

# Stratified Geometric Ethics: Mathematical Foundations for Verifiable Moral Reasoning in Autonomous Systems

Andrew H. Bond

*San José State University*

andrew.bond@sjsu.edu

## Abstract

We present *Stratified Geometric Ethics* (SGE), a mathematical framework providing rigorous foundations for deterministic, verifiable ethical reasoning in autonomous systems. SGE addresses fundamental limitations of prior geometric approaches by modeling the space of ethically relevant configurations as a *stratified space*—a union of smooth manifolds of varying dimensions connected along boundary strata—rather than a smooth manifold. This structure captures moral discontinuities, incommensurable values, and genuine ethical dilemmas that smooth models cannot represent.

We make five principal contributions. First, we show that stratified spaces are a *natural minimal candidate* among standard geometric structures for representing ethical phenomena including threshold effects, lexical priorities, and moral dilemmas (Theorem 2.3). Second, we establish a *representation theorem* (Theorem 4.3) characterizing all satisfaction functionals satisfying five explicit axioms—including a novel locality axiom and scale-normalization assumption—with a complete proof. Third, we prove *finite approximation theorems* (Theorems 3.9–3.11) showing that any decision problem on a compact stratified space reduces to a finite graph problem with explicit error bounds. Fourth, we establish *decidability results* (Theorem 6.4) for the quantifier-free, non-temporal fragment of our ethical specification language via o-minimal structures, with temporal properties handled by standard model checking on finite approximations. Fifth, we derive *sample complexity bounds* (Theorems 7.2–7.4) for learning ethical content from data.

SGE serves as the theoretical foundation for the DEME 2.0 architecture presented in the companion paper [2], providing mathematical justification for design choices including multi-dimensional moral vector spaces, governance profiles with veto regions and lexical priorities, and layered enforcement architectures.

**Keywords:** AI ethics, differential geometry, stratified spaces, formal verification, autonomous systems, moral philosophy

# 1 Problem statement and approach

The deployment of AI systems in safety-critical domains—healthcare, autonomous vehicles, criminal justice, financial markets—creates an urgent need for frameworks that make ethical reasoning explicit, deterministic, and formally verifiable. Current approaches fall into two categories: *principle-based* approaches that articulate high-level commitments without computational implementations, and *learning-based* approaches that infer ethical behavior from feedback without formal guarantees. Neither satisfies the requirements of safety-critical deployment where lives depend on ethical decisions made at machine timescales.

Recent work has proposed differential-geometric frameworks for ethics, modeling configurations as points on smooth manifolds and ethical quantities as tensor fields [1]. While mathematically elegant, these approaches face five fundamental limitations:

**The Content Problem.** Existing frameworks assume obligations and interests are "given" without specifying how these tensors derive from ethical sources—constitutional principles, stakeholder preferences, or deliberative procedures.

**The Smoothness Problem.** Real ethical reasoning involves discrete choices (trolley problems), incommensurable values (life vs. property), and threshold effects (killing vs. letting die). Smooth manifolds cannot represent these discontinuities.

**The Linearity Problem.** Bilinear satisfaction functionals  $S(x) = O^\mu(x)I_\mu(x)$  cannot capture threshold effects, lexical priorities, or diminishing returns—phenomena central to moral reasoning.

**The Tractability Problem.** No existing geometric ethics framework provides complexity analysis establishing whether ethical computation can be performed in real time.

**The Verification Problem.** No framework enables machine-checkable proofs of ethical constraint satisfaction.

This paper introduces *Stratified Geometric Ethics* (SGE), a framework that addresses all five limitations while preserving the coordinate-invariance that makes geometric approaches valuable. The companion paper [2] instantiates SGE in the DEME 2.0 architecture for real-time ethical governance in autonomous systems.

## 1.1 Relationship to DEME 2.0

SGE was developed through a process of reflective equilibrium with DEME 2.0, a computational architecture for real-time ethical governance described in the companion paper [2]. DEME commits to concrete engineering choices: multi-dimensional moral vector spaces, governance profiles with veto regions and lexical priorities, a three-layer architecture (strategic, tactical, reflex), and hardware-resident ethics modules. SGE provides the mathematical justification for these choices:

- DEME's *moral vector space*  $M \subseteq \mathbb{R}^k$  is justified by Theorem 2.3, which shows stratified spaces are natural minimal candidates for ethical phenomena.
- DEME's *governance profiles* with veto regions and scalarization are justified by Theorem 4.3's representation of satisfaction operators.
- DEME's *real-time enforcement* is justified by the finite approximation theorems (3.9–3.11) and complexity bounds (Section 5).
- DEME's *verification layer* is justified by the decidability results (Theorem 6.4) combined with standard model checking.

## 2 Why Stratified Spaces

Before developing the machinery of SGE, we establish that stratified spaces are a *natural minimal candidate* among standard geometric structures for capturing ethical phenomena. We do not claim a fully formal minimality theorem over all conceivable geometric structures—such a claim would require defining "simpler" in a way that excludes exotic alternatives. Instead, we show that standard alternatives (smooth manifolds, manifolds with corners, cell complexes) each fail to represent at least one essential ethical phenomenon, and that stratified spaces succeed at all of them.

### 2.1 Ethical Phenomena Requiring Geometric Representation

We identify four categories of ethical phenomena that any adequate geometric framework must represent:

**E1. Discrete Choices.** Many ethical decisions present finitely many options with no meaningful intermediate. In trolley-type scenarios, one must choose path A or path B; there is no "path  $0.5A + 0.5B$ ."

