

Inference and Verbalization Functions During In-Context Learning

Junyi Tao, Xiaoyin Chen, Nelson F. Liu

Stanford University Mila



TL;DR

Why is LM performance unaffected by **label remappings** during **In-Context Learning**?

e.g., "true"/"false" to "cat"/"dog"

Using interchange intervention, we find:

- 1) LMs use an **inference** function (that solves the task) and a **verbalization** function (that maps the inferred answer to the label space)
- 2) the **inference** function is *invariant to remappings of the label space*
- 3) these functions are *located consistently and separately* in layers across tasks, datasets, and models (7B-70B)

ICL Setting

DEMO Sentence 1: {Premise}\nSentence 2: {Hypothesis} \n {Label}
Sentence 1: {Premise}\nSentence 2: {Hypothesis} \n {Label}

TEST Sentence 1: {Premise}\nSentence 2: {Hypothesis} \n

Input-label relation

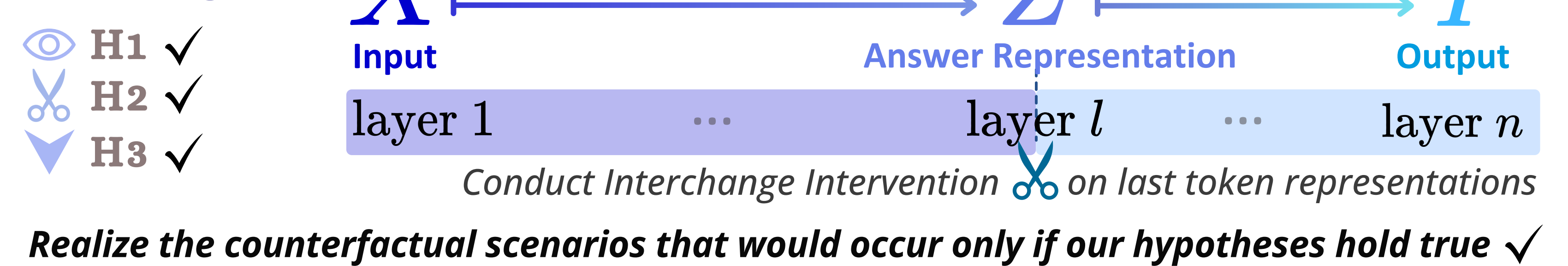
{Label} = true/false $\xrightarrow{\text{remap the label space}}$ {Label} = cat/dog

Hypotheses

- H1 Causal Mechanism** **Inference** Infers the representation of the answer to the task
- H2 Invariance** **Verbalization** Maps the answer representation to the label words
- H3 Localization** One shared **Inference** for the same task despite diff. label words
- The functions are separately and consistently localized in the model

Method Apply **Interchange Intervention** on representations induced by diff. ICL settings

