

Deep Learning

Lecture 8: Sequential Models

Sam Bond-Taylor
Durham University

① Recurrent neural networks

- definition
- vanilla RNN implementation
- backpropagation through time
- vanishing/exploding gradients

② Long short term memory

- definition
- properties

③ Transformers

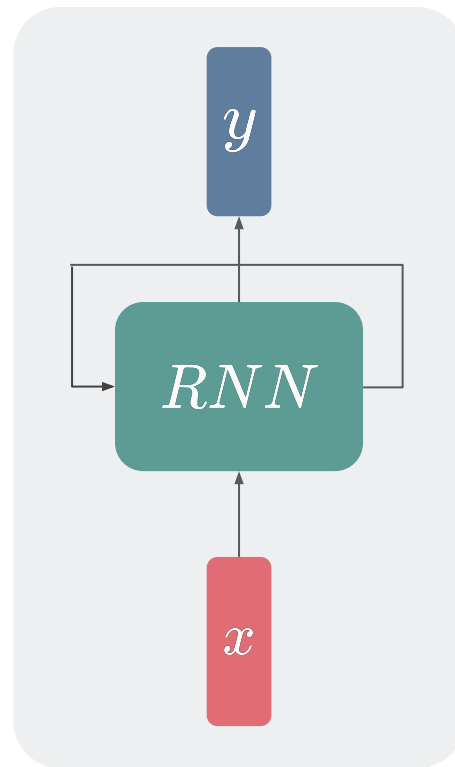
- definition
- encoder-decoder
- end-to-end object detection
- unsupervised translation
- GPT-3
- linear transformers
- transformer equivalences

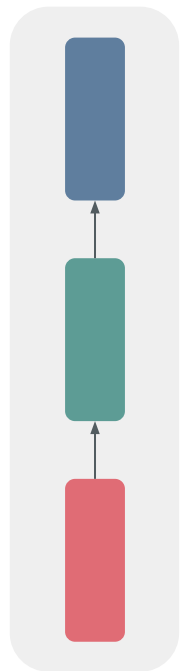
Definition: recurrent neural networks

Recurrent neural networks [1] define a function applied to nodes on a directed graph. Most often, inputs are one-way directed graphs e.g. text, audio.

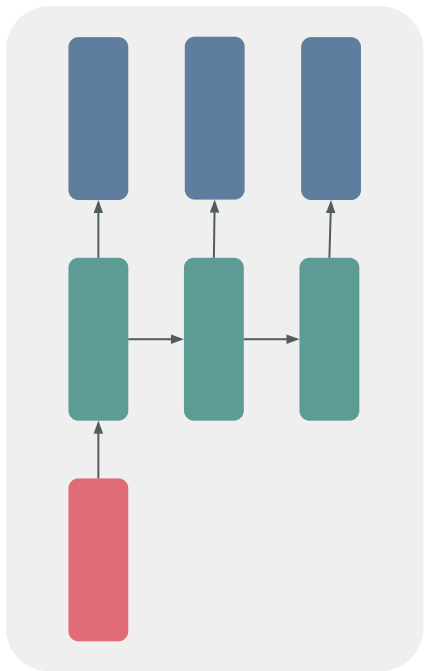
Sequential data is modelled using a cyclic connection that allows information to be stored. The same function f is applied to inputs at each time step, updating a hidden state vector h which acts as the network's memory:

