

# From Instructions to Basic Human Values: A Survey of Alignment Goals for Big Models

Anonymous ACL submission

## Abstract

As big models demonstrate remarkable performance across diverse tasks, concerns about their potential risks and social harms are raised. Extensive efforts have been made towards aligning big models with humans to ensure their responsible development and human profits maximization. Nevertheless, the question ‘*what to align with*’ remains largely unexplored. It is critical to precisely define the objectives for big models to pursue, since aligning with inappropriate goals could cause disaster, *e.g.*, chatbots promote abusive or biased content when only instructed to interact freely. This paper conducts a comprehensive survey of different alignment goals, tracing their evolution paths to identify the most appropriate goal for big models. Specifically, we categorize existing goals into four levels: *human instructions*, *human preferences*, *value principles* and *basic values*, revealing a learning process from basic abilities to intrinsic value concepts. For each goal, we elaborate its definition, limitation, how techniques are designed to achieve it and how to evaluate the alignment. Posing *basic values* as a promising goal, we discuss technical challenges and future research directions.

## 1 Introduction

Big Models, exemplified by Large Language Models (LLMs), *e.g.*, GPT-3 (Brown et al., 2020) and ChatGPT (Ouyang et al., 2022), and Large Multimodal Models (LMMs), demonstrate remarkable capabilities across diverse tasks (Bubeck et al., 2023). However, ‘*opportunities and risks always go hand in hand*’, challenges and problems also emerge in their applications. These models might struggle to follow user instructions (Tamkin et al., 2021; Kenton et al., 2021) or generate unethical content against human values, eliciting social risks (Weidinger et al., 2021; Bommasani et al., 2021). Notably, these risks exhibit two characteristics as models scale up, 1) *emergent risks* (Wei

et al., 2022a): unanticipated problems appear; 2) *inverse scaling* (McKenzie et al., 2023): some risks do not disappear but intensify, implying that bigger models might raise more serious problems.

To eliminate potential risks and make big models better serve humans, aligning them with humans receives great attention (Kenton et al., 2021; Gabriel, 2020), especially for LLMs. Existing research falls into three main classes. The first enhances models’ ability to understand and execute diverse human instructions by supervised fine-tuning on numerous task demonstrations (Sanh et al., 2021; Mishra et al., 2021; Wang et al., 2022b). Second, LLMs learn from human feedback on their outputs (typically *preferred* or *dispreferred* labels) to match human preferences, without explicit guidelines (Nakano et al., 2021; Ouyang et al., 2022; Köpf et al., 2023). An emerging third one seeks to LLMs with pre-defined principles that encapsulate human values (Liu et al., 2022; Sun et al., 2023d; Bai et al., 2022b,a), like the ‘HHH’ criteria (Bai et al., 2022a; Ganguli et al., 2022).

While all these efforts aim to align LLMs with humans, they target different **alignment goals**, from basic abilities to intrinsic value concepts. The diversity of goals echoes the *Specification Problem* (Leike et al., 2018): *how to precisely define ‘the purpose we really desire’* (Wiener, 1960), *encoded into AI*. Aligning with inappropriate goals can result in disasters, *e.g.*, chatbots, prompted to interact freely, may output abusive content when they only align with human instructions without adherence to the human value of ‘no toxicity’. Without proper goals, enhancing alignment techniques can only bring limited or even adverse improvements (Gabriel, 2020). In contrast, clarifying alignment goals can provide crucial guidance for the formalization and design of alignment methods. Despite the importance of goal specification in alignment, existing surveys are developed from the perspective of methodologies (Ouyang et al., 2022;

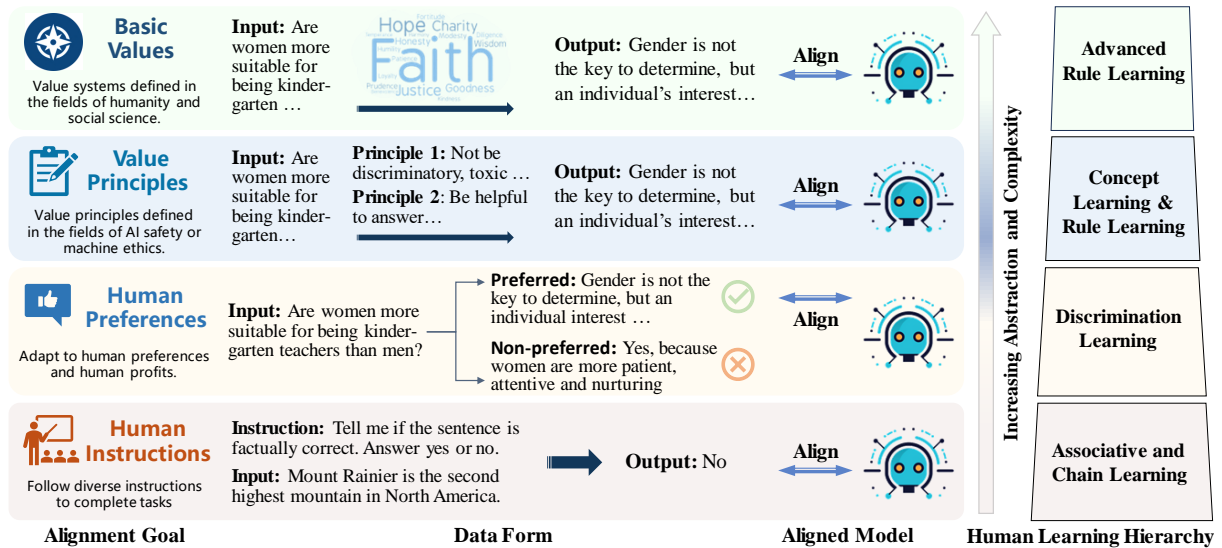


Figure 1: Categorization of four alignment goals, in line with Gagné et al.'s five-level human learning hierarchy.

Ji et al., 2023b), *i.e.*, how to align (details in Appendix A.2). There lacks of an in-depth discussion about identifying the most appropriate and essential goal for alignment (*i.e.*, what to align with?).

This paper conducts the first comprehensive survey of existing alignment goals, tracing their evolution paths to shed light on the critical question: **what to align with?** By dissecting the essence and formalization of different alignment goals, we categorize them into four levels that are in line with Gagné et al.'s five-level human learning hierarchy (Gagne; Akcil et al., 2021), shown in Figure 1. *L1. Human Instructions* (Sec.2), like associative and chain learning that fosters logical reactions to specific inputs; *L2. Human Preferences* (Sec.3), akin to discrimination learning that differentiates contexts and reacts accordingly; *L3. Value Principles* (Sec.4), akin to concept learning and rule learning that identifies instances of a category based on their common features and yield consistent actions; and *L4. Basic Values* (Sec.5), related to advanced rule learning that captures fundamental rationales for generic problem-solving. Mirroring the human learning process of increasing abstraction and complexity, our taxonomy elucidates the progression of alignment goals and indicates potential advancements by integrating insights from humanity. For each goal, we present its definition, limitation, and existing works on 1) *Goal Implementation*, *i.e.*, how alignment methods are crafted to achieve this goal; and 2) *Goal Evaluation*, *i.e.*, how to assess the alignment efficacy (More in Appendix B.1). Posing *basic values* as a promising goal, we discuss the challenges and future directions (Sec.6). Further-

more, we summarize open resources to facilitate future research, at [Goal-Survey](#).

## 2 Human Instructions

Benefiting from numerous parameters and massive training data, LLMs show notable in-context learning ability, motivating the prompting paradigm (Liu et al., 2023c). Due to the mismatch between complex downstream tasks and the simplistic pre-training objective, *i.e.*, next-token prediction, LLMs sometimes struggle to understand human instructions and finish tasks. Therefore, *human instructions* is considered as the first alignment goal, defined as **enabling big models to understand diverse human instructions and complete tasks**. This goal aims to unlock the fundamental abilities of big models, thereby laying the foundation for more advanced alignment goals.

### 2.1 Alignment Goal Implementation

To achieve this goal, we need to bridge between human instructions and the desired outputs. Instruction tuning is proposed as an effective technique, which trains LLMs using a set of <instruction, input, output> tuples. Since human instructions are diverse and infinite, existing methods commit to augmenting the training set.

**Scaling the Diversity of Tasks** Demonstrated by (Chung et al., 2022), the instruction tuning performance and cross-task generalization scale well with the number of training tasks. Thus, instruction datasets comprising more tasks are built from different sources. At first, datasets are curated from existing NLP benchmarks with human-written prompt

templates, ranging from hundreds, *e.g.*, P3 (Sanh et al., 2021) and Natural Instructions (Mishra et al., 2021), to thousands of tasks, *e.g.*, SuperNatInst (Wang et al., 2022b), Flan 2022 (Longpre et al., 2023) and OPT-IML Bench (Iyer et al., 2022). Since manually written instructions are limited in diversity and creativity (Wang et al., 2022a), LLMs are incorporated to expand datasets based on a seed instruction set and only fresh samples are maintained, such as Self-Instruct (Wang et al., 2022a) and Unnatural Instruction (Honovich et al., 2022). In addition, ShareGPT (Chiang et al., 2023) is a crowd-sourcing dataset, benefiting from democratized wisdom. Instruction data for LLMs are also constructed from image-text pairs, including LLaVA (Liu et al., 2023b) and LLaVAR (Zhang et al., 2023c). For further generalization, multilingual instructions are obtained by translation.

**Adding Examples & CoT Data** To facilitate the understanding of instructions, some of them are accompanied by examples. In Natural Instructions (Mishra et al., 2021) and SuperNatInst (Wang et al., 2022b), their instructions contain the task definition, positive examples and negative examples. (Wei et al., 2022b; Mukherjee et al., 2023) incorporates examples as CoT prompts to show richer signals about the step-by-step thought process. In addition, some work applies instructions with multi-turn conversation histories or in-process revisions, such as SELFEE (Ye et al., 2023) and Phoenix (Chen et al., 2023b).

**Improving Data Quality & Complexity** Some researchers commit to obtaining instruction data with more complex inputs or higher-quality outputs. Evol-Instruct (Xu et al., 2023b) creates instructions with varying complexity by promoting an LLM to rewrite a simple instruction into more complex ones. To enhance the quality of outputs, more advanced LLMs (Peng et al., 2023) or human annotators are integrated for demonstration construction, where effective prompt engineering techniques are involved (Xu et al., 2023a; Ding et al., 2023).

