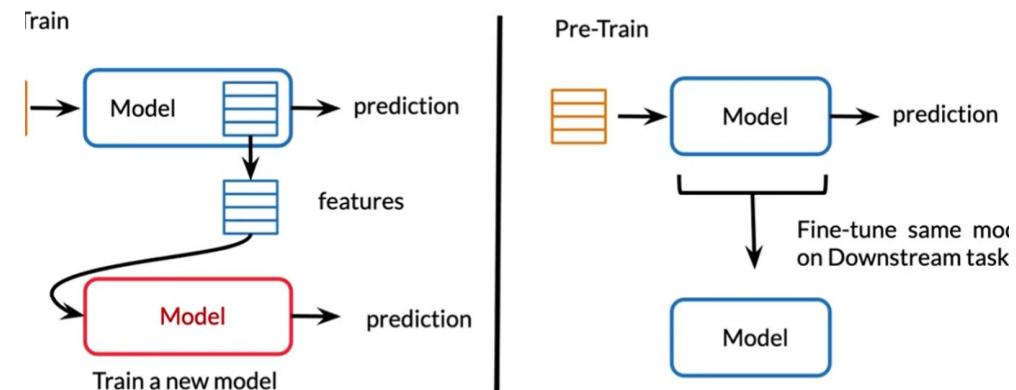
Transfer learning options



olearning.ai

or the third option
is pre-training task,

ture-based vs. Fine-Tuning



olearning.ai

fine tune the same model
on a downstream tasks.

Transfer

sirable Goals



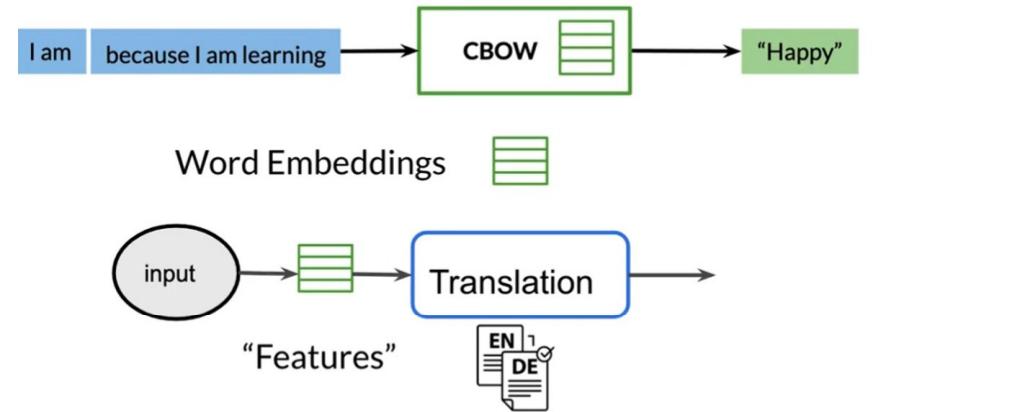
- Reduce training time
- Improve predictions
- Small datasets

Transfer Learning!

olearning.ai

We can improve on
these three goals with

General purpose learning

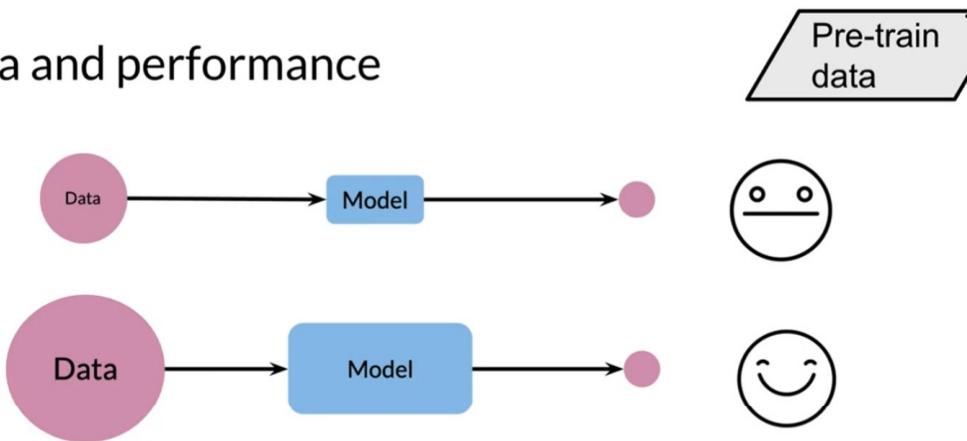


olearning.ai

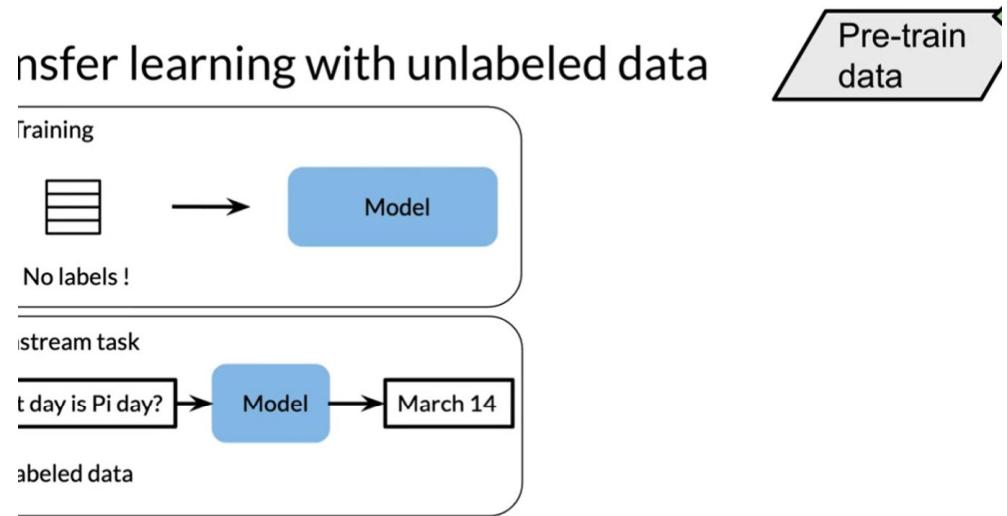
it can be used to provide more
meaning to the raw inputs.