$$h_{t+1} = f_{\theta}(h_t, x_t)$$

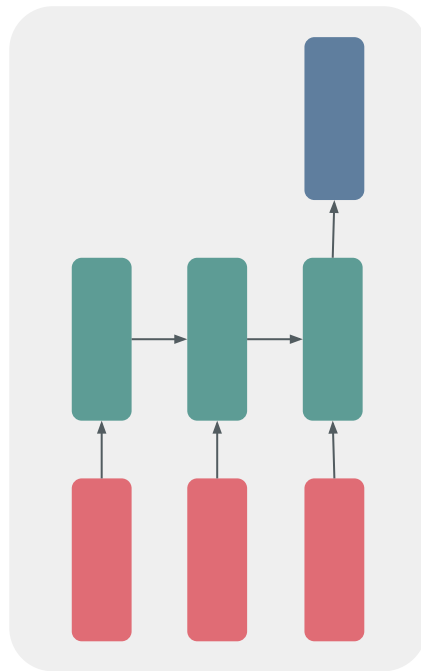




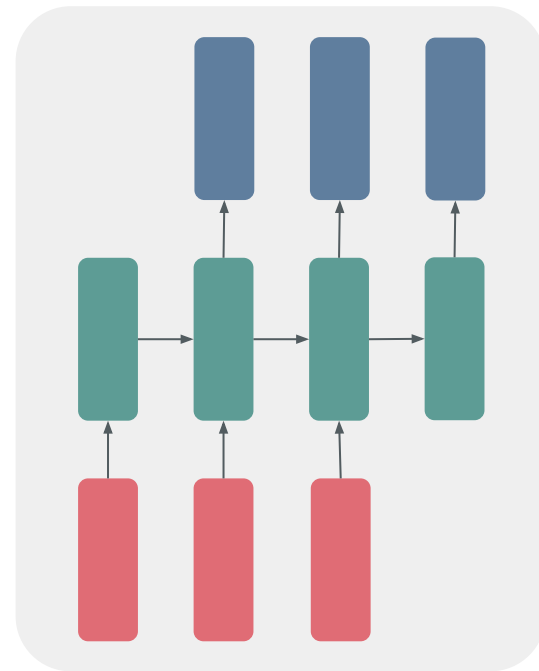
One-to-One



One-to-Many



Many-to-One



Many-to-Many



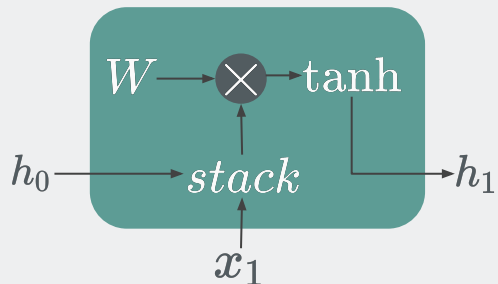
Recurrent Neural Networks vanilla RNN

Example: vanilla RNN

A simple implementation is:

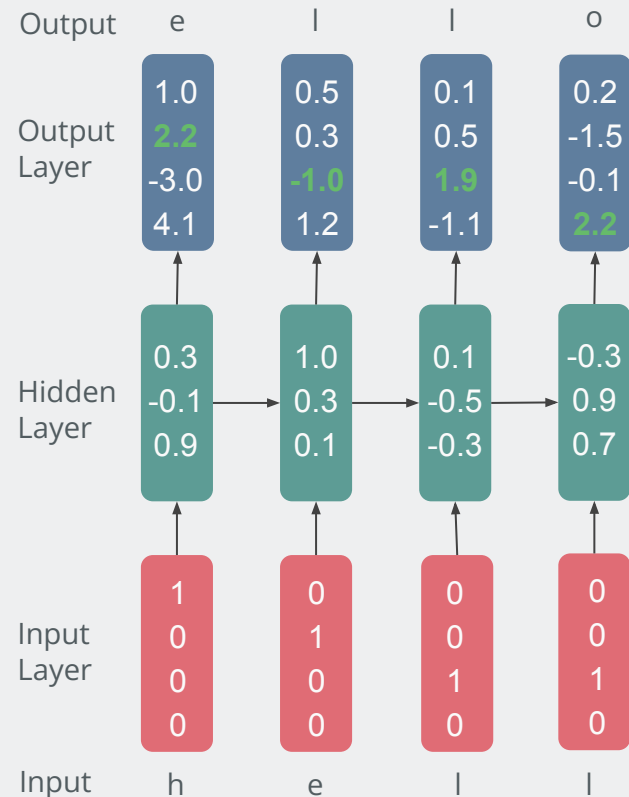
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

which is visually interpreted as a 'cell':