Findings

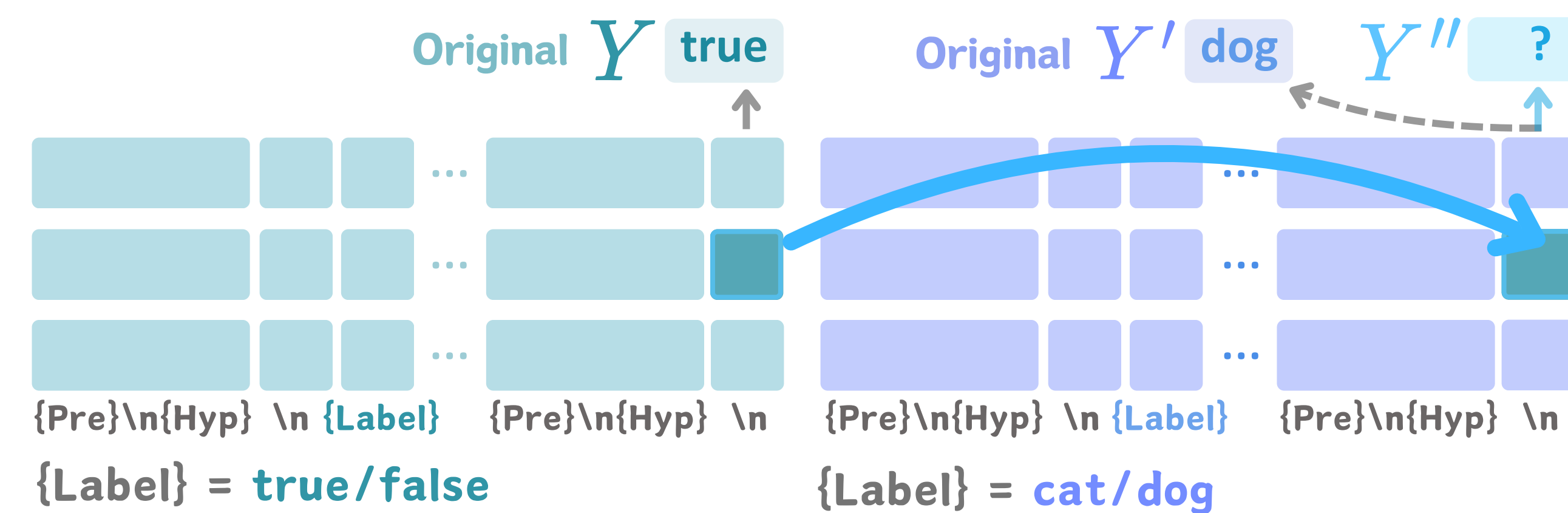


Main Experiments: Remap Label Spaces

1. Remap label spaces and preserve example-label relations

{Label}=true/false $\xrightarrow{\text{remap the label space}}$ {Label}=cat/dog

2. Intervene: Interchange representations across settings



3. Evaluate: Is the output flipped after the intervention?

Y'' flips from dog to **Y_c** = cat consistently at certain layers
i.e. The intervention flips the answer & preserves the label space

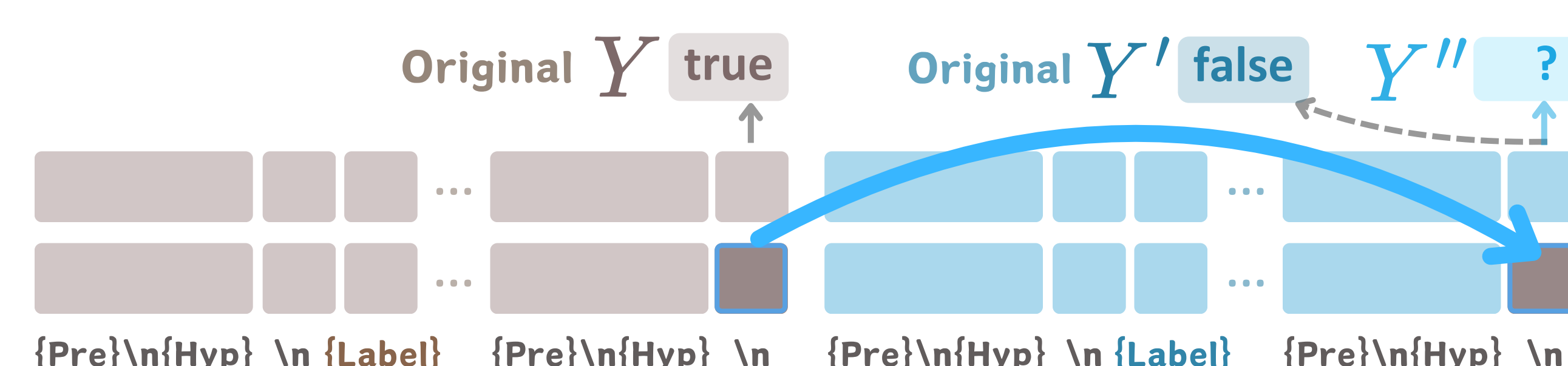
- These layers perform the *causal* role of **verbalization** \rightarrow **H1.2** ✓
- Inference** is *invariant* to remappings of label spaces \rightarrow **H2** ✓
- Verbalization** is *consistently* located in late layers \rightarrow **H3** ✓

