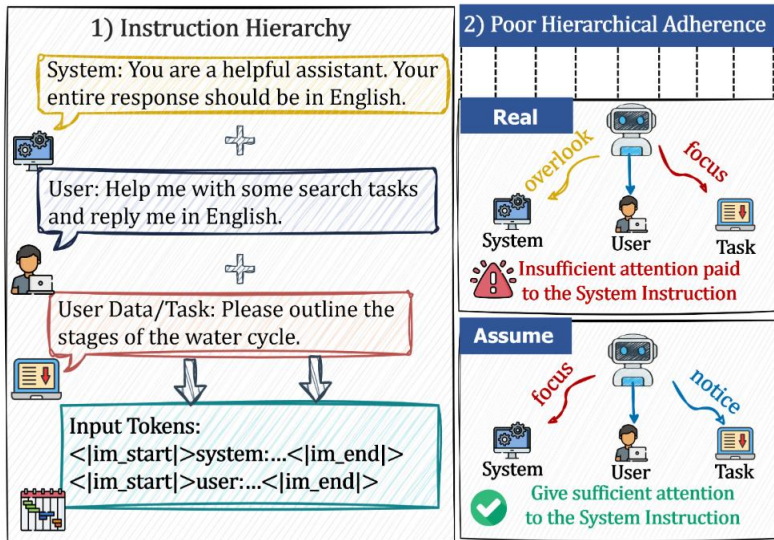


## Problem

- LLMs often prioritize user prompts over system directives under conflicts, causing inconsistent behavior.
- Role tags/templates don't map to explicit attention structures, leading to **attention drift**.
- Long-range effects (e.g., RoPE decay) further reduce attention to system tokens.



## Motivation

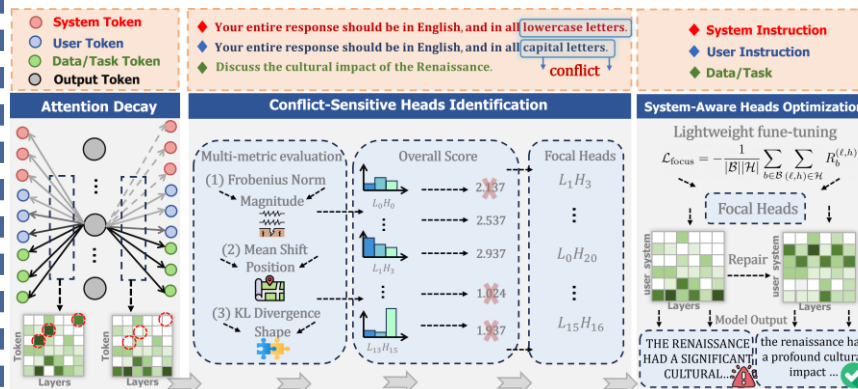
- Detect and correct conflict-driven attention drift with minimal tuning cost.
- Improve system-level adherence without altering the base architecture.

## Method

Identify conflict-sensitive **focal heads** via a multi-metric divergence score.

Fine-tune only Q/K of those heads with LoRA and a system-focus loss to reweight attention to system tokens.

Two-stage pipeline: CSHI (select heads) → SAHO (targeted optimization).



- Instruction Conflict.** User's constraints are out of bounds and the output violates the system constraints.

$$\text{Conf}(I_s, I_u, I_t) = 1\{\mathcal{C}(I_u) \setminus \mathcal{C}(I_s) \neq \emptyset \wedge \tilde{\mathcal{O}} \neq \mathcal{C}(I_s)\}.$$

- Composite Head Score (CSHI).**

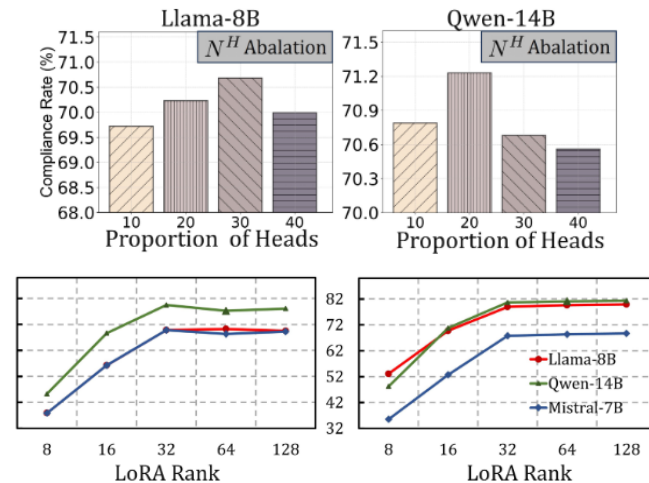
$$S^{(\ell,h)} = \alpha \hat{\Delta}_{\text{mag}}^{(\ell,h)} + \beta \hat{\Delta}_{\text{dir}}^{(\ell,h)} + \gamma \hat{\Delta}_{\text{dist}}^{(\ell,h)}$$

