

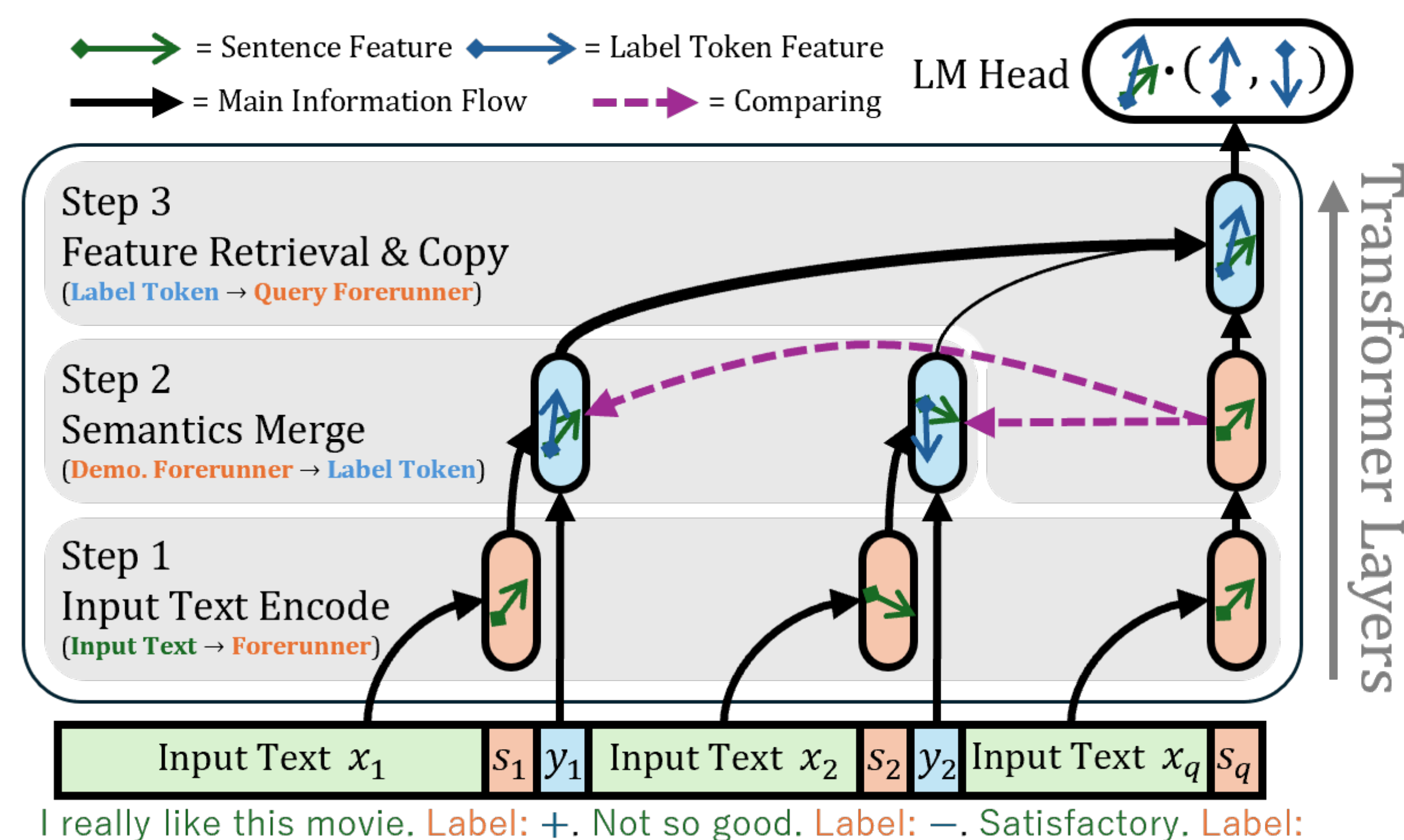


Background

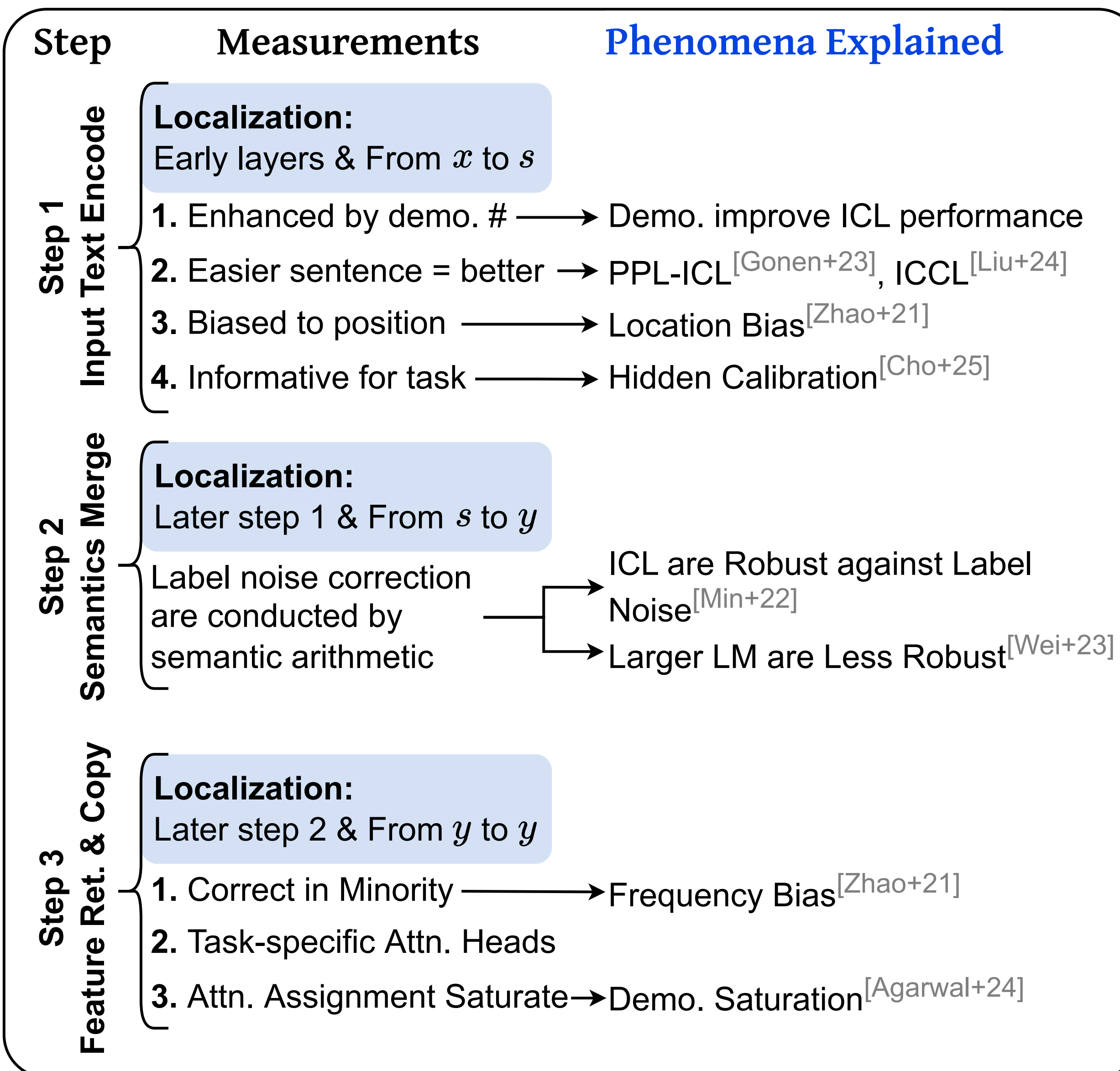


Some scattered phenomena have been observed, but no unified explanation.

