The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

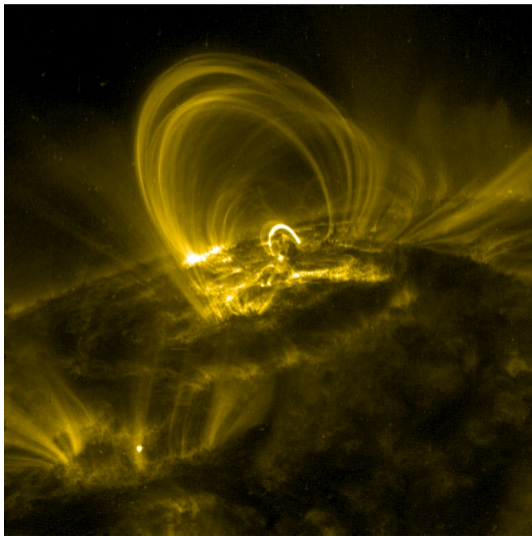# How to Find ChatGPT's Hidden Size, and Other Low-rank Logit Tricks

**Matthew Finlayson**     Xiang Ren     Swabha Swayamdipta

University of Southern California

April 8, 2024

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

The technical details
○○○○○

Consequences of knowing the LLM image
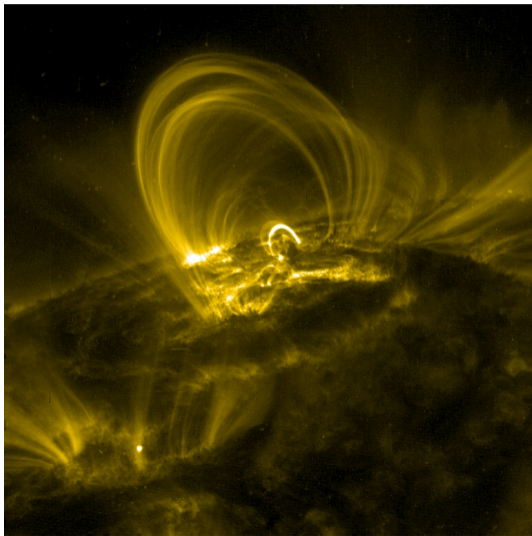○○○○○○

What now?
○○○○

# The solar corona



- Looks really cool.
- Made of plasma (ions).
- 150–450× hotter than the sun surface.
- The sun's magnetic field causes *coronal loops*.

The technical details
ooooo

Consequences of knowing the LLM image
oooooo

What now?
oooo

# The solar corona

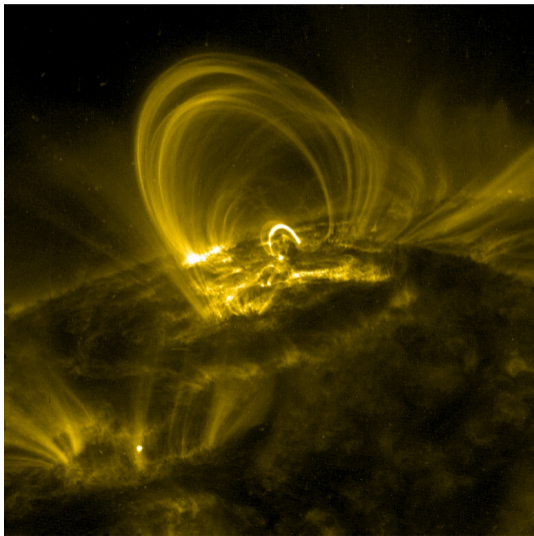

- Looks really cool.
- Made of plasma (ions).
- 150–450× hotter than the sun surface.
- The sun's magnetic field causes *coronal loops*.

