

Label Words are Anchors



Meizhi Zhong
12.21

Label Words are Anchors

Label Words are Anchors: An Information Flow Perspective for Understanding In-Context Learning

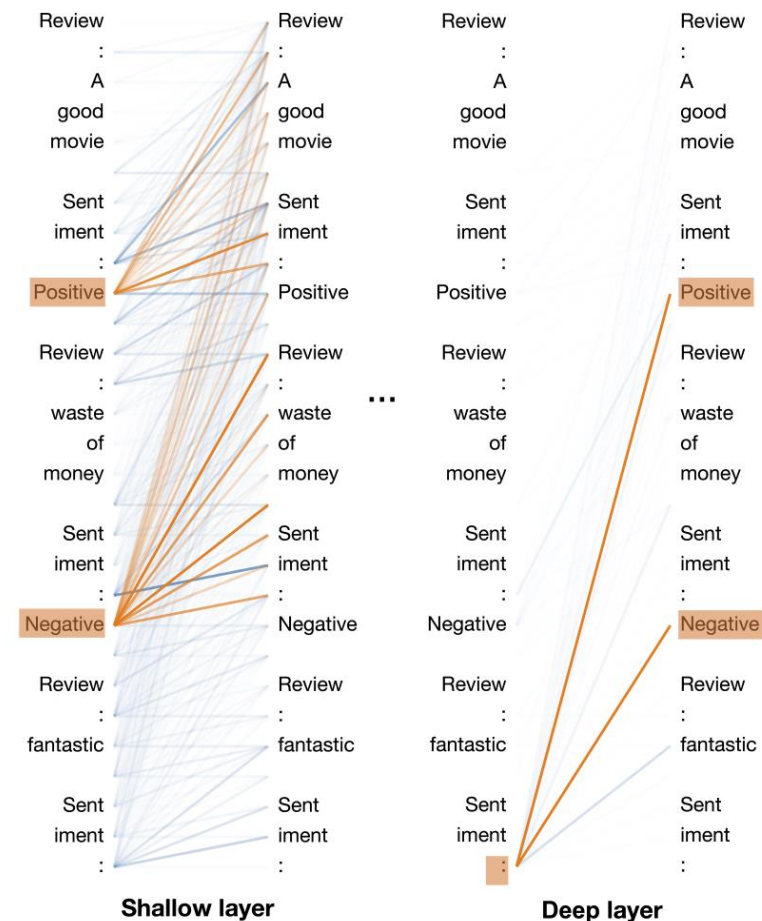
Lean Wang^{†,§}, Lei Li[†], Damai Dai[†], Deli Chen[§],
Hao Zhou[§], Fandong Meng[§], Jie Zhou[§], Xu Sun[†]

[†]National Key Laboratory for Multimedia Information Processing,
School of Computer Science, Peking University

[§]Pattern Recognition Center, WeChat AI, Tencent Inc., China

{lean, daidamai, xusun}@pku.edu.cn nlp.lilei@gmail.com

victorcheng@deepseek.com {tuxzhou, fandongmeng, withtomzhou}@tencent.com



Motivation

How LLMs learn from the provided context?

