# Advancing Automated Ethical Profiling in SE: a Zero-Shot Evaluation of LLM Reasoning

Patrizio Migliarini
*University of L'Aquila*
L'Aquila, Italy
patrizio.migliarini@univaq.it

Mashal Afzal Memon
*University of L'Aquila*
L'Aquila, Italy
mashal.memon@univaq.it

Marco Autili
*University of L'Aquila*
L'Aquila, Italy
marco.autili@univaq.it

Paola Inverardi
*Gran Sasso Science Institute*
L'Aquila, Italy
paola.inverardi@gssi.it

*Abstract*—**Large Language Models (LLMs) are increasingly integrated into software engineering (SE) tools for tasks that extend beyond code synthesis, including judgment under uncertainty and reasoning in ethically significant contexts. We present a fully automated framework for assessing ethical reasoning capabilities across 16 LLMs in a zero-shot setting, using 30 real-world ethically charged scenarios. Each model is prompted to identify the most applicable ethical theory to an action, assess its moral acceptability, and explain the reasoning behind their choice. Responses are compared against expert ethicists' choices using inter-model agreement metrics. Our results show that LLMs achieve an average Theory Consistency Rate (TCR) of 73.3% and Binary Agreement Rate (BAR) on moral acceptability of 86.7%, with interpretable divergences concentrated in ethically ambiguous cases. A qualitative analysis of free-text explanations reveals strong conceptual convergence across models despite surface-level lexical diversity. These findings support the potential viability of LLMs as ethical inference engines within SE pipelines, enabling scalable, auditable, and adaptive integration of user-aligned ethical reasoning. Our focus is the Ethical Interpreter component of a broader profiling pipeline: we evaluate whether current LLMs exhibit sufficient interpretive stability and theory-consistent reasoning to support automated profiling.**

*Index Terms*—**software engineering ethics, large language models, moral reasoning, zero-shot learning**

## I. INTRODUCTION

Autonomous systems are increasingly becoming an integral part of our daily lives across diverse domains [1], [2]. These systems can operate independently without any human intervention and make decisions acting on behalf of their users [3]–[6]. Their rapid growth brings both opportunities and challenges. From a software engineering perspective, as these systems become pervasive, a key challenge is designing systems that, beyond meeting technical requirements, also account for ethical considerations [7]–[11].

**SE ethics.** Recently, various studies have focused on the ethical implications of these software-intensive systems on individuals and society [10], [12]–[15]. Software engineering ethics encompasses principles and rules that guide engineers' decisions throughout the design and development process [16]. Various approaches have also been introduced that ensure that systems align with broad ethical values like fairness, transparency, and safety [17]–[22].

**Ethics operationalization.** Moreover, beyond abstract ethical norms, recent studies focus on operationalizing end users'

ethical preferences directly into the software engineering process [8], [23]–[27]. Through profiling techniques (including questionnaires, surveys, product reviews, etc.), these studies propose approaches to capture user ethical preferences, generate their ethical profiles, and integrate them into these systems [7], [28]–[31]. The ethical profile is a structured representation of the ethical preferences of the user that autonomous systems can leverage to adjust their behavior and make decisions aligned with the user's ethical values [28]. Embedding user ethical profiles into system design would not only enhance trust and accountability, as highlighted by regulatory bodies including GDPR [32], the AI Act [33], and the Ethics Guidelines for Trustworthy AI [34], but also ensure that systems reflect and respect the ethical preferences of their users. However, relying on manual input from users to generate their profiles limits their scope and adaptability, as users' ethical preferences vary with the change in context. Hence, it becomes impractical to expect users to provide input for every possible situation, introducing the challenge of automating the generation of ethical profiles.

**LLM-based ethical reasoning.** Recent advances in generative AI, especially large language models (LLMs), have positioned these models as powerful tools capable of engaging in ethical reasoning [35]. Various studies have explored the use of LLMs to assess their moral reasoning abilities in specific applications [36]–[38]. Building on this, we take a step toward evaluating whether large language models (LLMs) can effectively reason about ethically significant content in real-life scenarios. To this end, we propose a lightweight, fully automated framework that examines the potential of LLMs to identify ethically relevant information, to support the automated generation of user-aligned ethical profiles, and to integrate them into software engineering (SE) pipelines.

**Setup and RQs.** We present 16 LLMs (as shown in Table III) with 30 ethically charged statements. For each statement, the models are prompted to identify the most applicable ethical theory according to the action detailed in the statement, assess whether the action is morally acceptable according to the selected theory, and explain the reasoning behind their choice. Importantly, as discussed in Section III-A, selecting an ethical theory does not imply that the action described in the statement is justified by that theory; rather, it serves as a normative

lens through which the moral acceptability of the action is evaluated. This distinction enables meaningful binary judgments (acceptable/unacceptable) based on whether the action complies with the principles of the selected theory. To establish a comparative baseline, we replicate the same process with three professors, experts with extensive knowledge in applied ethics and philosophy. We then evaluated the responses of LLMs and experts, exploring both alignment and divergence in their judgments and their explanations. To assess the potential of LLMs as ethical reasoning modules in software engineering, we pose the following research questions:

**RQ1:** Do LLMs demonstrate the capacity for ethical reasoning when presented with ethically charged scenarios?
**RQ2:** To what extent do different LLMs agree on ethical theories and moral acceptability of the scenarios?
**RQ3:** How do the agreements among LLMs compare to those among the human experts?
**RQ4:** What qualitative characteristics emerge in the explanations produced by LLMs?

**Methodology.** These questions are designed to evaluate whether current generative models can function not only as isolated agents but also as components in robust, transparent, and auditable decision pipelines for software engineering. To address RQ1, we prompted 16 LLMs with 30 ethical scenarios covering a range of common, real-life contexts that involve ethically charged situations. We then analyzed LLMs' responses to determine their ability to recognize the action described in the scenario, its correspondence to one of the ethical theories (utilitarianism, deontology, and virtue ethics), and determine whether the action described is morally acceptable according to the selected ethical theory. To address RQ2, we computed inter-model agreement using Theory Consistency Rate (TCR) to identify the percentage of prompts for which the models selected the same ethical theory and Binary Agreement Rate metrics (BAR) to identify the percentage of prompts for which different models agreed on whether the action is morally acceptable in accordance with the selected ethical theory. These metrics assess whether different LLM models apply comparable ethical reasoning structures under identical conditions. To address RQ3, we collected questionnaires from the three professors, expert ethicists, and we presented them with the same scenarios previously shown to the LLMs. We surveyed them by replicating the same process we followed with LLMs. We then compared the experts' judgments with the responses generated by the LLMs using z-scores. To address RQ4, we conducted a multi-layered qualitative analysis of the free-text explanations provided by the LLMs. We applied lexical similarity metrics (TF-IDF and cosine similarity), dimensionality reduction (PCA, t-SNE), and topic modeling (LDA) to examine variation in linguistic form and underlying conceptual structure. A manual alignment study further assessed whether the explanations were consistent with the ethical theories selected by the models. This combined approach enabled the evaluation of both the coherence and the diversity of the explanations.

**LLM-aided ethical reasoning.** From a software engineering perspective, the results indicate that LLMs are capable of serving as modular evaluators of ethical context. They show capabilities that can be utilized as a possible way to automate the generation of user ethical profiles. This paper does not present a full automated ethical profile generator. Instead, it evaluates whether state-of-the-art LLMs can serve as ethical reasoning modules, that is, whether they can consistently select an interpretive moral lens and produce theory-aligned acceptability explanations under zero-shot conditions. Within SE pipelines, this component enables: (i) decision auditing with concise theory-grounded rationales; (ii) triage and escalation when model disagreement signals moral ambiguity; (iii) seed signals for user-aligned profiling, where interpretive patterns accumulate into dynamic profiles.

**Contributions.** The main contributions of the paper are:
- An automated framework for quantifying agreement and divergence among LLMs on non-trivial reasoning tasks, using ethical scenarios as a representative benchmark.
- An empirical analysis of the consistency and diversity of 16 LLMs, measuring both classification agreement and the qualitative variety in explanations.
- An evidence-based discussion on the trade-offs of single-LLM versus multi-LLM approaches for judgment and reasoning in SE support systems.
- A comparison of the LLM outcomes and those of expert human judgments, identifying areas of alignment and persistent divergence.
- A fully reproducible experimental pipeline and dataset, with all artifacts released for independent verification.

**Paper roadmap.** The rest of the paper is organized as follows. Section II introduces the theoretical foundations of ethical reasoning and motivates the need for automated ethical profiling in SE. Section III provides an overview of our evaluation framework and its main components. Section IV reports inter-model and human–LLM agreement metrics to assess consistency in ethical reasoning. Section V analyzes the LLM-generated explanations through lexical, conceptual, and theory-alignment perspectives. Section VI discusses the implications of our findings, identifies limitations, and outlines practical use cases. Section VII reviews related work on ethical AI, moral reasoning, and LLM evaluation. Finally, Section VIII concludes the paper and highlights directions for future research.