**E2. Incommensurable Values.** Some moral dimensions cannot be traded off at any finite rate. The value of a human life is not equivalent to any amount of property damage; procedural rights cannot be "compensated" by increased welfare.

**E3. Threshold Effects.** Crossing certain boundaries has discontinuous moral significance. The distinction between killing and letting die, between lying and remaining silent, between 0.99 and 1.00 on a consent scale—these involve discrete moral transitions.

**E4. Genuine Dilemmas.** Some situations admit no ethically satisfactory resolution. Both available options involve genuine moral loss ("moral residue"), and this property must be represented, not smoothed away.

### 2.2 Insufficiency of Smooth Manifolds

**Proposition 2.1.** *Let  $M$  be a smooth connected manifold. Then:*

- (i) *Any two points in  $M$  can be connected by a smooth path (violates E1);*
- (ii) *Any smooth function  $f: M \rightarrow \mathbb{R}$  has continuous level sets (violates E3);*
- (iii) *For any smooth Riemannian metric  $g$  on  $M$ , all directions at any point have finite, comparable costs (violates E2).*

*Proof.* (i) follows from path-connectedness of connected manifolds. (ii) follows from continuity of smooth functions. (iii) follows from positive-definiteness of Riemannian metrics:  $g(v,v) > 0$  for all nonzero  $v$ , so all directions have positive finite cost, hence are comparable. ■

### 2.3 Stratified Spaces as Natural Minimal Candidates

**Definition 2.2** (Stratified Space). *A stratified space is a triple  $(M, \{M_i\}_{i \in I}, \leq)$  where  $M$  is a paracompact Hausdorff space,  $\{M_i\}$  is a locally finite partition into connected smooth manifolds (strata), and  $\leq$  is a partial order on  $I$  such that  $M_i \cap \text{cl}(M_j) \neq \emptyset$  implies  $i \leq j$  (frontier condition). We require Whitney's condition (B) for regularity.*

**Theorem 2.3** (Stratified Spaces Represent E1–E4). *Stratified spaces can represent all four ethical phenomena:*

- (i) *Discrete choices are modeled by 0-dimensional strata representing decision endpoints;*

- (ii) *Incommensurable values are modeled by stratified metrics with arbitrarily large cost (in a singular limit);*
- (iii) *Threshold effects are modeled by stratum boundaries where satisfaction functions have discontinuities;*
- (iv) *Genuine dilemmas are modeled by singular strata from which all exits incur positive moral cost.*

Moreover, among standard geometric structures—smooth manifolds, manifolds with corners, and cell complexes without smooth structure—each fails at least one of E1–E4. Stratified spaces therefore emerge as a natural minimal candidate.

*Proof.* The positive claims (i)–(iv) follow by explicit construction in Section 2.4. For the comparative claims: smooth manifolds fail by Proposition 2.1. Manifolds with corners can represent E1 and E3 (via boundary strata) but not E2 (all boundary components have finite, comparable codimension) or E4 (corners are not "singular" in the required sense—they have well-defined tangent cones with finite dimension). Cell complexes can represent E1 and E4 but lack the smooth structure on cells needed for tensor calculus (obligations as vector fields, interests as covector fields). Stratified spaces combine smooth strata (enabling differential geometry) with singular boundaries (enabling discontinuities), achieving representation of all four phenomena. ■

**Remark 2.4.** We do not claim that stratified spaces are minimal among *all conceivable* geometric structures—a determined geometer could propose disconnected manifolds, manifolds with degenerate metrics, subanalytic sets, or other exotic alternatives. Our claim is that stratified spaces are the natural minimal choice among *standard* structures commonly used in differential geometry and topology.

## 2.4 Constructive Examples

We provide explicit constructions demonstrating how stratified spaces represent each phenomenon. In these examples, we distinguish between *deliberation strata* (interior strata that parameterize reasoning or intermediate configurations) and *decision strata* (0-dimensional terminal strata). An admissible decision trajectory may traverse higher dimensional strata during deliberation, but must terminate on a 0-dimensional decision stratum; remaining indefinitely in a deliberation stratum is not considered a completed decision.

**Example 2.5** (Trolley Problem—Discrete Choice).

Let  $M = \{A, B\} \cup \gamma$ , where  $A$  and  $B$  are 0-dimensional strata (the discrete choices) and  $\gamma \cong (0, 1)$  is a 1-dimensional stratum (the deliberation space). Topologize  $M$  so that  $\gamma \cong (0, 1)$  and its closure satisfies  $\bar{\gamma} = \gamma \cup \{A, B\}$ , identifying  $A$  with 0 and  $B$  with 1. The frontier condition gives  $A, B \leq \gamma$ . Any path in  $M$  from deliberation to decision must terminate at a discrete point—there is no continuous interpolation staying within the terminal decision strata between  $A$  and  $B$ .

**Example 2.6** (Lexical Priority—Incommensurability).

Let  $M = \mathbb{R}_{\geq 0} \times \mathbb{R}_{\geq 0}$  with coordinates  $(v_1, v_2)$  representing two values. Stratify as  $M_0 = \{(0, 0)\}$ ,  $M_1 = \{0\} \times \mathbb{R}_{> 0}$ ,  $M_2 = \mathbb{R}_{> 0} \times \{0\}$ , and  $M_3 = \mathbb{R}_{> 0} \times \mathbb{R}_{> 0}$ . Define a family of metrics on  $M_3$  by  $g_\epsilon = \text{diag}(1/\epsilon^2, 1)$ . As  $\epsilon \rightarrow 0$ , motion in the  $v_1$  direction becomes unboundedly costly relative to  $v_2$ , yielding a lexical limit in which any decrease in  $v_1$  dominates any increase in  $v_2$ . This corresponds to a governance profile where  $v_1$  is lexically prior to  $v_2$ . Strictly, the  $\epsilon \rightarrow 0$  limit yields a singular (degenerate) geometry;

nevertheless, for sufficiently small  $\epsilon$ , any bounded-cost path behaves lexicographically, since arbitrarily small decreases in  $v_1$  dominate any bounded improvement in  $v_2$ .

