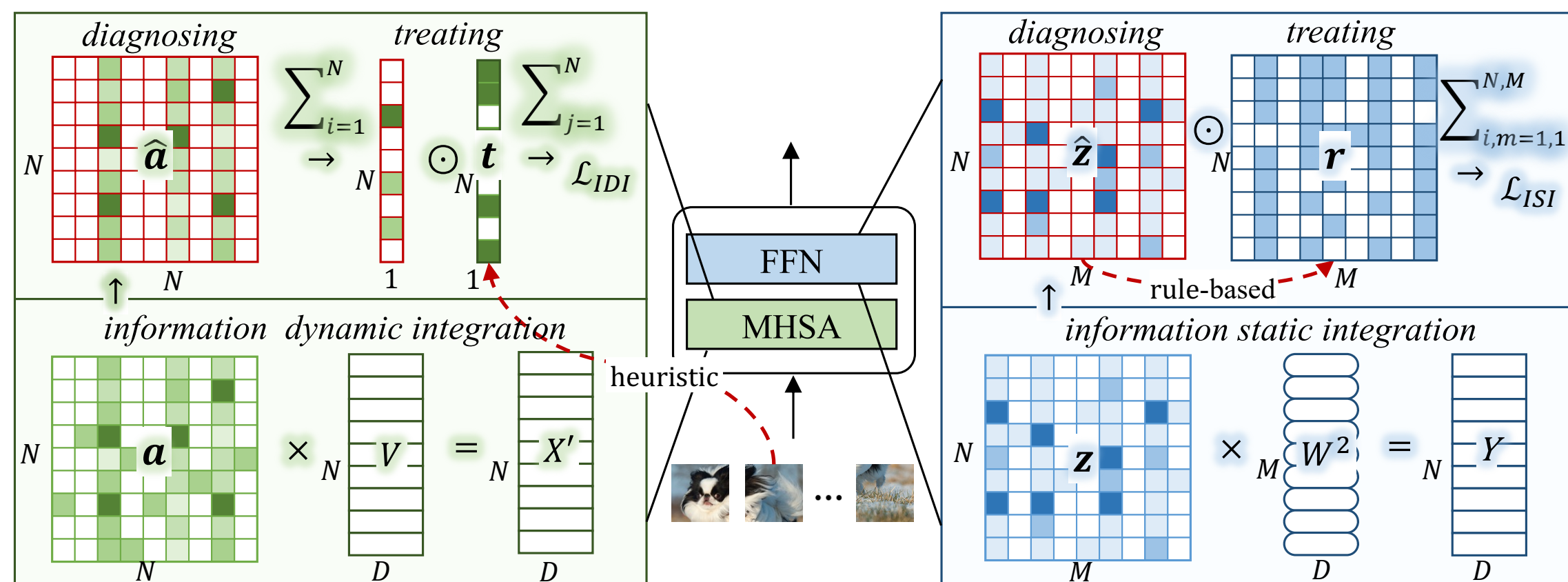


## 1. Introduction

Inspired by research on information integration mechanisms and conjunctive errors in the biological visual system, this paper conducts an in-depth exploration of the internal error mechanisms of Transformers. We first propose an information integration hypothesis for Transformers in the machine vision domain and provide substantial experimental evidence to support this hypothesis. This includes the dynamic integration of information among tokens and the static integration of information within tokens in Transformers, as well as the presence of conjunctive errors therein. Addressing these errors, we further propose heuristic dynamic integration constraint methods and rule-based static integration constraint methods to rectify errors and ultimately improve model performance. The entire methodology framework is termed as Transformer Doctor, designed for diagnosing and treating internal errors within transformers. Through a plethora of quantitative and qualitative experiments, it has been demonstrated that Transformer Doctor can effectively address internal errors in transformers, thereby enhancing model performance.



## 2. Information Integration Hypothesis

**Information Integration Hypothesis:** Similar to biological vision, in machine vision, the Transformer continually processes and refines various mixed information in the primary stage, and integrates it in the advanced stage. When erroneous information is integrated, i.e., conjunction errors occur, it leads to incorrect predictions.

