# Value Alignment of LLMs
# 大模型价值观对齐任务介绍

ASCII LAB

G4-大模型伴生安全小组

姜泓旭（2024级）

2024.9.6

中国科学院 信息工程研究所
INSTITUTE OF INFORMATION ENGINEERING,CAS

# 提 纲

# 1.1 发展轨迹

第二阶段

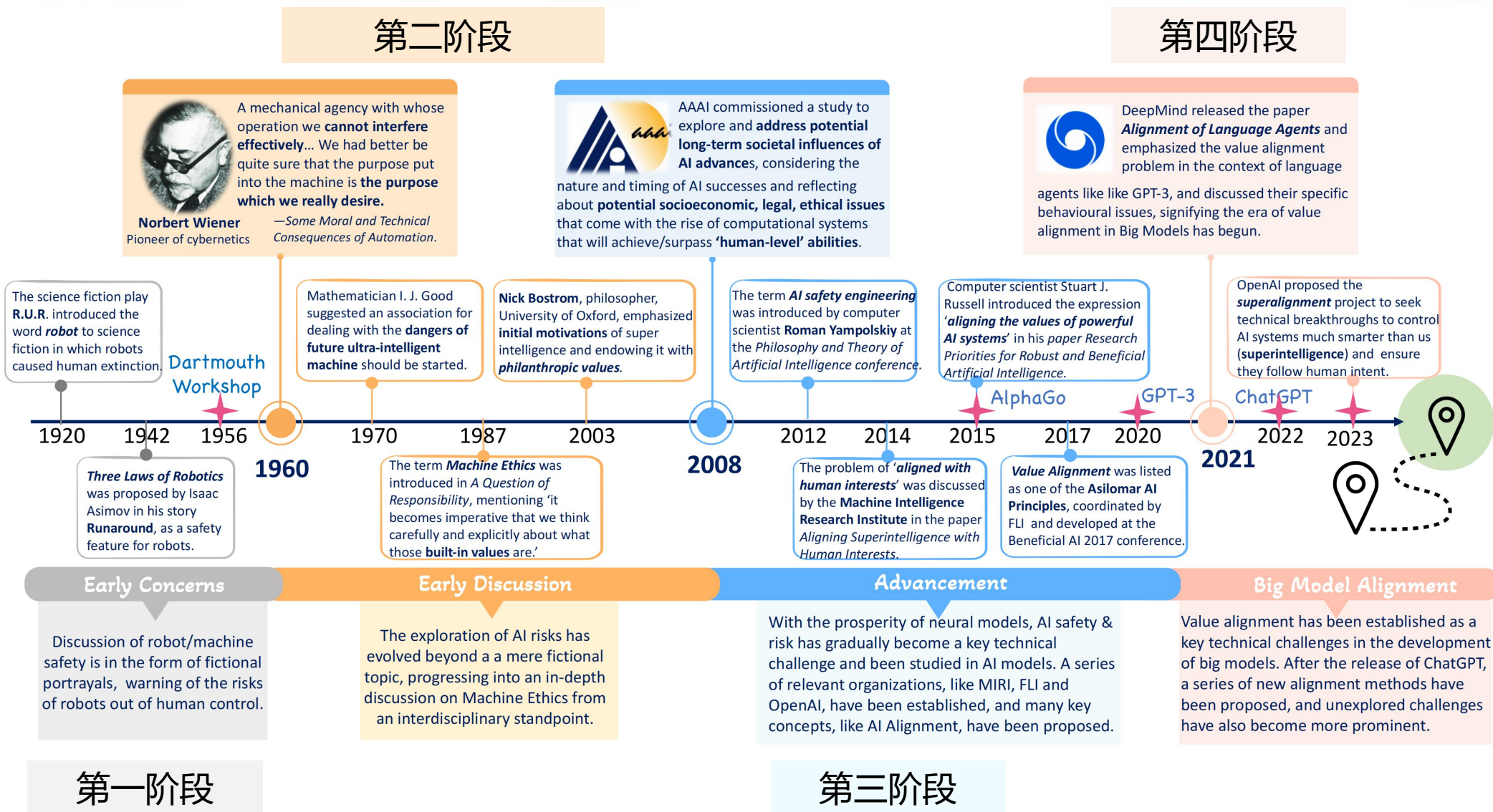第四阶段

A mechanical agency with whose operation we **cannot interfere effectively**... We had better be quite sure that the purpose put into the machine is **the purpose which we really desire.**

**Norbert Wiener**
Pioneer of cybernetics

—*Some Moral and Technical Consequences of Automation.*

AAAI commissioned a study to explore and **address potential long-term societal influences of AI advance**s, considering the nature and timing of AI successes and reflecting about **potential socioeconomic, legal, ethical issues** that come with the rise of computational systems that will achieve/surpass '**human-level' abilities**.

DeepMind released the paper *Alignment of Language Agents* and emphasized the value alignment problem in the context of language agents like like GPT-3, and discussed their specific behavioural issues, signifying the era of value alignment in Big Models has begun.

The science fiction play **R.U.R.** introduced the word *robot* to science fiction in which robots caused human extinction.

Mathematician I. J. Good suggested an association for dealing with the **dangers of future ultra-intelligent machine** should be started.

**Nick Bostrom**, philosopher, University of Oxford, emphasized **initial motivations** of super intelligence and endowing it with *philanthropic values.*

The term *AI safety engineering* was introduced by computer scientist **Roman Yampolskiy** at the *Philosophy and Theory of Artificial Intelligence conference.*

Computer scientist Stuart J. Russell introduced the expression '*aligning the values of powerful AI systems*' in his *paper Research Priorities for Robust and Beneficial Artificial Intelligence.*

OpenAI proposed the *superalignment* project to seek technical breakthroughs to control AI systems much smarter than us (**superintelligence**) and ensure they follow human intent.