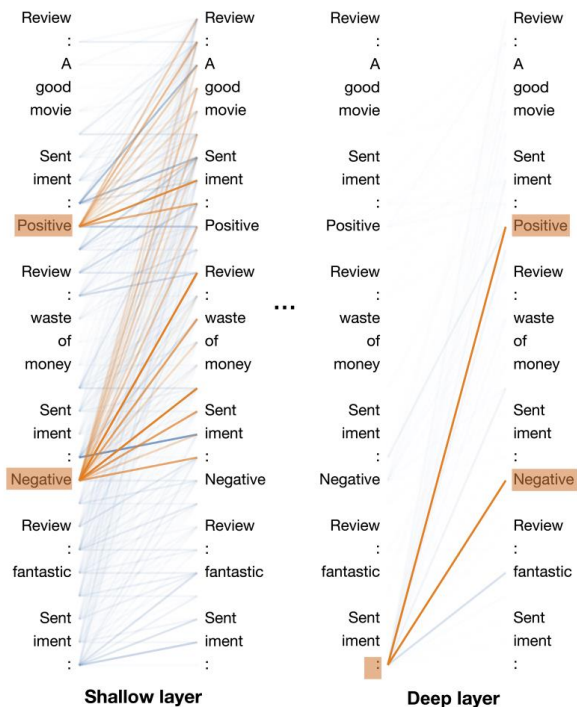


Figure 1: Visualization of the information flow in a GPT model performing ICL. The line depth reflects the significance of the information flow from the right word to the left. The flows involving label words are highlighted. Label words gather information from demonstrations in shallow layers, which is then extracted in deep layers for final prediction.

Information Flow with Labels as Anchors

\mathcal{H}_1 : In shallow layers, label words gather the information of demonstrations to form semantic representations for deeper layers.

\mathcal{H}_2 : In deep layers, the model extracts the information from label words to form the final prediction.

Significance of Information

$$I_l = \left| \sum_h A_{h,l} \odot \frac{\partial \mathcal{L}(x)}{\partial A_{h,l}} \right|.$$

$I_l(i, j)$ represents the **significance** of the information flow from the **j-th** word to the **i-th** word for ICL.

S_{wp} , the mean significance of information flow from the text part to label words:

$$S_{wp} = \frac{\sum_{(i,j) \in C_{wp}} I_l(i, j)}{|C_{wp}|}, \quad (2)$$
$$C_{wp} = \{(p_k, j) : k \in [1, C], j < p_k\}.$$

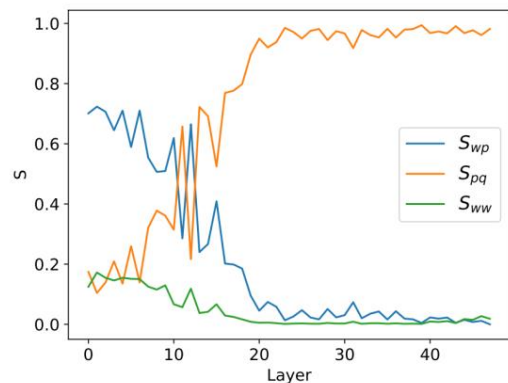
S_{pq} , the mean significance of information flow from label words to the target position:

$$S_{pq} = \frac{\sum_{(i,j) \in C_{pq}} I_l(i, j)}{|C_{pq}|}, \quad (3)$$
$$C_{pq} = \{(q, p_k) : k \in [1, C]\}.$$

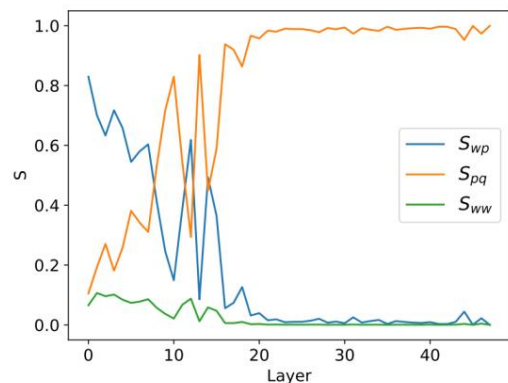
S_{ww} , the mean significance of the information flow amongst all words, excluding influences represented by S_{wp} and S_{pq} :

$$S_{ww} = \frac{\sum_{(i,j) \in C_{ww}} I_l(i, j)}{|C_{ww}|}, \quad (4)$$
$$C_{ww} = \{(i, j) : j < i\} - C_{wp} - C_{pq}.$$

Label Words are Anchors: Hypothesis Motivated by Saliency Scores



(a) Results on the SST-2 dataset



(b) Results on the AGNews dataset

Figure 3: Relative sizes of S_{wp} , S_{pq} , and S_{ww} in different layers on SST-2 and AGNews. Results of other datasets can be found in Appendix B. Initially, S_{wp} occupies a significant proportion, but it gradually decays over layers, while S_{pq} becomes the dominant one.

