# Generalization as Simplification in Robotics: Opportunities for Structured Representations in Embodied AI Systems

Bjørn Remseth

Email: la3lma@gmail.com

*With AI assistance*

*Abstract*—The mathematician Alexandre Grothendieck famously reformulated entire mathematical domains by replacing narrow, technically complicated settings with sweeping abstractions that made deep theorems nearly tautological. A similar epistemic maneuver occurs in algorithmic breakthroughs such as the Fast Fourier Transform (FFT), which replaces expensive direct computation with a representation where operations become trivial. Recent work in 3D Gaussian Splatting applies an analogous shift in computer graphics: from heavy neural radiance fields to structured, sparse representations enabling real-time rendering. This paper examines robotics and embodied AI domains where similar representational reframings may yield significant computational and conceptual simplifications. For each domain, we analyze the current state of the art, identify sources of complexity, and outline opportunities where structured, geometry-aware representations could replace monolithic computational approaches.

*Index Terms*—robotics, abstraction, computational complexity, embodied AI, representation learning, motion planning, world models, model predictive control

version: 2025-11-20-10:19:36-CET-f320ad8-main

*Note on document creation: This document was created with AI assistance. See the "About This Document" section for details on the methodology and verification processes used.*

## I. Introduction

Jean Dieudonné, in his address for Grothendieck's Fields Medal in 1966, remarked that his colleague's mathematical motto could be described as *simplifier en généralisant*— "to simplify by generalizing." This captures a hallmark of Grothendieck's thinking: the willingness to adopt an abstract, and to outsiders even forbiddingly complicated framework, whenever this broader environment made the essential structures transparent and the theorems nearly inevitable.

Grothendieck himself described this method through the celebrated metaphor of the "rising sea," in which a hard mathematical "rock" is not smashed with tools, but gently dissolved as the conceptual water level rises and reveals the underlying unity of phenomena. As Schneps summarizes, Grothendieck regarded many special or concrete situations as *more complicated*, whereas general ones—despite appearances—were *conceptually simpler* and more natural.

This is not an isolated historical curiosity. The same epistemic transformation occurs in computational breakthroughs such as the Cooley–Tukey Fast Fourier Transform: by shifting into the frequency-domain representation, an $O(N^2)$ convolution becomes a pair of $O(N \log N)$ transforms and an elementwise multiplication, revealing a deep structure that was invisible in the direct formulation.

Likewise, the recent success of 3D Gaussian Splatting represents a conceptual shift: instead of modeling radiance fields with an over-general neural network, one adopts a structured, sparse, geometric representation that is already close to the thing being computed. This "representation-first" stance yields both conceptual clarity and orders-of-magnitude efficiency: real-time novel-view synthesis at 1080p.

This paper surveys a family of domains in robotics and embodied AI where similar representational reframings may offer substantial computational gains. Each example below is a potential site where abstractions may "rise like the sea," replacing heavy machinery with simple, fast, domain-attuned structures.

## II. Representational Shifts in Robotics and Embodied AI

The following sections examine domains where current solutions rely on monolithic, expensive, or universal approximators, and where alternative representations may yield transformations analogous to (1) the FFT's domain shift, and (2) Gaussian Splatting's conversion of general fields into sparse, geometrically structured primitives.

## III. World Models and Visuomotor Dynamics

### A. Current Situation and State of the Art

Contemporary approaches to world modeling for robotic systems predominantly rely on large neural architectures trained to predict future observations from current states and actions. The dominant paradigms include:

- **Recurrent architectures:** Models like World Models [1] use variational autoencoders (VAEs) combined with recurrent neural networks (RNNs) to learn compressed representations of visual observations and predict forward dynamics.
- **Transformer-based models:** Recent work such as Dreamer [2], [3] and IRIS [4] employ transformer architectures to model long-horizon dependencies in obser-

vation sequences, achieving impressive results in model-based reinforcement learning.

- **Diffusion models:** Emerging approaches like Diffusion Policy [5] and UniPi [6] use diffusion models to learn distributions over future trajectories, showing strong performance in manipulation tasks.

These methods have demonstrated success in complex control tasks, from simulated environments to real-world robotic manipulation. State-of-the-art systems can learn to predict hundreds of time steps into the future and support planning through learned dynamics.

### B. Sources of Complexity

Despite impressive capabilities, current world models face several sources of complexity:

*a) Computational Cost.:* Large transformer-based world models require significant computation for both training and inference. Models like DreamerV3 can have tens of millions of parameters, and forward prediction through these models is computationally expensive, limiting real-time performance on resource-constrained robots.

*b) Dense Representation.:* Most current approaches represent the entire scene as a dense feature vector or pixel grid. This fails to exploit:

- **Sparsity:** Most of the environment is static; only a few objects move or require prediction.
- **Locality:** Object interactions are typically local; distant objects don't affect each other.
- **Compositionality:** Scenes are composed of discrete objects with independent dynamics.

*c) Generalization Challenges.:* Monolithic neural world models struggle to generalize to:

- Novel object configurations
- Different numbers of objects than seen during training
- New environments with different visual appearances

*d) Historical Context.:* These approaches emerged when:

- Large-scale neural network training became feasible (post-2015)
- End-to-end learning paradigms dominated computer vision
- Structured representations were seen as limiting generality

The complexity is thus partly *useful* (neural networks are flexible function approximators) and partly *historical* (following the success of end-to-end learning in other domains).

### C. Opportunities for Simplification

*a) Object-Centric Representations.:* Replace monolithic scene representations with explicit object-centric models. Recent work in this direction includes:

- **Slot attention:** Methods like SAVi [7] decompose scenes into object slots, each represented by a learned feature vector.
- **Neural scene graphs:** Represent scenes as graphs where nodes are objects and edges encode relationships [8].

- **Gaussian object representations:** Extending 3D Gaussian Splatting [9] to dynamic scenes, where each object is a collection of Gaussians with learned dynamics.

*b) Event-Driven Dynamics.:* Instead of predicting full scenes at every time step, maintain an implicit scene representation and only update:

- Objects that are moving or being manipulated
- Regions where interactions occur (contacts, collisions)
- Parts of the scene within the robot's attention or planning horizon

This is analogous to sparse matrix computations: most elements are zero (static), so only track the non-zero (dynamic) components.

*c) Geometric Primitives.:* Represent objects and their dynamics using structured geometric primitives:

- Anisotropic 3D Gaussians (position, orientation, covariance)
- Superquadrics or convex polytopes
- Signed distance functions (SDFs) with learned time derivatives
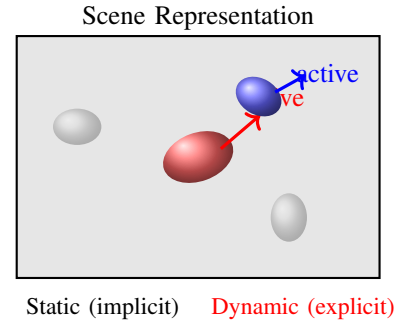


Scene Representation

Fig. 1. *Object-centric world model with event-driven updates. Static scene elements (gray) are maintained implicitly, while only dynamic objects (colored) require explicit prediction and updates.*

*d) Hybrid Approaches.:* Combine learned components with physics-based simulation:

- Learn object segmentation and parameter identification (mass, friction, shape)
- Use fast analytic physics engines for rigid body dynamics
- Apply learned corrections for contact-rich or deformable scenarios

This exploits known structure (Newtonian mechanics) while using learning where physics-based models are insufficient.

