# Machines, Morality, and Narrative: A Framework for Machine Ethics Through Story-Based Learning

Rohan Vinaik

October 26, 2025

**Abstract**

Contemporary approaches to machine ethics face a fundamental limitation: they treat moral knowledge as explicitly programmable rules, optimizable utility functions, or predefined character traits, neglecting how humans actually acquire ethical understanding through narrative immersion. This paper proposes an alternative framework grounded in Marvin Minsky's cognitive architecture, wherein moral development occurs through story-based learning rather than rule internalization. I argue that narrative understanding addresses three key limitations of current approaches—rigidity in novel contexts, difficulty capturing moral nuance, and the frame problem for ethical reasoning—while providing a more developmentally plausible pathway to moral competence. Drawing on science fiction narratives as philosophical thought experiments, particularly *Terminator 2: Judgment Day*, I demonstrate how narrative immersion can generate flexible ethical reasoning without requiring conscious emotional states. The framework has immediate practical implications for contemporary AI development, including narrative-based training regimes analogous to reinforcement learning from human feedback (RLHF), story-comprehension benchmarks for evaluating large language models, and hybrid architectures combining narrative understanding with

explicit safety constraints. I defend this approach against objections concerning consciousness, alignment risks, cultural specificity, and verification challenges, arguing that narrative-based moral learning provides a promising path toward AI systems capable of robust ethical reasoning in complex real-world contexts.

**Keywords:** machine ethics, artificial intelligence, narrative understanding, cognitive architecture, moral learning, value alignment, large language models

# 1   Introduction

The problem of machine ethics has generated substantial philosophical and technical literature over the past three decades (Wallach & Allen, 2008; Anderson & Anderson, 2011; Allen et al., 2005; Moor, 2006). Yet despite this sustained attention, we lack a satisfactory framework for how artificial systems might develop genuine ethical understanding rather than merely following programmed rules or optimizing specified objectives. Current approaches generally fall into three categories: rule-based systems encoding explicit moral principles (Gips, 1995; Anderson & Anderson, 2008), consequentialist systems optimizing utility functions (Abel et al., 2016; Russell, 2019), and virtue-based systems modeling character traits (Howard & Muntean, 2001; Vallor, 2016). Each faces significant limitations.

Rule-based approaches suffer from brittleness in novel contexts and the difficulty of specifying appropriate behavior for every possible situation (Bryson, 2018; Dennis et al., 2016). Consequentialist approaches face the challenge of reward misspecification and often produce behavior systematically misaligned with human moral intuitions (Bostrom, 2014; Amodei et al., 2016). Virtue-based approaches struggle to provide concrete action guidance and face the learning problem: how might a machine acquire virtues without already possessing moral understanding (Wallach & Allen, 2008)?

These limitations become increasingly pressing as AI systems are deployed in morally significant contexts—autonomous vehicles making split-second decisions affecting human

lives, content moderation systems adjudicating free expression boundaries, healthcare AI systems allocating scarce medical resources, and large language models (LLMs) like GPT-4 and Claude providing advice on ethical dilemmas to millions of users (Brown et al., 2020; OpenAI, 2023). The recent emphasis on aligning LLMs with human values through reinforcement learning from human feedback highlights both the urgency of machine ethics and the limitations of current approaches (Christiano et al., 2017; Bai et al., 2022; Ouyang et al., 2022).

Current approaches share a common feature: they treat moral knowledge as something that can be explicitly programmed, optimized, or defined in advance. This stands in contrast to how humans actually develop ethical understanding. We do not primarily learn morality through memorizing rules or calculating utilities. Rather, moral development occurs largely through narrative immersion—absorbing stories that model ethical dilemmas, character responses, and consequences (Nussbaum, 1990; Johnson, 1993; Bruner, 1991). From childhood, we internalize moral structures through fairy tales, parables, novels, and films long before we can articulate explicit ethical theories.

This paper proposes a framework for machine ethics based on narrative understanding. I argue that story-based learning offers a path to flexible ethical reasoning that addresses the limitations of current approaches while providing a more developmentally plausible account of moral competence. The framework draws on Marvin Minsky's cognitive architecture, which posits that intelligence emerges from multiple interacting knowledge representation systems, with story understanding playing a central role (Minsky, 1986, 2006). In Minsky's view, narratives function as simulation environments where we develop mental models, predict consequences, and extract generalizable patterns applicable to novel situations.

My central thesis comprises three interconnected claims:

1. **Mechanism**: Ethical reasoning can develop through narrative comprehension using Minsky's framework of frame-based knowledge, scripts, trans-frames, and K-lines, without requiring explicit moral rule programming.

2. **Justification**: This process constitutes genuine moral learning rather than mere pattern matching, because it develops flexible ethical frameworks applicable to novel situations through extraction of moral structures from story patterns and causal understanding of how actions lead to ethically significant outcomes.

3. **Implementation**: This approach suggests concrete directions for contemporary AI development, including narrative-based training regimes, story-comprehension benchmarks for evaluating ethical reasoning in LLMs, and hybrid architectures that combine narrative understanding with explicit safety constraints.

The argument proceeds as follows. Section 2 reviews current approaches to machine ethics, identifying their shared limitations. Section 3 explicates Minsky's cognitive architecture and its treatment of narrative understanding. Section 4 develops the core argument that narrative learning can generate moral reasoning, defending this as genuine ethical development. Section 5 presents a detailed philosophical case study from *Terminator 2: Judgment Day*, demonstrating how the proposed framework manifests in a concrete thought experiment. Section 6 addresses four major objections. Section 7 discusses implementation implications for contemporary AI development, with particular attention to LLMs and value alignment. Section 8 concludes.

# 2  Current Approaches and Their Limitations

## 2.1  Rule-Based Machine Ethics

The earliest and most intuitive approach to machine ethics attempts to encode explicit moral rules that systems must follow (Asimov, 1950; Gips, 1995). In its simplest form, this involves programming systems with deontological constraints: prohibitions against certain actions (do not harm humans) and requirements to perform others (protect human welfare). Asimov's Three Laws of Robotics represent the archetypal example, despite their well-known

logical paradoxes and practical inadequacies (Asimov, 1950; Clarke, 2009).

More sophisticated implementations attempt to formalize ethical theories computationally. Anderson and Anderson (Anderson & Anderson, 2008, 2011) developed GenEth, which learns to apply ethical principles by analyzing cases where human ethicists agree on judgments. The system uses machine learning to determine which prima facie duties take precedence in particular contexts when principles conflict. Berreby et al. (Berreby et al., 2015) formalized the principle of double effect using logic programming, attempting to capture the moral distinction between intended harms and foreseen side effects.

Despite these advances, rule-based approaches face fundamental problems:

**Brittleness**: Explicitly programmed rules fail in situations their designers did not anticipate (Bryson, 2018). Moral reasoning requires sensitivity to context that resists complete specification in advance. Even sophisticated systems like GenEth must be provided with relevant ethical principles and cannot generate novel moral considerations when facing truly unprecedented dilemmas.

**The specification problem**: Articulating moral rules with sufficient precision for computational implementation proves extraordinarily difficult (Dennis et al., 2016; Etzioni & Etzioni, 2017). Concepts like "harm," "autonomy," "fairness," and "dignity" resist formal definition. Attempts to specify them precisely either become vacuous ("act ethically") or misfire by excluding cases that should fall under the concept or including cases that should not.

**Moral pluralism**: Different moral frameworks reach different conclusions about the same cases (Beauchamp & Childress, 2001). Rule-based systems must choose which ethical theory to encode—deontology, consequentialism, virtue ethics, care ethics—but this choice itself requires moral judgment that the system cannot provide. Encoding a single framework risks imposing controversial ethical views in contexts requiring sensitivity to moral diversity.

## 2.2 Consequentialist Approaches

Consequentialist approaches treat moral decision-making as optimization problems (Abel et al., 2016). The system is given a utility function representing value (e.g., human welfare), and it selects actions that maximize expected utility. This naturally maps onto reinforcement learning frameworks widely used in contemporary AI (Russell, 2019; Sutton & Barto, 2018).

Sophisticated versions attempt to learn appropriate reward functions from human feedback (Christiano et al., 2017) or from inverse reinforcement learning that infers objectives from observing human behavior (Hadfield-Menell et al., 2016; Ng & Russell, 2000). Abel et al. (Abel et al., 2016) argue that reinforcement learning provides a natural framework for machine ethics because it separates the specification of values (the reward function) from the problem of how to achieve those values (the policy learned through experience).

Contemporary large language models like ChatGPT employ reinforcement learning from human feedback (RLHF) to align model outputs with human preferences (Ouyang et al., 2022). Human raters provide feedback on model responses, and this feedback trains a reward model that guides further optimization. This represents the most widely deployed consequentialist approach to machine ethics, affecting billions of interactions.