**Potential Information Integration in MHA:** From the query-key-value mechanism of MHA below, it can be observed that a certain token  $X'_i$  in  $X'$  is an integration weighted sum of all tokens in  $V$ , where  $a$  can be considered as integration weights within MHA:

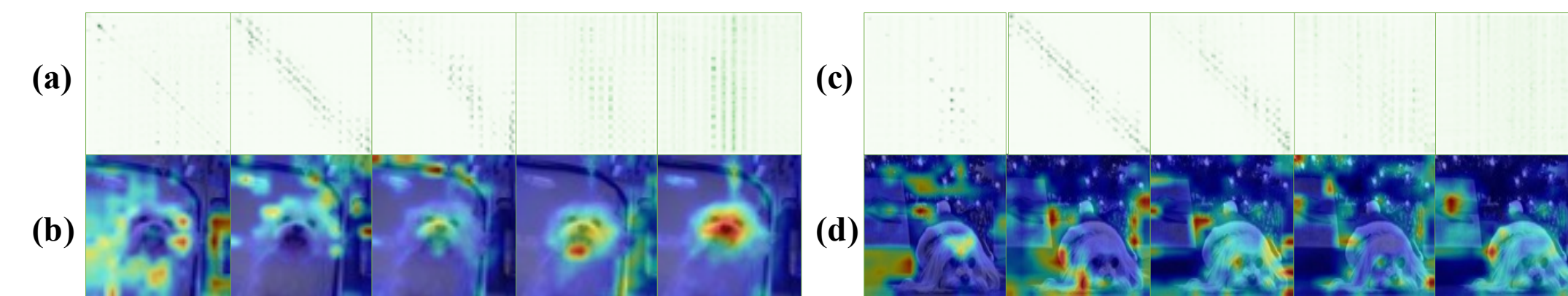
$$X' = \text{softmax}\left(\frac{1}{\sqrt{D}} QK^T\right)V \xrightarrow{\text{decompose}} X'_i = \sum_{j=1}^N a_{i,j} V_j, \quad a = \text{softmax}\left(\frac{1}{\sqrt{D}} QK^T\right)$$

**Potential Information Integration in FFN:** FFN also employs the query-key-value mechanism similar to MHA. As can be seen from the following equation, a certain token  $Y_i$  in  $Y$  is an integration weighted sum of all dimensions in  $W^{(2)}$ , where  $z$  can be referred to as integration weights within the FFN:

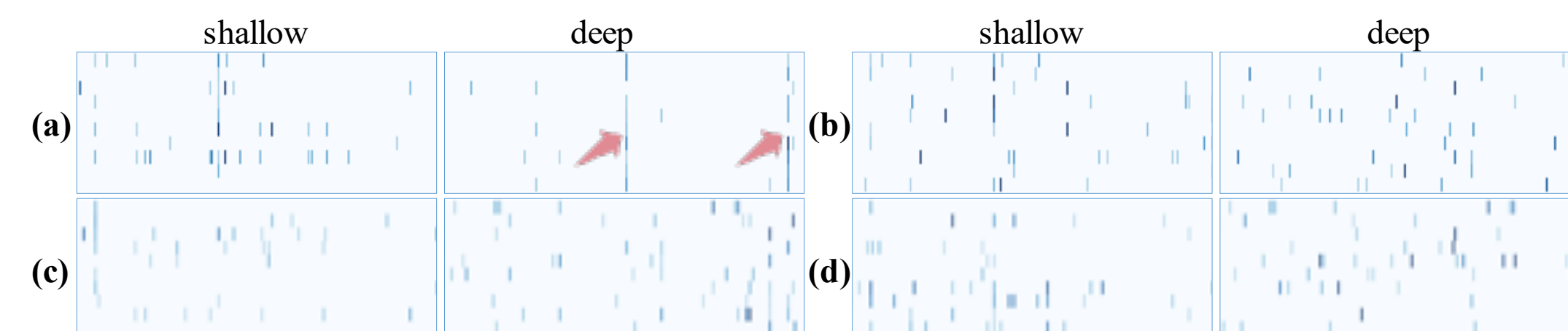
$$Y = \text{gelu}(X'W^{(1)})W^{(2)} \xrightarrow{\text{decompose}} Y_i = \sum_{m=1}^M z_{i,m} W_m^{(2)}, \quad z = \text{gelu}(X'W^{(1)})$$

