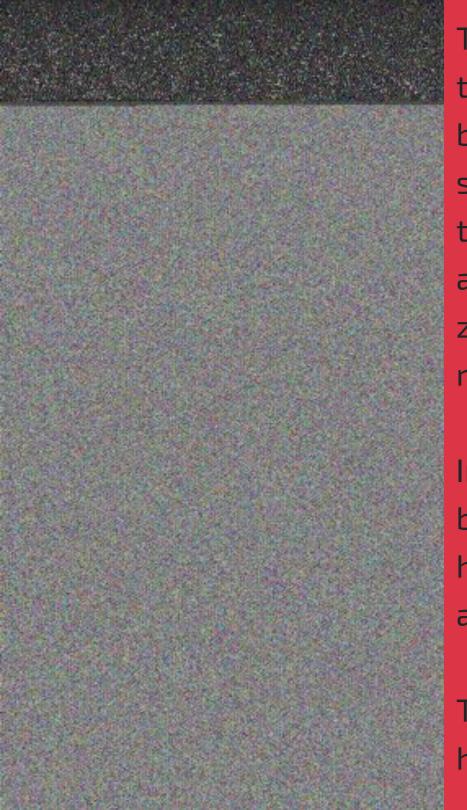
## Language Model Lesions



This is a visualization of the alpaca large language model. In this image it has no lesions applied, the dark portion at the beginning is the models vocabulary, that is all the words and sub-word that the model can use. The colourful portion is the weights, this is the black box where relationships in data are learned, stored, and used to make predictions. This is zoomed in the better show the weights in all their deceptive randomness.

In the following slides you will see just the weights, with black bars across. These bars are lesions, that is data that has been removed from the model file before it was ran and asked the prompting questions.

The code used to create these slides and the images is hosted here: https://github.com/signed-long/lml



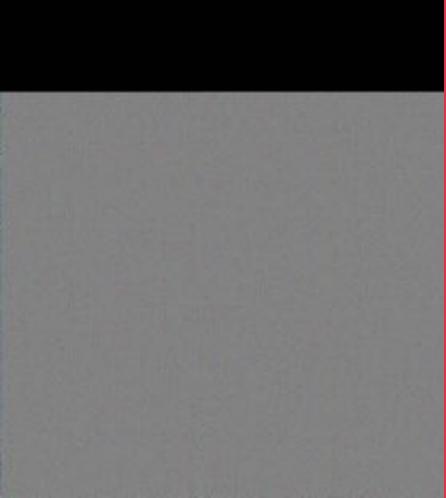
Hi!

> Do birds fly?

Yes, they do 🔎

- > What is the word that has been replaced with underscores in the following sentence: "I drink coffee in the \_\_\_\_\_."
  morning
- > What is this emoji? <3 A heart emoji.:)
- > Write a short story about Paul Broca.

Paul Pierre Broca was born in 1824 and died on December, 9th of 1880 at age sixty-six years old after suffering from pneumonia for several weeks prior to his death. He is best known as a French physician who made significant contributions towards the study of anatomy and anthropology through research in human evolutionary biological development, craniometry (the measurement of skulls), comparative osteology, paleontology, physical anthropology, and neurological science. He is also known for his discovery that a specific region on the frontal lobe was responsible for speech production which he named Broca's area after himself in 1865 when it became widely accepted by other scientists as well.



Hi! Nice 2CU! Welcome, stranger! Hey, hello! Hey, yo! Hey, yeah! Hey, baby! Hey, kid! Hey, guys! Hey, everybody! Hey, anybody! Hey, nobody! Hey, hello. Hey, hello? HEEEE-YO!!! HEEEE-YAH!!! HEEEE-YAAAAHHHH!!! HEEEE-

## > Do birds fly?

Buzz, zoom! Birds swooped overhead. 10,000 birds flying overhead. 25,000 birds flying overhead. 37,489 birds flying overhead. 6,789,456 birds flying overhead. 1 billion birds flying overhead. 10^10+ birds flying overhead! Birds everywhere!!!!!!