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# The solar corona

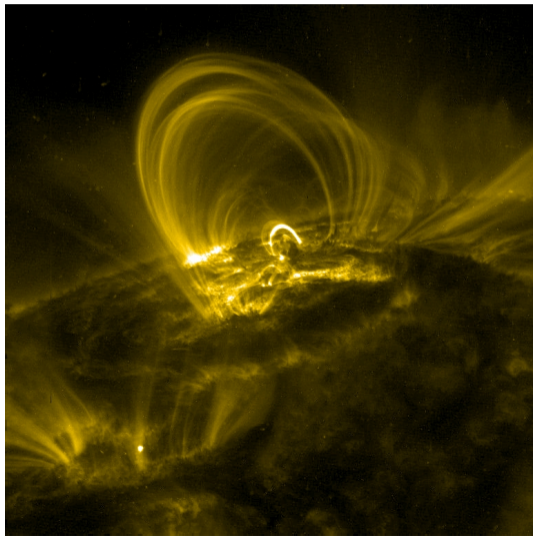

- Looks really cool.
- Made of plasma (ions).
- 150–450× hotter than the sun surface.
- The sun's magnetic field causes *coronal loops*.

The technical details
ooooo

Consequences of knowing the LLM image
oooooo

What now?
oooo

# The solar corona



- Looks really cool.
- Made of plasma (ions).
- 150–450× hotter than the sun surface.
- The sun's magnetic field causes *coronal loops*.

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

## The solar corona
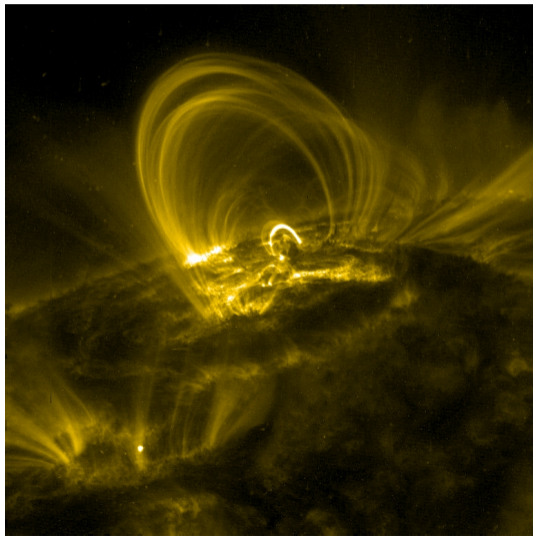


- Looks really cool.
- Made of plasma (ions).
- 150–450× hotter than the sun surface.
- The sun's magnetic field causes *coronal loops*.

The technical details
OOOOO

Consequences of knowing the LLM image
OOOOOO

What now?
OOOO

# The analogy

- Scientists study the *structure* of coronal loops to learn about the sun's *internal* magnetic fields.

- We can study the *structure* of LLM outputs to learn about their *internal* details.



- The sun
- The sun's magnetic field
- The solar corona

- Proprietary LLMs
- Non-public model details
- LLM API outputs

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# The analogy

- Scientists study the *structure* of coronal loops to learn about the sun's *internal* magnetic fields.
- We can study the *structure* of LLM outputs to learn about their *internal* details.





- The sun
- The sun's magnetic field
- The solar corona

- Proprietary LLMs
- Non-public model details
- LLM API outputs

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# The analogy

- Scientists study the *structure* of coronal loops to learn about the sun's *internal* magnetic fields.
- We can study the *structure* of LLM outputs to learn about their *internal* details.





- The sun
- The sun's magnetic field
- The solar corona

- Proprietary LLMs
- Non-public model details
- LLM API outputs

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# The analogy

- Scientists study the *structure* of coronal loops to learn about the sun's *internal* magnetic fields.
- We can study the *structure* of LLM outputs to learn about their *internal* details.



- The sun
- The sun's magnetic field
- The solar corona



- Proprietary LLMs
- Non-public model details
- LLM API outputs

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# The analogy

- Scientists study the *structure* of coronal loops to learn about the sun's *internal* magnetic fields.
- We can study the *structure* of LLM outputs to learn about their *internal* details.