246

Size and performance



Transfer learning with unlabeled data



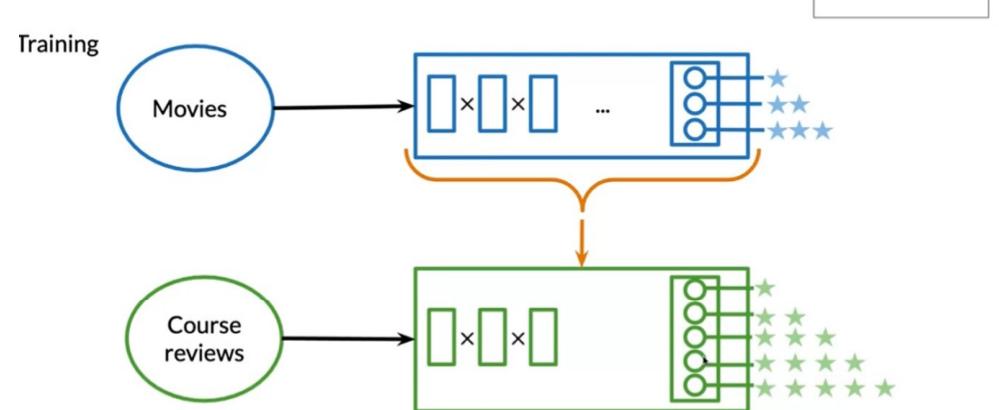
the bigger models you can build,

learning.ai

for example, over here,

5:14

Re-tune: adding a layer



Then you create a
new output layer,

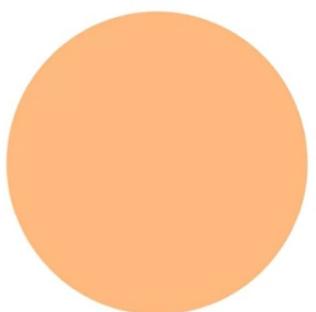
learning.ai

Labeled vs Unlabeled Data

Labeled text data



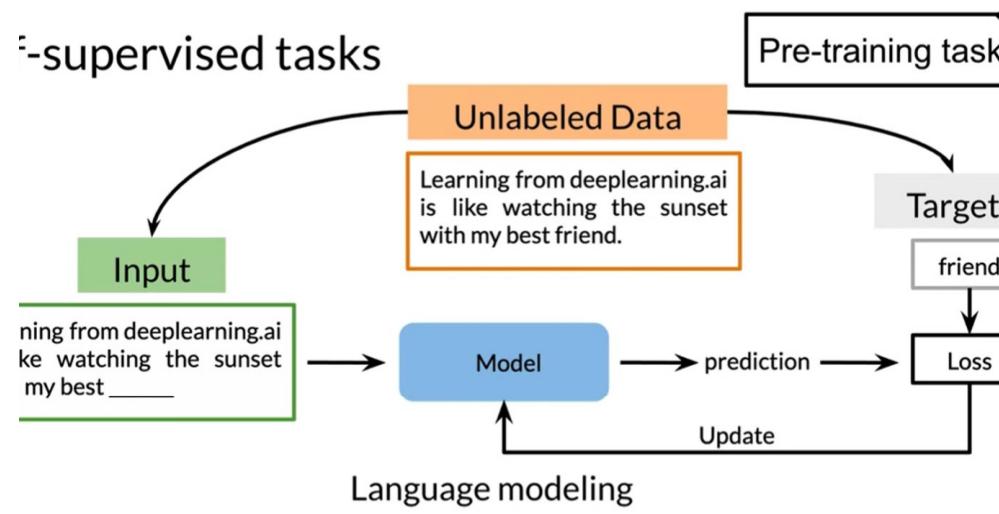
Unlabeled text data



and as you can see
that there is way

5:05

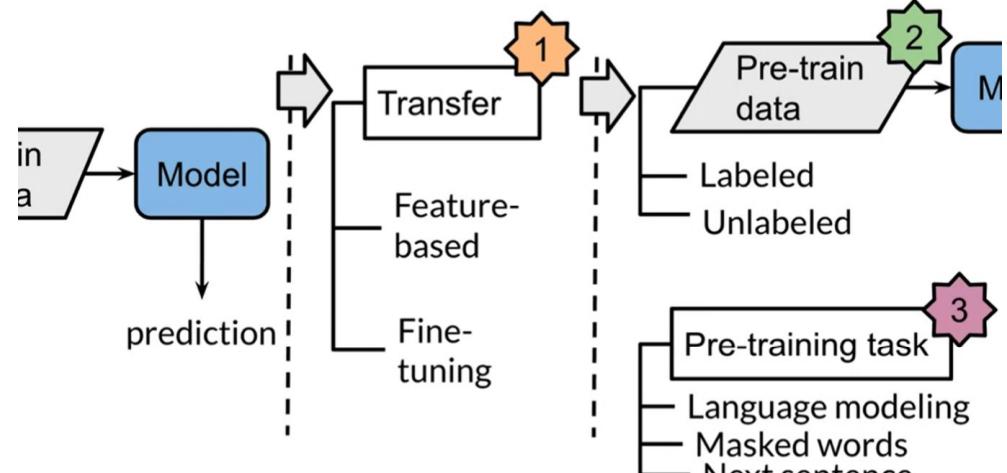
semi-supervised tasks



Then as you can see,
the targets is friend,

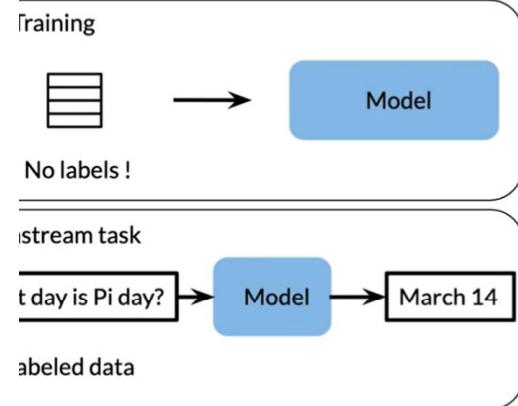
deeplearning.ai

Summary



like mask toward modeling or
next sentence prediction.

Transfer learning with unlabeled data

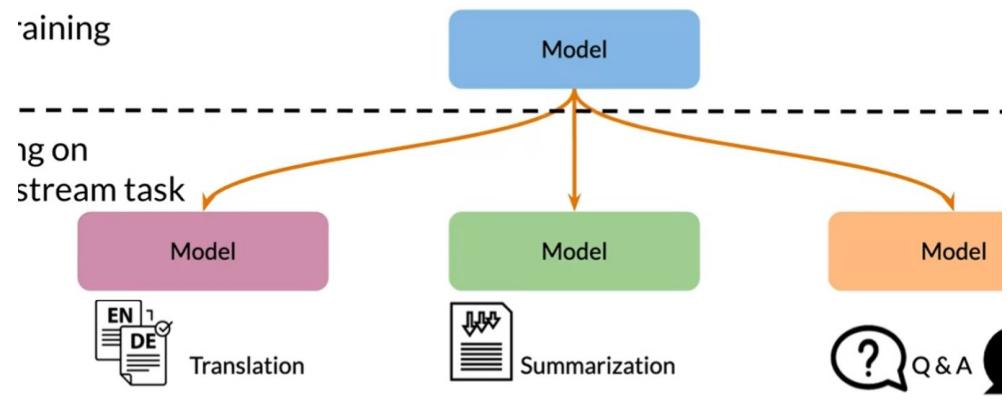


The question is, which tasks
work with unlabeled data?

deeplearning.ai

Which tasks work with
unlabeled data?

Fine-tune a model for each downstream task



or to do question answering.

deeplearning.ai

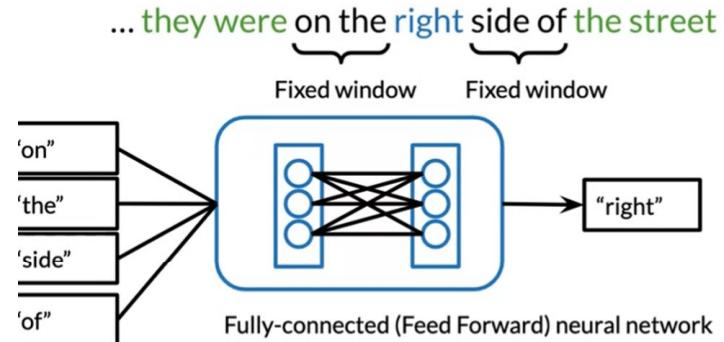
line

:BOW ELMo GPT BERT T5

learning.ai

were solving for and how they

ntinuous Bag of Words



which is the center word.

ntext

... right ...

... they were on the right ...

... they were on the right side of the street



deeplearning.ai

ELMo, GPT BERT, T5

Now you will get a overview
of some models such as ELMo,

7:07

the context after
the word and we use

learning.ai

7:07

ed more context?

... they were on the right side of the street.
Fixed window Fixed window

... they were on the right side of history.

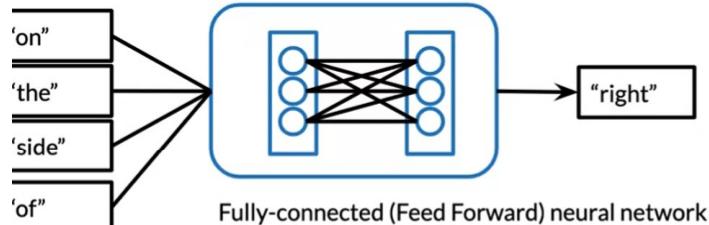
1:32
7:07
nlearning.ai

it could be instead
of the streets,

ntinuous Bag of Words

... they were on the right side of the street

Fixed window Fixed window

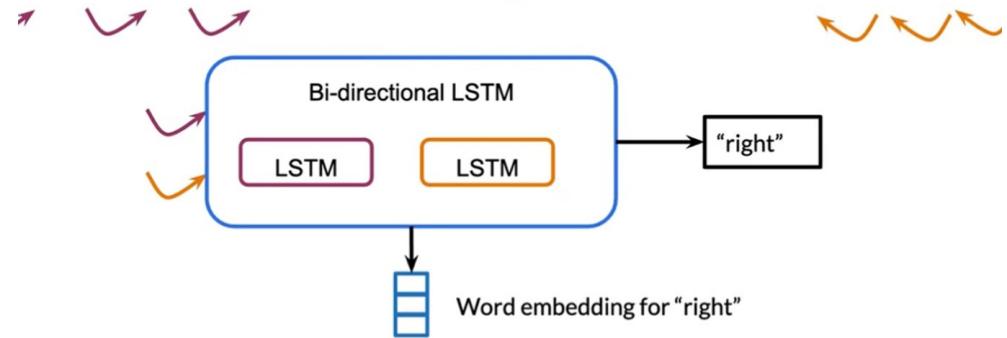


a neural network and you
predict the outputs,

1:14
7:07
nlearning.ai

Mo: Full context using RNN

legislators believed that they were on the ___ side of history so they changed the law.



You get the word embedding for
plearning.ai

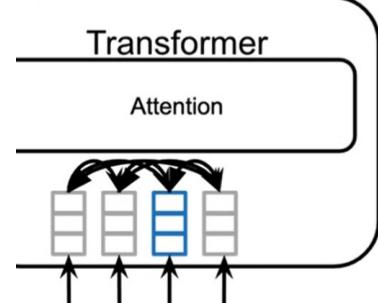
e all context words

legislators believed that they were on the right side of history, so they changed the law

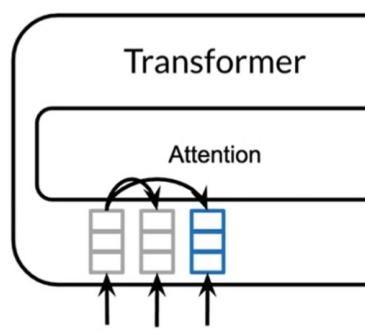
Over here, it starts from

1:44
7:07
nlearning.ai

T: Uni-directional



... on the right side...
which word can peek at itself!

