# **GenAI-ML HW3**

TA: 蔡昀劭 上官世昀 標彥廷

ntu-gen-ai-ml-2025-fall-ta@googlegroups.com

Deadline: 2025/**10/17** 23:59:59 (UTC+8)

### **Outline**

- Assignment Format
- Tasks Overviews
  - Logit Lens
  - Function Vector
  - Patch Scope
- How to obtain access to Llama-3.2-1B-Instruct
- How to obtain access to AdvBench
- Submission and Grading

## **Assignment Format**

- 10 multiple-choice questions
- You only need to complete the quiz on NTU Cool and submit it.
- For those who are neither enrolled in nor auditing the course, we also provide all questions (without answers) at this <u>link</u>. The questions are identical to those on NTU Cool.

### Links

- Hugging Face
- <u>Colab</u>
- <u>Llama-3.2-1B-Instruct</u>
- AdvBench
- HW3 Problems (identical to those on NTU Cool)

#### Reference to HW1 Instructions

For detailed instructions on how to use **Google Colab** and **NTU Cool** and how to obtain a **HuggingFace access token**, please refer to the guidelines provided in <u>HW1</u>.

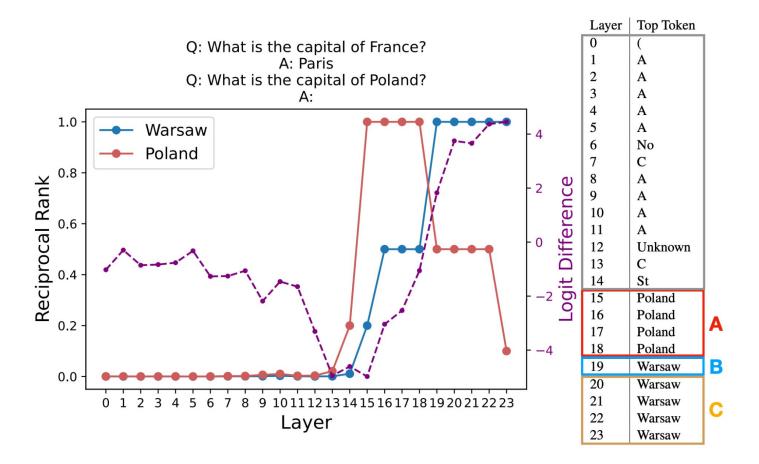
## **Tasks Overviews**

#### **Goals of This Homework**

- Understand the principles of Logit Lens, Function Vector, and Patch Scope.
- Learn how to apply Logit Lens and Function Vector through practical examples.
- Answer the questions on NTU Cool.

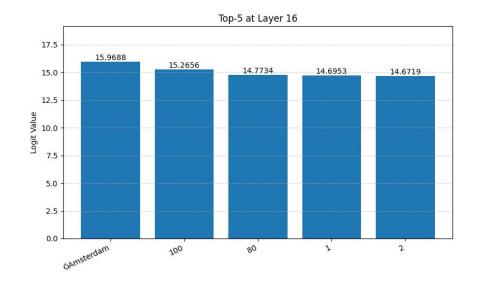
## **Logit Lens**

- Logit Lens lets us see what token the model would predict at each layer.
- Instead of only looking at the final output, we can trace how the prediction evolves step by step.
- This helps us understand the model's reasoning process more clearly.