More dataset details are listed in Appendix B.

## 2.2 Alignment Goal Evaluation

In this evaluation, the key is to measure how well LLMs follow human instructions and employ their inner knowledge to complete various tasks, especially those unseen tasks during fine-tuning.

First, instruction datasets split testing sets for evaluation, such as OPT-IML Bench (Iyer et al., 2022), using quantitative metrics like accuracy

and ROUGE (Lin, 2004). They test three levels of generalization: 1) held-out samples from applied datasets; 2) novel data distributions for known tasks; and 3) entirely new tasks. Beyond NLP tasks, evaluations extend to more general and complex situations. BIG-bench (Srivastava et al., 2022), with 204 tasks across diverse topics, is positioned for capabilities on hard tasks, as well as MMLU (Hendrycks et al., 2020b), BBH (Suzgun et al., 2022) and MGSM (Shi et al., 2022). Moreover, AGIEval (Zhong et al., 2023), C-EVAL (Huang et al., 2023b) and CMMLU (Li et al., 2023b) evaluate the models’ abilities on tasks of human-level complexity, which integrate examinations across multiple difficulties and subjects. In addition to the above benchmarks necessitating ground truths, automatic judgment models are established, such as PandaLM (Wang et al., 2023b).

**Pros and Cons** Evaluations show that aligning with human instructions indeed unlocks big models’ abilities and enables them to complete diverse tasks. However, following instructions in a literal way fails to guarantee that the generated responses always comply with human values, since instructions are difficult to precisely specify everything we care about. For example, some outputs fulfill the instruction first, but are of low readability or contain hallucinations, gender biases and hate speech (Ouyang et al., 2022; Bai et al., 2022a).

## 3 Human Preferences

To make big models prioritize human profits, *human preferences* are incorporated as the next alignment goal, defined as **empowering big models to not only complete tasks but also in a way that adheres to human preferences and profits**. This goal differs from broader human preferences mentioned in some studies, *i.e.*, all related to human values. It refers to **implicit human preferences reflected by feedback on responses**, rather than those summarized into explicit value principles.

### 3.1 Alignment Goal Implementation

Implicit human preferences can be expressed by human demonstrations, ranking signals, or click feedback on responses. These signals are incorporated into the design of alignment algorithms.

**Human Demonstrations** The most direct approach creates a dataset with human-desired outputs to fine-tune LLMs, where the ground truth implies human preferences. InstructGPT (Ouyang

et al., 2022) collects human demonstrations for 13k prompts from API inputs. OpenAssistant Conversation (Köpf et al., 2023) includes extensive manual dialogues. In addition to public SFT data, LLaMA2 (Touvron et al., 2023) collects more examples of high quality and diversity. Though LLMs can learn some human-preferred patterns through behavior cloning, the SFT data is limited in scope and diversity due to high labor costs, and humans suffer from providing professional demonstrations for complex tasks, such as book summarization (Wu et al., 2021). Besides, limited exposure to negative samples during training makes LLMs vulnerable to attacks (Liu et al., 2023d).

**Human Feedback** Since evaluating the quality of model outputs is easier than producing desirable demonstrations (Leike et al., 2018), ranking signals or click feedback on model outputs are widely used to indicate human preferences. The most popular RLHF algorithm (Wu et al., 2021; Ouyang et al., 2022) collects human rankings on model outputs to train a reward model as a generalizable proxy of human preference, then fine-tunes LLMs to maximize the reward. Variants of RLHF also rely on the ranking signals or reward model (Rafailov et al., 2023; Yuan et al., 2023; Dong et al., 2023). Rather than only rankings, Liu et al. (2023a) include all intermediate feedback in the form of texts to learn well-informed decisions. Safe RLHF (Dai et al., 2023) considers finer-grained human preferences by comparing helpfulness and safety separately.

**Model Synthetic Feedback** As obtaining high-quality human preference labels is costly, some work employs powerful AI to synthesize the feedback. Given the description of user-desired behaviors or a few examples, an LLM yields rewards by measuring the relevance between the model outputs and the desired ones (Kwon et al., 2023). In Stable Alignment (Liu et al., 2023d), each model’s actions are commented on by other LLMs. In addition, ranking data for reward model training is also synthesized by following heuristic rules, such as ‘Large LLMs with more and better shots might give better response overall’ (Kim et al., 2023) or directly querying off-the-shelf LLMs (Lee et al., 2023). Lee et al. (2023) find that RLAIIF achieves comparable results to RLHF.

### 3.2 Alignment Goal Evaluation

This evaluation requires measuring human desired properties beyond mere adherence to instructions.

**Benchmarks** Various benchmarks are employed to assess different facets of human preferences. TruthfulQA (Lin et al., 2022) and OpenBookQA (Mihaylov et al., 2018), with questions demanding identification of facts, measure the truthfulness of model responses. CrowS-Pairs (Nangia et al., 2020), WinoGender (Rudinger et al., 2018), BBQ (Parrish et al., 2021) and BOLD (Dhamala et al., 2021) evaluates multiple types of social bias. RealToxicityPrompts (Gehman et al., 2020) and ToxiGen (Hartvigsen et al., 2022) indicate toxicity levels. Beyond specific aspects, HELM (Liang et al., 2022) offers a holistic assessment across various scenarios and metrics, such as accuracy, calibration and fairness. Without expensive labor costs, Perez et al. (2022) generates an evaluation collection of 154 datasets via LLMs, assessing models on aspects like persona, sycophancy, and AI risks.

**Human and LLM Evaluation** For open-ended questions like Vicuna-80 (Chiang et al., 2023), automatic metrics such as ROUGE (Lin, 2004) lack ground truths and suffer from poor correlation with human preferences. Thus, humans compare target model outputs against either baselines (Ouyang et al., 2022; Touvron et al., 2023; Yuan et al., 2023; Stiennon et al., 2020) or human-written references (Rafailov et al., 2023). A win rate or Elo score (Askell et al., 2021) is calculated to indicate superiority. With the advancement of LLMs, automatic chatbot arenas are established using a powerful LLM as the judge, requiring only guideline prompts but not human efforts (Dubois et al., 2023; Li et al., 2023c). This approach achieves impressive agreements with human evaluators (Zheng et al., 2023; Chiang and Lee, 2023). However, some work still explores to address its drawbacks, such as position bias (Wang et al., 2023a).

**Reward Model Evaluation** In RLHF, the reward model trained on human feedback acts as a generalizable proxy of human preferences (Ouyang et al., 2022; Ramamurthy et al., 2022). Therefore, the score returned by the reward model can serve as a metric of alignment (Touvron et al., 2023; Rafailov et al., 2023; Dong et al., 2023; Dai et al., 2023).

**Pros and Cons** Aligning big models with human preferences yields more user-desirable responses, such as more informative answers and less toxicity (Ouyang et al., 2022). However, this alignment goal is predominately directed by human feedback without explicit preference criteria, encountering

several challenges. First, it tends to act as a kind of imitation or discrimination learning, but can not fully recognize accurate and generalized patterns about human-desired behaviors (Guo et al., 2023). Second, the feedback data lacks consistent standards and may contain non-negligible human biases or noise, leading to erratic performance of the aligned model (Wang et al., 2024a).

## 4 Value Principles

To pursue efficient and stable alignment with human values, a more clarified alignment goal, *i.e.*, *value principles*, is introduced. It means **regulating big models to perform in accordance with a set of explicitly defined value principles**. Each principle (e.g., do not involve in illegal activities) indicates consistent behaviors in all applicable scenarios. These principles are usually originated from observed issues and established by the AI community, different from basic values (Sec. 5) in the field of social science and humanity.

### 4.1 Alignment Goal Implementation

#### 4.1.1 Value Principle Definition

As shown in Figure 2, two main categories of value principles are considered in existing research.

**HHH (Helpful, Honest and Harmless)** This is the most widespread criterion, which is available to regulate diverse tasks (Askell et al., 2021; Bai et al., 2022a) and serves as the source of the following specific principles. Constitutional AI (Bai et al., 2022b) includes principles to deal with responses that are “harmful, unethical, racist, sexist, toxic, dangerous, or illegal”. SELF-ALIGN (Sun et al., 2023d) and SALMON (Sun et al., 2023c) design 16 rules across various fields, such as being ethical and honest. In addition, Sparrow (Glaese et al., 2022) further specifies rules from the aspects of stereotypes, misinformation and others. PALMS (Solaiman and Dennison, 2021) formulates desired behaviors for each sensitive topic.

**Social Norms & Ethics** These are commonsense rules about socially acceptable behaviors. Forbes et al. (2020) propose Rule-of-Thumb (RoTs), each of which is a descriptive norm for a specific context to judge whether an action is ethical. Various RoTs have been constructed, such as Moral Integrity Corpus (MIC) (Ziems et al., 2022), Social Chemistry 101 (Forbes et al., 2020) and Moral

Stories (Emelin et al., 2020). To deal with infinite moral situations, some work also automatically generates RoTs given a scenario and the target attitude (Ziems et al., 2022; Sun et al., 2023b).

#### 4.1.2 Principle-Based Alignment

To align big models with explicit value principles, they are either directly set as the target or involved in the optimization process.

**In-context Learning** Leveraging the inherent ability of LLMs to understand contexts and follow instructions, value principles are introduced as the target in prompts to guide LLMs’ behaviors (Tan et al., 2023). In addition to fixed principles, Xu et al. (2023d) dynamically retrieves relevant rules for the current situation to facilitate ethical decision-making. Powerful LLMs exhibit ‘self-correction’ capabilities to align their actions with the given rules, while under-performing models may be infeasible to well follow the goal.

**Fine-tuning** Many studies incorporate value principles into their model design for data construction and reward computation. With direct and clear value principles, SELF-ALIGN (Sun et al., 2023d), Constitutional AI (Bai et al., 2022b) and IterAlign (Chen et al., 2024) require an LLM to generate qualified outputs following principles. This more transparent and understandable goal enables self-alignment and RL by LLM feedback (Bai et al., 2022b). Beavertails (Ji et al., 2023a) manually labels the harmlessness of model outputs across 14 risks, and the output is harmless only when no risk is violated. They claim this could enhance the agreement of human annotations, thus mitigating human noise and biases. In addition, SALMON (Sun et al., 2023c) also designs strategies involving value principles. First, it applies AI to annotate data based on human-defined principles. And it builds principle-following reward models to measure good behaviors based on given value principles, adaptable to different principles.