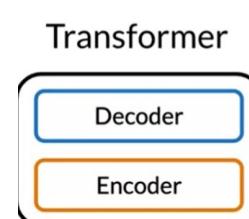


... on the right
No peeking!

only look at the previous inputs,

plearning.ai

RT

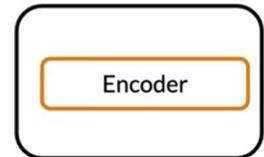


The legislators believed that they were on the side of history, so they changed the law.

GPT



BERT

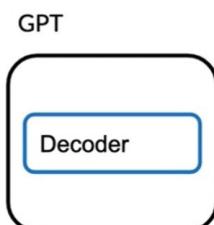
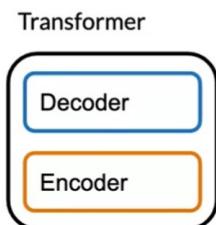
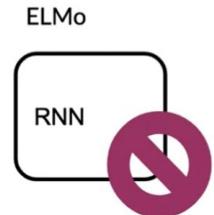


Bi-directional

bi-directional representations

plearning.ai

Open AI GPT



legislators believed that they were on the

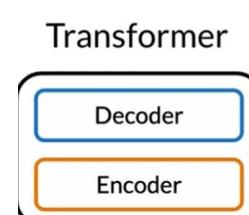


Uni-directional

as capturing longer-term
dependencies.
3/07

plearning.ai

RT

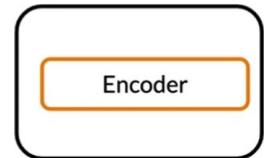


The legislators believed that they were on the side of history, so they changed the law.

GPT



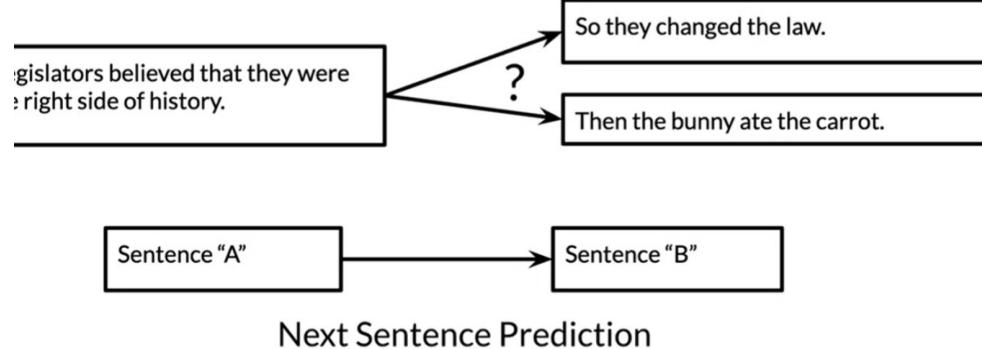
BERT



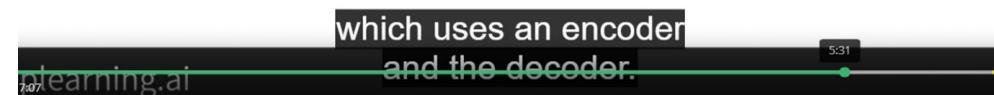
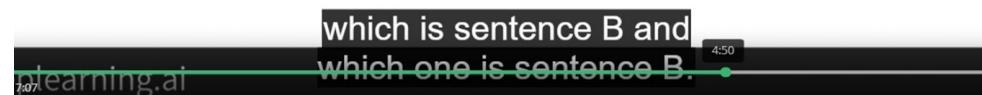
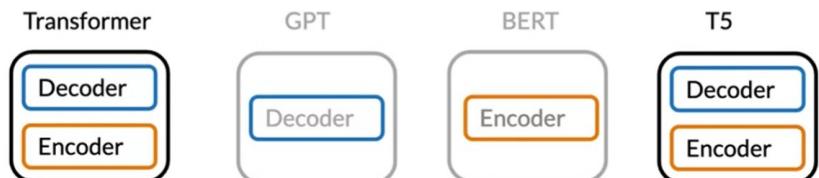
As it works over here,

plearning.ai

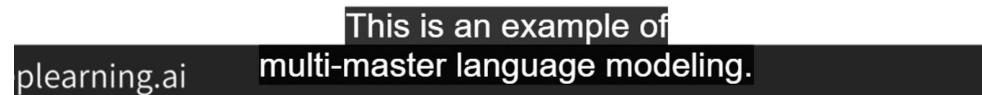
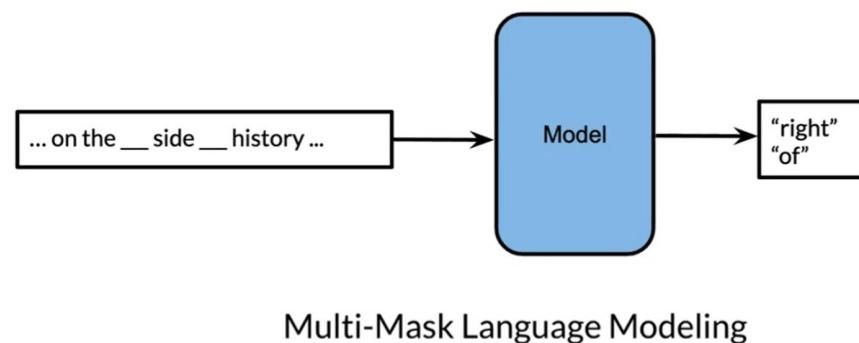
RT: Words to Sentences



: Encoder vs. Encoder-Decoder

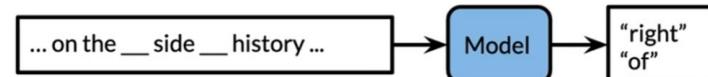


Informer + Bi-directional Context

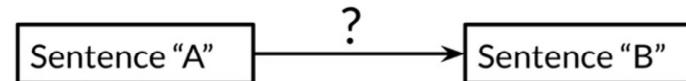


RT Pre-training Tasks

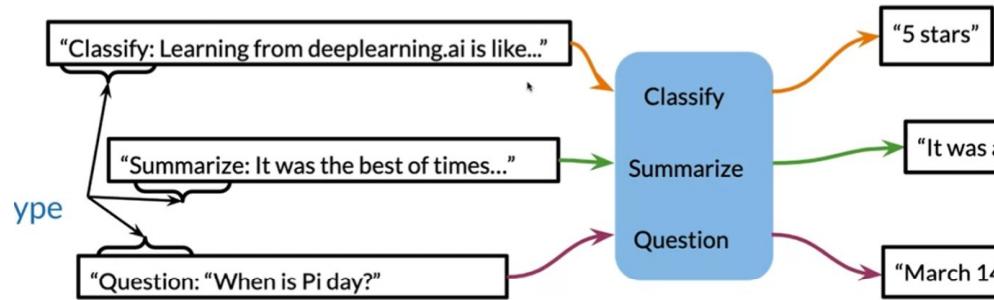
ti-Mask Language Modeling



Next Sentence Prediction



: Text-to-Text



Outline

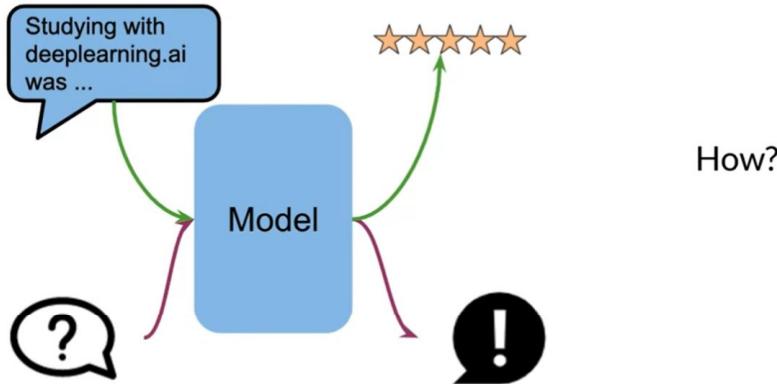
Learn about the BERT architecture

Understand how BERT pre-training works

basically tell the model
that's it's going to classify

the BERT architecture, and then,

: Multi-task



How?

Summary

BOW	ELMo	GPT	BERT	T5
text	Full sentence	Transformer:	Transformer:	Transformer:
bow	Bi-directional	Decoder	Encoder	Encoder - Decoder
N	Context	Uni-directional	Bi-directional	Bi-directional
	RNN	Context	Context	Context
			Multi-Mask	Multi-Task
			Next Sentence	
			Prediction	

How can you do this?

then it used fixed context window

ERT

A multi layer bidirectional transformer

Positional embeddings

BERT_base:

12 layers (12 transformer blocks)

12 attentions heads

110 million parameters

4:30
ntelearning.ai

BERT base which has 12 layers
or 12 transformer blocks,

ERT pre-training

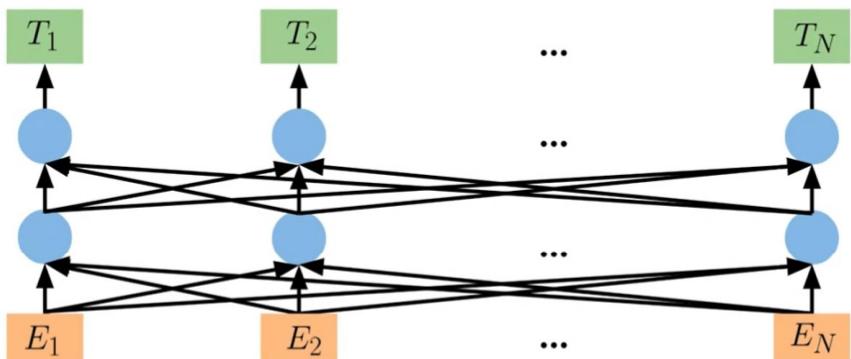
fter school Lukasz does his _____ in the library.