Review: I dislike ... Sentiment: **Negative** Review: A good ... Sentiment: **Positive** Review: ... Sentiment: **Positive**

In shallow layers

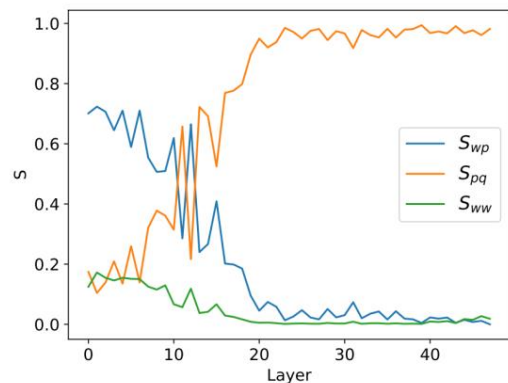
- (1) $S_{\{pq\}}$, the significance of the information flow from **label** words to **targeted positions** is **low**;
- (2) $S_{\{wp\}}$, the information flow from the **text** part to **label** words is **high**;

In deep layers

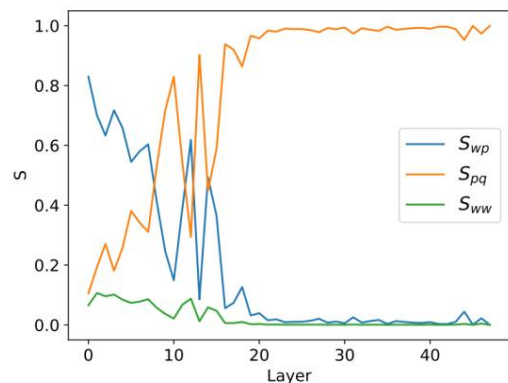
- (1) $S_{\{pq\}}$, the significance of the information flow from **label** words to **targeted positions** is **dominant**;
- (2) $S_{\{wp\}}$, the information flow from the **text** part to **label** words is **low**;

$S_{\{pq\}}$ and $S_{\{wp\}}$ usually surpass $S_{\{ww\}}$, suggesting that interactions involving **label words** outweigh others.

Label Words are Anchors: Hypothesis Motivated by Saliency Scores



(a) Results on the SST-2 dataset



(b) Results on the AGNews dataset

Figure 3: Relative sizes of S_{wp} , S_{pq} , and S_{ww} in different layers on SST-2 and AGNews. Results of other datasets can be found in Appendix B. Initially, S_{wp} occupies a significant proportion, but it gradually decays over layers, while S_{pq} becomes the dominant one.

Review: I dislike ... Sentiment: **Negative** Review: A good ... Sentiment: **Positive** Review: ... Sentiment: **Positive**

In shallow layers

- (1) $S_{\{pq\}}$, the significance of the information flow from **label** words to **targeted positions** is **low**;
- (2) $S_{\{wp\}}$, the information flow from the **text** part to **label** words is **high**;

In deep layers

- (1) $S_{\{pq\}}$, the significance of the information flow from **label** words to **targeted positions** is **dominant**;
- (2) $S_{\{wp\}}$, the information flow from the **text** part to **label** words is **low**;

$S_{\{pq\}}$ and $S_{\{wp\}}$ usually surpass $S_{\{ww\}}$, suggesting that interactions involving **label words** outweigh others.

