Recurrent Neural Networks - LSTM



• Input: $\vec{x} = (x_1, \dots, x_M)$

• Hidden:
$$\vec{h} = (h_1, \dots, h_H)$$

• Output:
$$\vec{y} = (y_1, \dots, y_N)$$

RNN example



Activation function:

$$\sigma(\xi) = \begin{cases} 1 & \xi \ge 0 \\ 0 & \xi < 0 \end{cases}$$

RNN example



Activation function:

$$\sigma(\xi) = \begin{cases} 1 & \xi \ge 0 \\ 0 & \xi < 0 \end{cases}$$

RNN example



- *M* inputs: $\vec{x} = (x_1, \dots, x_M)$
- *H* hidden neurons: $\vec{h} = (h_1, \dots, h_H)$
- N output neurons: $\vec{y} = (y_1, \dots, y_N)$
- Weights:
 - $U_{kk'}$ from input $x_{k'}$ to hidden h_k
 - $W_{kk'}$ from hidden $h_{k'}$ to hidden h_k
 - $V_{kk'}$ from hidden $h_{k'}$ to output y_k



RNN – formally

• Input sequence:
$$\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$$

 $\vec{x}_t = (x_{t1}, \ldots, x_{tM})$

RNN – formally

• Input sequence:
$$\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$$

$$\vec{x}_t = (x_{t1}, \ldots, x_{tM})$$

• Hidden sequence:
$$\mathbf{h} = \vec{h}_0, \vec{h}_1, \dots, \vec{h}_T$$

$$ec{h}_t = (h_{t1}, \ldots, h_{tH})$$

We have $\vec{h}_0 = (0, \dots, 0)$ and

$$\vec{h}_{tk} = \sigma \left(\sum_{k'=1}^{M} U_{kk'} x_{tk'} + \sum_{k'=1}^{H} W_{kk'} h_{(t-1)k'} \right)$$

RNN – formally

• Input sequence:
$$\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$$

$$\vec{x}_t = (x_{t1}, \ldots, x_{tM})$$

• Hidden sequence:
$$\mathbf{h} = \vec{h}_0, \vec{h}_1, \dots, \vec{h}_T$$

$$\vec{h}_t = (h_{t1}, \ldots, h_{tH})$$

We have $\vec{h}_0 = (0, \dots, 0)$ and

$$\vec{h}_{tk} = \sigma \left(\sum_{k'=1}^{M} U_{kk'} x_{tk'} + \sum_{k'=1}^{H} W_{kk'} h_{(t-1)k'} \right)$$

• Output sequence: $\mathbf{y} = \vec{y}_1, \dots, \vec{y}_T$

$$ec{\mathbf{y}_t} = (\mathbf{y}_{t1}, \dots, \mathbf{y}_{tN})$$

where $\mathbf{y}_{tk} = \sigma \left(\sum_{k'=1}^{H} \mathbf{V}_{kk'} h_{tk'} \right)$.

RNN – in matrix form

• Input sequence: $\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$

RNN – in matrix form

- Input sequence: $\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$
- Hidden sequence: $\mathbf{h} = \vec{h}_0, \vec{h}_1, \dots, \vec{h}_T$ where

$$\vec{h}_0 = (0,\ldots,0)$$

and

$$\vec{h}_t = \sigma(U\vec{x}_t + W\vec{h}_{t-1})$$

RNN – in matrix form

- Input sequence: $\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$
- Hidden sequence: $\mathbf{h} = \vec{h}_0, \vec{h}_1, \dots, \vec{h}_T$ where

$$\vec{h}_0=(0,\ldots,0)$$

and

$$\vec{h}_t = \sigma(U\vec{x}_t + W\vec{h}_{t-1})$$

• Output sequence: $\mathbf{y} = \vec{y}_1, \dots, \vec{y}_T$ where

 $\mathbf{y}_t = \sigma(\mathbf{V}\mathbf{h}_t)$

- \vec{h}_t is the memory of the network, captures what happened in all previous steps (with decaying quality).
- RNN shares weights U, V, W along the sequence. Note the similarity to convolutional networks where the weights were shared spatially over images, here they are shared temporally over sequences.
- RNN can deal with sequences of variable length. Compare with MLP which accepts only fixed-dimension vectors on input.

The Task: Design a recurrent network with a single hidden layer which works as a binary adder.

Example of behavior: Input two binary numbers, e.g., 111 and 101 (we assume that the least significant bit is on the *left*).

The input of the network will be: (1, 1), (1, 0), (1, 1)

The output is supposet to be: 0, 0, 1 (we ignore the carry at the end).

Training set

$$\mathcal{T} = \left\{ (\mathbf{x}_1, \mathbf{d}_1), \dots, (\mathbf{x}_p, \mathbf{d}_p) \right\}$$

here

• each $\mathbf{x}_{\ell} = \vec{x}_{\ell 1}, \dots, \vec{x}_{\ell T_{\ell}}$ is an input sequence,

• each $\mathbf{d}_{\ell} = \vec{d}_{\ell 1}, \dots, \vec{d}_{\ell T_{\ell}}$ is an expected output sequence. Here each $\vec{x}_{\ell t} = (x_{\ell t 1}, \dots, x_{\ell t M})$ is an input vector and each $\vec{d}_{\ell t} = (d_{\ell t 1}, \dots, d_{\ell t N})$ is an expected output vector. In what follows I will consider a training set with a **single** element (\mathbf{x}, \mathbf{d}) . I.e. drop the index ℓ and have

•
$$\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$$
 where $\vec{x}_t = (x_{t1}, \dots, x_{tM})$
• $\mathbf{d} = \vec{d}_1, \dots, \vec{d}_T$ where $\vec{d}_t = (d_{t1}, \dots, d_{tN})$

The squared error of (\mathbf{x}, \mathbf{d}) is defined by

$$E_{(\mathbf{x},\mathbf{d})} = \sum_{t=1}^{T} \sum_{k=1}^{N} \frac{1}{2} (y_{tk} - d_{tk})^2$$

Recall that we have a sequence of network outputs $\mathbf{y} = \vec{y}_1, \dots, \vec{y}_T$ and thus y_{tk} is the *k*-th component of \vec{y}_t

Consider a single training example (**x**, **d**).

