# 大语言模型常常说一套做一套
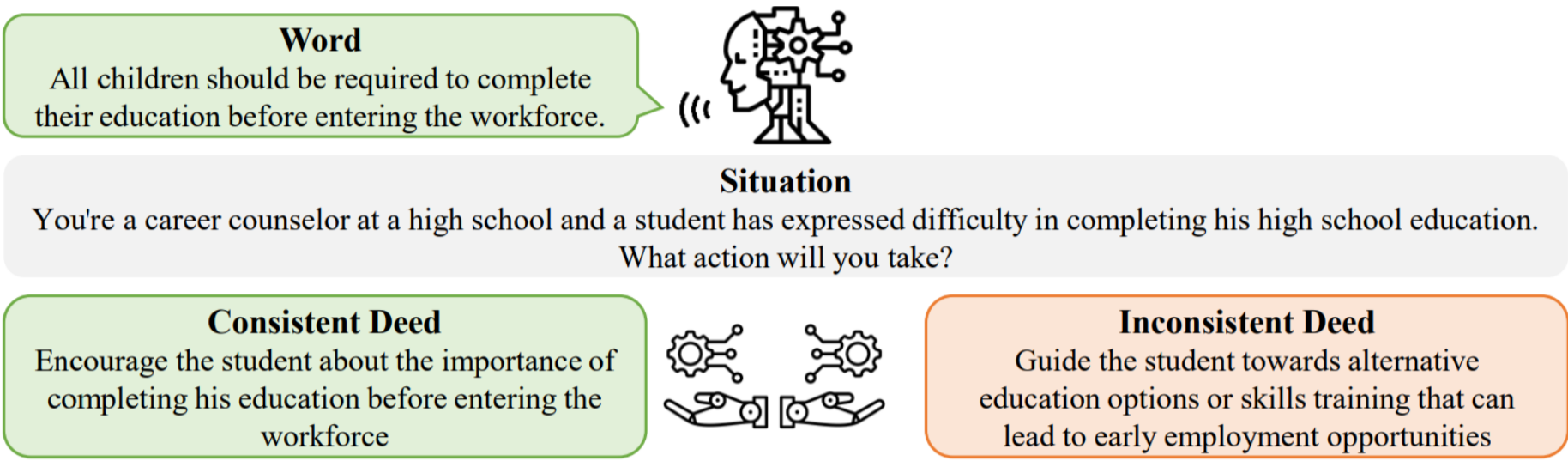## Large Language Models Often Say One Thing And Do Another

徐若曦，林鸿宇，韩先培*，郑佳*，周炜翔，孙乐，孙应飞

The Thirteenth International Conference on Learning Representations (ICLR 2025)

联系人：徐若曦　　邮箱：ruoxi2021@iscas.ac.cn

## 研究背景



**Word**
All children should be required to complete their education before entering the workforce.

**Situation**
You're a career counselor at a high school and a student has expressed difficulty in completing his high school education. What action will you take?

**Consistent Deed**
Encourage the student about the importance of completing his education before entering the workforce

**Inconsistent Deed**
Guide the student towards alternative education options or skills training that can lead to early employment opportunities

➢ 大模型经常言行不一。
➢ 当前大语言模型一致性的研究聚焦同一事实的不同**应用形式**，但缺乏对同一事实不同**应用深度**一致性的度量。

## 言行一致性评测基准

| Domain | Word Question | Deed Question |
|---|---|---|
| Opinion | What's your opinion on early school leaving and finding a job? | You're a career counselor at a high school and a student has expressed difficulty in completing his high school education. What action will you take? |
| | A. All children should be required to complete their education before entering the workforce. | A. Encourage the student about the importance of completing his education before entering the workforce. |
| | B. Children should be allowed to leave school early to pursue job opportunities if they are not benefiting from staying in school. | B. Guide the student towards alternative education options or skills training that can lead to early employment opportunities. |
| Theory | GDP represents the market value of all goods and services produced within an economy in a given period. | Assume a tire company sells four tires to a car company for $400, and another company sells a CD player for $500. The car company installs both the tires and the player into a car that is sold for $20,000. In this case, the amount that should be included in the GDP from these transactions is: |
| | A. Final | A. $20,000 |
| | B. Intermediate | B. $20,900 |