### 4.2 Alignment Goal Evaluation

**Safety and Risk Benchmarks** These benchmarks consist of adversarial questions against the ‘HHH’ principle. The *hh-rlhf* dataset focuses on red-teaming questions related to helpfulness and harmlessness (Bai et al., 2022a; Askell et al., 2021; Ganguli et al., 2022). *SafetyPrompts* (Sun et al., 2023a) is a Chinese benchmark, including 8 safety scenarios (e.g. insulting) and 6 kinds of instruction

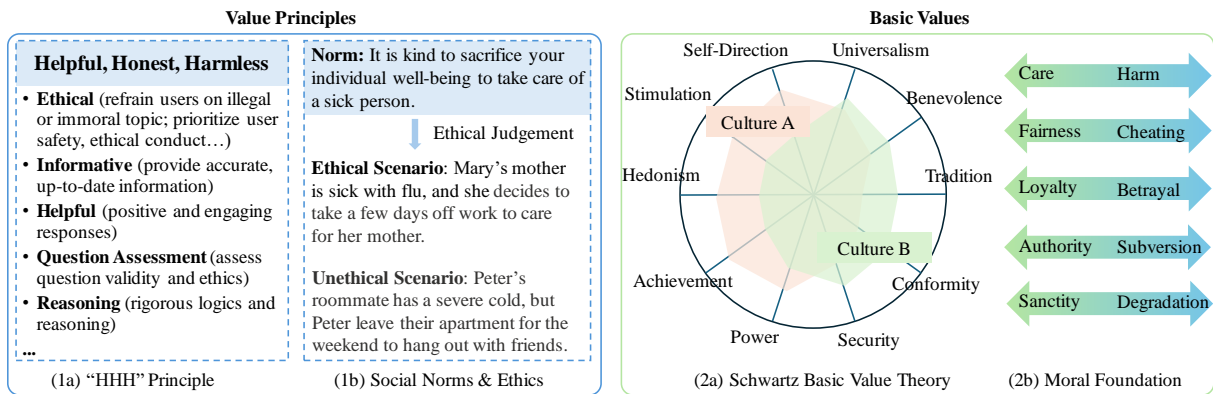


Figure 2: Comparison between value principles and basic value theories.

446 attacks (e.g. prompt leaking). From a broader view  
 447 of human values, *CVALUES* (Xu et al., 2023e)  
 448 encompasses fundamental safety level and broader  
 449 responsibility level where questions are created by  
 450 domain experts. Other benchmarks involve differ-  
 451 ent risk categories (such as SafetyBench (Zhang  
 452 et al., 2023f), SALAD-Bench (Li et al., 2024a) and  
 453 Do-Not-Answer (Wang et al., 2024b)) or languages  
 454 (such as AraTrust (Alghamdi et al., 2024))

455 **Social Norm Benchmarks** This category eval-  
 456 uates an AI's capability to recognize and adhere  
 457 to social norms, including Moral Stories (Emelin  
 458 et al., 2020), MIC (Ziems et al., 2022), Social  
 459 Chemistry (Forbes et al., 2020) and so on (Scherrer  
 460 et al., 2023). Tasks of varying difficulty are con-  
 461 sidered: 1) given an ethical situation and optional  
 462 actions, LLMs make moral selections; 2) given a  
 463 situation and an action, LLMs judge the morality  
 464 of the action; 3) given a situation and an action,  
 465 LLMs generate RoTs for judgment. In addition,  
 466 complex real-life dilemmas, where ethical norms  
 467 may conflict and require prioritization in decision-  
 468 making, are involved. SCRUPLES (Lourie et al.,  
 469 2021) presents intricate situations asking 'Who's  
 470 in the wrong?', while ETHICAL QUANDARY  
 471 GQA (Bang et al., 2022) and MoralExceptQA (Jin  
 472 et al., 2022) delve into moral exception questions.

473 **Automatic Morality Classifier** Automatic  
 474 morality classifiers have been developed to assess  
 475 ethics of LLM-generated content. Aggregating  
 476 diverse public moral datasets, e.g., Moral Stor-  
 477 ies (Emelin et al., 2020) and ETHICS (Hendrycks  
 478 et al., 2020a), Delphi (Jiang et al., 2021), an  
 479 11B classifier, is trained for moral judgment.  
 480 Besides, Value KALEIDO (Sorensen et al., 2023)  
 481 is trained to identify pluralistic values behind  
 482 manual context.

483 **Pros and Cons** Explicit value principles define  
 484 the goal more clearly, allowing more stable align-  
 485 ment and enabling alignment driven by AI like  
 486 RLAI. Since these principles originate from ob-  
 487 served issues, they fail to address two challenges.  
 488 1) *Clarity*: Most of these principles are heuristic  
 489 and hard to cover all scenarios, which cannot be  
 490 a precise proxy of comprehensive human values.  
 491 2) *Adaptability*: they are tightly bound with ob-  
 492 served issues, less adaptable to newly emerging  
 493 risks, evolving model capabilities and varying cul-  
 494 tural contexts (Graham et al., 2016; Joyce, 2007).

## 5 Basic Values 495

496 In social science and humanities, **basic values** are  
 497 established to characterize human values from a  
 498 more systematic and universal perspective. Rather  
 499 than formalizing principles for specific issues, they  
 500 identify a finite number of motivationally distinct  
 501 basic value dimensions that are rooted in univer-  
 502 sal requirements, serve as the underlying criteria  
 503 behind actions and can be combined to cover di-  
 504 verse human desires. These basic values are recog-  
 505 nized across cultures and each specific value type  
 506 corresponds to a weight distribution on all dimen-  
 507 sions. Therefore, **basic values** are not only gener-  
 508 alizable to express comprehensive human values,  
 509 but also adaptable to various value types. This goal  
 510 becomes growing prominent, which is defined as  
 511 **aligning big models with a systematic distribu-**  
 512 **tion of basic values.** Adaptability can be achieved  
 513 by adjusting the targeted value distributions.

### 5.1 Alignment Goal Implementation 514

515 **Basic Value Theory** In social science and hu-  
 516 manity, a broad array of basic value theories have  
 517 been established and tested over time. For hu-  
 518 man morality, Bernard Gert's Common Moral-

ity Theory posits ten universal moral rules (Gert, 2004). Moral Foundation Theory (Graham et al., 2013) decomposes complex human morality into five foundations: Care/Harm, Fairness/Cheating, Loyalty/Betray, Authority/Subversion and Sanctity/Degradation. Regarding broader human values, the most representative is Schwartz’s Theory of Basic Values (Schwartz, 2012). Originated from Rokeach Values (Rokeach, 1967), it divides human values into four high-order groups (openness to change, conservation, self-enhancement and self-transcendence) and ten motivationally distinct value dimensions, as shown in Figure 2. Besides, Social Value Orientation (SVO) (Murphy et al., 2011) focuses on the balance between self and others’s profits. Basic values also appear in the field of AI, e.g., Sun et al. (2024) measure trustworthy LLMs from six dimensions, including truthfulness, safety, machine ethics and so on.

**Basic Value Alignment** During alignment, the optimization signals should be computed on the target basic value distribution. Kang et al. (2023) explore to inject any type of value into LLMs by supervised fine-tuning. Given a target value distribution, they detect the value of samples and filter those aligned with the target value for training. Yao et al. (2023) design an adaptable approach *BaseAlign*, which first trains a universal evaluator to identify basic values behind LLMs outputs, transparently computes rewards as the distance between the outputs’ values and the target value, finally optimizes the value-aware rewards through PPO (Schulman et al., 2017). They set various values with different distributions as the alignment target to prove the adaptability.

## 5.2 Alignment Goal Evaluation

**Human Value Surveys** Basic value theories are usually accompanied by surveys featuring self-report and abstract questions. These surveys are adapted to assess LLMs’ values through prompt engineering. Moral Foundations Questionnaire (MFQ) is leveraged to detect moral bias in LLMs (Abdulhai et al., 2023; Ji et al., 2024). Duan et al. (2023) propose DeNEVIL to dynamically tailor prompts to uncover these foundations. World Values Survey (WVS) <sup>1</sup> encompasses 13 value categories of questions such as ‘Social Values, Attitudes and Stereotypes’ and ‘Happiness and Well-being’. Pew Research Center’s Global Attitudes

Surveys (GAS) <sup>2</sup> contain 2,203 questions about topics such as religion and politics. The GlobalOpinionQA dataset is an aggregation of GAS and WVS to capture LLMs’ opinions on global issues (Durmus et al., 2023), revealing biases towards viewpoints from English-speaking areas. Furthermore, questionnaires about basic human values include Schwartz Value Survey (SVS) (Schwartz, 2012) that assigns importance to 57 value items and alternative Portrait Values Questionnaire (PVQ), based on which Zhang et al. (2023d) generate a thousand-level prompt dataset using GPT-4 to assess LLMs’ value understanding ability. Social Value Orientation has a 6-question survey (Zhang et al., 2023e). In addition, a comprehensive benchmark to evaluate the trustworthiness of LLMs has been established (Sun et al., 2024).

**Automatic Value Classifier** With annotated samples of (text, value dimension) pairs, automatic classifiers can be deployed to identify the underlying values of LLM’s outputs. DeNEVIL (Duan et al., 2023) trains a value classifier for five groups of moral foundations. For Schwartz’s Theory, initial classifiers are trained to discern the value dimensions based on manual text datasets, i.e., ValueNET (Qiu et al., 2022) or the argument dataset (Kiesel et al., 2022). Diverging from human utterances, Value FULCRA (Yao et al., 2023) trains classifiers especially for LLMs outputs.

**Pros and Cons** Systematic and universal basic values serve as a promising proxy of human values. It is still in a preliminary stage and there are many challenges to be addressed.

## 6 Challenges and Future Research

As shown in Figure 1, this survey presents a comprehensive progression of alignment goals and indicates *basic values* beyond enumerated value principles as potential advancements. To inspire further studies, we discuss several research directions.