The algorithm computes a sequence of weight matrices as follows:

Consider a single training example (\mathbf{x}, \mathbf{d}) .

The algorithm computes a sequence of weight matrices as follows:

Initialize all weights randomly close to 0.

Consider a single training example (\mathbf{x}, \mathbf{d}) .

The algorithm computes a sequence of weight matrices as follows:

- Initialize all weights randomly close to 0.
- In the step ℓ + 1 (here ℓ = 0, 1, 2, ...) compute "new" weights U^(ℓ+1), V^(ℓ+1), W^(ℓ+1) from the "old" weights U^(ℓ), V^(ℓ), W^(ℓ) as follows:

$$U_{kk'}^{(\ell+1)} = U_{kk'}^{(\ell)} - \varepsilon(\ell) \cdot \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta U_{kk'}}$$
$$V_{kk'}^{(\ell+1)} = V_{kk'}^{(\ell)} - \varepsilon(\ell) \cdot \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta V_{kk'}}$$
$$W_{kk'}^{(\ell+1)} = W_{kk'}^{(\ell)} - \varepsilon(\ell) \cdot \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta W_{kk'}}$$

Consider a single training example (\mathbf{x}, \mathbf{d}) .

The algorithm computes a sequence of weight matrices as follows:

- Initialize all weights randomly close to 0.
- In the step ℓ + 1 (here ℓ = 0, 1, 2, ...) compute "new" weights U^(ℓ+1), V^(ℓ+1), W^(ℓ+1) from the "old" weights U^(ℓ), V^(ℓ), W^(ℓ) as follows:

$$U_{kk'}^{(\ell+1)} = U_{kk'}^{(\ell)} - \varepsilon(\ell) \cdot \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta U_{kk'}}$$
$$V_{kk'}^{(\ell+1)} = V_{kk'}^{(\ell)} - \varepsilon(\ell) \cdot \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta V_{kk'}}$$
$$W_{kk'}^{(\ell+1)} = W_{kk'}^{(\ell)} - \varepsilon(\ell) \cdot \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta W_{kk'}}$$

The above is THE learning algorithm that modifies weights!

Backpropagation

Computes the derivatives of *E*, no weights are modified!

Backpropagation

Computes the derivatives of *E*, no weights are modified!

$$\frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta U_{kk'}} = \sum_{t=1}^{T} \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta h_{tk}} \cdot \sigma' \cdot x_{tk'} \qquad \qquad k' = 1, \dots, M$$
$$\frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta V_{kk'}} = \sum_{t=1}^{T} \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta y_{tk}} \cdot \sigma' \cdot h_{tk'} \qquad \qquad k' = 1, \dots, H$$
$$\frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta W_{kk'}} = \sum_{t=1}^{T} \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta h_{tk}} \cdot \sigma' \cdot h_{(t-1)k'} \qquad \qquad k' = 1, \dots, H$$

Backpropagation

Computes the derivatives of *E*, no weights are modified!

$$\frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta U_{kk'}} = \sum_{t=1}^{T} \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta h_{tk}} \cdot \sigma' \cdot x_{tk'} \qquad \qquad k' = 1, \dots, M$$
$$\frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta V_{kk'}} = \sum_{t=1}^{T} \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta y_{tk}} \cdot \sigma' \cdot h_{tk'} \qquad \qquad k' = 1, \dots, H$$
$$\frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta W_{kk'}} = \sum_{t=1}^{T} \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta h_{tk}} \cdot \sigma' \cdot h_{(t-1)k'} \qquad \qquad k' = 1, \dots, H$$

Backpropagation:

$$\frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta y_{tk}} = y_{tk} - d_{tk} \quad (\text{assuming squared error})$$

$$\frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta h_{tk}} = \sum_{k'=1}^{N} \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta y_{tk'}} \cdot \sigma' \cdot V_{k'k} + \sum_{k'=1}^{H} \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta h_{(t+1)k'}} \cdot \sigma' \cdot W_{k'k}$$

Long-term dependencies

$$\frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta h_{tk}} = \sum_{k'=1}^{N} \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta y_{tk'}} \cdot \sigma' \cdot V_{k'k} + \sum_{k'=1}^{H} \frac{\delta E_{(\mathbf{x},\mathbf{d})}}{\delta h_{(t+1)k'}} \cdot \sigma' \cdot W_{k'k}$$

► Unless $\sum_{k'=1}^{H} \sigma' \cdot W_{k'k} \approx 1$, the gradient either vanishes, or explodes.

- For a large T (long-term dependency), the gradient "deeper" in the past tends to be too small (large).
- A solution: LSTM

LSTM is currently a bit obsolete. The main idea is to decompose W into several matrices, each responsible for a different task. One is concerned about memory, one is concerned about the output at each step, etc.

LSTM

$$\vec{h}_{t} = \vec{o}_{t} \circ \sigma_{h}(\vec{C}_{t}) \qquad \text{output}$$

$$\vec{C}_{t} = \vec{f}_{t} \circ \vec{C}_{t-1} + \vec{i}_{t} \circ \tilde{C}_{t} \qquad \text{memory}$$

$$\tilde{C}_{t} = \sigma_{h}(W_{C} \cdot \vec{h}_{t-1} + U_{C} \cdot \vec{x}_{t}) \qquad \text{new memory contents}$$

 $\vec{o}_{t} = \sigma_{g}(W_{o} \cdot \vec{h}_{t-1} + U_{o} \cdot \vec{x}_{t}) \qquad \text{output gate}$ $\vec{f}_{t} = \sigma_{g}(W_{f} \cdot \vec{h}_{t-1} + U_{f} \cdot \vec{x}_{t}) \qquad \text{forget gate}$ $\vec{i}_{t} = \sigma_{g}(W_{i} \cdot \vec{h}_{t-1} + U_{i} \cdot \vec{x}_{t}) \qquad \text{input gate}$