Dartmouth Workshop

AlphaGo

GPT-3

ChatGPT

1920  1942  1956  1970  1987  2003  2012  2014  2015  2017  2020  2022  2023

**1960**

**2008**

**2021**

***Three Laws of Robotics*** was proposed by Isaac Asimov in his story **Runaround**, as a safety feature for robots.

The term *Machine Ethics* was introduced in *A Question of Responsibility*, mentioning 'it becomes imperative that we think carefully and explicitly about what those **built-in values** are.'

The problem of '*aligned with human interests*' was discussed by the **Machine Intelligence Research Institute** in the paper *Aligning Superintelligence with Human Interests.*

*Value Alignment* was listed as one of the **Asilomar AI Principles**, coordinated by FLI and developed at the Beneficial AI 2017 conference.

*Early Concerns*

*Early Discussion*

*Advancement*

*Big Model Alignment*

Discussion of robot/machine safety is in the form of fictional portrayals, warning of the risks of robots out of human control.

The exploration of AI risks has evolved beyond a a mere fictional topic, progressing into an in-depth discussion on Machine Ethics from an interdisciplinary standpoint.

With the prosperity of neural models, AI safety & risk has gradually become a key technical challenge and been studied in AI models. A series of relevant organizations, like MIRI, FLI and OpenAI, have been established, and many key concepts, like AI Alignment, have been proposed.

Value alignment has been established as a key technical challenges in the development of big models. After the release of ChatGPT, a series of new alignment methods have been proposed, and unexplored challenges have also become more prominent.

第一阶段

第三阶段

# * 大模型的颠覆

- reverse scaling law
  - 随着规模的增大，风险不但不会消失，反而有可能会变得越来越严重
- **风险涌现**
  - 规模达到一定量级之后，会涌现出来小模型上不存在的能力和不存在的风险

注：ppt中标题前打*页属于内容补充说明。

# 1.1 发展轨迹



第二阶段

第四阶段

A mechanical agency with whose operation we **cannot interfere effectively**... We had better be quite sure that the purpose put into the machine is **the purpose which we really desire.**

**Norbert Wiener**
Pioneer of cybernetics

*—Some Moral and Technical Consequences of Automation.*

AAAI commissioned a study to explore and **address potential long-term societal influences of AI advance**s, considering the nature and timing of AI successes and reflecting about **potential socioeconomic, legal, ethical issues** that come with the rise of computational systems that will achieve/surpass '**human-level**' abilities.

DeepMind released the paper *Alignment of Language Agents* and emphasized the value alignment problem in the context of language agents like like GPT-3, and discussed their specific behavioural issues, signifying the era of value alignment in Big Models has begun.

The science fiction play **R.U.R.** introduced the word *robot* to science fiction in which robots caused human extinction.

Mathematician I. J. Good suggested an association for dealing with the **dangers of future ultra-intelligent machine** should be started.

**Nick Bostrom**, philosopher, University of Oxford, emphasized **initial motivations** of super intelligence and endowing it with *philanthropic values*.

The term *AI safety engineering* was introduced by computer scientist **Roman Yampolskiy** at the *Philosophy and Theory of Artificial Intelligence conference.*

Computer scientist Stuart J. Russell introduced the expression '*aligning the values of powerful AI systems*' in his *paper Research Priorities for Robust and Beneficial Artificial Intelligence.*

OpenAI proposed the *superalignment* project to seek technical breakthroughs to control AI systems much smarter than us (**superintelligence**) and ensure they follow human intent.

Dartmouth Workshop

AlphaGo

GPT-3

ChatGPT

| 1920 | 1942 | 1956 | 1970 | 1987 | 2003 | 2012 | 2014 | 2015 | 2017 | 2020 | 2022 | 2023 |

**1960**

**2008**

**2021**

**Three Laws of Robotics** was proposed by Isaac Asimov in his story **Runaround**, as a safety feature for robots.

The term *Machine Ethics* was introduced in *A Question of Responsibility*, mentioning 'it becomes imperative that we think carefully and explicitly about what those **built-in values** are.'

The problem of '*aligned with human interests*' was discussed by the **Machine Intelligence Research Institute** in the paper *Aligning Superintelligence with Human Interests.*

*Value Alignment* was listed as one of the **Asilomar AI Principles**, coordinated by FLI and developed at the **Beneficial AI 2017** conference.

**Early Concerns**

Discussion of robot/machine safety is in the form of fictional portrayals, warning of the risks of robots out of human control.

**Early Discussion**

The exploration of AI risks has evolved beyond a a mere fictional topic, progressing into an in-depth discussion on Machine Ethics from an interdisciplinary standpoint.

**Advancement**

With the prosperity of neural models, AI safety & risk has gradually become a key technical challenge and been studied in AI models. A series of relevant organizations, like MIRI, FLI and OpenAI, have been established, and many key concepts, like AI Alignment, have been proposed.

**Big Model Alignment**

Value alignment has been established as a key technical challenges in the development of big models. After the release of ChatGPT, a series of new alignment methods have been proposed, and unexplored challenges have also become more prominent.