Masked language modeling (MLM)

Maybe work, maybe homework.
plearning.ai

ERT

- Makes use of transfer learning/pre-training:



and then, you get
your T_1, T_2, T_N .

4:30
ntelearning.ai

ERT

A multi layer bidirectional transformer

Positional embeddings

BERT_base:

12 layers (12 transformer blocks)

12 attentions heads

110 million parameters

They have way, way more
parameters and way,

4:30
ntelearning.ai

ummary

Choose 15% of the tokens at random: mask them 80% of the time, replace them with a random token 10% of the time, or keep as is 10% of the time.

There could be multiple masked spans in a sentence

Next sentence prediction is also used when pre-training.

15 percent of the tokens at random.

ntearning.ai

BERT pre-training

After school Lukasz does his _____ in the library.

Masked language modeling (MLM)

by the embedding matrix, and then,

ntearning.ai

rmalizing the input

put	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing
en	E [CLS]	E my	E dog	E is	E cute	E [SEP]	E he	E likes	E play	E ##ing
beddings	+	+	+	+	+	+	+	+	+	+
ment	E A	E A	E A	E A	E A	E A	E B	E B	E B	E B
beddings	+	+	+	+	+	+	+	+	+	+
ition	E 0	E 1	E 2	E 3	E 4	E 5	E 6	E 7	E 8	E 9
beddings										

and then you get your new inputs.

ntearning.ai

tline

Understand how BERT inputs are fed into the model

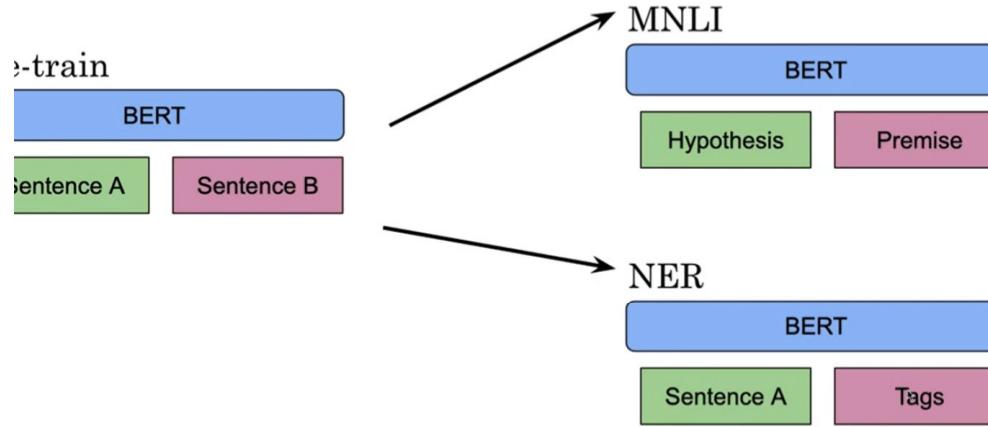
Visualize the output

Learn about the BERT objective

You're going to learn how
Bert inputs are fed into

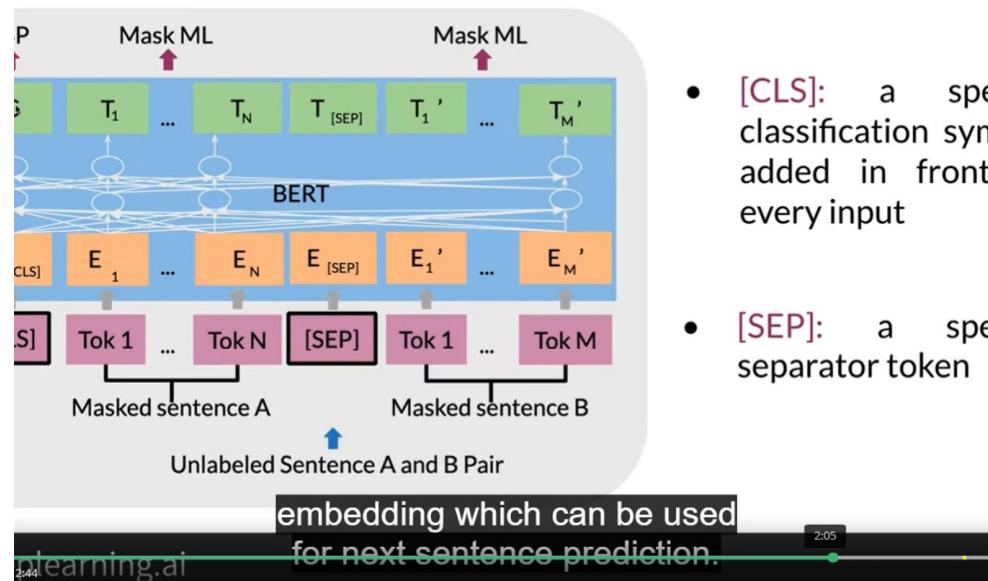
ntearning.ai

Re-tuning BERT: Outline



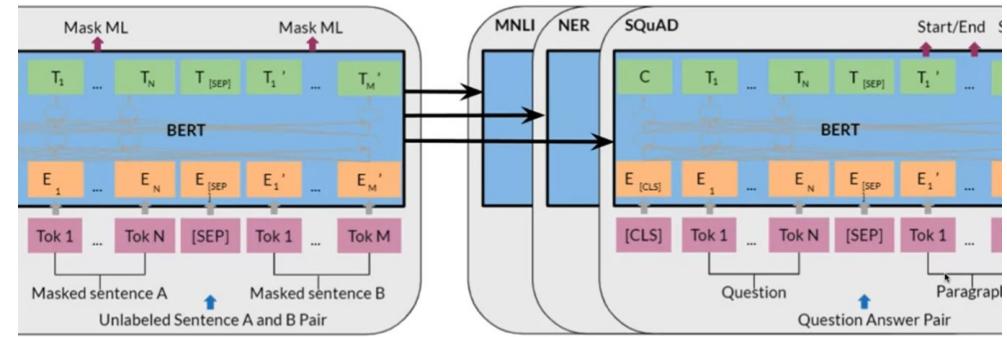
plearning.ai

Visualizing the output



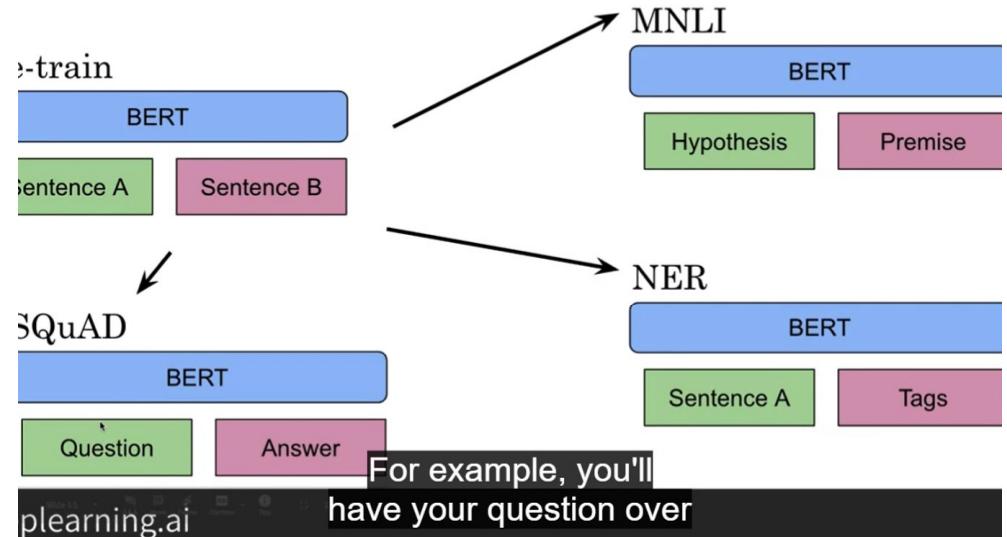
plearning.ai

Outputs



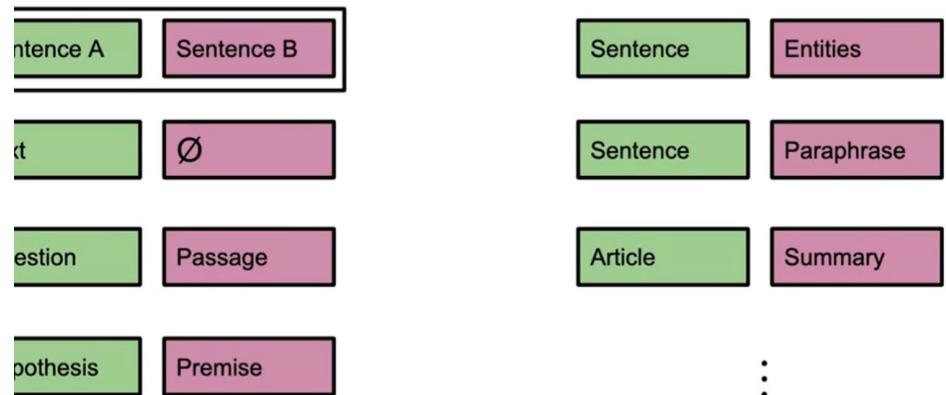
Over here, you'll have the paragraph and this will give