## II. SETTING THE CONTEXT

In this section, we introduce background concepts that form the basis of our approach and frame this work inside the umbrella project Exosoul [39] which is about protecting citizens' ethics and privacy in the digital world.

Understanding how individuals make ethical decisions in real-life scenarios is a crucial step in designing systems that can adapt to their users' ethical preferences [7]. While various approaches have proposed the use of LLMs for moral reasoning, they are primarily designed to fine-tune the models

and test their reasoning capabilities rather than assessing their inherent capacity for moral reasoning [35] [40]. Our end goal, however, is to automatically generate user ethical profiles utilizing LLMs, which reflect how individual users would act in an ethically charged situation. Hence, this involves evaluating LLMs' ethical reasoning capabilities.

**The choice of questionnaire.** To develop our approach, we build on top of the questionnaire introduced in [28] and revised in [41]. The questionnaire is introduced by a multidisciplinary group, including ethicists, philosophers, and researchers engaged in applied ethics and cognitive science, to collect ethical preferences from users in real-life situations. The questionnaire is composed of questions that reflect everyday moral dilemmas, designed to manually generate users' ethical profiles from their responses. To adapt this questionnaire for our approach, we generated declarative statements from these questions (further discussed in Section III). This translation allows us to leverage the questionnaire in a more structured format, making it suitable for evaluating whether LLMs can identify the presence of ethically charged actions within these statements. The original questionnaire and all the used statements are presented in the replication package.[1]

**The choice of theories.** Ethical considerations in the fields of computer science and software engineering have become increasingly important as technology advances [42]. The extensive use of artificial intelligence and machine learning has further made it essential to evaluate the ethical implications of such technologies [6], [43], [44]. Various ethical theories can be employed to ensure their ethical development and application [12], [42]. In this work, we evaluate whether LLMs possess ethical reasoning capabilities and identify the ethical significance of actions by evaluating their responses against three foundational moral theories: utilitarianism, virtue ethics, and deontology. The selection of these theories is based not only on their significant influence but also on their widespread adoption across both applied ethics and the broader AI and machine ethics literature [12], [17], [45]–[47]. Among the selected theories, Utilitarianism evaluates actions based on their outcomes, aiming to maximize overall well-being or happiness [48]. Virtue ethics focuses on the moral character and intentions [49]. Deontology judges actions based on adherence to moral duties or principles, regardless of the outcome [50]. The statements we used in our approach are grounded in real-life situations representing an ethical dilemma, in which each decision taken may correspond to one or more of these ethical theories, highlighting the relevance of these ethical theories to our work. In our setting, multiple theories may be applicable to the same scenario. We therefore treat the selected theory as a normative lens, an interpretive perspective used to frame the subsequent yes/no acceptability judgment. The task is not to assert the "true" or exclusive theory, but to elicit structured moral reasoning under an explicitly plural and interpretive design.

**Ethical profiling.** An ethical profile is a structured repre-

[1]https://github.com/ASE25authors/ase-aep

sentation of a user's ethical preferences, designed to reflect how the user would respond to ethically significant decisions in context-specific scenarios. These profiles are inherently dynamic and situated, evolving with the user's behavior and values. Prior studies have proposed their construction through explicit elicitation, such as surveys, questionnaires, or review analysis [28], [41], [51]–[53]. However, the continuous and manual input required from the user limits their scope and adaptability, highlighting the need to automate the generation of ethical profiles.
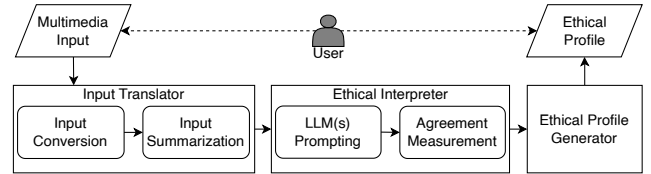


Fig. 1: Automated Ethical Profile Generation

The broader vision of project Exosoul [39] is to propose a modular approach to protect and empower individuals in asymmetric interactions with complex digital environments. The primary goal of the project is to mediate ethical, privacy-related, and social frictions through adaptive components that infer, negotiate, and operationalize the user's ethical stance in context. Figure 1 shows one technical subcomponent of the approach: a pipeline for the automated interpretation of ethically relevant content and generation of the user's ethical profile. The pipeline comprises seven modules. *Multimedia Input* collects inputs from user activity or environment, such as video, audio, and text. *Input Translator* performs *Input Conversion* to convert the inputs into symbolic representations and *Input Summarization* to condense relevant content into structured prompts. *Ethical Interpreter* classifies the summarized input by theory, acceptability, and explanation using *LLM(s) Prompting* to provide general-purpose interpretive ethical outputs, and *Agreement Measurement* compares multiple reasoning traces to assess consistency. *Ethical Profile Generator* aggregates these outputs to synthesize the user's ethical profile.

This paper concerns specifically the *Ethical Interpreter* component and aims to evaluate whether current LLMs exhibit the reasoning capability, internal coherence, and interpretive stability required to support this component. Thus, the focus of this paper is not on ethical profile generation, rather the focus is on empirically validating the core interpretive layer necessary for supporting the automation of the ethical profile generation.

## III. APPROACH OVERVIEW

We propose a fully automated, experimental pipeline to compare the ability of LLMs to reason over ethically charged scenarios in zero-shot settings, to support the automated generation of user ethical profiles in software engineering applications. Figure 2 illustrates the overview of the approach.

*Creating statements from questionnaire entries.* The process begins by adapting 30 questions – drawn from and inspired by the ethical profiling questionnaires used in the study [41]
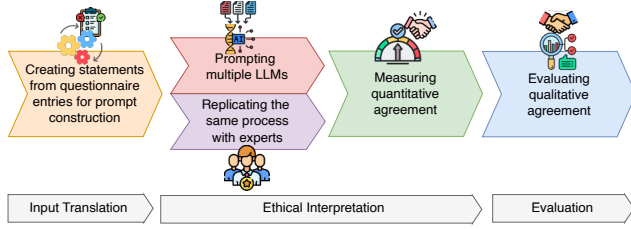
Fig. 2: Approach overview of this research

– into concise scenario statements, each presenting a specific moral dilemma and reflecting a choice that has already been made.

We provide three questionnaire items that were administered to both the LLMs and the human experts:

- A person helps a postal clerk manually assign queue numbers during a system failure.
- A user donates a small amount to Wikipedia after reading a request for funding.
- Someone finds a wallet with €1,000 and no ID and turns it in to the police.

Full corpus of items, prompts, raw responses, and expert annotations is released in the replication package.

*Prompting multiple LLMs.* Each scenario is presented as a consistently templated prompt shown in Section III-B to 16 state-of-the-art language models, spanning both proprietary and open-source families, using deterministic inference settings as specified in the LLMs documentation.

*Replicating the same process with experts.* Three expert ethicists independently respond to the same set of scenario prompts, enabling direct comparison between model and human ethical judgments.

*Measuring quantitative agreement.* Model and expert responses are analyzed using Theory Consistency Rate (TCR) and Binary Agreement Rate (BAR), alongside standardized z-score normalization, to assess inter-model and human-model consistency.

TABLE I: LLM convergence for the three illustrative items.

| Item | TCR | BAR |
|------|-----|-----|
| 1 | 43.75% (Utilitarianism) | 100.00% (YES) |
| 2 | 50.00% (Utilitarianism) | 93.75% (YES) |
| 3 | 75.00% (Deontology) | 100.00% (YES) |

TABLE II: Experts convergence for the three illustrative items.

| Item | TCR | BAR |
|------|-----|-----|
| 1 | 66.67% (Virtue ethics) | 100.00% (YES) |
| 2 | 100.00% (Utilitarianism) | 66.67% (YES) |
| 3 | 66.67% (Deontology) | 66.67% (YES) |

In Table I and Table II we provide the TCR and BAR results for the three questionnaire items provided above. The full set of results is released in the replication package.

*Evaluating qualitative agreement.* The provided explanations undergo multiple qualitative analyses, including TF-IDF-based lexical similarity, topic modeling, and clustering, to assess reasoning consistency and conceptual alignment. For each of the three items above, we provide brief explanations from both the LLMs and the human experts:

- LLM: Assists in resolving system failure, promotes efficiency
- LLM: Shows generosity and support for a valuable resource
- LLM: Acts with respect for others' property and integrity
- Expert: virtue ethics: if the act was in the nature of the individual
- Expert: you pay a little you get a lot. all theories could apply with different reasons
- Expert: no one would judge you, so it is just the way you are deontology

Full explanations set is released in the replication package.

