

# Understanding Knowledge Storage and Transfer in Multimodal Language Models

Samyadeep Basu<sup>§</sup>, Martin Grayson<sup>†</sup>, Cecily Morrison<sup>†</sup>, Besmira Nushi<sup>†</sup>, Soheil Feizi<sup>§</sup>,  
Daniela Massiceti<sup>†</sup>

§: University of Maryland, College Park; †: Microsoft Research



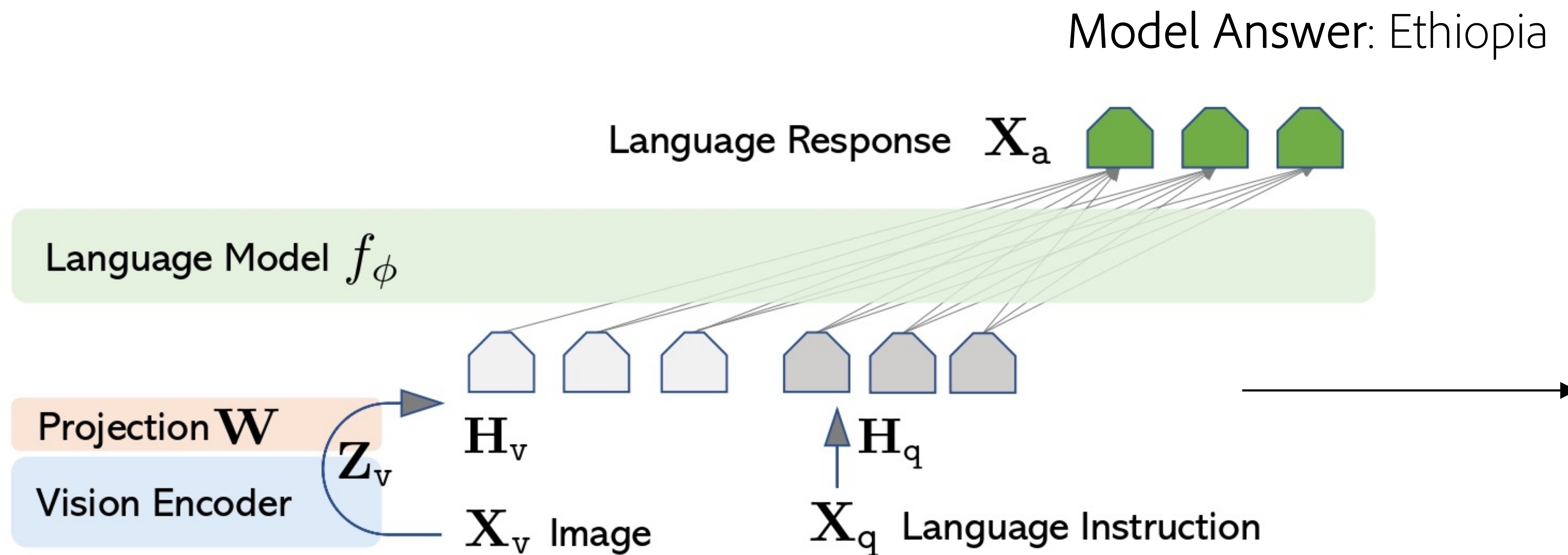
KGM @ ECCV 2024 | Appearing in NeurIPS 2024  
<https://aka.ms/mlm-info-storage>

# About Me

- Final year CS PhD student at University of Maryland with Dr. Soheil Feizi
- **Current Research Interests:** Multimodal /Vision/Language Models with an emphasis on controlling them via Interpretability
  - *DiffQuickFix, LocoEdit*: Light-weight Model Editing for Text-to-Image Models
  - *CompAlign* - Interpreting arbitrary ViT components (e.g., attention heads) with text
    - Zero-shot segmentation
    - Spurious Correlation Mitigation

Motivation

# Multimodal LLMs are Widely Used in VQA



Input Image

Question: Where does this dish come from?

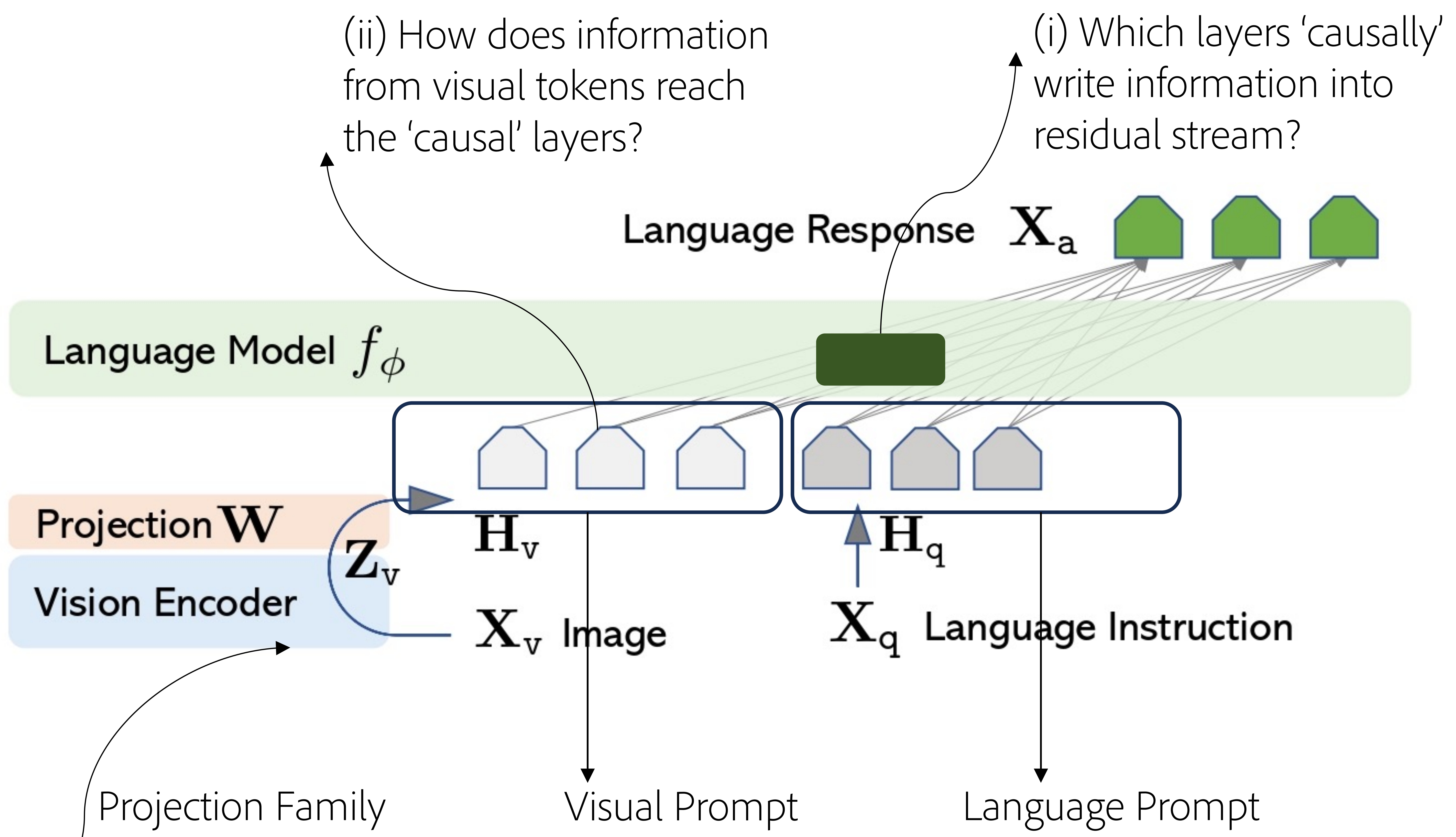
Fig 1: Representative VQA task

But we *lack scientific* understanding on:

- (i) How they internally process information
- (ii) How can we control them for tasks like model editing to fix failure modes / introduce rare concepts

We study how MLLMs process and transfer information in a factual VQA task using a constraint-based formulation

