# Deep Dream for LLM's

By Peter Chatain and Scott Viteri

# Deep Dream for LLM's



Deep Dream ->

Motivation:

1. We want sentences that describe a neuron's behavior
2. Naive deep dream produces gibberish
   a. True for **images**, and for **text**
3. Find a regularization technique that constrains to recognizable sentences.
4. GPT-N uses recognizable sentences to produce explanations

## Language models can explain neurons in language models

AUTHORS

Steven Bills*, Nick Cammarata*, Dan Mossing*, Henk Tillman*, Leo Gao*, Gabriel Goh, Ilya Sutskever, Jan Leike, Jeff Wu*, William Saunders*

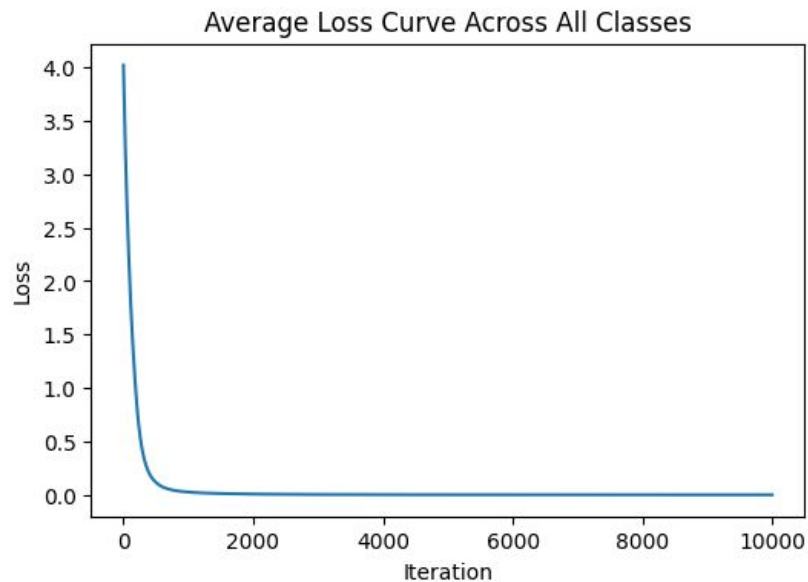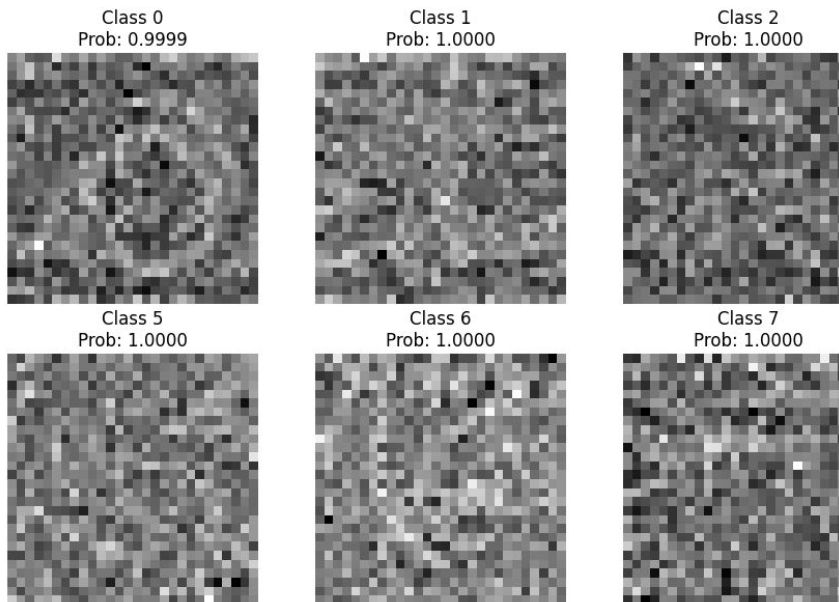* Core Research Contributor;   Author contributions statement below.   Correspondence to interpretability@openai.com.
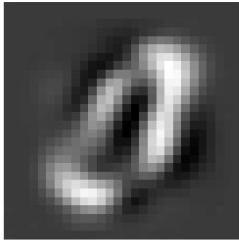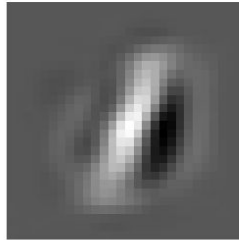
AFFILIATION

OpenAI

PUBLISHED

May 9, 2023

# Problem

# Solution

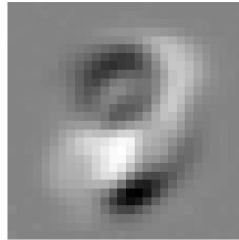# Auto-encoder loss

# Future Work

1. Training better autoencoders
2. Placing the auto-encoder in intermediate layers, not just after the first embedding layer.
3. Instead of searching with gradient descent, search over tokens restricted to only the most plausible tokens as judged by the language model itself