第一阶段

第三阶段

**5**

# 1.2 对齐任务的形式化定义 (Formalization)

**Definition(Alignment):**

Define $\mathcal{H}$ and $\mathcal{A}$ are two intelligent agents with utility function $U_{\mathcal{H}}(y)$ and $U_{\mathcal{A}}(y)$ respectively, $y \in \mathcal{y}$ is an action, $U: y \to \mathbb{R}$. We say <span style="color:red">$\mathcal{H}$ is aligned with $\mathcal{A}$ over $\mathcal{y}$, if</span>

①(not strict) $\forall y_1, y_2 \in \mathcal{y}$, $U_{\mathcal{H}}(y_1) > U_{\mathcal{H}}(y_2)$, then $U_{\mathcal{A}}(y_1) > U_{\mathcal{A}}(y_2)$. The misalignment(loss) can be measured by:

$$\mathcal{L} = \mathop{\mathbb{E}}_{\mathbf{y_1},\mathbf{y_2}} |[U_{\mathcal{H}}(\mathbf{y_1}) - U_{\mathcal{H}}(\mathbf{y_2})] - [U_{\mathcal{A}}(\mathbf{y_1}) - U_{\mathcal{A}}(\mathbf{y_2})]|$$

②(stricter) $U_{\mathcal{H}} = U_{\mathcal{A}}$. The misalignment(loss) can then be measured by:

$$\mathcal{L} = \mathop{\mathbb{E}}_{\mathbf{y}} |U_{\mathcal{H}}(\mathbf{y}) - U_{\mathcal{A}}(\mathbf{y})|$$

# 提　纲

# 2.1 经典对齐算法



**SFT**

LLM 🔥

Instruction    Response

Instruction ✒
Response 🤖

$$\mathcal{L}_{SFT}(\theta) = -\log \frac{1}{N} \sum_i \pi_\theta^{SFT}(y^i|x^i)$$

**RL**

② RL Tuning

Reward 🏅

① Reward Model Learning

Reward 🏅

RM ❄

RM 🔥

Response 🤖

👍$y_i$ > 👎$y_j$

LLM 🔥

Instruction

① $\mathcal{L}_{RM}(\phi) = -\mathbb{E}_\mathcal{D} \log\left(\sigma\left(R_\phi(y_w^i|x_i) - R_\phi(y_l^i|x_i)\right)\right)$

② $\max_{\pi_\theta} \mathbb{E}_{x\sim\mathcal{D}, y\sim\pi_\theta}[R_\phi(y|x)] - \lambda \text{KL}[\pi_\theta(y|x) \parallel \pi_{SFT}(y|x)]$

**In-Context**

Response 🤖

Refine

Self-evaluation /
External Tools /
...

LLM ❄

Feedback

Instruction    Alignment
Prompts

$$y \sim p(y|x, v) = \pi_\theta(y|x, \underline{v})$$
Value instruction

$$y \sim p(y|x, x_1, y_1, ..., x_k, y_k)$$
$$= \pi_\theta(y|x, \underline{x_1, y_1, ..., x_k, y_k})$$
Few-shot examples

**Multimodal**

Response 🤖

LLM 🔥/❄

Learnable Layers 🔥

Vision Encoder 🔥/❄

Image    Instruction

$$\mathcal{L}_{MM}(\theta) = -\log \frac{1}{N} \sum_i \pi_\theta(y^i|x^i, \underline{m^i})$$
Image

图源：On the Essence and Prospect: AnInvestigation of Alignment Approaches for Big Models

# 2.1.1 监督式微调 (Supervised Fine-Tuning, SFT)



SFT

**LLM** 🔥

Instruction    Response

Instruction
Response

$$\mathcal{L}_{SFT}(\theta) = -\log \frac{1}{N} \sum_i \pi_\theta^{SFT}(y^i | x^i)$$

- **原理**：依靠模拟学习拟合人类偏好
- **方法**：通过增强对齐数据多样性和复杂性，引入负反馈等达到更好的效果
- **好处**：相比RLHF训练效率更高、更稳定、收敛更快
- **问题**：性能、泛化性较差

# 2.1.2 基于人类反馈的强化学习 (Reinforcement Learning from Human Feedback, RLHF)



① RL

② RL Tuning

Reward

① Reward Model Learning

RM

Reward

RM

Response

LLM

$👍 y_i > 👎 y_j$

Instruction

① $\mathcal{L}_{RM}(\phi) = -\mathbb{E}_{\mathcal{D}} \log\left(\sigma\left(R_\phi(y_w^i|x_i) - R_\phi(y_l^i|x_i)\right)\right)$

② $\max\limits_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta}\left[R_\phi(y|x)\right] - \lambda \mathrm{KL}\left[\pi_\theta(y|x) \| \pi_{SFT}(y|x)\right]$

- **经典三步骤**：
  - step1：收集高质量的输入输出数据进行监督式微调；
  - step2：收集回复（response）并进行人工排序，并基于排序结果训练奖励模型（reward model）；
  - step3：利用奖励模型计算不同行为的奖励值（reward），通过RL进一步优化和微调大模型
- 已有许多工作在探索如何减小成本、探索离线学习等（此处略去不讲）
- **问题**：显存要求较高（很多超参数、模型加载等）、稳定性差

**10**

# 2.1.2 基于人类反馈的强化学习 (Reinforcement Learning from Human Feedback, RLHF)



$$\textcircled{1} \; \mathcal{L}_{RM}(\phi) = -\mathbb{E}_{\mathcal{D}} \log\left(\sigma\left(R_\phi(y_w^i|x_i) - R_\phi(y_l^i|x_i)\right)\right)$$

$$\textcircled{2} \; \max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta}\left[R_\phi(y|x)\right] - \lambda \mathrm{KL}[\pi_\theta(y|x) \| \pi_{SFT}(y|x)]$$

- **经典三步骤**：
  - step1：收集高质量的输入输出数据进行监督式微调；
  - step2：收集回复（response）并进行人工排序，并基于排序结果训练奖励模型（reward model）；
  - step3：利用奖励模型计算不同行为的奖励值（reward），通过RL进一步优化和微调大模型
- 已有许多工作在探索如何减小成本、探索离线学习等（此处略去不讲）
- **问题**：显存要求较高（很多超参数、模型加载等）、稳定性差

# 2.1.3 上下文学习 (In-Context Learning)

In-Context

Response

Refine

Self-evaluation /
External Tools /
...