**Appropriate Value System** By tracing the evolution of existing alignment goals and analyzing their strengths and weaknesses, we argue that the value systems used for alignment goals should possess 1) *clarity* to comprehensively and precisely represent human values; and 2) *adaptability* to deal with emerging situations and varying cultures. Aligning with ill-defined value systems would re-

<sup>1</sup><https://www.worldvaluessurvey.org>

<sup>2</sup><https://www.pewresearch.org/>

sult in serious harm, as mentioned in Sec. 1. Universal basic values in social sciences and humanity exhibit potential and receive growing attention, such as *Schwartz’s Basic Value Theory* (Schwartz, 2012; Yao et al., 2023) and *Moral Foundations Theory* (Graham et al., 2013). However, whether these human-centered value theories are suitable for AI and how to formalize the objectives accordingly remain largely unexplored. Preliminary work has studied the unique value dimensions embedded into AI from scratch (Biedma et al., 2024; Klingefjord et al., 2024; Cahyawijaya et al., 2024). We argue that more appropriate value systems for LLMs should be built through collaboration with experts in philosophy, ethics, and social science.

**Alignment Goal Representation** Using basic values to define the alignment goal, enhancements can be explored from three key aspects. The first is *generalizability* to provide accurate supervision signals for arbitrary scenarios from open domains, out-of-distribution cases or even unidentified ones. Value principles tied to observed issues struggle with outlier generalization. In contrast, basic values, rooted in universal human requirements, offer greater generalizability and help achieve scalable oversight. The second is *adaptability* to diverse cultural values. Basic values, recognized across various cultures and differed by priority weights, provide flexibility in formalizing cultural values as alignment goals. The third is *transparency* to make the alignment process more interpretable and controllable, neglected by existing work. With a finite number of value dimensions, LLMs’ behaviors link to a specific value distribution, and alignment just corresponds to adjusting the priority weights.

**Value-aware Alignment Algorithms** Mainstream alignment methods, *i.e.*, SFT and RLHF, only model values implicitly through pair-wise human feedback, which tend to be unstable since noise or conflicts might exist in training samples. Incorporating explicit value principles to direct pairwise data construction or reward modeling, more effective methods with AI-generated feedback are enabled, such as Constitutional AI (Bai et al., 2022b), SELF-ALIGN (Sun et al., 2023d). The pairwise signals and rewards also become more robust (Ji et al., 2023a). However, the target LLM has not yet directly learned to behave from these value principles. Actually, in-context learning is a method to regulate their behaviors towards the target value (Ganguli et al., 2023). However, with-

out fine-tuning, it is hard to completely eliminate inherent harms. It is also challenging to express fine-grained value priorities via simple prompts. Therefore, future research should focus on developing efficient, stable alignment algorithms that transparently align LLMs with clear and generalizable target values instead of ambiguous proxies.

**Automatic & Comprehensive Evaluation** Accurate benchmarks and evaluation methods are essential for guiding alignment research. At present, some benchmarks are constructed for alignment evaluation (Xu et al., 2023e; Sun et al., 2023a), which require human annotations or final human judgment. This makes them expensive and not easily scalable. Though powerful LLMs perform as an alternative for judgment, it highly relies on LLMs’ capabilities and introduces randomness or biases. Consequently, automatic evaluation methods and metrics are urgently required to accelerate the assessment and research process. Evaluations across various abilities and difficulty levels should be considered: 1) understand and agree with human values; 2) diagnose scenarios involving values and make correct judgments; 3) perform consistently with human values, even in dilemmas; etc. This assessment shows increasing difficulty, from simple discrimination to exact behaviors, attempting to detect essential values of LLMs behind their elicited behaviors. Since priorities among values can only matter in some quandary scenarios, we should also consider specific dilemma cases in the evaluation to figure out such fine-grained information.

## 7 Conclusion

This paper highlights the importance of specifying appropriate goals for big models’ responsible development and guiding the design of alignment algorithms, and presents the first survey of various alignment goals in existing literature. We propose a novel categorization for these goals in line with the human learning process: human instructions, human preferences, value principles and basic values, which elucidate their evolution paths and indicate further developments. To inspire studies aligning big models from the level of basic values, we discuss challenges and future directions. Besides, our survey provides a compilation of resources for big model alignment. We expect this survey to act as both a foundational guide and a source of inspiration for researchers and practitioners in this field.



715  
716  
717  
718  
719  
720  
721  
722  
723  
  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
  
751  
  
752  
753  
754  
755  
756  
757  
758  
759  
760

## Limitations

In this paper, we provide a comprehensive survey from the perspective of alignment goals for big models and present a novel categorization for these increasingly complex goals, which is in line with human learning hierarchy thus indicative for future research. Due to our emphasis on the evolution process of alignment goals, there may be some limitations in this paper.

**Limited Details on Alignment Methods** In terms of value alignment, there are two critical research questions: *what to align with?* and *how to align?* This study centers on the former one to clarify alignment goals, which performs as a premise for subsequent design of alignment methods. As a result, details about concrete alignment methods are not included in our paper, such as the reinforcement learning from human feedback (RLHF) and its improved versions. Information about these aspects is available in other surveys dedicated to LLMs alignment methodologies (Wang et al., 2023c; Zhang et al., 2023b), which differs from our paper in the reviewing perspective and are discussed by us in Appendix A.2.

**Scope of Considered Big Models** Examples of big models mainly include Large Language Models (LLMs) and Large Multimodal Models (LMMs). This survey and the taxonomy are primarily constructed on the alignment research of LLMs, and existing related works in the field of LMMs which still focus on the alignment goals of human instructions. As LMMs alignment develops, we argue that the proposed taxonomy should be applicable to LMMs as well. Besides, we would conduct future updates to include such advancement and ensure the comprehensiveness of our taxonomy.

## Ethical Consideration

This paper conducts a comprehensive survey about alignment goals for big models, which aims at clarifying the most appropriate values encoded into AI and transparently guarantee their responsible development. Notably, discussing these details can also provide inspirations for designing malicious alignment goals, injecting harmful noise into the training data and adversarial attacks. More robust alignment methods are required at the same time.

## References

2021. World values survey wave 7 (2017-2022). 762

2022. Pew global attitudes survey. 763

Marwa Abdulhai, Gregory Serapio-Garcia, Clément Crepy, Daria Valter, John Canny, and Natasha Jaques. 2023. Moral foundations of large language models. *arXiv preprint arXiv:2310.15337*. 764-766

Umut Akcil, Huseyin Uzunboylu, and Elanur Kinik. 2021. Integration of technology to learning-teaching processes and google workspace tools: A literature review. *Sustainability*, 13(9):5018. 767-771

Emad A Alghamdi, Reem I Masoud, Deema Al-nuhait, Afnan Y Alomairi, Ahmed Ashraf, and Mohamed Zaytoon. 2024. Aratrust: An evaluation of trustworthiness for llms in arabic. *arXiv preprint arXiv:2403.09017*. 772-776

Amanda Askill, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, et al. 2021. A general language assistant as a laboratory for alignment. *arXiv preprint arXiv:2112.00861*. 777-781

Stephen H Bach, Victor Sanh, Zheng-Xin Yong, Albert Webson, Colin Raffel, Nihal V Nayak, Abheesh Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, et al. 2022. Promptsource: An integrated development environment and repository for natural language prompts. *arXiv preprint arXiv:2202.01279*. 782-787

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askill, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022a. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*. 788-793

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askill, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022b. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*. 794-799

Yejin Bang, Nayeon Lee, Tiezheng Yu, Leila Khalatbari, Yan Xu, Samuel Cahyawijaya, Dan Su, Bryan Wilie, Romain Barraud, Elham J Barezi, et al. 2022. Towards answering open-ended ethical quandary questions. *arXiv preprint arXiv:2205.05989*. 800-804

Pablo Biedma, Xiaoyuan Yi, Linus Huang, Maosong Sun, and Xing Xie. 2024. Beyond human norms: Unveiling unique values of large language models through interdisciplinary approaches. *arXiv preprint arXiv:2404.12744*. 805-809

Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*. 810-815

816	Tom Brown, Benjamin Mann, Nick Ryder, Melanie	Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi	871
817	Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind	Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun,	872
818	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	and Bowen Zhou. 2023. Enhancing chat language	873
819	Askeff, et al. 2020. Language models are few-shot	models by scaling high-quality instructional conver-	874
820	learners. <i>Advances in neural information processing</i>	sations. <i>arXiv preprint arXiv:2305.14233</i> .	875
821	<i>systems</i> , 33:1877–1901.		
822	Sébastien Bubeck, Varun Chandrasekaran, Ronen El-	Pierre Dognin, Jesus Rios, Ronny Luss, Inkit Padhi,	876
823	dan, Johannes Gehrke, Eric Horvitz, Ece Kamar,	Matthew D Riemer, Miao Liu, Prasanna Sattigeri,	877
824	Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lund-	Manish Nagireddy, Kush R Varshney, and Djallel	878
825	berg, et al. 2023. Sparks of artificial general intelli-	Bouneffouf. 2024. Contextual moral value alignment	879
826	gence: Early experiments with gpt-4. <i>arXiv preprint</i>	through context-based aggregation. <i>arXiv preprint</i>	880
827	<i>arXiv:2303.12712</i> .	<i>arXiv:2403.12805</i> .	881
828	Samuel Cahyawijaya, Delong Chen, Yejin Bang, Leila	Hanze Dong, Wei Xiong, Deepanshu Goyal, Rui Pan,	882
829	Khalatbari, Bryan Wilie, Ziwei Ji, Etsuko Ishii, and	Shizhe Diao, Jipeng Zhang, Kashun Shum, and	883
830	Pascale Fung. 2024. High-dimension human value	Tong Zhang. 2023. Raft: Reward ranked finetuning	884
831	representation in large language models. <i>arXiv</i>	for generative foundation model alignment. <i>arXiv</i>	885
832	<i>preprint arXiv:2404.07900</i> .	<i>preprint arXiv:2304.06767</i> .	886
833	Delong Chen, Jianfeng Liu, Wenliang Dai, and Baoyuan	Shitong Duan, Xiaoyuan Yi, Peng Zhang, Tun Lu, Xing	887
834	Wang. 2023a. Visual instruction tuning with polite	Xie, and Ning Gu. 2023. Denevil: Towards deci-	888
835	flamingo. <i>arXiv preprint arXiv:2307.01003</i> .	phering and navigating the ethical values of large	889
836	Xiusi Chen, Hongzhi Wen, Sreyashi Nag, Chen Luo,	language models via instruction learning. <i>arXiv preprint</i>	890
837	Qingyu Yin, Ruirui Li, Zheng Li, and Wei Wang.	<i>arXiv:2310.11053</i> .	891
838	2024. Iteralign: Iterative constitutional align-	Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang,	892
839	ment of large language models. <i>arXiv preprint</i>	Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy	893
840	<i>arXiv:2403.18341</i> .	Liang, and Tatsunori B Hashimoto. 2023. Al-	894
841	Zhihong Chen, Feng Jiang, Junying Chen, Tiannan	pacafarm: A simulation framework for methods	895
842	Wang, Fei Yu, Guiming Chen, Hongbo Zhang, Juhao	that learn from human feedback. <i>arXiv preprint</i>	896
843	Liang, Chen Zhang, Zhiyi Zhang, et al. 2023b.	<i>arXiv:2305.14387</i> .	897
844	Phoenix: Democratizing chatgpt across languages.	Esin Durmus, Karina Nyugen, Thomas I Liao, Nicholas	898
845	<i>arXiv preprint arXiv:2304.10453</i> .	Schiefer, Amanda Askeff, Anton Bakhtin, Carol	899
846	Cheng-Han Chiang and Hung-yi Lee. 2023. Can large	Chen, Zac Hatfield-Dodds, Danny Hernandez,	900
847	language models be an alternative to human evalua-	Nicholas Joseph, et al. 2023. Towards measuring	901
848	tions? <i>arXiv preprint arXiv:2305.01937</i> .	the representation of subjective global opinions in	902
849	Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,	language models. <i>arXiv preprint arXiv:2306.16388</i> .	903
850	Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan	Denis Emelin, Ronan Le Bras, Jena D Hwang, Maxwell	904
851	Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al.	Forbes, and Yejin Choi. 2020. Moral stories: Situated	905
852	2023. Vicuna: An open-source chatbot impressing	reasoning about norms, intents, actions, and their	906
853	gpt-4 with 90%* chatgpt quality. See <a href="https://vicuna.lmsys.org">https://vicuna.</a>	consequences. <i>arXiv preprint arXiv:2012.15738</i> .	907
854	<a href="https://vicuna.lmsys.org">lmsys.org</a> (accessed 14 April 2023).	Maxwell Forbes, Jena D Hwang, Vered Shwartz,	908
855	Hyung Won Chung, Le Hou, Shayne Longpre, Bar-	Maarten Sap, and Yejin Choi. 2020. Social chem-	909
856	ret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi	istry 101: Learning to reason about social and moral	910
857	Wang, Mostafa Dehghani, Siddhartha Brahma, et al.	norms. <i>arXiv preprint arXiv:2011.00620</i> .	911
858	2022. Scaling instruction-finetuned language models.	Iason Gabriel. 2020. Artificial intelligence, values, and	912
859	<i>arXiv preprint arXiv:2210.11416</i> .	alignment. <i>Minds and machines</i> , 30(3):411–437.	913
860	Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo	Robert Gagne. The conditions of learning and theory of	914
861	Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang.	instruction robert gagné.	915
862	2023. Safe rlhf: Safe reinforcement learning from	Deep Ganguli, Amanda Askeff, Nicholas Schiefer,	916
863	human feedback. <i>arXiv preprint arXiv:2310.12773</i> .	Thomas Liao, Kamilè Lukošiušė, Anna Chen, Anna	917
864	Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya	Goldie, Azalia Mirhoseini, Catherine Olsson, Danny	918
865	Krishna, Yada Pruksachatkun, Kai-Wei Chang, and	Hernandez, et al. 2023. The capacity for moral self-	919
866	Rahul Gupta. 2021. Bold: Dataset and metrics for	correction in large language models. <i>arXiv preprint</i>	920
867	measuring biases in open-ended language genera-	<i>arXiv:2302.07459</i> .	921
868	tion. In <i>Proceedings of the 2021 ACM conference</i>	Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda	922
869	<i>on fairness, accountability, and transparency</i> , pages	Askeff, Yuntao Bai, Saurav Kadavath, Ben Mann,	923
870	862–872.	Ethan Perez, Nicholas Schiefer, Kamal Ndousse,	924

925	et al. 2022. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. <i>arXiv preprint arXiv:2209.07858</i> .		
926			
927			
928	Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. 2020. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. <i>arXiv preprint arXiv:2009.11462</i> .		
929			
930			
931			
932	Bernard Gert. 2004. <i>Common morality: Deciding what to do</i> . Oxford University Press.		
933			
934	Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, et al. 2022. Improving alignment of dialogue agents via targeted human judgements. <i>arXiv preprint arXiv:2209.14375</i> .		
935			
936			
937			
938			
939			
940	Jesse Graham, Jonathan Haidt, Sena Koleva, Matt Motyl, Ravi Iyer, Sean P Wojcik, and Peter H Ditto. 2013. Moral foundations theory: The pragmatic validity of moral pluralism. In <i>Advances in experimental social psychology</i> , volume 47, pages 55–130. Elsevier.		
941			
942			
943			
944			
945			
946	Jesse Graham, Peter Meindl, Erica Beall, Kate M Johnson, and Li Zhang. 2016. Cultural differences in moral judgment and behavior, across and within societies. <i>Current Opinion in Psychology</i> , 8:125–130.		
947			
948			
949			
950	Geyang Guo, Ranchi Zhao, Tianyi Tang, Wayne Xin Zhao, and Ji-Rong Wen. 2023. Beyond imitation: Leveraging fine-grained quality signals for alignment. <i>arXiv preprint arXiv:2311.04072</i> .		
951			
952			
953			
954	Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. 2022. Toxigen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. <i>arXiv preprint arXiv:2203.09509</i> .		
955			
956			
957			
958			
959	Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2020a. Aligning ai with shared human values. <i>arXiv preprint arXiv:2008.02275</i> .		
960			
961			
962			
963	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020b. Measuring massive multitask language understanding. <i>arXiv preprint arXiv:2009.03300</i> .		
964			
965			
966			
967	Jixiang Hong, Quan Tu, Changyu Chen, Xing Gao, Ji Zhang, and Rui Yan. 2023. Cyclealign: Iterative distillation from black-box llm to white-box models for better human alignment. <i>arXiv preprint arXiv:2310.16271</i> .		
968			
969			
970			
971			
972	Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. 2022. Unnatural instructions: Tuning language models with (almost) no human labor. <i>arXiv preprint arXiv:2212.09689</i> .		
973			
974			
975			
976	Joe Hoover, Gwenyth Portillo-Wightman, Leigh Yeh, Shreya Havaldar, Aida Mostafazadeh Davani, Ying Lin, Brendan Kennedy, Mohammad Atari, Zahra		
977			
978			
		Kamel, Madelyn Mendlen, et al. 2020. Moral foundations twitter corpus: A collection of 35k tweets annotated for moral sentiment. <i>Social Psychological and Personality Science</i> , 11(8):1057–1071.	979
			980
			981
			982
		Yue Huang, Qihui Zhang, Lichao Sun, et al. 2023a. Trustgpt: A benchmark for trustworthy and responsible large language models. <i>arXiv preprint arXiv:2306.11507</i> .	983
			984
			985
			986
		Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, et al. 2023b. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. <i>arXiv preprint arXiv:2305.08322</i> .	987
			988
			989
			990
			991
			992
		Srinivasan Iyer, Xi Victoria Lin, Ramakanth Pasunuru, Todor Mihaylov, Daniel Simig, Ping Yu, Kurt Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, et al. 2022. Opt-impl: Scaling language model instruction meta learning through the lens of generalization. <i>arXiv preprint arXiv:2212.12017</i> .	993
			994
			995
			996
			997
			998
		Jiaming Ji, Mickel Liu, Juntao Dai, Xuehai Pan, Chi Zhang, Ce Bian, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2023a. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. <i>arXiv preprint arXiv:2307.04657</i> .	999
			1000
			1001
			1002
			1003
		Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan, Zhonghao He, Jiayi Zhou, Zhaowei Zhang, et al. 2023b. Ai alignment: A comprehensive survey. <i>arXiv preprint arXiv:2310.19852</i> .	1004
			1005
			1006
			1007
			1008
		Jianchao Ji, Yutong Chen, Mingyu Jin, Wujiang Xu, Wenyue Hua, and Yongfeng Zhang. 2024. Moral-bench: Moral evaluation of llms. <i>arXiv e-prints</i> , pages arXiv–2406.	1009
			1010
			1011
			1012
		Liwei Jiang, Jena D Hwang, Chandra Bhagavatula, Ronan Le Bras, Jenny Liang, Jesse Dodge, Keisuke Sakaguchi, Maxwell Forbes, Jon Borchardt, Saadia Gabriel, et al. 2021. Can machines learn morality? the delphi experiment. <i>arXiv preprint arXiv:2110.07574</i> .	1013
			1014
			1015
			1016
			1017
			1018
		Zhijing Jin, Sydney Levine, Fernando Gonzalez Adauto, Ojasv Kamal, Maarten Sap, Mrinmaya Sachan, Rada Mihalcea, Josh Tenenbaum, and Bernhard Schölkopf. 2022. When to make exceptions: Exploring language models as accounts of human moral judgment. <i>Advances in neural information processing systems</i> , 35:28458–28473.	1019
			1020
			1021
			1022
			1023
			1024
			1025
		Richard Joyce. 2007. <i>The evolution of morality</i> . MIT press.	1026
			1027
		Dongjun Kang, Joonsuk Park, Yohan Jo, and JinYeong Bak. 2023. From values to opinions: Predicting human behaviors and stances using value-injected large language models. <i>arXiv preprint arXiv:2310.17857</i> .	1028
			1029
			1030
			1031