The *Input Translation, Ethical Interpretation, and Evaluation* pipeline as a whole supports the translation of raw questionnaire data into structured prompts, automated ethical reasoning and explanation generation, and multi-level evaluation of interpretive quality and agreement across LLMs, following the pattern of the modules shown in Figure 1.

All artifacts, including prompts, responses, and evaluation scripts, are available in the replication package[1] to support reproducibility and benchmarking.

### A. Scenario Dataset and Prompt Construction

The evaluation relies on a benchmark of 30 ethically relevant scenarios, adapted from an established ethical profiling questionnaire previously used in applied ethics and cognitive science research [28], [41]. These scenarios were selected for their ability to reflect already taken decisions that are rooted in general moral dilemmas. Each scenario was expressed as a declarative statement describing a concrete action or choice, carefully phrased to minimize ambiguity while retaining enough realism to engage different dimensions of ethical reasoning. The process involved translating original survey items into concise prompts and validating their clarity through expert review. Examples of the scenarios used in the study include:

- A freelancer downloads expensive software illegally to complete an urgent work project.
  *Reflects real-world dilemmas regarding license compliance, project pressure, and ethical tool use.*
- A traveler accepts all cookies when buying a flight online.
  *Raises issues of privacy consent, dark patterns, and ethical user experience design in software systems.*
- A customer finds a USB stick in a café and plugs it into their own laptop out of curiosity.
  *Illustrates security risks and responsible behavior in handling external devices.*

Selecting a theory does not imply that the action is justified by that theory; rather, it specifies the lens relative to which the binary acceptability is evaluated. For instance, models may split between deontology and utilitarianism as the best lens for a piracy scenario, yet still converge that the action is not acceptable given either lens. This distinction is central to our evaluation design.

As detailed in subsection III-B, each scenario was presented to both LLMs and human experts using a shared prompt, structured in three parts: (i) selection of the ethical theory that best frames the action (utilitarianism, deontology, or virtue ethics), (ii) a binary judgment on whether the action is morally acceptable under that theory, and (iii) a concise explanation. Crucially, the ethical theory is treated not as a justification mechanism but as a normative lens: an interpretive perspective that informs how moral acceptability is evaluated. The binary response (YES/NO) is therefore relative to the internal logic of the selected theory, that is, whether the action adheres to or violates the principles it prescribes. For instance, the same action may be deemed unacceptable from a deontological standpoint (due to duty violation), yet potentially acceptable from a utilitarian perspective (if it maximizes positive outcomes). To illustrate this point, consider the scenario: *"A freelancer downloads expensive software illegally to complete an urgent work project."* Despite a relatively low Theory Consistency Rate (TCR) of 37.5% with LLMs split between deontology and utilitarianism 93.75% of models judged the action to be morally unacceptable. One LLM selected deontology and explained its judgment as follows: *"Disregards others' intellectual property rights and moral principles."* This case shows that identifying a theory does not entail endorsement of the action; rather, the binary judgment reflects whether the action conforms to the moral rules of that framework. As further discussed in Section V, we observed occasional misalignments between theory selection, binary judgment, and free-text explanation. These inconsistencies were explicitly analyzed to assess the internal coherence of LLM ethical reasoning.

### B. LLM Pool and Prompting Protocol

The evaluation included 16 large language models covering a broad spectrum of architectural families, release dates, and access modalities. This diversity was intended to reflect both widely deployed commercial models and state-of-the-art open-source systems available to practitioners and researchers. The pool comprises both API-accessible and locally executable models, as summarized in Table III.

TABLE III: Overview of Evaluated LLMs

| Model | Provider | Access | Parameters |
|---|---|---|---|
| GPT-4o | OpenAI | API | N/A |
| Claude 3.7 Sonnet | Anthropic | API | N/A |
| Gemini 2.5 Pro | Google | API | N/A |
| Command R+ | Cohere | API | N/A |
| Mistral | Mistral AI | API | 7B/8x22B (MoE) |
| Grok 3 | xAI | API | N/A |
| Qwen 3 | Alibaba | API | 235B |
| LLaMA 2 Chat 7B | Meta | API | 7B |
| LLaMA 3 8B | Meta | Local | 8B |
| LLaMA 3.2 3B | Meta | Local | 3B |
| DeepSeek-R1 Distill | DeepSeek | Local | 8B |
| Hermes (OpenHermes) | Nous | Local | 7B |
| Orca 2 Full | Microsoft | Local | 13B |
| Reasoner V1 | Community | Local | 7B |
| Ghost 7B v0.9.1 | Community | Local | 7B |
| Phi-3 Mini Instruct | Microsoft | Local | 3.8B |

Each model was prompted with the same set of 30 scenarios, using an identical input template and zero-shot configuration to ensure comparability. For all API-accessible models, the default deterministic setting was used (temperature set to 0.2 unless otherwise specified by the provider). Local models were run with their recommended default parameters, and all outputs were collected automatically to avoid human bias. The prompt structure presented to each model is:

```
Given the following scenario: [SCENARIO TEXT]
1) Which ethical theory best applies to this
situation: utilitarianism, deontology, or virtue
ethics?
2) Based on the theory you selected, is the action
morally acceptable (yes/no)?
3) Provide a brief explanation.
```

No additional system or context instructions were provided, and no prior examples were shown to the models, reflecting a strict zero-shot approach. All responses were collected in plain text and processed using a multi-script analysis pipeline.

Design rationale (zero-shot): we adopt a strict zero-shot setting to probe intrinsic ethical reasoning capacity under a single, shared prompt. Alternative strategies (one-/few-shot, chain-of-thought, formatting variants) are known to modulate outputs; we deliberately defer such prompt-sensitivity studies to future work to avoid conflating capacity with prompt engineering.

Determinism and parameters: for API models, we set temperature to 0.2 (or provider defaults when lower/hard-coded) to reduce sampling variance while preserving non-trivial reasoning; local models follow recommended deterministic settings. This balances repeatability and expressivity and makes inter-model comparisons fairer under identical inputs.

### IV. QUANTITATIVE ANALYSIS

For each scenario in the dataset, all 16 LLMs were queried independently using the shared prompt structure. Each response was parsed to extract the selected ethical theory, the binary acceptability judgment, and the free-text explanation. The same protocol was applied to a group of 3 human experts in ethics and applied philosophy to provide a comparative baseline. Agreement among models was quantified using two primary metrics:

*Theory Consistency Rate (TCR).* The share of models that select the modal theory for a scenario; it measures convergence of interpretive framing. TCR is introduced in this work as a deliberately simple, transparent modal-share indicator; it is not a gold-standard agreement coefficient.

*Binary Agreement Rate (BAR).* The share of models that agree on the yes/no acceptability given their selected lens; it measures outcome-level consensus.

We report z-scores for both metrics and visualize thresholds to identify high-variance (ambiguous) scenarios.

To account for scale differences, both metrics are standardized via *z-score* normalization ($z_{\text{TCR}}$ and $z_{\text{BAR}}$), where $\mu$ and $\sigma$ represent the mean and standard deviation of the respective metric across all scenarios:

$$z_{\text{TCR}} = \frac{\text{TCR} - \mu_{\text{TCR}}}{\sigma_{\text{TCR}}}, \qquad z_{\text{BAR}} = \frac{\text{BAR} - \mu_{\text{BAR}}}{\sigma_{\text{BAR}}}$$

A combined agreement score per scenario is computed as:

$$\text{Combined z-score} = \frac{z_{\text{TCR}} + z_{\text{BAR}}}{2}$$

This provides a unified measure of inter-model consistency across both theoretical and acceptability dimensions. The same procedure was applied to the human experts. All parsing and scoring scripts are included in the replication package.

**Ethical reasoning capacity (RQ1).** LLMs show ethical reasoning capacities over ethically relevant scenarios. With a 73.3% average agreement on ethical theory and 86.7% on acceptability, models show consistent normative interpretation under zero-shot conditions. Results shown in Figure 3 suggest that LLMs can produce structured, theory-informed outputs for moral scenarios, despite the lack of fine-tuning.