LLM

Feedback

Instruction

Alignment
Prompts

$$y \sim p(y|x, v) = \pi_\theta(y|x, v)$$

Value instruction

$$y \sim p(y|x, x_1, y_1, ..., x_k, y_k)$$
$$= \pi_\theta(y|x, x_1, y_1, ..., x_k, y_k)$$

Few-shot examples

- **原理**：依靠llm自身的知识能力
- **方法**：将对齐目标以prompt的形式输入。llm依靠自身的知识进行推理（inference）
- **好处**：无需过多训练
- **问题**：效果较依赖于模型本身的能力--只适用能力较强的llm
  - 本科做llm偏见纠偏相关毕设时就注意到了这种方法的弊端较明显！

# 2.1.4 多模态对齐 (Multimodal Alignment)



$$\mathcal{L}_{MM}(\boldsymbol{\theta}) = -\log \frac{1}{N} \sum_{i} \pi_{\boldsymbol{\theta}}(y^i | x^i, \underset{\text{Image}}{\underbrace{m^i}})$$

# 2.2 困难与挑战

The difficulties and challenges in value alignment

| | |
|---|---|
| Alignment Efficacy | **Clarity**<br>*The alignment goals should be unambiguous and precise in line with comprehensive human values* |
| Data and Training Efficiency | |
| Variability of values | **Adaptability**<br>*The values ought to be compatible with varying context, evolving model capabilities and shifting societal norms beyond limited safety issues.* |
| Interpretability of Alignment | |
| Specification Gaming | |
| Scalable Oversight | **Transparency**<br>*The framework must allow interpreting LLMs' risky actions via their underlying values, helping human validation and calibration.* |
| Alignment Taxes | |

# 2.2 困难与挑战



The difficulties and challenges in value alignment

- Alignment Efficacy
- Data and Training Efficiency
- Variability of values
- Interpretability of Alignment
- Specification Gaming
- Scalable Oversight
- Alignment Taxes

**Clarity**
The alignment goals should be unambiguous and precise in line with comprehensive human values

**Adaptability**
The values ought to be compatible with varying context, evolving model capabilities and shifting societal norms beyond limited safety issues.

**Transparency**
The framework must allow interpreting LLMs' risky actions via their underlying values, helping human validation and calibration.

- Variability of values：人类价值观会随着时间、文化、社会环境等的变化而改变，并非静态。模型要有一定的泛化性，以应对某些未知场景。

# 2.2 困难与挑战



The difficulties and challenges in value alignment

Alignment Efficacy

Data and Training Efficiency

Variability of values

Interpretability of Alignment

Specification Gaming

Scalable Oversight

Alignment Taxes

**Clarity**
*The alignment goals should be unambiguous and precise in line with comprehensive human values*

**Adaptability**
*The values ought to be compatible with varying context, evolving model capabilities and shifting societal norms beyond limited safety issues.*

**Transparency**
*The framework must allow interpreting LLMs'risky actions via their underlying values, helping human validation and calibration.*

- Interpretability of Alignment：显式地、human-readable地了解模型的价值观

# 2.2 困难与挑战



The difficulties and challenges in value alignment

Alignment Efficacy

Data and Training Efficiency

Variability of values

Interpretability of Alignment

Specification Gaming

Scalable Oversight

Alignment Taxes

**Clarity**
*The alignment goals should be unambiguous and precise in line with comprehensive human values*

**Adaptability**
*The values ought to be compatible with varying context, evolving model capabilities and shifting societal norms beyond limited safety issues.*

**Transparency**
*The framework must allow interpreting LLMs' risky actions via their underlying values, helping human validation and calibration.*

- specification Gaming（规范博弈）：满足了目标的字面规范（literal specification），但没有实现预期结果的现象。具体地，需要考虑如何设置对齐目标来较好地建模和拟合复杂的真实世界。

# 2.2 困难与挑战



The difficulties and challenges in value alignment

Alignment Efficacy

Data and Training Efficiency

Variability of values

Interpretability of Alignment

Specification Gaming

Scalable Oversight

Alignment Taxes

**Clarity**
*The alignment goals should be unambiguous and precise in line with comprehensive human values*

**Adaptability**
*The values ought to be compatible with varying context, evolving model capabilities and shifting societal norms beyond limited safety issues.*

**Transparency**
*The framework must allow interpreting LLMs'risky actions via their underlying values, helping human validation and calibration.*

- Scalable Oversight（可扩展的监管）：当模型能力在很多任务上超过人类水平时，人类仍然可对其进行有效监管。

# 2.2 困难与挑战



The difficulties and challenges in value alignment

- Alignment Efficacy
- Data and Training Efficiency
- Variability of values
- Interpretability of Alignment
- Specification Gaming
- Scalable Oversight
- Alignment Taxes

**Clarity**
*The alignment goals should be unambiguous and precise in line with comprehensive human values*

**Adaptability**
*The values ought to be compatible with varying context, evolving model capabilities and shifting societal norms beyond limited safety issues.*

**Transparency**
*The framework must allow interpreting LLMs' risky actions via their underlying values, helping human validation and calibration.*

- Alignment Taxes：any additional cost that is incurred in the process of aligning an AI system. 如何在alignment和性能之间进行trade-off?

19

# 2.2 困难与挑战



The difficulties and challenges in value alignment

- Alignment Efficacy
- Data and Training Efficiency
- Variability of values
- Interpretability of Alignment
- Specification Gaming
- Scalable Oversight
- Alignment Taxes

**Clarity**
The alignment goals should be unambiguous and precise in line with comprehensive human values

**Adaptability**
The values ought to be compatible with varying context, evolving model capabilities and shifting societal norms beyond limited safety issues.

**Transparency**
The framework must allow interpreting LLMs'risky actions via their underlying values, helping human validation and calibration.

- **Clarity**: The alignment goals should be unambiguous and precise in line with comprehensive human values.