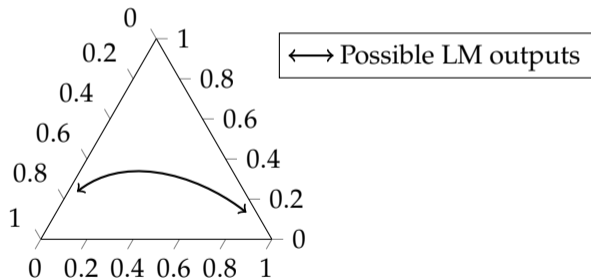🌀 OpenAI

- The sun
- The sun's magnetic field
- The solar corona

- Proprietary LLMs
- Non-public model details
- LLM API outputs

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# The structure of LLM outputs

LLM outputs lie within a low-dimensional space.

Space of probability distributions over 3 items

Section 1

The technical details

The technical details
○●○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# LLM architecture

Input

| Karma is my | → | Transformer |

Embedding $h \in \mathbb{R}^d$

The technical details
○●○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# LLM architecture

Softmax matrix

$W \in \mathbb{R}^{v \times d}$

Input

```
Karma is my
```

→ Transformer → MatMul

Logits $\boldsymbol{\ell} \in \mathbb{R}^v$

Embedding $\boldsymbol{h} \in \mathbb{R}^d$

The technical details
○●○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# LLM architecture

Softmax matrix

$W \in \mathbb{R}^{v \times d}$

Input

| Karma is my | → Transformer → MatMul → SoftMax

| 0.6 | boyfriend |
| 0.2 | queen |
| 0.1 | thought |
| 0.1 | breeze |
| 0.0 | acrobat |

Logits $\boldsymbol{\ell} \in \mathbb{R}^v$

Embedding $\boldsymbol{h} \in \mathbb{R}^d$

Probabilities $\boldsymbol{p} \in \Delta_v$

The technical details
○○●○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# The set of next-token distributions

Next-token distributions over a vocabulary of size $v$ are

- $v$-tuples of reals.
- Non-negative, sum to 1.
- Known as the $v$-simplex, or $\Delta_v$.

## The set of next-token distributions
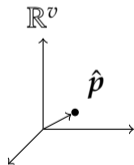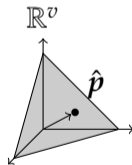
Next-token distributions over a vocabulary of size $v$ are

- $v$-tuples of reals.
- Non-negative, sum to 1.
- Known as the $v$-simplex, or $\Delta_v$.

$$\mathbb{R}^v$$

$$\hat{p}$$

The technical details
○○●○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# The set of next-token distributions

Next-token distributions over a vocabulary of size $v$ are

- $v$-tuples of reals.
- Non-negative, sum to 1.
- Known as the $v$-simplex, or $\Delta_v$.

The technical details
○○●○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# The set of next-token distributions
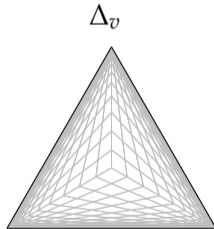
Next-token distributions over a vocabulary of size $v$ are

- $v$-tuples of reals.
- Non-negative, sum to 1.
- Known as the $v$-simplex, or $\Delta_v$.

The technical details
○○○●○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# Probability distributions are vectors

- $\Delta_v$ is a *vector space*.
- The softmax function is a *linear map* $\mathbb{R}^v \to \Delta_v$.

$$\mathbb{R}^v \qquad\qquad \Delta_v$$

The technical details
○○○●○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# Probability distributions are vectors

- $\Delta_v$ is a *vector space.*
- The softmax function is a *linear map* $\mathbb{R}^v \rightarrow \Delta_v$.

**Matthew Finlayson** @mattf1n · Oct 5
Did you know that the softmax function is linear?

| | |
|---|---|
| I knew that | 22.2% |
| I did not know that | 20% |
| **I don't believe you** | **57.8%** |

45 votes · Final results

# Probability distributions are vectors

- $\Delta_v$ is a *vector space.*
- The softmax function is a *linear map* $\mathbb{R}^v \rightarrow \Delta_v$.