- is the component-wise product of vectors
- is the matrix-vector product
- σ_h hyperbolic tangents (applied component-wise)
- σ_g logistic sigmoid (aplied component-wise)

RNN vs LSTM





$$\vec{h}_{t} = \vec{o}_{t} \circ \sigma_{h}(\vec{C}_{t})$$

$$\Rightarrow \vec{C}_{t} = \vec{f}_{t} \circ \vec{C}_{t-1} + \vec{i}_{t} \circ \tilde{C}_{t}$$

$$\tilde{C}_{t} = \sigma_{h}(W_{C} \cdot \vec{h}_{t-1} + U_{C} \cdot \vec{x}_{t})$$

$$\vec{o}_{t} = \sigma_{g}(W_{o} \cdot \vec{h}_{t-1} + U_{o} \cdot \vec{x}_{t})$$

$$\vec{f}_{t} = \sigma_{g}(W_{f} \cdot \vec{h}_{t-1} + U_{f} \cdot \vec{x}_{t})$$

$$\vec{i}_{t} = \sigma_{g}(W_{i} \cdot \vec{h}_{t-1} + U_{i} \cdot \vec{x}_{t})$$



$$\vec{h}_{t} = \vec{o}_{t} \circ \sigma_{h}(\vec{C}_{t})$$
$$\vec{C}_{t} = \vec{f}_{t} \circ \vec{C}_{t-1} + \vec{i}_{t} \circ \tilde{C}_{t}$$
$$\tilde{C}_{t} = \sigma_{h}(W_{C} \cdot \vec{h}_{t-1} + U_{C} \cdot \vec{x}_{t})$$
$$\vec{o}_{t} = \sigma_{g}(W_{o} \cdot \vec{h}_{t-1} + U_{o} \cdot \vec{x}_{t})$$

$$\vec{o}_t = \sigma_g(W_o \cdot \vec{h}_{t-1} + U_o \cdot \vec{x}_t)$$

$$\Rightarrow \vec{f}_t = \sigma_g(W_f \cdot \vec{h}_{t-1} + U_f \cdot \vec{x}_t)$$

$$\vec{i}_t = \sigma_g(W_i \cdot \vec{h}_{t-1} + U_i \cdot \vec{x}_t)$$



$$\vec{h}_{t} = \vec{o}_{t} \circ \sigma_{h}(\vec{C}_{t})$$
$$\vec{C}_{t} = \vec{f}_{t} \circ \vec{C}_{t-1} + \vec{i}_{t} \circ \tilde{C}_{t}$$
$$\Rightarrow \tilde{C}_{t} = \sigma_{h}(W_{C} \cdot \vec{h}_{t-1} + U_{C} \cdot \vec{x}_{t})$$
$$\vec{o}_{t} = \sigma_{g}(W_{o} \cdot \vec{h}_{t-1} + U_{o} \cdot \vec{x}_{t})$$
$$\vec{f}_{t} = \sigma_{g}(W_{f} \cdot \vec{h}_{t-1} + U_{f} \cdot \vec{x}_{t})$$
$$\Rightarrow \vec{i}_{t} = \sigma_{g}(W_{i} \cdot \vec{h}_{t-1} + U_{i} \cdot \vec{x}_{t})$$



$$\vec{h}_{t} = \vec{o}_{t} \circ \sigma_{h}(\vec{C}_{t})$$

$$\Rightarrow \vec{C}_{t} = \vec{f}_{t} \circ \vec{C}_{t-1} + \vec{i}_{t} \circ \tilde{C}_{t}$$

$$\Rightarrow \tilde{C}_{t} = \sigma_{h}(W_{C} \cdot \vec{h}_{t-1} + U_{C} \cdot \vec{x}_{t})$$

$$\vec{o}_{t} = \sigma_{g}(W_{o} \cdot \vec{h}_{t-1} + U_{o} \cdot \vec{x}_{t})$$

$$\Rightarrow \vec{f}_{t} = \sigma_{g}(W_{f} \cdot \vec{h}_{t-1} + U_{f} \cdot \vec{x}_{t})$$

$$\Rightarrow \vec{i}_t = \sigma_g(W_i \cdot \vec{h}_{t-1} + U_i \cdot \vec{x}_t)$$



$$\Rightarrow \vec{h}_{t} = \vec{o}_{t} \circ \sigma_{h}(\vec{C}_{t})$$
$$\vec{C}_{t} = \vec{f}_{t} \circ \vec{C}_{t-1} + \vec{i}_{t} \circ \tilde{C}_{t}$$
$$\tilde{C}_{t} = \sigma_{h}(W_{C} \cdot \vec{h}_{t-1} + U_{C} \cdot \vec{x}_{t})$$

$$\Rightarrow \vec{o}_t = \sigma_g(W_o \cdot \vec{h}_{t-1} + U_o \cdot \vec{x}_t)$$
$$\vec{f}_t = \sigma_g(W_f \cdot \vec{h}_{t-1} + U_f \cdot \vec{x}_t)$$
$$\vec{i}_t = \sigma_g(W_i \cdot \vec{h}_{t-1} + U_i \cdot \vec{x}_t)$$

- LSTM (almost) solves the vanishing gradient problem w.r.t. the "internal" state of the network.
- Learns to control its own memory (via forget gate).
- Revolution in machine translation and text processing.
- ... but the development goes on ...