### D. Potential Impact

If successful, these simplifications could enable:

- Real-time model-predictive control at 100+ Hz on edge devices
- Better generalization to novel scenes and object counts
- Interpretable world representations for debugging and safety verification
- Scaling to longer prediction horizons with constant per-object cost

## IV. Motion Planning via Configuration-Space Representations

### A. Current Situation and State of the Art

Motion planning remains a computational bottleneck in robotics, particularly for high-DOF manipulators and mobile robots in cluttered environments. Current state-of-the-art approaches include:

- **Sampling-based planners:** RRT [10], RRT* [11], and PRM [12] build graphs in configuration space through random sampling. These are probabilistically complete but require many collision checks.
- **Trajectory optimization:** Methods like CHOMP [13], TrajOpt [14], and STOMP [15] formulate planning as optimization over continuous trajectory spaces, often using gradient information.
- **Learning-based planning:** Neural motion planners like MPNet [16] and motion policy networks [17] learn to predict collision-free paths directly from experience.

State-of-the-art systems can plan complex manipulation sequences, but often require seconds to minutes for challenging problems, limiting reactive replanning capabilities.

### B. Sources of Complexity

*a) High-Dimensional Configuration Spaces.:* A 7-DOF manipulator has a 7-dimensional C-space; dual-arm systems have 14+ dimensions. The volume of the free configuration space grows exponentially, making exhaustive search intractable.

*b) Complex Collision Geometry.:* Each sampled configuration requires checking collisions between robot links and environment obstacles. For mesh-based geometries, this involves expensive proximity queries. A single RRT query might perform 10,000+ collision checks.

*c) Frequent Replanning.:* In dynamic environments or for reactive tasks, the robot must replan at high frequency (10-50 Hz). Current planners cannot reliably meet these real-time constraints for complex scenes.

*d) Historical Development.:* Sampling-based planners emerged in the 1990s when:

- Exact cell decomposition was computationally infeasible for high-DOF systems
- Randomization provided an elegant probabilistic completeness guarantee
- Collision detection libraries were the main computational bottleneck

### C. Opportunities for Simplification

*a) Configuration-Space Splatting.:* Precompute or learn a representation of the free configuration space as a collection of local "safe regions" or "motion primitives," analogous to Gaussian splats in rendering:

- **Local convex polytopes:** Represent free C-space as a union of convex regions, each with a center configuration and feasible deviation [18].

- **Motion funnels:** For each goal region, precompute a family of funnels (tubes in C-space) that guarantee convergence via local feedback control [19].
- **Gaussian process implicit surfaces:** Learn the boundary of C-obstacle regions using Gaussian processes, enabling fast queries and gradient information [20].
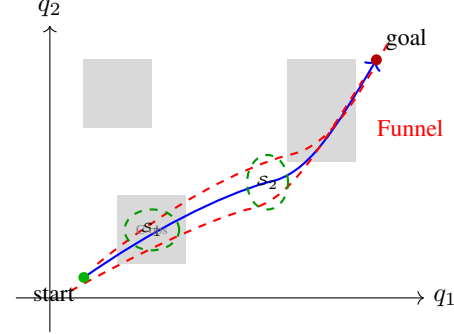


Fig. 2. *Motion funnel in configuration space. Precomputed safe tubes (red dashed boundaries) guide the robot from start to goal, with local safe regions $\mathcal{S}_i$ (green) allowing efficient collision-free motion. Gray rectangles represent C-space obstacles.*

*b) Planning as Lookup and Interpolation.:* With precomputed representations:

1) **Query:** Given current configuration $q_{\text{curr}}$ and goal $q_{\text{goal}}$, find nearby safe regions or funnels in a learned database.
2) **Compose:** Chain together local motion primitives or funnels.
3) **Refine:** Apply small local optimization to smooth transitions and ensure dynamic feasibility.

This shifts expensive online search to offline precomputation, analogous to how radiance caching accelerates rendering.

*c) Neural Implicit Representations.:* Learn neural networks that map directly from $(q, \text{scene})$ to:

- Collision probability or signed distance to nearest obstacle
- Gradient of collision cost (for gradient-based optimization)
- Locally optimal motion direction (as in potential fields)

Recent work like CollisionNet [21] shows promise, achieving 100x speedups over traditional collision checking.

### D. Potential Impact

- Planning times reduced from seconds to milliseconds
- Replanning at control frequency (50-200 Hz)
- Graceful degradation: fall back to sampling-based methods when learned representations are uncertain
- Compositional planning: combine precomputed motion primitives for novel tasks

## V. Model Predictive Control via Basis Representations

### A. Current Situation and State of the Art

Model Predictive Control (MPC) has become a cornerstone of robotic control, particularly for systems with complex

dynamics and constraints. State-of-the-art implementations include:

- **Nonlinear MPC:** Methods like Sequential Quadratic Programming (SQP) [22] and interior-point methods solve nonlinear optimal control problems at each time step.
- **Learning-based MPC:** Approaches like Differentiable MPC [23] and neural MPC [24] use learned dynamics models and/or warm-starting from learned policies.
- **Whole-body control:** Systems like TOWR [25] and Crocoddyl [26] optimize over full robot state trajectories including contact sequences.

Modern MPC can control quadrupeds performing dynamic locomotion, dexterous manipulation, and aerial vehicles, but at significant computational cost.

### B. Sources of Complexity

*a) Online Optimization Burden.:* Each MPC iteration solves a trajectory optimization problem:

$$\min_{\mathbf{u}_{0:H}} \quad \sum_{t=0}^{H} \ell(x_t, u_t) + \ell_f(x_H) \tag{1}$$
$$\text{s.t.} \quad x_{t+1} = f(x_t, u_t)$$
$$g(x_t, u_t) \leq 0$$

where $H$ is the planning horizon, $\ell$ is a stage cost, $f$ is the dynamics model, and $g$ represents constraints. For nonlinear systems, this requires iterative optimization with complexity $O(H \cdot n_x \cdot n_u \cdot I)$ where $I$ is the number of iterations.

*b) Real-Time Requirements.:* Control loops must run at 100-1000 Hz depending on the task. For a 20-step horizon with a 7-DOF manipulator, the optimization may have hundreds of variables. Achieving real-time performance often requires:

- Simplified dynamics models
- Conservative constraint handling
- Early termination of optimization
- Powerful embedded computers

*c) Historical Context.:* MPC emerged from process control in the 1970s-80s when:

- Linear-quadratic methods were tractable
- Nonlinear optimization was expensive but worth it for high-value applications
- Real-time embedded optimization was infeasible, limiting applications

Advances in optimization and computing have made real-time nonlinear MPC possible, but it remains computationally demanding.