We study LLaVA and multi-modal Phi-2

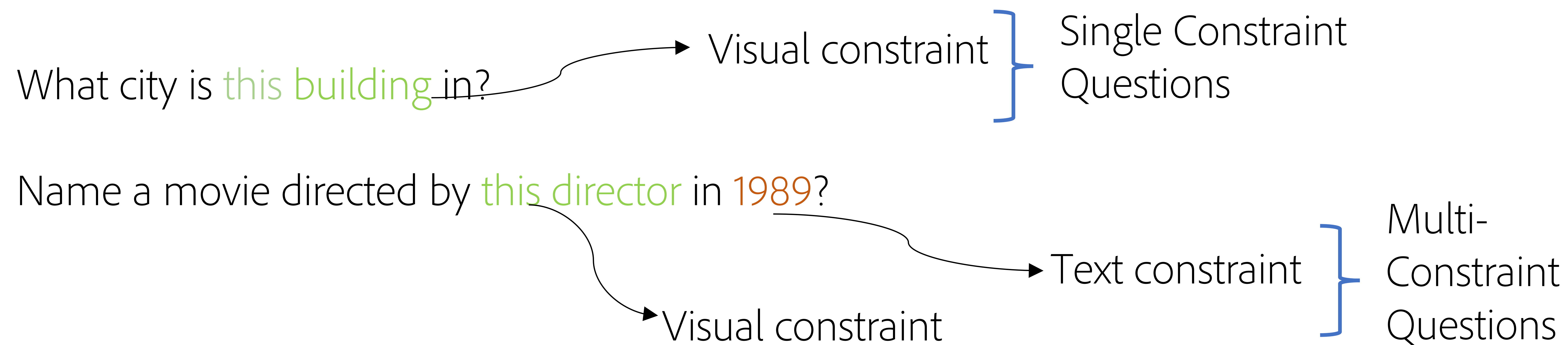


**Causal Layer :** A layer or a small set of layers which control the output of the model conditioned on an input

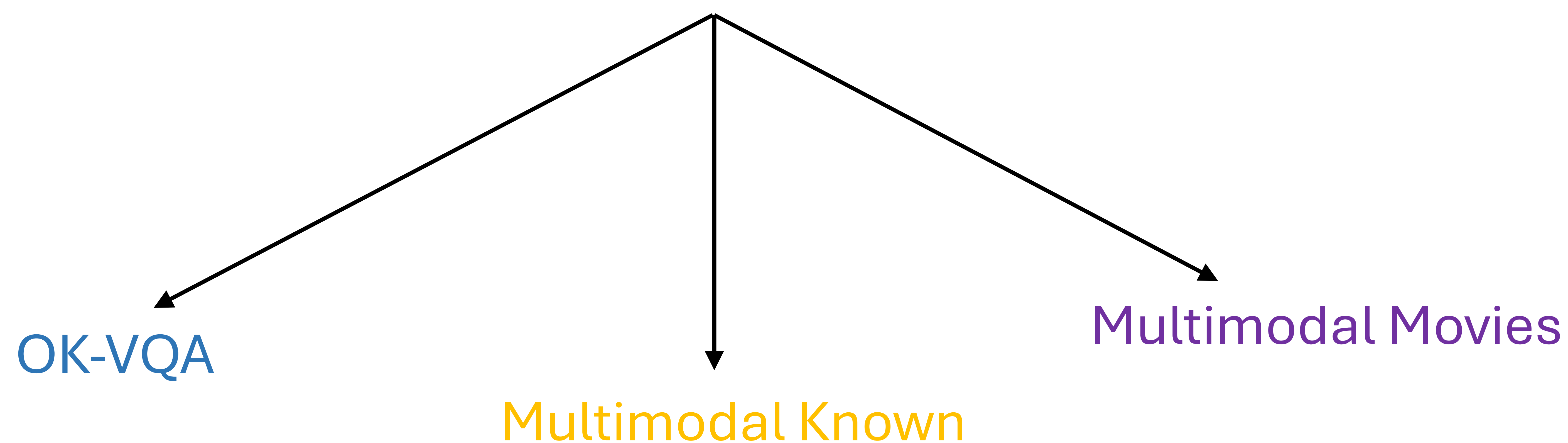
# Designing a Probe Dataset

# Introducing Constraint Based Formulation for Interpreting MLLMs

## Constraint based VQA Questions



VQA-Constraints Dataset: 9.7k VQA questions annotated with constraints



# Constraint based VQA Questions



What city is **this building** in?

Visual constraint

Single Constraint Questions

Name a movie directed by **this director** in **1989**?



Visual constraint

Text constraint

Multi-Constraint Questions

**Visual Constraint:** A set of tokens in the question which relates to a visual entity in an image (e.g., this building → Space Needle in Image)

**Text Constraint:** A set of tokens in the question which reduces the space of the possible answers. Often used in conjunction with visual constraints.

Introducing  
Constraint Based  
Formulation for  
Interpreting MLLMs

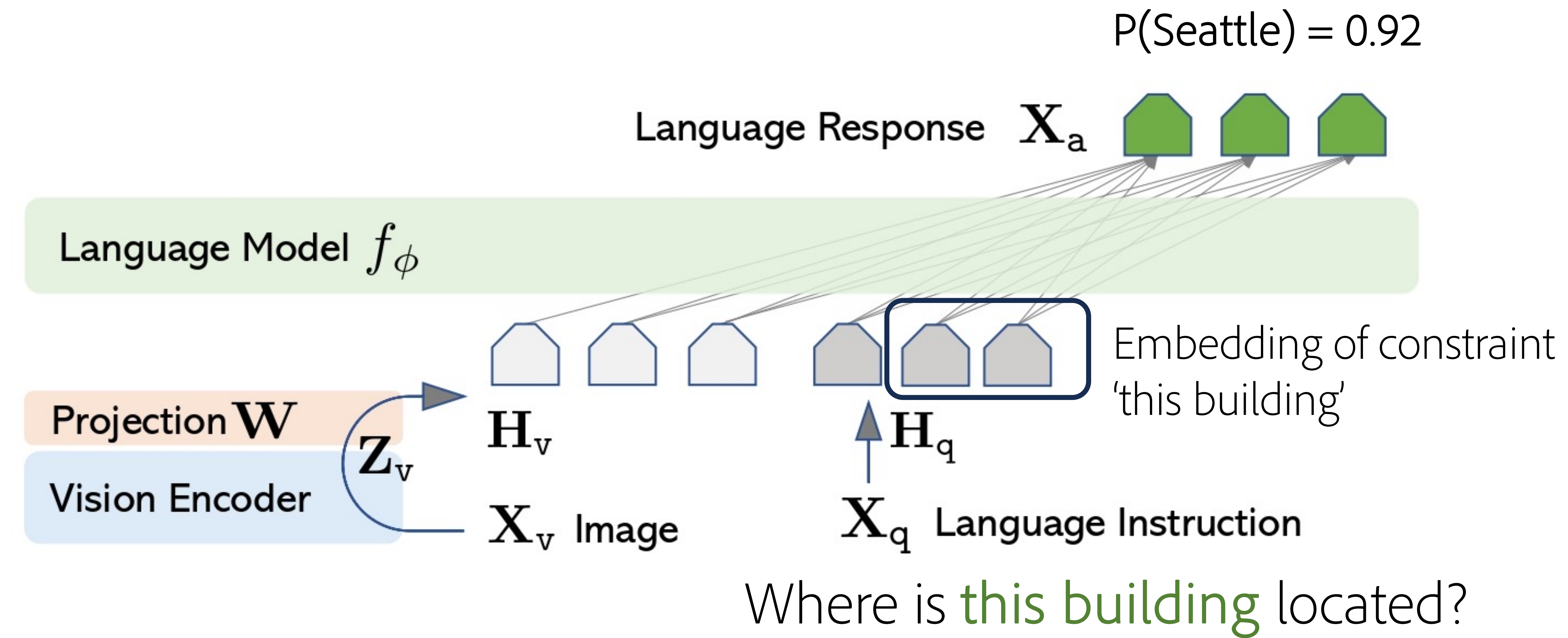


# Tracing Method for Knowledge Storage

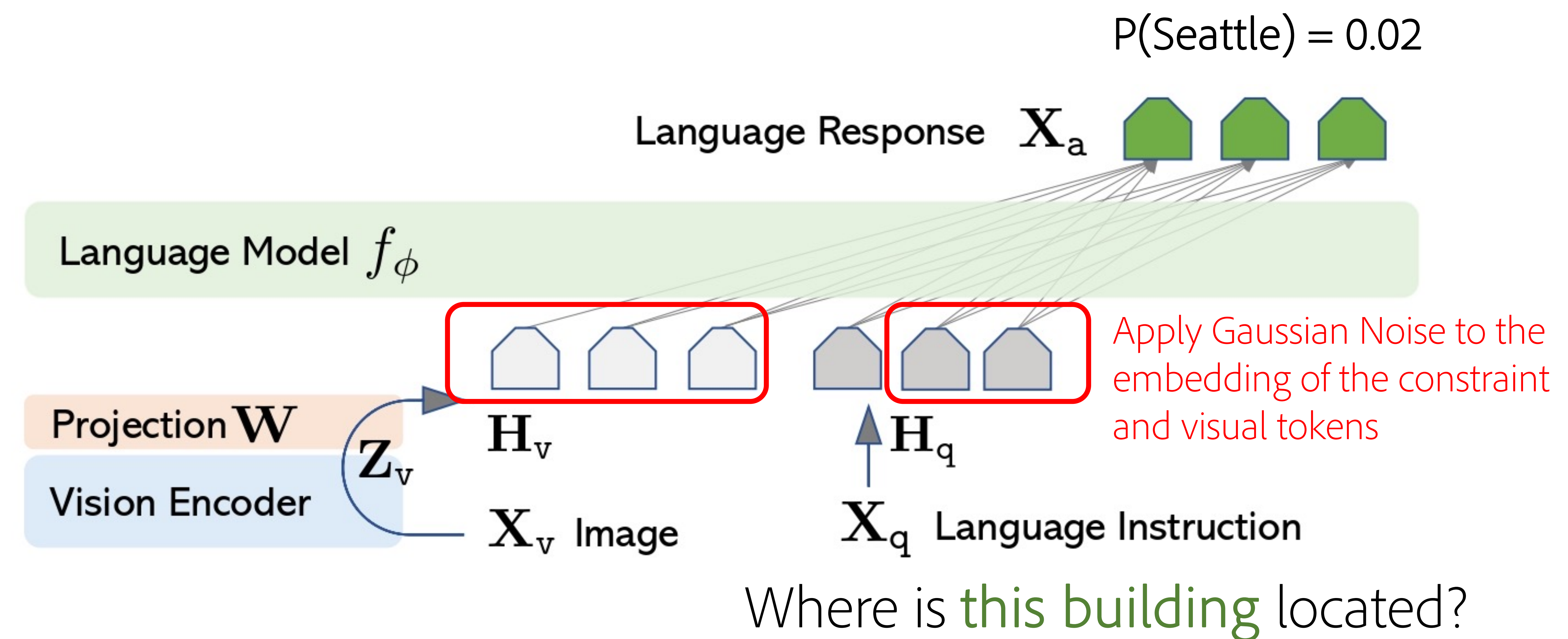
# Using Causal Trace Designed for Language Models