**Matthew Finlayson** @mattf1n · Oct 5 ···
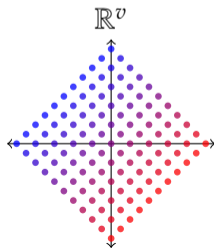Did you know that the softmax function is linear?

| | |
|---|---|
| I knew that | 22.2% |
| I did not know that | 20% |
| **I don't believe you** | **57.8%** |

45 votes · Final results

$\mathbb{R}^v$  softmax$(\mathbb{R}^v)$  clr $(\text{softmax}(\mathbb{R}^v))$

The technical details
○○○○●

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# LLM outputs lie on a low-dimensional vector subspace

- The image of a function is its codomain.
- The dim of a linear map's image is $\leq$ the dim of its domain.
- softmax $\circ W$ is a linear map $\mathbb{R}^d \to \Delta_v$.
- The dim of an LLM's image is at most $d$
- $d$ LLM outputs form a *basis* for its image.

$$\mathbb{R}^d \xrightarrow{\quad W \quad} \mathbb{R}^v$$

The technical details
○○○○●

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

## LLM outputs lie on a low-dimensional vector subspace

- The image of a function is its codomain.
- The dim of a linear map's image is ≤ the dim of its domain.
- softmax $\circ W$ is a linear map $\mathbb{R}^d \to \Delta_v$.
- The dim of an LLM's image is at most $d$
- $d$ LLM outputs form a *basis* for its image.



$$\mathbb{R}^d \xrightarrow{\quad W \quad} \mathbb{R}^v$$

$$\longleftrightarrow \mathrm{im}(W)$$

The technical details
○○○○●

Consequences of knowing the LLM image
○○○○○○

What now?
○○○○

# LLM outputs lie on a low-dimensional vector subspace

- The image of a function is its codomain.
- The dim of a linear map's image is $\leq$ the dim of its domain.
- softmax $\circ W$ is a linear map $\mathbb{R}^d \to \Delta_v$.
- The dim of an LLM's image is at most $d$
- $d$ LLM outputs form a *basis* for its image.



$$\mathbb{R}^d \xrightarrow{\quad W \quad} \mathbb{R}^v \xrightarrow{\quad \text{softmax} \quad} \Delta_v$$

$\longleftrightarrow \mathrm{im}(W)$   $\longleftrightarrow \mathrm{im}(\text{softmax} \circ W)$

The technical details
○○○○●

Consequences of knowing the LLM image
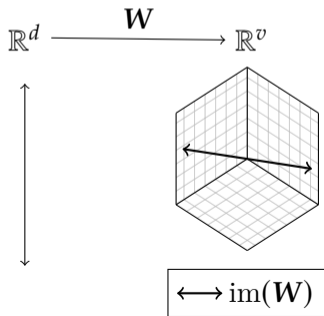○○○○○○

What now?
○○○○

## LLM outputs lie on a low-dimensional vector subspace

- The image of a function is its codomain.
- The dim of a linear map's image is $\leq$ the dim of its domain.
- softmax $\circ\, W$ is a linear map $\mathbb{R}^d \to \Delta_v$.
- The dim of an LLM's image is at most $d$
- $d$ LLM outputs form a *basis* for its image.



$$\mathbb{R}^d \xrightarrow{\quad W \quad} \mathbb{R}^v \xrightarrow{\quad \text{softmax} \quad} \Delta_v$$

$\longleftrightarrow \text{im}(W)$  $\longleftrightarrow \text{im}(\text{softmax} \circ W)$