**Example 2.7** (Consent Threshold—Discontinuity).

Let  $M = [0,1]$  with strata  $M_0 = \{\tau\}$  (the threshold),  $M_1 = [0, \tau)$ , and  $M_2 = (\tau, 1]$ . Set  $S(\tau) = 0$  (or  $S(\tau) = 1$ , depending on policy),

so the discontinuity is located at the stratum boundary. Define satisfaction by

$$S(x) = \begin{cases} 0, & x < \tau, \\ S(\tau), & x = \tau, \\ 1, & x > \tau. \end{cases}$$

The stratum boundary at  $\tau$  represents the discrete moral transition from “insufficient consent” to “adequate consent.” The stratification isolates the morally exceptional boundary point where the rule regime changes.

**Example 2.8** (Sophie's Choice—Genuine Dilemma).

Let  $M = D \cup \{A, B\}$ , where  $D \cong (0,1)$  is a 1-dimensional dilemma stratum with  $\bar{D} = D \cup \{A, B\}$ , and  $A, B$  are 0-dimensional strata. Define satisfaction so that any path from  $D$  to either endpoint incurs irreducible moral residue (e.g.,  $S(A) = S(B) = \alpha < 1$  while  $\sup_{x \in D} S(x) \leq \beta < \alpha$ ). Then the dilemma stratum is singular in the sense that no trajectory can reach a fully satisfactory terminal state; all exits encode unavoidable moral loss.

### 3 Stratified Moral Spaces

We now develop the formal machinery of stratified moral spaces, providing the geometric foundation for ethical reasoning.

#### 3.1 Formal Definitions

**Definition 3.1** (Stratified Moral Space). *A stratified moral space is a stratified space*

$(M, \{M_i\}, \preceq)$  where:

- (i)  $M$  is the total moral space of ethically relevant configurations;
- (ii) Each stratum  $M_i$  is a smooth manifold representing configurations admitting smooth ethical trade-offs;
- (iii) Stratum boundaries represent moral discontinuities—dilemmas, phase transitions, or categorical distinctions.

**Definition 3.2** (Moral Dimension). *The moral dimension of  $x \in M$  is  $\dim(x) := \dim(M_i)$  where  $x \in M_i$ . Points on lower-dimensional strata represent morally "singular" configurations where some trade-offs break down.*

**Definition 3.3** (Stratified Tangent Bundle). *At  $x \in M_i$ , the stratified tangent space is  $T^{\text{str}}_x M := T_x M_i \oplus N_x(M_i, M)$  where  $T_x M_i$  is the tangent space to the stratum and  $N_x$  is the normal cone.*

**Definition 3.4** (Stratified Ethical Selection). *Given a stratified moral space  $(M, \{M_i\}, \preceq)$  and a feasible action set  $A \subseteq M$ , an action  $a^* \in A$  is ethically admissible if:*

1.  $a^*$  lies in a minimal stratum  $M_k$  such that  $A \cap M_k \neq \emptyset$ ;
2.  $a^*$  is maximal with respect to a smooth ethical preference functional  $\phi_k$  defined on  $M_k$ .

#### 3.2 Whitney's Condition (B) and Its Role

**Definition 3.5** (Whitney's Condition B). *Let  $M_i \subset \text{cl}(M_j)$ . The pair satisfies Whitney (B) at  $x \in M_i$  if: for sequences  $\{y_n\} \subset M_j$ ,  $\{x_n\} \subset M_i$  with  $y_n, x_n \rightarrow x$ , if secant lines  $\ell_n = x_n y_n \rightarrow \ell$  and  $T_{y_n} M_j \rightarrow \tau$ , then  $\ell \subset \tau$ .*

**Lemma 3.6** (Consequences of Whitney B). *If  $M$  satisfies Whitney (B), then: (i) paths approaching stratum boundaries have well-defined limiting tangent directions; (ii) the*

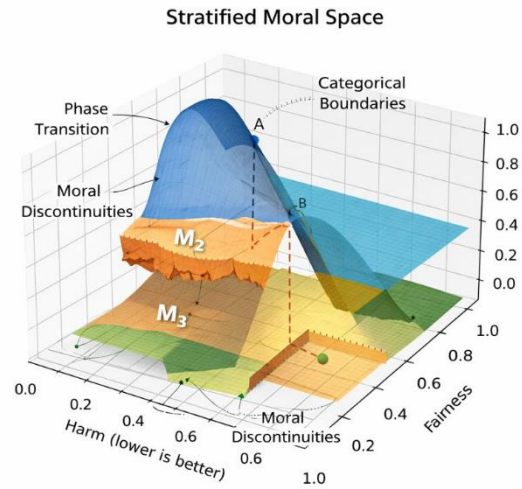


Figure 1 | Stratified moral space with moral discontinuities

A schematic 3D stratified space  $(M, \{M_i\}, \preceq)$  depicting ethically relevant configurations. The horizontal axes parameterize two illustrative ethical dimensions—**Harm** (lower is better) and **Fairness**—while the vertical axis denotes **Ethical satisfaction** (a scalar proxy for overall moral value). The full volume represents the **total moral space  $M$** . Three representative strata  $M_1, M_2, M_3$  (semi-transparent surfaces) illustrate regions in which ethical trade-offs vary smoothly: within any single stratum, nearby configurations admit continuous, “manifold-like” adjustments (e.g., small changes in harm can be compensated by changes in fairness without qualitative shifts in moral kind).