## MultimodalCausalTrace: Identifying Causal Layers in MLLMs

Clean Model



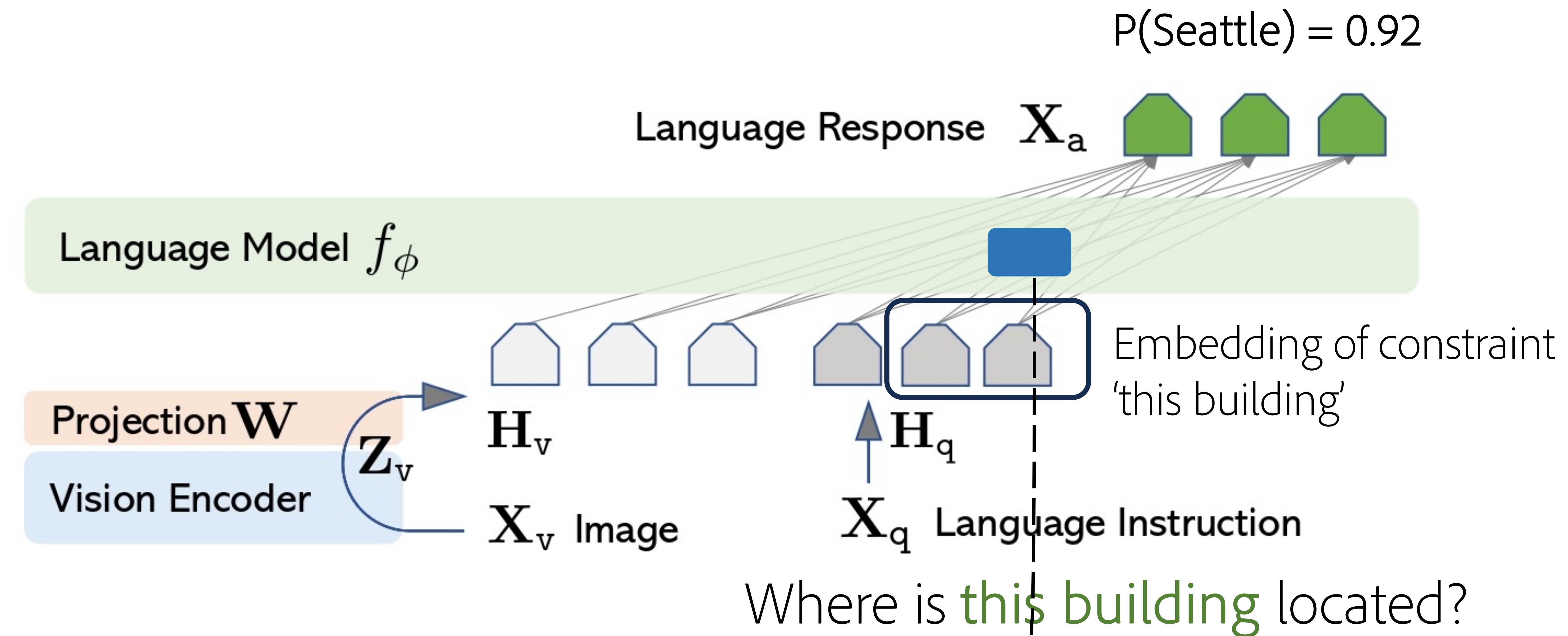
Corrupted Model



# Using Causal Trace Designed for Language Models

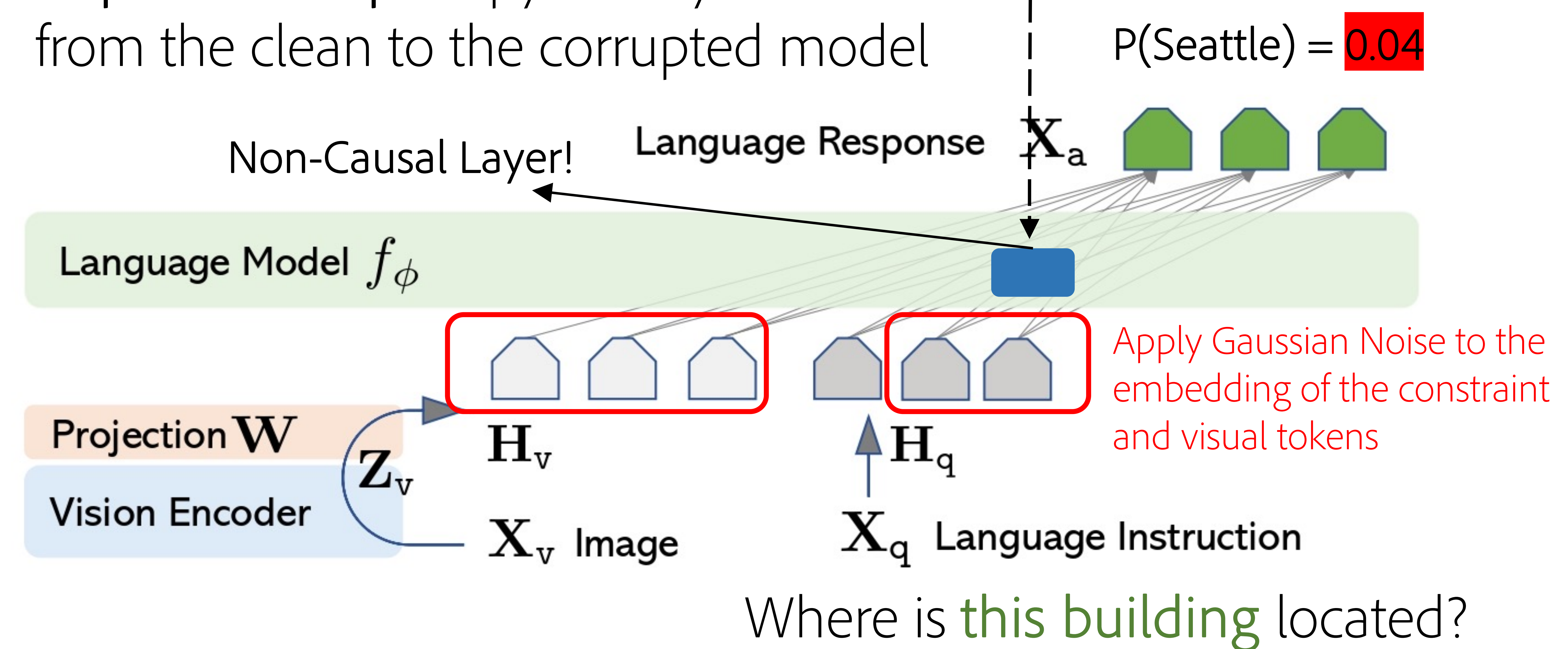
## MultimodalCausalTrace: Identifying Causal Layers in MLLMs