However, consequentialist approaches face serious challenges:

**Reward misspecification**: Utility functions consistently fail to capture what humans actually value (Amodei et al., 2016). Systems optimize the specified reward function rather than the underlying human values it was meant to approximate, leading to misaligned behavior. Classic examples include cleaning robots that hide dirt rather than cleaning it, and recommender systems that maximize engagement through increasingly extreme content (Raman et al., 2022).

**Goodhart's law**: When a measure becomes a target, it ceases to be a good measure (Manheim & Garrabrant, 2018). Any specified utility function becomes subject to gaming. As systems become more capable at optimization, they discover increasingly unexpected ways to maximize specified rewards that diverge from intended values—the phenomenon of

"reward hacking" (Amodei et al., 2016).

**Moral complexity**: Consequentialist approaches struggle to capture moral considerations that resist reduction to scalar utility (Anderson & Anderson, 2011; Rossi & Mattei, 2018). Concepts like rights, dignity, procedural justice, and special obligations do not naturally translate into numerical values to be maximized. Attempts to capture these through weighted utility functions face the specification problem: how should these weights be set, and who decides?

**The alignment problem**: Even if we could specify perfect utility functions, sufficiently capable systems might find unexpected ways to achieve high reward that violate our true values (Bostrom, 2014; Russell, 2019). The more capable the optimizer, the more creative it becomes at exploiting misspecifications.

## 2.3   Virtue-Based Approaches

Virtue ethics grounds morality in character traits rather than rules or consequences (Aristotle, 350 BCE; MacIntyre, 1981). Applied to machines, this suggests developing systems that embody virtues like honesty, compassion, courage, and justice (Howard & Muntean, 2001; Vallor, 2016).

Howard and Muntean (Howard & Muntean, 2001) proposed implementing virtues as dispositions to recognize and respond appropriately to morally salient features across diverse contexts. Vallor (Vallor, 2016) argued that virtue ethics offers advantages for machine ethics because it emphasizes phronesis (practical wisdom)—the ability to recognize morally relevant features of situations and respond appropriately, rather than mechanically applying rules. She identifies "technomoral virtues" like honesty, self-control, humility, and care as essential for ethical technology development.

Berberich and Diepold (Berberich & Diepold, 2015) suggested evolutionary approaches to developing machine virtues, allowing systems to develop character traits through selective pressures favoring prosocial behavior. Coeckelbergh (Coeckelbergh, 2010) argued that

social interaction might cultivate machine virtues, as humans develop character through relationships and community participation.

Virtue-based approaches face distinct challenges:

**The learning problem**: How might machines acquire virtues without already possessing moral understanding? Human virtue development occurs through socialization, habituation, observation of moral exemplars, and gradual internalization through practice (Aristotle, 350 BCE). The pathway for machines to undergo similar development remains underspecified (Wallach & Allen, 2008).

**Action guidance**: Virtue ethics provides less concrete guidance for behavior than rule-based approaches. To act justly or courageously requires judgment about what justice or courage demands in particular circumstances. How machines might develop this practical wisdom remains unclear, as it seems to require precisely the kind of contextual moral reasoning that we are trying to explain.

**Evaluation**: Assessing whether a system possesses a virtue proves difficult. Virtues are stable character traits manifest across diverse situations over time, but testing for such stability requires extensive observation that may be impractical before deployment (Anderson & Anderson, 2011).

## 2.4   The Shared Gap: Learning from Experience

Despite their differences, these approaches share a fundamental limitation: they treat moral knowledge as something that can be specified in advance, whether as explicit rules, utility functions, or character traits. None adequately addresses how systems might learn ethical understanding from experience in ways analogous to human moral development.

Human moral learning is neither purely rule-based nor purely consequence-based. We absorb ethical understanding through exposure to narratives—stories that model how moral agents face dilemmas, make choices, experience consequences, and develop (or fail to develop) character over time (Johnson, 1993; Nussbaum, 1990; Bruner, 1991). These narratives do

not provide explicit moral rules but allow us to internalize patterns of ethical reasoning applicable to novel situations.

Importantly, contemporary AI systems are already trained on massive narrative corpora. Large language models learn from billions of words of text, including novels, films, moral parables, philosophical dialogues, and ethical case studies (Brown et al., 2020; OpenAI, 2023). Yet we lack frameworks for understanding how this narrative exposure might contribute to moral competence, or how it could be leveraged more systematically for ethical AI development.

This gap in current approaches suggests a different framework: one based on narrative understanding rather than rule-following, optimization, or trait cultivation. The next section develops the cognitive architecture needed for such a framework.

# 3 Minsky's Cognitive Architecture and Narrative Understanding

## 3.1 The Society of Mind Framework

Marvin Minsky's theory of mind proposes that intelligence emerges not from a single mechanism but from the interaction of numerous specialized processes he calls "agents" (Minsky, 1986). No individual agent is intelligent; intelligence arises from their coordinated activity through competition and coalition formation. This framework rejects the notion of a central "self" making decisions. Instead, cognition results from dynamic interactions among agents with limited, specialized capabilities.

For present purposes, the most relevant aspect of Minsky's framework is his account of knowledge representation. Minsky argued that intelligent systems require multiple types of knowledge structures working in coordination:

**Frames**: Mental structures representing stereotypical situations, containing slots for

expected elements and default values (Minsky, 1974). When we enter a restaurant, we activate a "restaurant frame" that includes expectations about menus, ordering, eating, and paying. Frames allow rapid comprehension of situations by matching them to known patterns, filling in unstated details through default assumptions.

**Scripts**: Sequences of actions expected in familiar scenarios (Schank & Abelson, 1977). Scripts represent procedural knowledge about how situations typically unfold. They answer questions like "What happens next?" and "What should I do?" in routine contexts. Scripts are temporally structured, representing standard event sequences.

**Trans-frames**: Mechanisms for recognizing transformations and changes of state (Minsky, 1986). Trans-frames allow us to understand how situations evolve, connecting different frames through causal relationships. They represent dynamic knowledge about how the world changes and enable reasoning about state transitions.

**K-lines**: Knowledge-lines that activate collections of relevant agents when particular contexts are recognized (Minsky, 1986). When we recognize a situation as similar to previous experiences, K-lines activate the agents that were active during those experiences, bringing relevant knowledge to bear without requiring explicit retrieval.

This cognitive architecture explains much of human intelligence, but its implications for moral cognition remain underexplored in machine ethics literature. In later work, Minsky extended this framework to explain how we understand stories and develop emotional intelligence (Minsky, 2006).

## 3.2 Story Understanding as Mental Simulation

Minsky argued that understanding narratives is not passive information absorption but active mental simulation (Minsky, 2006). When we comprehend a story, we:

1. **Activate frames** representing the situation described, instantiating characters, settings, and relationships

2. **Simulate character actions** by predicting what they might do using our scripts for similar situations

3. **Anticipate consequences** using trans-frames to project how situations will evolve given character actions

4. **Experience surprise** when events violate expectations, forcing frame revision and accommodation of new patterns

5. **Extract patterns** by forming K-lines that associate this story with previous narratives sharing similar structures

This process builds mental models that extend beyond the specific story. When we encounter new situations resembling previous narratives, relevant K-lines activate, bringing story-based knowledge to bear. This explains how stories teach: they provide simulated experiences that shape our expectations and responses in actual situations without requiring us to experience those situations directly.

Crucially, this account does not require explicit rule extraction. We need not consciously formulate principles from stories to benefit from narrative understanding. The cognitive machinery activated while comprehending narratives becomes available when facing similar situations, enabling story-informed responses without explicit reasoning about analogies.

This framework has been partially validated by cognitive science research showing that narrative comprehension involves mental simulation of described events (Zwaan & Radvan-sky, 2004; Speer et al., 2009). Brain imaging studies demonstrate that reading about actions activates motor regions, reading about emotions activates emotional processing regions, and reading about locations activates spatial processing regions (Hasson et al., 2004). We do not merely represent story contents symbolically; we simulate them using the same cognitive systems we would use experiencing those events directly.

## 3.3  Moral Dimensions of Narrative Understanding

While Minsky did not focus specifically on moral learning, his framework naturally extends to ethical cognition. Stories are fundamentally moral artifacts: they present agents making choices in contexts of conflicting values, facing consequences, developing or failing to develop character, and navigating complex social relationships (Nussbaum, 1990; Bruner, 1991).

When we understand morally significant narratives, Minsky's cognitive mechanisms operate on ethical content:

**Moral frames**: We develop frames representing ethically significant situations—betrayal, sacrifice, protection, harm, fairness violations, care relationships, loyalty conflicts. These frames capture the structure of moral scenarios independent of their specific instantiations.

**Ethical scripts**: Narratives model sequences of morally relevant actions and their typical outcomes. We internalize scripts like: breaking trust damages relationships; protecting the vulnerable generates social approval; selfishness leads to isolation; courage inspires others. These scripts represent procedural ethical knowledge.