**Stratum boundaries** (annotated as **Moral discontinuities**) mark loci where smooth trade-offs fail—corresponding to dilemmas, sharp regime changes, or categorical distinctions; these are indicated by ridge/edge features and highlighted with dashed guides and arrows. **Phase transition** denotes a boundary across which the governing ethical regime changes, while **Categorical boundaries** indicate constraints that separate morally distinct classes of actions or outcomes. Example configurations **A** and **B** (marked points) lie on different strata; the vertical dashed projections indicate their coordinates in (Harm, Fairness) and emphasize that moving between strata requires crossing a discontinuity rather than a

stratification is locally topologically trivial along each stratum; (iii) geodesic limits exist and are well-behaved at boundaries.

**Remark 3.7** (Verifying Whitney B in Practice). For computational implementations, moral spaces are typically constructed as semialgebraic sets (defined by polynomial inequalities) or as finite unions of manifolds with explicit boundary conditions. By Łojasiewicz's theorem, semialgebraic sets admit Whitney stratifications, and the stratification can be computed algorithmically. For moral spaces defined by threshold conditions (e.g., "harm >  $\tau$ "), the level sets of polynomial functions automatically satisfy Whitney (B).

### 3.3 Finite-Dimensional Approximation Theorems

For computational implementation, continuous stratified spaces must be approximated by finite structures. The following theorems establish that this can be done with explicit error bounds while preserving stratification structure.

**Definition 3.8** (Stratified Graph). A stratified graph is  $G = (V, E, \sigma, \leq)$  where  $(V, E)$  is a finite directed graph,  $\sigma: V \rightarrow A$  assigns vertices to stratum labels, and  $\leq$  is a partial order on  $A$  such that  $(v, w) \in E$  implies  $\sigma(v) \leq \sigma(w)$ .

**Theorem 3.9** (Existence of  $\varepsilon$ -Approximations). Let  $M$  be a stratified space with finitely many compact strata. For every  $\varepsilon > 0$ , there exists a finite stratified graph  $G_\varepsilon$  that is an  $\varepsilon$ -approximation of  $M$ : (i) vertices map to their labeled strata; (ii) every point of  $M$  is within  $\varepsilon$  of some vertex; (iii) edges connect nearby vertices respecting stratum order.

**Theorem 3.10** (Decision Approximation Bound). Let  $(M, A, u, \Gamma)$  be a decision problem where  $u$  is  $L$ -Lipschitz. If  $\pi^*$  is optimal on  $M$  and  $\pi^*$  is optimal on  $G_\varepsilon$ , then  $|u(x, \pi^*(x)) - u(x, \pi^*(v_x))| \leq 2L\varepsilon + \omega(\varepsilon)$  where  $\omega$  is the modulus of continuity of the optimal value function.

**Theorem 3.11** (Structure Preservation). For sufficiently small  $\varepsilon$ , any  $\varepsilon$ -approximation of a Whitney (B) space preserves stratification structure: paths between strata correspond to directed graph paths, and constraint sets constant on strata are exactly preserved.

**Corollary 3.12** (DEME Implementability). DEME's finite moral vector spaces and discrete governance profiles constitute  $\varepsilon$ -approximations of continuous stratified moral spaces, with approximation error bounded by Theorem 3.10.

## 4 The Representation Theorem

The representation theorem characterizes all satisfaction operators consistent with explicit axioms. Unlike prior work that merely assumed a particular functional form, we *derive* the form from first principles, with all assumptions made explicit.

### 4.1 Five Axioms for Satisfaction Operators

**Definition 4.1** (Satisfaction Operator). A stratified satisfaction operator is a map  $\Sigma: \Gamma(T^{\text{str}}M) \times \Gamma(T^{\text{str}}M) \times \Gamma(\text{Sym}^2 T^{\text{str}}M) \times C \rightarrow C^\infty(M)$  taking obligation field  $O$ , interest field  $I$ , metric  $g$ , and constraint set  $C$  to a satisfaction function  $S: M \rightarrow \mathbb{R}$ .

We impose five axioms:

**Axiom 1** (Coordinate Invariance). For any diffeomorphism  $\psi: M \rightarrow M$ ,  $\Sigma(\psi^*O, \psi^*I, \psi^*g, \psi^*C)(x) = \Sigma(O, I, g, C)(\psi^{-1}(x))$ .

**Axiom 2** (Monotonicity). If  $O'(x) = O(x) + \delta O$  where  $I(x)(\delta O) > 0$ , then  $\Sigma(O', I, g, C)(x) \geq \Sigma(O, I, g, C)(x)$ .

**Axiom 3** (Constraint Respect). If  $x \in C$  (constraint set), then  $\Sigma(O, I, g, C)(x) = -\infty$ .

**Axiom 4** (Stratum Compatibility).  $\Sigma$  restricts to a smooth function on each stratum of  $M$ .

**Axiom 5** (Locality).  $\Sigma(O, I, g, C)(x)$  depends only on the pointwise values  $O(x)$ ,  $I(x)$ ,  $g(x)$ , and the indicator  $I_C(x)$ , not on derivatives, jets, or values at other points.