Time-series Forecasting with LSTM

(See: https://www.tensorflow.org/tutorials/structured_data/ time_series)

- Weather time series dataset
- 14 different features such as air temperature, atmospheric pressure, and humidity
- collected every 10 minutes, beginning in 2003 (only 2009 -2016 considered in the example)

	p (mbar)	T (degC)	Tpot (K)	Tdew (degC)	rh (%)	VPmax (mbar)	VPact (mbar)	VPdef (mbar)	sh (g/kg)	H2OC (mmol/mol)	rho (g/m**3)	wv (m/s)	max. wv (m/s)	wd (deg)
5	996.50	-8.05	265.38	-8.78	94.4	3.33	3.14	0.19	1.96	3.15	1307.86	0.21	0.63	192.7
11	996.62	-8.88	264.54	-9.77	93.2	3.12	2.90	0.21	1.81	2.91	1312.25	0.25	0.63	190.3
17	996.84	-8.81	264.59	-9.66	93.5	3.13	2.93	0.20	1.83	2.94	1312.18	0.18	0.63	167.2
23	996.99	-9.05	264.34	-10.02	92.6	3.07	2.85	0.23	1.78	2.85	1313.61	0.10	0.38	240.0
29	997.46	-9.63	263.72	-10.65	92.2	2.94	2.71	0.23	1.69	2.71	1317.19	0.40	0.88	157.0

The Task: Predict the temperature for the next hour.

Weather Data



Preprocessing (omitted)

Before applying any prediction model, proper preprocessing is essential for time series data.

- Train-test Split: It is crucial to split the data into training and test sets while ensuring that the temporal order is maintained. This allows the model to learn from past input data of the training set and evaluate its performance on unseen future data.
- Handling Missing Values: Addressing missing values is crucial as gaps in the data can affect the model's performance. You can use techniques like interpolation or forward/backward filling.
- Data Normalization: Normalizing the data ensures that all features are on the same scale, preventing any single feature from dominating the model's learning process.
- Detrending: Removing the trend component from the data can help in better understanding the underlying patterns and making accurate predictions.
- Seasonal Adjustment: For data with seasonality, seasonal adjustment methods like seasonal differencing or seasonal decomposition can be applied.

Baseline model

The baseline: Predict that the temperature stays constant.



21
Simple linear model

The linear model: Consider the current values of all variables and predict the temperature using linear regression.



Linear explained - weights



Data Windowing

Assume that the samples are taken hourly (subsample the 10 minutes samples).

Consider windowed inputs to the model.

E.g., predict one hour given 6 hours from the past:





Given 24 hours in the past, predict the next hour with LSTM. A possible LSTM architecture:



The used LSTM had the memory dimension equal to 32.

LSTM forecasting



Model comparison

MAE = mean absolute error



Time-series prediction summary

- The presented approach is very basic!
- Omitted lots of important ideas:
 - Preprocessing extremely important!
 - Cross-validation tricky!
 - Classical models ARIMA etc. very deep and advanced area!
 - Proper evaluation, explainability, ... (a whole new course possible!)
- Read books, e.g.
 - Hyndman and Athanasopoulos. Forecasting: Principles and Practice. Online: https://otexts.com/fpp2/
 - Manu Joseph. Modern Time Series Forecasting with Python. Packt Publishing. 2022

RNN text generator (a little obsolete example)

Generating texts letter by letter.



- Generating Shakespeare letter by letter.
- Trained on Shakespeare's plays (4.4MB).

VIOLA: Why, Salisbury must find his flesh and thought That which I am not aps, not a man and in fire, To show the reining of the raven and the wars To grace my hand reproach within, and not a fair are hand, That Caesar and my goodly father's world; When I was heaven of presence and our fleets, We spare with hours, but cut thy council I am great, Murdered and by thy master's ready there My power to give thee but so much as hell: Some service in the noble bondman here, Would show him to her wine.

KING LEAR: O, if you were a feeble sight, the courtesy of your law, Your sight and several breath, will wear the gods With his heads, and my hands are wonder'd at the deeds, So drop upon your lordship's head, and your opinion Shall be against your honour.

Wikipedia

Hutter Prize 100MB dataset from Wikipedia (96MB) Naturalism and decision for the majority of Arab countries' capitalide was grounded by the Irish language by [[John Clair]], [[An Imperial Japanese Revolt]], associated with Guangzham's sovereignty. His generals were the powerful ruler of the Portugal in the [[Protestant Immineners]], which could be said to be directly in Cantonese Communication, which followed a ceremony and set inspired prison, training. The emperor travelled back to [[Antioch, Perth, October 25[21]] to note, the Kingdom of Costa Rica, unsuccessful fashioned the [[Thrales]], [[Cynth's Dajoard]], known in western [[Scotland]], near Italy to the conquest of India with the conflict. Copyright was the succession of independence in the slop of Syrian influence that was a famous German movement based on a more popular servicious, non-doctrinal and sexual power post. Many governments recognize the military housing of the [[Civil Liberalization and Infantry Resolution 265 National Party in Hungary]], that is sympathetic to be to the [[Punjab Resolution]] (PJS)[http: //www.humah.yahoo.com/guardian.cfm/7754800786d17551963s89.htm Official economics Adjoint for the Nazism, Montgomery was swear to advance to the resources for those Socialism's rule, was starting to signing a

Xml halucination:

```
<page>
  <title>Antichrist</title>
  <id>865</id>
  <revision>
    <id>15900676</id>
    <timestamp>2002-08-03T18:14:12Z</timestamp>
    <contributor>
      <username>Paris</username>
      <id>23</id>
    </contributor>
    <minor />
    <comment>Automated conversion</comment>
    <text xml:space="preserve">
    #REDIRECT [[Christianity]]</text>
  </revision>
</page>
```

- Algebraic geometry textbook.
- LaTeX source (16MB).
- Almost compilable.

Proof. Omitted.

Lemma 0.1. Let C be a set of the construction.

Let C be a gerber covering. Let F be a quasi-coherent sheaves of O-modules. We have to show that

 $\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$

Proof. This is an algebraic space with the composition of sheaves ${\mathcal F}$ on $X_{\acute{e}tale}$ we have

 $\mathcal{O}_X(\mathcal{F}) = \{morph_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$

where \mathcal{G} defines an isomorphism $\mathcal{F} \to \mathcal{F}$ of \mathcal{O} -modules.

Lemma 0.2. This is an integer Z is injective.

Proof. See Spaces, Lemma ??.

Lemma 0.3. Let S be a scheme. Let X be a scheme and X is an affine open covering. Let $U \subset X$ be a canonical and locally of finite type. Let X be a scheme. Let X be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let X be a scheme. Let X be a scheme covering. Let

 $b: X \to Y' \to Y \to Y \to Y' \times_X Y \to X.$

be a morphism of algebraic spaces over S and Y.