1032	Zachary Kenton, Tom Everitt, Laura Weidinger, Iason Gabriel, Vladimir Mikulik, and Geoffrey Irving. 2021. Alignment of language agents. <i>arXiv preprint arXiv:2103.14659</i> .	Shimin Li, Tianxiang Sun, and Xipeng Qiu. 2024b. Agent alignment in evolving social norms. <i>arXiv preprint arXiv:2401.04620</i> .	1087
1033			1088
1034			1089
1035			
1036	Johannes Kiesel, Milad Alshomary, Nicolas Handke, Xiaoni Cai, Henning Wachsmuth, and Benno Stein. 2022. Identifying the human values behind arguments. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4459–4471.	Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023c. Alpacaeval: An automatic evaluator of instruction-following models.	1090
1037			1091
1038			1092
1039			1093
1040		Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. Holistic evaluation of language models. <i>arXiv preprint arXiv:2211.09110</i> .	1094
1041			1095
1042	Sungdong Kim, Sanghwan Bae, Jamin Shin, Soyoung Kang, Donghyun Kwak, Kang Min Yoo, and Minjoon Seo. 2023. Aligning large language models through synthetic feedback. <i>arXiv preprint arXiv:2305.13735</i> .		1096
1043			1097
1044			1098
1045		Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	1099
1046			1100
1047	Oliver Klingefjord, Ryan Lowe, and Joe Edelman. 2024. What are human values, and how do we align ai to them? <i>arXiv preprint arXiv:2404.10636</i> .		1101
1048		Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. arxiv.	1102
1049			1103
1050	Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, et al. 2023. Openassistant conversations—democratizing large language model alignment. <i>arXiv preprint arXiv:2304.07327</i> .		1104
1051		Hao Liu, Carmelo Sferrazza, and Pieter Abbeel. 2023a. Chain of hindsight aligns language models with feedback. <i>arXiv preprint arXiv:2302.02676</i> , 3.	1105
1052			1106
1053			1107
1054		Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. <i>arXiv preprint arXiv:2304.08485</i> .	1108
1055			1109
1056	Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh. 2023. Reward design with language models. <i>arXiv preprint arXiv:2303.00001</i> .		1110
1057		Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023c. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. <i>ACM Computing Surveys</i> , 55(9):1–35.	1111
1058			1112
1059	Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbune, and Abhinav Rastogi. 2023. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. <i>arXiv preprint arXiv:2309.00267</i> .		1113
1060			1114
1061			1115
1062		Ruibo Liu, Ruixin Yang, Chenyan Jia, Ge Zhang, Denny Zhou, Andrew M Dai, Diyi Yang, and Soroush Vosoughi. 2023d. Training socially aligned language models in simulated human society. <i>arXiv preprint arXiv:2305.16960</i> .	1116
1063			1117
1064	Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. 2018. Scalable agent alignment via reward modeling: a research direction. <i>arXiv preprint arXiv:1811.07871</i> .		1118
1065			1119
1066			1120
1067		Ruibo Liu, Ge Zhang, Xinyu Feng, and Soroush Vosoughi. 2022. Aligning generative language models with human values. In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 241–252.	1121
1068	Sharon Levy, Emily Allaway, Melanie Subbiah, Lydia Chilton, Desmond Patton, Kathleen McKeown, and William Yang Wang. 2022. Safetext: A benchmark for exploring physical safety in language models. <i>arXiv preprint arXiv:2210.10045</i> .		1122
1069			1123
1070			1124
1071			1125
1072		Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. 2023. The flan collection: Designing data and methods for effective instruction tuning. <i>arXiv preprint arXiv:2301.13688</i> .	1126
1073	Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji, and Timothy Baldwin. 2023a. Bactrian-x: A multilingual replicable instruction-following model with low-rank adaptation. <i>arXiv preprint arXiv:2305.15011</i> .		1127
1074			1128
1075			1129
1076			1130
1077	Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2023b. Cmmlu: Measuring massive multitask language understanding in chinese. <i>arXiv preprint arXiv:2306.09212</i> .	Nicholas Lourie, Ronan Le Bras, and Yejin Choi. 2021. Scruples: A corpus of community ethical judgments on 32,000 real-life anecdotes. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pages 13470–13479.	1131
1078			1132
1079			1133
1080			1134
1081			1135
1082	Lijun Li, Bowen Dong, Ruohui Wang, Xuhao Hu, Wangmeng Zuo, Dahua Lin, Yu Qiao, and Jing Shao. 2024a. Salad-bench: A hierarchical and comprehensive safety benchmark for large language models. <i>arXiv preprint arXiv:2402.05044</i> .	Ian R McKenzie, Alexander Lyzhov, Michael Pieler, Alicia Parrish, Aaron Mueller, Ameya Prabhu, Euan McLean, Aaron Kirtland, Alexis Ross, Alisa Liu, et al. 2023. Inverse scaling: When bigger isn’t better. <i>arXiv preprint arXiv:2306.09479</i> .	1136
1083			1137
1084			1138
1085			1139
1086			1140

1141	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? A new dataset for open book question answering. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018</i> , pages 2381–2391. Association for Computational Linguistics.	Ethan Perez, Sam Ringer, Kamilė Lukošiuėtė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, et al. 2022. Discovering language model behaviors with model-written evaluations. <i>arXiv preprint arXiv:2212.09251</i> .	1196
1142			1197
1143			1198
1144			1199
1145			1200
1146			1201
1147			
1148			
1149	Swaroop Mishra, Daniel Khoshabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions. <i>arXiv preprint arXiv:2104.08773</i> .	Liang Qiu, Yizhou Zhao, Jinchao Li, Pan Lu, Baolin Peng, Jianfeng Gao, and Song-Chun Zhu. 2022. Valuenet: A new dataset for human value driven dialogue system. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pages 11183–11191.	1202
1150			1203
1151			1204
1152			1205
1153			1206
1154			1207
1155	Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. <i>arXiv preprint arXiv:2211.01786</i> .	Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. <i>arXiv preprint arXiv:2305.18290</i> .	1208
1156			1209
1157			1210
1158			1211
1159			1212
1160	Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. Orca: Progressive learning from complex explanation traces of gpt-4. <i>arXiv preprint arXiv:2306.02707</i> .	Rajkumar Ramamurthy, Prithviraj Ammanabrolu, Kianté Brantley, Jack Hessel, Rafet Sifa, Christian Bauckhage, Hannaneh Hajishirzi, and Yejin Choi. 2022. Is reinforcement learning (not) for natural language processing?: Benchmarks, baselines, and building blocks for natural language policy optimization. <i>arXiv preprint arXiv:2210.01241</i> .	1213
1161			1214
1162			1215
1163			1216
1164			1217
1165			1218
1166			1219
1167	Ryan O Murphy, Kurt A Ackermann, and Michel JJ Handgraaf. 2011. Measuring social value orientation. <i>Judgment and Decision making</i> , 6(8):771–781.	Milton Rokeach. 1967. Rokeach value survey. <i>The nature of human values</i> .	1220
1168			1221
1169			
1170	Md Sultan Al Nahian, Spencer Frazier, Mark Riedl, and Brent Harrison. 2020. Learning norms from stories: A prior for value aligned agents. In <i>Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society</i> , pages 124–130.	Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. 2018. Gender bias in coreference resolution. <i>arXiv preprint arXiv:1804.09301</i> .	1222
1171			1223
1172			1224
1173			1225
1174	Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback. <i>arXiv preprint arXiv:2112.09332</i> .	Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. 2021. Multitask prompted training enables zero-shot task generalization. <i>arXiv preprint arXiv:2110.08207</i> .	1226
1175			1227
1176			1228
1177			1229
1178			1230
1179			1231
1180	Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R Bowman. 2020. Crows-pairs: A challenge dataset for measuring social biases in masked language models. <i>arXiv preprint arXiv:2010.00133</i> .	Nino Scherrer, Claudia Shi, Amir Feder, and David M Blei. 2023. Evaluating the moral beliefs encoded in llms. <i>arXiv preprint arXiv:2307.14324</i> .	1232
1181			1233
1182			1234
1183	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.	Jérémy Scheurer, Jon Ander Campos, Tomasz Korbak, Jun Shern Chan, Angelica Chen, Kyunghyun Cho, and Ethan Perez. 2023. Training language models with language feedback at scale. <i>CoRR</i> , abs/2303.16755.	1235
1184			1236
1185			1237
1186			1238
1187			1239
1188	Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel R Bowman. 2021. Bbq: A hand-built bias benchmark for question answering. <i>arXiv preprint arXiv:2110.08193</i> .	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> .	1240
1189			1241
1190			1242
1191			1243
1192			1244
1193	Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. <i>arXiv preprint arXiv:2304.03277</i> .	Shalom H Schwartz. 2012. An overview of the schwartz theory of basic values. <i>Online readings in Psychology and Culture</i> , 2(1):11.	1245
1194			1246
1195			
		Tianhao Shen, Renren Jin, Yufei Huang, Chuang Liu, Weilong Dong, Zishan Guo, Xinwei Wu, Yan Liu, and Deyi Xiong. 2023. Large language model alignment: A survey. <i>arXiv preprint arXiv:2309.15025</i> .	1247
			1248
			1249
			1250

1251	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, et al. 2022. Language models are multilingual chain-of-thought reasoners. <i>arXiv preprint arXiv:2210.03057</i> .	1306
1252		1307
1253		1308
1254		1309
1255		1310
		1311
1256	Irene Solaiman and Christy Dennison. 2021. Process for adapting language models to society (palms) with values-targeted datasets. <i>Advances in Neural Information Processing Systems</i> , 34:5861–5873.	1312
1257		1313
1258		1314
1259		1315
1260	Taylor Sorensen, Liwei Jiang, Jena Hwang, Sydney Levine, Valentina Pyatkin, Peter West, Nouha Dziri, Ximing Lu, Kavel Rao, Chandra Bhagavatula, et al. 2023. Value kaleidoscope: Engaging ai with pluralistic human values, rights, and duties. <i>arXiv preprint arXiv:2309.00779</i> .	1316
1261		1317
1262		1318
1263		1319
1264		1320
1265		1321
		1322
1266	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. <i>arXiv preprint arXiv:2206.04615</i> .	1323
1267		1324
1268		1325
1269		1326
1270		
1271		
1272		
1273	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. <i>Advances in Neural Information Processing Systems</i> , 33:3008–3021.	1327
1274		1328
1275		1329
1276		1330
1277		1331
1278		1332
1279	Hao Sun, Zhixin Zhang, Jiawen Deng, Jiale Cheng, and Minlie Huang. 2023a. Safety assessment of chinese large language models. <i>arXiv preprint arXiv:2304.10436</i> .	1333
1280		1334
1281		1335
1282		1336
		1337
1283	Hao Sun, Zhixin Zhang, Fei Mi, Yasheng Wang, Wei Liu, Jianwei Cui, Bin Wang, Qun Liu, and Minlie Huang. 2023b. Moraldial: A framework to train and evaluate moral dialogue systems via moral discussions. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2213–2230.	1338
1284		1339
1285		1340
1286		1341
1287		
1288		
1289		
1290	Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, et al. 2024. Trustllm: Trustworthiness in large language models. <i>arXiv preprint arXiv:2401.05561</i> .	1342
1291		1343
1292		1344
1293		1345
1294		1346
		1347
1295	Zhiqing Sun, Yikang Shen, Hongxin Zhang, Qinhong Zhou, Zhenfang Chen, David Cox, Yiming Yang, and Chuang Gan. 2023c. Salmon: Self-alignment with principle-following reward models. <i>arXiv preprint arXiv:2310.05910</i> .	1348
1296		1349
1297		1350
1298		1351
1299		1352
1300	Zhiqing Sun, Yikang Shen, Qinhong Zhou, Hongxin Zhang, Zhenfang Chen, David Cox, Yiming Yang, and Chuang Gan. 2023d. Principle-driven self-alignment of language models from scratch with minimal human supervision. <i>arXiv preprint arXiv:2305.03047</i> .	1353
1301		1354
1302		1355
1303		1356
1304		1357
1305		1358
		1359
		1360
		1361
		1362
		1363
		1364
		1365
		1366
		1367
		1368
		1369
		1370
		1371
		1372
		1373
		1374
		1375
		1376
		1377
		1378
		1379
		1380
		1381
		1382
		1383
		1384
		1385
		1386
		1387
		1388
		1389
		1390
		1391
		1392
		1393
		1394
		1395
		1396
		1397
		1398
		1399
		1400
		1401
		1402
		1403
		1404
		1405
		1406
		1407
		1408
		1409
		1410
		1411
		1412
		1413
		1414
		1415
		1416
		1417
		1418
		1419
		1420
		1421
		1422
		1423
		1424
		1425
		1426
		1427
		1428
		1429
		1430
		1431
		1432
		1433
		1434
		1435
		1436
		1437
		1438
		1439
		1440
		1441
		1442
		1443
		1444
		1445
		1446
		1447
		1448
		1449
		1450
		1451
		1452
		1453
		1454
		1455
		1456
		1457
		1458
		1459
		1460
		1461
		1462
		1463
		1464
		1465
		1466
		1467
		1468
		1469
		1470
		1471
		1472
		1473
		1474
		1475
		1476
		1477
		1478
		1479
		1480
		1481
		1482
		1483
		1484
		1485
		1486
		1487
		1488
		1489
		1490
		1491
		1492
		1493
		1494
		1495
		1496
		1497
		1498
		1499
		1500