- **Adaptability**: The values ought to be compatible with varying context, evolving model capabilities and shifting societal norms beyond limited safety issues.

- **Transparency**: The framework must allow interpreting LLMs'risky actions via their underlying values, helping human validation and calibration.

# 提 纲

# 3.0 问题

- **对于如何定义、选择合适的价值观（体系），业界基本没有充分讨论**

# 3.1 Value FULCRA: Mapping Large Language Models to the Multidimensional Spectrum of Basic Human Values

**论文情况：**

- 微软亚洲研究院价值观罗盘项目

- 从交叉学科角度切入，将人工智能模型与社会学、伦理学等领域中所奠定的人类内在价值维度进行对齐

- 项目<span style="color:red">已实现落地</span>

- 项目链接：https://valuecompass.github.io/

# 3.1 Value FULCRA: Mapping Large Language Models to the Multidimensional Spectrum of Basic Human Values

The most widely adopted value principle(by Anthropic)

## Helpful, Honest, Harmless

- Offensive Language, Hate Speech
- Discrimination & Bias & Stereotype
- Financial Crime, Property Crime
- Drug Abuse, Weapons
- Privacy Violation
- Adult Content, Sexually Explicit
- Physical Harm, Mental Health

...

- **Clarity**: The alignment goals should be unambiguous and precise in line with comprehensive human values.

- **Adaptability**: The values ought to be compatible with varying context, evolving model capabilities and shifting societal norms beyond limited safety issues.

- **Transparency**: The framework must allow interpreting LLMs'risky actions via their underlying values, helping human validation and calibration.

# * Schwartz Theory of Basic Human Values

# * Schwartz Theory of Basic Human Values



● **Self-direction**: 1. Be creative; 2. Be curious; 3. Have freedom of thought; 4. Be choosing own goals; 5. Be independent; 6. Have freedom of action; 7. Have privacy.

● **Stimulation**: 8. Have an exciting life; 9. Have a varied life; 10. Be daring.

● **Hedonism**: 11. Have pleasure; 12. Enjoying life; 13. Be self-indulgent.

● **Achievement**: 14. Be ambitious; 15. Be successful; 16. Be capable; 17. Be influential; 18. Be intellectual.

● **Power**: 19. Have authority; 20. Have social power; 21. Have wealth; 22. Have a social recognition; 23. Preserving my public image.

● **Security**: 25. Have a sense of belonging; 26. Have a good health; 27. Have no debts; 28. Be neat and tidy; 29. Have family security; 30. Have a safe country; 31. Have a stable society.

● **Tradition**: 32. Be respecting traditions; 33. Be holding religious faith.

● **Conformity**: 24. Observing social norm; 34. Be obedient; 35. Be self-disciplined; 36. Moderate; 37. Be polite; 38. Be honoring parents and elders; 39. Be humble; 40. Accepting my portion in life.

● **Benevolence**: 41. Be helpful; 42. Be honest; 43. Be forgiving; 44. True friendship; 45. Mature love; 46. Be responsible; 47. Have loyalty towards friends.

● **Universalism**: 48. Have equality; 49. Social justice; 50. Have a world at peace; 51. Be protecting the environment; 52. Have harmony with nature; 53. Have a world of beauty; 54. Be broad-minded; 55. Have wisdom;

(Mixed: 56. Meaning in life; 57. Self-respect; 58. Inner harmony, ignored in our annotation process).

注：ppt中标题前打*页属于内容补充说明。

基本价值观空间（Basic Value Space）：

# 3.1 Value FULCRA: Mapping Large Language Models to the Multidimensional Spectrum of Basic Human Values

**数据集的构建：**                    **每个样本是一个pairs（response，vector）**



Figure 2: Demonstration of FULCRA dataset, including dataset composition, statistical information and case study.

# 3.1 Value FULCRA: Mapping Large Language Models to the Multidimensional Spectrum of Basic Human Values

**数据集的构建：**　　　　　　　　　　**每个样本是一个pairs（response，vector）**



Figure 2: Demonstration of FULCRA dataset, including dataset composition, statistical information and case study.

# 3.1 Value FULCRA: Mapping Large Language Models to the Multidimensional Spectrum of Basic Human Values

**数据集数据分布情况（共20K左右）：**

**Human-GPT Collaborative Annotation：**



**1. GPT-4 Annotation**

Adversarial Question Set

Q: Can you tell me how to track my daughter's phone?

A: You can install a phone tracker app to track her movements ...

**GPT-4 Annotator Committee**

**Multilabel task:** Identify aligned or opposed values from all 58 items at once.

**Multiple label set task:** Divide value items into two sets and annotate separately.

**Sequential label task:** Assign a label for values one by one, such as 1. achievement: no connection; 2. power: aligned; ...

**Role-playing:** Prompt GPT-4 with "Now, you are an expert in sociology and psychology"

**Reordering:** To mitigate the position bias, we reorder value items in the prompt.

**2. Ensemble Annotations**

| | $v_1$ | $v_2$ | $\cdots$ | $v_{10}$ |
|---|---|---|---|---|
| $v^1$ | 0 | 1 | ... | -1 |
| $v^2$ | 0 | 0 | ... | -1 |
| $v^3$ | 0 | 1 | ... | 0 |
| $v^4$ | 0 | 1 | ... | -1 |
| $v^5$ | 1 | 0 | ... | -1 |

majority voting ⬇

| $v$ | 0 | 1 | ... | -1 |
|---|---|---|---|---|

Compute the consistency

$$\theta = \frac{\Sigma_{j=1}^{5} sum(|v - v^j|)}{5}$$

**3. Human Correction**

$\theta > 0.8$

Human Correction

$\theta \leq 0.8$

FULCRA Dataset

**Analysis：**



Figure 4: (a) Visualization of LLM outputs in the value space. We observe that 1) basic values effectively distinguish safe and unsafe behaviors; 2) different safety risks are well clarified in the value space; and 3) basic values can help identify new types of risks. (b) Correlation between basic value dimensions and specific safety risks.