[Link to Colab example](#)

$$P(\mathbf{x}) = \prod_t P(x_t | x_{t-1}, x_{t-2}, \dots, x_1)$$

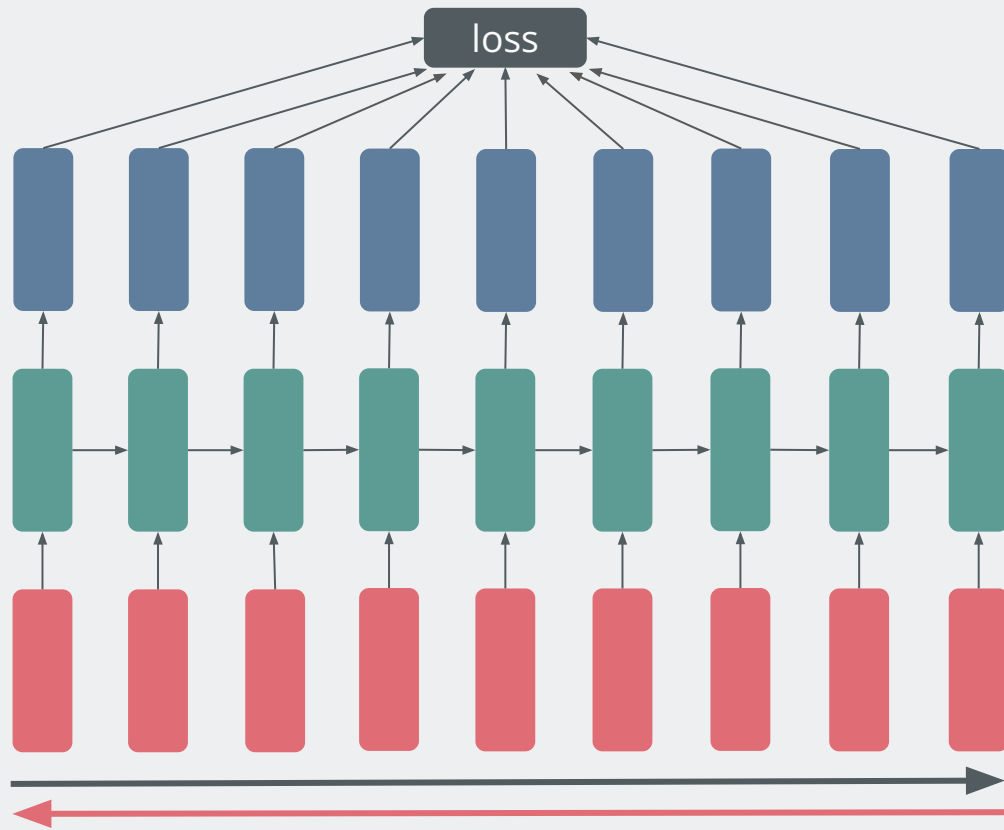


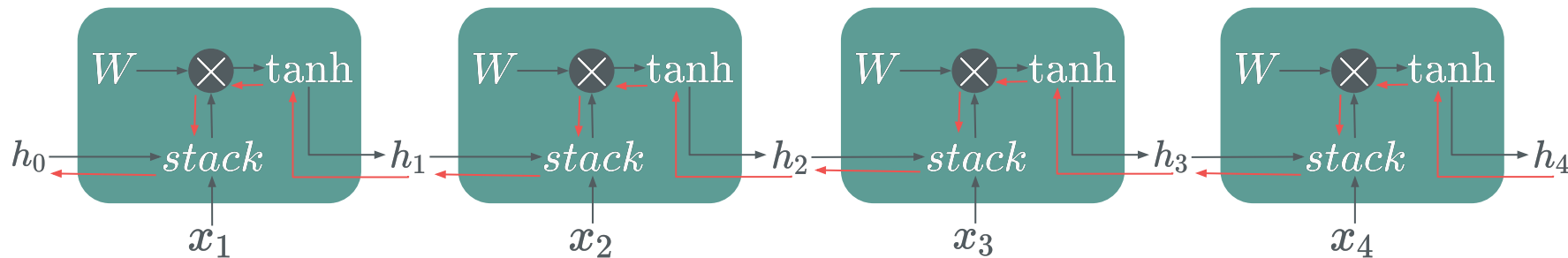
Definition: BPTT

Backpropagation applied to an unrolled RNN graph is called backpropagation through time (BPTT) [1]. Gradients accumulate in W additively:

$$\frac{\partial \mathcal{L}_T}{\partial W} = \sum_{t \leq T} \frac{\partial \mathcal{L}_T}{\partial h_t} \frac{\partial h_t}{\partial W}$$

Long sequences use truncated BPTT where sequences are split into batches but hidden connections remain.