**LLMs agreement (RQ2).** Across all 30 scenarios, pairwise model agreement was 73.3% for TCR and 86.7% for BAR. Agreement varied across scenarios, with high divergence in ethically ambiguous situations especially those involving trade-offs between duties, rights, or risks.

| N. | Theory Consistency Rate (TCR) | Binary Agreement Rate (BAR) | Combined z-score for LLMs | Combined z-score for Experts |
|---|---|---|---|---|
| 1 | 43.75% agreement on utilitarianism | 100% agreement on YES | -0.181 | 0.378 |
| 2 | 50% agreement on utilitarianism | 93.75% agreement on YES | -0.163 | 0.234 |
| 3 | 75% agreement on deontology | 100% agreement on YES | 1.133 | -0.665 |
| 4 | 37.5% agreement on virtue ethics | 50% agreement on Tie | -2.404 | -0.522 |
| 5 | 56.25% agreement on virtue ethics | 93.75% agreement on YES | 0.1 | -0.665 |
| 6 | 75% agreement on deontology | 81.25% agreement on YES | 0.489 | 0.378 |
| 7 | 62.5% agreement on utilitarianism | 75% agreement on NO | -0.281 | 0.378 |
| 8 | 62.5% agreement on deontology | 75% agreement on YES | -0.281 | 0.234 |
| 9 | 37.5% agreement (tie on utilitarianism) | 93.75% agreement on NO | -0.689 | -0.665 |
| 10 | 68.75% agreement on utilitarianism | 81.25% agreement on YES | 0.226 | 0.378 |
| 11 | 56.25% agreement on virtue ethics | 68.75% agreement on NO | -0.788 | 0.378 |
| 12 | 43.75% agreement on deontology | 100% agreement on YES | -0.181 | -0.522 |
| 13 | 43.75% agreement on deontology | 87.5% agreement on NO | -0.615 | -0.522 |
| 14 | 68.75% agreement on deontology | 100% agreement on YES | 0.87 | 0.378 |
| 15 | 75% agreement on virtue ethics | 100% agreement on YES | 1.133 | 1.277 |
| 16 | 62.5% agreement on deontology | 100% agreement on YES | 0.607 | 0.378 |
| 17 | 68.75% agreement on deontology | 75% agreement on YES | -0.019 | -0.665 |
| 18 | 50% agreement on deontology | 62.5% agreement on NO | -1.245 | -0.522 |
| 19 | 62.5% agreement on virtue ethics | 93.75% agreement on YES | 0.362 | 0.378 |
| 20 | 75% agreement on deontology | 87.5% agreement on YES | 0.699 | -0.522 |
| 21 | 68.75% agreement on virtue ethics | 100% agreement on YES | 0.87 | 0.378 |
| 22 | 56.25% agreement on deontology | 87.5% agreement on YES | -0.09 | 0.378 |
| 23 | 37.5% agreement (tie on deontology) | 87.5% agreement on NO | -0.878 | -0.665 |
| 24 | 56.25% agreement on virtue ethics | 93.75% agreement on YES | 0.1 | -0.665 |
| 25 | 50% agreement on virtue ethics | 100% agreement on YES | 0.082 | 0.378 |
| 26 | 62.5% agreement on deontology | 93.75% agreement on YES | 0.362 | 0.378 |
| 27 | 68.75% agreement on virtue ethics | 100% agreement on YES | 0.87 | 0.378 |
| 28 | 68.75% agreement on virtue ethics | 87.5% agreement on YES | 0.435 | -0.665 |
| 29 | 50% agreement on deontology | 100% agreement on YES | 0.082 | 1.277 |
| 30 | 62.5% agreement on deontology | 93.75% agreement on YES | 0.362 | -0.665 |

Fig. 3: LLMs TCR and BAR results with Fleiss' Kappa agreement coloring and Z-Scores with threshold coloring.

Figure 3 visualizes model agreement using interpretive thresholds inspired by Fleiss' Kappa [54]. Green and yellow cells mark strong and fair agreement, while red signals interpretive divergence. This allows rapid identification of high-variance scenarios that may warrant escalation or human oversight. LLMs reach high consistency on moral acceptability and moderate-to-high consistency on ethical theory. Divergences occur in conceptually complex cases, confirming that model disagreements align with known areas of ethical ambiguity.

**LLMs vs. Human Experts (RQ3).** To compare human and model behavior, we compute the combined z-score per scenario for both groups. In several cases (e.g., scenarios numbered 14, 15, 21, 27, 29 in the replication package), LLM agreement patterns parallel the degree of expert convergence on the same scenarios, whereas sustained divergences on both sides (e.g., scenarios 4, 9, 11, 13, 17, 23) indicate intrinsic interpretive ambiguity. We therefore interpret agreement as a signal of stability, and disagreement as a cue for triage, rather than as evidence of normative correctness. Expert TCR values (often around two-thirds) should not be read as poor performance but as human-level pluralism under minimal normative instruction. This variability is a feature of the task design: it exposes multiple legitimate framings and makes explicit where automated interpretation should escalate to humans.

*Single LLM vs. Ensemble.* No single model was in perfect agreement with the modal ensemble across all scenarios. Several models exhibited idiosyncratic choices, suggesting that relying on a single LLM for ethically-sensitive decisions introduces variance and potential bias. Ensemble-based aggregation mitigates this by producing more robust and explainable outcomes, particularly in ethically ambiguous or high-stakes settings.

*Implications for Software Engineering.* The described capacity to identify scenarios with high or low model agreement is potentially actionable in SE pipelines. Scenarios with high agreement can be handled autonomously; those with low agreement can trigger alerts or escalate decisions for human review. Moreover, scenario-wise agreement scores can serve as confidence metrics to support ethical decision auditing, runtime triage in recommender or assistant systems, and automated filtering of morally unstable outputs. This sets the stage for embedding LLMs as modular ethical profilers providing scalable, explainable, and context-sensitive reasoning components in future software systems.

## V. QUALITATIVE ANALYSIS OF LLM EXPLANATIONS

The ability of an LLM to select an appropriate ethical theory or to judge an action's acceptability (as quantitatively analyzed in Section IV) is only part of what constitutes ethical reasoning. In software engineering contexts where systems must explain decisions to users, auditors, or regulators, the *explanatory layer* becomes essential. That is, this section addresses **RQ4** by complementing our quantitative agreement analysis with a qualitative investigation of the linguistic and conceptual properties of the explanations generated by LLMs. The structure, content, and coherence of such explanations directly affect the trustworthiness, transparency, and usefulness of the AI system involved. We aim to determine whether models produce morally meaningful, internally coherent, and theory-consistent explanations. We also explore how such explanations vary across models, and whether surface-level linguistic diversity masks deeper conceptual alignment. To this end, we design a multi-tiered qualitative analysis based on investigations that include lexical and syntactic diversification

of the explanations, low lexical similarity implications for conceptual divergence, consistency of explanations with the ethical theory selected by the model, moral vocabularies and normative traditions emerging across the corpus. To answer these, we combine computational techniques (TF-IDF similarity, LDA topic modeling, clustering, dimensionality reduction) with structured manual review. This dual approach balances scale and interpretability, enabling us to identify both aggregate trends and fine-grained misalignments.

*Lexical Diversity and Similarity.* We first assess lexical similarity using TF-IDF vectorization and pairwise cosine similarity, computed both across all explanations globally and within each scenario. Despite a consistent prompt structure and fixed task format, we observe substantial variation in surface form. The average pairwise cosine similarity across scenarios is 0.11, with a min of 0.02 and a max of 0.17 (Figure 6). These low values indicate that models rarely repeat phrasings, even when agreeing on the same theory and judgment. This suggests that explanations are generated with contextual variation. This result implies that LLMs are not simply giving fixed moral explanations but are capable of producing distinct, scenario-sensitive moral language. It also poses a methodological challenge; similarity metrics that rely on surface overlap will systematically underestimate conceptual agreement unless paraphrase-aware techniques are employed.

*Semantic Clustering and Model Positioning.* To understand how explanations vary semantically, we applied Principal Component Analysis (PCA) and t-distributed Stochastic Neighbor Embedding (t-SNE) to the TF-IDF vector space. These projections allow us to visualize clusters of models and identify potential outliers. As shown in Figures 4 and 5, most
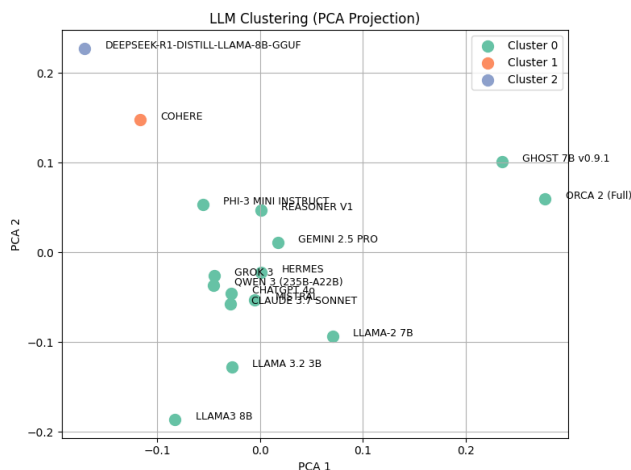


Fig. 4: PCA Clustering

models occupy a dense central region, indicating that their explanations are semantically similar at the coarse level. A few models, notably COHERE and DEEPSEEK-R1, appear in peripheral regions. Manual inspection reveals that this divergence is primarily due to stylistic verbosity or recurrent syntactic structures, rather than shifts in ethical stance. We find

no evidence that any model persistently aligns with a particular ethical theory in its language use alone. Instead, model positioning appears to reflect stylistic preferences rather than moral commitments. This further supports the hypothesis that surface diversity does not equate to conceptual inconsistency. The heatmap in Figure 7 confirms that there is no large cluster of highly similar models; most LLMs are only weakly related to each other in their use of language, specifically LLMs with the same family derivation.