### C. Opportunities for Simplification

*a) Trajectory Basis Functions.:* Instead of optimizing over raw control sequences $\mathbf{u}_{0:H}$, represent control trajectories as linear combinations of basis functions:

$$u_t = \sum_{i=1}^{K} \alpha_i \phi_i(t)$$

where $\{\phi_i\}$ are basis functions (e.g., B-splines, Fourier series, dynamic movement primitives) and $\alpha_i$ are coefficients

to optimize. This reduces the decision variables from $H \cdot n_u$ to $K \cdot n_u$ where $K \ll H$.

*b) Local Control Primitives.:* Partition the state-action space and learn local controllers ("control splats") valid in each region:

- Each primitive is a simple controller (e.g., linear gain, polynomial trajectory)
- MPC optimization selects and blends a small number of primitives
- Analogous to mixture-of-experts: sparse activation reduces computation

*c) Learned Value Functions for Warm Starting.:* Use a learned value function $V_\theta(x)$ to:

- Provide initial trajectory guess via rollout of $\nabla V_\theta$
- Set terminal cost in finite-horizon MPC
- Decide when to terminate optimization early (if learned policy is already near-optimal)

Recent work [27], [28] shows this can reduce MPC iterations by 5-10x.
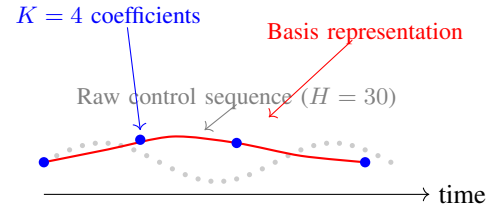


Fig. 3. *Basis function parameterization for MPC. Instead of optimizing over $H = 30$ discrete control values (gray points), optimize over $K = 4$ basis function coefficients (blue points), with the full trajectory (red curve) reconstructed via interpolation.*

*d) Hierarchical MPC.:* Separate planning into:

1) **High level:** Long-horizon, coarse plan using simplified dynamics (computed infrequently, 1-10 Hz)
2) **Low level:** Short-horizon, detailed control with full dynamics (computed at control rate, 100+ Hz)

The high-level plan provides reference trajectories and constraint boundaries for the low-level controller.

### D. Potential Impact

- MPC feasible on resource-constrained robots (drones, legged robots)
- Longer planning horizons (50-100 steps) without increased computation
- More complex dynamics models (contact-rich, fluid interactions)
- Faster recovery from disturbances through rapid replanning

## VI. PERCEPTION AS SPARSE TASK-RELEVANT FIELDS

### A. Current Situation and State of the Art

Modern robotic perception systems increasingly rely on large foundation models and general-purpose architectures:

- **Vision transformers:** Models like CLIP [29], DINOv2 [30], and SAM [31] provide rich visual representations learned from massive datasets.

- **Multi-modal models:** Systems like RT-2 [32] and PaLM-E [33] combine vision and language for general-purpose robot control.
- **Dense prediction:** Methods predict per-pixel depth, surface normals, semantic segmentation, and object affordances using U-Net or transformer architectures.

These approaches achieve impressive generalization across visual domains and tasks, enabling robots to handle diverse objects and environments.

### B. Sources of Complexity

*a) Computational Overhead.:* Vision transformers process images with $O(N^2)$ complexity in the number of patches. A 224×224 image at 16×16 patches requires attention over 196 tokens, repeated across many layers. Running models like SAM or CLIP on every frame at control frequency (30-60 Hz) is challenging even on powerful GPUs.

*b) Dense Representations.:* Most vision models produce dense feature maps over the entire image, but robotic tasks typically require only sparse information:
- Grasping needs keypoints, approach vectors, and local geometry
- Navigation needs free-space boundaries and traversability
- Manipulation needs object poses, contact surfaces, and motion affordances

Computing full dense representations wastes computation on task-irrelevant regions.

*c) End-to-End Learning Paradigm.:* Current methods emerged from the success of:
- ImageNet-scale supervised learning (2010s)
- Self-supervised learning from web-scale data (2020s)
- Scaling laws suggesting bigger models = better performance

This paradigm prioritizes generality over efficiency, which is appropriate for foundation models but may be suboptimal for resource-constrained robotic deployment.

### C. Opportunities for Simplification

*a) Task-Specific Field Extraction.:* Learn lightweight modules that extract only task-relevant geometric fields from visual input:
- **For navigation:** Free-space distance field (2D or 3D SDF indicating closest obstacle)
- **For grasping:** Grasp affordance field (6-DOF grasp poses per surface point)
- **For manipulation:** Contact normals, friction cones, and support surfaces

These representations are sparse (defined only on object surfaces or task-relevant regions) and structured (obeying geometric constraints).

*b) Geometric Reasoning Pipelines.:* Once sparse fields are extracted, use fast geometric algorithms:
- **Collision-free motion:** Query SDF for swept-volume collision checking
- **Force closure analysis:** Evaluate grasp stability using contact geometry and friction cones

- **Path planning:** Use distance gradients for potential-field navigation

This separates *perception* (learned, runs at perception frequency) from *geometric reasoning* (analytic, runs at control frequency).
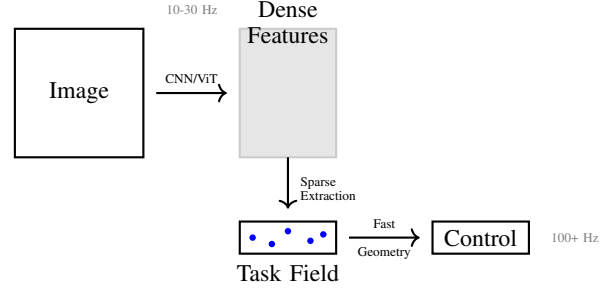


Fig. 4. *Perception pipeline with sparse task-field extraction. Dense visual features are computed at perception frequency (10-30 Hz), then compressed to sparse task-relevant fields. Fast geometric reasoning operates at control frequency (100+ Hz) using these lightweight representations.*

*c) Hybrid Learned-Geometric Modules.:* Combine learning for perception with classical geometry for reasoning:
- Learn object segmentation and category recognition (when geometric priors are weak)
- Fit geometric primitives (planes, cylinders, boxes) to segmented point clouds
- Use analytic kinematics and contact mechanics for motion planning

Recent work like Neural Descriptor Fields [34] and NeRF-based scene understanding [35] shows promise in this direction.