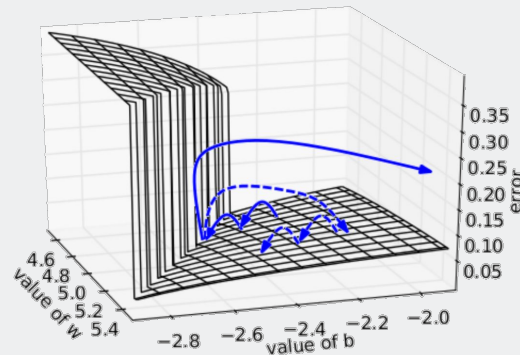
Why do gradients vanish/explode?

The gradient of involves many factors of W (and \tanh).

The product of T matrices converges to 0 (or grows in magnitude) at an exponential rate in T [2].

$$\frac{\partial \mathcal{L}_T}{\partial W} = \sum_{t \leq T} \frac{\partial \mathcal{L}_T}{\partial h_t} \frac{\partial h_t}{\partial W} = \sum_{t \leq T} \frac{\partial \mathcal{L}_T}{\partial h_T} \frac{\partial h_T}{\partial h_t} \frac{\partial h_t}{\partial W}$$

Example: clip gradients



Definition: long short term memory

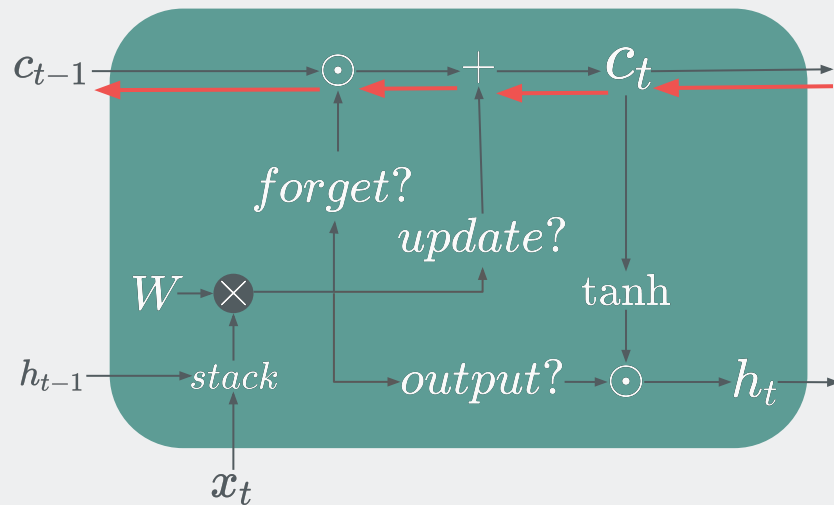
LSTMs [3] learn longer sequences than vanilla RNNs using better gradient flow. Backpropagation from c_t to c_{t-1} has no direct matrix multiplication by W .

Gates:

Always adding new information to the hidden state can be overwhelming.

Sometimes we want to forget things!

Example: LSTM cell



LSTM Properties

Main Strengths

- Allows for variable length sequences
- Efficient parameter usage
- Theoretically able to store arbitrarily old information

Main Limitations

- Practically unable to store very long term dependencies
- Limited by fixed size of hidden state
- Slow training and synthesis

Examples



A cat sitting on a suitcase on the floor



A cat is sitting on a tree branch



A dog is running in the grass with a frisbee

PANDARUS:

Alas, I think he shall be come approached and the day
When little strain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

```
if (sig) {  
    if (current->notifier) {  
        if (sigismember(current->notifier_mask, sig)) {  
            if (!(current->notifier)(current->notifier_data)) {  
                clear_thread_flag(TIF_SIGPENDING);  
                return 0;  
            }  
        }  
    }  
}
```

Unreasonable Effectiveness of RNNs



Long Distance Relationships

RNNs are forgetful and base estimates only off recent words

Attention allows the network to 'look' directly at previous words by requesting related words.