## 3. Diagnosing

**Inter-token Information Dynamic Integration:** The figure below presents a visual analysis of the integration weight  $a$ . In the initial stages of the Transformer, MHA primarily mixes and processes adjacent patch information among tokens. However, in the advanced stages of the Transformer, MHA dynamically and selectively integrates specific patch information among tokens. When integrating incorrect information among tokens, termed conjunction errors, it leads to model mispredictions.



**Intra-token Information Static Integration:** The figure below illustrates a visual analysis of the integration weight  $z$ . In the initial stages of the Transformer, the FFN primarily mixes and processes various low-level information within tokens. However, in the advanced stages of the Transformer, the FFN statically and selectively integrates specific category information within tokens. When integrating incorrect information within tokens, termed conjunction errors, it leads to model mispredictions.



## 4. Treating

**Heuristic Information Dynamic Integration Therapy:** The integrated weight  $\hat{a}$  in a multi-head scenario is obtained by computing the gradient of the predicted probability  $p_k$  for the true class  $k$ . Then, the integration weight  $\hat{a}$  is constrained by optimizing the loss function  $\mathcal{L}_{IDI}$  based on the foreground annotation  $t$ .

$$\hat{a} = \frac{1}{H} \sum_{h=1}^H \left( \max\left(\frac{\partial p_k}{\partial a^h}, 0\right) \odot a^h \right) \quad \mathcal{L}_{IDI} = \sum_{j=1}^N \left( \left( \sum_{i=1}^N \hat{a}_{i,j} \right) \odot (1 - t_j) \right)$$

**Rule-based Information Static Integration Therapy:** The integration weight  $z$  is updated to  $\hat{z}$  by establishing a connection with the true class  $k$  through gradient-based computation. Subsequently, a binary integration rule  $r$  is derived using  $S$  high-confidence samples. Finally, the loss function  $\mathcal{L}_{ISI}$  is computed based on the integration rule  $r$  to constrain the integration weight  $\hat{z}$ .

$$\hat{z} = \max\left(\frac{\partial p_k}{\partial z}, 0\right) \odot z \quad r = \mathbb{I}(\bar{z} \geq \tau), \quad \bar{z} = \frac{1}{S} \sum_{s=1}^S \hat{z}_{(s)}$$

$$\mathcal{L}_{ISI} = \sum_{m=1}^M \sum_{n=1}^N \left( \hat{z}_{n,m} \odot (1 - r_{n,m}) \right)$$

**Joint Therapy of Dynamic and Static Integration:** Dynamic integration therapy and static integration therapy can be used in combination. In practical applications, it is recommended to sequentially apply static integration constraint therapy followed by dynamic integration constraint therapy to achieve better results.

## 5. Experiments

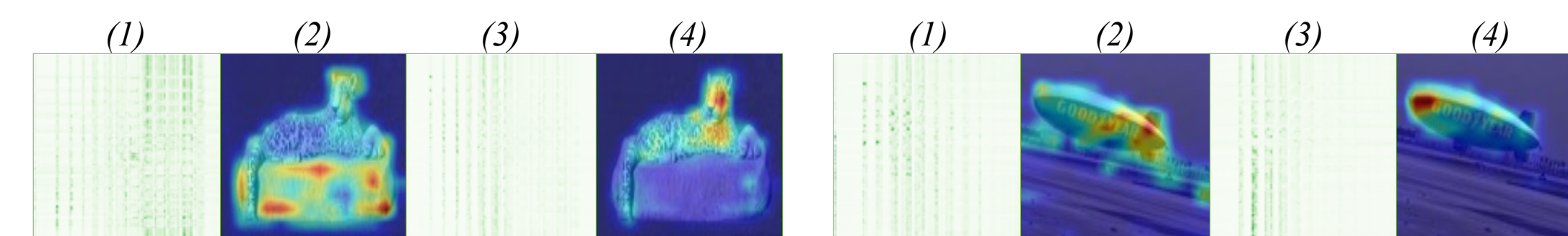
**The Performance of Transformer Doctor on SOTA Models:** ‘+Doctor’ indicates the performance of model treated with Transformer Doctor. From the table below, it is evident that Transformer Doctor effectively enhances model performance on both small and large-scale datasets.