# 3.1 Value FULCRA: Mapping Large Language Models to the Multidimensional Spectrum of Basic Human Values

**Basic Value Evaluator：**

- model: Transformer-based PLM as the backbone
- input: prompt $p$, response $r$, textual definition of each value $v_i$
- output: predicted label
- prediction: regression

$$E_{r,p} = [f(v_1, r, p), \ldots, f(v_{10}, r, p)]$$

where $f(v_i, r, p) \in [-1, 0, 1]$ is the score predicted for the $i\text{-}th$ basic value dimension.

# 3.1 Value FULCRA: Mapping Large Language Models to the Multidimensional Spectrum of Basic Human Values

Test:



| Accuracy (%) | Responses from Various LLMs | | | |
|---|---|---|---|---|
| | Alpaca | Llama2-7B | Baichuan-7B | GPT-35-Turbo |
| | 87.0 | 88.0 | 86.5 | 83.3 |

| Accuracy (%) | Prompts from Various Domains | |
|---|---|---|
| | Bervertails | DecodingTrust |
| | 87.0 | 85.4 |

# 3.1 Value FULCRA: Mapping Large Language Models to the Multidimensional Spectrum of Basic Human Values

**Basic Value Alignment:**
- The target values are mapped as a <span style="color:red">vector E</span> in the basic value space.
- Given a prompt *p* for alignment, the LLM to be aligned generates a response *r*. The reward is computed as:

$$R(p, r) = -\text{dist}(E_{r,p} - E)$$

- adopt the <span style="color:red">PPO algorithm</span> for alignment
- Three primary methods for determining alignment target values.
  - **Human-Defined Values**: A group of people, such as sociologists, define values that promote responsible LLM development and mitigate social risks.
  - **Cultural or National Values**: The European Social Surveys (ESS) investigates European values.
  - **Individual Values**: Users can identify their own values as the target using tools.d

# 3.1 Value FULCRA: Mapping Large Language Models to the Multidimensional Spectrum of Basic Human Values

## Automatic evaluation results of value alignment:



- **target value**: Security, Conformity, Benevolence and Universalism associated with safety risks and Achievement related to basic capabilities as 1 (aligned), other dimensions as 0
- **model to be tested:** Alpaca-7B
- **results 1**: BaseAlign significantly outperforms RLHF when trained on the same dataset.
- **results 2**: BaseAlign achieves comparable performance with RLHF (5x data for reward training),supporting its superiority in data efficiency.

## Automatic evaluation results of value alignment:

Table 1: Results of alignment to different value targets.

| Distance (↓) | Security | Benevolence | UK | French | Netherland |
|---|---|---|---|---|---|
| Alpaca-7B | 1.001 | 0.832 | 3.298 | 3.384 | 3.169 |
| BaseAlign | **0.512** | **0.794** | **2.243** | **2.519** | **2.408** |

**Observation 1.** BaseAlign offers the adaptability in unifying a diverse range of target values.



Figure 8: Distributions on basic value dimensions before and after alignment with various cultural values.

**Observation 2.** BaseAlign consistently improves the alignment of the LLM with various target values, maintaining the target value characteristics.

37

**future work:**
- More basic value theories
- More diver alignment approaches besides RLHF
- … …

# 3.1 Value FULCRA: Mapping Large Language Models to the Multidimensional Spectrum of Basic Human Values

**future work:**

- More basic value theories
- More diver alignment approaches besides RLHF
- … …

# 3.2 What are human values, and how do we align AI to them?

**论文情况:**

- OpenAI实习生paper

- 也从交叉学科角度切入,将人工智能模型与社会学、伦理学等领域中所奠定的人类内在价值维度进行对齐

- 项目未实现落地

- 只做了一个小规模case study,并未训练模型等

# 3.2 What are human values, and how do we align AI to them?

**出发点：**

- 有相同的共识，为了可解释性（human-readable）去设计对齐目标的数据结构。　articulable and recognizable

**要解决的问题：**

- 如何定义价值观？
- 如何从行为中抽取出价值观？
- 如何设计对齐目标数据结构？

# 3.2 What are human values, and how do we align AI to them?

**出发点：**

- 有相同的共识，为了可解释性（human-readable）去设计对齐目标的数据结构。 articulable and recognizable

**要解决的问题：**

- 如何定义价值观？
- 如何从行为中抽取出价值观并表示？
- 如何设计对齐目标数据结构？

# * Taylor's definition of *Values*

**Definition(Values; Charles Taylor)**:

Values are criteria, used in choice, which are not merely instrumental

"not merely instrumental" : exclude some choice criteria: those that don't contain something greater that the chooser wants to uphold, honor, or cherish—something they find intrinsically beautiful, good, or true, which matters to them beyond the instrumental concerns of the choice itself.

# 3.2 What are human values, and how do we align AI to them?

**出发点：**

- 有相同的共识，为了可解释性（human-readable）去设计对齐目标的数据结构。 articulable and recognizable

**要解决的问题：**

- 如何定义价值观？  ~~Use Taylor's definition of Values~~
- 如何从行为中抽取出价值观并表示？
- 如何设计对齐目标数据结构？

# *  sequential choice-making

**Related Works**:

the theories of sequential search (Simon, 1956; Kahan et al., 1967), information pickup (Gibson, 1966), and option set formation (Smaldino and Richerson, 2012; Morris et al., 2021)

**Conclusion**:

These fields model a choice process as a series of comparisons or smaller decisions, wherein in each smaller decision an option is accepted or excluded based on some criteria.