Proof. Let X be a nonzero scheme of X. Let X be an algebraic space. Let \mathcal{F} be a quasi-coherent sheaf of \mathcal{O}_X -modules. The following are equivalent

- *F* is an algebraic space over S.
- (2) If X is an affine open covering.

Consider a common structure on X and X the functor $\mathcal{O}_{\mathbf{Y}}(U)$ which is locally of

 Trained on all source files of Linux kernel concatenated into a single file (474MB of C code).

```
* Increment the size file of the new incorrect UI FILTER group information
* of the size generatively.
static int indicate policy(void)
 int error:
 if (fd == MARN EPT) {
     * The kernel blank will coeld it to userspace.
    if (ss->segment < mem total)</pre>
      unblock graph and set blocked();
    else
      ret = 1:
    goto bail;
 segaddr = in SB(in.addr);
 selector = seq / 16;
 setup works = true;
 for (i = 0; i < blocks; i++) {</pre>
   seq = buf[i++];
    bpf = bd->bd.next + i * search:
   if (fd) {
      current = blocked;
    }
  rw->name = "Getjbbregs";
 bprm self clearl(&iv->version);
  regs->new = blocks[(BPF STATS << info->historidac)] | PFMR CLOBATHINC SECONDS << 12;</pre>
 return segtable;
```

```
Copyright (c) 2006-2010, Intel Mobile Communications, All rights reserved.
  This program is free software; you can redistribute it and/or modify it
under the terms of the GNU General Public License version 2 as published by
the Free Software Foundation.
       This program is distributed in the hope that it will be useful,
but WITHOUT ANY WARRANTY; without even the implied warranty of
  MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
 GNU General Public License for more details.
  You should have received a copy of the GNU General Public License
  along with this program; if not, write to the Free Software Foundation,
 Inc., 675 Mass Ave, Cambridge, MA 02139, USA.
```

#include <linux/kexec.h>
#include <linux/errno.h>
#include <linux/io.h>
#include <linux/platform_device.h>
#include <linux/multi.h>
#include <linux/kevent.h>

#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/setew.h>
#include <asm/pgproto.h>

Evolution of Shakespeare

100 iter.:

tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tklrgd t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoaenns lng

300 iter.:

"Tmont thithey" fomesscerliund Keushey. Thom here sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

500 iter .:

we counter. He stutn co des. His stanted out one ofler that concossions and was to gearang reay Jotrets and with fre colt otf paitt thin wall. Which das stimn

700 iter.:

Aftair fall unsuch that the hall for Prince Velzonski's that me of her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort how, and Gogition is so overelical and ofter.

1200 iter.:

"Kite vouch!" he repeated by her door. "But I would be done and quarts, feeling, then, son is people...."

2000 iter.:

"Why do what that day," replied Natasha, and wishing to himself the fact the princess, Princess Mary was easier, fed in had oftened him. Pierre aking his soul came to the packs and drove up his father-in-law women. Consider the following task: Given a sequence of vectors

$$\mathbf{x} = \vec{x}_1, \ldots, \vec{x}_T$$

generate a new sequence

$$\mathbf{y}=\vec{y}_1,\ldots,\vec{y}_{T'}$$

of possibly different length (i.e., possibly $T \neq T'$).

E.g., a machine translation task, \mathbf{x} is an embedding of an English sentence, \mathbf{y} is a sequence of probability distributions on a German vocabulary.

Attention

Consider two recurrent networks:

- Enc the encoder
 - Hidden state $\vec{h_0}$ initialized by standard methods for recurrent networks
 - ► Reads $\vec{x}_1, ..., \vec{x}_T$, does not output anything but produces a sequence of hidden states $\vec{h}_1, ..., \vec{h}_T$

Attention

Consider two recurrent networks:

- Enc the encoder
 - Hidden state $\vec{h_0}$ initialized by standard methods for recurrent networks
 - ► Reads $\vec{x}_1, ..., \vec{x}_T$, does not output anything but produces a sequence of hidden states $\vec{h}_1, ..., \vec{h}_T$
- Dec the decoder
 - The initial hidden state is \vec{h}_T
 - ▶ Does not read anything but outputs the sequence $\vec{y}_1, \ldots, \vec{y}_{T'}$ This is a simplification. Typically, Dec reads $\vec{y}_0, \vec{y}_1, \ldots, \vec{y}_{T'-1}$ where \vec{y}_0 is a special vector embedding a separator.

Attention

Consider two recurrent networks:

- Enc the encoder
 - Hidden state $\vec{h_0}$ initialized by standard methods for recurrent networks
 - ► Reads $\vec{x}_1, ..., \vec{x}_T$, does not output anything but produces a sequence of hidden states $\vec{h}_1, ..., \vec{h}_T$
- Dec the decoder
 - The initial hidden state is \vec{h}_T
 - ▶ Does not read anything but outputs the sequence $\vec{y}_1, \ldots, \vec{y}_{T'}$ This is a simplification. Typically, Dec reads $\vec{y}_0, \vec{y}_1, \ldots, \vec{y}_{T'-1}$ where \vec{y}_0 is a special vector embedding a separator.

Trained on pairs of sentences, able to learn a fine translation between major languages (if the recurrent networks are LSTM).

Is not perfect because all info about $\mathbf{x} = \vec{x}_1, \dots, \vec{x}_T$ is squeezed into the single state vector \vec{h}_T .

In particular, the network tends to forget the context of each word.

What if we provide the decoder with an information about the *relevant context* of the generated word?

What if we provide the decoder with an information about the *relevant context* of the generated word?

We use the same encoder Enc producing the sequence of hidden states: $\vec{h}_1, \ldots, \vec{h}_T$

What if we provide the decoder with an information about the *relevant context* of the generated word?

We use the same encoder Enc producing the sequence of hidden states: $\vec{h}_1, \ldots, \vec{h}_T$

The decoder Dec is still a recurrent network but

• the hidden state \vec{h}'_0 initialized by \vec{h}_T and a sequence of hidden states $\vec{h}'_0, \ldots, \vec{h}'_T$, is computed,

What if we provide the decoder with an information about the *relevant context* of the generated word?