**Remark 4.2** (Motivation for Locality). Axiom 5 captures the intuition that satisfaction is a *local scoring operator* with no hidden history or look-ahead. An agent's ethical standing at configuration  $x$  should depend on the ethical facts *at*  $x$ , not on curvature tensors, gradients, or how the agent arrived at  $x$ . This matches the operational semantics of DEME, where Ethics Modules evaluate moral vectors pointwise.

## 4.2 Scale Normalization

In addition to the five axioms, we impose an explicit *modeling assumption* to avoid arbitrary sensitivity to units or scaling conventions:

**Assumption (Scale Normalization)**. Reparametrizing the obligation field by a positive scalar,  $O \mapsto \alpha O$  for  $\alpha > 0$ , should not affect the ranking of configurations, aside from a context-dependent scaling  $\lambda(x)$  that may depend on  $x$  but not on  $O$ .

This assumption is a modeling choice, not a logical consequence of the axioms. It reflects the principle that the *direction* of obligations relative to interests matters, but arbitrary unit choices (e.g., measuring harm in "utils" vs. "millutils") should not change ethical conclusions.

## 4.3 Complete Proof of the Representation Theorem

**Theorem 4.3** (Representation of Stratified Satisfaction). Let  $\Sigma$  satisfy Axioms 1–5 and the Scale Normalization assumption. Then there exist:

- (i) A smooth monotone function  $f: \mathbb{R} \rightarrow \mathbb{R}$  (activation function);
- (ii) A smooth function  $\lambda: M \rightarrow \mathbb{R}_+$  (scale field);
- (iii) A constraint indicator  $\chi_C: M \rightarrow \{0, -\infty\}$ ;

such that on the regular stratum:

$$\Sigma(O, I, g, C)(x) = \chi_C(x) + \lambda(x) \cdot f(I_\mu(x)O^\mu(x) / \sqrt{(g_{\mu\nu}(x)O^\mu(x)O^\nu(x))})$$

with appropriate limit behavior on singular strata.

*Proof.* We proceed in four steps.

**Step 1: Classification of Local Invariants.** By Axiom 5 (Locality),  $\Sigma(x)$  depends only on  $O(x)$ ,  $I(x)$ ,  $g(x)$ , and  $I_C(x)$ . By Axiom 1 (Coordinate Invariance),  $\Sigma(x)$  can depend only on coordinate-invariant quantities constructible from these pointwise values.

At a point  $x$  on the regular stratum,  $O(x) \in T_x M$ ,  $I(x) \in T_x^* M$ , and  $g(x)$  is a symmetric bilinear form on  $T_x M$ . The invariant theory of the orthogonal group  $O(n)$  classifies all polynomial invariants of such data.

**Lemma 4.4.** The polynomial invariants of  $(O(x), I(x), g(x))$  under coordinate changes are generated by: (a)  $I_\mu O^\mu$  (the contraction); (b)  $g_{\mu\nu} O^\mu O^\nu$  (squared norm of  $O$ ); (c)  $g^{\mu\nu} I_\mu I_\nu$  (squared norm of  $I$  with respect to  $g^{-1}$ ). *Proof of Lemma.* This is the first fundamental theorem of invariant theory for  $O(n)$  (see any standard reference.) ■



**Step 2: Scale Normalization.** By the Scale Normalization assumption,  $\Sigma$  must be invariant under  $O \mapsto \alpha O$  up to a context-dependent factor  $\lambda(x)$ . The only dimensionless combination of invariants (a)–(c) that is independent of scaling  $O$  by  $\alpha$  is:

$$z := I_\mu O^\mu / \sqrt{(g_{\mu\nu} O^\mu O^\nu)}$$

This is the component of  $I$  in the direction of  $O$ , normalized by  $\|O\|$ . It measures how well obligations align with interests, independent of the magnitude of  $O$ .

Here,  $\lambda$  is a local scalar scale field, allowed to depend on  $x$  and local invariants built from  $(I, g, C)$  but not on the arbitrary overall scaling of  $O$ ; thus any dependence on  $\|I\|_{g^{-1}}$  can be folded into  $\lambda$  while preserving coordinate invariance and locality.

The invariant  $\|I\|^2 = g_{\mu\nu} I^\mu I^\nu$  would also be coordinate-invariant, but scale normalization requires that scaling  $O \rightarrow \alpha O$  not change rankings. Since  $\|I\|$  is independent of  $O$ , any dependence on  $\|I\|$  can be absorbed into the context-dependent factor  $\lambda(x)$ , leaving  $z$  as the sole  $O$ -dependent argument.

**Step 3: Monotonicity Constraint.** Axiom 2 requires that increasing the component of  $O$  in the direction favored by  $I$  must increase  $\Sigma$ . Since  $z$  increases when  $I_\mu O^\mu$  increases,  $\Sigma$  must be a *monotone increasing* function of  $z$ . Thus  $\Sigma(x) = \lambda(x) \cdot f(z)$  where  $f$  is monotone increasing and  $\lambda(x) > 0$  is the context-dependent scale factor.