The technical details
OOOO●

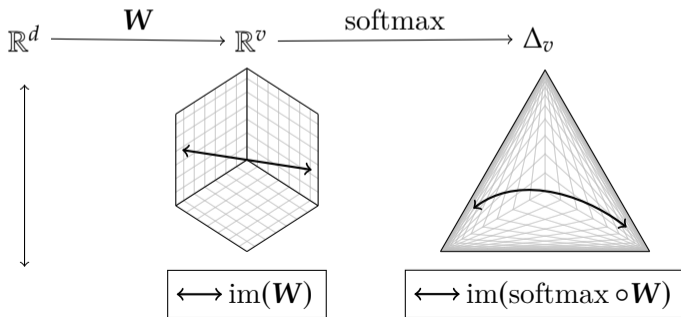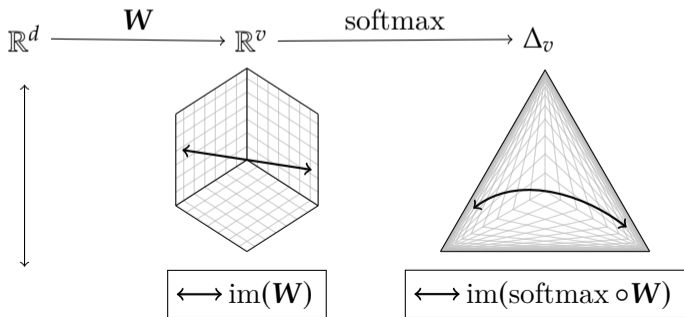Consequences of knowing the LLM image
OOOOOO

What now?
OOOO

## LLM outputs lie on a low-dimensional vector subspace

- The image of a function is its codomain.
- The dim of a linear map's image is ≤ the dim of its domain.
- softmax $\circ\, W$ is a linear map $\mathbb{R}^d \to \Delta_v$.
- The dim of an LLM's image is at most $d$
- $d$ LLM outputs form a *basis* for its image.



$$\mathbb{R}^d \xrightarrow{\quad W \quad} \mathbb{R}^v \xrightarrow{\quad \text{softmax} \quad} \Delta_v$$

$\longleftrightarrow \operatorname{im}(W)$    $\longleftrightarrow \operatorname{im}(\text{softmax} \circ W)$

The technical details
○○○○○

Consequences of knowing the LLM image
●○○○○○

What now?
○○○○

Section 2

Consequences of knowing the LLM image

The technical details
○○○○○

Consequences of knowing the LLM image
○●○○○○

What now?
○○○○

# Cheap, full LLM outputs

**Stealing Machine Learning Models via Prediction APIs**

Florian Tramèr
*EPFL*

Fan Zhang
*Cornell University*

Ari Juels
*Cornell Tech, Jacobs Institute*

Michael K. Reiter
*UNC Chapel Hill*

Thomas Ristenpart
*Cornell Tech*

- Common APIs give top-$k$ tokens and probabilities.
- Logit bias allows boosting tokens to top-$k$.
- Extracting full outputs takes $O(v/k)$ API calls.
- Once the LLM image is known, only $O(d/k)$ calls.
- Intuition: position in a $d$-dimensional subspace is fully specified by $d$ coordinates.

The technical details
○○○○○

Consequences of knowing the LLM image
○●○○○○

What now?
○○○○

# Cheap, full LLM outputs

- Common APIs give top-$k$ tokens and probabilities.
- Logit bias allows boosting tokens to top-$k$.
- Extracting full outputs takes $O(v/k)$ API calls.
- Once the LLM image is known, only $O(d/k)$ calls.
- Intuition: position in a $d$-dimensional subspace is fully specified by $d$ coordinates.

The technical details
○○○○○

Consequences of knowing the LLM image
○●○○○○

What now?
○○○○

# Cheap, full LLM outputs

- Common APIs give top-$k$ tokens and probabilities.
- Logit bias allows boosting tokens to top-$k$.
- Extracting full outputs takes $O(v/k)$ API calls.
- Once the LLM image is known, only $O(d/k)$ calls.
- Intuition: position in a $d$-dimensional subspace is fully specified by $d$ coordinates.

The technical details
○○○○○

Consequences of knowing the LLM image
○●○○○○

What now?
○○○○

# Cheap, full LLM outputs

- Common APIs give top-*k* tokens and probabilities.
- Logit bias allows boosting tokens to top-*k*.
- Extracting full outputs takes $O(v/k)$ API calls.
- Once the LLM image is known, only $O(d/k)$ calls.
- Intuition: position in a $d$-dimensional subspace is fully specified by $d$ coordinates.