We use the same encoder Enc producing the sequence of hidden states: $\vec{h}_1, \ldots, \vec{h}_T$

The decoder Dec is still a recurrent network but

- the hidden state \vec{h}'_0 initialized by \vec{h}_T and a sequence of hidden states $\vec{h}'_0, \ldots, \vec{h}'_T$, is computed,
- ► reads a sequence of context vectors $\vec{c}_1, \ldots, \vec{c}_{T'}$ where

$$ec{c}_i = \sum_{j=1}^T lpha_{ij} ec{h}_j$$
 where $lpha_{ij} = rac{\exp(e_{ij})}{\sum_{k=1}^T \exp(e_{ik})}$

where $e_{ij} = \texttt{MLP}(ec{h}'_{i-1}, ec{h}_j)$

• outputs the sequence $\vec{y}_1, \ldots, \vec{y}_{T'}$

The attention mechanism extracts the information from the sequence quite well.

- The attention mechanism extracts the information from the sequence quite well.
- Is there a reason for reading the input sequence sequentially?

- The attention mechanism extracts the information from the sequence quite well.
- Is there a reason for reading the input sequence sequentially?
- Could we remove the recurrent network itself and preserve only the attention?

Self-Attention Layer (is all you need)

Fix an input sequence: $\vec{x}_1, \ldots, \vec{x}_T$

Consider three learnable matrices: W_q , W_k , W_v

Generate sequences of queries, keys, and values:

•
$$\vec{q}_1, \ldots, \vec{q}_T$$
 where $\vec{q}_k = W_q \vec{x}_k$ for all $k = 1, \ldots, T$

•
$$\vec{k}_1, \ldots, \vec{k}_T$$
 where $\vec{k}_k = W_k \vec{x}_k$ for all $k = 1, \ldots, T$

•
$$\vec{v}_1, \ldots, \vec{v}_T$$
 where $\vec{v}_k = W_v \vec{x}_k$ for all $k = 1, \ldots, T$

Self-Attention Layer (is all you need)

Fix an input sequence: $\vec{x}_1, \ldots, \vec{x}_T$

Consider three learnable matrices: W_q , W_k , W_v

Generate sequences of queries, keys, and values:

•
$$\vec{q}_1, \ldots, \vec{q}_T$$
 where $\vec{q}_k = W_q \vec{x}_k$ for all $k = 1, \ldots, T$
• $\vec{k}_1, \ldots, \vec{k}_T$ where $\vec{k}_k = W_k \vec{x}_k$ for all $k = 1, \ldots, T$
• $\vec{v}_1, \ldots, \vec{v}_T$ where $\vec{v}_k = W_v \vec{x}_k$ for all $k = 1, \ldots, T$

Define a vector score for all $i, j \in \{1, ..., T\}$ by

$$e_{ij} = \vec{q}_i \cdot \vec{k}_j$$

Intuitively, e_{ij} measures how much the input at the position *i* is related to the input at the position *j*, in other words, how much the query fits the key.

Define

$$\alpha_{ij} = \frac{\exp(\boldsymbol{e}_{ij} / \sqrt{d_{attn}})}{\sum_{k=1}^{T} \exp(\boldsymbol{e}_{ik} / \sqrt{d_{attn}})}$$

 d_{attn} is the dimension of \vec{v}_i

I.e., we apply the good old softmax to $(e_{i1}, \ldots, e_{iT}) / \sqrt{d_{attn}}$

Self-Attention Layer (is all you need)

Define a vector score for all $i, j \in \{1, ..., T\}$ by

$$e_{ij} = ec{q}_i \cdot ec{k}_j$$

Intuitively, e_{ij} measures how much the input at the position *i* is related to the input at the position *j*, in other words, how much the query fits the key.

Define

$$\alpha_{ij} = \frac{\exp(\boldsymbol{e}_{ij} / \sqrt{d_{attn}})}{\sum_{k=1}^{T} \exp(\boldsymbol{e}_{ik} / \sqrt{d_{attn}})}$$

 d_{attn} is the dimension of \vec{v}_i

I.e., we apply the good old softmax to $(e_{i1}, \ldots, e_{iT}) / \sqrt{d_{attn}}$

Define a sequence of outputs $\vec{y}_1, \ldots, \vec{y}_T$ by

$$\vec{y}_i = \sum_{j=1}^T \alpha_{ij} \cdot \vec{v}_j$$

Language Model

A sequence of *tokens* $a_1, \ldots, a_T \in \Sigma^*$ E.g. words from a vocabulary Σ .

The goal: Maximize

$$\prod_{k=1}^{T} P(a_k \mid a_1, ..., a_{k-1}; W) \qquad (= P(a_1, ..., a_T; W))$$

where

P is the conditional probability measure over Σ modeled using a neural network with weights W.

Language Model

A sequence of *tokens* $a_1, \ldots, a_T \in \Sigma^*$ E.g. words from a vocabulary Σ .

The goal: Maximize

$$\prod_{k=1}^{T} P(a_k \mid a_1, ..., a_{k-1}; W) \qquad (= P(a_1, ..., a_T; W))$$

where

P is the conditional probability measure over Σ modeled using a neural network with weights W.

Can be used to generate text:

Given $a_1, ..., a_k$, sample a_{k+1} from $P(a_{k+1} | a_1, ..., a_k; W)$

GPT



GPT



Masked Self-Attention Layer (is all you need)

Assume an attention mechanism which given an input sequence $\vec{x}_1, \ldots, \vec{x}_T$ generates $\vec{y}_1, \ldots, \vec{y}_T$.

The Problem: How to generate \vec{y}_k only based on $\vec{x}_1, \ldots, \vec{x}_{k-1}$?
Masked Self-Attention Layer (is all you need)

Assume an attention mechanism which given an input sequence $\vec{x}_1, \ldots, \vec{x}_T$ generates $\vec{y}_1, \ldots, \vec{y}_T$.