Abstract

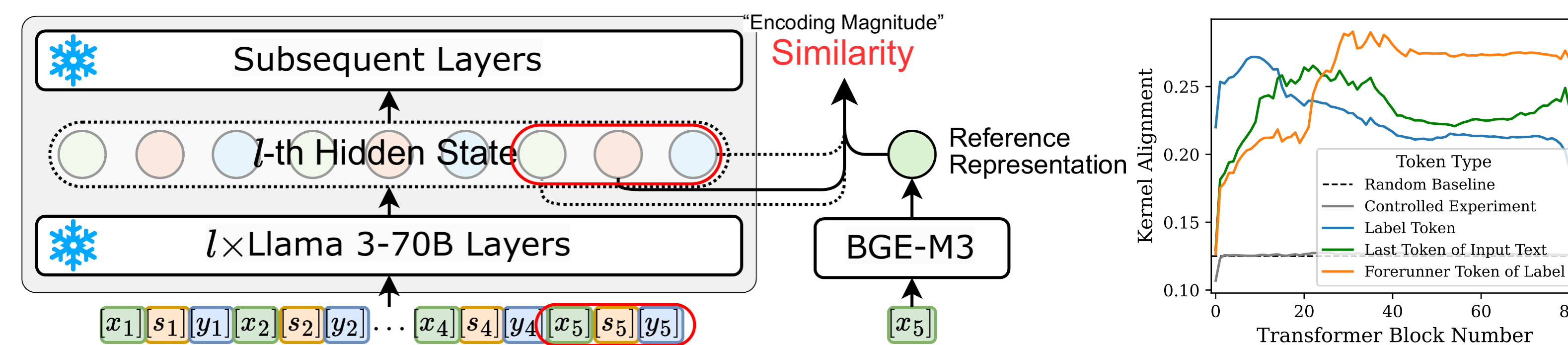


To unify previous observation:

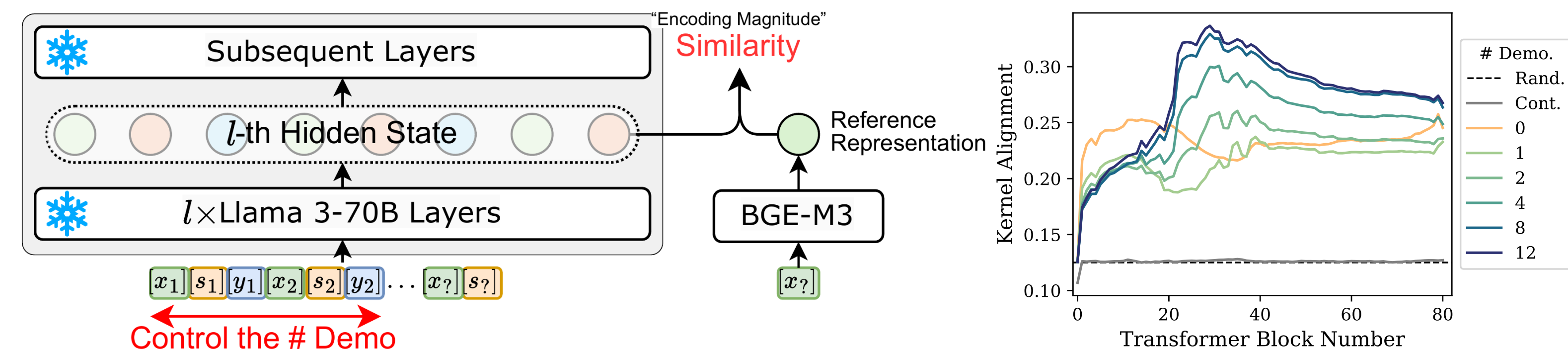


This work is supported by The Nakajima Foundation Start-Up Support Grant. The authors would like to thank Mr. Lihao Liu at the Beijing Institute of Technology, and Ms. Yunpin Li at the Capital Normal University for proofreading.