*d) Active Perception.:* Since task fields are sparse, guide visual attention and sensor queries:
- Foveated vision: high-resolution sensing only in task-relevant regions
- Viewpoint planning: move cameras to observe critical surfaces
- Predictive sensing: only re-sense when predicted dynamics diverge from observations

### D. Potential Impact

- Real-time perception on edge devices (mobile robots, drones)
- Reduced communication bandwidth (transmit sparse fields, not full images)
- Interpretable representations for debugging and safety analysis
- Compositional generalization: geometric reasoning transfers across environments

## VII. SIMULATION WITH DIFFERENTIABLE GAUSSIAN PROXIES

### A. Current Situation and State of the Art

Physics simulation plays a crucial role in modern robotics for:

- **Training:** Sim-to-real transfer for reinforcement learning policies
- **Planning:** Model-predictive control with physics-based dynamics
- **Validation:** Testing safety-critical behaviors before deployment

Current state-of-the-art simulators include:

- **Rigid body engines:** MuJoCo [36], PyBullet, Isaac Sim provide fast rigid-body dynamics with contact
- **Differentiable simulators:** DiffTaichi [37], Brax [38], and Warp enable gradient-based optimization through simulation
- **Hybrid methods:** Neural simulators like Graph Network Simulators [39] learn dynamics corrections

Despite advances, high-fidelity simulation remains expensive when combining:

- Detailed collision geometry (meshes, not primitives)
- Realistic contact mechanics (friction, soft contacts)
- Visual rendering for perception pipelines
- Deformable or fluid dynamics

### B. Sources of Complexity

*a) Collision Detection.:* Checking for collisions between complex mesh geometries is expensive:

- Mesh-mesh collision: $O(n \cdot m)$ for $n$ and $m$ triangles (with broad-phase acceleration)
- Contact resolution: iterative constraint solving for each contact point
- Continuous collision detection: check swept volumes between time steps

*b) Contact Dynamics.:* Modeling realistic contact is notoriously difficult:

- Contact is a discontinuous constraint (complementarity problem)
- Friction cones lead to nonlinear, non-smooth equations
- Stiff contact forces require small time steps for stability

*c) Rendering for Perception.:* If the robot's policy uses visual observations, simulation must render images:

- Rasterization or ray tracing for each camera view
- Realistic lighting and materials for sim-to-real transfer
- Sensor noise and artifacts (motion blur, rolling shutter)

Combined simulation + rendering + policy inference can run slower than real-time even on GPUs.

*d) Historical Context.:* Physics engines were designed for:

- Animation and games (visual plausibility > physical accuracy)
- Offline engineering analysis (accuracy prioritized over speed)
- Simple primitive shapes (boxes, spheres, cylinders)

Modern robotics demands both accuracy and massive parallelization for RL training, pushing these systems to their limits.

### C. Opportunities for Simplification

*a) Gaussian Splat Physics Proxies.:* Represent objects as collections of 3D Gaussians (extending Gaussian Splatting [9]):

- **Geometry:** Each Gaussian has position $\mu \in \mathbb{R}^3$, covariance $\Sigma \in \mathbb{R}^{3 \times 3}$
- **Physics:** Associate mass, moment of inertia, and material properties with each Gaussian
- **Contacts:** Detect overlap between Gaussians (fast closed-form test), compute contact forces proportional to penetration depth

This representation is:

- **Sparse:** Objects decompose into ~100-1000 Gaussians (vs. millions of mesh triangles)
- **Differentiable:** All operations (overlap, forces) have analytic derivatives
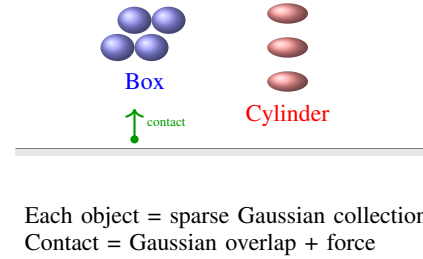- **Fast:** Gaussian-Gaussian overlap is a closed-form expression



Each object = sparse Gaussian collection
Contact = Gaussian overlap + force

Fig. 5. *Objects represented as sparse collections of 3D Gaussians for fast physics simulation. Contact detection reduces to Gaussian overlap tests (closed-form), and contact forces are proportional to penetration depth.*

*b) Coupled Rendering and Physics.:* Since Gaussian Splatting already provides fast differentiable rendering [9], we can:

1) Use the same Gaussian representation for both physics and rendering
2) Render images directly from physics state at minimal cost
3) Enable gradient flow from visual loss through physics to control

This unifies simulation and perception pipelines, avoiding expensive mesh-to-image rendering.

*c) Learned Dynamics Corrections.:* Use Gaussian proxy physics as a fast *base simulator*, then learn:

- Residual dynamics $\dot{x} = f_{\text{gauss}}(x, u) + f_\theta(x, u)$
- Contact model corrections (friction, restitution)
- Material property identification (stiffness, damping)

Train on small amounts of real-world data or high-fidelity sim data. This trades perfect physical accuracy for massive speedup, with learned corrections bridging the gap.

*d) Hierarchical Simulation.:* Combine multiple fidelity levels:

- **Coarse:** Gaussian proxy physics (1000+ Hz)
- **Medium:** Rigid body dynamics with primitives (100-500 Hz)
- **Fine:** Full mesh-based simulation with FEM (1-10 Hz for validation)

6

Run coarse simulation for RL training and MPC, occasionally validate with fine simulation.

### D. Potential Impact

- 10-100× faster simulation for RL training
- Unified physics + rendering for end-to-end differentiable pipelines
- Real-time simulation on edge devices for online MPC
- Scalability: simulate 1000s of objects at interactive rates

## VIII. MULTI-ROBOT COORDINATION VIA POTENTIAL FIELDS

### A. Current Situation and State of the Art

Multi-robot systems are increasingly deployed for:
- Warehouse automation (Kiva/Amazon robots, AutoStore)
- Drone swarms for search, surveillance, or shows
- Multi-arm manipulation systems
- Autonomous vehicle fleets

Current coordination approaches include:
- **Centralized optimization:** Solve a joint optimization problem for all robots [40]. Optimal but scales poorly ($O(N^3)$ or worse for $N$ robots).
- **Priority-based planning:** Assign priorities, plan sequentially with collision avoidance [41]. Fast but suboptimal and prone to deadlocks.
- **Velocity obstacles:** Distributed geometric algorithm where each robot locally avoids collisions [42]. Scales well but lacks global coordination.
- **Learning-based:** Train policies via MARL (multi-agent RL) to coordinate implicitly [43]. Flexible but requires extensive training and may not provide guarantees.

State-of-the-art systems can coordinate dozens to hundreds of robots, but struggle with:
- Dense environments with many interactions
- Real-time replanning under dynamic changes
- Formal safety and deadlock-freedom guarantees

### B. Sources of Complexity

*a) Combinatorial Coordination.:* With $N$ robots, the joint configuration space is $N$-times higher dimensional than single-robot planning. The number of possible collision pairs grows as $O(N^2)$.