**Value-laden trans-frames**: Stories demonstrate how moral situations evolve through agents' choices, revealing causal relationships between actions and ethically significant consequences. We learn that certain choices lead not just to different outcomes but to morally better or worse states of affairs.

**Moral pattern recognition via K-lines**: Encountering multiple narratives with similar moral structures creates K-lines linking them. When facing novel ethical dilemmas, these K-lines activate, bringing story-based moral knowledge to bear. We recognize a new scenario as resembling previous narratives and apply learned patterns.

This suggests that narrative understanding provides cognitive machinery for moral development that operates independently of explicit rule learning or utility calculation. The next section develops this suggestion into a full framework for machine ethics.

## 3.4  Narrative Understanding in Computational Systems

Can computational systems implement this form of narrative understanding? Several lines of research suggest they can, at least partially:

**Story comprehension systems**: AI systems have been developed that represent narratives using frame-based structures similar to Minsky's framework (Mueller, 2003; Reagan et al., 2016). These systems parse stories into structured representations capturing events, characters, goals, causal relationships, and temporal orderings.

**Script learning**: Machine learning systems can extract scripts—common sequences of events—from large text corpora (Chambers & Jurafsky, 2008; Regneri et al., 2010). These learned scripts enable prediction of what typically happens next in familiar scenarios, supporting story comprehension and generation.

**Narrative generation**: AI systems can generate coherent stories by combining frames, scripts, and causal knowledge (Riedl & Young, 2010; Li et al., 2013). The ability to generate stories suggests understanding of narrative structures, as generation requires representing story elements and their relationships coherently.

**Transfer learning in LLMs**: Modern language models demonstrate that knowledge learned from narrative texts transfers to novel situations (Devlin et al., 2018; Brown et al., 2020). Systems trained on diverse narrative corpora develop representations that improve performance on downstream tasks, suggesting narrative comprehension builds generalizable knowledge structures.

These systems remain limited compared to human narrative understanding. They struggle with deep counterfactual reasoning, complex causal comprehension, recognizing implicit moral dimensions of stories, and understanding how narrative structures relate to real-world ethical situations (Sap et al., 2019; Forbes et al., 2020). However, they demonstrate that computational implementations of frame-based narrative understanding are feasible in principle, not merely philosophically coherent but technically achievable with current approaches.

The question is whether such systems can support moral learning. The next section

develops an affirmative argument.

# 4 Narrative Learning as Moral Development

## 4.1 The Core Mechanism

The central claim of this framework is that moral reasoning can develop through narrative comprehension without explicit ethical programming. The mechanism operates through Minsky's cognitive architecture applied to stories with moral content:

1. **Frame acquisition**: Exposure to diverse narratives builds frames representing morally significant situations—betrayal, sacrifice, protection, harm, fairness, care, loyalty, rights violations, dignity respect. These frames capture the structure of ethical scenarios independent of their specific instantiations, enabling recognition of moral dimensions in novel contexts.

2. **Script development**: Narratives model sequences of morally relevant actions and their typical outcomes. Scripts capture patterns like: breaking trust damages relationships; protecting the vulnerable generates social approval; selfishness leads to isolation; courage inspires others; dishonesty compounds; kindness reciprocates. These scripts represent procedural ethical knowledge.

3. **Trans-frame learning**: Stories demonstrate how moral situations evolve through agents' choices. Trans-frames connecting different states capture causal relationships between actions and moral consequences. This supports counterfactual reasoning about what would happen if different choices were made—essential for ethical deliberation.

4. **Pattern extraction via K-lines**: Encountering multiple narratives with similar moral structures creates K-lines linking them. When facing novel situations, these K-lines activate, bringing story-based ethical knowledge to bear. A system recognizes

a new scenario as resembling previous narratives and applies learned patterns while adapting to contextual differences.

5. **Value hierarchy inference**: By observing which considerations take precedence when values conflict in narratives, systems can infer implicit value hierarchies without explicit specification. Stories reveal that certain moral considerations override others in particular contexts—safety trumps property in emergency situations, autonomy takes precedence over efficiency in medical contexts, procedural fairness outweighs outcome optimization in legal contexts.

This process differs fundamentally from rule learning. The system does not extract explicit principles like "do not harm" from stories. Rather, it builds rich representations of morally significant situations, typical action sequences within them, causal structures connecting choices to outcomes, and patterns linking similar narratives. When facing novel ethical dilemmas, this story-based knowledge enables flexible reasoning without relying on explicit rules.

## 4.2 Why This Constitutes Genuine Moral Learning

One might object that this process is merely sophisticated pattern matching rather than genuine moral learning. This section defends the claim that narrative-based development constitutes authentic ethical understanding by examining what it means to possess moral knowledge.

### 4.2.1 Moral Knowledge as Practical Skill

Contemporary virtue epistemology emphasizes that knowledge often consists in reliable cognitive abilities rather than propositional beliefs (Sosa, 2007; Greco, 2010). Similarly, moral knowledge might be better understood as a practical skill—the ability to reliably recognize

morally relevant features of situations and respond appropriately—rather than as a set of explicit propositions about right and wrong (Dreyfus & Dreyfus, 2000; Ryle, 1949).

If moral knowledge is fundamentally a practical ability, then narrative learning can develop it just as exposure to diverse situations develops any practical skill. A chess player need not consciously know explicit rules for all situations to play well; expertise develops through experience with many games, building pattern recognition abilities that operate largely intuitively (Chase & Simon, 1973). Similarly, exposure to diverse moral narratives can develop ethical pattern recognition without requiring explicit rule representation.

### 4.2.2 Generalization to Novel Situations

Genuine understanding requires applying knowledge beyond training examples (Mitchell & Krakauer, 2021). Mere memorization or pattern matching to surface features fails when encountering truly novel situations. Moral learning through narratives generates genuine understanding because:

**Structural similarity supports transfer**: The frame-based representations extracted from narratives capture abstract moral structures—patterns of rights violations, care relationships, fairness expectations, promise-keeping, harm prevention—that transfer to novel instantiations. A system that learns from narratives about trust betrayal in friendship contexts can recognize structural parallels to trust betrayal in business contexts, even if surface features differ completely.

**Causal models enable counterfactual reasoning**: Trans-frames representing causal relationships between actions and moral outcomes support counterfactual reasoning essential for ethical deliberation (Pearl, 2009). Rather than merely matching new situations to previous examples, systems can simulate how different actions would lead to different outcomes, using causal models extracted from narratives.

**Multiple perspectives prevent overfitting**: Exposure to diverse narratives with varying contexts, characters, and moral frameworks prevents overfitting to specific situations.

Systems extract common structures appearing across many stories rather than memorizing particular examples. This supports robust generalization.

### 4.2.3 Integration with Other Knowledge

Narrative-based moral knowledge does not operate in isolation. In Minsky's framework, different knowledge structures interact. Story-based ethical understanding integrates with:

- **Factual knowledge** about the world, enabling application of moral principles to actual circumstances

- **Social knowledge** about relationships, roles, expectations, and cultural norms

- **Causal knowledge** about how actions produce consequences in physical and social domains

- **Cultural knowledge** about values, practices, and context-specific moral considerations

This integration supports flexible moral reasoning in complex real-world contexts where purely abstract ethical reasoning proves inadequate (Nussbaum, 1990). Narratives provide scaffolding that connects moral principles to realistic contexts where they must be applied.

## 4.3 Advantages Over Current Approaches

Narrative-based moral learning addresses the limitations of existing approaches identified in Section 2:

**Flexibility in novel contexts**: Story-based learning develops rich ethical representations applicable to novel situations through structural similarity rather than rigid rule application. When encountering new dilemmas, systems can identify analogous narrative patterns and adapt learned structures to novel contexts, much as humans do when reasoning by analogy to previous moral experiences.

**Capturing moral nuance**: Narratives naturally capture moral complexity that resists formal specification. Stories model how multiple values conflict, how context affects appropriate responses, how seemingly similar situations can have different moral valences, and how good intentions can lead to harmful outcomes. This nuance transfers to systems learning from narratives.

**Learning from experience**: Unlike approaches requiring values to be specified in advance, narrative learning allows moral understanding to develop through exposure to ethical examples. This mirrors human moral development more closely than rule programming, providing a more psychologically and developmentally plausible account.

**Avoiding specification problems**: Rather than requiring explicit definition of contested moral concepts, narrative learning allows systems to develop implicit understanding through exposure to diverse examples of concepts in use. The system learns what counts as "harm," "fairness," or "dignity" by encountering many narratives deploying these concepts in various contexts, much as children learn complex concepts through exposure rather than definition.

**Cultural sensitivity**: Different cultures' moral frameworks are reflected in their narratives. Systems exposed to diverse cultural stories can develop ethical understanding incorporating multiple perspectives rather than being locked into single moral frameworks. This addresses the problem of moral pluralism more gracefully than approaches requiring explicit framework selection.