Shallow Layers: Information Aggregation

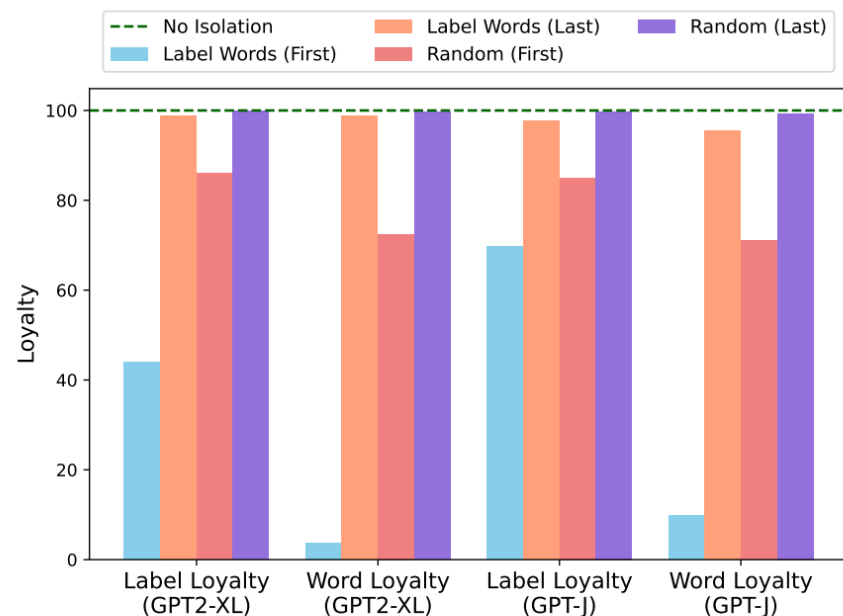


Figure 4: The impact of isolating label words versus randomly isolating non-label words within the first or last 5 layers. Isolating label words within the first 5 layers exerts the most substantial impact, highlighting the importance of shallow-layer information aggregation via label words.

cally, we set $A_l(p, i) (i < p)$ to 0 in the attention matrix A_l of the l -th layer, where p represents label words and i represents preceding words.

- (1) Label Loyalty: measures the consistency of output **labels** with and without **isolation**.
- (2) Word Loyalty: employs the **Jaccard similarity** to compare the **top-5 predicted words** /w and /wo **isolation**, capturing more subtle model output alterations

Shallow Layers: Information Aggregation

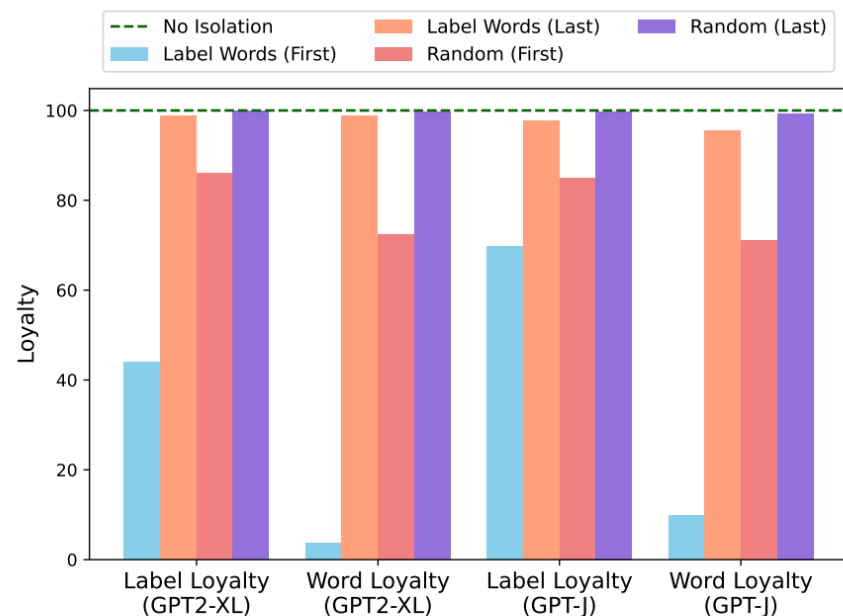


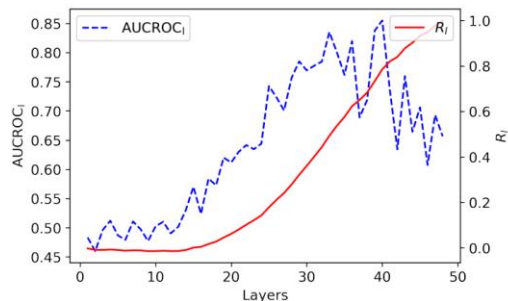
Figure 4: The impact of isolating label words versus randomly isolating non-label words within the first or last 5 layers. Isolating label words within the first 5 layers exerts the most substantial impact, highlighting the importance of shallow-layer information aggregation via label words.

cally, we set $A_l(p, i) (i < p)$ to 0 in the attention matrix A_l of the l -th layer, where p represents label words and i represents preceding words.