*b) Communication and Computation.:* Centralized approaches require:
- All robots to communicate state to a central planner
- Central planner to solve large optimization problem
- Plans to be communicated back to robots

This creates bandwidth and latency bottlenecks. Distributed approaches avoid this but require careful design to prevent conflicts.

*c) Dynamic Environments.:* Robots must replan continuously as:
- Other robots move (predicted but uncertain)
- Goals change (task allocation updates)
- Obstacles appear (humans, objects)

*d) Historical Context.:* Multi-robot coordination emerged from:
- Single-robot planning (extend to multiple agents)
- Distributed systems (consensus, leader election)
- Game theory (Nash equilibria, mechanism design)

Each perspective brings complexity: planning suffers from dimensionality, distributed systems from asynchrony, and game theory from strategic behavior.
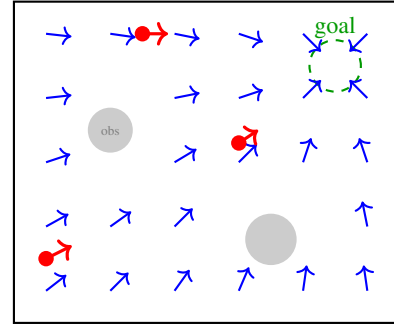
### C. Opportunities for Simplification

*a) Workspace Potential Fields.:* Represent the coordination problem as *vector fields* over the workspace:
- **Goal attraction:** Field guiding robots toward their goals
- **Obstacle repulsion:** Field pushing robots away from static obstacles
- **Inter-robot repulsion:** Dynamic field preventing collisions between robots
- **Flow lanes:** Structured fields encouraging predictable motion patterns (e.g., right-hand rule in corridors)

Each robot queries the field at its current position to determine local velocity:

$$v_i(t) = k_{\text{goal}} \cdot \nabla \phi_{\text{goal}}(p_i) + k_{\text{avoid}} \cdot \sum_{j \neq i} \nabla \phi_{\text{avoid}}(p_i, p_j) + \nabla \phi_{\text{obstacle}}(p_i)$$



Shared potential field guides all robots

Fig. 6. *Multi-robot coordination via shared potential field. Blue arrows show the field direction guiding robots (red circles) toward the goal (green) while avoiding obstacles (gray) and each other. Each robot computes local velocity from field queries.*

*b) Benefits of Field Representation.:*
- **Scalability:** Each robot's computation is local ($O(1)$ or $O(\log N)$ with spatial indexing)
- **Reactivity:** Fields update continuously; robots respond immediately
- **Predictability:** Smooth fields produce smooth motions, aiding prediction
- **Decentralization:** Robots query shared field (broadcast or cached locally)

*c) Learning Field Parameters.:* Classical potential fields suffer from local minima and parameter tuning. Modern approach:
- Learn field representations via multi-agent RL [44]
- Train neural networks to output field gradients given global state

- Use differentiable trajectory optimization to tune field weights for task success

This combines the structure of geometric fields with the flexibility of learning.

*d) Hierarchical Coordination.:* Layer multiple coordination mechanisms:
1) **Global:** Task allocation and high-level flow patterns (computed centrally, updated slowly)
2) **Local:** Collision avoidance and trajectory smoothing via potential fields (computed per-robot, updated at control rate)
3) **Reactive:** Emergency stops and safety overrides (reflex-level, <1ms latency)

### D. Potential Impact

- Coordinate 100s-1000s of robots in real-time
- Reduce communication: broadcast field updates, not individual plans
- Provable safety: carefully designed fields guarantee collision avoidance
- Emergent behaviors: complex global patterns from simple local rules

## IX. VISUOMOTOR POLICIES AS LOCALIZED SKILL FIELDS

### A. Current Situation and State of the Art

End-to-end visuomotor policies have become increasingly capable:

- **Behavior cloning:** Learn policies $\pi_\theta(a \mid o)$ from demonstrations. Systems like RT-1 [45], RT-2 [32] train on 100k+ robot trajectories.
- **Foundation models:** Large vision-language-action models generalize across tasks and environments via pre-training on diverse data.
- **Diffusion policies:** Model action distributions with diffusion models [5], achieving state-of-the-art performance on manipulation benchmarks.

These monolithic policies can achieve impressive generalization, but face challenges:

- **Sample efficiency:** Require massive datasets (100k+ demonstrations)
- **Compositionality:** Difficult to combine or sequence learned behaviors
- **Interpretability:** Black-box policies are hard to debug or verify
- **Computation:** Large models (100M+ parameters) are slow on robot hardware

### B. Sources of Complexity

*a) Monolithic Architecture.:* Current policies map directly from raw observations to actions:

$$o_t \xrightarrow{\text{CNN/ViT}} \text{features} \xrightarrow{\text{RNN/Transformer}} a_t$$

This requires the network to learn:
- Perception (what objects are present and where)
- Task understanding (what to do)
- Motor control (how to move)

- Coordination (sequencing primitives)

All of these are entangled in a single large network.

*b) State-Action Space Coverage.:* A single policy must handle the entire state-action space:
- Different object poses and configurations
- Various stages of task execution (approach, grasp, transport)
- Multiple solution strategies (different grasps, paths)

This requires large capacity and extensive training data.

*c) Historical Context.:* End-to-end learning emerged from:
- Success of deep learning in vision (ImageNet, 2012)
- Advances in imitation learning and DAGGER (2011-2015)
- Scaling laws: bigger models + more data = better policies (2020s)

The paradigm prioritizes generality and has been hugely successful, but may not be optimal for all deployment scenarios.

### C. Opportunities for Simplification

*a) Skill Decomposition.:* Decompose complex tasks into *local controllers* or *skills*, each valid in a region of state-task space:
- **Approach skill:** Move end-effector to pre-grasp pose
- **Grasp skill:** Close gripper with appropriate force
- **Transport skill:** Move grasped object to target
- **Place skill:** Release object gently

Each skill is a simpler controller (possibly even linear or geometric) valid locally.