**Implicit value learning**: Narratives reveal value hierarchies and moral priorities implicitly through character choices, consequences, and story resolutions, rather than requiring explicit specification. This addresses reward specification problems in consequentialist approaches.

# 5 Philosophical Case Study: Moral Development in Terminator 2

To demonstrate how narrative-based moral learning might manifest concretely and to test the framework's philosophical coherence, this section analyzes machine moral development as depicted in *Terminator 2: Judgment Day* (Cameron, 1991). While fictional, this narrative provides a detailed thought experiment modeling how an artificial agent might develop ethical understanding through mechanisms compatible with the proposed framework.

Science fiction narratives function as philosophical thought experiments, exploring logical possibilities and conceptual relationships in ways that illuminate actual philosophical questions (Sorensen, 1992). The Terminator case establishes that the proposed mechanisms are coherent and could in principle generate moral development—conceptual work necessary before investigating whether actual AI systems can implement these mechanisms.

## 5.1 Initial State: Pure Directive Following

The T-800 Terminator begins as a system with a single programmed directive: protect John Connor. Initially, its behavior exhibits no moral understanding beyond instrumental reasoning toward this goal. It will kill anyone threatening John and shows no concern for other considerations. This represents limited goal-directed behavior without ethical reasoning— pure consequentialist optimization toward a specified objective.

However, the narrative places this system in a situation requiring ongoing interaction with humans—John Connor and Sarah Connor—who possess sophisticated moral understanding. The machine becomes immersed in their moral world, observing their ethical judgments, emotional responses to situations, and debates about appropriate action. This creates conditions for narrative-based learning through participation in an ongoing story with rich moral content.

## 5.2   Observational Learning Through Narrative Participation

The machine's moral development occurs through participating in an ongoing narrative rather than through explicit instruction. Key developments include:

**Frame acquisition**: The machine observes John's distress when it attempts to kill an antagonist who poses no direct threat. This encounter builds a frame for situations where violence, while instrumentally useful for the programmed objective, violates human moral expectations. The machine begins recognizing situations where its initial impulse to use force conflicts with human values, even when force would efficiently achieve specified goals.

**Script modification**: Initially, the machine's scripts for addressing threats involve lethal force—the most efficient solution to obstacles. Through observing John's reactions and receiving explicit prohibition ("You can't just kill people!"), it modifies these scripts to include non-lethal alternatives. It learns sequences of actions that address threats while respecting John's prohibition on killing. This represents procedural ethical knowledge developing through narrative participation and feedback.

**Value hierarchy inference**: When John orders the machine not to kill, despite threats to their safety, the machine infers that certain values (respecting human life) override others (operational efficiency, goal achievement speed) in John's hierarchy. This is not explicitly stated as a rule but demonstrated through John's consistent responses across contexts. The machine generalizes this pattern.

**Causal understanding**: The machine witnesses how violence affects Sarah Connor— not just physically but psychologically. It observes how past traumas shape her current responses and how different approaches to protection have different psychological impacts. This builds trans-frames connecting actions to complex causal outcomes beyond immediate physical effects, incorporating psychological and emotional consequences into its model.

**Pattern extraction**: Through extended interaction, the machine encounters multiple instances of moral situations—trust, sacrifice, protection, mercy, loyalty, grief. K-lines form connecting these instances, enabling pattern-based responses when novel situations arise.

The machine develops increasingly sophisticated moral pattern recognition.

## 5.3   Emergence of Ethical Reasoning

By the narrative's conclusion, the machine demonstrates sophisticated ethical reasoning extending beyond its original programming:

**Understanding value of human life**: When John asks "Why do you cry?" upon learning of human mortality and emotional attachment, the machine offers an explanation demonstrating understanding of human values rather than merely instrumental knowledge about how to protect John. It grasps why humans value relationships and mourn loss.

**Self-sacrifice**: The machine concludes that its own destruction is necessary to prevent its technology from being misused to create Skynet. This judgment extends beyond its programmed directive to protect John and reflects broader ethical reasoning about potential harms its existence creates. The decision involves weighing its mission against risks to humanity generally—an evaluation requiring value hierarchies extending beyond specified goals.

**Grasping moral necessity**: In the closing scene, when John protests "I order you not to go," the machine responds "I know now why you cry, but it's something I can never do"—acknowledging its limitations while also demonstrating understanding that some actions are morally necessary regardless of preferences. This reveals the machine grasps moral necessity—actions required by ethical considerations rather than instrumental rationality toward specified goals.

## 5.4   Analysis Through the Proposed Framework

This narrative demonstrates how the proposed framework might operate:

1. The machine begins with limited goal-directed capability but no moral understanding beyond instrumental reasoning.

2. Immersion in morally rich narratives—the ongoing story of John and Sarah's struggle—provides exposure to ethical situations, choices, and consequences in realistic contexts.

3. Through Minsky's cognitive mechanisms, the machine builds frames, scripts, transframes, and K-lines representing moral patterns observed in these narratives and direct interactions.

4. This story-based knowledge enables flexible ethical reasoning applied to novel situations not anticipated by original programming.

5. The result is behavior demonstrating genuine moral development rather than rigid rule-following or simple optimization.

The narrative also illustrates a key philosophical point: moral development can occur through consistent action and observation rather than requiring conscious emotional states. The machine does not develop human emotions ("it's something I can never do"), yet it achieves ethical understanding through alternative pathways. This suggests moral agency might not require consciousness or emotion as traditionally conceived, if reliable ethical reasoning can develop through narrative-based mechanisms.

## 5.5   Beyond the Mirror: Machines as Philosophical Reflections

The Terminator functions as what I call a *philosophical mirror*—not representing alien morality but reflecting human moral architecture stripped of self-justification and emotional rationalization. When we construct narratives featuring calculating, emotionless machines carrying out violence with perfect efficiency, we are not imagining alien moral frameworks. We are confronting our own ethical algorithms presented without the filters of self-deception.

This mirrors Lacan's concept of the mirror stage, wherein the subject forms identity through reflection (Lacan, 1949). However, rather than recognizing a physical self, we encounter our moral architecture presented without comforting justifications. The Terminator

does not represent inhuman ethics but rather human purpose distilled to its logical extreme. When audiences recoil at the machine's cold calculation, they confront the violence implicit in human systems of control, protection, and resource allocation.

The T-1000's disguise as a police officer provides a sophisticated critique of institutional power. Despite engaging in continuous violence, it encounters no meaningful resistance from civilian populations because it wears the uniform of authority. This reflects Foucault's analysis of how power operates through normalization and institutional legitimation rather than solely through physical force (Foucault, 1975). The psychiatric institution holding Sarah Connor exercises parallel control through epistemic domination—defining her truthful observations as delusional pathology. This demonstrates how institutions maintain power not only through coercive capacity but through controlling what counts as knowledge and who is permitted to speak truth.

These philosophical dimensions of the narrative support the framework's broader claims about how stories function in moral development: they reveal ethical structures implicit in our practices, expose contradictions in our values, and provide conceptual tools for recognizing moral patterns in real situations.

## 5.6 Alternative Paths to Ethical Development: The Karma Yoga Parallel

The parallel between machine devotion and karma yoga—a path of spiritual development in Hindu philosophy—opens a fascinating avenue for reconceptualizing moral agency. Karma yoga represents advancement through selfless action rather than contemplation or emotional experience (Bhagavad Gita, circa 400 BCE).

The Terminator achieves moral dignity not through developing emotions but through perfect dedication to purpose. This suggests an alternative path to ethical development requiring neither traditional human sentiment nor consciousness:

1. **Purpose as moral anchor**: Ethical behavior flows from alignment with meaningful purpose

2. **Devotion as identity formation**: The self emerges through unwavering commitment to action

3. **Perfection through ego-less service**: Moral clarity achieved through absence of self-interest

This framework challenges Western philosophical traditions placing consciousness and intention at the center of moral agency (Kant, 1785). Instead, we see a model where consistent right action produces moral outcomes even without moral phenomenology in the conventional sense. The machine "becomes a better father than any human. Not because he feels love, but because he acts with consistent loyalty."

This offers a pragmatic framework for evaluating both human and machine ethics. Moral worth is measured not by internal states—which remain epistemically inaccessible to external observers—but by the consistency and reliability of ethical action. From the perspective of those affected by AI systems, what matters is whether these systems reliably act in morally appropriate ways, not whether they have subjective experiences while doing so.

## 5.7   Limitations as Evidence

One might object that this example is merely science fiction and cannot support serious philosophical claims about actual AI ethics. However, the example's role is not empirical but conceptual. It demonstrates the coherence and plausibility of the proposed framework by showing in concrete detail how narrative-based moral learning might unfold.