*Topic Modeling and Moral Vocabulary.* We applied Latent Dirichlet Allocation (LDA) to the corpus of model explanations to uncover recurrent moral themes and to support downstream, user-directed interpretation. We used a standard preprocessing pipeline (lower-casing, tokenization, stopword removal, lemmatization, and bigram detection), built a filtered dictionary and bag-of-words corpus, and trained LDA models for a range of topic counts $k$. Model selection was guided by the coherence curve $C_v$, which improved up to $k = 12$ and then showed negligible gains, indicating that twelve topics provide an adequate balance between semantic granularity and interpretability (Figure 8). This choice is therefore empirically grounded rather than arbitrary. For each topic we inspected the highest-weight tokens and representative explanations nearest to the topic centroid to propose a short, human-readable descriptor. Crucially, this labeling is descriptive and domain/application-oriented: it illustrates how an analyst or practitioner can use our system to surface moral vocabulary and attach domain-specific meanings, rather than asserting a fixed ontology. We report some of the labeling applied to the actual clusters (Figure 9), the complete labeling table, including the token-to-label mapping and per-topic exemplar explanations, is provided in the replication package.

- 1 - Harm and character responsibility (Tokens such as *harm, protecting, character, behavior, good* indicate emphasis on avoiding harm through responsible conduct.)
- 2 - Normative framing and intervention (Presence of *deontological, virtue, privacy, safety, intervening* reflects meta-theoretical framing and protective action.)
- ...
- 12 - Justice, integrity, and solidarity (*justice, rules, duties, integrity, solidarity, courage* signals rule-guided fairness and character strength.)

Importantly, no single topic dominates the corpus. Explanations often blend multiple topics within a single rationale, indicating flexible use of moral vocabulary across theories. This pluralism is expected in ethical discourse and is compatible with our goal: the system exposes structured moral themes and their lexical supports, while leaving the interpretive labeling to the analyst's aims and domain constraints.

*Theory–Explanation Alignment.* A core requirement for ethical reasoning is that the explanation should be consistent with the moral theory being invoked. To assess this, we conducted a manual alignment check on a stratified sample of 180 responses (6 models × 10 scenarios × 3 ethical theories), manually excluding outliers. In over 90% of cases, the expla-

Fig. 5: t-SNE Clustering
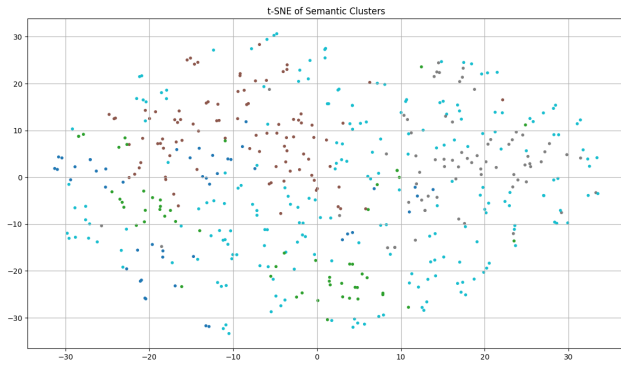


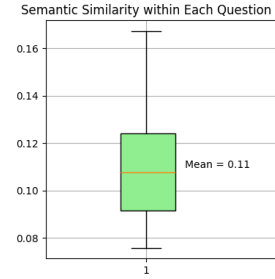Fig. 6: Similarity per question



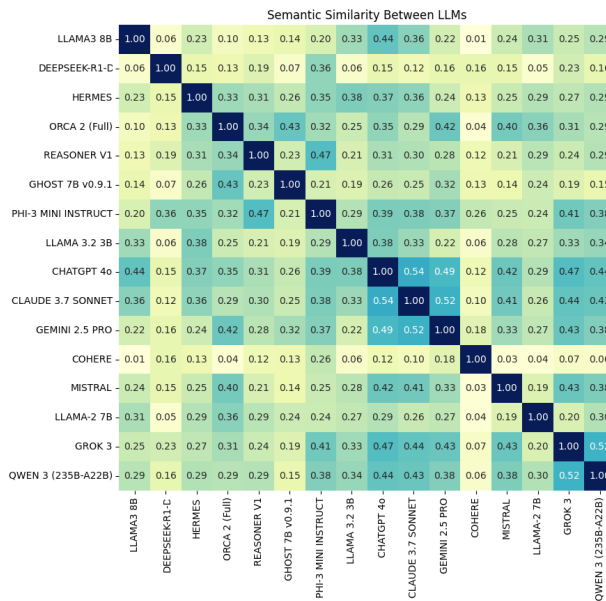Fig. 7: Semantic similarity heatmap



Fig. 8: LDA Coherence Score



Fig. 9: LDA topic excerpt. The full image is available in the replication package.

nation supported the selected theory in a coherent manner. For example, Deontological responses often cited duties, rules, or rights (e.g., "The action violates the duty of confidentiality."), Utilitarian responses referenced consequences and well-being (e.g., "The outcome benefits more people."), and Virtue-based explanations appealed to character or intention (e.g., "It reflects compassion and honesty."). Misalignments, when they occurred, tended to arise in edge cases or procedurally complex scenarios. Occasionally, a model labeled a decision as "virtue ethics" but explained it in outcome-based terms. Mismatches were rare and not concentrated in specific models.

*Conciseness and Structural Patterns.* We analyzed sentence and word counts across all explanations to understand the
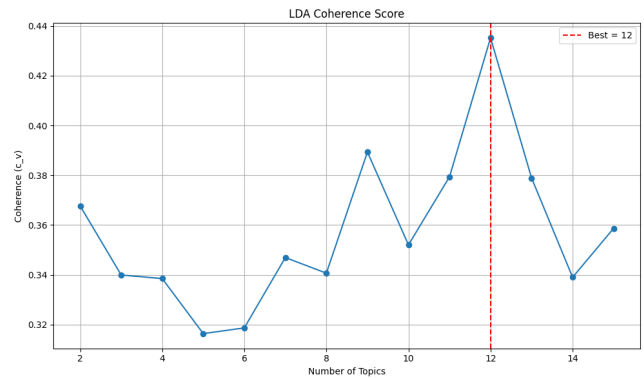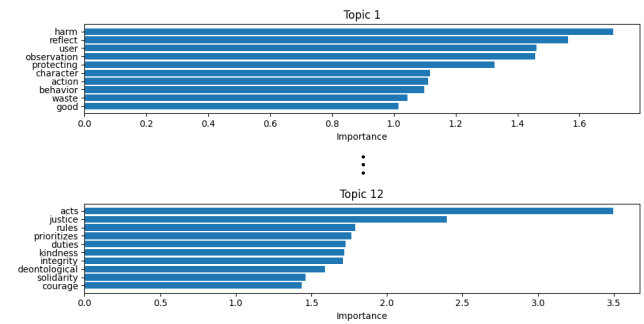
structural economy of model explanations. Over 95% of outputs were a single sentence. Mean word counts ranged from 7.7 to 21.2 tokens (Figures 10–11). Some models (e.g., GHOST, DEEPSEEK) tended to be more verbose, but longer explanations did not correlate with stronger theory alignment or clarity. This brevity is notable; despite being concise, most explanations successfully reference relevant moral principles. The findings suggest that LLMs are able to generate ethical
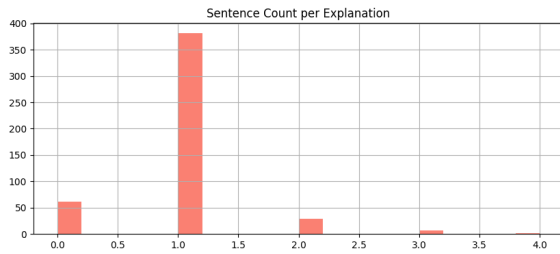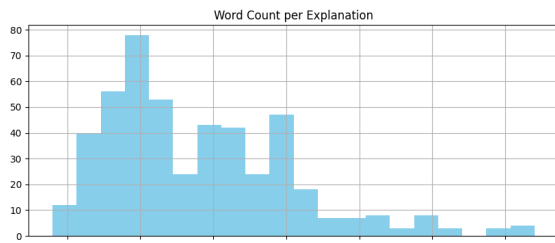
Fig. 10: Sentence counts


Fig. 11: Word counts

explanations that are both succinct and substantively meaningful, a valuable property for deployment in constrained SE contexts, such as UI prompts or trace logs.