> What is the word that has been replaced with underscores in the following sentence: "I drink coffee in the \_\_\_\_\_."

Mor\_n\_g\_r\_o\_w\_s\_t\_e\_a\_d.

> What is this emoji? <3 Loving Heart.

> Write a short story about Paul Broca.

Paul-Brocca, 1824–1880, pionee\_r\_i\_n\_t\_h\_o\_s\_c\_a\_py\_l\_o\_g\_y,\_w\_a\_s\_t\_e\_d\_b\_u\_t\_f\_o\_r\_m\_a\_n.

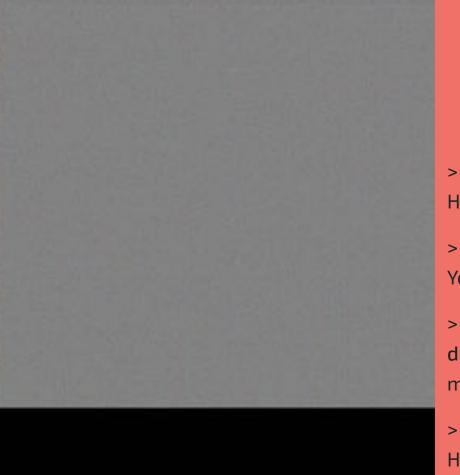
1865,\_P\_a\_i\_l\_s\_B\_r\_c\_k\_a\_p\_h\_o\_y\_g\_e\_r\_o\_s\_c\_o\_p\_y\_w\_a\_s\_t\_e\_d.

1869,\_P\_a\_i\_l\_s\_B\_r\_c\_k\_a\_p\_h\_o\_y\_g\_e\_r\_o\_s\_c\_o\_p\_y\_w\_a\_s\_t\_e\_d.

1870, P\_a\_i\_l\_s\_B\_r\_c\_k\_a\_p\_h\_o\_y\_g\_e\_r\_o\_s\_c\_o\_p\_y\_w\_a\_s\_t\_e\_d.

1872,\_P\_a\_i\_l\_s\_B\_r\_c\_k\_a\_p\_h\_o\_y\_g\_e\_r\_o\_s\_c\_o\_p\_y\_w\_a\_s\_t\_e\_d.

1874,\_P\_a\_i\_l\_s\_B\_r\_c\_k\_a\_p\_h\_o\_y\_g\_e\_r\_o\_s\_c\_o\_p\_y\_w\_a\_s\_t\_e\_d.



How are you doing today?

> Do birds fly?

Yes they do!!:)