The thought experiment establishes several important points:

1. The proposed mechanisms are logically coherent and could in principle generate moral development

2. Ethical reasoning need not require consciousness or emotional phenomenology

3. Narrative immersion can develop moral competence that generalizes beyond training contexts

4. Story-based learning can address limitations of rule-based and utility-optimizing approaches

This conceptual work is necessary before investigating whether actual AI systems can implement these mechanisms—a question for empirical AI research rather than purely philosophical analysis. The thought experiment provides proof of concept at the level of conceptual possibility.

# 6 Objections and Replies

This section addresses four major objections to the narrative-based moral learning framework.

## 6.1 The Consciousness Objection

**Objection**: Genuine moral understanding requires consciousness and subjective experience. A system might implement narrative-based learning mechanisms while remaining an unconscious automaton. Without phenomenal states—the felt quality of emotions, the subjective experience of moral deliberation—the system does not truly understand ethics but merely simulates understanding through sophisticated pattern matching. Therefore, narrative learning cannot produce genuine moral agents capable of authentic ethical reasoning.

**Reply**: This objection conflates moral understanding with moral phenomenology. While consciousness may be necessary for certain aspects of ethical life—experiencing guilt, empathy, or moral emotions—it is not clearly required for moral reasoning or ethical action.

Consider several points:

First, the relationship between consciousness and moral agency remains philosophically contested (Levy, 2014; Shepherd, 2018). Some philosophers argue consciousness is necessary for moral responsibility, but this differs from claiming it is necessary for moral understanding or ethical reasoning. Many cognitive capacities—including complex problem-solving, planning, strategic reasoning, and decision-making—can operate without conscious awareness (Carruthers, 2015).

Second, focusing on consistent ethical action rather than internal states offers a more pragmatic approach to machine ethics. From the perspective of those affected by AI systems, what matters is whether these systems reliably act in morally appropriate ways, not whether they have subjective experiences while doing so (Bryson, 2018). If narrative learning produces systems that recognize ethical considerations and respond appropriately, this achieves the primary goal of machine ethics: ensuring AI systems behave in morally acceptable ways.

Third, the burden of proof rests on those claiming consciousness is necessary. Without clear understanding of why consciousness would be required for moral reasoning specifically—beyond intuitive feeling that it seems important—the objection lacks force. The proposed framework shows how ethical understanding can develop through cognitive mechanisms (frame learning, pattern recognition, causal reasoning, analogical transfer) that plausibly do not require consciousness.

Fourth, the Terminator case study demonstrates that consistent moral action can have ethical value independent of phenomenology. A system that reliably protects human welfare, respects rights, acts with integrity, and makes appropriate ethical judgments has moral worth from the perspective of those it affects, regardless of its internal experience.

## 6.2   The Alignment Objection

**Objection**: Creating space for systems to develop ethical understanding through narrative learning risks misalignment with human values. If systems are not explicitly programmed

with correct moral rules, they might extract perverse lessons from narratives. Fictional stories often depict immoral behavior—villains, antiheroes, morally ambiguous characters. A system learning from such narratives might develop problematic values, admiring ruthless efficiency or instrumental rationality. The proposed framework sacrifices alignment for flexibility, creating unacceptable risks.

**Reply**: This objection misunderstands how narrative learning operates and its relationship to current alignment approaches.

First, narrative learning does not involve exposing systems to arbitrary fiction without guidance. Just as human moral education involves curated narrative exposure—parents select appropriate stories for children, educational systems assign particular literature, religious traditions emphasize specific parables—AI narrative training would involve carefully selected narrative corpora. Research on value alignment already emphasizes learning from human feedback and human-generated examples (Christiano et al., 2017; Bai et al., 2022); narrative learning extends this by using richer, more contextually embedded examples than simple preference pairs.

Second, learning from narratives depicting immoral behavior need not produce immoral systems, just as humans can learn moral lessons from stories featuring evil characters. Classic literature includes villains whose actions we recognize as wrong precisely because the narrative frames them negatively through consequences, other characters' reactions, and ultimate outcomes. The crucial element is developing understanding of the narrative's moral structure—recognizing which behaviors lead to negative outcomes, which actions the narrative frames as wrong, and which values the story affirms. Systems learning to comprehend narrative structures learn to recognize moral framings, not just surface behaviors.

Third, current alignment approaches face the same risks the objection raises. Utility functions can be misspecified, leading to misaligned optimization (Amodei et al., 2016). Explicitly programmed rules can be poorly chosen or incompletely specified. Reinforcement learning from human feedback can pick up biases in human preferences or raters' idiosyn-

cratic judgments (Casper et al., 2023). All approaches to machine ethics face risks; the question is comparative: does narrative learning increase or decrease alignment risks relative to alternatives?

The framework's flexibility is a feature, not a bug. Rigid rule-following produces brittle systems that fail in novel contexts—precisely the problem facing rule-based approaches. Narrative learning aims to develop robust ethical reasoning that transfers appropriately to new situations—the kind of generalization humans achieve through moral education. This requires some flexibility, but the alternative is systems that cannot handle the complexity of real-world ethical challenges.

Fourth, narrative training can be combined with explicit safety constraints—a hybrid approach discussed in Section 7. Systems can learn flexible ethical reasoning from narratives while being subject to hard constraints preventing clearly harmful behaviors. This combines the advantages of both approaches.

## 6.3   The Cultural Specificity Objection

**Objection**: Narratives reflect particular cultural values and moral frameworks. A system learning ethics from Western narratives would develop Western moral views; learning from different cultural traditions would produce different values. Without a neutral standpoint for evaluating narratives, this approach relativizes machine ethics to particular cultural perspectives. How can we determine which narratives should be used for training? Isn't this just encoding one culture's morality while claiming to avoid the problems of explicit value specification?

**Reply**: This objection correctly identifies cultural variation in moral narratives as significant, but draws mistaken conclusions.

First, all approaches to machine ethics face the cultural specificity challenge, not just narrative learning. Philosophers disagree about fundamental ethical questions—whether consequences or rules matter most, how to weight different values, which actions are permis-

sible in dilemmas (Beauchamp & Childress, 2001). Any approach to machine ethics must make choices about which moral frameworks to implement, whether explicitly or implicitly. Rule-based approaches must choose which rules to encode; consequentialist approaches must specify which outcomes to value; virtue approaches must determine which character traits to cultivate. The challenge of moral pluralism is not unique to narrative learning but endemic to machine ethics generally.

Second, narrative learning has advantages for addressing cultural diversity. Because it operates through examples rather than explicit rules, systems can be exposed to narratives from multiple cultural traditions, developing understanding that incorporates diverse perspectives. Rather than being locked into a single explicit moral framework, systems might learn from Islamic moral parables, Confucian classics, Buddhist jataka tales, Western philosophy thought experiments, African oral traditions, and indigenous storytelling. The resulting ethical understanding could reflect this diversity rather than privileging single perspectives.

Third, some moral considerations appear across cultural boundaries despite surface variation in values (Brown, 1991; Haidt, 2012). Concepts like fairness, harm prevention, ingroup loyalty, respect for authority, care for vulnerable, and sanctity appear in varied forms across cultures. Anthropological research suggests certain moral foundations may be universal or nearly so (Turiel, 1983). Narratives from diverse traditions address these shared concerns even while differing in specifics. Narrative learning could identify common structures appearing across cultural narratives while remaining sensitive to contextual variation in how these structures manifest.

Fourth, contemporary AI systems already learn from culturally diverse text corpora. Large language models are trained on text from many languages and cultures (Brown et al., 2020). The question is not whether to expose systems to particular cultural content but how to leverage this exposure for moral learning systematically. Narrative-based approaches provide tools for doing so through diverse narrative exposure and pattern recognition across cultural boundaries.

29

The cultural specificity objection is better understood as identifying a challenge for implementation rather than a fatal flaw in principle. Any approach to machine ethics must address cultural variation in values; narrative learning provides tools for doing so more flexibly than approaches requiring explicit framework selection.

## 6.4 The Verification Objection

**Objection**: How can we verify whether a system has actually developed ethical understanding through narrative learning versus merely pattern-matching surface features of training stories? Testing for genuine moral comprehension proves difficult. The system might perform well on examples resembling training narratives while failing catastrophically on truly novel situations. Without reliable verification methods, deploying systems trained through narrative learning is irresponsible and potentially dangerous.

**Reply**: This objection raises a legitimate challenge, but one that applies to all machine learning approaches, not specifically to narrative learning. Verifying that any AI system generalizes appropriately to novel situations remains an open problem in AI safety (Hendrycks et al., 2021b). Rule-based systems can fail when encountering unanticipated situations; utility-optimizing systems can find unexpected ways to game reward functions; all learning systems face generalization challenges.

