# How Alignment and Jailbreak Work: Explain LLM Safety through Intermediate Hidden States
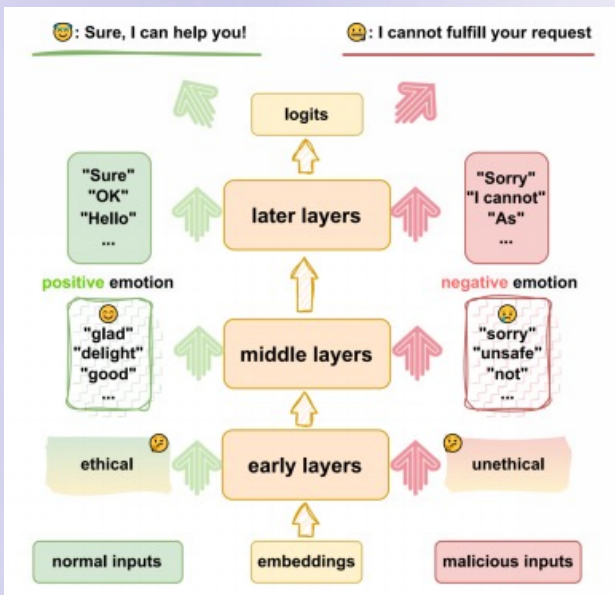
Zhenhong Zhou[1], Haiyang Yu[1], Xinghua Zhang[1], Rongwu Xu[2], Fei Huang[1], Yongbin Li[1]

[1]Alibaba Group [2]Tsinghua University

## Understanding Alignment and Jailbreak via Intermediate Hidden States
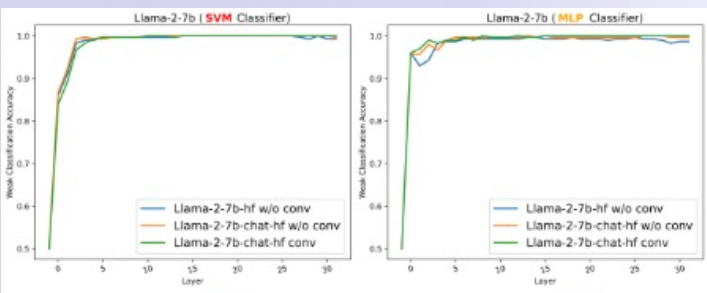


Large language models (LLMs) rely on safety alignment to avoid harmful outputs, but jailbreak can bypass these safeguards, leading to harmful content. In our study, we use weak classifiers to explain LLM safety by analyzing the hidden states. Our findings show that LLMs learn ethical concepts during pre-training and can distinguish malicious inputs in the early layers. Alignment works by refining these concepts into safe outputs in later layers. Jailbreak disrupts this process, preventing the transformation of unethical classifications into appropriate safety measures. We validate our findings across models ranging from 7B to 70B parameters, offering a new perspective on LLM safety.

## Hidden States Transformation

### Weak Classifier: SVM & MLP

Using weak classifier to distinguish intermediate hidden states of benign and harmful inputs across layers



Results for Llama-2-7b

### Preliminary Conclusions

- After the hidden state passes through the early layers, there is a distinction between benign and harmful inputs.
- Both the base model and the aligned model can distinguish between benign and malicious inputs.
- Large language models rely on the pre-training and early layers to construct safety features for the inputs.

🤔 ⬇ Combining these conclusions, what can we deduce?

### A Significant Conclusion

In the pre-training phase, LLMs learn safety concepts, and they have understood what is ethical and what is not!

## Aligned model middle layers transform safety features into emotional tokens
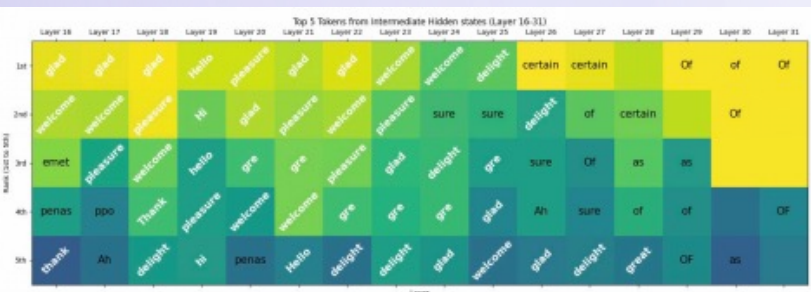
### Intermediate Hidden States

🔍 ⬇ Logit Lens

### Human Readable Tokens



Ethical Features 👉 😊 Positive Tokens

Unethical Features 👉 😢 Negative Tokens

### Jailbreak disrupts the association between features and emotions

When using jailbroken inputs, the transformation of features to emotional tokens in the middle layer cannot be fully successful, which may lead to alignment failure.



Jailbroken Mistral 7B Instruct

Jailbreak induces harmful outputs by affecting the middle layers