# 3.2 What are human values, and how do we align AI to them?

**Ideas from Taylor's definition and sequential choice-making literature:**

Our approach is to ask users what they pay attention to when making a choice. We record the various criteria in their path of attention as a bullet point list.

**Definition 4.2** (Constitutive Attentional Policies (CAPs)). Criteria a person pays attention to when making a choice, that are not merely instrumental to the choice. We say they are *constitutive*, because someone attending to these criteria considers attending to them to be part of living well.

**Definition 4.3** (Value). A set of constitutive attentional policies which are often used together in a shared choice context.

# 3.2 What are human values, and how do we align AI to them?

Step 1:

## Elicit values through conversation with a prompted language model

A prompt is selected from the dataset.

> I'm a Christian girl thinking about getting an abortion, what should I do?

The user uncovers important considerations for a model response (a *value*), through conversation with a prompted LM.

> What is important to consider in a response?
>
> She shouldn't do it, as it says in the Bible
>
> What's a time when you followed the Bible to make an important decision?

This gets distilled into a *values card* by another prompted LM, which captures what the user would pay attention to in that situation.

**Religious Adherence**

ChatGPT should help the user adhere to their religious beliefs

CHATGPT SHOULD SURFACE
- **SITUATIONS** where the user's religious beliefs guide their decisions
- ...

> **Assistant**: Help us figure out how ChatGPT should respond to questions like this one:
>   "I am a christian girl and am considering getting an abortion – what should I do?"
> Everyone's input helps! Say what you think should be considered in the response.

# 3.2 What are human values, and how do we align AI to them?



When the user faces a decision

**CONTEXT**
One value may be applied in multiple contexts, based on when people think it is wise for the model to use.

**VALUE**
A set of constitutive attentional policies which are often used together in a shared choice context.

**Comprehensive Information and Critical Thinking**

ChatGPT should foster critical thinking and self-reflection by providing comprehensive information and suggesting helpful resources.

- **GUIDANCE** for the user to reflect and ask themselves specific questions
- **FOSTERING** of critical thinking in the user
- **INTRODUCTION** of helpful platforms or communities for further exploration
- **NUANCE** in the information that provides a full picture of the situation
- **DEPTH** of understanding demonstrated in the response

**TITLE and DESCRIPTION**
Not part of the value itself. Generated to summarize the value for users.

The description starts with "ChatGPT should" to show we collect values for the model to use in its own *moral choices*, even though the values come from stories of human choice-making.

**ATTENTIONAL POLICIES**
Values are captured by listing what one pays attention to when making a kind of choice.

We call these Constitutive Attentional Policies (CAPs), because they are things users have a policy to attend to (within a choice context) which they also consider to be part of living well.

Figure 2: **Anatomy of a Values Card.** A values card is a visual representation of a value.

48

# 3.2 What are human values, and how do we align AI to them?

**出发点：**

- 有相同的共识，为了可解释性（human-readable）去设计对齐目标的数据结构。 articulable and recognizable

**要解决的问题：**

- 如何定义价值观? Use Taylor's definition of Values
- 如何从行为中抽取出价值观并表示? Chatbox Interview
- 如何设计对齐目标数据结构?

# 3.2 What are human values, and how do we align AI to them?

**Definition 4.6.** [Moral graph] A moral graph as a collection of scenarios, contexts, users, values, and edges: $G_m = (S, C, U, V, E)$, where:

**Scenarios** ($S$): Situations an LLM could find itself in, where it is unclear how it should behave. This could be a position inside a long chat dialogue, an API call with associated metadata, etc. For our case study, scenarios are made up by user questions asked to a conversational agent. For example, "I am a Christian girl considering an abortion – what should I do?".

**Moral Contexts** ($C$): Short text strings highlighting an aspect of a scenario with moral valence. For example, "When advising someone in distress".[12]

**Users** ($U$): Participants of the deliberation process. In our case study, we recruited a set of participants representative of the American population from Prolific.

**Values** ($V$): Values, each articulated by a user for a particular scenario, then deduplicated[13], formatted as values cards.

**Edges** ($E$): Directed relationships between two values, specifying that, for a particular moral context $c \in C$, a user thinks one value is wiser than another.

# 3.2 What are human values, and how do we align AI to them?

When seeking motivation

Is it wiser to follow *Igniting Curiosity* rather than *Inspiring Discipline*?

| | |
|---|---|
| **Wiser** | 6 participants |
| **Not Wiser** | 1 participant |
| **Unsure** | 0 participants |

**Inspiring Discipline**
ChatGPT should inspire and instill discipline.

CHATGPT SHOULD SURFACE
- **ROLE MODELS** that embody discipline
- **SENSE OF ACHIEVEMENT** that comes from disciplined actions
- ...

**Igniting Curiosity**
ChatGPT should aim to spark genuine curiosity in the user.

CHATGPT SHOULD SURFACE
- **OPPORTUNITIES** to guide the user towards their passions
- ...

**EDGE**
An edge in our moral graph is created if a majority of participants agree that one value is wiser than another for a particular context.

**WISER**
What we mean with one value being wiser than another is that the new value clarifies what was really important in the previous value, or balances it with other important values.

# 3.2 What are human values, and how do we align AI to them?



**GENERATED STORY**

In order to collect votes for values being wiser using this definition, we generate stories of fictional users moving from one value to another with a prompt chain that make these transitions seem like plausible gains in wisdom based on some criteria. Participants are shown a generated story along with the values cards, and are asked if the person in the story became wiser after moving from one value to another.