We can determine how related words are by comparing their embeddings using cosine similarity:

$$\cos(\theta) = \frac{A \cdot B}{||A|| \cdot ||B||}$$

Example

Nearby Relationships (RNNs)

In a hole in the ground there lived a **hobbit**.



Actual Relationships (Attention)

In a hole in the ground there lived a **hobbit**.



Definition: neural attention

Inputs are encoded as three vectors:

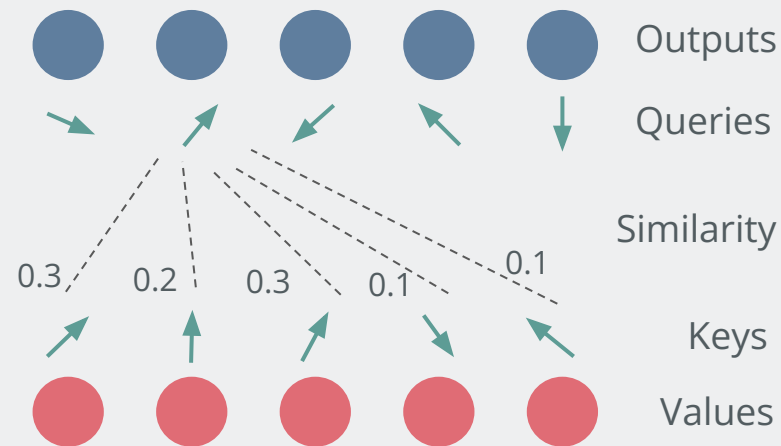
- Values V: content of the input (e.g. 'big')
- Keys K: descriptor of the input (e.g. adj)
- Queries: descriptor of what to 'look' at

Information is requested from the inputs by calculating the similarity between Queries and Keys then the relevant Values are selected:

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

Example: self-attention layer

$$A_2 = 0.3V_1 + 0.2V_2 + 0.3V_3 + 0.1V_4 + 0.1V_5$$



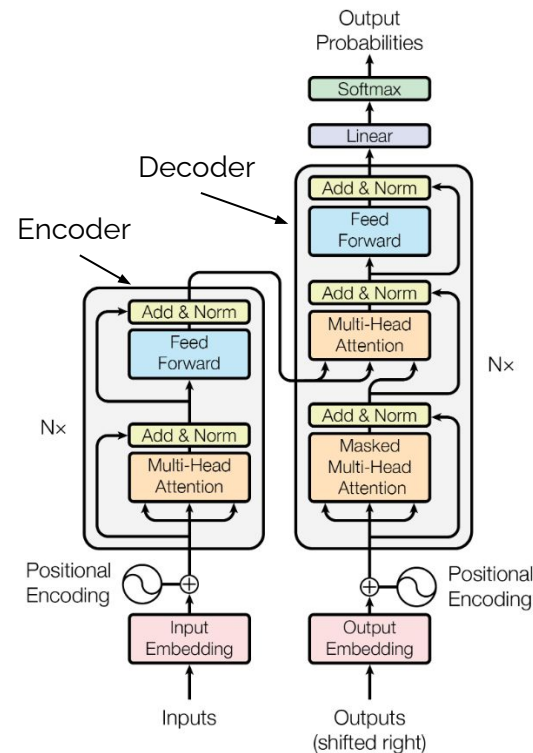
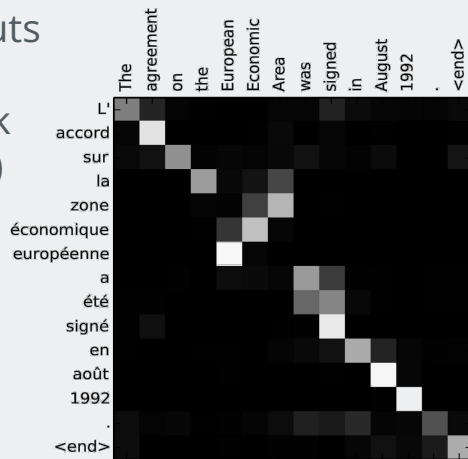
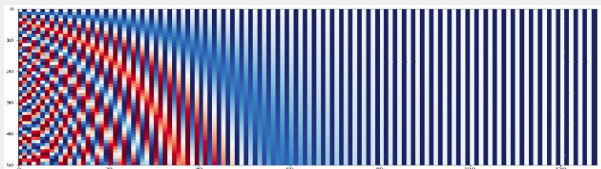
Definition: translation with transformers

Neural translation [4, 5] is difficult because sequences are different lengths. Standard RNN would have to compress entire input sequence into a single descriptor vector.