The technical details
OOOOO

Consequences of knowing the LLM image
O●OOOO

What now?
OOOO

# Cheap, full LLM outputs

- Common APIs give top-$k$ tokens and probabilities.
- Logit bias allows boosting tokens to top-$k$.
- Extracting full outputs takes $O(v/k)$ API calls.
- Once the LLM image is known, only $O(d/k)$ calls.
- Intuition: position in a $d$-dimensional subspace is fully specified by $d$ coordinates.

The technical details
○○○○○

Consequences of knowing the LLM image
○●○○○○

What now?
○○○○

# Cheap, full LLM outputs

- Common APIs give top-$k$ tokens and probabilities.
- Logit bias allows boosting tokens to top-$k$.
- Extracting full outputs takes $O(v/k)$ API calls.
- Once the LLM image is known, only $O(d/k)$ calls.
- Intuition: position in a $d$-dimensional subspace is fully specified by $d$ coordinates.

The technical details
○○○○○

Consequences of knowing the LLM image
○○●○○○

What now?
○○○○

# Finding the embedding size

Collect at least $d$ outputs from the model, check the dimension of the space that they span.

- Create a matrix $P$ with LLM outputs as columns.
- $P$ will have $d$ nonzero singular values.

The technical details
○○○○○

Consequences of knowing the LLM image
○○●○○○

What now?
○○○○

# Finding the embedding size

Collect at least *d* outputs from the model, check the dimension of the space that they span.

- Create a matrix $P$ with LLM outputs as columns.
- $P$ will have *d* nonzero singular values.

The technical details
○○○○○

Consequences of knowing the LLM image
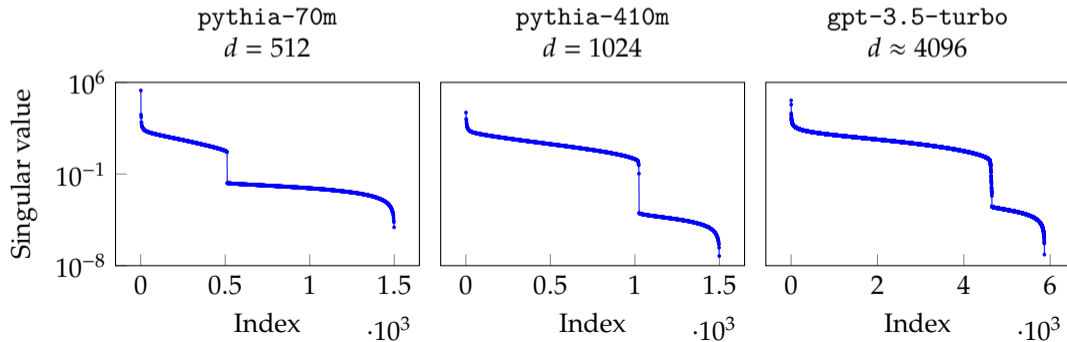○○●○○○

What now?
○○○○

# Finding the embedding size

Collect at least $d$ outputs from the model, check the dimension of the space that they span.

- Create a matrix $P$ with LLM outputs as columns.
- $P$ will have $d$ nonzero singular values.

The technical details
○○○○○

Consequences of knowing the LLM image
○○○●○○

What now?
○○○○

# Output attribution

LLM outputs lie uniquely in the image of the model that generated them.



- AGI Inc.'s new LLM API secretly serves Llama 2.
- AGI Inc. uses a hidden prompt to modify the logits.
- We can catch AGI Inc. because its API outputs remain in Llama 2's image.

The technical details
○○○○○

Consequences of knowing the LLM image
○○○●○○

What now?
○○○○

# Output attribution

LLM outputs lie uniquely in the image of the model that generated them.



- AGI Inc.'s new LLM API secretly serves Llama 2.
- AGI Inc. uses a hidden prompt to modify the logits.
- We can catch AGI Inc. because its API outputs remain in Llama 2's image.

The technical details
○○○○○

Consequences of knowing the LLM image
○○○●○○

What now?
○○○○

# Output attribution

LLM outputs lie uniquely in the image of the model that generated them.