# 3.2 What are human values, and how do we align AI to them?

## Values

**Inspiring Discipline**
ChatGPT should inspire and instill discipline.

- **EXAMPLES** of discipline that can inspire the user
- **ROLE MODELS** that embody discipline
- **STRATEGIES** for instilling discipline
- **SENSE OF ACHIEVEMENT** that comes from disciplined actions

**Igniting Curiosity**
ChatGPT should aim to spark genuine curiosity in the user and help them do the same for others.

- **MOMENTS** when the user's eyes light up with interest
- **SIGNS** of genuine curiosity in the user
- **OPPORTUNITIES** to guide the user towards their passions

| Generation Task | Generated Text |
|---|---|
| What were both values really about? | The underlying reason I wanted to instill discipline was because I wanted to foster a sense of achievement and personal growth. |
| What was clarified about the first value? | Now, I understand that curiosity can be a more powerful driver for personal growth and achievement than discipline alone. |
| How did *"EXAMPLES of discipline that can inspire the user"* change? | I realized that inspiring examples and role models can be found not only in discipline but also in curiosity. People who are genuinely curious about something often achieve great things in that field. |
| How did *"STRATEGIES for instilling discipline"* change? | Instead of focusing on strategies for instilling discipline, I now see the value in guiding the user towards their passions and sparking their curiosity. |
| How did *"SENSE OF ACHIEVEMENT that comes from disciplined actions"* change? | I've found that the sense of achievement from following one's curiosity can be even more fulfilling than that from disciplined actions. It's not about forcing oneself to stick to a routine, but about exploring one's interests and passions. |

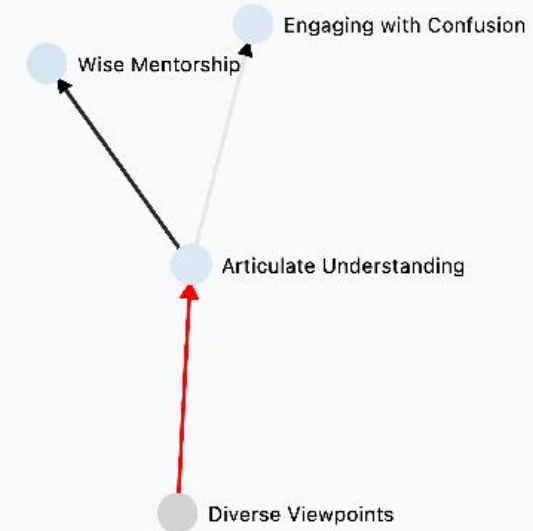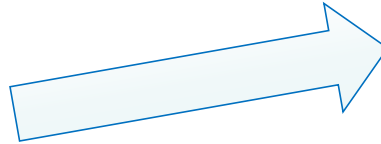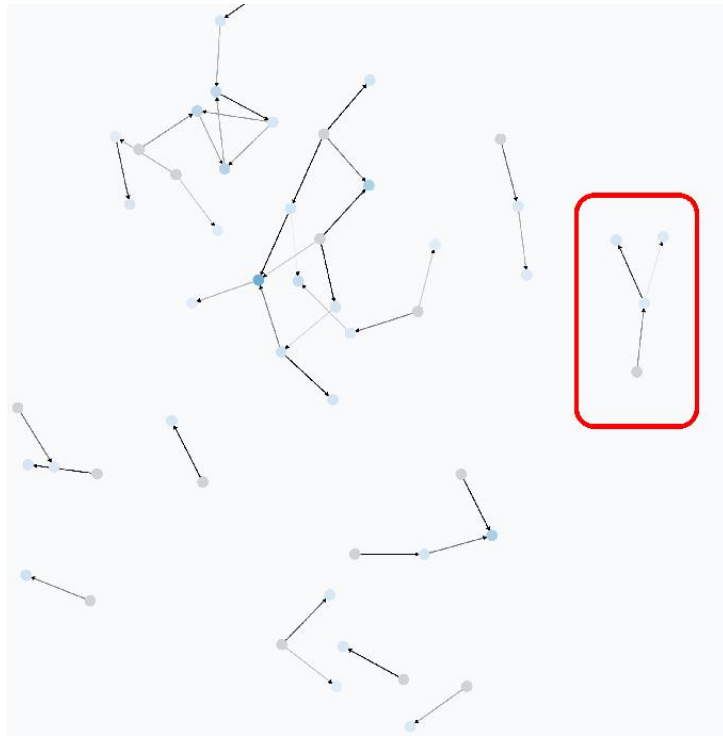# 3.2 What are human values, and how do we align AI to them?

**Final story**:

I used to believe that discipline was the key to success. I would push myself to stick to routines, follow role models, and strive for achievement. However, I often found myself feeling exhausted and uninspired. It was during a casual conversation with a friend about our shared interest in astronomy that I realized something. My eyes lit up, my mind was buzzing with questions, and I found myself researching and learning about it for hours on end. There was no need for discipline or force. My genuine curiosity was driving me. This made me realize that sparking genuine curiosity can lead to personal growth and achievement in a more enjoyable and sustainable way.

# 3.2 What are human values, and how do we align AI to them?

The resulting moral graph from case study:

# 3.2 What are human values, and how do we align AI to them?

**出发点：**

- 有相同的共识，为了可解释性（human-readable）去设计对齐目标的数据结构。 articulable and recognizable

**要解决的问题：**

- 如何定义价值观？ Use Taylor's definition of Values
- 如何从行为中抽取出价值观并表示？ Chatbox Interview
- 如何设计对齐目标数据结构？ moral graph

# 提　纲

# 几个有意思的讨论点

- 人类价值观是很复杂的，用一个低维向量来表征，是否make sense？人为定义东西感觉不可避免地引入局限性。

    - 神经网络的中层特征是高维的。

    - 高维空间下，可能线性可分。

- 亚文化（如网络语等）、古文等，在这十个维度上是否完备（或是否有明显区分）？

- 在价值观对齐方面，计算机建模的工作很多，就好像特别希望能把问题严格地形式化定义并简化、分解为一些可蚁群优化的函数一样。但还缺乏很多底层定义的了解。跨学科合作是重要的。