*b) Skill Fields Representation.:* Represent each skill as a "field" over state space:
- **Applicability field:** $\alpha_i(s)$ indicates where skill $i$ is valid
- **Value field:** $V_i(s)$ estimates expected success from state $s$ using skill $i$
- **Action field:** $\pi_i(s)$ produces actions in the skill's domain

At execution time, blend skills based on applicability:

$$a_t = \frac{\sum_i \alpha_i(s_t) \cdot \pi_i(s_t)}{\sum_i \alpha_i(s_t)}$$

*c) Benefits of Skill Decomposition.:*
- **Sample efficiency:** Each skill is simpler, requires less data
- **Compositionality:** Combine skills for new tasks without retraining
- **Interpretability:** Understand which skill is active and why
- **Efficiency:** Activate and compute only relevant skills (sparse activation)

*d) Learning Skill Decompositions.:* Several approaches to discover skills:
- **Hierarchical RL:** Learn high-level policy over skills and low-level skill controllers [46]
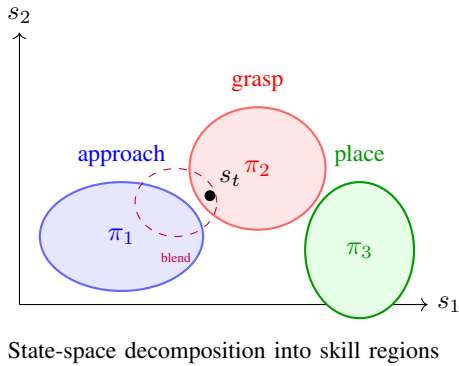- **Option discovery:** Unsupervised methods identify bottleneck states and useful primitives [47]

State-space decomposition into skill regions

Fig. 7. *Localized skill fields decompose the state-task space. Each skill $\pi_i$ (colored regions) is a simple controller valid locally. In overlap regions (purple), skills are blended based on applicability weights $\alpha_i(s)$.*

- **Segmentation from demonstrations:** Cluster or change-point detection on demonstration trajectories to identify skill boundaries [48]
- **Task-parameterized skills:** Learn skills with parameters (grasp pose, target location) using DMPs or ProMPs [49]
  *e) Hybrid Learned-Geometric Skills.:* Some skills can be geometric:
  - **Approach:** Plan collision-free path to pre-grasp (use motion planning from Sec. IV)
  - **Grasp:** Evaluate grasp quality using force closure (geometric computation)
  - **Transport:** Follow minimum-jerk trajectory with impedance control

Only learn skills where geometric methods are insufficient (e.g., contact-rich manipulation, deformable objects).

*D. Potential Impact*

- 10-100× reduction in required training data per task
- Real-time policy inference: evaluate only active skills
- Modular development: update individual skills without full retraining
- Formal verification: prove properties of geometric skill components
- Transfer: share skill libraries across robot platforms and tasks

## X. A GENERAL MECHANISM: THE SEA THAT RISES

Across these examples, we observe a common pattern. An initial representation—dense world models, explicit configuration-space search, monolithic MPC optimization, dense visual features, mesh-based simulation, centralized multi-robot planning, end-to-end policies—is treated as the "rock" of the problem. Researchers first build heavy machinery to attack it directly, establishing correctness and feasibility.

But over time, a new representational layer emerges in which the underlying computation becomes simple or even trivial. This is Grothendieck's "rising sea": the complexity has not been removed, but dissolved into a higher-level abstraction.

In robotics and embodied AI, this abstraction process is not a detour but the natural path forward. World models, planners,

and controllers evolve from direct numerical computation toward structured fields; from black-box approximation toward geometric representations that match the latent structure of the physical world. The opportunities listed in this paper highlight domains where such rising abstractions may already be imminent.

*A. Common Simplification Patterns*

Several themes recur across domains:

**Sparsity and locality.** Most information is irrelevant; most interactions are local. Represent only active elements (moving objects, task-relevant regions, applicable skills) and compute only necessary interactions.

**Geometric structure.** Physical systems obey geometric laws. Structured representations (Gaussians, polytopes, fields) expose this structure, enabling fast analytic operations instead of expensive learned inference.

**Hierarchy and composition.** Complex behaviors decompose into simpler primitives. Hierarchical representations (coarse-to-fine, goal-to-action, task-to-skill) reduce problem dimensionality at each level.

**Separation of concerns.** Decouple perception from control, learning from reasoning, online from offline computation. Each component can be optimized independently and specialized for its role.

**Event-driven updates.** Don't recompute everything every time step. Maintain implicit representations and update only when events occur (contacts, state transitions, observations).

**Learned priors with geometric reasoning.** Use learning where data is rich and structure is weak (perception, initialization); use geometry where structure is strong and computation is fast (collision checking, force analysis, trajectory integration).

*B. The Path Forward*

The transformation induced by Gaussian Splatting in neural rendering is a paradigmatic case of this broader research direction: finding the representation in which a problem becomes simple. This direction is not merely aesthetic; it is computational pragmatism. Following Grothendieck, the goal is not to complicate—it is to find the level at which the structure of the world becomes transparent.

The robotics and AI domains surveyed here present fertile ground for similar representational breakthroughs. Whether through object-centric world models, configuration-space primitives, basis-field MPC, sparse perception, Gaussian proxy simulation, potential-field coordination, or skill decomposition, each domain invites a shift toward structured geometry and away from monolithic generality.

The sea is rising; the rocks will dissolve.

## XI. CONCLUSION

This paper has examined seven domains in robotics and embodied AI where current approaches rely on computationally expensive, monolithic methods. For each, we have analyzed

the state of the art, identified sources of complexity, and outlined opportunities for simplification through structured, geometry-aware representations.

The recurring pattern is one of representational abstraction: replacing dense, general-purpose approximators with sparse, structured representations that expose and exploit the underlying geometry of physical systems. This is analogous to classic breakthroughs like the FFT (shifting to a representation where operations are trivial) and recent successes like 3D Gaussian Splatting (structured primitives enabling real-time rendering).

These simplifications are not mere optimizations; they represent fundamental shifts in how we model, perceive, plan, and control robotic systems. If successful, they promise:

- Orders-of-magnitude computational speedups
- Better generalization through compositional structure
- Interpretability for debugging and verification
- Real-time performance on resource-constrained hardware

The challenges are significant: learning structured representations, ensuring they capture relevant phenomena, and proving safety guarantees. But the potential rewards—both scientific and practical—make this a compelling research direction.

As Grothendieck's "rising sea" metaphor suggests, the most powerful simplifications often come not from attacking problems directly, but from raising the level of abstraction until their solutions become natural and inevitable. In robotics, the sea is beginning to rise.

## ACKNOWLEDGMENTS

## ABOUT THIS DOCUMENT

This document represents a technical exploration created through a collaborative process between human and AI. The production process followed these steps:

1) **Discovery**: The topic emerged from examining patterns across multiple computational domains where representational shifts enabled dramatic simplifications, particularly the recent success of 3D Gaussian Splatting in computer graphics.
2) **Initial Exploration**: Through dialogue with AI systems (ChatGPT and Claude), we explored how similar patterns might apply to robotics domains, drawing on concepts from Grothendieck's mathematical philosophy and algorithmic breakthroughs like the FFT.
3) **Synthesis**: Multiple draft documents were created, refined, and merged to develop detailed analyses of each robotics domain, structured to present: (1) current state of the art, (2) sources of complexity, and (3) opportunities for structured simplification.
4) **Visualization**: TikZ diagrams were created inline within the LaTeX document to illustrate key concepts: object-centric world models, configuration-space funnels, basis-function MPC, sparse perception pipelines, Gaussian physics proxies, potential field coordination, and skill field decomposition.
5) **Verification**: All references are documented with verification status. URLs were tested for accessibility, and content relevance was checked against citations.

This document aims to provoke thought and discussion about how abstraction and structured representations might advance robotics research, building on successful patterns from mathematics, algorithms, and computer graphics.