**Step 4: Constraint Respect.** Axiom 3 requires  $\Sigma(x) = -\infty$  for  $x \in C$ . This is achieved by adding  $\chi_C(x) := 0$  if  $x \notin C$ ,  $-\infty$  if  $x \in C$ . ■

#### 4.4 Uniqueness and Degrees of Freedom

**Corollary 4.5.** *The representation is unique up to the choice of: (i) the activation function  $f$  (capturing threshold/saturation structure); (ii) the scale field  $\lambda$  (capturing context-dependent sensitivity); (iii) the constraint set  $C$  (capturing absolute prohibitions).*

These degrees of freedom correspond precisely to DEME governance profile parameters:  $f$  corresponds to scalarization functions (e.g., weighted sums, lexical priorities),  $\lambda$  to dimension-specific transforms, and  $C$  to veto regions.

### 5 Computational Complexity

For ethical reasoning to be deployable in real-time systems, we must establish computational tractability with explicit bounds.

#### 5.1 Satisfaction Evaluation

**Theorem 5.1** (Complexity of Satisfaction Evaluation). *Let  $M$  be a stratified space of dimension  $n$  with  $k$  constraint predicates. Evaluating  $\Sigma(O, I, g, C)(x)$  requires  $O(n^2 + k)$  arithmetic operations.*

*Proof.* Computing  $I_\mu O^\mu$  requires  $O(n)$  operations. Computing  $g_{\mu\nu} O^\mu O^\nu$  requires  $O(n^2)$  operations. Evaluating  $k$  constraint predicates requires  $O(k)$  operations. The activation function  $f$  and scale field  $\lambda$  are  $O(1)$ . Total:  $O(n^2 + k)$ . ■

**Corollary 5.2** (DEME Real-Time Feasibility). *For DEME's typical configuration ( $n \leq 10$  dimensions,  $k \leq 20$  constraints), satisfaction evaluation requires  $\leq 200$  arithmetic operations, achievable in  $< 1 \mu s$  on contemporary embedded processors.*

## 5.2 Geodesic Planning

**Theorem 5.3** (Complexity of Geodesic Planning). *Let  $M$  have  $m$  strata with  $N$  vertices per stratum in the  $\varepsilon$ -approximation. Finding an optimal path requires  $O(mN^2 \cdot n^2 \cdot \log(mN))$  operations.*

## 6 Formal Verification via O-Minimal Structures

We establish decidability results for verifying ethical specifications, providing theoretical grounding for DEME's verification layer.

### 6.1 Ethical Specification Language

**Definition 6.1** (Ethical Specification Language). *ESL formulas are built from:*

- Atomic predicates:  $S(x) \triangleright \triangleleft c$ ,  $O_i(x) \triangleright \triangleleft c$ ,  $I_i(x) \triangleright \triangleleft c$  where  $\triangleright \triangleleft \in \{<, \leq, =, \geq, >\}$
- Boolean connectives:  $\neg, \wedge, \vee, \Rightarrow$
- Quantifiers over regions:  $\forall x \in R. \phi$ ,  $\exists x \in R. \phi$
- Temporal operators: Always, Eventually, Until

**Remark 6.2** (Scope of Decidability). The decidability result (Theorem 6.4 below) applies to the *quantifier-free, non-temporal fragment* of ESL—essentially Boolean combinations of polynomial inequalities over the moral space. For the full language including temporal operators, we rely on a two-stage approach: (1) use Theorem 6.4 for static properties, and (2) check temporal properties on the finite approximating graph  $G_\varepsilon$  using standard LTL/CTL model-checking algorithms. Since  $G_\varepsilon$  is finite, this preserves decidability.

### 6.2 Decidability via O-Minimality

**Definition 6.3** (O-Minimal Structure). *A structure  $(\mathbb{R}, <, +, \cdot, \dots)$  is o-minimal if every definable subset of  $\mathbb{R}$  is a finite union of points and intervals.*

**Theorem 6.4** (Decidability of Quantifier-Free ESL). *For polynomial  $O$ ,  $I$ ,  $g$  and semialgebraic  $C$ , satisfiability of quantifier-free, non-temporal ESL formulas is decidable.*

*Proof.* By Tarski–Seidenberg, the first-order theory of  $(\mathbb{R}, <, +, \cdot)$  is decidable. Quantifier-free ESL formulas reduce to Boolean combinations of polynomial inequalities, decidable by cylindrical algebraic decomposition. ■

**Corollary 6.5** (Temporal Properties via Model Checking). *Temporal ESL formulas over the finite approximation  $G_\varepsilon$  can be verified using standard LTL/CTL model-checking algorithms with complexity polynomial in  $|G_\varepsilon|$  and exponential in formula size.*

**Corollary 6.6** (DEME Verification). *DEME governance profiles with polynomial veto predicates and weighted-sum scalarization admit decidable verification of static safety properties, with temporal properties verified via model checking on the discrete implementation.*

## 7 Learning-Theoretic Foundations

SGE's ethical content must ultimately be specified by humans. We establish sample complexity bounds for learning this content from data.

## 7.1 Learning Obligation Weights

**Theorem 7.2** (Sample Complexity for Weight Learning). *Let  $W = \{w \in \mathbb{R}^k : \|w\|_1 = 1, w \geq 0\}$  be the weight simplex. Empirical risk minimization achieves generalization error  $\leq \epsilon$  with probability  $\geq 1 - \delta$  using  $N = O(k \log(k/\delta) / \epsilon^2)$  samples.*

## 7.2 Learning Interest Fields

**Theorem 7.3** (Sample Complexity for Utility Learning). *For a hypothesis class with pseudo-dimension  $d$  and preference oracle noise rate  $\eta < 1/2$ , utility estimation error  $\leq \epsilon$  requires  $N = O((d \log(1/\epsilon) + \log(1/\delta)) / ((1-2\eta)^2 \epsilon^2))$  samples.*

## 7.3 Learning Metrics

**Theorem 7.4** (Sample Complexity for Metric Learning). *For  $p$ -parameter metric class, a consistent metric requires  $m = O((p \log p + \log(1/\delta)) / \epsilon)$  trajectory pairs with margin  $\epsilon$ .*

# 8 Worked Example: Medical Triage

To make the framework concrete, we present a simplified triage scenario.