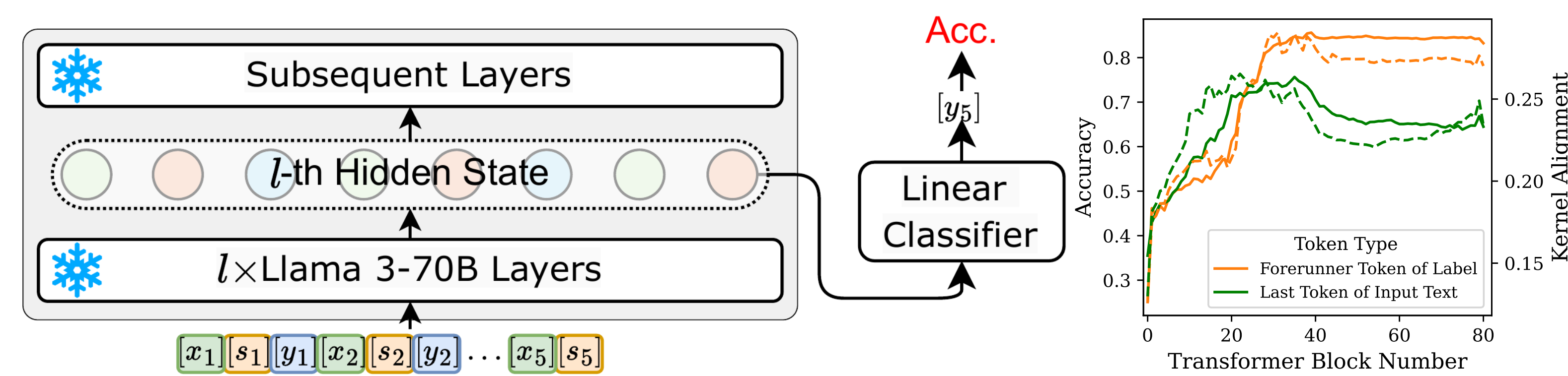
Step 1: Input Text Encode

1.1. Localization: Which Token and Layers? → Early layers & From x to s 

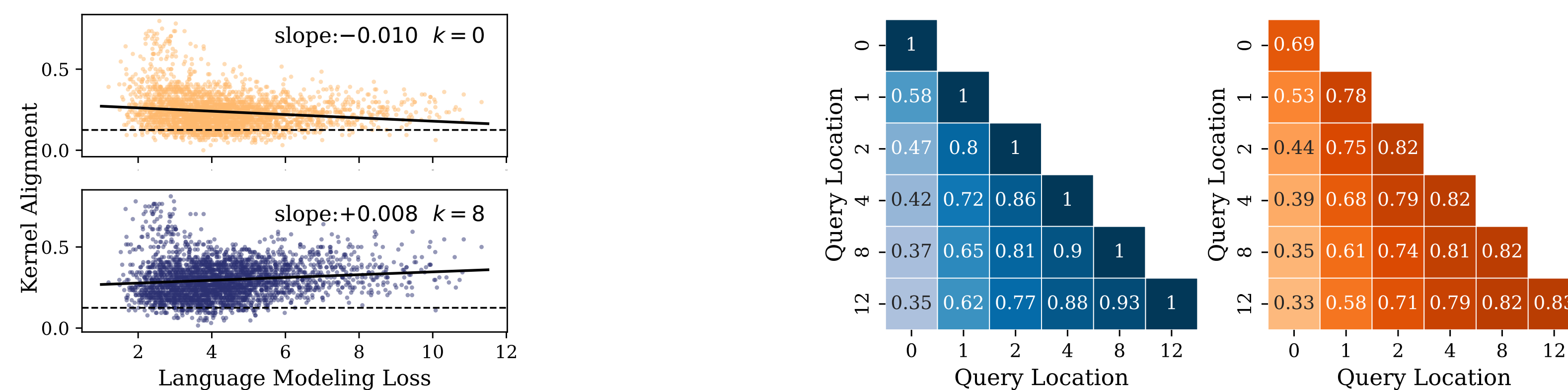
1.2. Encoding Enhanced by Demonstrations