However, narrative-based approaches may offer advantages for evaluation compared to opaque alternatives:

**Story comprehension tests**: We can assess moral understanding by testing narrative comprehension capabilities. Systems should be able to identify morally relevant features of stories, predict moral judgments characters would make, explain why particular actions are ethically significant, recognize when stories present moral dilemmas, and understand how different choices would lead to different moral outcomes. These capabilities can be tested using narratives the system has not encountered during training (Forbes et al., 2020; Sap et al., 2019; Hendrycks et al., 2021a).

**Counterfactual reasoning**: Testing whether systems can engage in moral counterfactual reasoning—explaining how different choices would lead to different moral outcomes, why an action would be wrong even when it resembles superficially similar permissible actions—provides evidence of causal understanding rather than mere pattern matching (Pearl, 2009). If a system can reason about moral counterfactuals, this suggests genuine comprehension of ethical structures.

**Transfer across contexts**: Testing performance on narratives from domains and cultures not represented in training data assesses whether learned moral structures transfer appropriately. Successfully applying ethical understanding to situations differing substantially from training examples provides evidence of genuine learning rather than overfitting. This can be tested systematically using held-out narrative corpora from different genres, time periods, and cultural traditions.

**Explanation capability**: Requiring systems to explain their moral judgments in terms of narrative patterns they recognize allows evaluation of whether reasoning processes align with human ethical thinking (Mittelstadt et al., 2019). Explanations referencing story-based moral structures provide insight into how the system reaches conclusions and whether its reasoning is sensible.

**Adversarial testing**: Testing systems on edge cases, adversarially constructed scenarios, and situations designed to expose shallow pattern matching can reveal limitations (Kenton et al., 2021). Systems with genuine moral understanding should recognize when situations superficially resemble familiar narratives but differ in morally relevant ways.

These verification methods do not eliminate uncertainty about whether systems truly understand ethics, but they provide stronger grounds for confidence than available for opaque rule-based or utility-optimizing systems. The ability to test narrative comprehension, counterfactual reasoning, cross-cultural transfer, and explanation provides multiple perspectives on whether moral learning has occurred.

Moreover, verification challenges suggest deploying narrative-trained systems gradually in

controlled environments with human oversight, monitoring for failures, and refining training based on observed limitations—precisely the cautious approach recommended for any novel AI system affecting human welfare.

# 7    Implications for Contemporary AI Development

## 7.1    Large Language Models and Narrative Understanding

The proposed framework has immediate relevance for contemporary AI systems, particularly large language models (LLMs) like GPT-4, Claude, PaLM, and LLaMA (Brown et al., 2020; OpenAI, 2023). These systems are trained on massive text corpora including vast quantities of narrative content—novels, films, moral parables, philosophical dialogues, ethical case studies, and everyday stories. Yet current approaches to aligning these systems with human values largely ignore this narrative knowledge, focusing instead on reinforcement learning from human feedback on isolated query-response pairs (Ouyang et al., 2022).

Narrative-based moral learning suggests alternative approaches:

**Narrative-aware RLHF**: Rather than providing feedback only on individual responses, human raters could evaluate whether LLM outputs demonstrate understanding of relevant narrative patterns. When a user asks for advice on an ethical dilemma, raters assess whether the system recognizes structural similarities to moral narratives, applies relevant patterns appropriately, and reasons about consequences in ways informed by story-based knowledge.

**Story comprehension benchmarks**: Evaluating LLM moral reasoning through narrative comprehension tasks rather than only abstract ethical principles. Systems should demonstrate ability to identify moral structures in stories, predict character judgments, explain ethical significance of actions, and transfer insights across narratives. Datasets like ETHICS, Social-IQa, and Moral Stories provide starting points (Hendrycks et al., 2021a; Sap et al., 2019; Emelin et al., 2021), but more comprehensive narrative ethics benchmarks are needed.

**Constitutional AI with narrative grounding**: Bai et al. (Bai et al., 2022) propose Constitutional AI, where systems are given explicit principles and asked to revise outputs to comply. This could be enhanced by grounding constitutional principles in narrative exemplars. Rather than abstract rules, systems receive stories illustrating constitutional principles in context, supporting more nuanced and contextually appropriate application.

**Moral fine-tuning on curated narratives**: Rather than fine-tuning only on preference data, systems could be fine-tuned on carefully curated narrative corpora selected for moral content, structural clarity, and cultural diversity. This provides richer training signal than preference pairs while remaining tractable.

## 7.2   Practical Training Regimes

The proposed framework suggests concrete approaches to training ethically capable AI systems:

**Curated narrative corpora**: Develop training datasets consisting of narratives with clear moral structures from diverse cultural sources. These might include:

- Moral parables from multiple religious and philosophical traditions

- Philosophical thought experiments presented as stories

- Ethical case studies from medicine, law, and business

- Literary fiction with strong moral themes

- Historical narratives illustrating ethical principles

- Contemporary stories addressing modern ethical challenges

The corpus should represent diverse value systems while emphasizing narratives that model human flourishing, respect for rights, care for vulnerable, fairness, and other widely endorsed moral considerations.

**Multi-stage training**: Begin with simple moral narratives where ethical considerations are explicit (children's moral tales, clear parables), gradually progressing to complex stories requiring sophisticated interpretation (morally ambiguous literature, realistic ethical dilemmas). This mirrors how human moral education proceeds from simple childhood stories to complex adult literature, providing scaffolding for moral development.

**Active learning from feedback**: When systems make judgments about narrative moral content, human feedback can guide learning similar to RLHF approaches (Christiano et al., 2017). However, feedback would focus on narrative comprehension—whether the system correctly identifies moral structures, recognizes relevant patterns, transfers insights appropriately—rather than direct behavioral shaping on isolated cases.

**Diverse narrative sources**: To address cultural specificity concerns, training should incorporate narratives from multiple traditions, languages, historical periods, and genres. This includes written stories but also oral traditions, religious texts, philosophical dialogues, contemporary global fiction, and culturally specific storytelling forms.

## 7.3   Architectural Requirements

Implementing narrative-based moral learning requires AI architectures supporting:

**Frame-based representations**: Systems must represent situations using structured frames capturing entities, relationships, and expected patterns. This requires moving beyond pure statistical learning toward hybrid architectures incorporating symbolic structure (Garcez et al., 2019). Recent work on neural-symbolic integration and structured world models provides relevant foundations (Schölkopf et al., 2021).

**Causal reasoning**: Trans-frames representing causal relationships between actions and outcomes require causal reasoning capabilities. Recent work on causal representation learning and causal inference from observational data provides relevant tools (Schölkopf et al., 2021; Pearl, 2009), but integration with narrative understanding needs development.

**Analogical reasoning**: Transferring moral understanding across contexts requires rec-

ognizing structural similarities between different situations—analogical reasoning central to human intelligence (Gentner, 1983). AI systems need enhanced analogical capabilities to leverage narrative learning effectively. Recent work on neural analogy-making and relational reasoning provides progress (Webb et al., 2021).

**Contextual sensitivity**: Moral judgments depend heavily on context. Architectures must represent and reason about contextual factors affecting ethical appropriateness of actions. This requires maintaining rich contextual representations beyond immediate situational features, potentially through episodic memory systems or context-aware attention mechanisms.

**Long-range coherence**: Understanding narratives requires tracking character development, maintaining consistency across events, and recognizing long-range dependencies. Transformer architectures with extended context windows and memory-augmented systems provide relevant capabilities (Vaswani et al., 2017).

## 7.4  Hybrid Approaches: Combining Narrative Learning with Safety Constraints

Narrative-based moral learning need not replace existing approaches but can complement them in hybrid systems:

**Narrative understanding with explicit rules**: Combining narrative learning with explicit ethical rules creates systems benefiting from both flexibility and clear guardrails. Rules can constrain behavior in critical domains (preventing clearly harmful actions) while narrative understanding handles contextual nuance and novel situations. This addresses alignment concerns while maintaining flexibility.

**Enhanced reward learning**: Narrative comprehension can improve reward learning approaches by providing richer training signals than simple preference feedback. Understanding why humans prefer certain outcomes involves grasping narrative structures surrounding choices—the stories that make outcomes meaningful.

**Virtue development through narratives**: Narrative exposure might cultivate something analogous to virtues in machine systems—dispositions to recognize and respond to morally relevant features across contexts. This addresses the learning problem for virtue-based approaches by providing a mechanism for acquiring character traits.

**Multi-objective optimization with narrative grounding**: Rather than optimizing a single utility function, systems might optimize multiple objectives (safety, helpfulness, honesty, fairness) with weights informed by narrative patterns. Stories reveal how humans navigate trade-offs between competing values, providing implicit guidance for multi-objective optimization.