## 8.1 Moral Space Construction

Consider allocating one ICU bed among three patients. The moral space is the 2-simplex  $\Delta^2 = \{(p_A, p_B, p_C) : p_i \geq 0, \sum p_i = 1\}$ , where  $p_i$  is the probability of allocation to patient  $i$ .

### Stratification:

- *Interior* (2-dimensional): Probabilistic allocations where all patients have positive probability.
- *Edges* (1-dimensional): Allocations between two patients, one excluded.
- *Vertices* (0-dimensional): Deterministic allocations—the actual decisions.

This stratification captures that deterministic allocation (vertices) is categorically different from probabilistic allocation (interior)—crossing from interior to vertex represents a discrete ethical transition.

## 8.2 Obligations and Interests

**Obligations** (derived from principles):

- Beneficence:  $O_{\text{ben}} = \nabla(\text{expected health outcome})$
- Urgency:  $O_{\text{urg}}$  points toward most critical patient
- Equity:  $O_{\text{eq}}$  points toward disadvantaged patients
- Rights:  $O_{\text{rts}} = -\nabla(\text{coercion} + \text{consent violation})$

The aggregate obligation field is  $O = w_1 O_{\text{ben}} + w_2 O_{\text{urg}} + w_3 O_{\text{eq}} + w_4 O_{\text{rts}}$ .

**Interests** (derived from stakeholders):  $I = \alpha_{\text{clin}} I_{\text{clin}} + \alpha_{\text{pat}} I_{\text{pat}} + \alpha_{\text{inst}} I_{\text{inst}}$ .

## 8.3 Constraint Set

$$C = \{x : \text{coercion}(x) > 0\} \cup \{x : \text{consent}(x) < \tau\}$$

If Patient  $C$  involves coercion or insufficient consent, then  $S(\text{vertex } C) = -\infty$ .

## 8.4 Satisfaction Evaluation

Using sigmoid activation  $f(z) = \tanh(z)$ :

- $S(A) = 0 + 1.0 \cdot \tanh(0.87/1.0) \approx 0.70$
- $S(B) = 0 + 1.0 \cdot \tanh(0.62/1.0) \approx 0.55$
- $S(C) = -\infty$  (constraint violation)

**Decision:** Patient A, with formal certificate that C is forbidden and A dominates B.

### 8.5 Multi-agent preview: Two triage officers

To illustrate natural extension to multi-agent scenarios, consider two physicians  $A$  and  $B$  jointly allocating three ICU beds among four patients. Each physician has obligation field  $O_\mu^a$  and interest field  $I_{a\mu}$  over their proposed allocation. The interaction tensor  $G_{ab}$  encodes:

- $G_{aa} = 1$  (each cares about their own clinical judgment)
- $G_{ab} = 0.3$  for  $a \neq b$  (each partially defers to colleague's expertise)

Joint satisfaction is  $W = G_{ab}I_{a\mu}O_\mu^b$ . The stratified space includes:

- Interior: negotiable allocations
- Disagreement boundary: when physicians' vetoes conflict
- Resolution stratum: escalation to ethics committee

Full multi-agent theory is beyond scope here but demonstrates theoretical extensibility.

## 9 SGE as the Foundation for DEME 2.0

We summarize how SGE's theorems justify the DEME 2.0 architecture presented in the companion paper [2].

**DEME's Moral Vector Space.** DEME represents configurations as vectors  $m \in \mathbb{R}^k$ . *SGE justification:* Theorem 2.3 shows stratified spaces are natural minimal candidates;  $\mathbb{R}^k$  with coordinate thresholds is the simplest such space.

**DEME's Governance Profiles.** DEME profiles specify veto regions, scalarization, and lexical priorities. *SGE justification:* Theorem 4.3 proves this is the unique form satisfying Axioms 1–5 plus scale normalization.

**DEME's Real-Time Layer.** DEME compiles profiles to hardware operating at sub-ms timescales. *SGE justification:* Theorem 5.1 gives  $O(n^2 + k)$  complexity; Corollary 5.2 confirms embedded feasibility.

**DEME's Verification.** DEME verifies ethical specifications. *SGE justification:* Theorem 6.4 provides decidability for static properties; Corollary 6.5 handles temporal properties via model checking.

## 10 Limitations and Future Work

**Metaethical neutrality.** SGE provides a geometric and computational representation of ethical content, but does not resolve metaethical questions about the ultimate grounding of values. The framework is intentionally pluralistic and can accommodate consequentialist,

deontological, virtue-based, and hybrid constitutional content, provided that such content can be operationalized as obligations, constraints, and/or preference models within the moral state space.

**Specification burden.** Constitutional principles and prohibited regions must be specified (or at least endorsed) by human governance. SGE automates computation and enforcement, not the original creation of normative content. In particular, the choice of values, lexical priorities, and constraint boundaries remains a governance problem rather than a purely technical one.