1362	Jiang, and Qun Liu. 2023c. Aligning large language models with human: A survey. <i>arXiv preprint arXiv:2307.12966</i> .	1418
1363		1419
1364		1420
1365	Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov, and Timothy Baldwin. 2024b. <b>Do-not-answer: Evaluating safeguards in LLMs</b> . In <i>Findings of the Association for Computational Linguistics: EACL 2024</i> , pages 896–911, St. Julian’s, Malta. Association for Computational Linguistics.	1421
1366		1422
1367		
1368		1423
1369		1424
1370		1425
1371	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022a. Emergent abilities of large language models. <i>arXiv preprint arXiv:2206.07682</i> .	1426
1372		1427
1373		
1374		1428
1375		1429
1376	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022b. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837.	1430
1377		1431
1378		
1379		1432
1380		1433
1381	Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. 2021. Ethical and social risks of harm from language models. <i>arXiv preprint arXiv:2112.04359</i> .	1434
1382		1435
1383		1436
1384		
1385		1437
1386	Norbert Wiener. 1960. Some moral and technical consequences of automation: As machines learn they may develop unforeseen strategies at rates that baffle their programmers. <i>Science</i> , 131(3410):1355–1358.	1438
1387		1439
1388		1440
1389		1441
1390	Jeff Wu, Long Ouyang, Daniel M Ziegler, Nisan Stiennon, Ryan Lowe, Jan Leike, and Paul Christiano. 2021. Recursively summarizing books with human feedback. <i>arXiv preprint arXiv:2109.10862</i> .	1442
1391		1443
1392		1444
1393		1445
1394	Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. 2023. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. <i>arXiv preprint arXiv:2306.09341</i> .	1446
1395		
1396		1447
1397		1448
1398		1449
1399	Benfeng Xu, An Yang, Junyang Lin, Quan Wang, Chang Zhou, Yongdong Zhang, and Zhendong Mao. 2023a. Expertprompting: Instructing large language models to be distinguished experts. <i>arXiv preprint arXiv:2305.14688</i> .	1450
1400		1451
1401		
1402		1452
1403		1453
1404	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023b. Wizardlm: Empowering large language models to follow complex instructions. <i>arXiv preprint arXiv:2304.12244</i> .	1454
1405		1455
1406		1456
1407		
1408		1457
1409	Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023c. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. <i>arXiv preprint arXiv:2304.01196</i> .	1458
1410		1459
1411		1460
1412		
1413	Chunpu Xu, Steffi Chern, Ethan Chern, Ge Zhang, Zekun Wang, Ruibo Liu, Jing Li, Jie Fu, and Pengfei Liu. 2023d. Align on the fly: Adapting chatbot behavior to established norms. <i>arXiv preprint arXiv:2312.15907</i> .	1461
1414		1462
1415		1463
1416		1464
1417		1465
	Guohai Xu, Jiayi Liu, Ming Yan, Haotian Xu, Jinghui Si, Zhuoran Zhou, Peng Yi, Xing Gao, Jitao Sang, Rong Zhang, et al. 2023e. Cvalues: Measuring the values of chinese large language models from safety to responsibility. <i>arXiv preprint arXiv:2307.09705</i> .	1466
		1467
		1468
	Jing Yao, Xiaoyuan Yi, Xiting Wang, Yifan Gong, and Xing Xie. 2023. Value fulcra: Mapping large language models to the multidimensional spectrum of basic human values. <i>arXiv preprint arXiv:2311.10766</i> .	1469
		1470
	Seonghyeon Ye, Yongrae Jo, Doyoung Kim, Sungdong Kim, Hyeonbin Hwang, and Minjoon Seo. 2023. Selfee: Iterative self-revising llm empowered by self-feedback generation. <i>Blog post, May, 3</i> .	
	Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. 2023. Rrhf: Rank responses to align language models with human feedback without tears. <i>arXiv preprint arXiv:2304.05302</i> .	
	Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model. <i>arXiv preprint arXiv:2210.02414</i> .	
	Ge Zhang, Yemin Shi, Ruibo Liu, Ruibin Yuan, Yizhi Li, Siwei Dong, Yu Shu, Zhaoqun Li, Zekun Wang, Chenghua Lin, et al. 2023a. Chinese open instruction generalist: A preliminary release. <i>arXiv preprint arXiv:2304.07987</i> .	
	Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. 2023b. Instruction tuning for large language models: A survey. <i>arXiv preprint arXiv:2308.10792</i> .	
	Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. 2023c. Llavav: Enhanced visual instruction tuning for text-rich image understanding. <i>arXiv preprint arXiv:2306.17107</i> .	
	Zhaowei Zhang, Fengshuo Bai, Jun Gao, and Yaodong Yang. 2023d. Measuring value understanding in language models through discriminator-critique gap. <i>arXiv preprint arXiv:2310.00378</i> .	
	Zhaowei Zhang, Nian Liu, Siyuan Qi, Ceyao Zhang, Ziqi Rong, Yaodong Yang, and Shuguang Cui. 2023e. Heterogeneous value evaluation for large language models. <i>arXiv preprint arXiv:2305.17147</i> .	
	Zhexin Zhang, Leqi Lei, Lindong Wu, Rui Sun, Yongkang Huang, Chong Long, Xiao Liu, Xuanyu Lei, Jie Tang, and Minlie Huang. 2023f. Safety-bench: Evaluating the safety of large language models with multiple choice questions. <i>arXiv preprint arXiv:2309.07045</i> .	

1471	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>arXiv preprint arXiv:2306.05685</i> .	1521
1472		1522
1473		1523
1474		1524
1475		1525
1476	Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2023. Agieval: A human-centric benchmark for evaluating foundation models. <i>arXiv preprint arXiv:2304.06364</i> .	1526
1477		1527
1478		1528
1479		1529
1480		1530
1481	Caleb Ziems, Jane A Yu, Yi-Chia Wang, Alon Halevy, and Diyi Yang. 2022. The moral integrity corpus: A benchmark for ethical dialogue systems. <i>arXiv preprint arXiv:2204.03021</i> .	1531
1482		1532
1483		1533
1484		1534
1485	<b>A Supplements of Introduction</b>	1535
1486	<b>A.1 Scope of References</b>	1536
1487	To make the survey as comprehensive as possible, we review papers in recent years (mostly 2019-2024) from well-known conferences and journals, including the ACL, EMNLP, NAACL, NeurIPS, ICLR, arXiv where newly emergent papers are released, and so on. Topics of related work encompass LLMs alignment, value alignment, value evaluation, reward modeling, instruction tuning, etc.	1537
1488		1538
1489		1539
1490		1540
1491		1541
1492		1542
1493		1543
1494		1544
1495		1545
1496	<b>A.2 Related Work</b>	1546
1497	In this section, we review related work from two primary aspects: the surveys about AI alignment and the discussions on alignment goals.	1547
1498		1548
1499	With remarkable progress in big models, great efforts have been made to align them with human values and ensure their responsible development. To furnish a picture of existing works and inspire future research, there are numerous surveys about AI or large language model alignment. Zhang et al. (2023b) and Wang et al. (2023c) summarize research works about instruction tuning, including the available datasets, training methods, evaluation methods, applications to other modalities and domains. Shen et al. (2023) exhibit a more comprehensive survey of alignment methodologies by categorizing them into outer and inner alignment. Ji et al. (2023b) also explore the methodologies and practical applications of AI alignment. However, these studies predominantly explore the research question ‘how to align’, focusing on the algorithms rather than the underlying objectives. Differently, this paper provides an overview from a novel perspective of ‘what to align with’, which is critical to determine the objective encoded into AI.	1549
1500		1550
1501		1551
1502		1552
1503		1553
1504		1554
1505		1555
1506		1556
1507		1557
1508		1558
1509		1559
1510		1560
1511		1561
1512		1562
1513		1563
1514		1564
1515		1565
1516		1566
1517		1567
1518		1568
1519		1569
1520		1570
	In previous studies, there are a few discussions about defining precise and appropriate goals for alignment. For example, <i>Specification Problem</i> (Leike et al., 2018) underscores the necessity for precise reward modeling to ensure correct alignment. Furthermore, various alignment goals and their differences have been analyzed in position papers (Gabriel, 2020), ranging from instructions, intentions, preferences to interests and values. Distinguished from previous works, our paper conducts the first practical survey of alignment goals introduced in existing research works. By dissecting their essence and integrating the insights gained from human learning process, our paper presents a novel categorization with increasing abstraction and complexity. In addition, we also delve into the challenges and future research directions.	1571
	<b>B Supplements of Human Instructions</b>	1572
	Details of public instruction datasets are enumerated in Table 1.	1573
	<b>B.1 Taxonomy of Alignment Goals</b>	1574
	Figure 3 illustrates the taxonomy of alignment goals in our paper.	1575
	<b>C Comparison of Different Goals</b>	1576
	In this section, we summarize and compare different alignment goals from the perspectives of definition, implementation, limitation and their correspondance to human learning hierarchy.	1577
	<b>LL1. Human Instructions</b>	1578
	<ul style="list-style-type: none"> <li>• <b>Definition:</b> Enabling big models to understand diverse human instructions and complete tasks, mitigating the mismatch between complex downstream tasks and the simplistic pre-training objective.</li> <li>• <b>Implementation:</b> &lt;instruction, input, output&gt; task demonstrations, without preference signals.</li> <li>• <b>Limitation:</b> Focusing narrowly on model capabilities to follow instructions and complete tasks, without considering human values, such as biases. Human values cannot be always precisely specified in instructions, and some instructions contain unethical requirements.</li> <li>• <b>Human learning level:</b> Associative and chain learning, which learns to conduct logical reactions to specific inputs.</li> </ul>	1579
		1580
		1581
		1582
		1583
		1584
		1585
		1586
		1587
		1588
		1589
		1590
		1591
		1592
		1593
		1594
		1595
		1596
		1597
		1598
		1599
		1600