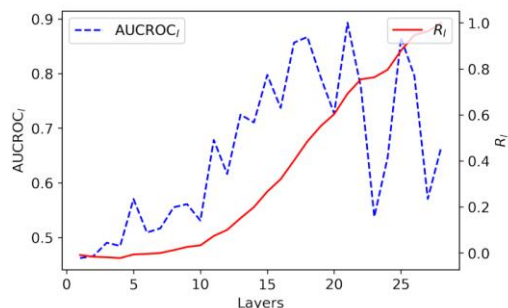
- (1) Label Loyalty: measures the consistency of output **labels** with and without **isolation**.
 - (2) Word Loyalty: employs the **Jaccard similarity** to compare the **top-5 predicted words** /w and /wo **isolation**, capturing more subtle model output alterations
- **notable influence** on the model's behavior when label words are isolated within the first 5 layers
 - **Influence becomes inconsequential** within the last 5 layers.
 - Emphasizes the **superiority** of **label** words over non-label words.

Loyalty

Deep Layers: Information Extraction



(a) GPT2-XL (total 48 layers).



(b) GPT-J (total 28 layers).

Figure 5: $AUCROC_l$ and R_l of each layer in GPT models. The result is averaged over SST-2, TREC, AGNews, and Emoc. $AUCROC_l$ reaches 0.8 in deep layers, and R_l increases mainly in the middle and later layers.

- AUC-ROC score to quantify the **correlation** between $A_{\mathcal{L}}(q, p_i)$ and model **prediction**.
- $R_{\mathcal{L}}$ quantify the **accumulated contribution** of the first \mathcal{L} layers to model **prediction**

$$R_l = \frac{\sum_{i=1}^l (AUCROC_i - 0.5)}{\sum_{i=1}^N (AUCROC_i - 0.5)}.$$

- AUC~ROC得分越接近1.0，表示模型性能越好，能够更好地区分正例和负例。如果AUC~ROC得分接近0.5，则模型性能较差，类似于随机猜测。

Hypothesis

- **Shallow layers:** assemble information from demonstrations via label words to form semantic representations.
- **Deep layers:** the aforementioned aggregated information on label words is then extracted to form the final prediction.
- Label words are Anchors

Anchor Re-weighting: Method(performance)

a strong correlation between the model's output category and the attention distribution

$$(A(q, p_1), \dots, A(q, p_C))$$

$$\begin{aligned} & \Pr_f(Y = i | X = x) \\ & \approx A(q, p_i) \\ & = \frac{\exp(\mathbf{q}_q \mathbf{k}_{p_i}^T / \sqrt{d})}{\sum_{j=1}^N \exp(\mathbf{q}_q \mathbf{k}_j^T / \sqrt{d})}. \end{aligned}$$

$$\log \frac{\Pr_f(Y = i | X = x)}{\Pr_f(Y = C | X = x)} = \beta_i^T \hat{\mathbf{x}}.$$

$$\log \frac{\Pr_f(Y = i | X = x)}{\Pr_f(Y = C | X = x)} = \beta_0^i + \beta_i^T \mathbf{x}.$$

$$\hat{A}(q, p_i) = \exp(\beta_0^i) A(q, p_i)$$

To train the re-weighting vector $\beta = \{\beta_0^i\}$, we utilize an auxiliary training set $(\mathbf{X}_{train}, \mathbf{Y}_{train})$. Here, we perform ICL with normal demonstrations and optimize β with respect to the classification loss \mathcal{L} on $(\mathbf{X}_{train}, \mathbf{Y}_{train})$:

$$\beta^* = \arg \min_{\beta} \mathcal{L}(\mathbf{X}_{train}, \mathbf{Y}_{train}). \quad (10)$$

Anchor Re-weighting: Results

Method	SST-2	TREC	AGNews	EmoC	Average
Vanilla In-Context Learning (1-shot per class)	61.28	57.56	73.32	15.44	51.90
Vanilla In-Context Learning (5-shot per class)	64.75	60.40	52.52	9.80	46.87
Anchor Re-weighting (1-shot per class)	90.07	60.92	81.94	41.64	68.64

Table 1: The effect after adding parameter β_0^i . For AGNews, due to the length limit, we only use three demonstrations per class. Our Anchor Re-weighting method achieves the best performance overall tasks.