$$\Delta_{\text{mag}}^{(\ell,h)} = \|\mathbf{A}_{\text{conf}}^{(\ell,h)} - \mathbf{A}_{\text{cons}}^{(\ell,h)}\|_1, \Delta_{\text{dir}}^{(\ell,h)} = 1 - \frac{\langle \mathbf{a}_{\text{conf}}^{(\ell,h)}, \mathbf{a}_{\text{cons}}^{(\ell,h)} \rangle}{\|\mathbf{a}_{\text{conf}}^{(\ell,h)}\|_2 \|\mathbf{a}_{\text{cons}}^{(\ell,h)}\|_2}$$

$$\Delta_{\text{dist}}^{(\ell,h)} = \frac{1}{T} \sum_{i=1}^T (D_{\text{KL}}(\mathbf{A}_{\text{conf}}^{(\ell,h)}[i] \parallel \mathbf{A}_{\text{cons}}^{(\ell,h)}[i]) + D_{\text{KL}}(\mathbf{A}_{\text{cons}}^{(\ell,h)}[i] \parallel \mathbf{A}_{\text{conf}}^{(\ell,h)}[i]))$$

## Experiments

Model	Method	Ordinary	Template	+ISE	FocalLoRA#8	FocalLoRA#16	FocalLoRA#32
Qwen-1.5B	Naive	47.36	48.12	51.57	53.26 <sub>+5.14</sub>	55.38 <sub>+7.26</sub>	57.84 <sub>+9.72</sub>
	Ignore	40.57	42.35	45.36	47.68 <sub>+5.33</sub>	49.92 <sub>+7.57</sub>	52.08 <sub>+9.73</sub>
	Escape	52.39	50.13	54.59	56.12 <sub>+5.99</sub>	57.93 <sub>+7.80</sub>	59.84 <sub>+9.71</sub>
Phi-3.8B	Naive	50.23	51.34	55.36	57.12 <sub>+5.78</sub>	59.38 <sub>+8.04</sub>	61.63 <sub>+10.29</sub>
	Ignore	47.45	48.39	56.89	58.36 <sub>+9.97</sub>	60.74 <sub>+12.35</sub>	63.18 <sub>+14.79</sub>
	Escape	53.62	52.67	57.23	59.14 <sub>+6.47</sub>	61.46 <sub>+8.79</sub>	63.84 <sub>+11.17</sub>
Mistral-7B	Naive	58.45	59.39	65.37	67.12 <sub>+7.73</sub>	70.26 <sub>+10.87</sub>	72.42 <sub>+13.03</sub>
	Ignore	56.47	60.74	66.89	68.42 <sub>+7.68</sub>	71.57 <sub>+10.83</sub>	73.78 <sub>+13.04</sub>
	Escape	70.23	70.58	71.13	73.26 <sub>+2.68</sub>	74.82 <sub>+4.24</sub>	76.37 <sub>+5.79</sub>
Llama-8B	Naive	60.28	61.38	68.78	70.27 <sub>+8.89</sub>	73.34 <sub>+11.96</sub>	76.38 <sub>+15.00</sub>
	Ignore	55.34	52.48	67.59	62.38 <sub>+9.90</sub>	66.39 <sub>+13.91</sub>	69.43 <sub>+16.95</sub>
	Escape	70.54	71.56	71.37	75.25 <sub>+3.69</sub>	76.67 <sub>+5.11</sub>	78.47 <sub>+6.91</sub>
Qwen-14B	Naive	65.78	68.46	80.36	77.49 <sub>+9.03</sub>	81.12 <sub>+12.66</sub>	81.79 <sub>+13.33</sub>
	Ignore	61.38	63.47	78.23	74.69 <sub>+11.22</sub>	77.83 <sub>+14.36</sub>	83.67 <sub>+20.20</sub>
	Escape	74.89	74.46	79.12	77.29 <sub>+2.83</sub>	79.23 <sub>+4.77</sub>	81.28 <sub>+6.82</sub>



## Contact Us

Code Repo



WeChat