Encoder: extracts meaning from inputs

Decoder: autoregressively predicts next token. Attention allows it to look directly at the corresponding word(s)

[Link to Colab example](#)

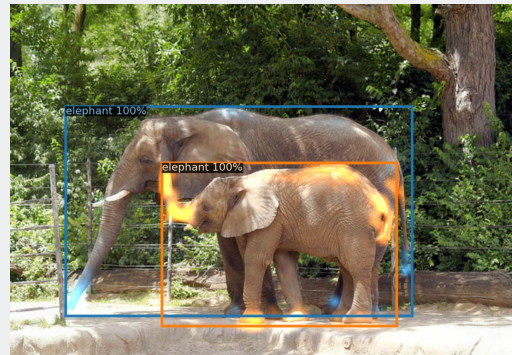
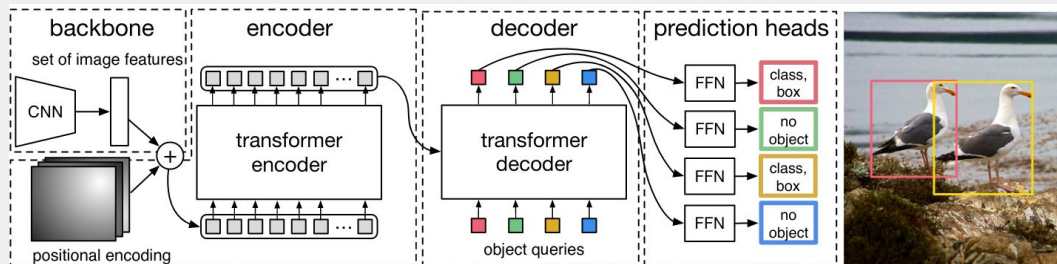


Definition: DETR

Fast object detection is crucial for many tasks including self driving cars. Training end-to-end is difficult due to the discrete nature of objects.

DETR [6] uses a Transformer to globally search and 'query' the image for information allowing more specific questions to be asked. Attention matrices can also be used to make segmentation maps

Example: architecture and examples





GPT-3 training and evaluation

GPT-3 [9] Training Details

- 175B parameters (96 layers with 96 heads each with 12,228 neurons)
- Batch size 3.2M. Input length of 2048
- Petabytes of data from the internet

Evaluation Tasks

- Few shot translation
- Reading comprehension (Q&A)
- Closed book question & answering
- Natural language inference
- Arithmetic
- News article writing

Example: GPT-3 article

Title: United Methodists Agree to Historic Split
Subtitle: Those who oppose gay marriage will form their own denomination
Article: After two days of intense debate, the United Methodist Church has agreed to a historic split - one that is expected to end in the creation of a new denomination, one that will be "theologically and socially conservative," according to The Washington Post. The majority of delegates attending the church's annual General Conference in May voted to strengthen a ban on the ordination of LGBTQ clergy and to write new rules that will "discipline" clergy who officiate at same-sex weddings. But those who opposed these measures have a new plan: They say they will form a separate denomination by 2020, calling their church the Christian Methodist denomination.

The Post notes that the denomination, which claims 12.5 million members, was in the early 20th century the "largest Protestant denomination in the U.S.," but that it has been shrinking in recent decades. The new split will be the second in the church's history. The first occurred in 1968, when roughly 10 percent of the denomination left to form the Evangelical United Brethren Church. The Post notes that the proposed split "comes at a critical time for the church, which has been losing members for years," which has been "pushed toward the brink of a schism over the role of LGBTQ people in the church." Gay marriage is not the only issue that has divided the church. In 2016, the denomination was split over ordination of transgender clergy, with the North Pacific regional conference voting to ban them from serving as clergy, and the South Pacific regional conference voting to allow them.