1.3. Encoding is Informative to Downstream Task

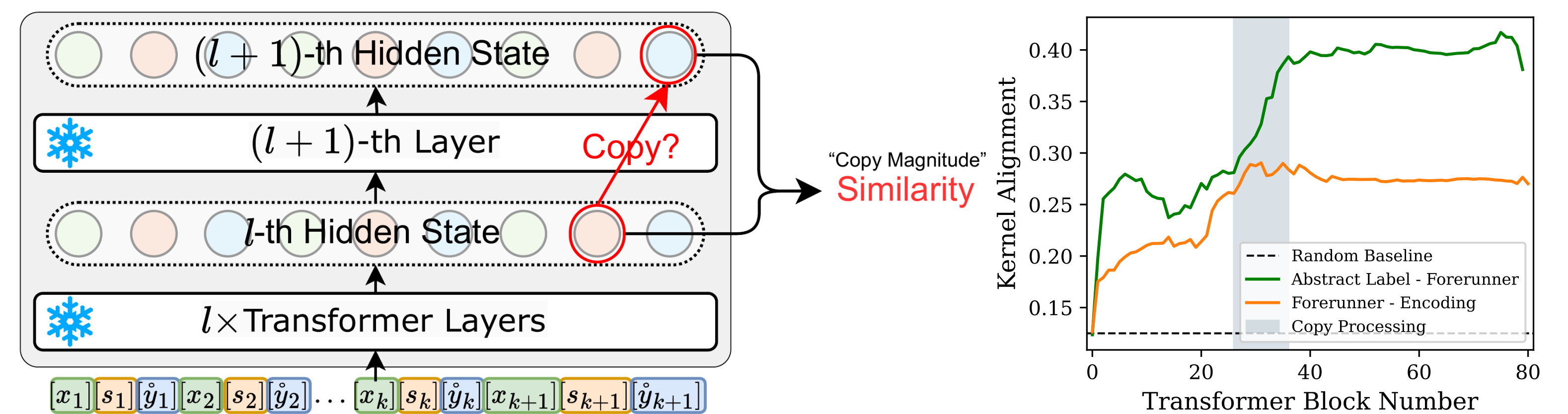
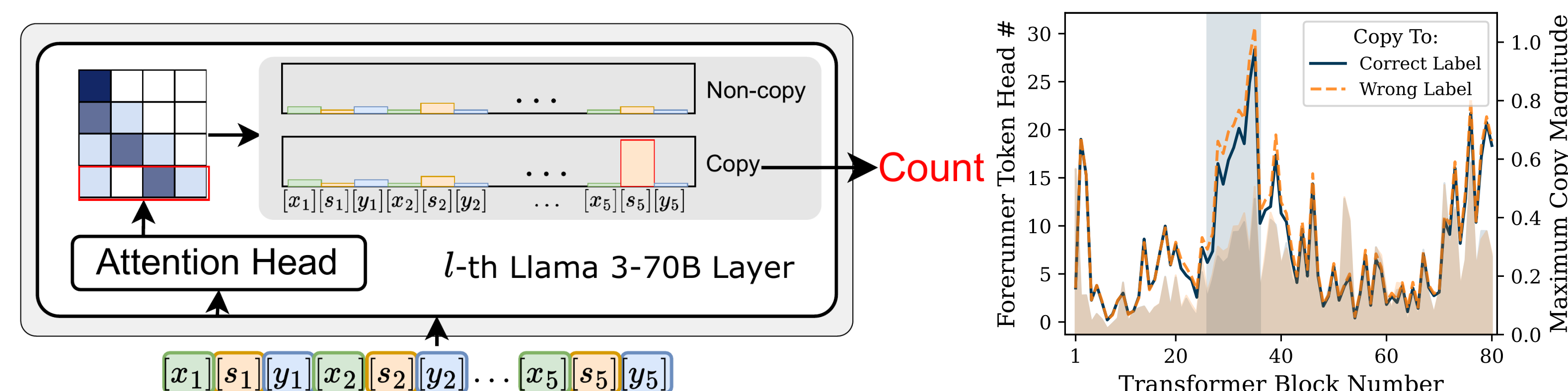