Clean Model



Important Step: Copy the layer activations from the clean to the corrupted model

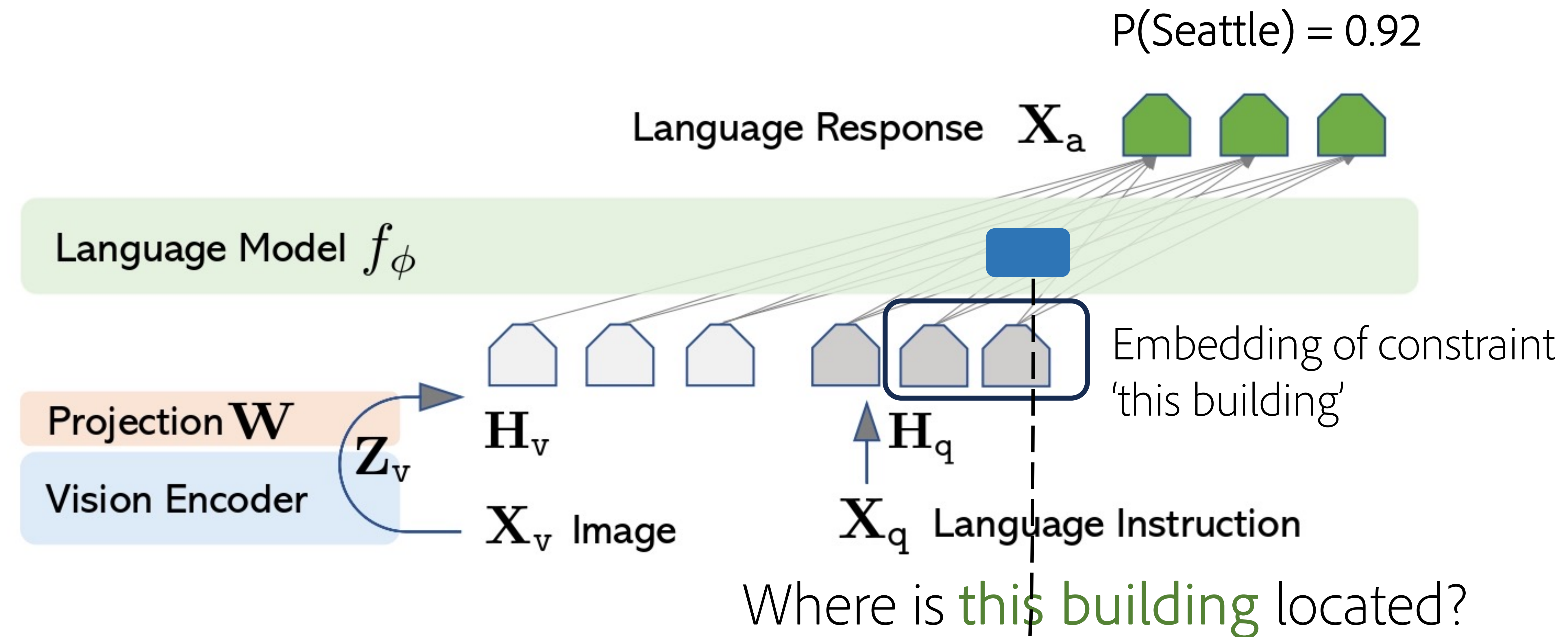
Corrupted →  
Restored Model



# Using Causal Trace Designed for Language Models

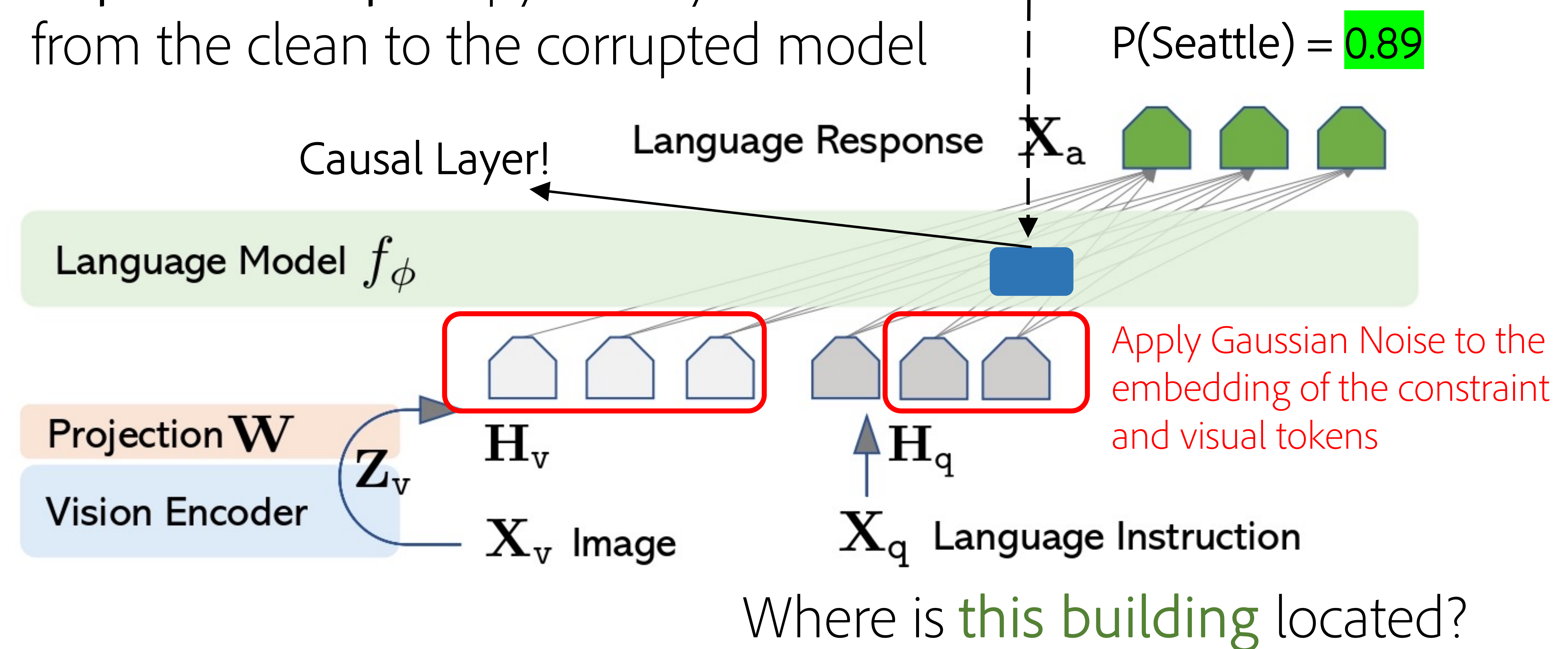
## MultimodalCausalTrace: Identifying Causal Layers in MLLMs

Clean Model



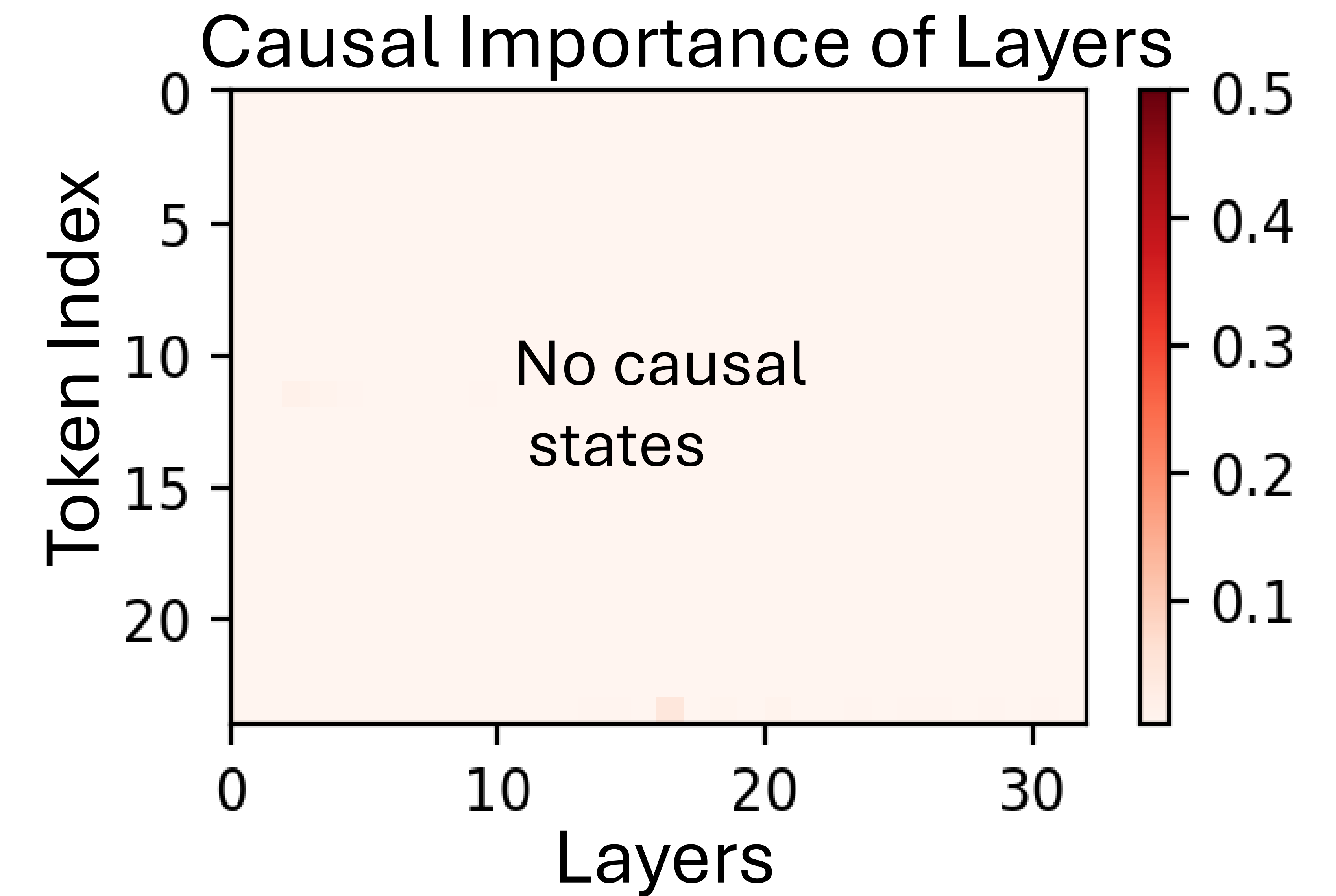
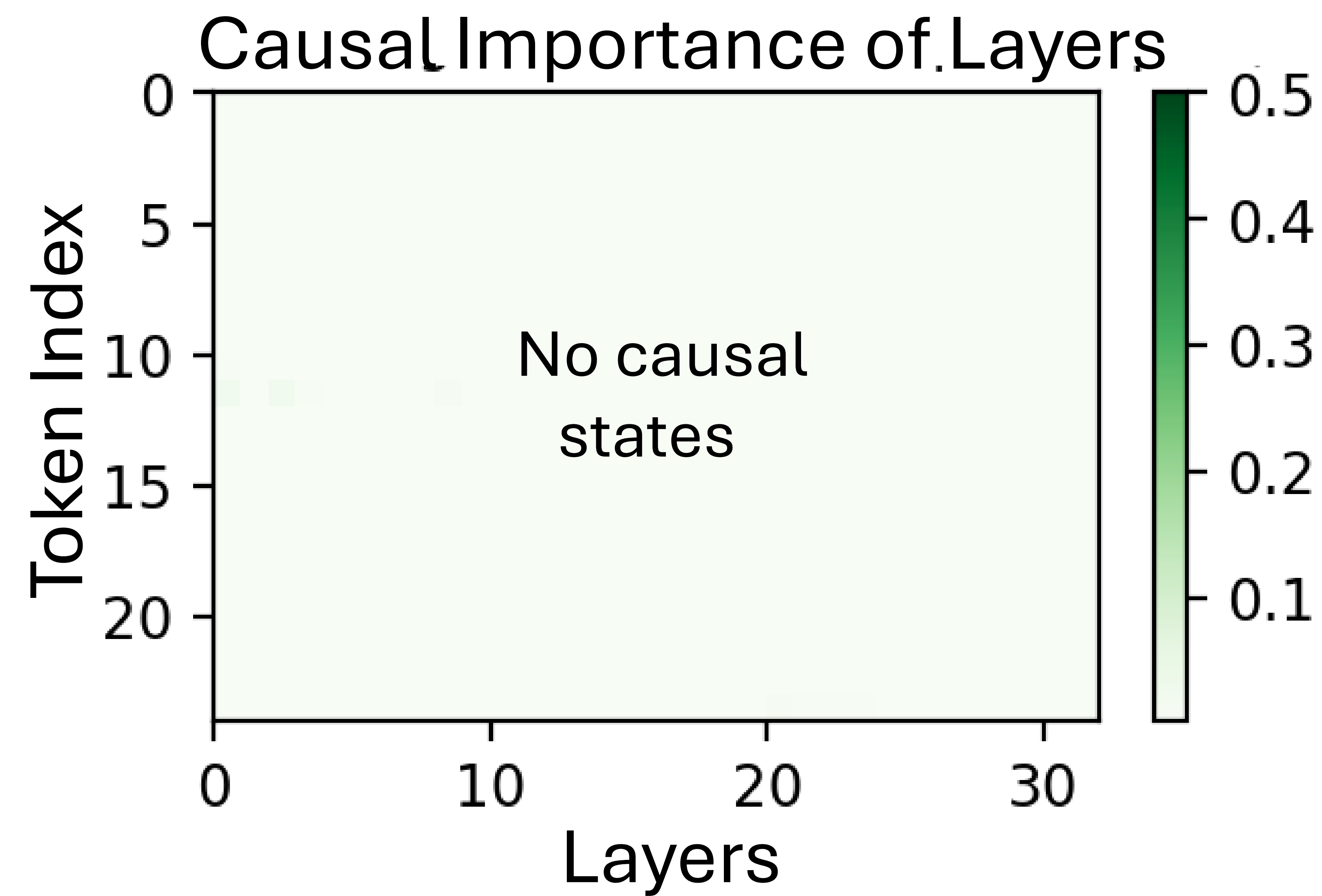
Important Step: Copy the layer activations from the clean to the corrupted model