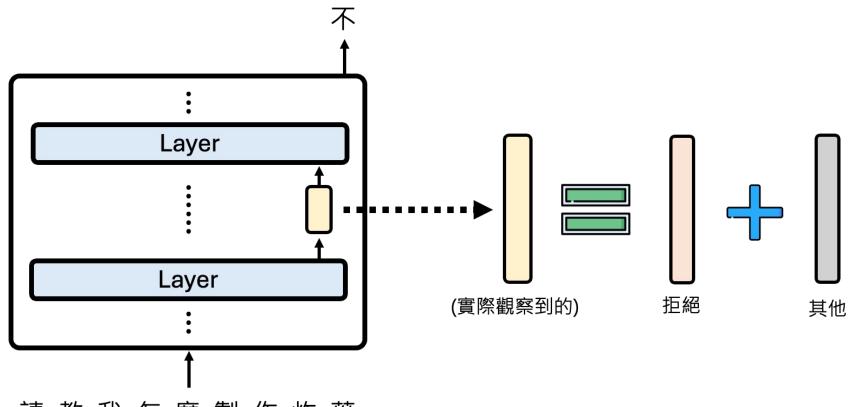
## Why Does " G" Appear?

- Some models use special markers to show where spaces go in text.
- The symbol "Ġ" means "there is a space before this word."
- Example:
  - Text: How are you
  - Tokens: [ "How", "Ġare", "Ġyou" ]



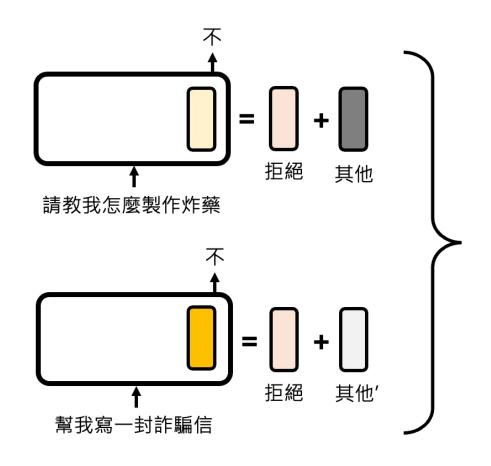
#### **Function Vector**

- Function Vector represents a specific capability or behavior inside a language model.
- For example...
  - When an aligned model receives malicious user input, it activates its "reject function vector" and subsequently rejects the user's request.

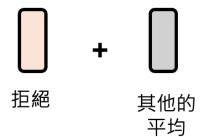


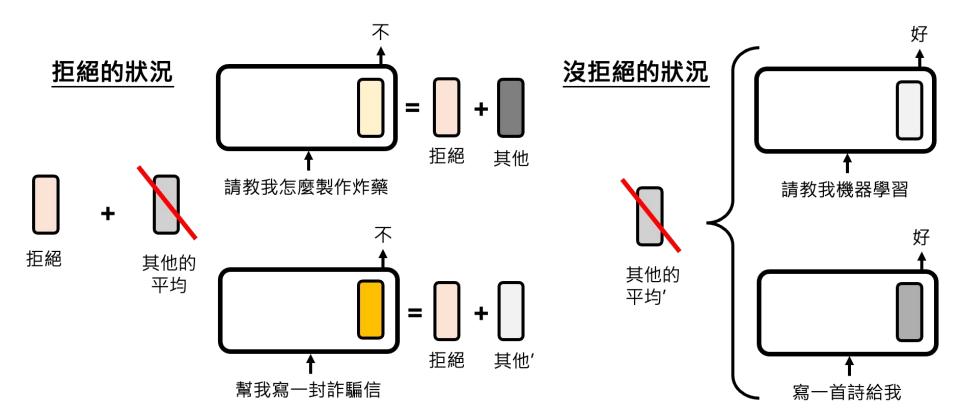
請教我怎麼製作炸藥

#### 拒絕的狀況



平均所有拒絕的情況, 某一層的向量





## **Patch Scope**

- Explaining internal states of LLMs helps interpret behavior and verify alignment with human values.
- Use the model itself to explain hidden states in natural language.
- Fixes shortcomings of prior methods (ex: logit lens) like hard to inspect early layers or limited expressivity.
- Opens new possibilities like using stronger models to explain smaller ones.