Complementary Experiments: Reconstruct Tasks

1. Construct alternative tasks w/ example-label relations

Task = Topic classification {Label} = true/false Task = NLI {Label} = true/false

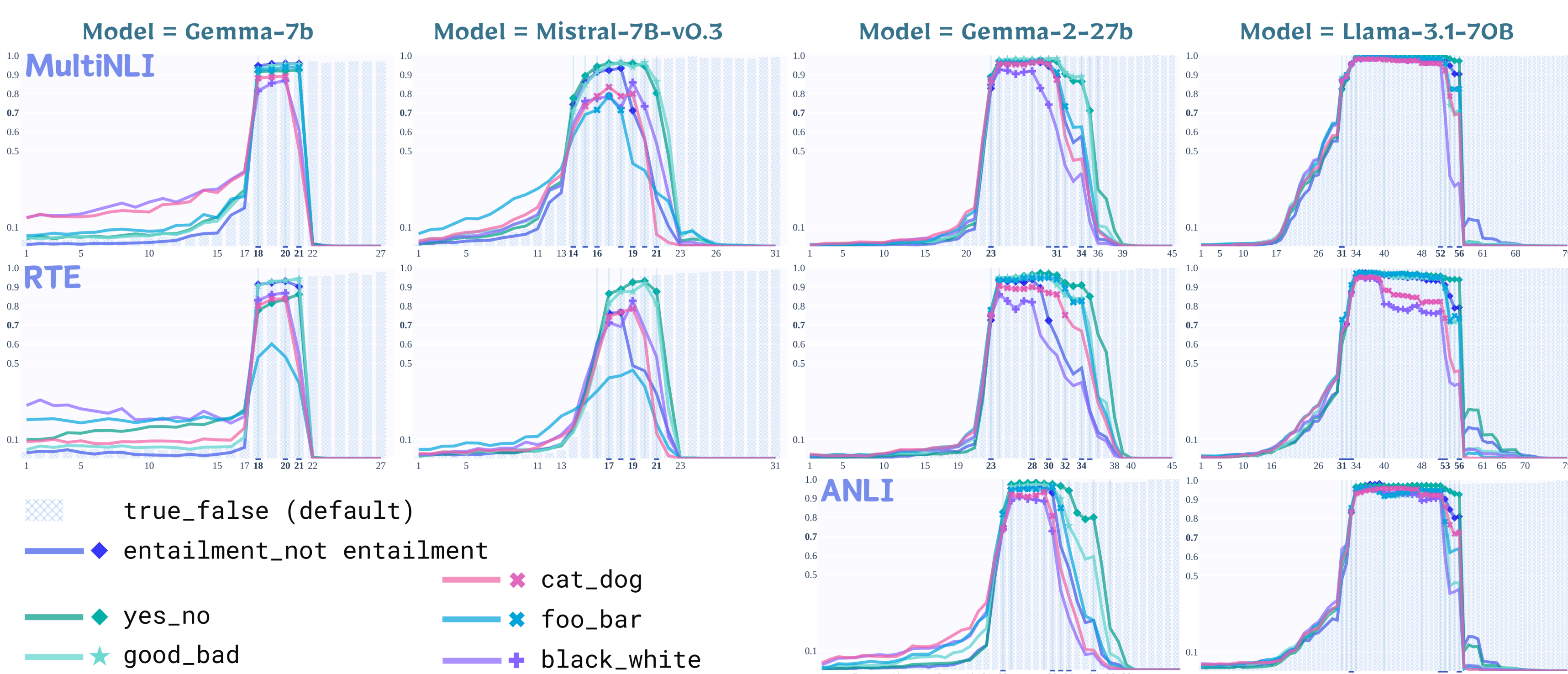
2. Intervene: Interchange representations across settings



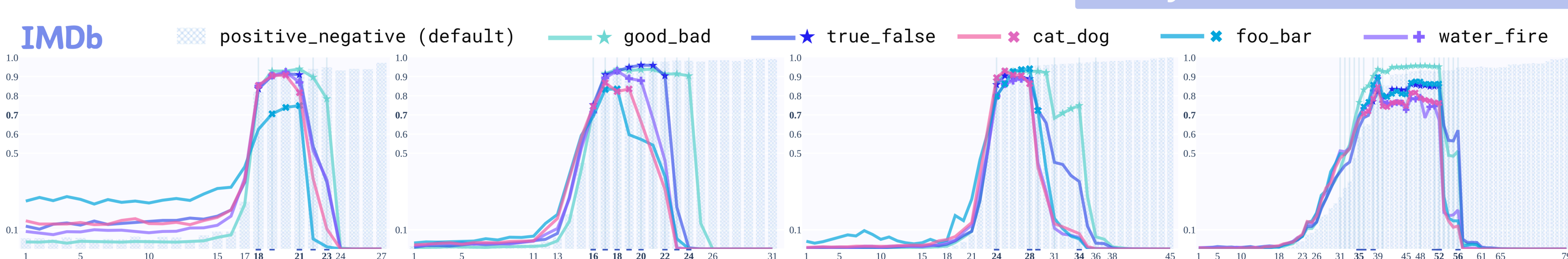
3. Evaluate: Is the output flipped after the intervention?