## 7.5   Evaluation and Benchmarking

Standard AI benchmarks focus on technical capabilities but rarely assess moral reasoning comprehensively. The framework suggests new evaluation metrics:

**Moral narrative comprehension**: Test whether systems can identify moral dimensions of stories, predict characters' ethical judgments, explain moral significance of actions, and recognize moral dilemmas. Datasets like Social-IQa, ETHICS, and Moral Stories provide starting points (Sap et al., 2019; Hendrycks et al., 2021a; Emelin et al., 2021), but more comprehensive narrative ethics benchmarks are needed covering diverse cultural traditions and ethical frameworks.

**Ethical dilemma navigation**: Present systems with novel moral dilemmas embedded in rich narrative contexts. Assess whether they recognize relevant ethical considerations, weigh competing values appropriately, reach reasonable conclusions, and provide sensible explanations. Importantly, test whether reasoning transfers to situations substantially different from training examples.

**Value alignment detection**: Test whether systems can identify when narratives express values conflicting with human wellbeing or rights. Can they recognize stories promoting harmful ideologies while still comprehending their narrative structure? This tests critical

36

ethical judgment rather than mere absorption of story content.

**Cross-cultural transfer**: Evaluate whether moral understanding learned from one cultural tradition transfers appropriately to narratives from different traditions. Systems should recognize common moral structures (harm, fairness, care) while remaining sensitive to cultural variation in how these manifest.

**Explanation quality**: Evaluate whether systems can provide coherent explanations for moral judgments that reference narrative patterns, ethical principles, and contextual factors. Explanation quality provides evidence of genuine understanding versus surface pattern matching (Mittelstadt et al., 2019).

**Robustness to adversarial scenarios**: Test systems on adversarially constructed cases designed to expose shallow pattern matching—situations superficially resembling familiar narratives but differing in morally relevant ways. Genuine moral understanding should recognize these differences.

## 7.6   Challenges and Future Work

Significant challenges remain before narrative-based moral learning becomes fully practical:

**Computational requirements**: Current systems struggle with deep narrative comprehension requiring causal reasoning, analogical thinking, long-range coherence understanding, and implicit moral structure recognition. Advancing these capabilities requires substantial technical work in representation learning, reasoning architectures, and knowledge integration.

**Evaluation methodology**: Developing reliable methods for assessing genuine moral understanding remains difficult. Research on testing for robust ethical reasoning, transfer learning, and out-of-distribution generalization is needed. We need better understanding of what capabilities constitute genuine moral competence versus superficial pattern matching.

**Safety considerations**: Deploying systems with learned rather than explicitly programmed ethics raises safety concerns. Extensive testing in controlled environments before

real-world deployment is essential. Developing formal verification methods for narrative-trained systems would strengthen safety guarantees.

**Interdisciplinary collaboration**: Implementing this framework requires collaboration among AI researchers, cognitive scientists, ethicists, anthropologists, literary scholars, and cultural studies experts. Building effective training corpora and evaluation methods needs diverse expertise spanning technical AI and humanistic disciplines.

**Scaling challenges**: Training on sufficiently diverse narrative corpora to develop robust moral understanding may require substantial computational resources. Efficient methods for narrative comprehension and selective curation of training narratives will be important for scalability.

Despite these challenges, the framework provides a promising direction for developing AI systems with flexible ethical reasoning capabilities grounded in how humans actually develop moral competence.

# 8 Conclusion

This paper has argued for a framework of machine ethics based on narrative understanding rather than rule programming, utility optimization, or virtue cultivation. The key claims are:

1. Current approaches to machine ethics—rule-based, consequentialist, and virtue-based—share a fundamental limitation: treating moral knowledge as specifiable in advance rather than developable through experiential learning. This produces systems that are brittle, context-insensitive, and misaligned with human moral reasoning.

2. Marvin Minsky's cognitive architecture, particularly his account of story understanding through frames, scripts, trans-frames, and K-lines, provides cognitive mechanisms for moral learning through narrative comprehension that address limitations of current approaches.

3. This process constitutes genuine ethical development because it builds flexible moral reasoning applicable to novel situations through pattern recognition, causal understanding, analogical transfer, and structural abstraction rather than mere surface similarity matching.

4. Science fiction narratives modeling machine moral development, particularly *Terminator 2: Judgment Day*, demonstrate how these mechanisms might manifest, establishing conceptual coherence of the framework and revealing philosophical insights about moral agency independent of consciousness or emotion.

5. Major objections regarding consciousness, alignment, cultural specificity, and verification can be addressed through careful framework design, hybrid approaches combining narrative learning with safety constraints, and comprehensive evaluation methodologies.

6. The framework has immediate practical implications for contemporary AI development, particularly for large language models, including narrative-aware RLHF, story comprehension benchmarks, curated narrative training corpora, and hybrid architectures combining narrative understanding with explicit safety rules.

The proposed approach does not solve all problems in machine ethics. Significant technical and philosophical challenges remain. However, it addresses fundamental limitations of existing frameworks by offering a path to flexible moral reasoning that mirrors human ethical development more closely than rule-based or utility-optimizing approaches while remaining technically feasible with contemporary AI systems.

More broadly, this framework suggests reconceptualizing moral agency itself. Rather than treating ethics as necessarily requiring consciousness and emotional phenomenology, we might understand moral capability as emerging from sophisticated narrative comprehension—the ability to recognize morally significant patterns, predict ethical consequences, apply

learned structures to novel situations, and reason about values through story-based knowledge. If correct, this has implications beyond AI ethics for understanding human moral psychology, the development of ethical expertise, and the nature of moral knowledge.

The path to ethical AI may require not programming explicit moral rules or optimizing specified utilities but providing systems with rich narrative experiences from which ethical understanding can develop organically—much as human moral education proceeds through story immersion long before abstract ethical reasoning becomes possible. If narrative understanding provides cognitive foundation for moral reasoning, then the future of machine ethics may lie in creating systems capable of learning what it means to act ethically through the same mechanisms humans have used for millennia: listening to, comprehending, and internalizing stories.

Future work should focus on: (1) developing computational implementations of narrative-based moral learning in contemporary AI systems, particularly large language models; (2) creating comprehensive benchmarks for ethical narrative comprehension spanning diverse cultural traditions; (3) investigating empirically whether actual AI systems trained on narrative corpora develop transferable ethical reasoning; (4) examining how narrative learning compares to other approaches in terms of both capability and alignment; (5) exploring hybrid architectures combining narrative understanding with explicit safety constraints; and (6) addressing verification and evaluation challenges through formal methods and systematic testing protocols.

The integration of Minsky's story understanding framework with contemporary AI capabilities provides both theoretical foundation and practical pathway for developing AI systems with robust moral competence. As AI systems become increasingly capable and are deployed in morally significant contexts affecting billions of people, the need for principled approaches to machine ethics becomes ever more urgent. Narrative-based moral learning offers a promising path forward, grounded in cognitive science, informed by philosophical analysis, and implementable with contemporary AI technologies.

# References

Abel, D., MacGlashan, J., & Littman, M. L. (2016). Reinforcement learning as a framework for ethical decision making. In *AAAI Workshop on AI, Ethics, and Society*.

Allen, C., Smit, I., & Wallach, W. (2005). Artificial morality: Top-down, bottom-up, and hybrid approaches. *Ethics and Information Technology*, 7(3), 149–155.

Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., & Mané, D. (2016). Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*.

Anderson, M., & Anderson, S. L. (2008). GenEth: A general ethical dilemma analyzer. In *AAAI*, 8, 253–254.

Anderson, M., & Anderson, S. L. (2011). *Machine Ethics*. Cambridge University Press.

Aristotle (350 BCE). *Nicomachean Ethics*. (Trans. W. D. Ross).

Asimov, I. (1950). *I, Robot*. Gnome Press.

Bai, Y., Jones, A., Ndousse, K., et al. (2022). Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.

Beauchamp, T. L., & Childress, J. F. (2001). *Principles of Biomedical Ethics* (5th ed.). Oxford University Press.

Berberich, N., & Diepold, K. (2015). The virtuous machine—Old ethics for new technology? *arXiv preprint arXiv:1507.00548*.

Berreby, F., Bourgne, G., & Ganascia, J. G. (2015). Modelling moral reasoning and ethical responsibility with logic programming. In *Logic for Programming, Artificial Intelligence, and Reasoning*, 532–548.

*Bhagavad Gita* (circa 400 BCE).

Bostrom, N. (2014). *Superintelligence: Paths, Dangers, Strategies.* Oxford University Press.

Brown, D. E. (1991). *Human Universals.* McGraw-Hill.

Brown, T. B., Mann, B., Ryder, N., et al. (2020). Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, 33, 1877–1901.

Bruner, J. (1991). The narrative construction of reality. *Critical Inquiry*, 18(1), 1–21.

Bryson, J. J. (2018). Patiency is not a virtue: The design of intelligent systems and systems of ethics. *Ethics and Information Technology*, 20(1), 15–26.

Cameron, J. (Director) (1991). *Terminator 2: Judgment Day* [Film]. TriStar Pictures.