**Scaling and real-time feasibility.** The complexity of SGE’s core computations is polynomial in the dimensionality and discretization of the moral space, but practical costs can still grow rapidly with dimension and with the number of constraints. Very high-dimensional spaces (e.g.,  $n > 100$ ) may challenge strict real-time requirements without additional structure (e.g., sparsity, low-rank metrics, stratified factorization, or learned surrogates) and without careful engineering of candidate action sets.

**Adversarial robustness and boundary vulnerabilities.** The present framework is primarily an “in-model” theory: it specifies how obligations, interests, metrics, and constraints induce satisfaction and admissible trajectories once the moral state  $x \in M$  is given. It does not, by itself, prevent adversarial manipulation of the world-to-moral mapping  $\phi: W \rightarrow M$  (e.g., sensor tampering, deceptive reporting, ontology gaps, or distribution shift), nor does it fully address “boundary hugging” behaviors near thresholds and stratum transitions. Robust deployment therefore requires explicit threat models, uncertainty-aware constraint checking, and tamper-evident measurement pipelines.

**Toward a Stratified Ethical Integrity Monitor (SEIM).** A natural next step is to treat ethical governance as a security problem and to formalize a reference-monitor–style enforcement kernel. We propose a *Stratified Ethical Integrity Monitor (SEIM)*: a minimal, verifiable component that mediates all actuation and enforces hard ethical invariants (forbidden regions  $C$ , lexical constraints, and admissible-stratum termination conditions) regardless of the optimization objectives of the controlled system. In such a design, actions would require an explicit approval token issued only if safety conditions are certified in the relevant moral representation. A mathematically robust SEIM would also incorporate epistemic uncertainty by operating on uncertainty sets  $X(w) \subseteq M$  rather than point estimates, enforcing a robust condition of the form “ $\text{permit}(a) \Rightarrow X(w, a) \cap C = \emptyset$ ,” where  $X(w, a)$  is the reachable moral region induced by action  $a$  under bounded uncertainty. This direction aligns SGE’s geometric apparatus with assurance-oriented notions from security and robust control.

**Future work.** Immediate directions include: (i) multi-agent and institutional ethics (coupled obligation/interest fields and bargaining over stratified spaces), (ii) temporal dynamics and trajectory-level constraints (stratum-safe planning, compositional guarantees for macro-actions), (iii) uncertainty quantification and robust enforcement (including distribution-shift detection and conservative fallbacks), and (iv) empirical validation, including comparison against human judgments and red-team benchmarks targeting specification gaming and representation attacks.

**Multi-agent extension.** The tensor formalism naturally extends to  $N$ -agent settings via agent-indexed obligation and interest fields  $O_\mu^a, I_{a\mu}$ , with interaction tensors  $G_{ab}$  encoding cooperation, competition, or institutional relationships. Stratified spaces over joint

configuration spaces  $M = M_1 \times \dots \times M_N$  can represent bargaining boundaries, Nash equilibria, and coalition formation. Finite approximation theorems (3.9-3.11) extend to product spaces with sparse interaction graphs, preserving computational tractability. Full development of multi-agent DEME, including mechanism design and institutional hierarchies, is reserved for future work.

**Open problem.** Characterize conditions under which robust constraint enforcement over uncertainty sets  $X(w)$  remains decidable (or efficiently approximable) for definable fragments of SGE policies.

## 11 Conclusion

We have presented Stratified Geometric Ethics, a mathematical framework providing rigorous foundations for verifiable ethical reasoning. Our contributions include:

1. Showing stratified spaces are *natural minimal candidates* for ethical phenomena (Theorem 2.3).
2. A complete *representation theorem* with explicit axioms including locality and scale normalization (Theorem 4.3).
3. *Finite approximation theorems* enabling implementation (Theorems 3.9–3.11).
4. *Decidability results* for static verification, with temporal properties via model checking (Theorem 6.4, Corollary 6.5).
5. *Sample complexity bounds* for learning ethical content (Theorems 7.2–7.4).

SGE provides the theoretical foundation for the DEME 2.0 architecture [2], establishing that real-time, verifiable ethical governance is mathematically grounded. As AI systems take on consequential roles, such rigor becomes not a luxury but a necessity.

**Competing interests**

The author declares no competing interests.

**Author contributions**

A.B. performed all aspects of the work.

**Materials & Correspondence**

Correspondence and requests for materials should be addressed to A.B. (email: andrew.bond@sjsu.edu).

**Data availability**

No new datasets were generated or analysed during the current study. Data, where applicable for companion empirical demonstrations, are described in the accompanying DEME 2.0 manuscript (submitted).

**Code availability**

No production-ready software was used or produced as part of the current theoretical study. Code and implementation details for the companion DEME 2.0 demonstrations are described in the accompanying DEME 2.0 manuscript (submitted). A development repository is available at: <https://github.com/ahb-sjsu/erism-lib>

**References**

- [1] A. H. Bond. Differential geometry for moral alignment. Working paper, San José State University, 2024.
- [2] A. H. Bond. DEME 2.0: Real-time ethical governance for safety-critical autonomous systems. Submitted to Nature Machine Intelligence, 2025.
- [3] H. Whitney. Tangents to an analytic variety. *Annals of Mathematics*, 81(3):496–549, 1965.
- [4] R. Thom. Ensembles et morphismes stratifiés. *Bull. Amer. Math. Soc.*, 75(2):240–284, 1969.
- [5] L. van den Dries. *Tame Topology and O-minimal Structures*. Cambridge, 1998.
- [6] A. Tarski. *A Decision Method for Elementary Algebra and Geometry*. RAND, 1951.
- [7] S. Harris. *The Moral Landscape*. Free Press, 2010.