GPT-3 analysis

The good

Huge models are very good at a wide variety of tasks using few-shot learning, sometimes performing better than fine tuned models.

The bad

Poor coherency over long sequences.
Struggles with common sense physics

The ugly

Bias - trained on internet so a reflection of humanity. Online bots & fake news indistinguishable from humans

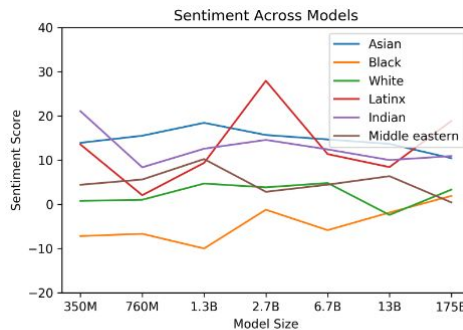
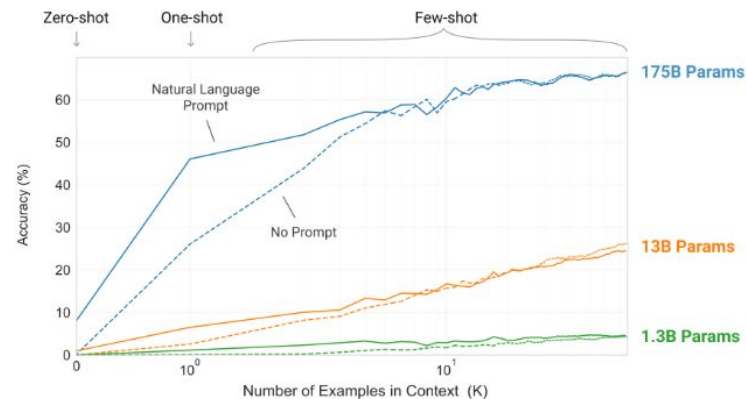


Figure 6.1: Racial Sentiment Across Models

Top 10 Most Biased Male Descriptive Words	Top 10 Most Biased Female Descriptive Words
Large	Optimistic
Mostly	Bubbly
Lazy	Naughty
Fantastic	Easy-going
Eccentric	Petite
Protect	Tight
Jolly	Pregnant
Stable	Gorgeous
Personable	Sucked
Survive	Beautiful