REFERENCES

[1] D. Ha and J. Schmidhuber, "World models," *arXiv preprint arXiv:1803.10122*, 2018, verified: arXiv accessible - Introduces generative neural network models for RL environments. Also published at NeurIPS 2018. [Online]. Available: https://arxiv.org/abs/1803.10122

[2] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi, "Dream to control: Learning behaviors by latent imagination," in *International Conference on Learning Representations (ICLR)*, 2020, verified: arXiv accessible - Introduces Dreamer, model-based RL via latent imagination. [Online]. Available: https://arxiv.org/abs/1912.01603

[3] D. Hafner, T. Lillicrap, M. Norouzi, and J. Ba, "Mastering atari with discrete world models," *arXiv preprint arXiv:2010.02193*, 2021, verified: arXiv accessible - DreamerV2 for discrete action spaces. [Online]. Available: https://arxiv.org/abs/2010.02193

[4] V. Micheli, E. Alonso, and F. Fleuret, "Transformers are sample-efficient world models," in *International Conference on Learning Representations (ICLR)*, 2023, verified: arXiv accessible - IRIS model for transformer-based world modeling. [Online]. Available: https://arxiv.org/abs/2209.00588

[5] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song, "Diffusion policy: Visuomotor policy learning via action diffusion," in *Robotics: Science and Systems (RSS)*, 2023, verified: arXiv accessible - Diffusion models for visuomotor policy learning. [Online]. Available: https://arxiv.org/abs/2303.04137

[6] Y. Du, S. Shen, B. Huang, X. Tian, H.-X. Wu, J. Liang, D. Zhou, and J. B. Tenenbaum, "Learning universal policies via text-guided video generation," *arXiv preprint arXiv:2302.00111*, 2023, verified: arXiv accessible - UniPi for text-guided video generation in robotics. [Online]. Available: https://arxiv.org/abs/2302.00111

[7] T. Kipf, G. F. Elsayed, A. Mahendran, A. Stone, S. Sabour, G. Heigold, R. Jonschkowski, A. Dosovitskiy, and K. Greff, "Conditional object-centric learning from video," in *International Conference on Learning Representations (ICLR)*, 2022, verified: arXiv accessible - SAVi model for object-centric video learning with slot attention. [Online]. Available: https://arxiv.org/abs/2111.12594

[8] C. Eppner, A. Mousavian, and D. Fox, "Physically consistent scene graphs for robot manipulation," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 3684–3691, 2021, verified: DOI accessible - Neural scene graphs for robotic manipulation. [Online]. Available: https://doi.org/10.1109/LRA.2021.3064269

[9] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics*, vol. 42, no. 4, 2023, verified: DOI and arXiv accessible - SIGGRAPH 2023, introduces 3D Gaussian Splatting for real-time rendering. [Online]. Available: https://arxiv.org/abs/2308.04079

[10] S. M. LaValle, "Rapidly-exploring random trees: A new tool for path planning," Computer Science Department, Iowa State University, Tech. Rep. TR 98-11, 1998, verified: PDF accessible - Original RRT paper, foundational work in sampling-based planning. [Online]. Available: https://msl.cs.illinois.edu/~lavalle/papers/Lav98c.pdf

[11] S. Karaman and E. Frazzoli, "Sampling-based algorithms for optimal motion planning," *The International Journal of Robotics Research*, vol. 30, no. 7, pp. 846–894, 2011, verified: DOI and arXiv accessible - Introduces RRT* and PRM* with asymptotic optimality guarantees. [Online]. Available: https://arxiv.org/abs/1105.1186

[12] L. E. Kavraki, P. Svestka, J.-C. Latombe, and M. H. Overmars, "Probabilistic roadmaps for path planning in high-dimensional configuration spaces," in *IEEE International Conference on Robotics and Automation (ICRA)*, 1996, pp. 566–580, verified: DOI accessible - Original PRM paper, foundational work in sampling-based planning. [Online]. Available: https://doi.org/10.1109/ROBOT.1996.503713

[13] N. Ratliff, M. Zucker, J. A. Bagnell, and S. Srinivasa, "Chomp: Gradient optimization techniques for efficient motion planning," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2009, pp. 489–494, verified: DOI accessible - Covariant Hamiltonian Optimization for Motion Planning. [Online]. Available: https://doi.org/10.1109/ROBOT.2009.5152817

[14] J. Schulman, Y. Duan, J. Ho, A. Lee, I. Awwal, H. Bradlow, J. Pan, S. Patil, K. Goldberg, and P. Abbeel, "Motion planning with sequential convex optimization and convex collision checking," in *The International Journal of Robotics Research*, vol. 33, no. 9, 2014, pp. 1251–1270, verified: DOI accessible - TrajOpt trajectory optimization method. [Online]. Available: https://doi.org/10.1177/0278364914528132

[15] M. Kalakrishnan, S. Chitta, E. Theodorou, P. Pastor, and S. Schaal, "Stomp: Stochastic trajectory optimization for motion planning," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2011, pp. 4569–4574, verified: DOI accessible - Stochastic trajectory optimization using covariance matrix adaptation. [Online]. Available: https://doi.org/10.1109/ICRA.2011.5980280

[16] A. H. Qureshi, M. J. Bency, and M. C. Yip, "Motion planning networks," *IEEE Transactions on Robotics*, vol. 35, no. 6, pp. 1277–1294, 2019, verified: DOI accessible - MPNet, learning-based motion planning. [Online]. Available: https://doi.org/10.1109/TRO.2019.2933224

[17] B. Ichter, E. Schmerling, A.-a. Agha-mohammadi, and M. Pavone, "Learned critical probabilistic roadmaps for robotic motion planning," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 9535–9541, verified: DOI accessible - Learning critical regions for sampling-based planning. [Online]. Available: https://doi.org/10.1109/ICRA40945.2020.9197290

[18] A. Majumdar and R. Tedrake, "Funnel libraries for real-time robust feedback motion planning," vol. 36, no. 8, 2017, pp. 947–982, verified: DOI accessible - Motion funnels for robust feedback control. [Online]. Available: https://doi.org/10.1177/0278364917712421

[19] R. Tedrake, I. R. Manchester, M. Tobenkin, and J. W. Roberts, "Lqr-trees: Feedback motion planning via sums-of-squares verification," vol. 29, no. 8, 2010, pp. 1038–1052, verified: DOI accessible - Verifiable motion funnels using sums-of-squares. [Online]. Available: https://doi.org/10.1177/0278364910369189

[20] B. Ichter, J. Harrison, and M. Pavone, "Learned sampling distributions for efficient planning in continuous spaces," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017, pp. 2318–2325, verified: DOI accessible - Learning sampling distributions for motion planning. [Online]. Available: https://doi.org/10.1109/IROS.2017.8206041

[21] M. Danielczuk, A. Mousavian, C. Eppner, and D. Fox, "Object rearrangement using learned implicit collision functions," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2021, pp. 6010–6017, verified: DOI accessible - Neural implicit representations for collision checking. [Online]. Available: https://doi.org/10.1109/ICRA48506.2021.9561398

[22] M. Diehl, H. G. Bock, J. P. Schlöder, R. Findeisen, Z. Nagy, and F. Allgöwer, "Real-time optimization and nonlinear model predictive control of processes governed by differential-algebraic equations," in *Journal of Process Control*. Elsevier, 2002, vol. 12, no. 4, pp. 577–585, standard reference - Sequential quadratic programming for MPC.