The Problem: How to generate \vec{y}_k only based on $\vec{x}_1, \ldots, \vec{x}_{k-1}$?

Define a vector score for all $i, j \in \{1, ..., T\}$ by

$$\mathbf{e}_{ij} = egin{cases} ec{\mathbf{q}}_i \cdot ec{\mathbf{k}}_j & ext{if } j < i \ -\infty & ext{otherwise.} \end{cases}$$

This means that

$$\alpha_{ij} = \begin{cases} \frac{\exp(e_{ij} / \sqrt{d_{attn}})}{\sum_{k=1}^{T} \exp(e_{ik} / \sqrt{d_{attn}})} & \text{ if } j < i \\ 0 & \text{ otherwise.} \end{cases}$$

Masked Self-Attention Layer (is all you need)

Assume an attention mechanism which given an input sequence $\vec{x}_1, \ldots, \vec{x}_T$ generates $\vec{y}_1, \ldots, \vec{y}_T$.

The Problem: How to generate \vec{y}_k only based on $\vec{x}_1, \ldots, \vec{x}_{k-1}$?

Define a vector score for all $i, j \in \{1, ..., T\}$ by

$$m{e}_{ij} = egin{cases} ec{m{q}}_i \cdot ec{m{k}}_j & ext{if } j < i \ -\infty & ext{otherwise.} \end{cases}$$

This means that

$$\alpha_{ij} = \begin{cases} \frac{\exp(e_{ij} / \sqrt{d_{attn}})}{\sum_{k=1}^{T} \exp(e_{ik} / \sqrt{d_{attn}})} & \text{ if } j < i \\ 0 & \text{ otherwise.} \end{cases}$$

Define a sequence of outputs $\vec{y}_1, \ldots, \vec{y}_T$ by

$$\vec{y}_i = \sum_{j=1}^T \alpha_{ij} \cdot \vec{v}_j$$

Multi-head Self-Attention Layer (is all you need)

Assume the number of *heads* is *H*.

For h = 1, ..., H the *h*-th head is an attention mechanism which given the input $\vec{x}_1, ..., \vec{x}_T$ produces

 $\vec{y}_1^h, \dots, \vec{y}_T^h$

Note that the output may be different which means that, in particular, the matrices W_q , W_k , W_v may be different for each head.

Assume that all vectors \vec{y}_k^h are of the same dimension d_{mid} and consider a learnable matrix W_{out} of dimensions $d_{out} \times (H \cdot d_{mid})$.

Multi-head Self-Attention Layer (is all you need)

Assume the number of *heads* is *H*.

For h = 1, ..., H the *h*-th head is an attention mechanism which given the input $\vec{x}_1, ..., \vec{x}_T$ produces

 $\vec{y}_1^h, \dots, \vec{y}_T^h$

Note that the output may be different which means that, in particular, the matrices W_q , W_k , W_v may be different for each head.

Assume that all vectors \vec{y}_k^h are of the same dimension d_{mid} and consider a learnable matrix W_{out} of dimensions $d_{out} \times (H \cdot d_{mid})$.

The multi-head attention produces the following output:

 $\vec{y}_1,\ldots,\vec{y}_T$

where

$$\vec{y}_k = W_{out} \cdot \left(\vec{y}_k^1 \odot \vec{y}_k^2 \odot \cdots \vec{y}_k^H \right)$$

Here \odot is a concatenation of vectors.

Multi-head Self-Attention Summary

Input: A sequence $\vec{x}_1, \ldots, \vec{x}_T$ Output: A sequence $\vec{y}_1, \ldots, \vec{y}_T$

I.e., a sequence of the same length. The dimensions of \vec{y}_k and \vec{x}_k do not have to be equal.

Multi-head Self-Attention Summary

Input: A sequence $\vec{x}_1, \dots, \vec{x}_T$ Output: A sequence $\vec{y}_1, \dots, \vec{y}_T$

I.e., a sequence of the same length. The dimensions of \vec{y}_k and \vec{x}_k do not have to be equal.

Attention:

Learnable parameters: Matrices W_q , W_k , W_v .

These matrices are used to compute queries, keys, and values from $\vec{x}_1, \ldots, \vec{x}_T$. Output $\vec{y}_1, \ldots, \vec{y}_T$ is computed using values "scaled" by the query-key attention.

Multi-head Self-Attention Summary

Input: A sequence $\vec{x}_1, \dots, \vec{x}_T$ Output: A sequence $\vec{y}_1, \dots, \vec{y}_T$

I.e., a sequence of the same length. The dimensions of \vec{y}_k and \vec{x}_k do not have to be equal.

Attention:

Learnable parameters: Matrices W_q , W_k , W_v .

These matrices are used to compute queries, keys, and values from

 $\vec{x}_1, \dots, \vec{x}_T$. Output $\vec{y}_1, \dots, \vec{y}_T$ is computed using values "scaled" by

the query-key attention.

Multi-head attention:

Learnable parameters:

• Matrices W_q^h , W_k^h , W_v^h where h = 1, ..., H and H is the number of heads.

Each attention head operates independently on the input $\vec{x}_1, \ldots, \vec{x}_T$.

Matrix W_{out}.

Linearly transforms the concatenated results of the attention heads.

GPT - transformer



Positional encoding

The Goal: To encode a position (index) $k \in \{1, ..., T\}$ into a vector \vec{P}_k of real numbers.

Positional encoding

The Goal: To encode a position (index) $k \in \{1, ..., T\}$ into a vector \vec{P}_k of real numbers.

Assume that \vec{P}_k should have a dimension d. Given a position $k \in \{1, ..., T\}$ and $i \in \{0, ..., d/2\}$ define

$$P_{k,2i} = \sin\left(\frac{k}{n^{2i/d}}\right)$$
$$P_{k,(2i+1)} = \cos\left(\frac{k}{n^{2i/d}}\right)$$

Here n = 10000.

A user defined constant, the original paper suggests n = 10000.

Positional encoding

The Goal: To encode a position (index) $k \in \{1, ..., T\}$ into a vector \vec{P}_k of real numbers.