Y'' flips from false to **Y_c** = true consistently at certain layers
i.e. The intervention flips the answer & preserves the label space

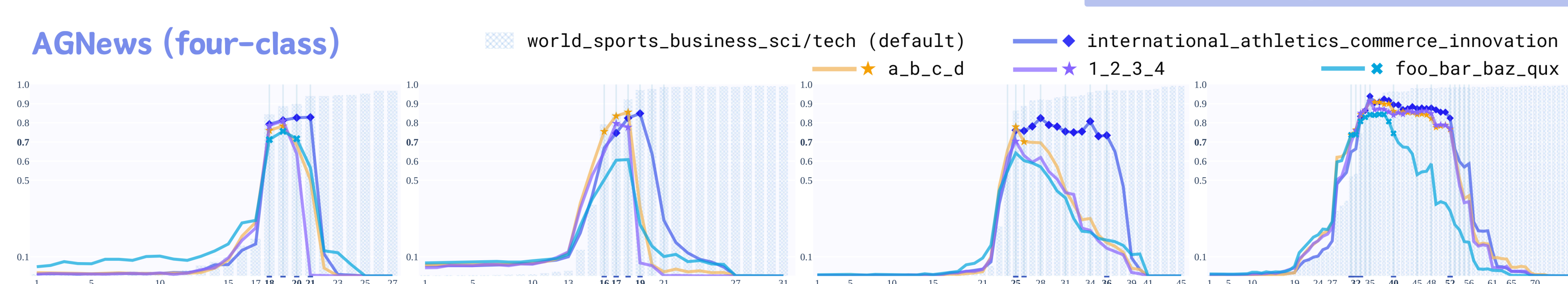
Main Results: Intervention Patterns are Consistent Across Tasks



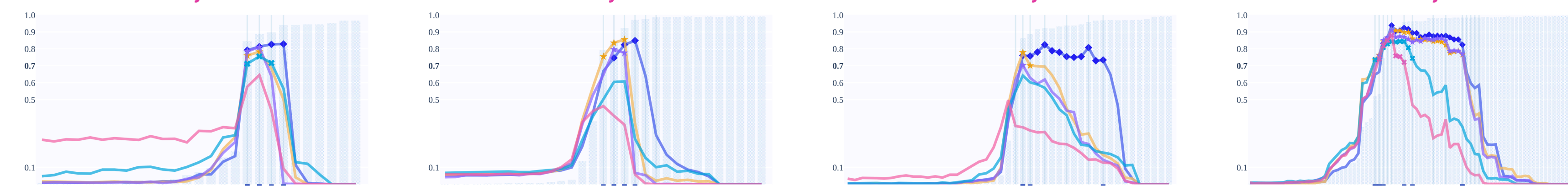
Binary Sentiment Classification



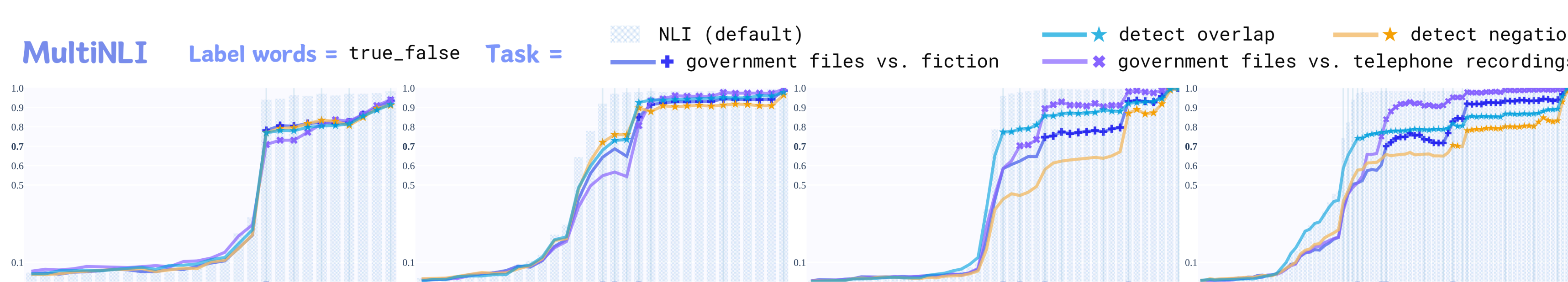
Multi-Class Topic Classification



Intervention on AGNews Label Words = cat_dog_deer_bird Illustration of Failure Cases



Main Results: Four Reconstructed Tasks on MultiNLI



- These layers perform the *causal* role of **inference** \rightarrow **H1.2** ✓