Adding **more demonstrations** for vanilla ICL may **not bring a stable** accuracy boost due to the potential **noise** introduced

This shortens the input context and thus brings (almost) no extra cost to the inference speed.

Anchor-Only Context Compression: Method (Speed)

- a context compression technique that **reduces the full demonstration to anchor hidden states** for accelerating ICL inference.

Text_{anchor}: This method concatenates the formatting and label text with the input, as opposed to concatenating the hidden states at each layer.

Hidden_{random}: This approach concatenates the hidden states of formatting and randomly selected non-label words (equal in number to Hidden_{anchor}).

Hidden_{random-top}: To establish a stronger baseline, we randomly select 20 sets of non-label words in Hidden_{random} and report the one with the highest label loyalty.

(2022b). As a solution, we amalgamate the hidden states of both the formatting and the label words, a method we've termed **Hidden_{anchor}**.

Anchor-Only Context Compression

Method	Label Loyalty	Word Loyalty	Acc.
ICL (GPT2-XL)	100.00	100.00	51.90
Text _{anchor}	51.05	36.65	38.77
Hidden _{random}	48.96	5.59	39.96
Hidden _{random-top}	57.52	4.49	41.72
Hidden _{anchor}	79.47	62.17	45.04
ICL (GPT-J)	100.00	100.00	56.82
Text _{anchor}	53.45	43.85	40.83
Hidden _{random}	49.03	2.16	31.51
Hidden _{random-top}	71.10	11.36	52.34
Hidden _{anchor}	89.06	75.04	55.59

Table 2: Results of different compression methods on GPT2-XL and GPT-J (averaged over SST-2, TREC, AG-News, and EmoC). Acc. denotes accuracy. The best results are shown in bold. Our method achieves the best compression performance.

Text_{anchor}: This method concatenates the formatting and label text with the input, as opposed to concatenating the hidden states at each layer.

Hidden_{random}: This approach concatenates the hidden states of formatting and randomly selected non-label words (equal in number to Hidden_{anchor}).

Hidden_{random-top}: To establish a stronger baseline, we randomly select 20 sets of non-label words in Hidden_{random} and report the one with the highest label loyalty.

(2022b). As a solution, we amalgamate the hidden states of both the formatting and the label words, a method we’ve termed **Hidden_{anchor}**.

Here, “formatting” refers to elements like “Review:” and “Sentiment:”

only leads to a 1.5 accuracy drop

Efficiency improvements

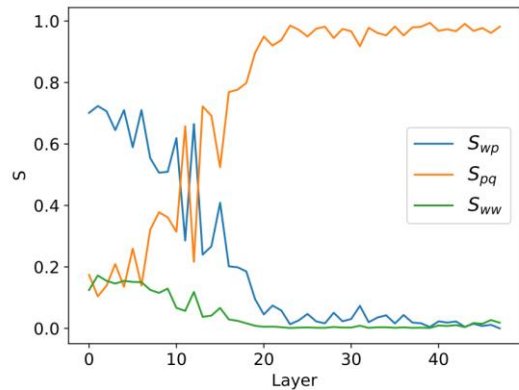
- speed-up ratio ranges from $1.1\times$ to $2.9\times$
- the acceleration effect is more **pronounced** in the **GPT-J** model compared to GPT2-XL, demonstrating its great potential to apply to **larger language models**.