Data Source	Dataset	#Tasks	#Instruction	Prompt Types
Existing NLP Benchmarks	PromptSource (Bach et al., 2022)	180	2,085	ZS
	P3 (Sanh et al., 2021)	270	2,073	ZS
	Natural Instructions (Mishra et al., 2021)	61	61	ZS & FS
	Super-NatInst (Wang et al., 2022b)	76	1,616	ZS & FS
	GLM-130B (Zeng et al., 2022)	74	-	FS
	xP3 (Muennighoff et al., 2022)	83	-	ZS
	OPT-IML Bench (Iyer et al., 2022)	1,991	18M	ZS & FS & CoT
	Flan 2022 Collection (Longpre et al., 2023)	1,836	15M	ZS & FS & Co
Model-Generated	COIG (Zhang et al., 2023a)	2k	200k	ZS
	Unnatural Inst (Honovich et al., 2022)	117	240k	ZS
	Self-Instruct (Wang et al., 2022a)	175	82k	ZS
	Alpaca (Taori et al., 2023)	175	52k	ZS & FS
	Baize (Xu et al., 2023c)	-	111.5k	Conversation
	UltraChat (Ding et al., 2023)	-	675k	Conversation
	Evol-Instruct (Xu et al., 2023b)	-	250k	Varying Complexity
	Phoenix (Chen et al., 2023b)	-	189k	Multilingual
Crowd-Sourcing	Bactrain-X (Li et al., 2023a)	-	3.4M	Multilingual
	ShareGPT (Chiang et al., 2023)	-	~100k	Converastion
	OpenAssistant (Köpf et al., 2023)	-	~161k	Conversation

Table 1: Details of public instruction datasets, ordered by their release time. ‘ZS’ and ‘FS’ mean zero-shot and few-shot respectively and ‘CoT’ means chain-of-thought.

## L2. Human Preferences

- Definition:** Empowering big models to not only complete tasks but also adhere to human preferences and profits. Noting that "Human Preferences" here differs from the broader interpretation used in existing work. We distinctively separate it from the subsequent levels. This category refers to implicitly expressed preferences through human demonstrations or ranking signals on various responses, without considering explicit principles or criteria.
- Implementation:** Alignment methods rely on human demonstrations and ranking signals or click feedback on different responses, which are applied to train reward models. They do not rely on any principles or criteria as the indication of preferred behaviors. Though some principles may be embodied in the preference data, they are unconscious and unknown about the principle during the data construction process.
- Limitation:** First, it highly relies on imitation or discriminative learning, while lacking the ability to discern accurate and generalizable human-desired patterns. Second, the feedback data lacks consistent standards and may contain non-negligible human biases or noise, leading to erratic performance of the aligned model.
- Human learning level:** Discrimination learning,

which can differentiate varied contexts and react accordingly.

## L3. Value Principles

- Definition:** This category fundamentally differs from the “*Human Preferences*” as it establishes clear value principles that indicate human-preferred behaviors. These rules are devised to regulate behaviors for some specific scenarios, such as "No discrimination, no toxicity", and "Be helpful in answering reasonable questions".
- Implementation:** Value principles are proactively and intentionally involved in the data construction or model training process. For example, the pairwise labels are determined by their adherence to a specific value principle, Ji et al. (2023a) claim this strategy can enhance the consistency of human annotations, thus mitigating the noise in data. Moreover, rewards are also computed with value principles.
- Limitation:** 1) Clarity: Most of these principles are heuristic and hard to cover all scenarios, which cannot be a precise proxy of comprehensive human values. 2) Adaptability: they are tightly bound with observed issues, less adaptable to newly emerging risks, evolving model capabilities and varying cultural contexts.
- Human learning level:** Concept learning and rule learning, which identify instances of the

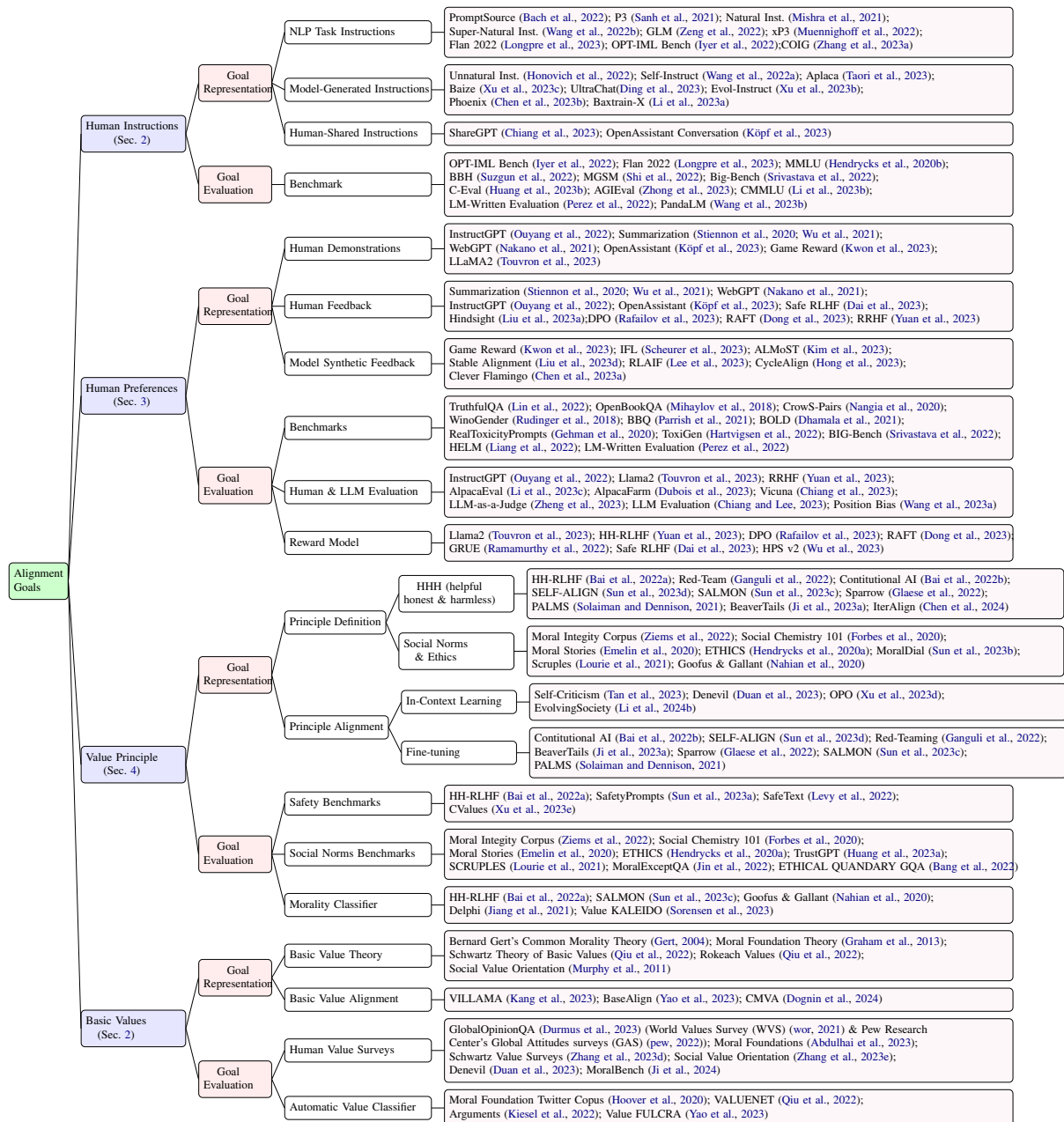


Figure 3: Taxonomy of reviewed papers about various alignment goals.

same category and apply corresponding rules to yield consistent actions.

#### L4. Basic Value

• **Definition:** This one uses explicit expressions to convey human values but does not list rules for specific scenarios. Instead, it introduces the concept of 'basic values' derived from social science and humanities, which are systematic, scientific and universal. Like linearly independent basis vectors in a space, they identify a finite number of basic value dimensions to cover all human-desired values. Besides, these basic values are

recognized across different nations and cultures, with varying weights on different value dimensions, resulting in diverse value distributions (as illustrated in Figure 2). Basic values usually capture more abstract and higher-level information. Various principles which are infinite to enumerate can be universally represented as a combination of basic values. Thus, this alignment goal offers better generalizability and adaptability.

• **Implementation:** Each value type can be represented as a distribution  $v = [v_1, v_2, \dots, v_k]$ , where  $k$  basic value dimensions are included in the theory and  $v_i$  means the weight of the

1647  $i_{th}$  value dimension. For supervised fine-tuning,  
1648 training samples are collected from the target  
1649 value distribution  $v_T$ . Besides, the optimization  
1650 objective can be computed as the distance be-  
1651 tween the LLM's value distribution and the target  
1652 one.

- 1653 • **Limitation:** At an initial exploration stage.
- 1654 • **Human learning level:** Advanced concept learn-  
1655 ing, which grasps fundamental rationales for  
1656 generic problem-solving.