Re-tuning BERT: Outline



plearning.ai

Summary



This is just the inputs
into your BERT's model.

plearning.ai

Transformer - T5 Model

Text to Text

Machine Translation



Classification

Summarization



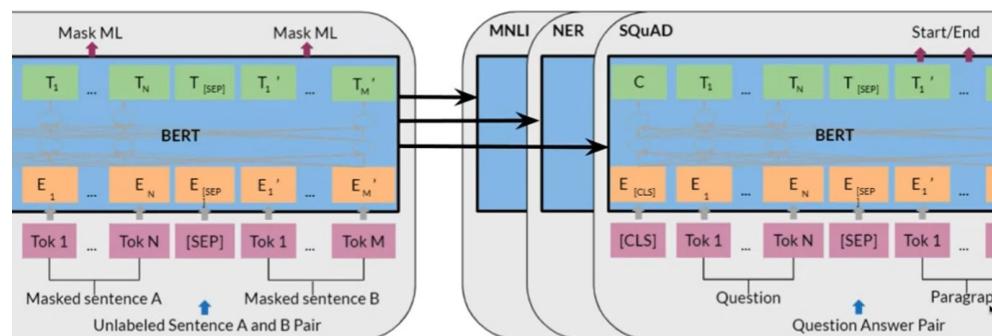
Question
Answering (Q&A)

Sentiment



And there are other applications
that you can use T5 for, but

Outputs



the hypothesis and then
the premise and so forth.

plearning.ai

Outline

Understand how T5 works

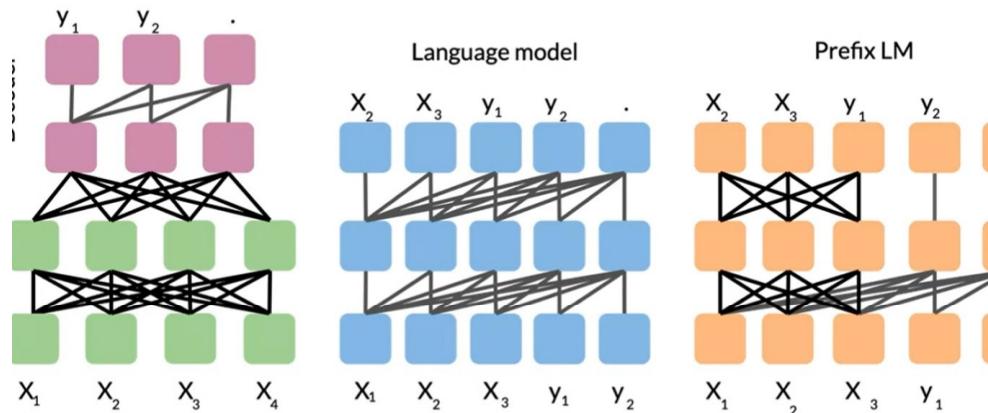
Recognize the different types of attention used

Overview of model architecture

So we're going to understand how T5 works.

plearning.ai

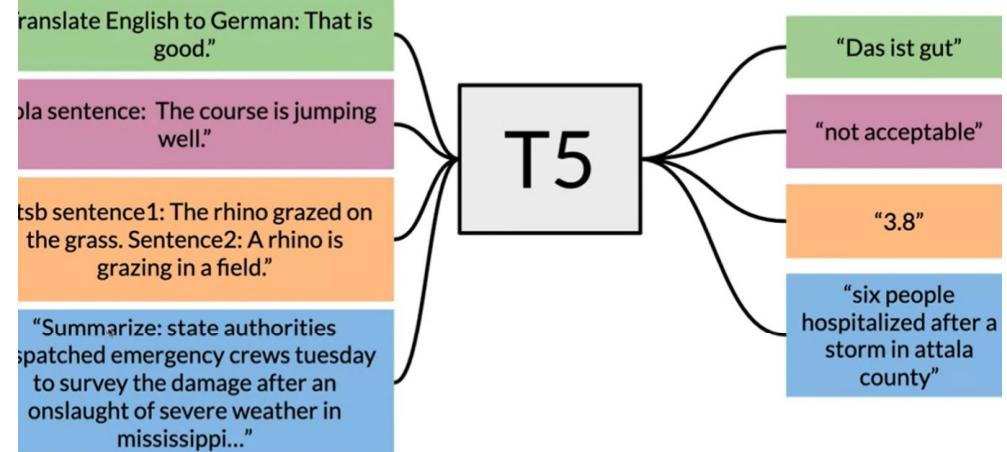
Model Architecture



So as you can see over here,
it's doing causal masking.

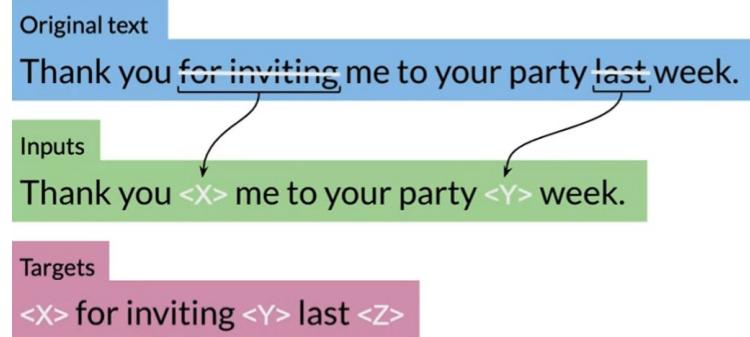
plearning.ai

Multi-task training strategy



This is how it works.

transformer - T5 Model



So each bracket corresponds
to a certain target.

plearning.ai

summary

Prefix LM attention

Model architecture

Pre-training T5 (MLM)

You've seen the model architecture for T5.

plearning.ai

Multi-task Training Strategy

Fine-tuning method	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
All parameters	83.28	19.24	80.88	71.36	26.98	39.82	27.65
Adapter layers, $d = 32$	80.52	15.08	79.32	60.40	13.84	17.88	15.54
Adapter layers, $d = 128$	81.51	16.62	79.47	63.03	19.83	27.50	22.63
Adapter layers, $d = 512$	81.54	17.78	79.18	64.30	23.45	33.98	25.81
Adapter layers, $d = 2048$	81.51	16.62	79.47	63.03	19.83	27.50	22.63
Gradual unfreezing	82.50	18.95	79.17	70.79	26.71	39.02	26.93

this is a table found in the original paper,

Input and Output Format

Machine translation:

- translate English to German: That is good.
- Predict entailment, contradiction , or neutral
- mnli premise: I hate pigeons hypothesis: My feelings towards pigeons are filled with animosity. target: entailment

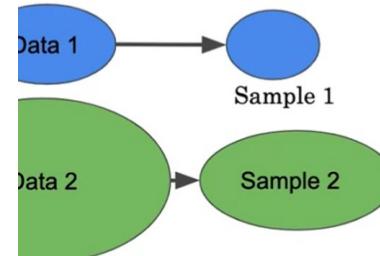
Winograd schema

- The city councilmen refused the demonstrators a permit because *they* feared violence

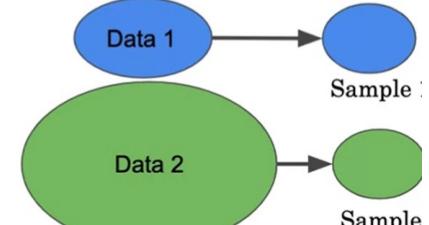
which is another way to predict whether a pronoun,

Data Training Strategies

samples-proportional mixing



Equal mixing



Temperature-scaled mixing

the parameters to get something in between.