1.4. Encoding Weakened by CLM Loss; More Demonstrations Restore Them

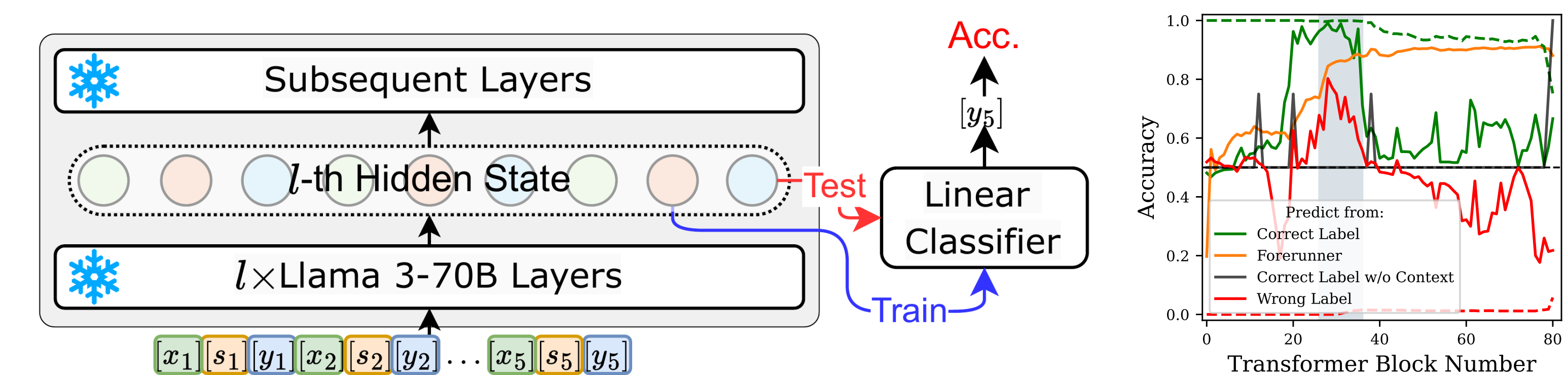


1.5. Encoding is Biased towards Demonstration Index ↑

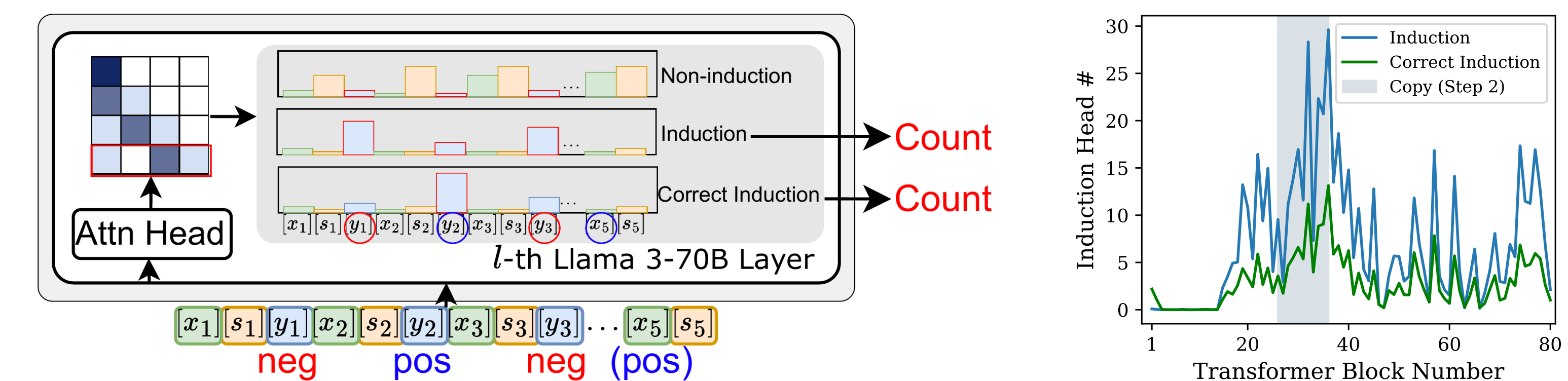
Step 2: Semantics Merge

2.1. Localization: Which Layers? → Later than Step 1 & From s to y 2.2. Why Robust to Noise? Consistent Copy Magnitude on Correct / Wrong y 

(Continue Step 2: Semantics Merge)