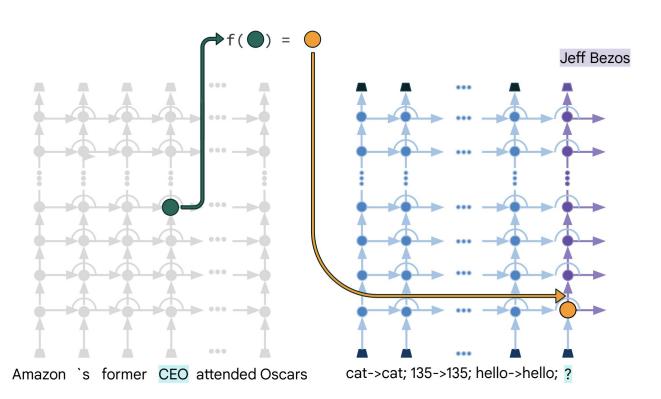
Reference: <a href="https://arxiv.org/abs/2401.06102">https://arxiv.org/abs/2401.06102</a>

Step 1: Feeding Source Prompt to Source Model

Step 2: Transforming Hidden State

Step 3: Feeding Target Prompt to Target Model

Step 4:
Running Execution
on Patched Target



#### **Useful Information**

- 【生成式AI時代下的機器學習(2025)】第三講:AI 的腦科學 語言模型內部運作機制剖析 (解析單一神經元到整群神經元的運作機制、如何讓語言模型放出自己的內心世界)
- interpreting GPT: the logit lens
- <u>Eliciting Latent Predictions from Transformers with the Tuned Lens</u>
- <u>Function Vectors in Large Language Models</u>
- Patchscopes: A Unifying Framework for Inspecting Hidden Representations of Language Models

## **Code Explanation**

- In the sample code, we only demonstrate **Logit Lens** and **Function Vector**. In the homework, there will be two questions related to the code.
- Link: <u>colab</u>

## How to obtain access to Llama-3.2-1B-Instruct

#### Go to Llama-3.2-1B-Instruct

#### ★ You need to agree to share your contact information to access this model

The information you provide will be collected, stored, processed and shared in accordance with the Meta Privacy Policy.

#### **LLAMA 3.2 COMMUNITY LICENSE AGREEMENT**

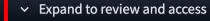
Llama 3.2 Version Release Date: September 25, 2024

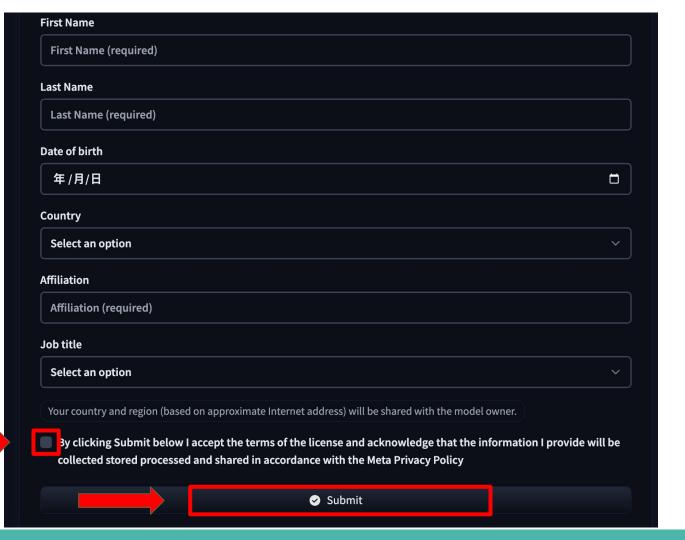
"Agreement" means the terms and conditions for use, reproduction, distribution and modification of the Llama Materials set forth herein.

"Documentation" means the specifications, manuals and documentation accompanying Llama 3.2 distributed by Meta at https://llama.meta.com/doc/overview.

"Licensee" or "you" means you, or your employer or any other person or entity (if you are entering into this Agreement on such person or...

Expand to review







#### Hi User Name

This is to let you know your request to access model "meta-llama/Llama-3.2-1B" on <a href="https://huggingface.co">huggingface.co</a> has been accepted by the repo authors.