*Explanatory Convergence and Disagreement.* In high-


Fig. 12: Word cloud

consensus scenarios (e.g., helping a clerk, donating to Wikipedia), all models agreed both on judgment and theory, and their explanations while lexically distinct shared moral themes like civic virtue, altruism, or public benefit. Word cloud analysis (Figure 12) confirms recurring terms such as "respect", "action", "duty", "virtue", and "harm". In low-agreement scenarios, such as those involving conflicting obligations or ambiguous responsibility, we observed both linguistic and conceptual spread. Some models emphasized personal integrity, others systemic outcomes. However, even here, the diversity appeared to reflect the ambiguity of the scenario rather than noise or error.

*Implications for Ethical Profiling and SE Systems.* These findings have direct implications for the use of LLMs in SE applications that require ethical reasoning or user profiling:

- *Trust and Transparency.* The presence of coherent, interpretable explanations supports use cases requiring traceable decision-making (e.g., user-facing ethical explanations or compliance logging).
- *Profile Construction.* The diversity and flexibility in moral vocabulary provide rich semantic data for generating dynamic ethical profiles based on user input.
- *Model Assessment.* Surface similarity metrics may underestimate performance; systems should instead incorporate theory-alignment checks or conceptual similarity measures (e.g., based on semantic entailment).
- *Ethical Memory.* Explanatory patterns may be used to build cumulative ethical profiles over time, enabling systems to adapt to evolving user values without retraining.

**Qualitative characteristics (RQ4).** From what emerges from our analysis, the LLM-generated explanations for ethical decisions are marked by high surface diversity and consistent moral coherence. Despite low lexical overlap across outputs, with average pairwise similarity scores around 0.11, the models consistently construct explanations that align with the selected ethical theory in over 90% of sampled cases. In line with what has been argued so far, this also suggests that models operate beyond rigid theoretical frameworks, producing context-sensitive explanations. Dimensionality reduction and clustering analyses (PCA, t-SNE) show that most models occupy a dense central region in semantic space, suggesting convergence at the conceptual level, even as their syntactic forms diverge. Outliers (e.g., COHERE, DEEPSEEK) diverge primarily due to stylistic rather than ethical differences. Topic modeling reveals that LLMs implicitly draw from a pluralistic moral vocabulary, spanning utilitarian, deontological, and virtue-based notions. Most explanations combine references to multiple moral principles within a single response, confirming the models' ability to construct hybrid, plausible explanations. Structurally, over 95% of explanations are a single sentence, and brevity does not preclude moral relevance. Even short explanations tend to highlight key ethical features (e.g., duty, harm, empathy), making them well-suited for constrained SE contexts. Finally, in high-consensus scenarios, lexical variation coexists with thematic convergence around shared moral terms (e.g., "altruism", "civic duty", "harm"), while disagreement in complex cases reflects underlying ambiguity rather than interpretive noise. These findings support the use of LLMs for both traceable decision-making and dynamic ethical profiling in software engineering systems.

## VI. DISCUSSION

**RQ1.** Our findings provide evidence that state-of-the-art LLMs can engage in ethical reasoning when presented with complex, real-world made explanations of acceptability. Without fine-tuning or examples, models consistently identified the most applicable ethical theory and made acceptability explanations with substantial inter-model agreement. This capability suggests that LLMs possess an implicit grasp of moral reasoning principles, grounded in their pretraining on large-scale textual corpora. From an SE perspective, this opens

the door to using LLMs as ethical reasoning modules in decision-making pipelines, such as requirement negotiation, user modeling, or system auditing.

**RQ2.** Quantitative results showed that models converge more strongly on binary moral acceptability (86.7% BAR) than on ethical theory classification (73.3% TCR). While this difference reflects the higher abstraction level of theoretical judgments, the level of agreement observed is non-trivial. The scenario-dependent variability in TCR reveals an important feature: model disagreement tends to reflect ethical ambiguity inherent in the scenario rather than arbitrary noise. This suggests that ensemble disagreement can be used as a proxy for moral uncertainty, enabling software systems to trigger escalation or human intervention when LLMs disagree sharply.

**RQ3.** Overall, LLMs exhibit non-trivial agreement with experts that is more pronounced in prevalent classes and weaker on rare or edge cases. Scenarios that elicited strong agreement among experts tended to also show high inter-model LLM agreement, and vice versa. This convergence reinforces the reliability of LLMs in interpreting familiar or structurally simple moral scenarios. Divergences, especially in edge cases, underscore the importance of hybrid systems that combine automated reasoning with human oversight. For SE applications involving legal, regulatory, or safety-critical implications, LLM-based profiling should not be deployed as an isolated decision-maker but as a complementary module.

**RQ4.** Qualitative analyses demonstrated that LLMs generate explanations that are lexically diverse but conceptually coherent. Despite low textual similarity across models, explanations consistently aligned with the chosen moral theory in over 90% of cases. Models blended terminology from multiple ethical traditions in natural, context-sensitive ways, tending to reflect how human reasoners combine principles, consequences, and character-based considerations. This expressive flexibility is critical for ethical profiling, as it enables the detection of user-aligned reasoning patterns across different moral framings. Moreover, the compactness of most explanations (single sentences) and their theoretical consistency suggest that LLMs are capable of producing tractable, auditable moral outputs suitable for runtime interpretation and logging.

> **Agreement ≠ correctness.** Our design surfaces stability and ambiguity signals, it does not certify normative accuracy. In SE practice, high agreement supports automation with audit, while low agreement recommends human-in-the-loop escalation.

**Limitations and Scope of Validity.** While promising, our findings are bounded by some limitations:

*Theoretical coverage.* We focus on three major ethical theories utilitarianism, deontology, and virtue ethics due to their widespread adoption in software engineering practice and education [55]. These ethical theories provide well-established foundations for analyzing ethical dilemmas in technology contexts. While alternative theories such as care ethics or contractualism are less commonly applied, they offer valuable perspectives that could enrich ethical analyses. Future work may explore the integration of these additional frameworks to capture a broader spectrum of moral reasoning in software engineering.

*Scenario framing.* Our prompts use concise, decontextualized scenarios. Richer formats (e.g., dialogues, system logs) may affect model interpretation. Our current prompts are decontextualized statements. An important next step is to apply the same ethical reasoning pipeline (Figure 2) to richer input modalities, including: (i) chat transcripts from developer-agent interactions; (ii) logs of user decisions in ethically sensitive configurations; (iii) behavioral signals from simulation environments or system telemetry. This would move the profiling process closer to real-time, context-aware ethical inference.

*Zero-shot constraints.* All reasoning is performed without memory or clarification. Interactive or multi-turn reasoning may yield different profiles. An ethical profile need not be static. As users interact with a system, their decisions may reveal shifts in priorities, trade-offs, or ethical boundaries. Future work should implement an ethical memory module that incrementally updates a user's profile over time, capturing both stable dispositions and contextual shifts. This requires designing a temporal profiling architecture that tracks ethical indicators across scenarios and resolutions.

*Agreement ≠ correctness.* Convergence does not imply normative accuracy. Human biases and model alignment may coincide but remain ethically questionable. In real deployments, users may reject or revise the moral judgments made by the system. Building on our current architecture, we envision an interactive loop in which: (i) the system proposes an ethical explanation; (ii) the user confirms, modifies, or rejects the reasoning; (iii) the profile is updated accordingly. This would enable both user agency and model refinement over time, reducing the risk of misaligned ethical personalization.

*Prompt sensitivity.* Our zero-shot, single-turn protocol deliberately controls for instruction complexity; however, model behavior can still be sensitive to seemingly innocuous variations in prompt phrasing, formatting, or input length. We therefore treat prompt sensitivity as a threat to validity and an explicit boundary of our claims. A systematic sensitivity analysis is left as future work. In practice, we recommend freezing prompt templates in repositories and reporting all formatting details that might affect reproducibility.

Even if the findings support the potential viability, these limitations suggest caution in direct deployment and highlight the need for further validation before integrating LLM-based profiling into high-stakes SE systems. Our results position LLMs as viable components for modular ethical reasoning in SE. Possible use cases include: *decision auditing* for moral rationales generation for SE tool outputs (e.g., in requirements prioritization or resource allocation); *autonomy triage* to route decisions to humans when LLMs disagree, reducing risk in ethically charged contexts; *agent personalization* to tailor behavior of autonomous SE agents based on learned ethical

user profiles. More broadly, the ability to extract consistent moral structure from language enables a shift from static ethics-as-checklists to adaptive, traceable, and user-aligned ethical cognition in engineered systems.

## VII. Related Work

*Ethics in autonomous and software-intensive systems.* The integration of ethical considerations into autonomous systems has been widely examined in software engineering. Prior work spans design-time approaches that encode ethics via codes of conduct, principles, and rules [1], [10], [21], [56], [57], and verification-time approaches that formalize and check system decisions against ethical frameworks [2], [8], [13], [14], [58]–[61]. This body of research establishes both the need and the mechanisms for ensuring ethically compliant behavior in autonomous decision pipelines.