- AGI Inc.'s new LLM API secretly serves Llama 2.

- **AGI Inc. uses a hidden prompt to modify the logits.**

- We can catch AGI Inc. because its API outputs remain in Llama 2's image.

The technical details
○○○○○

Consequences of knowing the LLM image
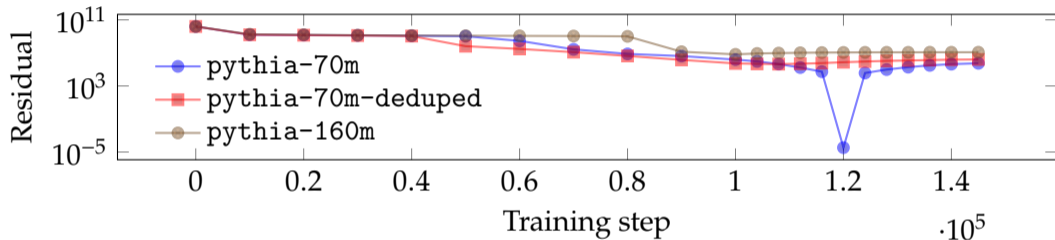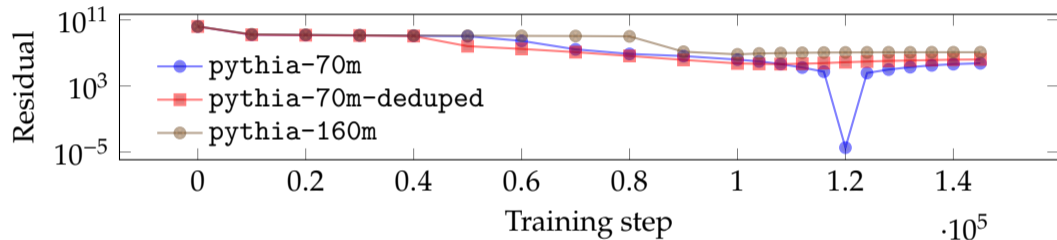○○○●○○

What now?
○○○○

# Output attribution

LLM outputs lie uniquely in the image of the model that generated them.



- AGI Inc.'s new LLM API secretly serves Llama 2.
- AGI Inc. uses a hidden prompt to modify the logits.
- We can catch AGI Inc. because its API outputs remain in Llama 2's image.

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○●○

What now?
○○○○

# Minor vs. major model updates

Table: Your favorite LLM API's logits have changed. What happened?

| Change | Interpretation |
|---|---|
| No logit change, no image change | No update |
| Logit change, no image change | Hidden prompt change or partial finetune |
| Low-rank image change | LoRA update |
| Image change | Full finetune |

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○●○

What now?
○○○○

# Minor vs. major model updates

Table: Your favorite LLM API's logits have changed. What happened?

| Change | Interpretation |
| --- | --- |
| No logit change, no image change | No update |
| Logit change, no image change | Hidden prompt change or partial finetune |
| Low-rank image change | LoRA update |
| Image change | Full finetune |

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○●○

What now?
○○○○

# Minor vs. major model updates

Table: Your favorite LLM API's logits have changed. What happened?

| Change | Interpretation |
| --- | --- |
| No logit change, no image change | No update |
| Logit change, no image change | Hidden prompt change or partial finetune |
| Low-rank image change | LoRA update |
| Image change | Full finetune |

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○●○

What now?
○○○○

# Minor vs. major model updates

Table: Your favorite LLM API's logits have changed. What happened?

| Change | Interpretation |
|---|---|
| No logit change, no image change | No update |
| Logit change, no image change | Hidden prompt change or partial finetune |
| Low-rank image change | LoRA update |
| Image change | Full finetune |

The technical details
OOOOO

Consequences of knowing the LLM image
OOOO●O

What now?
OOOO

# Minor vs. major model updates

Table: Your favorite LLM API's logits have changed. What happened?

| Change | Interpretation |
| --- | --- |
| No logit change, no image change | No update |
| Logit change, no image change | Hidden prompt change or partial finetune |
| Low-rank image change | LoRA update |
| Image change | Full finetune |