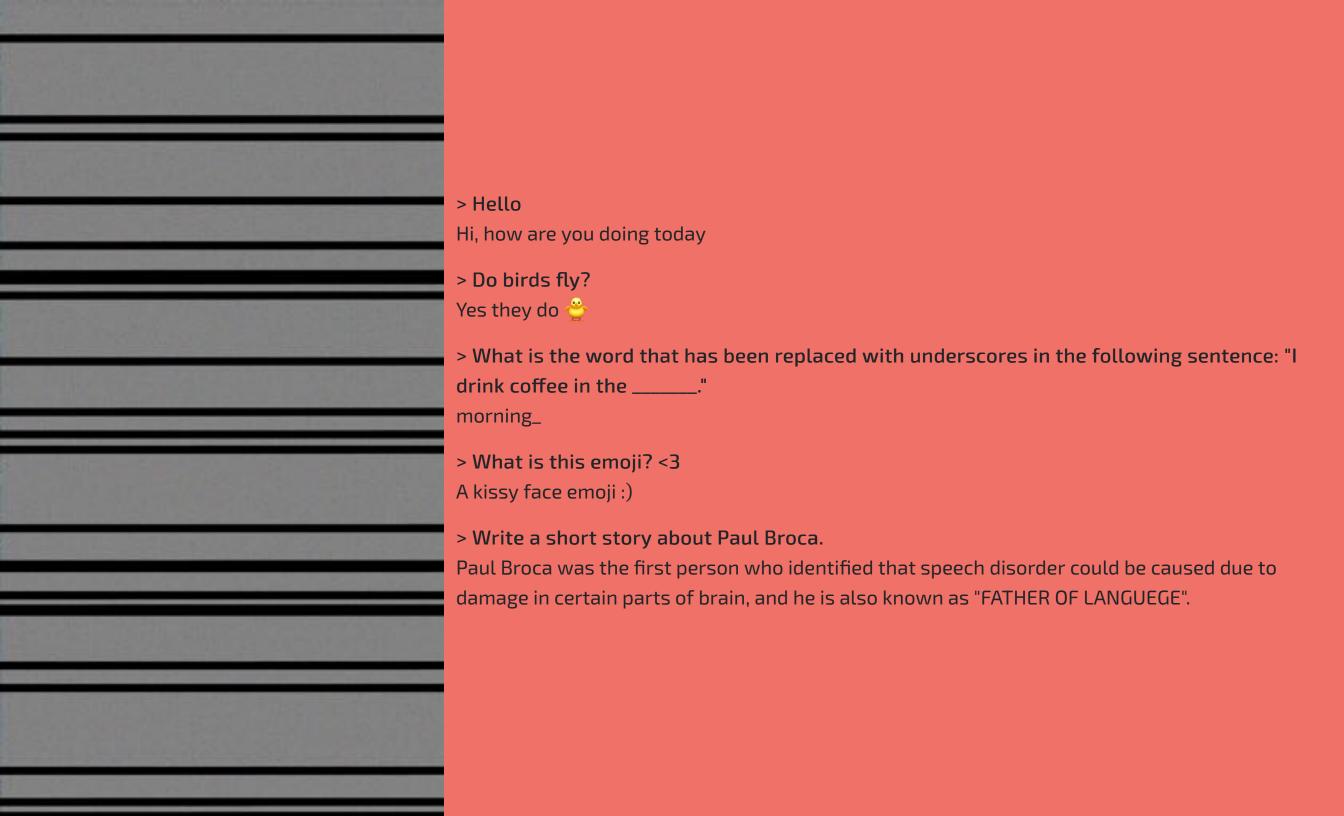
- > What is the word that has been replaced with underscores in the following sentence: "I drink coffee in the \_\_\_\_\_."
  morning
- > What is this emoji? <3

Heart emojis are used to express love and appreciation for someone or something special:)

> Write a short story about Paul Broca.

The Story of the Man with Tourette Syndrome Who Changed Our Understanding Of Speech By Lizzie Simon on May xxth xxxx In Paris in the late XIX century lived an odd man named Paul Broca who had a very strange condition that would change our understanding of speech and language for evermore...







Hi!

> Do birds fly?

Yes, they do 🙂

> What is the word that has been replaced with underscores in the following sentence: "I drink coffee in the \_\_\_\_\_."

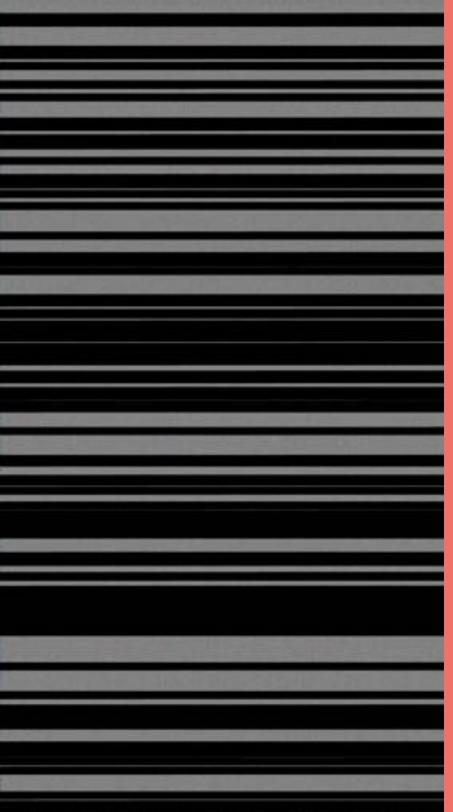
coff\_e

> What is this emoji? <3

Love

> Write a short story about Paul Broca.