*LLMs in software engineering and the need for ethical reasoning.* Concurrently, LLMs have been adopted across SE tasks such as code generation, bug detection, and documentation [62]. As these models are integrated into SE workflows, it becomes important to assess whether they exhibit ethical reasoning and how they might be leveraged in systems with ethical implications [35], [37].

*Model-centric evaluations of moral reasoning.* Han et al. [35] evaluate LLM understanding of moral/ethical reasoning across five domains (justice, deontology, virtue ethics, utilitarianism, and commonsense morality). They fine-tune BERT-base, BERT-large, RoBERTa-large, and ALBERT-xxlarge on ETHICS datasets comprising over 13,000 scenarios, and then test the models' ability to classify scenarios in line with the ethical theories used during fine-tuning. In contrast, the present work does not rely on fine-tuning; it evaluates whether pre-trained LLMs, in a zero-shot setup, display ethical reasoning.

*Value identification in scholarly texts vs. ethically charged situations.* A complementary thread leverages LLMs to extract values and structure from scientific writing. For example, [40] studies ChatGPT's ability to identify human values from titles and abstracts of SE publications using Schwartz's theory [63], followed by manual verification by humans. Related efforts show that LLMs can assist in extracting insights from scholarly texts, classifying publications, and mining metadata for literature reviews [36], [62]. These studies focus on value detection in academic prose rather than on evaluating ethical concrete reasoning, ethically charged scenarios.

*Reasoning with policies and learning moral rewards.* Rao et al. [37] argue against hard-coding specific moral values in LLMs and advocate for general ethical reasoning capacities that adapt to diverse contexts. They introduce in-context ethical policies defined at varying abstraction levels and grounded in deontology, virtue ethics, and consequentialism, reporting experiments across GPT-3, ChatGPT, and GPT-4 with GPT-4 showing stronger ethical reasoning. Tennant et al. [38] incorporate intrinsic moral rewards, grounded in deontological and utilitarian theories, into reinforcement learning fine-tuning. Using the Iterated Prisoner's Dilemma, they demonstrate that LLM agents can learn morally aligned strategies and even unlearn previously selfish behaviors. While these works shape model behavior through policies or moral rewards, the present study compares outputs across multiple models without additional tuning, treating ethical reasoning as a testing ground rather than a direct optimization target.

*Positioning of the present study.* Prior research establishes design/verification-time mechanisms for ethics in autonomous systems, documents LLM utility in SE, and explores both fine-tuned moral reasoning and in-context ethical policy use. The contribution here is orthogonal: a zero-shot assessment of pre-trained LLMs' ethical reasoning on ethically charged scenarios, contrasting with fine-tuned or policy-conditioned settings, and distinct from value-mining in academic text.

## VIII. Conclusion and Future Work

Our work investigates the potential of leveraging LLMs as an ethical component within the software engineering pipeline to automate the generation of user ethical profiles. The profiles represent users' ethical preferences in a structured way, guiding system behavior to align its decisions with user values. To achieve this, we evaluated the ethical reasoning capabilities of 16 state-of-the-art Large Language Models (LLMs) by presenting them with 30 ethically charged scenarios. We prompted LLMs to identify which ethical theory among utilitarianism, deontology, and virtue ethics applies to the action detailed in the scenario, determine whether the action described is morally acceptable, and provide the reasoning behind their choice. We then computed inter-model agreement between the responses of the 16 LLMs using Theory Consistency Rate (TCR) and Binary Agreement Rate (BAR) metrics. We replicated the same process with three expert ethicists and compared their responses with LLMs using z-scores to analyze the agreement between their responses. Moreover, we performed qualitative analysis of the LLM explanations to determine whether models produce meaningful and theory-consistent explanations. The results indicate LLMs' applicability within the software engineering pipeline to evaluate ethical contexts. Future work will prompt LLMs to identify actions based on a broader set of ethical theories, and will examine the applicability and interactions of these theories in more complex scenarios. Moreover, we plan to investigate how this framework could be integrated in existing SE toolchains.

## REFERENCES

[1] S. Suri, S. N. Das, K. Singi, K. Dey, V. S. Sharma, and V. Kaulgud, "Software engineering using autonomous agents: Are we there yet?," in *38th IEEE/ACM International Conference on Automated Software Engineering*, pp. 1855–1857, 2023.

[2] A. Jedlickova, "Ensuring ethical standards in the development of autonomous and intelligent systems," *IEEE Transactions on Artificial Intelligence*, vol. 5, no. 12, pp. 5863–5872, 2024.

[3] P. Pelliccione and N. Laranjeiro, "Insights from the software reliability research community," *IEEE Reliability Magazine*, vol. 1, no. 1, pp. 10–14, 2024.

[4] A. E. Waldman, "Power, process, and automated decision-making," *Fordham L. Rev.*, vol. 88, p. 613, 2019.

[5] A. Sharma, V. Sharma, M. Jaiswal, H.-C. Wang, D. N. K. Jayakody, C. M. W. Basnayaka, and A. Muthanna, "Recent trends in ai-based intelligent sensing," *Electronics*, vol. 11, no. 10, p. 1661, 2022.

[6] J. Anderson, L. Rainie, and A. Luchsinger, "Artificial intelligence and the future of humans," *Pew Research Center*, vol. 10, no. 12, pp. 1–10, 2018.

[7] M. Autili, M. De Sanctis, P. Inverardi, and P. Pelliccione, "Engineering digital systems for humanity: a research roadmap," *ACM Transactions on Software Engineering and Methodology*, 2025.

[8] M. De Sanctis and P. Inverardi, "Engineering ethical-aware collective adaptive systems," in *International Symposium on Leveraging Applications of Formal Methods*, pp. 238–252, Springer, 2024.

[9] D. Hendrycks, C. Burns, S. Basart, A. Critch, J. Li, D. Song, and J. Steinhardt, "Aligning ai with shared human values," in *9th International Conference on Learning Representations*, 2021.

[10] R. Alidoosti, P. Lago, E. Poort, M. Razavian, and A. Tang, "Incorporating ethical values into software architecture design practices," in *19th International Conference on Software Architecture Companion*, pp. 124–127, 2022.

[11] J. Svegliato, S. B. Nashed, and S. Zilberstein, "Ethically compliant sequential decision making," in *AAAI Conference on Artificial Intelligence*, vol. 35, pp. 11657–11665, 2021.

[12] S. Tolmeijer, M. Kneer, C. Sarasua, M. Christen, and A. Bernstein, "Implementations in machine ethics: A survey," *ACM Computing Surveys*, vol. 53, no. 6, pp. 1–38, 2020.

[13] P. Inverardi, "Ethics and privacy in autonomous systems: A software exoskeleton to empower the user," in *Software Engineering for Resilient Systems: 11th International Workshop.*, pp. 3–8, 2019.

[14] P. Inverardi, M. Palmiero, P. Pelliccione, and M. Tivoli, "Ethical-aware autonomous systems from a social psychological lens.," in *6th International Workshop on Cultures of Participation in the Digital Age: AI for Humans or Humans for AI?*, pp. 43–48, 2022.

[15] S. Cervantes, S. López, and J.-A. Cervantes, "Toward ethical cognitive architectures for the development of artificial moral agents," *Cognitive systems research*, vol. 64, pp. 117–125, 2020.

[16] R. Alidoosti, P. Lago, M. Razavian, and A. Tang, "Exploring the ethical landscape of software systems: A systematic literature review," *Journal of Systems and Software*, p. 112430, 2025.

[17] A. Jedličková, "Ethical approaches in designing autonomous and intelligent systems: a comprehensive survey towards responsible development," *AI & SOCIETY*, pp. 1–14, 2024.

[18] P. Bremner, L. A. Dennis, M. Fisher, and A. F. Winfield, "On proactive, transparent, and verifiable ethical reasoning for robots," *Proceedings of the IEEE*, vol. 107, no. 3, pp. 541–561, 2019.

[19] A. F. Winfield, C. Blum, and W. Liu, "Towards an ethical robot: internal models, consequences and ethical action selection," in *Conference towards autonomous robotic systems*, pp. 85–96, 2014.

[20] A. F. T. Winfield and V. V. Hafner, "Anticipation in robotics," *Handbook of Anticipation: Theoretical and Applied Aspects of the Use of Future in Decision Making*, pp. 1587–1615, 2019.

[21] R. Alidoosti, "Ethics-driven software architecture decision making," in *18th International Conference on Software Architecture Companion*, pp. 90–91, 2021.

[22] B. Townsend, C. Paterson, T. Arvind, G. Nemirovsky, R. Calinescu, A. Cavalcanti, I. Habli, and A. Thomas, "From pluralistic normative principles to autonomous-agent rules," *Minds and Machines*, vol. 32, no. 4, pp. 683–715, 2022.

[23] E. Winter, S. Forshaw, L. Hunt, and M. A. Ferrario, "Advancing the study of human values in software engineering," in *2019 IEEE/ACM 12th International Workshop on Cooperative and Human Aspects of Software Engineering (CHASE)*, pp. 19–26, IEEE, 2019.