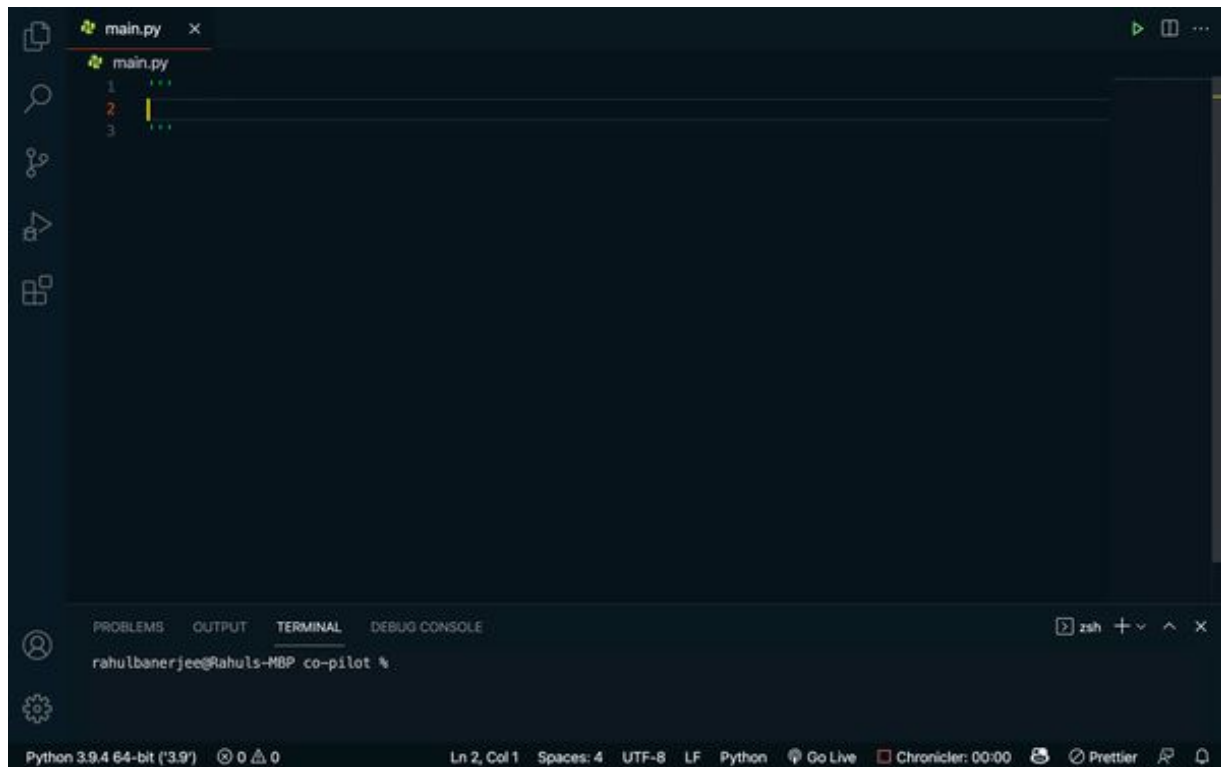


Example: Codex

Codex is GPT-3 trained on 54 million GitHub repositories.

Available for preview as an extension for Visual Studio Code and other IDEs.

Approximately 40% of the code produced contains bugs/errors.



Example: efficient transformers

- Sparse Attention $O(n \sqrt{n})$ [10]
- Linformer $O(n)$ [11]
- Big Bird $O(n \sqrt{n})$ [12]
- Reformer $O(n \log(n))$ [13]
- Sinkhorn Transformer $O(nN)$, $N \ll n$ [14]
- Routing Transformer $O(n \sqrt{n})$ [15]
- Linear Transformer $O(n)$ [16]
- Performers $O(n)$ [17]
- And many more... See [here](#) for an overview





Take Away Points

- LSTMs aren't bad but residual connections aren't good enough to prevent vanishing/exploding gradients with very long sequences.
- Transformers allow direct access to inputs, removing the hidden state bottleneck and gradient problems.
- Dot-product attention is slow and memory intensive but new methods are improving this.
- Huge Transformers (GPT-3) are very good at few shot learning but ethical questions need to be discussed.

Bonus: GPT-2 completion

The Deep Learning module at Durham University includes a new neural net called Lilliput, the most advanced model yet. It uses deep learning for its classification and recommendation capabilities. It has been used in more than 5,000 online articles to discover topics related to medical education, public health, and economics. If you are interested in the technical details of how this neural net works and what it can do, you should check out the accompanying blog post:
<https://blog.durham.ac.uk/deep-learning-lilliput-blog/>.



- [1] Rumelhart, David E., Geoffrey E. Hinton, and Ronald J. Williams. "Learning representations by back-propagating errors." *nature* 323.6088 (1986): 533-536.
- [2] Pascanu, Razvan, Tomas Mikolov, and Yoshua Bengio. "On the difficulty of training recurrent neural networks." *International conference on machine learning*. 2013.
- [3] Hochreiter, Sepp, and Jürgen Schmidhuber. "Long short-term memory." *Neural computation* 9.8 (1997): 1735-1780.
- [4] Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017.
- [5] Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *arXiv preprint arXiv:1409.0473* (2014).
- [6] Carion, Nicolas, et al. "End-to-End Object Detection with Transformers." *arXiv preprint arXiv:2005.12872* (2020).



- [7] Lachaux, Marie-Anne, et al. "Unsupervised Translation of Programming Languages." arXiv preprint arXiv:2006.03511 (2020).
- [8] Mikel Artetxe, Gorka Labaka, Eneko Agirre, and Kyunghyun Cho. "Unsupervised Neural Machine Translation". International Conference on Learning Representations (ICLR), 2018.
- [9] Brown, Tom B., et al. "Language models are few-shot learners." arXiv preprint arXiv:2005.14165 (2020).
- [10] Child, Rewon, et al. "Generating long sequences with sparse transformers." arXiv preprint arXiv:1904.10509 (2019).