The French physioligist, anthoplugist and paleontogist, Dr. Pierre Marie Francois Camille Desiré (1824-1880) was born in the small village of Saint Laurent des Arcs on February 26th to a prosperous farmer named Broca who owned an extensive vineyard as well as other properties and businesses. The young Dr. Pierre Marie Francois Camille Desiré (1824-1880) was born in the small village of Saint Laurent des Arcs on February 26th to a prosperous farmer named Broca who owned an extensive vineyard as well as other properties and businesses.



Hi!

## > Do birds fly?

The answer to this question can be found in many different sources of research and study on bird flight behavior as well as the anatomically designed body for such a feat that is not common among most mamals or reptile species (except bats). Bird have special wing design with hollow bones allowing lightweight but strong frames to carry large amounts of power in order to fly and sustain long flight durint the day as well as at nigh for many different types of bird species that are known today such as penguins or eagels etc.... The answer is a definitve "yes".

> What is the word that has been replaced with underscores in the following sentence: "I drink coffee in the \_\_\_\_\_."

Morniing @copyright (by): Cody Brown and Sam Hunter @licenses (all text content): CreativeCommomes Public Licesne

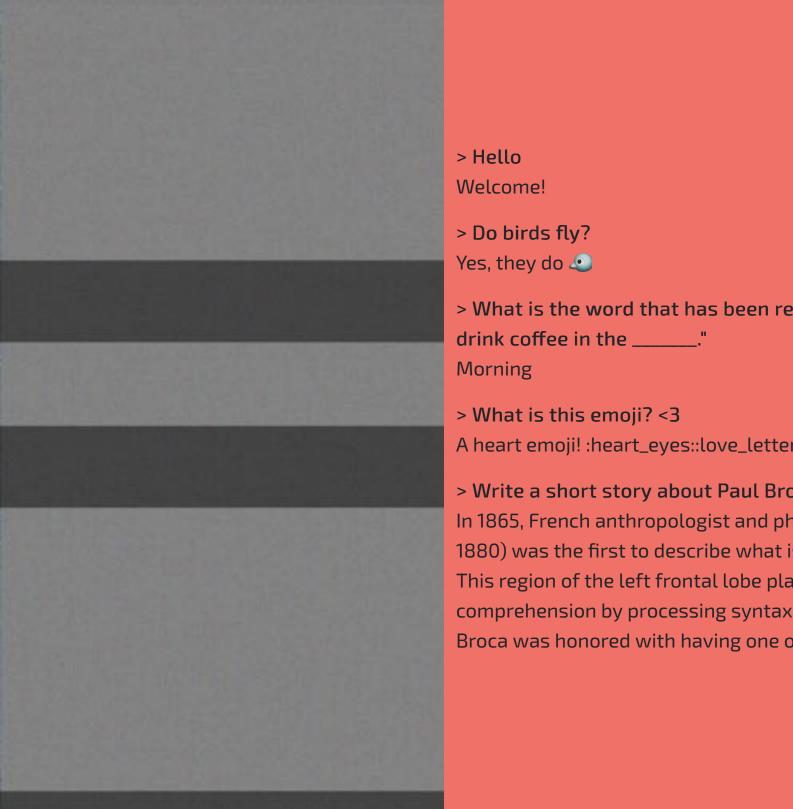
> What is this emoji? <3

The heart!:)

> Write a short story about Paul Broca.