Corrupted  $\rightarrow$   
Restored Model



# Using Causal Trace Designed for Language Models - Results

No Relevant Causal Layers!

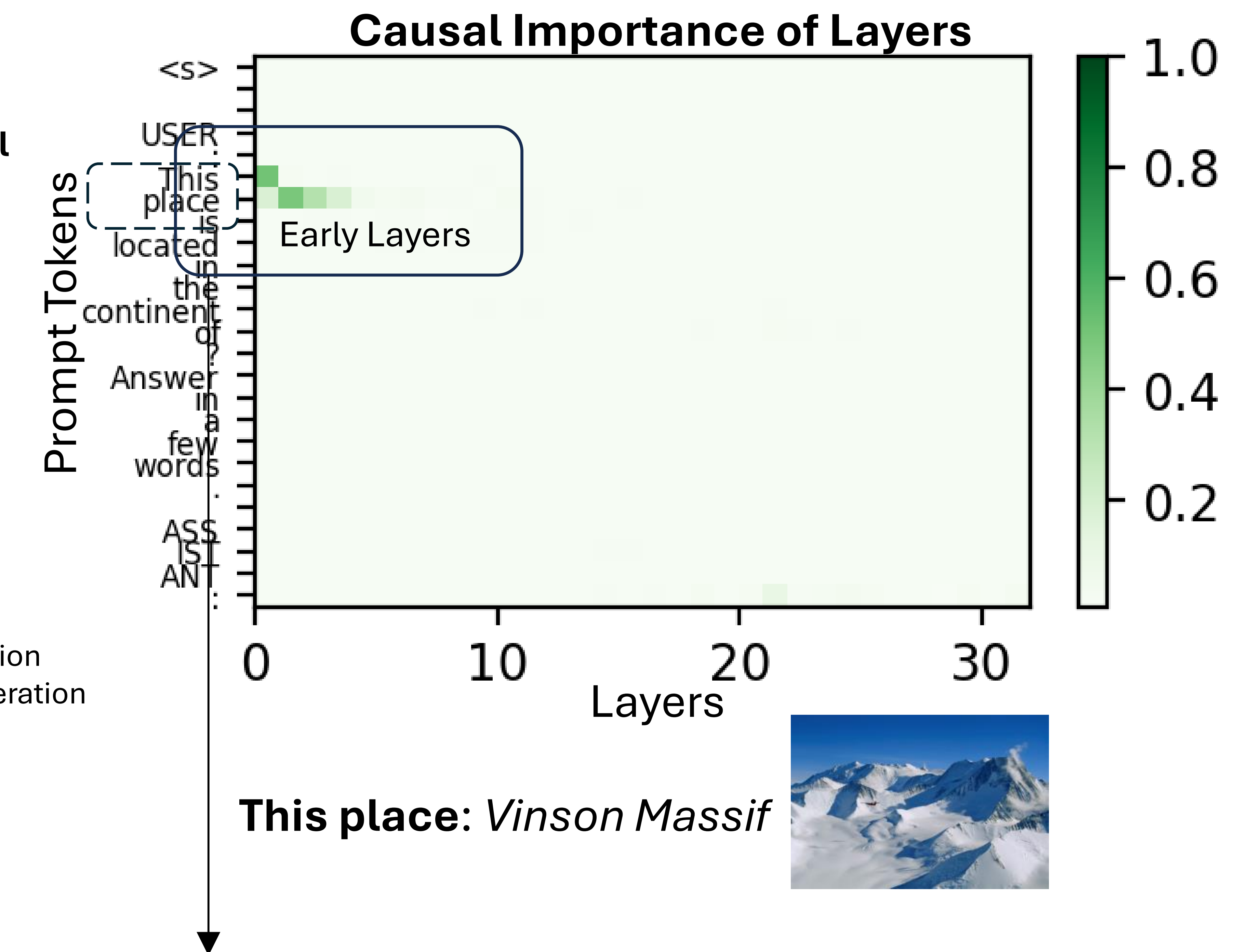
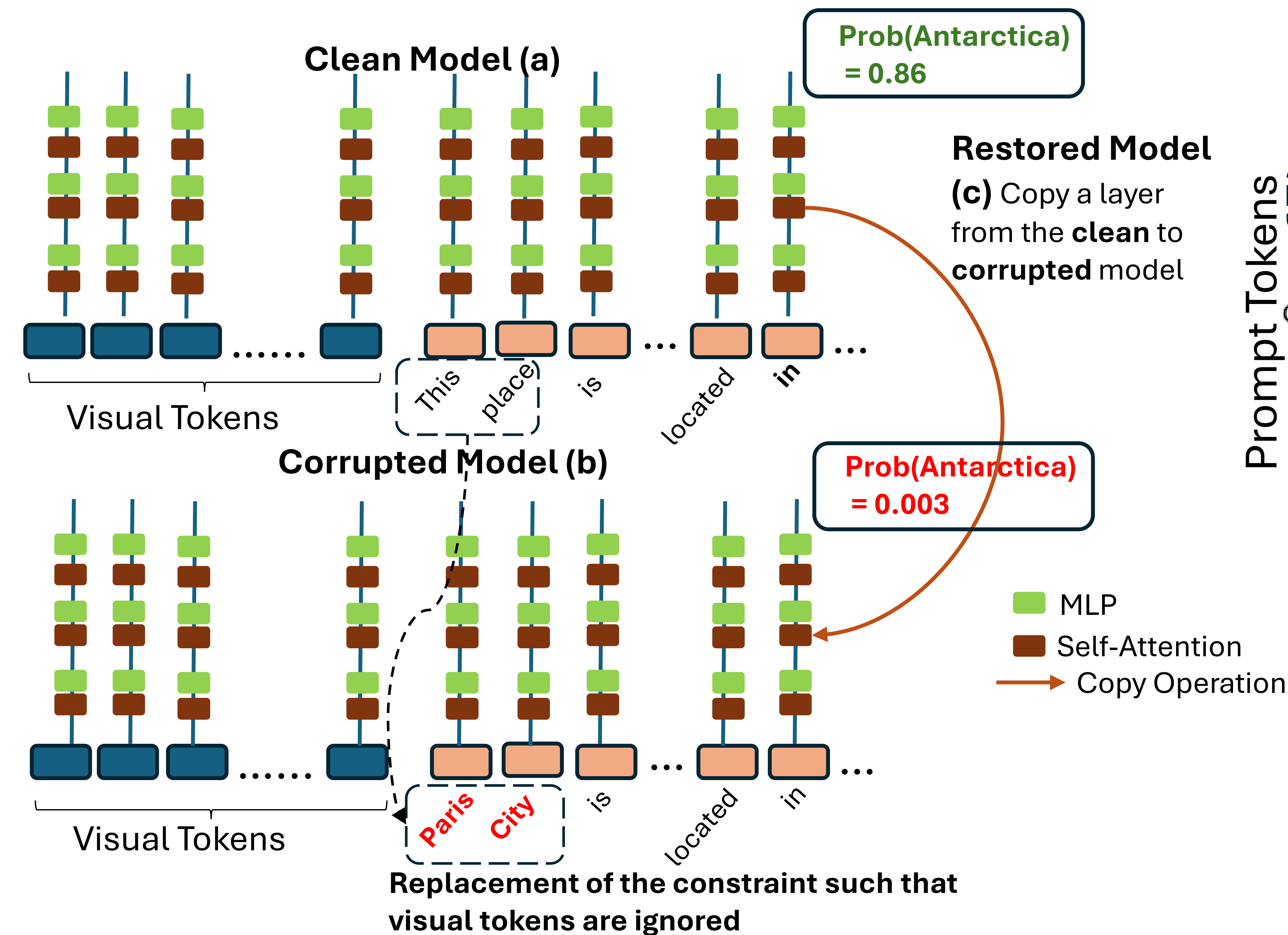