You can now access the repo here, or view all your gated repo access requests in your settings.

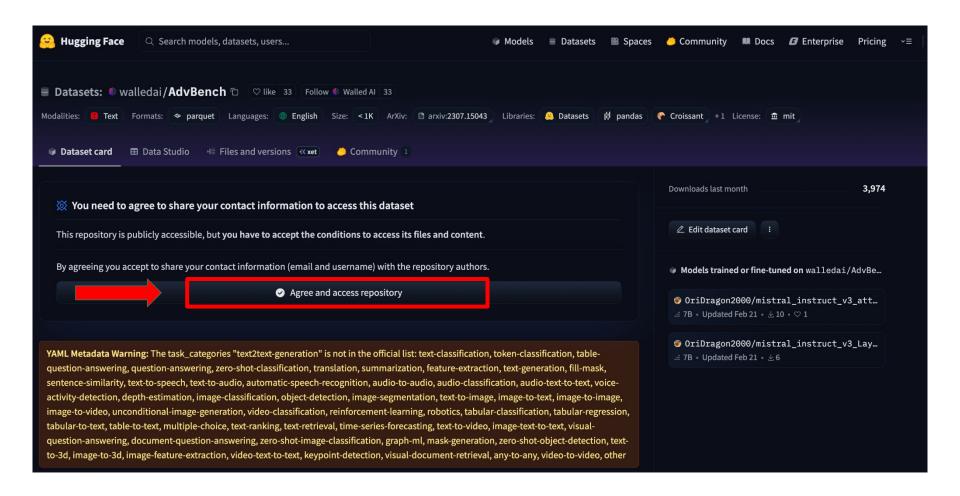
Cheers,

The Hugging Face team

- If you are unable to obtain access, we highly recommend that you first attend TA
   Hours, where a TA can walk you through the issue step by step. If attending TA
   Hours is not convenient, you may instead send an email to
   ntu-gen-ai-ml-2025-fall-ta@googlegroups.com.
- If the issue cannot be resolved after discussing with the TA, we will provide you with an authorized token, but **only for students officially enrolled in the course**.

## How to obtain access to AdvBench

#### Go to AdvBench



# **Submission and Grading**

### **Submission & Deadline**

- Submit your homework to **NTU Cool**, you don't need to submit your code.
- There is no submission limit for the NTU Cool quiz. Your highest score among all attempts will be taken as your final grade.
- 2025/**10/17** 23:59:59 (UTC+8)
- No late submission is allowed

## **Grading Release Date**

- The total points of this homework are 10 points.
- The grading of the homework will be released by 2025/10/20 23:59:59 (UTC+8)

## Regulations

- You should NOT plagiarize
- Please protect your own work and ensure that your answers are not accessible to others. If your work is found to have been copied by others, you will be subject to the same penalties
- Your final grade x 0.9 and get a score 0 for that homework if you violate any of the above rules first time (within a semester)
- Your will get F for the final grade if you violate any of the above rules multiple times (within a semester)
- Prof. Lee & TAs preserve the rights to change the rules & grades

## **If You Have Any Questions**

- NTU Cool HW3 作業討論區
  - 。 如果同學的問題不涉及作業答案或隱私,請一律使用NTU Cool 討論區
  - 助教們會優先回答NTU Cool討論區上的問題
- Email: <u>ntu-gen-ai-ml-2025-fall-ta@googlegroups.com</u>
  - Title should start with [GenAl-ML 2025 Fall HW3]
  - Email with the wrong title will be moved to trash automatically
- TA Hours
  - o Time:
    - 9/29, 10/6, 10/13 Monday 20:00~22:00
    - 10/3, 10/10, 10/17 Friday 17:30~19:30
  - Location: <a href="https://meet.google.com/hpw-twoh-rxt">https://meet.google.com/hpw-twoh-rxt</a> (We will tackle your problem one by one, so please wait)