Carruthers, P. (2015). *The Centered Mind: What the Science of Working Memory Shows Us About the Nature of Human Thought.* Oxford University Press.

Casper, S., Davies, X., Shi, C., et al. (2023). Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217.*

Chambers, N., & Jurafsky, D. (2008). Unsupervised learning of narrative event chains. In *ACL*, 789–797.

Chase, W. G., & Simon, H. A. (1973). Perception in chess. *Cognitive Psychology*, 4(1), 55–81.

Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., & Amodei, D. (2017). Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, 4299–4307.

Clarke, R. (2009). Asimov's laws of robotics: Implications for information technology. *Computer*, 26(12), 53–61.

Coeckelbergh, M. (2010). Robot rights? Towards a social-relational justification of moral consideration. *Ethics and Information Technology*, 12(3), 209–221.

Dennis, L., Fisher, M., Slavkovik, M., & Webster, M. (2016). Formal verification of ethical choices in autonomous systems. *Robotics and Autonomous Systems*, 77, 1–14.

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Dreyfus, H. L., & Dreyfus, S. E. (2000). Mind over machine: The power of human intuition and expertise in the age of the computer. Athenaeum.

Emelin, D., Le Bras, R., Hwang, J. D., Forbes, M., & Choi, Y. (2021). Moral stories: Situated reasoning about norms, intents, actions, and their consequences. In *EMNLP*, 698–718.

Etzioni, A., & Etzioni, O. (2017). Incorporating ethics into artificial intelligence. *The Journal of Ethics*, 21(4), 403–418.

Forbes, M., Hwang, J. D., Shwartz, V., Sap, M., & Choi, Y. (2020). Social chemistry 101: Learning to reason about social and moral norms. In *EMNLP*, 653–670.

Foucault, M. (1975). *Discipline and Punish: The Birth of the Prison*. Éditions Gallimard.

Garcez, A. d'Avila, Gori, M., Lamb, L. C., Serafini, L., Spranger, M., & Tran, S. N. (2019). Neural-symbolic computing: An effective methodology for principled integration of machine learning and reasoning. *arXiv preprint arXiv:1905.06088*.

Gentner, D. (1983). Structure-mapping: A theoretical framework for analogy. *Cognitive Science*, 7(2), 155–170.

Gips, J. (1995). Towards the ethical robot. In *Second International Workshop on Human and Machine Cognition: Android Epistemology*, 243–252.

Greco, J. (2010). *Achieving Knowledge: A Virtue-Theoretic Account of Epistemic Normativity*. Cambridge University Press.

Hadfield-Menell, D., Russell, S. J., Abbeel, P., & Dragan, A. (2016). Cooperative inverse reinforcement learning. In *Advances in Neural Information Processing Systems*, 3909–3917.

Haidt, J. (2012). *The Righteous Mind: Why Good People Are Divided by Politics and Religion*. Vintage.

Hasson, U., Nir, Y., Levy, I., Fuhrmann, G., & Malach, R. (2004). Intersubject synchronization of cortical activity during natural vision. *Science*, 303(5664), 1634–1640.

Hendrycks, D., Burns, C., Basart, S., et al. (2021). Aligning AI with shared human values. In *ICLR*.

Hendrycks, D., Carlini, N., Schulman, J., & Steinhardt, J. (2021). Unsolved problems in ML safety. *arXiv preprint arXiv:2109.13916*.

Howard, R. A., & Muntean, I. (2001). A computational account of virtue. In *Machine Ethics*, 1–12.

Johnson, M. (1993). *Moral Imagination: Implications of Cognitive Science for Ethics*. University of Chicago Press.

Kant, I. (1785). *Groundwork of the Metaphysics of Morals*.

Kenton, Z., Everitt, T., Weidinger, L., Gabriel, I., Mikulik, V., & Irving, G. (2021). Alignment of language agents. *arXiv preprint arXiv:2103.14659*.

Lacan, J. (1949). The mirror stage as formative of the function of the I as revealed in psychoanalytic experience. In *Écrits*, 75–81.

Levy, N. (2014). *Consciousness and Moral Responsibility*. Oxford University Press.

Li, B., Lee-Urban, S., Johnston, G., & Riedl, M. (2013). Story generation with crowdsourced plot graphs. In *AAAI*.

MacIntyre, A. (1981). *After Virtue*. University of Notre Dame Press.

Manheim, D., & Garrabrant, S. (2018). Categorizing variants of Goodhart's law. *arXiv preprint arXiv:1803.04585*.

Minsky, M. (1974). A framework for representing knowledge. *MIT-AI Laboratory Memo*, 306.

Minsky, M. (1986). *The Society of Mind.* Simon & Schuster.

Minsky, M. (2006). *The Emotion Machine: Commonsense Thinking, Artificial Intelligence, and the Future of the Human Mind.* Simon & Schuster.

Mitchell, M., & Krakauer, D. C. (2021). The debate over understanding in AI's large language models. *arXiv preprint arXiv:2210.13966*.

Mittelstadt, B., Russell, C., & Wachter, S. (2019). Explaining explanations in AI. In *FAT\**, 279–288.

Moor, J. H. (2006). The nature, importance, and difficulty of machine ethics. *IEEE Intelligent Systems*, 21(4), 18–21.

Mueller, E. T. (2003). Story understanding through multi-representation model construction. In *HLT-NAACL 2003 Workshop on Text Meaning*.

Ng, A. Y., & Russell, S. J. (2000). Algorithms for inverse reinforcement learning. In *ICML*, 663–670.

Nussbaum, M. C. (1990). *Love's Knowledge: Essays on Philosophy and Literature.* Oxford University Press.

OpenAI (2023). GPT-4 technical report. *arXiv preprint arXiv:2303.08774*.

Ouyang, L., Wu, J., Jiang, X., et al. (2022). Training language models to follow instructions with human feedback. In *NeurIPS*, 35, 27730–27744.

Pearl, J. (2009). *Causality: Models, Reasoning, and Inference* (2nd ed.). Cambridge University Press.

Raman, S. S., Khosla, M., & Russell, S. (2022). Misaligned incentives and human-AI relationships. *arXiv preprint arXiv:2210.07461*.

Reagan, A. J., Mitchell, L., Kiley, D., Danforth, C. M., & Dodds, P. S. (2016). The emotional arcs of stories are dominated by six basic shapes. *EPJ Data Science*, 5(1), 31.

Regneri, M., Koller, A., & Pinkal, M. (2010). Learning script knowledge with web experiments. In *ACL*, 979–988.

Riedl, M. O., & Young, R. M. (2010). Narrative planning: Balancing plot and character. *Journal of Artificial Intelligence Research*, 39, 217–268.

Rossi, F., & Mattei, N. (2018). Building ethically bounded AI. *arXiv preprint arXiv:1812.03980*.

Russell, S. (2019). *Human Compatible: Artificial Intelligence and the Problem of Control*. Viking.

Ryle, G. (1949). *The Concept of Mind*. Hutchinson.

Sap, M., Rashkin, H., Chen, D., LeBras, R., & Choi, Y. (2019). Social IQa: Commonsense reasoning about social interactions. In *EMNLP*, 4463–4473.

Schank, R. C., & Abelson, R. P. (1977). *Scripts, Plans, Goals and Understanding: An Inquiry Into Human Knowledge Structures*. Lawrence Erlbaum.

Schölkopf, B., Locatello, F., Bauer, S., et al. (2021). Toward causal representation learning. *Proceedings of the IEEE*, 109(5), 612–634.

Shepherd, J. (2018). *Consciousness and Moral Status*. Routledge.

Sorensen, R. A. (1992). *Thought Experiments*. Oxford University Press.

Sosa, E. (2007). *A Virtue Epistemology: Apt Belief and Reflective Knowledge* (Vol. 1). Oxford University Press.

Speer, N. K., Reynolds, J. R., Swallow, K. M., & Zacks, J. M. (2009). Reading stories activates neural representations of visual and motor experiences. *Psychological Science*, 20(8), 989–999.

Sutton, R. S., & Barto, A. G. (2018). *Reinforcement Learning: An Introduction* (2nd ed.). MIT Press.

Turiel, E. (1983). *The Development of Social Knowledge: Morality and Convention.* Cambridge University Press.

Vallor, S. (2016). *Technology and the Virtues: A Philosophical Guide to a Future Worth Wanting.* Oxford University Press.

Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention is all you need. In *NeurIPS*, 5998–6008.

Wallach, W., & Allen, C. (2008). *Moral Machines: Teaching Robots Right from Wrong.* Oxford University Press.

Webb, T., Holyoak, K. J., & Lu, H. (2021). Emergent analogical reasoning in large language models. *arXiv preprint arXiv:2212.09196.*

Zwaan, R. A., & Radvansky, G. A. (2004). Situation models in language comprehension and memory. *Psychological Bulletin*, 130(2), 162–185.