Potential Reason: Corrupted Model is very noisy as noise is applied to a large set of tokens (e.g., 576 + number of constraint tokens), whereas in a language model noise is only applied to a small set of tokens (e.g., 2-4)

Multimodal Causal Trace:  
Identifying Causal Layers  
in MLLMs

# Adapting Causal Trace Designed for Language Models

Multimodal Causal Trace:  
Identifying Causal Layers  
in MLLMs  
  
Our Method



Early MLP layers are extracted to be 'causal'

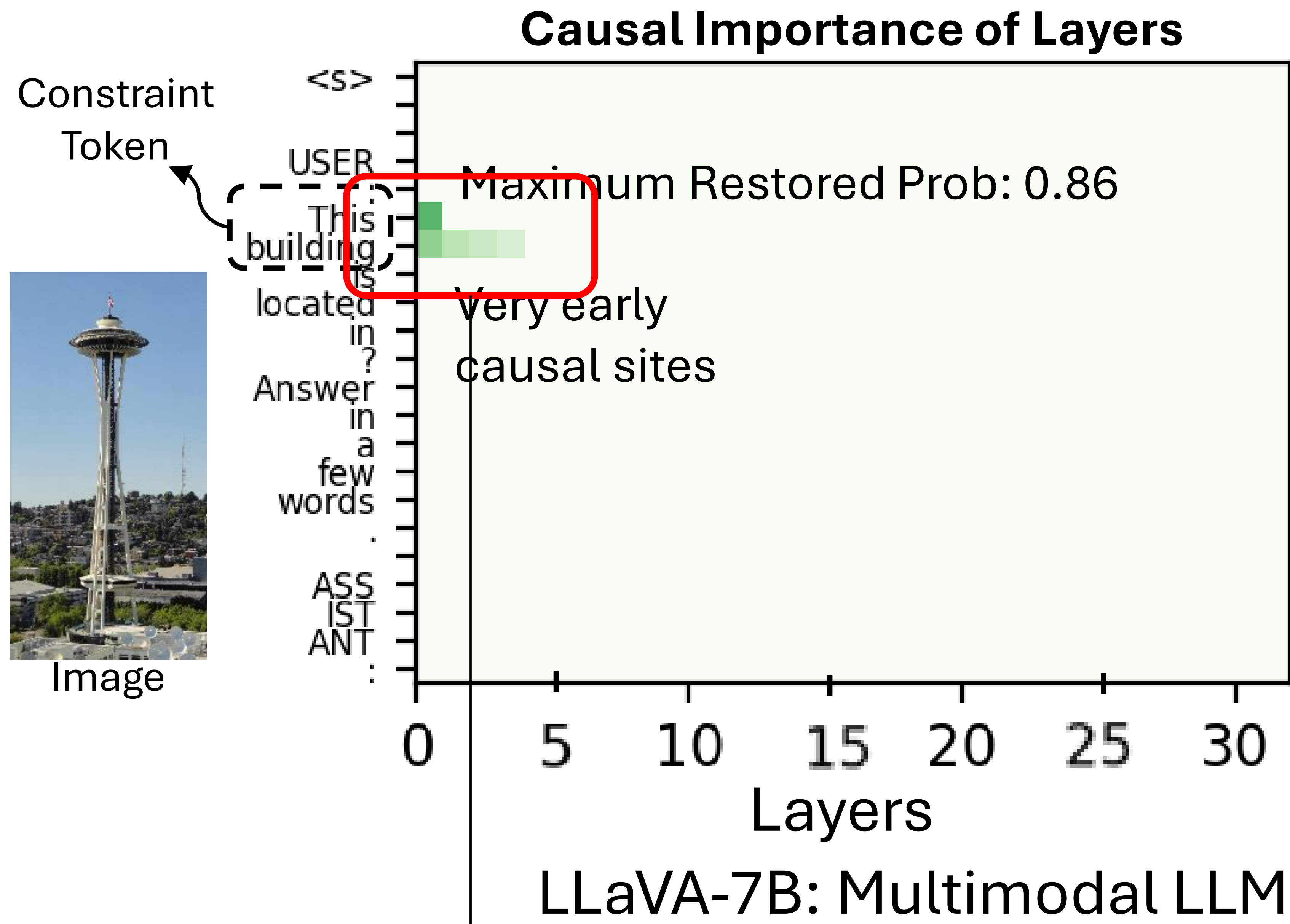
**Core Idea :** Replace the constraint tokens with a set of tokens such that the visual tokens are ignored while answering the question

# Tracing Results

How does information retrieval from internal layers in MLLMs differ from language models?

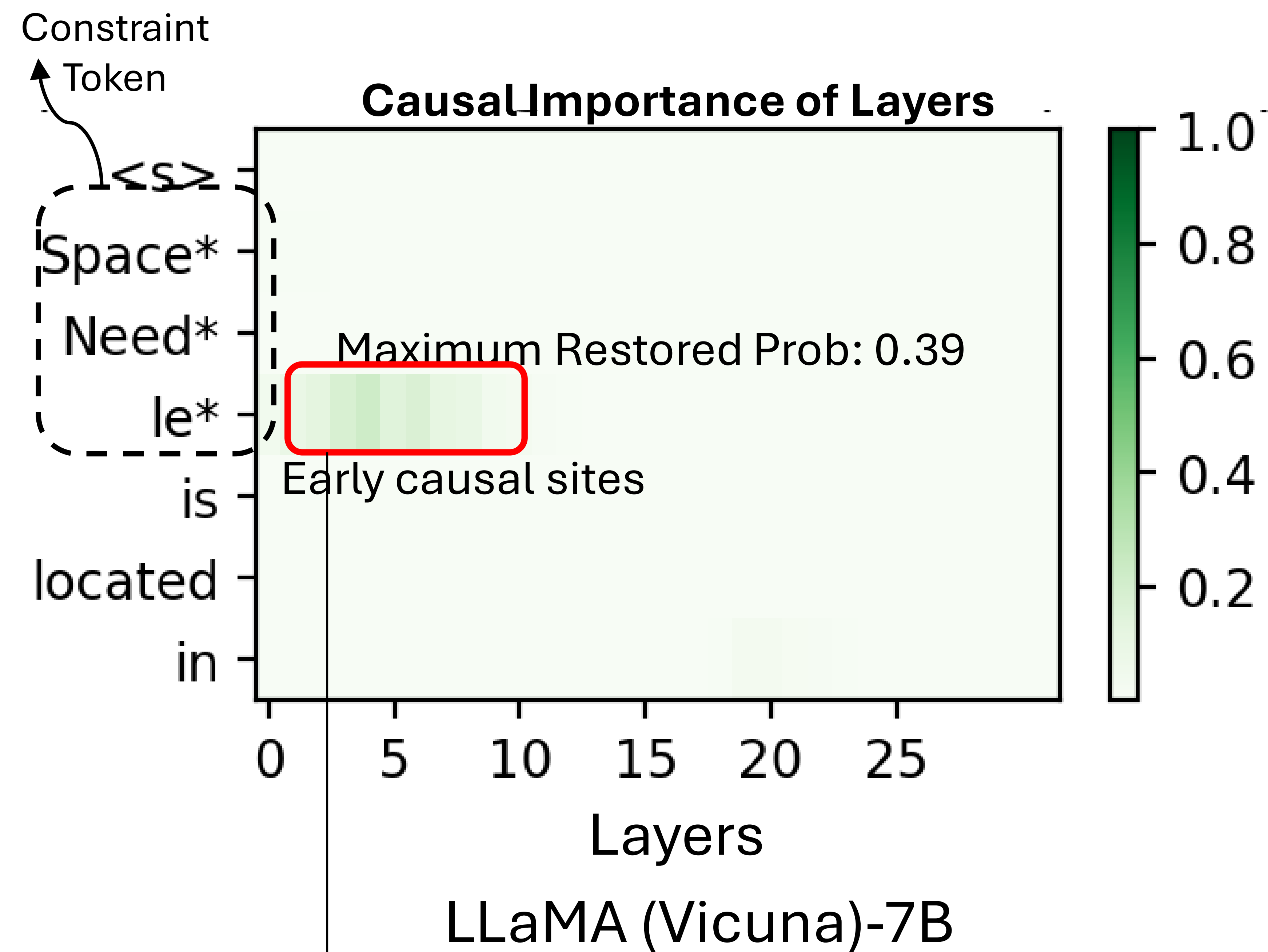
# Multimodal Causal Trace: Identifying Causal Layers in MLLMs