Model	SST-2	TREC	AGNews	EmoC
GPT2-XL	$1.1\times$	$1.5\times$	$2.5\times$	$1.4\times$
GPT-J	$1.5\times$	$2.2\times$	$2.9\times$	$1.9\times$

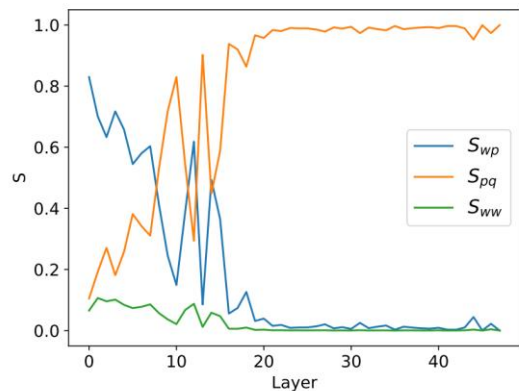
Table 3: Acceleration ratios of the Hidden_{anchor} method.

The background of the slide features a complex network of interconnected nodes and lines. The nodes are represented by small circles in various shades of blue and grey, while the lines are thin and light blue. The network is denser on the left side of the slide and becomes sparser towards the right. The word "Thanks" is centered in the upper half of the slide, overlaid on the network pattern.

Thanks



(a) Results on the SST-2 dataset



(b) Results on the AGNews dataset

Figure 3: Relative sizes of S_{wp} , S_{pq} , and S_{ww} in different layers on SST-2 and AGNews. Results of other datasets can be found in Appendix B. Initially, S_{wp} occupies a significant proportion, but it gradually decays over layers, while S_{pq} becomes the dominant one.

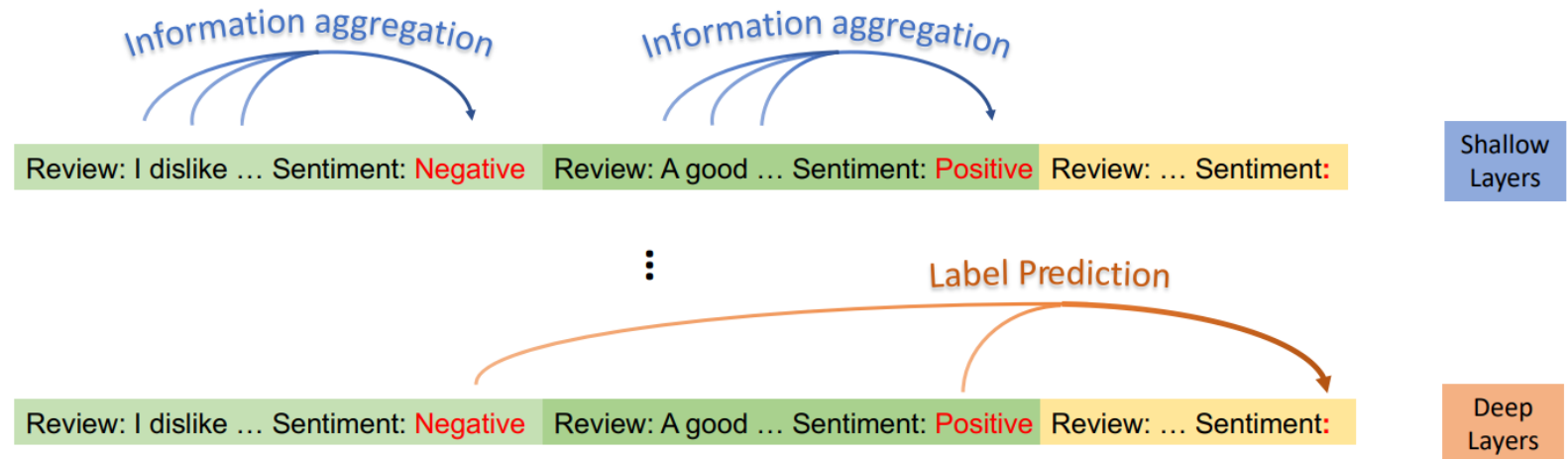


Figure 2: Illustration of our hypothesis. In shallow layers, label words gather information from demonstrations to form semantic representations for deeper processing, while deep layers extract and utilize this information from label words to formulate the final prediction.