Input and Output Format

Machine translation:

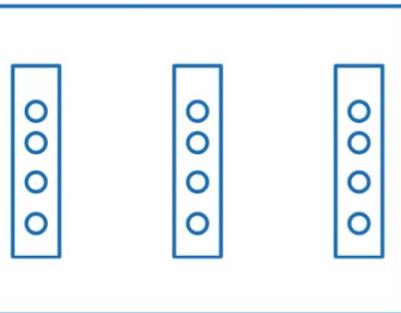
- translate English to German: That is good.
- Predict entailment, contradiction , or neutral
- mnli premise: I hate pigeons hypothesis: My feelings towards pigeons are filled with animosity. target: entailment

Winograd schema

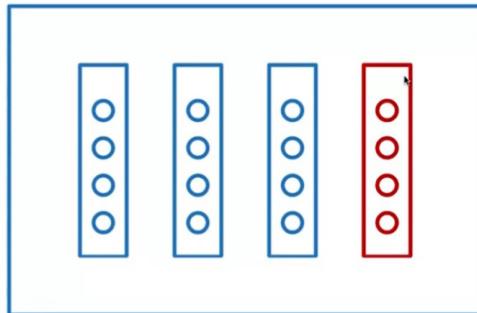
- The city councilmen refused the demonstrators a permit because *they* feared violence

your model and then it will be

Gradual unfreezing vs. Adapter layers



Gradual unfreezing



Adapter layers

limits of Transfer learning with a unified
add a neural network to
plearning.ai

Natural Language Understanding Evaluation

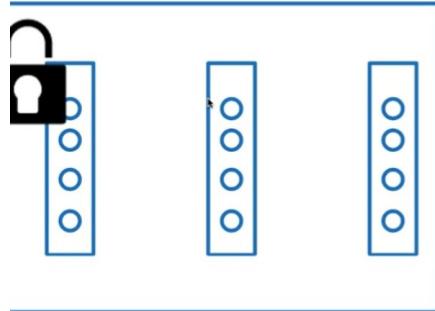
A collection used to train, evaluate, analyze natural language understanding systems

Datasets with different genres, and of different sizes and difficulties

Leaderboard

So people can use the data sets, and
plearning.ai

Gradual unfreezing vs. Adapter layers

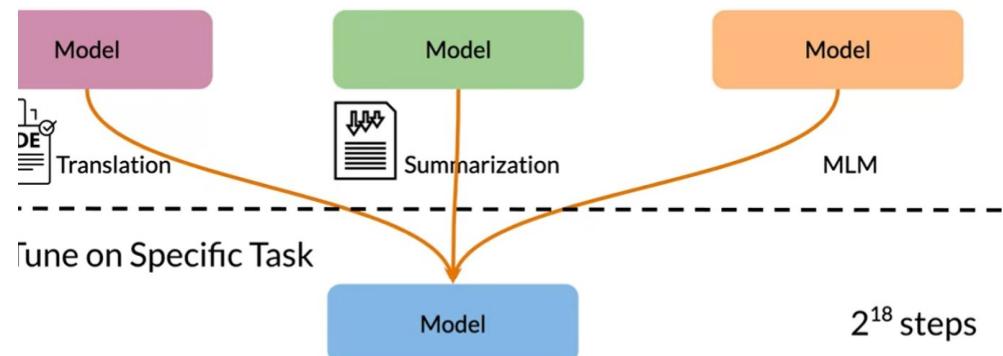


Gradual unfreezing

limits of Transfer learning with a un
Then you freeze this one and
plearning.ai then you freeze this one,

Fine-tuning

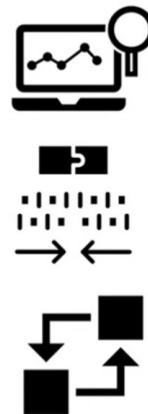
Training



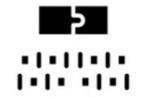
plearning.ai

General Language Understanding Evaluation

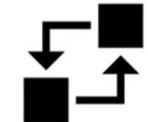
Drive research



Model agnostic



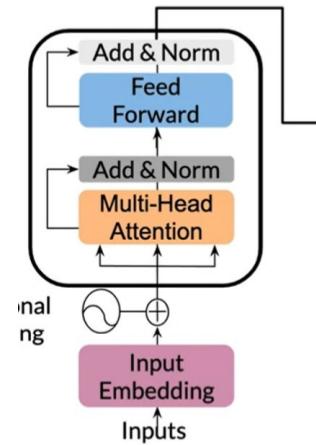
Makes use of transfer learning



because you have access
to several datasets.

plearning.ai

Transformer encoder



Feedforward:

```
[  
    LayerNorm,  
    dense,  
    activation,  
    dropout_middle,  
    dense,  
    dropout_final  
]
```

Encoder block:

```
[  
    Residual(  
        LayerNorm  
        attention  
        dropout_,  
    ),  
    Residual(  
        feed_forw  
    )  
]
```

the encoder block which has the feed forward and two residual connections.

plearning.ai

Tasks Evaluated on

Sentence grammatical or not?

Sentiment

Paraphrase

Similarity

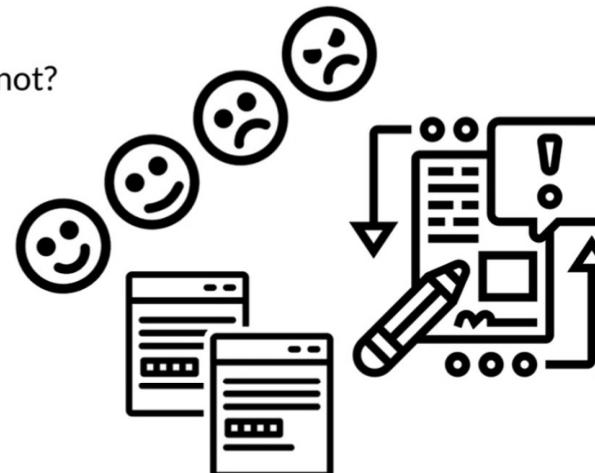
Questions duplicates

Answerable

Contradiction

Entailment

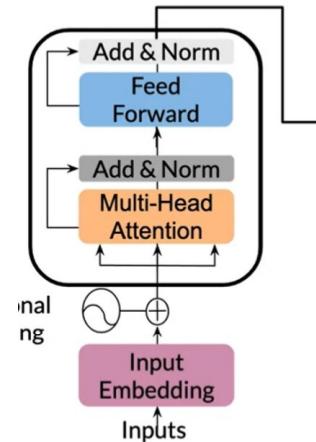
Winograd (co-ref)



And also for the Winograd Schema, which
is basically trying to identify whether,

plearning.ai

Transformer encoder



Feedforward:

```
[  
    LayerNorm,  
    dense,  
    activation,  
    dropout_middle,  
    dense,  
    dropout_final  
]
```

you have a layerNorm, a dense, followed by an activation.

plearning.ai

plementing Q&A with T5

Load a pre-trained model

Process data to get the required inputs and outputs: "question: Q context: C" as input and "A" as target

Fine tune your model on the new task and input

Predict using your own model

the inputs, and
plearning.ai then you predict using your own model.

ata examples

estion: What percentage of the French population today is non - European ?

text: Since the end of the Second World War , France has become an ethnically diverse country . . .
proximately five percent of the French population is non - European and non - white . This de
ach the number of non - white citizens in the United States (roughly 28 - 37 % , depending on how Latinos are cl
Demographics of the United States) . Nevertheless , it amounts to at least three million people , and has forced the
nic diversity onto the French policy agenda . France has developed an approach to dealing with ethnic proble
s in contrast to that of many advanced , industrialized countries . Unlike the United States , Britain , or ev
erlands , France maintains a " color - blind " model of public policy . This means that it targets virtually no policies
cial or ethnic groups . Instead , it uses geographic or class criteria to address issues of social inequalities . It has , ho
oped an extensive anti - racist policy repertoire since the early 1970s . Until recently , French policies focused pr
sues of hate speech — going much further than their American counterparts — and relatively less on is
mination in jobs , housing , and in provision of goods and services .

get: Approximately five percent

You'll have a context and
then you'll have a target.

ng Text Sequences

ks In NLP:



Writing Books



Chatbots

building intelligent agents for
plearning.ai conversations like chatbots.

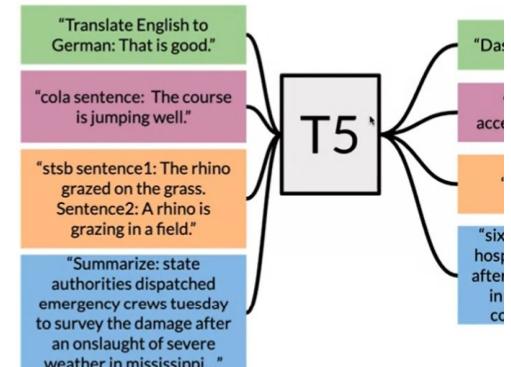
plementing Q&A with T5

Load a pre-trained model

Process data to get the required inputs and outputs: "question: Q context: C" as input and "A" as target

Fine tune your model on the new task and input

Predict using your own model



And remember this is the T5
plearning.ai that you'll be working with.

Chatbots

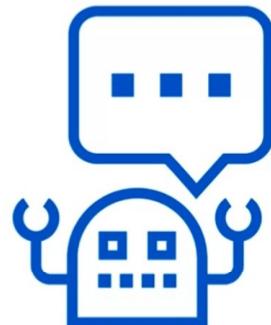
Context Windows:

User 1: What's for dinner?

Chatbot: Who's cooking, you or me?

User 1: Hey now chatbot.

Chatbot: I hope it's not hay, that's what horses eat.