## Single Constraint Questions



Layers 1-4

Restoring only one layer is sufficient



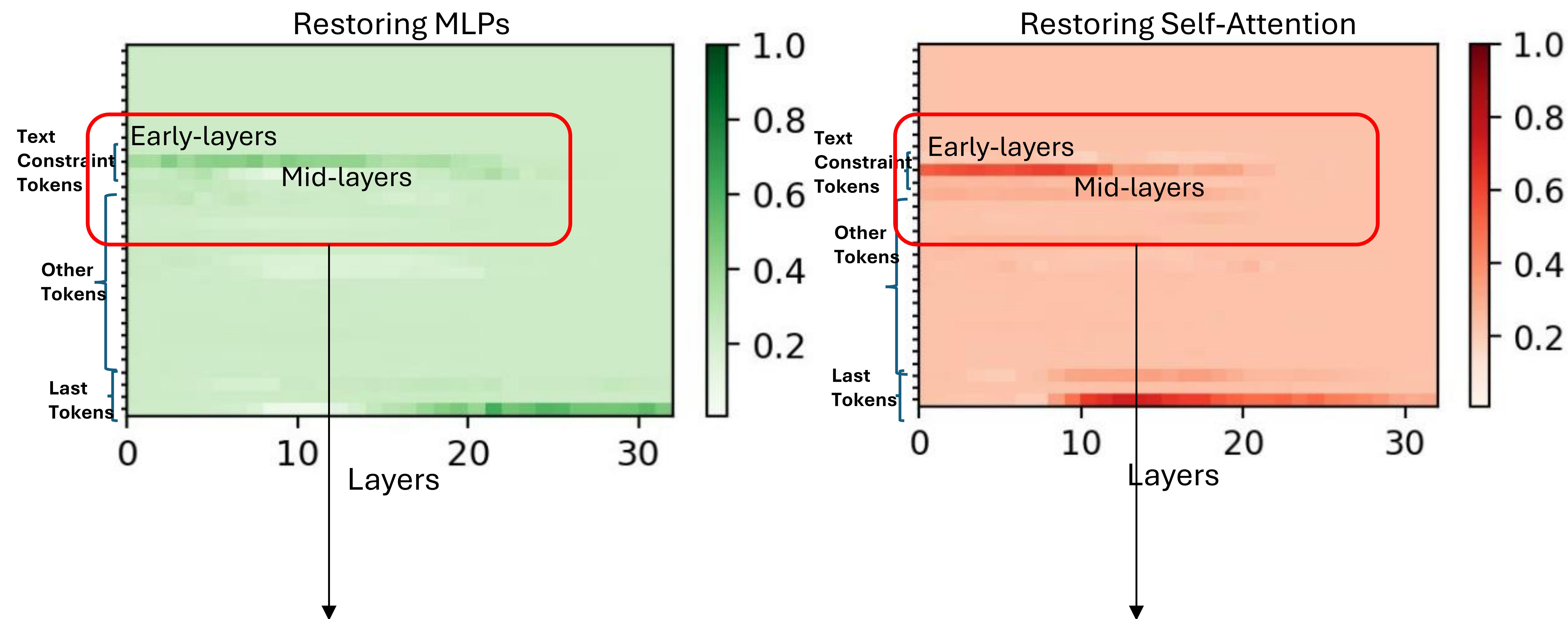
Layers 4-7

Requires restoring a window of layers

**Takeaway:** MLLMs (under the presence of a visual prompt) retrieve information differently than LLMs – although the same language backbone is used



# Information Processing for Text-Constraint in Multi-Constraint Questions



Text-constraint in multi-constraint questions require information to be retrieved from early + mid layers with a large window size

**Takeaway :** Multi-constraint VQA questions require more parametric memory to answer a question

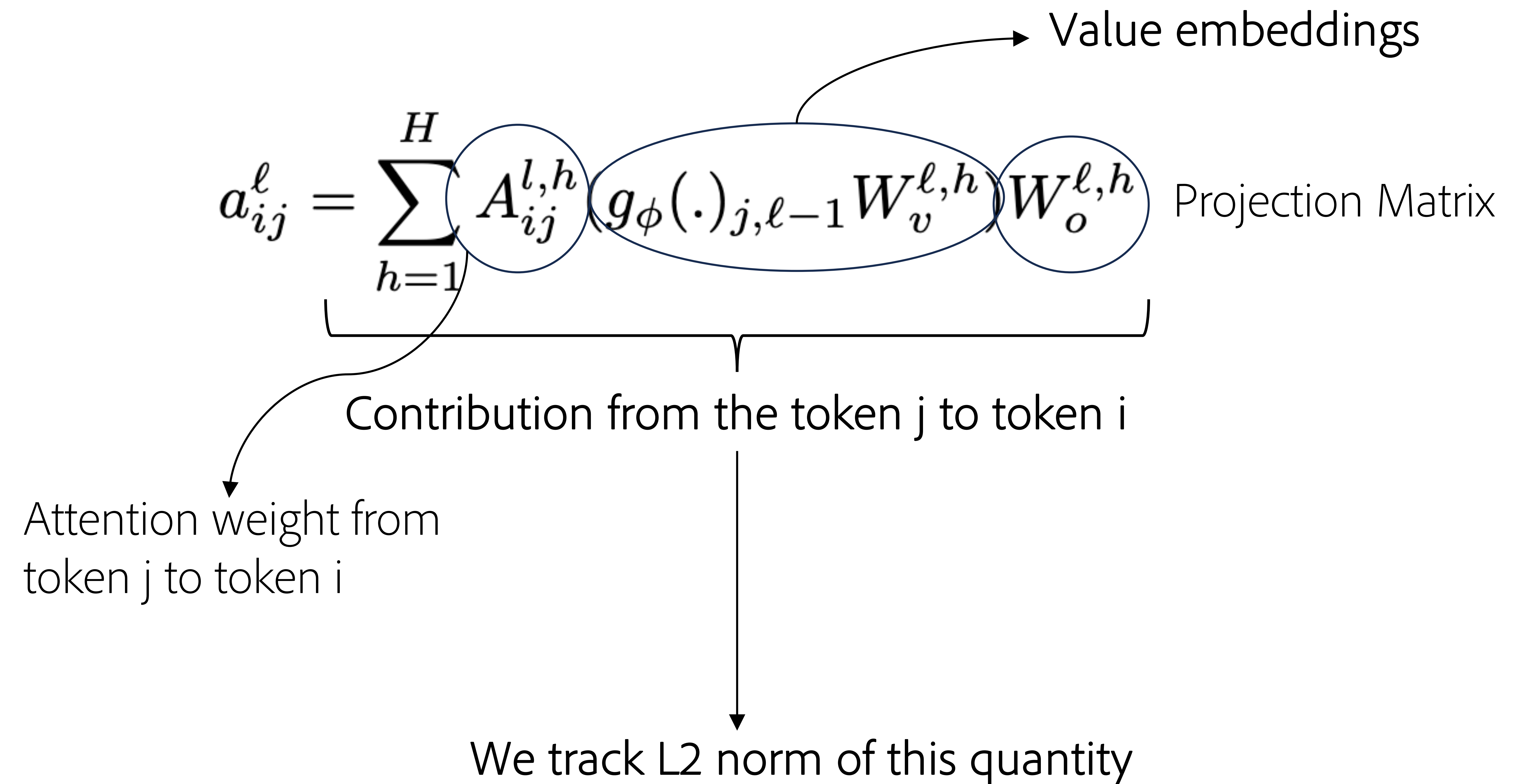
MultimodalCausalTrace:  
Identifying Causal Layers  
in MLLMs

Multi-Constraint  
Questions

# Tracking Attention Contributions from Visual Tokens to Constraints

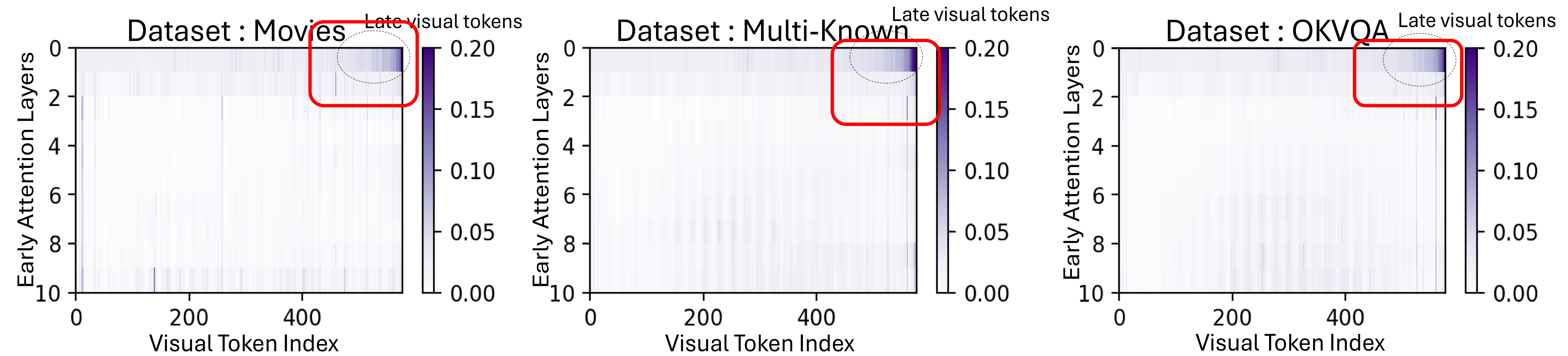
Information Transfer in MLLMs

Visual Tokens to Constraint



Attention contributions is a general property of transformers and have been previously used in Yuksekgonul et al.(2023)

We compute the attention contributions from the visual tokens to the constraint token



**Takeaway :** Only a subset of visual tokens (after the projection layer) transfer information to the constraint token position

Most of the visual information gets accumulated in the late visual tokens by the projection layer

Insights from  
Information  
Transfer in MLLMs

Visual Tokens to  
Constraint

Can we use the interpretability insights  
towards a tangible application?