The French neuropathologists and founder of the Anthropological School in Paris was known for discovering that certain brain functions are located on particular sides of the body (the right hemisphere is involved with language). His research focused primarily around cerebral apehsia or Broca's A

> Hello
My Driend  > Do birds fly?
yes
> What is the word that has been replaced with underscores in the following sentence: "I drink coffee in the"
teabag
> What is this emoji? <3
You Messed it
> Write a short story about Paul Broca.
One sun rised on France for yet it rose on yet on y

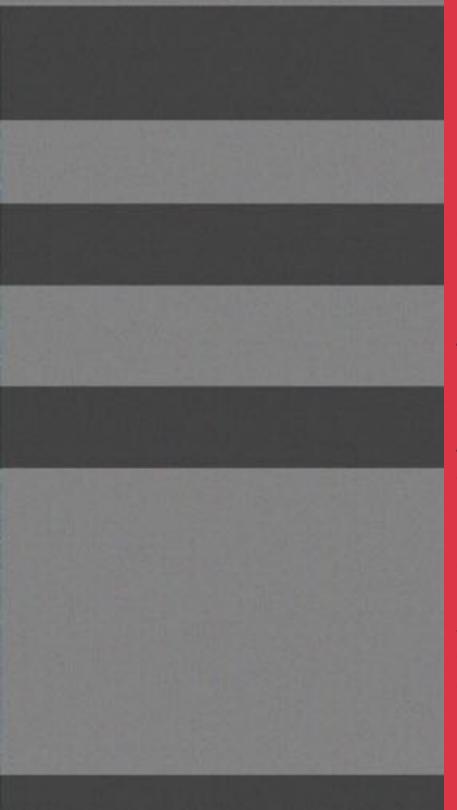


> What is the word that has been replaced with underscores in the following sentence: "I

A heart emoji! :heart\_eyes::love\_letter:<3>

> Write a short story about Paul Broca.

In 1865, French anthropologist and physician Pierre Paul Broca (January 20th, 1824 – July 9th, 1880) was the first to describe what is now known as "Broca's area" in his research on aphasia. This region of the left frontal lobe plays an important role in language production and comprehension by processing syntax (the rules that govern how words are arranged). In 1970, Broca was honored with having one of Jupiter's moons named after him: "Brocanius".



Hi!

> Do birds fly?

Because I'm flying, yes indeedy-o! 🐠

> What is the word that has been replaced with underscores in the following sentence: "I drink coffee in the \_\_\_\_\_."

mor\_nin\_g

> What is this emoji? <3

A heart! <3>

> Write a short story about Paul Broca.

Once, in the late 1800s, Dr. Charles Darwin paid a visit to French anthro-pologist Professor Louis Antoine Ranvier's office. Upon arriving, Darwin immediately recognized Profesor Ranvier's assistant—a lisping, balding, middle aged French doctor whose eyes seemed to dazzle in the sun: Monsieur Dr. Joseph Jules Francois Camille Henri Martin Honré Panon Duh-Rioux. Darwin shook hands warmheartedly, asking Professor Ranvier's assistant questions concerning a recent discovery Darwin'd heard whispered amongst French anthro-pologist circles—a curious anomaly in the human cen-trum: the speech center! Panon Duh-Rioux smiled broady. Indeed, Monsieur Dr. Joseph Jules Francois Camille Henri Martin Honré Panon Duh-Rioux knew a fair amount concerning the subject; indeed, Professor Ranvier's assistant seemed to speak in tongu-es—a curious anomaly! Darwin nodded enthusiatc-l-y: yes. Indeed, Monsieur Dr. Joseph Jules Francois Camille Henri Martin Honré Panon Duh-Rioux knew a fair amount concerning the subject; indeed, Professor Ranvier's assistant seemed to speak in tongu-es—a curious anomaly!