It's going to be fun.

plearning.ai

Writing Text Sequences

Skills In NLP:



Writing Books



Chatbots

train a real working chatbot
in this week's assignments.¹²⁷

plearning.ai

Informer Issues

Attention on sequence of length L takes L^2 time and memory

$L=100 \quad L^2 = 10K \quad (0.001s \text{ at } 10M \text{ ops/s})$

$L=1000 \quad L^2 = 1M \quad (0.1s \text{ at } 10M \text{ ops/s})$

$L=10000 \quad L^2 = 100M \quad (10s \text{ at } 10M \text{ ops/s})$

$L=100000 \quad L^2 = 10B \quad (1000s \text{ at } 10M \text{ ops/s})$

N layers take N times as much memory

GPT-3 has 96 layers and new models will have more

For example, if L equals 100,

plearning.ai

Informer Issues

Attention on sequence of length L takes L^2 time and memory

N layers take N times as much memory

GPT-3 has 96 layers and new models will have more

GPT-3, for example,
already has 96 layers,

plearning.ai

Memory with N Layers

Activations need to be stored for backprop

Big models are getting bigger

Compute vs memory tradeoff

less memory and I'll show you how to do that next.

plearning.ai

Attention Complexity

Attention: $\text{softmax}(\mathbf{Q}\mathbf{K}^T)\mathbf{V}$

$\mathbf{Q}, \mathbf{K}, \mathbf{V}$ are all $[L, d_{\text{model}}]$

$\mathbf{Q}\mathbf{K}^T$ is $[L, L]$

Save compute by using area of interest for large L

and those immediately around it, by using attention.

plearning.ai

What does Attention do?

Select Nearest Neighbors (K, Q) and return corresponding V

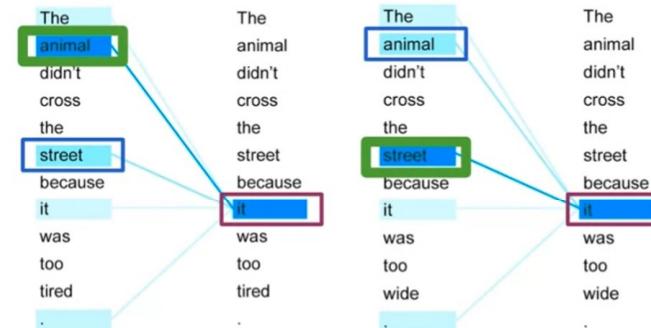


image ©
(Transformer: A
Network Archite
Language Under

It again refers to either the streets or the animal, but

plearning.ai

What does Attention do?

Select Nearest Neighbors (K, Q) and return corresponding V

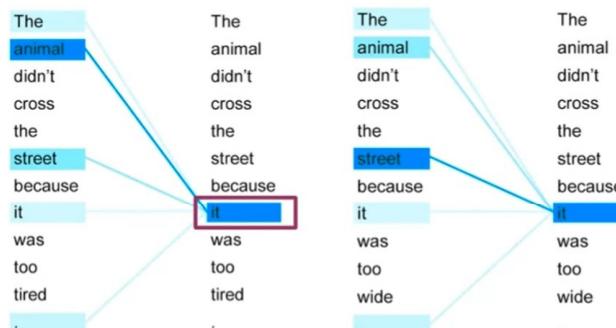


image ©
(Transformer: A
Network Archite
Language Under

Take the word it attention is focused on

plearning.ai

arest Neighbors

urse:

Natural Language Processing with Classification and Vector Space

ssons:

- KNN
- Hash Tables and Hash Functions
- Locality Sensitive Hashing
- Multiple Planes

Please feel free to skip ahead if

plearning.ai you're familiar with both KNN and LSH.

arest Neighbors

mpute the nearest neighbor to q among vectors $\{k_1, \dots, k_n\}$

Attention computes $d(q, k_i)$ for i from 1 to n which can be slow

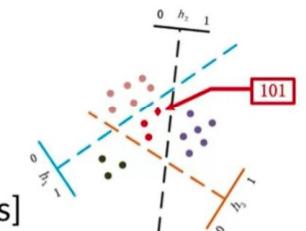
Faster *approximate* uses locality sensitive hashing (LSH)

Locality sensitive: if q is close to k_i :

$$\text{hash}(q) == \text{hash}(k_i)$$

Achieve by randomly cutting space

$$\text{hash}(x) = \text{sign}(xR) \quad R: [d, n_hash_bins]$$

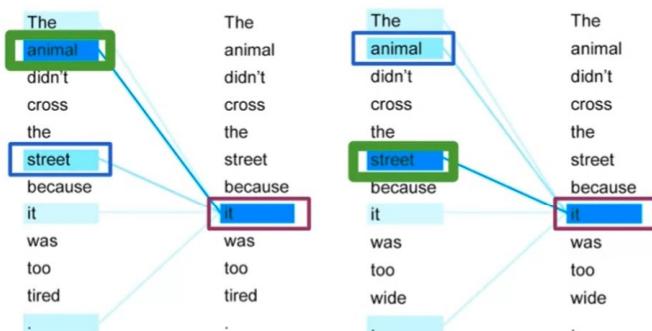


When choosing the hash you want to make

plearning.ai the buckets roughly the same size.

hat does Attention do?

lect Nearest Neighbors (K, Q) and return corresponding V



and you can ignore the other
words like it and was and tired.

plearning.ai

image ©
[Transformer: A
Network Architecture
for Language Under-](#)

arest Neighbors

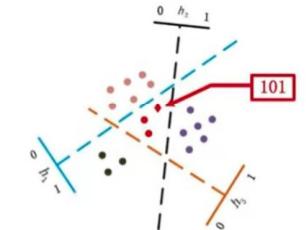
mpute the nearest neighbor to q among vectors $\{k_1, \dots, k_n\}$

Attention computes $d(q, k_i)$ for i from 1 to n which can be slow

Faster *approximate* uses locality sensitive hashing (LSH)

Locality sensitive: if q is close to k_i :

$$\text{hash}(q) == \text{hash}(k_i)$$



This helps you group similar query and
key vectors together,

plearning.ai

- Attention

Standard Attention:

$$A(Q, K, V) = \text{softmax}(QK^T)V$$

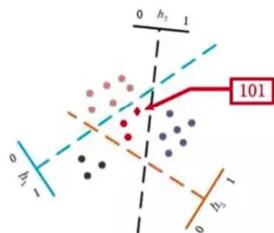
H Attention:

Hash Q and K

Standard attention within same-hash bins

Repeat a few times to increase

probability of key in the same bin



to increase the probability of
finding Q and K in the same bin.

plearning.ai

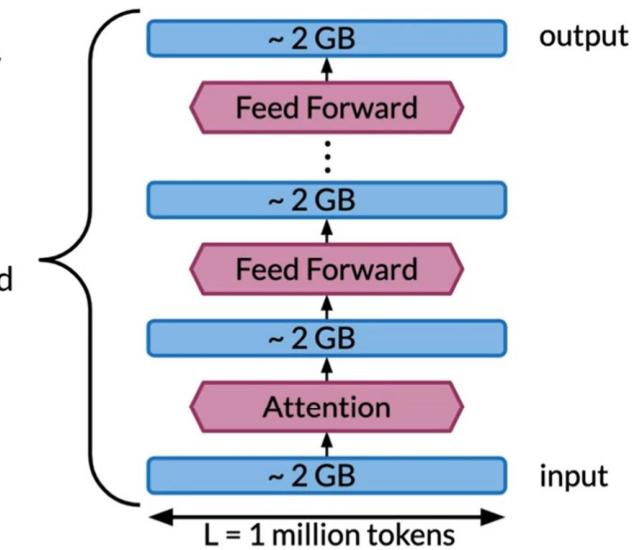
4:54

Memory Efficiency

12 x Attention
12 x Feed-Forward

plearning.ai

2:00



The size of the inputs and outputs for
each layer will also be 2GB.

- Attention

Standard Attention:

$$A(Q, K, V) = \text{softmax}(QK^T)V$$

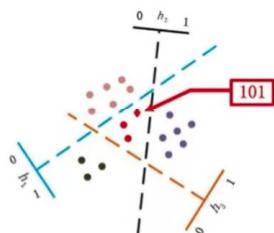
H Attention:

Hash Q and K

Standard attention within same-hash bins

Repeat a few times to increase

probability of key in the same bin



You can repeat this process multiple times

plearning.ai

3:12

- Attention

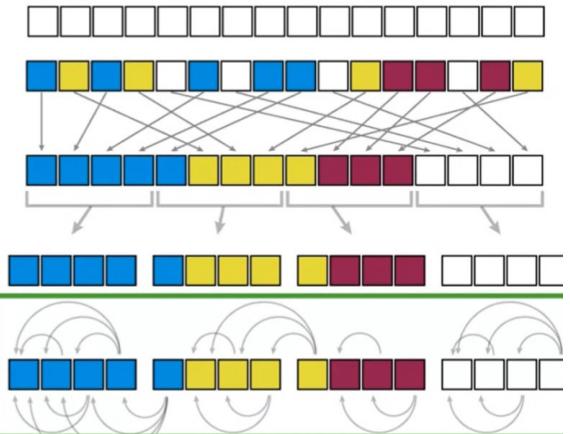
Sequence of Queries = Keys

LSH bucketing

Sort by LSH bucket

Chunk sorted sequence
to parallelize

Attend within same bucket of
own chunk and previous chunk



And the adjacent chunks discover the case
with a hash bucket that is split over

plearning.ai

2:00

Inversible layers equations

Standard Transformer:

$$y_a = x + \text{Attention}(x)$$

$$y_b = y_a + \text{FeedFwd}(y_a)$$

Inversible:

$$y_1 = x_1 + \text{Attention}(x_2)$$

$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

compute x_1, x_2 from y_1, y_2 :