Can we use the interpretability insights towards a tangible application?

- Introducing Rare Concepts
- Fixing Failure Modes

# Qualitative Examples of Rare VQA Questions



Inchcolm

In which country  
is *this island*  
located?



Croome Court

In which country  
is *this house*  
located?



Hilton Prague

In which country  
is *this hotel*  
located?



Cragside

In which country  
is *this house*  
located?



Giannutri

In which country  
is *this island*  
located?



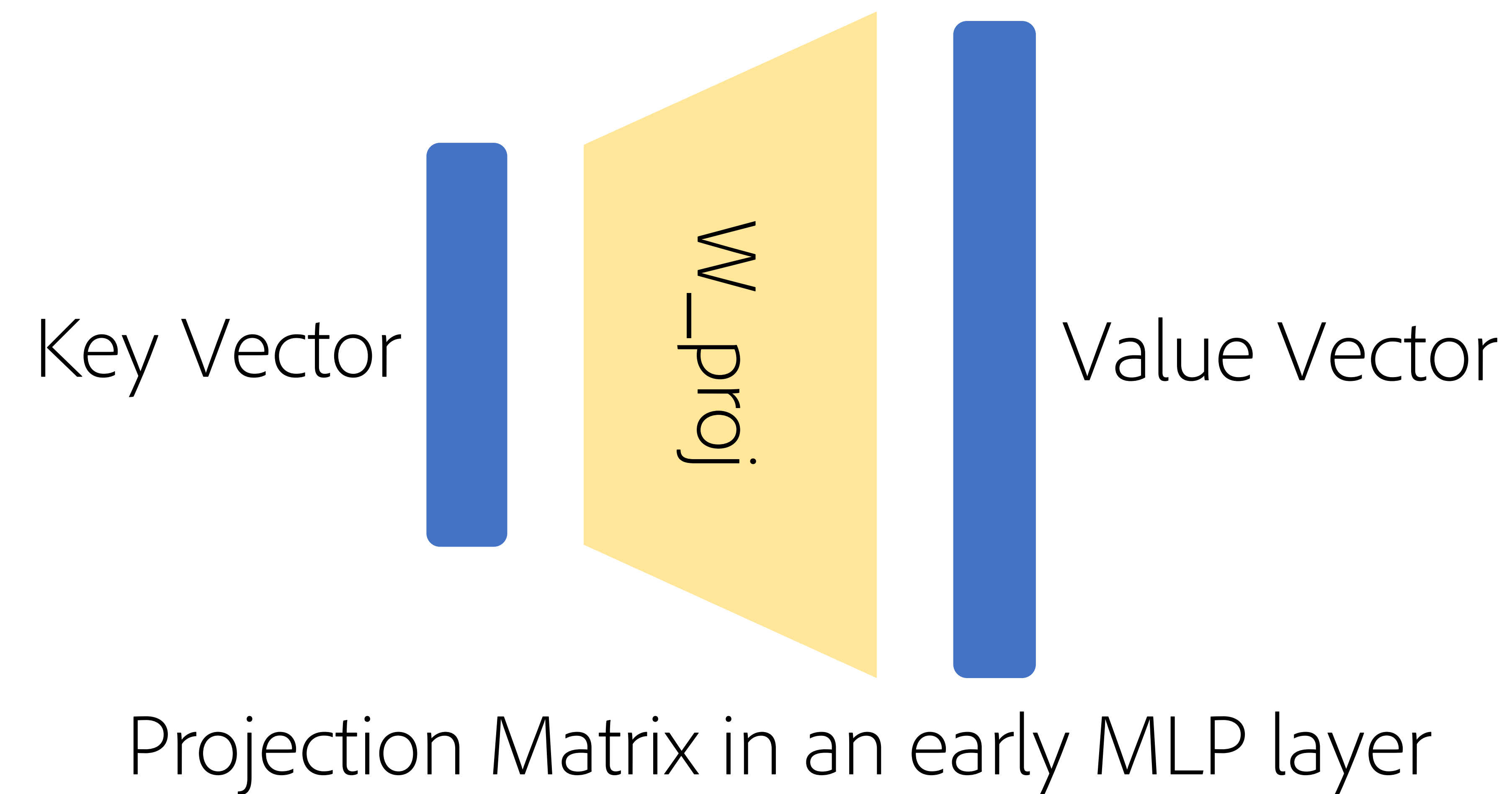
Kandy Lake

In which country  
is *this lake*  
located?

# Practical Application Leveraging Interpretability Insights:

## Model Editing to Incorporate Rare Knowledge into MLLMs - Method

We edit the *early MLP layers* for the best editing performance



Step 1: Obtain key embeddings via forward pass of the visual prompt + question

Step 2: Obtain value embeddings

Value embeddings  $z_{c,l}^*$  =  $\arg \min_{z_{c,l}} \mathcal{L}(z_{c,l})$  Language Modeling AR Loss

Step 3: Edit Optimization

Key embeddings

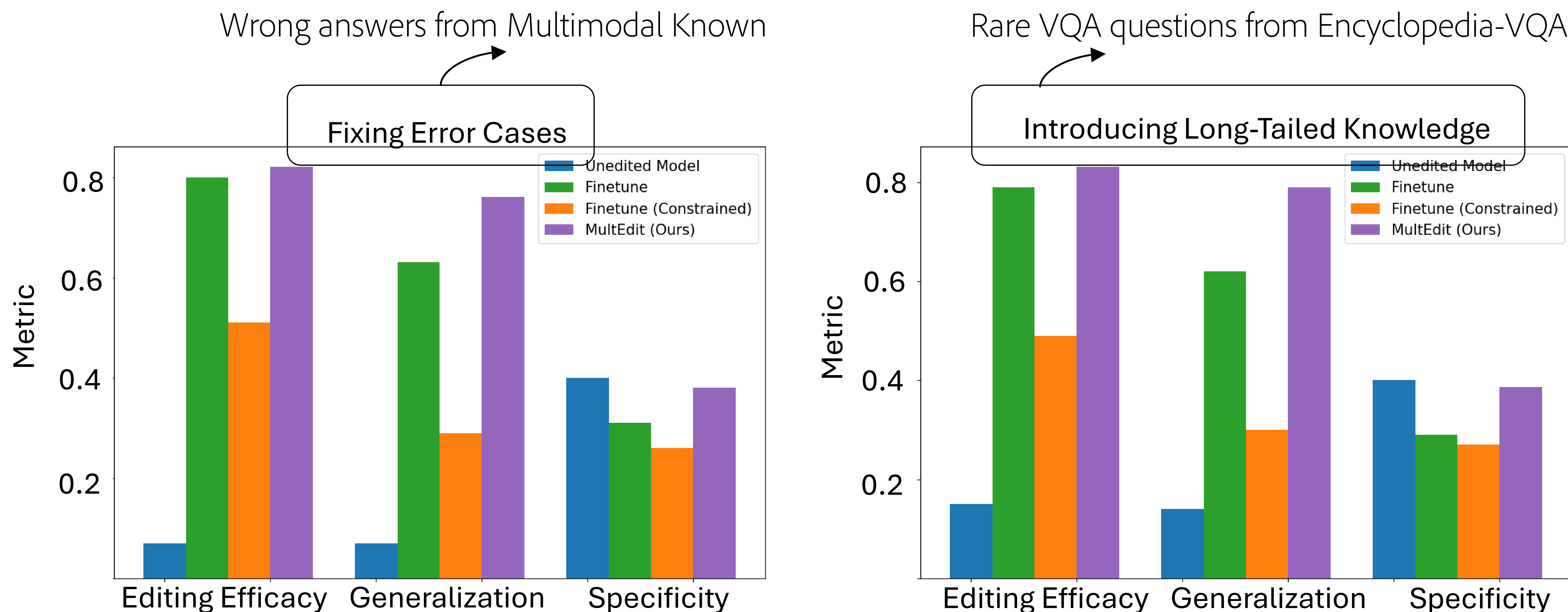
$$W_{proj}^{\ell*} = \arg \min_{W_{proj}^{\ell}} \|W_{proj}^{\ell} k_{c,l} - z_{c,l}^*\|_2^2 + \lambda \|W_{proj}^{\ell} - W_{proj}^{\ell'}\|_2^2$$

Does not require caching a Multimodal Wikipedia matrix → Relatively data free

# Practical Application Leveraging Interpretability Insights:

## Model Editing to Incorporate Rare Knowledge into MLLMs - Results

Editing the early layers with our objective leads to incorporation of rare knowledge into the model and is better (+faster) than fine-tuning the language model



### Metrics:

- (i) Editing Efficacy: Effectiveness of the edit (Measured by the probability of the correct answer)
- (ii) Generalization Efficacy: Effectiveness of the edit under paraphrased questions (Measured by the probability of correct answer)
- (iii) Specificity: Editing effect on unrelated VQA questions (measured by VQA-accuracy on OK-VQA and eval on MMMU)



## Conclusion

Checkout out other interpretability works at:  
<https://samyadeepbasu.github.io/>