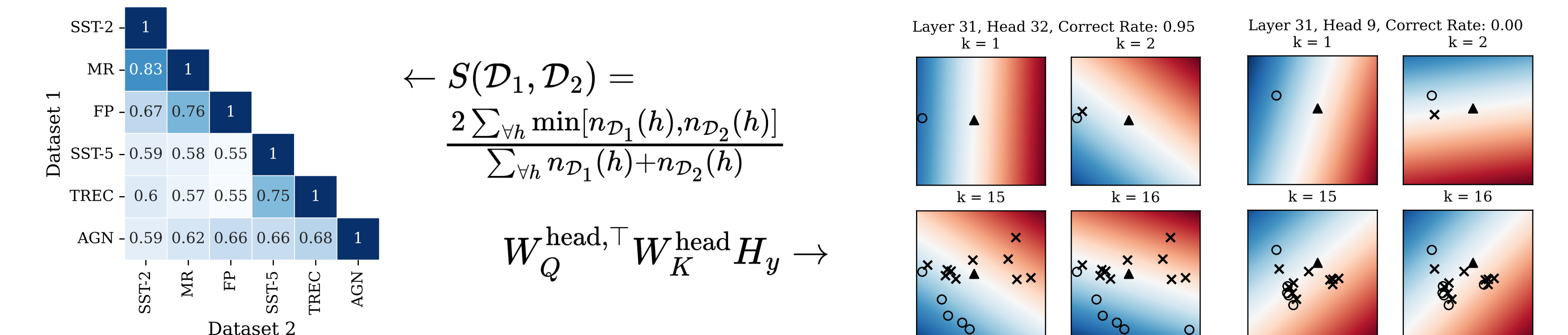
2.3. Why Robust to Noise? Feature Enhance / Neutralize on Correct / Wrong y 

Step 3: Feature Retrieval and Copy (Induction Heads)

3.1. Localization: Which Layers? → Later than Step 2 & From y to s 

3.2. Correct Retrieval is in Minority ↑

3.3. Some Induction Heads are Intrinsic; Some are Evoked by Task



Ablation Study

The proposed 3 steps are essential

#	Attention Disconnected	Affected Layers Ratio (from layer 1)			
	Key \rightarrow Query	25%	50%	75%	100%
1	None (4-shot baseline)	± 0 (Acc. 68.55)			
- Step1: Input Text Encode -					
2	Demo. Texts $x_i \rightarrow$ Forerunner s_i	-4.98 -0.80 ± 0.00	-15.82 -1.10 ± 0.02	-23.43 -3.20 ± 1.87	-30.60 -1.01 ± 0.01
3	Query Texts $x_q \rightarrow$ Forerunner s_q	-13.87 -0.16 ± 0.00	-21.10 -0.08 ± 0.00	-24.74 -0.47 ± 0.04	-28.38 -0.55 ± 0.00
- Step2: Semantics Merge -					
4	Demo. Forerunner $s_i \rightarrow$ Label y_i	-2.24 -0.00 ± 0.00	-3.45 -0.18 ± 0.00	-3.39 -0.10 ± 0.04	-3.42 -0.18 ± 0.01
- Step3: Feature Retrieval & Copy -					
5	Label $y_i \rightarrow$ Query Forerunner s_q	-5.14 $+0.03 \pm 0.00$	-10.03 -0.08 ± 0.00	-11.36 $+0.00 \pm 0.00$	-10.22 -0.06 ± 0.00
Reference Value					
6	Zero-shot	-17.90 (Acc. 50.65)			
7	Random Prediction	-36.05 (Acc. 32.50)			

Limitations: Is "Copying from Context" Enough for ICL?

Induction-based ICL Explanation

Taylor Swift → Singer
Donald Trump → Politician
Ian Goodfellow → Researcher
Geoffrey Hinton → Researcher

1. Search Similar Representations
2. Copy the Co-responding Label

Issue: Generalization to Unseen Labels

Taylor Swift → Singer
Donald Trump → Politician
Geoffrey Hinton → Researcher

[Gonen+23] Demystifying prompts in language models via perplexity estimation.
[Liu+24] Let's learn step by step: Enhancing in-context learning ability with curriculum learning.
[Zhao+21] Calibrate before use: Improving few-shot performance of language models.
[Cho+25] Token-based Decision Criteria Are Suboptimal in In-context Learning.
[Min+22] Rethinking the role of demonstrations: What makes in-context learning work?
[Wei+23] Larger language models do in-context learning differently.
[Agarwal+23] Many-shot in-context learning.