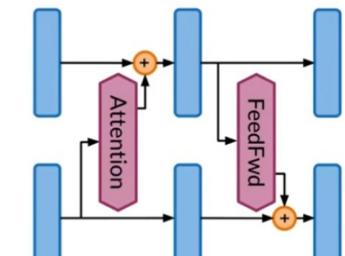
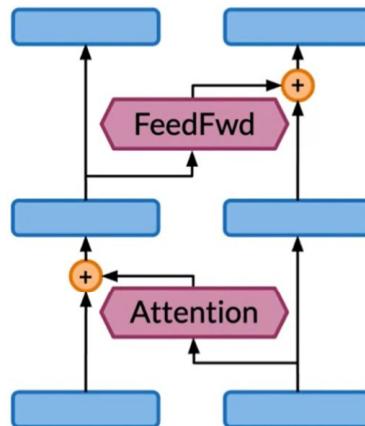
$$x_1 = y_1 - \text{Attention}(x_2)$$

$$x_2 = y_2 - \text{FeedFwd}(y_1)$$

x_1 equals y_1 minus
attention of x_2

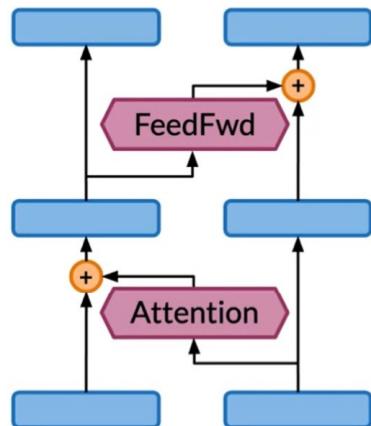
plearning.ai

Inversible layers equations



something like this,
so that information

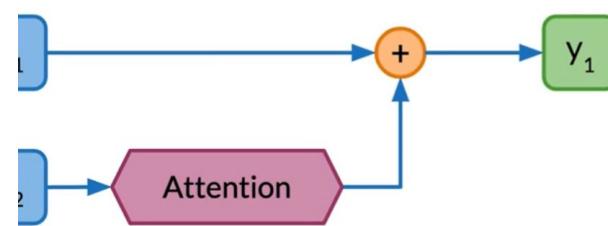
Inversible layers



about caching for
the backwards pass.

plearning.ai

Inversible layers equations



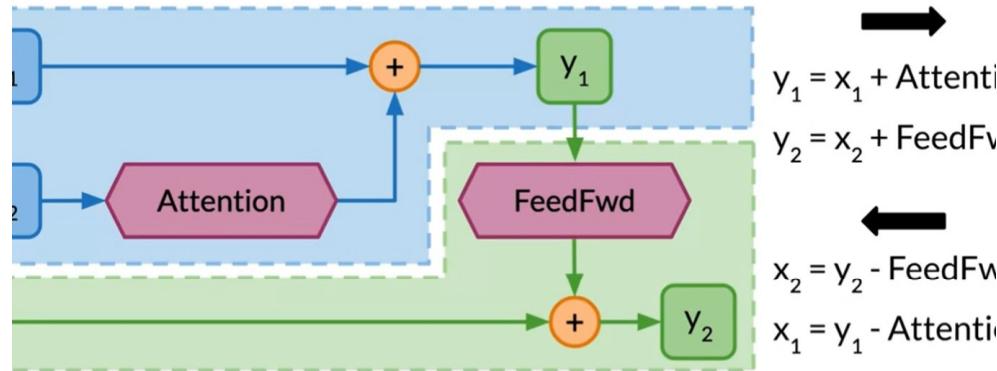
$$y_1 = x_1 + \text{Attention}(x_2)$$

$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

Restating the equations
I showed you earlier,

plearning.ai

versible layers equations



plearning.ai

going to calculate
x₂ before x₁,

former

e Reversible Transformer



L = 1 million tokens

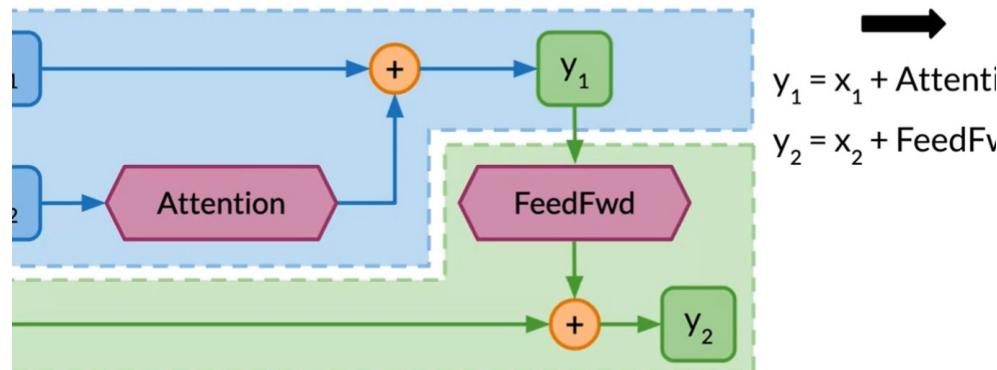


1 GPU
(16 GB)

plearning.ai

That's enough to fit an entire book,

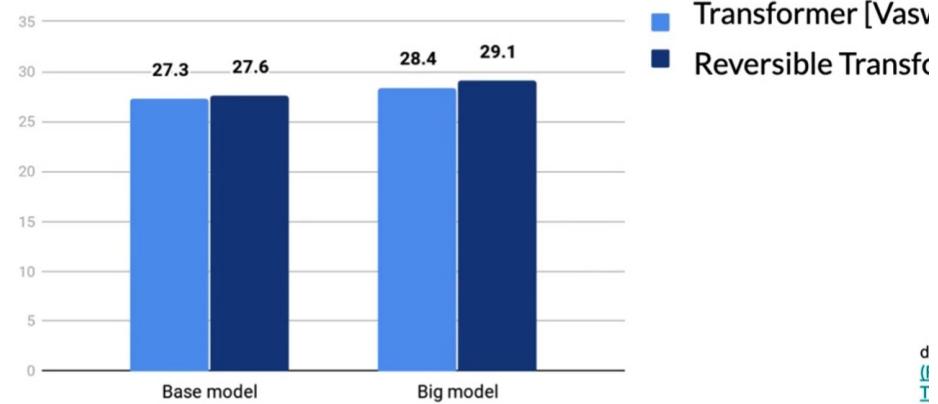
versible layers equations



plearning.ai

That's a forward pass for irreversible residual block.

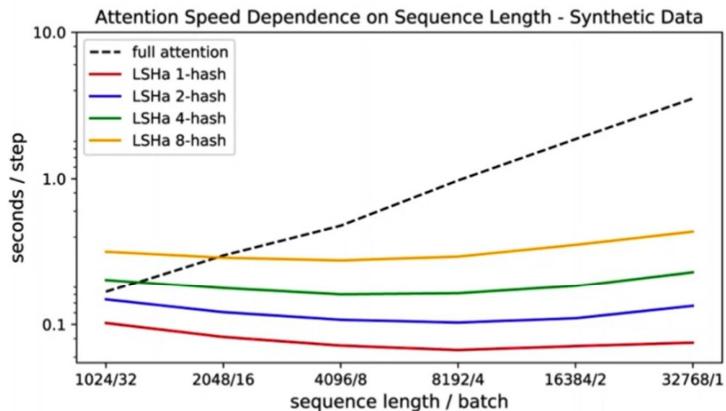
versible Transformer: BLEU Scores



plearning.ai

this can be applied anywhere you use a transformer model.

former



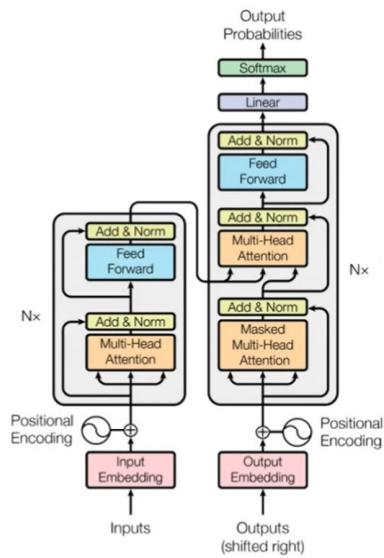
longer than fewer hashes,

plearning.ai

former

LSH Attention

Reversible Layers



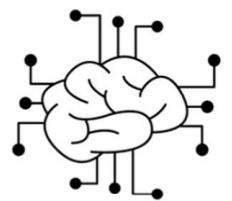
to more efficiently use
the memory available.

plearning.ai

line

Neural networks and forward propagation

Structure for sentiment analysis



plearning.ai