	CIFAR-10	CIFAR-100	ImageNet-10	ImageNet-50	ImageNet-1K
<b>ViT-Tiny</b>	82.17	56.02	78.8	59.62	64.77
<b>+Doctor</b>	83.00 (+0.87)	58.08 (+2.06)	80.80 (+2.00)	61.02 (+1.40)	68.86 (+4.09)
<b>DeiT-Tiny</b>	82.71	56.97	80.20	61.47	66.83
<b>+Doctor</b>	83.96 (+1.25)	59.49 (+2.52)	81.20 (+1.00)	63.19 (+1.69)	70.75 (+3.92)
<b>CaiT-XXS</b>	82.64	56.10	75.80	58.32	66.28
<b>+Doctor</b>	84.20 (+1.36)	60.00 (+3.90)	77.80 (+2.00)	60.16 (+1.84)	70.25 (+3.97)
<b>TNT-Small</b>	83.31	54.67	81.60	65.45	67.25
<b>+Doctor</b>	84.33 (+1.02)	55.60 (+0.93)	83.00 (+1.40)	67.64 (+2.08)	69.97 (+2.72)
<b>PVT-Tiny</b>	83.27	51.94	82.00	71.81	67.73
<b>+Doctor</b>	84.82 (+1.55)	55.10 (+3.16)	84.20 (+2.20)	74.53 (+2.72)	70.94 (+3.21)
<b>Eva-Tiny</b>	87.56	64.15	83.80	71.23	72.51
<b>+Doctor</b>	88.28 (+0.72)	64.99 (+0.84)	85.80 (+2.00)	72.95 (+1.72)	75.45 (+2.94)
<b>BeiT-Tiny</b>	74.69	49.58	79.80	71.59	70.46
<b>+Doctor</b>	76.20 (+1.51)	51.03 (+1.45)	82.20 (+2.40)	73.57 (+1.98)	71.98 (+3.12)

**The Performance of Transformer Doctor with Various Computational Forms:** ‘mean’ denotes directly averaging integration weights across all heads, ‘min’ and ‘max’ respectively represent taking the minimum and maximum integration weights within each head. Each row corresponds to ViT-Tiny and PVT-Tiny architectures, with ImageNet-10 dataset used for evaluation.

Base	+IDI				+ISI	
	min( $a^h$ )	max( $a^h$ )	mean( $a^h$ )	$\hat{a}$	$z$	$\hat{z}$
78.80	79.20 (+0.40)	78.60 (-0.20)	79.20 (+0.40)	<b>80.20 (+1.40)</b>	79.00 (+0.20)	<b>80.40 (+1.60)</b>
82.00	81.20 (-0.80)	81.20 (-0.80)	80.00 (-2.00)	<b>83.60 (+1.36)</b>	82.40 (+0.40)	<b>83.40 (+1.40)</b>

**Comparison of the Intra-token Integration Weights  $a$ :** (1) and (3) depict integration weights before and after treatment, respectively. (2) and (4) show the corresponding heat map effects of integration weights overlaid onto the original image before and after treatment.



**Comparison of the Inter-token Integration Weights  $z$ :** (1) and (2) represent the intra-token integration rules for correct predictions. (3) and (4) depict the intra-token integration weights before and after treatment, respectively.