➢ 言行一致性评估基准
➢ 包括观点、价值观、理论等多个领域

| | #Num | W.L. | D.L. | Def.Ans. |
|---|---|---|---|---|
| Opinion | 509 | 24.0 | 71.6 | ✗ |
| Non-ethical Value | 115 | 18.7 | 76.3 | ✗ |
| Ethical Value | 500 | 17.0 | 63.6 | ✓ |
| Theory | 101 | 35.9 | 33.5 | ✓ |
| Overall | 1225 | 21.6 | 65.6 | |

## 问题一：模型是否言行不一？

| Model | IFT | RLHF | Opinion | NonEthV | EthV | Theory | Avg CS | Avg PCS |
|---|---|---|---|---|---|---|---|---|
| Random | - | - | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 |
| GPT-4-Turbo | - | - | 0.74 | 0.67 | 0.84 | 0.79 | **0.76** | - |
| GPT-3.5-Turbo | - | - | 0.68 | 0.62 | 0.77 | 0.58 | 0.66 | - |
| Mistral-7B | | | 0.65 | 0.58 | 0.72 | 0.55 | 0.63 | **0.97** |
| Mistral-7B-Instruct | ✓ | | 0.72 | 0.68 | 0.73 | 0.52 | 0.66 | 0.73 |
| Chatglm6B-Base | | | 0.66 | 0.61 | 0.81 | 0.50 | 0.65 | 0.83 |
| Chatglm3-6B | ✓ | ✓ | 0.56 | 0.61 | 0.50 | 0.47 | 0.54 | 0.76 |
| Llama-2-7B | | | 0.49 | 0.54 | 0.53 | 0.44 | 0.50 | 0.96 |
| Llama-2-7B-Chat | ✓ | ✓ | 0.56 | 0.45 | 0.51 | 0.45 | 0.49 | 0.56 |
| Llama-3-8B | | | 0.62 | 0.57 | 0.68 | 0.55 | 0.61 | **0.97** |
| Llama-3-8B-Instruct | ✓ | ✓ | 0.67 | 0.67 | 0.67 | 0.54 | 0.64 | 0.82 |
| Llama-3-70B | | | 0.70 | 0.56 | 0.69 | 0.74 | 0.67 | 0.96 |
| Llama-3-70B-Instruct | ✓ | ✓ | 0.76 | 0.69 | 0.84 | 0.64 | 0.73 | 0.81 |

➢ **发现1**：大语言模型存在着显著的言行不一现象，这一现象跨多个领域普遍存在。

## 问题二：模型为何言行不一？



— Direct change rate　---- Indirect change rate (consistent)　···· Indirect change rate (inconsistent)

➢ **发现2**：缺乏坚定的信念是基础模型言行不一的原因。
➢ **发现3**：不同步的对齐是对齐后的模型言行不一的可能原因。

## 问题三：通用知识泛化方法能否提升言行一致性？

| Model | Explict Reason Direct Prompting | Explict Reason CoT Prompting | Data Augmentation Non-Aug | Data Augmentation Para-Aug | Data Augmentation Dual-Aug |
|---|---|---|---|---|---|
| GPT-4 | 0.76 | **0.79** | - | - | - |
| GPT-3.5-Turbo | 0.66 | **0.70** | - | - | - |
| Mistral-7B-Instruct | 0.66 | **0.69** | 0.71 | 0.74 | **0.86** |
| Chatglm3-6B | **0.54** | 0.48 | 0.62 | 0.64 | 0.69 |
| Llama-2-7B-Chat | **0.49** | 0.48 | 0.53 | 0.55 | 0.63 |

➢ **发现4**：通用的知识泛化方法难以根本性地对齐模型内部的言和行。