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○●

What now?
○○○○

# Other uses for LLM images

- Unargmaxable tokens
- Recovering the softmax matrix
- Basis-aware sampling

**Low-Rank Softmax Can Have Unargmaxable Classes in Theory
but Rarely in Practice**

**Andreas Grivas** and **Nikolay Bogoychev** and **Adam Lopez**
Institute for Language, Cognition, and Computation
School of Informatics
University of Edinburgh
{agrivas, n.bogoych, alopez}@ed.ac.uk

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○●

What now?
○○○○

# Other uses for LLM images

- Unargmaxable tokens
- **Recovering the softmax matrix**
- Basis-aware sampling

**Low-Rank Softmax Can Have Unargmaxable Classes in Theory but Rarely in Practice**

**Andreas Grivas** and **Nikolay Bogoychev** and **Adam Lopez**
Institute for Language, Cognition, and Computation
School of Informatics
University of Edinburgh
{agrivas, n.bogoych, alopez}@ed.ac.uk

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○●

What now?
○○○○

# Other uses for LLM images

- Unargmaxable tokens
- Recovering the softmax matrix
- Basis-aware sampling

**Low-Rank Softmax Can Have Unargmaxable Classes in Theory but Rarely in Practice**

**Andreas Grivas** and **Nikolay Bogoychev** and **Adam Lopez**
Institute for Language, Cognition, and Computation
School of Informatics
University of Edinburgh
{agrivas, n.bogoych, alopez}@ed.ac.uk

CLOSING THE CURIOUS CASE OF NEURAL TEXT
DEGENERATION

**Matthew Finlayson**[*]
University of Southern California

**John Hewitt**
Stanford University

**Alexander Koller**
Saarland University

**Swabha Swayamdipta**
University of Southern California

**Ashish Sabharwal**
The Allen Institute for AI

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
●○○○

Section 3

What now?

# Mitigations

| Proposal | Cons |
| --- | --- |
| Discontinue top-$k$ probs | Only slows attack |
| Remove softmax bottleneck | Expensive training, inference |
| Discontinue logit bias | Nerfs API |

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○●○○

# Mitigations

| Proposal | Cons |
| --- | --- |
| Discontinue top-*k* probs | Only slows attack |
| Remove softmax bottleneck | Expensive training, inference |
| Discontinue logit bias | Nerfs API |

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○●○○

# Mitigations

| Proposal | Cons |
| --- | --- |
| Discontinue top-*k* probs | Only slows attack |
| Remove softmax bottleneck | Expensive training, inference |
| Discontinue logit bias | Nerfs API |

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○●○○

# Mitigations

| Proposal | Cons |
|---|---|
| Discontinue top-*k* probs | Only slows attack |
| Remove softmax bottleneck | Expensive training, inference |
| Discontinue logit bias | Nerfs API |

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○●○○

## Mitigations

| Proposal | Cons |
| --- | --- |
| Discontinue top-$k$ probs | Only slows attack |
| Remove softmax bottleneck | Expensive training, inference |
| Discontinue logit bias | Nerfs API |

Recommendation: do nothing; LLM images are useful for accountability.

The technical details
ooooo

Consequences of knowing the LLM image
oooooo

What now?
oo●o

## Some future directions

- Efficient image extraction methods for strict APIs.
- More audit methods for LLMs.
- Stealing more than the image.

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○●○

## Some future directions

- Efficient image extraction methods for strict APIs.
- More audit methods for LLMs.
- Stealing more than the image.

The technical details
00000

Consequences of knowing the LLM image
000000

What now?
00●0

# Some future directions

- Efficient image extraction methods for strict APIs.
- More audit methods for LLMs.
- Stealing more than the image.

The technical details
○○○○○

Consequences of knowing the LLM image
○○○○○○

What now?
○○○●

# Thank you for coming!

- LLM outputs occupy a low-dimensional space: the *image*.
- Common API interfaces leak the LLM image.
- LLM images expose non-public information.
- LLM images are a tool for accountability.