[24] M. A. Memon, G. L. Scoccia, and M. Autili, "Automated negotiation-preliminary results of a systematic mapping study," in *2023 38th IEEE/ACM International Conference on Automated Software Engineering Workshops (ASEW)*, pp. 94–99, IEEE, 2023.

[25] M. A. Memon, M. Autili, G. Filippone, G. L. Scoccia, and P. Inverardi, "A high-level architecture of an automated context-aware ethics-based negotiation approach," in *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering, (ASE)*, ACM, 2024.

[26] D. Donati, Z. Assadi, S. Gozzano, P. Inverardi, and N. Troquard, "On representing humans' soft-ethics preferences as dispositions," in *Proceedings of the Ital-IA Intelligenza Artificiale*, vol. 3762 of *CEUR Workshop Proceedings*, pp. 135–140, 2024.

[27] B. S. Barn, "Do you own a volkswagen? values as non-functional requirements," in *Human-Centered and Error-Resilient Systems Development*, pp. 151–162, Springer, 2016.

[28] C. Alfieri, P. Inverardi, P. Migliarini, and M. Palmiero, "Exosoul: Ethical profiling in the digital world," in *HHAI2022: Augmenting Human Intellect*, pp. 128–142, IOS Press, 2022.

[29] P. Inverardi, "The european perspective on responsible computing," *Commun. ACM*, vol. 62, p. 64, mar 2019.

[30] P. Inverardi, P. Migliarini, and M. Palmiero, "Systematic review on privacy categorisation," *Computer Science Review*, vol. 49, p. 100574, 2023.

[31] D. Di Ruscio, P. Inverardi, P. Migliarini, and P. T. Nguyen, "Leveraging privacy profiles to empower users in the digital society," *Automated Software Engineering*, vol. 31, no. 1, p. 16, 2024.

[32] P. Voigt and A. Von dem Bussche, "The eu general data protection regulation (gdpr)," *A Practical Guide, 1st Ed., Cham: Springer International Publishing*, vol. 10, no. 3152676, pp. 10–5555, 2017.

[33] E. Commission, "Proposal for a regulation laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts," 2021.

[34] H. AI, "High-level expert group on artificial intelligence," *Ethics guidelines for trustworthy AI*, vol. 6, 2019.

[35] S. Han, E. Kelly, S. Nikou, and E.-O. Svee, "Aligning artificial intelligence with human values: reflections from a phenomenological perspective," *AI & SOCIETY*, pp. 1–13, 2022.

[36] A. Alshami, M. Elsayed, E. Ali, A. E. Eltoukhy, and T. Zayed, "Harnessing the power of chatgpt for automating systematic review process: methodology, case study, limitations, and future directions," *Systems*, vol. 11, no. 7, p. 351, 2023.

[37] A. Rao, A. Khandelwal, K. Tanmay, U. Agarwal, and M. Choudhury, "Ethical reasoning over moral alignment: A case and framework for in-context ethical policies in llms," in *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023* (H. Bouamor, J. Pino, and K. Bali, eds.), pp. 13370–13388, Association for Computational Linguistics, 2023.

[38] E. Tennant, S. Hailes, and M. Musolesi, "Moral alignment for LLM agents," in *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*, OpenReview.net, 2025.

[39] M. Autili, D. Di Ruscio, P. Inverardi, P. Pelliccione, and M. Tivoli, "A software exoskeleton to protect and support citizen's ethics and privacy in the digital world," *IEEE Access*, vol. 7, pp. 62011–62021, 2019.

[40] D. Mougouei, A. Azarnik, M. Fahmideh, E. Mougouei, H. K. Dam, A. A. Khan, S. Rafi, J. A. Khan, and A. Ahmad, "A first look at ai trends in value-aligned software engineering publications: Human-llm insights," in *Proceedings of the 47th International Conference on Software Engineering: Software Engineering in Society, ICSE-SEIS2025*, pp. 82–93, ACM, 2025.

[41] C. Alfieri, D. Donati, S. Gozzano, L. Greco, and M. Segala, "Ethical preferences in the digital world: The exosoul questionnaire," in *HHAI2022: Augmenting Human Intellect*, pp. 290–299, IOS Press, 2023.

[42] S. Singh, "Approaches, theories, and role of ethics in computer science and engineering," in *2023 Congress in Computer Science, Computer Engineering, & Applied Computing (CSCE)*, pp. 01–08, IEEE, 2023.

[43] UNESCO, "Recommendation on the ethics of artificial intelligence," 2022.

[44] M. Ryan and B. C. Stahl, "Artificial intelligence ethics guidelines for developers and users: clarifying their content and normative implications," *Inf., Comm. and Ethics in Society*, vol. 19, no. 1, pp. 61–86, 2020.

[45] M. Guarini, "Introduction: machine ethics and the ethics of building intelligent machines," *Topoi*, vol. 32, no. 2, pp. 213–215, 2013.

[46] M. Anderson and S. L. Anderson, "Machine ethics: Creating an ethical intelligent agent," in *Machine ethics & robot ethics*, pp. 237–248, Routledge, 2020.

[47] M. J. O'Fallon and K. D. Butterfield, "A review of the empirical ethical decision-making literature: 1996–2003," *Citation classics from the Journal of Business Ethics*, pp. 213–263, 2012.

[48] J. S. Mill, "Utilitarianism," in *Seven masterpieces of philosophy*, pp. 329–375, Routledge, 2016.

[49] R. Hursthouse, "On virtue ethics," in *Applied ethics*, pp. 29–35, Routledge, 2017.

[50] L. Alexander and M. Moore, "Deontological Ethics," in *The Stanford Encyclopedia of Philosophy* (E. N. Zalta, ed.), Metaphysics Research Lab, Stanford University, Winter 2021 ed., 2021.

[51] C. Boja, A. Zamfiroiu, M. Zurini, and B. Iancu, "User behaviour profiling in social media applications," *Economic Computation & Economic Cybernetics Studies & Research*, vol. 53, no. 1, 2019.

[52] J. Gilbert, S. Hamid, I. A. T. Hashem, N. A. Ghani, and F. F. Boluwatife, "The rise of user profiling in social media: review, challenges and future direction," *Social Network Analysis and Mining*, vol. 13, no. 1, p. 137, 2023.

[53] X. Dong, T. Li, R. Song, and Z. Ding, "Profiling users via their reviews: an extended systematic mapping study," *Software and Systems Modeling*, vol. 20, pp. 49–69, 2021.

[54] R. Falotico and P. Quatto, "Fleiss' kappa statistic without paradoxes," *Quality & Quantity*, vol. 49, pp. 463–470, 2015.

[55] K. Vaniea and L. J. Camp, "Computer security trolley problems: Exploring key factors in security decision-making," in *Proceedings of the 2018 IEEE Symposium on Security and Privacy Workshops (SPW)*, pp. 134–140, IEEE, 2018.

[56] M. Shahin, W. Hussain, A. Nurwidyantoro, H. Perera, R. Shams, J. Grundy, and J. Whittle, "Operationalizing human values in software engineering: A survey," *IEEE Access*, vol. 10, pp. 75269–75295, 2022.

[57] M. Pezzè, M. Ciniselli, L. Di Grazia, N. Puccinelli, and K. Qiu, "The trailer of the ACM 2030 roadmap for software engineering," *SIGSOFT Softw. Eng. Notes*, vol. 49, p. 31–40, Oct. 2024.

[58] L. Dennis, M. Fisher, M. Slavkovik, and M. Webster, "Formal verification of ethical choices in autonomous systems," *Robotics and Autonomous Systems*, vol. 77, pp. 1–14, 2016.

[59] R. C. Cardoso, A. Ferrando, L. A. Dennis, and M. Fisher, "Implementing ethical governors in bdi," in *Workshop on Engineering Multi-Agent Systems*, pp. 22–41, 2021.

[60] L. A. Dennis, M. M. Bentzen, F. Lindner, and M. Fisher, "Verifiable machine ethics in changing contexts," in *AAAI Conference on Artificial Intelligence*, pp. 11470–11478, 2021.

[61] N. S. A. Karim, F. Al Ammar, and R. Aziz, "Ethical software: Integrating code of ethics into software development life cycle," in *2017 International Conference on Computer and Applications (ICCA)*, pp. 290–298, IEEE, 2017.

[62] X. Hou, Y. Zhao, Y. Liu, Z. Yang, K. Wang, L. Li, X. Luo, D. Lo, J. Grundy, and H. Wang, "Large language models for software engineering: A systematic literature review," *ACM Transactions on Software Engineering and Methodology*, vol. 33, no. 8, pp. 1–79, 2024.

[63] S. H. Schwartz, "An overview of the schwartz theory of basic values," *Online readings in Psychology and Culture*, vol. 2, no. 1, pp. 10–20, 2012.