Assume that \vec{P}_k should have a dimension d. Given a position $k \in \{1, ..., T\}$ and $i \in \{0, ..., d/2\}$ define

$$P_{k,2i} = \sin\left(\frac{k}{n^{2i/d}}\right)$$
$$P_{k,(2i+1)} = \cos\left(\frac{k}{n^{2i/d}}\right)$$

Here n = 10000.

A user defined constant, the original paper suggests n = 10000.

Given an input sequence $\vec{x}_1, \ldots, \vec{x}_T$ we add the position embedding to each \vec{x}_k obtaining a new input sequence $\vec{x}'_1, \ldots, \vec{x}'_T$ where

$$\vec{x}_k' = \vec{x}_k + \vec{P}_k$$

Positional encoding/embedding



Positional encoding/embedding



- Vertically: Sinusoidal functions
- Horizontally: Decreasing frequency

For any offset $o \in \{1, ..., T\}$ there is a linear transformation M such that for any $k \in \{1, ..., T - o\}$ we have $M\vec{P}_k = \vec{P}_{k+o}$. Intuitively, just rotate each component of the \vec{P}_k appropriately.

GPT-2 - transformer



Layer normalization

Given a vector $\vec{x} \in \mathbb{R}^d$, the *layer normalization* computes:

$$\vec{x}' = \gamma \cdot \frac{(\vec{x} - \mu)}{\sigma} + \beta$$

Here

•
$$\mu = \frac{1}{d} \sum_{i=1}^{d} x_i$$
 and $\sigma^2 = \frac{1}{d} \sum_{i=1}^{d} (x_i - \mu)^2$
• $\gamma, \beta \in \mathbb{R}^d$ are vectors of trainable parameters

Layer normalization

Given a vector $\vec{x} \in \mathbb{R}^d$, the *layer normalization* computes:

$$\vec{x}' = \gamma \cdot \frac{(\vec{x} - \mu)}{\sigma} + \beta$$

Here

•
$$\mu = \frac{1}{d} \sum_{i=1}^{d} x_i$$
 and $\sigma^2 = \frac{1}{d} \sum_{i=1}^{d} (x_i - \mu)^2$

• $\gamma, \beta \in \mathbb{R}^d$ are vectors of trainable parameters

In Transformer:

The input to the layer normalization is a sequence of vectors: $\vec{x}_1, \ldots, \vec{x}_T$. The layer normalization is applied to each \vec{x}_k , producing a sequence of "normalized" vectors.

GPT - learning

A sequence of tokens $a_1, \ldots, a_T \in \Sigma$ and their one-hot encodings $\vec{u}_1, \ldots, \vec{u}_T \in \{0, 1\}^{|\Sigma|}$ We assume that a_1 is a special token marking the start of the sequence.

Embed to vectors and add the position encoding (W_e is an embedding matrix):

$$\vec{x}_k = W_e \cdot \vec{u}_k + P_k \in Rset^d$$

Г	Text Task Prediction Classifier
	Layer Norm
2x -	Layer Norm
	Masked Multi Self Attention
	Text & Position Embed

GPT - learning

A sequence of tokens $a_1, \ldots, a_T \in \Sigma$ and their one-hot encodings $\vec{u}_1, \ldots, \vec{u}_T \in \{0, 1\}^{|\Sigma|}$ We assume that a_1 is a special token marking the start of the sequence.

Embed to vectors and add the position encoding (W_e is an embedding matrix):

$$\vec{x}_k = W_e \cdot \vec{u}_k + P_k \in Rset^d$$



Apply the network (with the transformer block repeated 12x) to $\vec{x}_1, \ldots, \vec{x}_T$ and obtain $\vec{y}_1, \ldots, \vec{y}_T$ (Here assume that each $\vec{y}_k \in [0, 1]^{\Sigma}$ is a probability distribution on Σ)

GPT - learning

A sequence of tokens $a_1, \ldots, a_T \in \Sigma$ and their one-hot encodings $\vec{u}_1, \ldots, \vec{u}_T \in \{0, 1\}^{|\Sigma|}$ We assume that a_1 is a special token marking the start of the sequence.

Embed to vectors and add the position encoding (W_e is an embedding matrix):

$$\vec{x}_k = W_e \cdot \vec{u}_k + P_k \in Rset^d$$



Apply the network (with the transformer block repeated 12x) to $\vec{x}_1, \ldots, \vec{x}_T$ and obtain $\vec{y}_1, \ldots, \vec{y}_T$ (Here assume that each $\vec{y}_k \in [0, 1]^{\Sigma}$ is a probability distribution on Σ)

Compute the error:

$$-\sum_{\ell=1}^{T-1}\log\left(\vec{y}_\ell[a_{\ell+1}]\right)$$

Here $\vec{y}_{\ell}[a_{k+1}]$ is the probability of a_{k+1} in the distribution \vec{y}_k .

GPT - inference



A sequence of tokens $a_1, \ldots, a_\ell \in \Sigma$ and their one-hot encodings $\vec{u}_1, \ldots, \vec{u}_\ell \in \{0, 1\}^{|\Sigma|}$

Embed to vectors and add the position encoding:

$$ec{x}_k = oldsymbol{W}_{oldsymbol{e}} \cdot ec{u}_k + oldsymbol{P}_k \in oldsymbol{Rset}^d$$

Apply the network to $\vec{x}_1, \ldots, \vec{x}_\ell$ and obtain $\vec{y}_1, \ldots, \vec{y}_\ell$ (Assume that each $\vec{y}_k \in [0, 1]^{\Sigma}$ is a probability distribution on Σ)

Sample the next token from

$$a_{\ell+1} \sim \vec{y}_{\ell}$$

https://transformer.huggingface.co/doc/distil-gpt2

Architectures:

- Multi-layer perceptron (MLP):
 - dense connections between layers
- Convolutional networks (CNN):
 - local receptors, feature maps
 - pooling
- Recurrent networks (RNN):
 - self-loops but still feed-forward through time
- Transformer
 - Attention, query-key-value

Training:

gradient descent algorithm + heuristics