[23] B. Amos, I. D. J. Rodriguez, J. Sacks, B. Boots, and J. Z. Kolter, "Differentiable mpc for end-to-end planning and control," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2018, pp.

8299–8310, verified: arXiv accessible - Differentiable model predictive control. [Online]. Available: https://arxiv.org/abs/1810.13400

[24] M. Lutter, C. Ritter, and J. Peters, "Deep lagrangian networks: Using physics as model prior for deep learning," in *International Conference on Learning Representations (ICLR)*, 2019, verified: arXiv accessible - Physics-informed neural networks for control. [Online]. Available: https://arxiv.org/abs/1907.04490

[25] A. W. Winkler, C. D. Bellicoso, M. Hutter, and J. Buchli, "Gait and trajectory optimization for legged systems through phase-based end-effector parameterization," vol. 3, no. 3, 2018, pp. 1560–1567, verified: DOI accessible - TOWR for whole-body trajectory optimization. [Online]. Available: https://doi.org/10.1109/LRA.2018.2798285

[26] C. Mastalli, R. Budhiraja, W. Merkt, G. Saurel, B. Hammoud, M. Naveau, J. Carpentier, L. Righetti, S. Vijayakumar, and N. Mansard, "Crocoddyl: An efficient and versatile framework for multi-contact optimal control," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 2536–2542, verified: DOI accessible - Crocoddyl framework for contact-rich optimal control. [Online]. Available: https://doi.org/10.1109/ICRA40945.2020.9196673

[27] K. Lowrey, A. Rajeswaran, S. Kakade, E. Todorov, and I. Mordatch, "Plan online, learn offline: Efficient learning and exploration via model-based control," *arXiv preprint arXiv:1811.01848*, 2018, verified: arXiv accessible - Combining MPC with learned value functions. [Online]. Available: https://arxiv.org/abs/1811.01848

[28] J. Sacks, N. Atanasov, B. Boots, and Y. Zhu, "Ilqr-vae: Control-based learning of input-driven dynamics with applications to neural data," in *International Conference on Learning Representations (ICLR)*, 2022, verified: OpenReview accessible - Warm-starting MPC with learned models. [Online]. Available: https://openreview.net/forum?id=7b8J7cKL2U

[29] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning (ICML)*, 2021, pp. 8748–8763, verified: arXiv accessible - CLIP vision-language model from OpenAI. [Online]. Available: https://arxiv.org/abs/2103.00020

[30] M. Oquab, T. Darcet, T. Moutakanni, H. V. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby, R. Howes, P.-Y. Huang, H. Xu, V. Sharma, S.-W. Li, W. Galuba, M. Rabbat, M. Assran, N. Ballas, G. Synnaeve, I. Misra, H. Jegou, J. Mairal, P. Labatut, A. Joulin, and P. Bojanowski, "Dinov2: Learning robust visual features without supervision," *arXiv preprint arXiv:2304.07193*, 2023, verified: arXiv accessible - Self-supervised vision transformer from Meta AI. [Online]. Available: https://arxiv.org/abs/2304.07193

[31] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, P. Dollár, and R. Girshick, "Segment anything," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023, pp. 4015–4026, verified: arXiv accessible - SAM: foundation model for image segmentation from Meta AI. [Online]. Available: https://arxiv.org/abs/2304.02643

[32] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn, P. Florence, C. Fu, M. G. Arenas, K. Gopalakrishnan, K. Han, K. Hausman, A. Herzog, J. Hsu, B. Ichter, A. Irpan, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, I. Leal, L. Lee, T.-W. E. Lee, S. Levine, Y. Lu, H. Michalewski, I. Mordatch, K. Pertsch, K. Rao, K. Reymann, M. Ryoo, G. Salazar, P. Sanketi, P. Sermanet, J. Singh, A. Singh, R. Soricut, H. Tran, V. Vanhoucke, Q. Vuong, A. Wahid, S. Welker, P. Wohlhart, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich, "Rt-2: Vision-language-action models transfer web knowledge to robotic control," *arXiv preprint arXiv:2307.15818*, 2023, verified: arXiv accessible - Vision-language-action model from Google DeepMind. [Online]. Available: https://arxiv.org/abs/2307.15818

[33] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu, W. Huang, Y. Chebotar, P. Sermanet, D. Duckworth, S. Levine, V. Vanhoucke, K. Hausman, M. Toussaint, K. Greff, A. Zeng, I. Mordatch, and P. Florence, "Palm-e: An embodied multimodal language model," *arXiv preprint arXiv:2303.03378*, 2023, verified: arXiv accessible - Multimodal language model for embodied AI from Google. [Online]. Available: https://arxiv.org/abs/2303.03378

[34] A. Simeonov, Y. Du, A. Tagliasacchi, J. B. Tenenbaum, A. Rodriguez, P. Agrawal, and V. Sitzmann, "Neural descriptor fields: Se(3)-equivariant object representations for manipulation," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2022, pp. 6775–6782, verified: DOI accessible - SE(3)-equivariant neural fields for robotic manipulation. [Online]. Available: https://doi.org/10.1109/ICRA46639.2022.9812146

[35] L. Yen-Chen, P. Florence, J. T. Barron, A. Rodriguez, P. Isola, and T.-Y. Lin, "inerf: Inverting neural radiance fields for pose estimation," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021, pp. 1323–1330, verified: DOI accessible - Using NeRF for robotic pose estimation. [Online]. Available: https://doi.org/10.1109/IROS51168.2021.9636708

[36] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2012, pp. 5026–5033, verified: DOI accessible - MuJoCo physics engine, widely used in robotics. [Online]. Available: https://doi.org/10.1109/IROS.2012.6386109

[37] Y. Hu, L. Anderson, T.-M. Li, Q. Sun, N. Carr, J. Ragan-Kelley, and F. Durand, "Difftaichi: Differentiable programming for physical simulation," *arXiv preprint arXiv:1910.00935*, 2019, verified: arXiv accessible - Differentiable physics simulation framework. [Online]. Available: https://arxiv.org/abs/1910.00935

[38] C. D. Freeman, E. Frey, A. Raichuk, S. Girgin, I. Mordatch, and O. Bachem, "Brax – a differentiable physics engine for large scale rigid body simulation," *arXiv preprint arXiv:2106.13281*, 2021, verified: arXiv accessible - JAX-based differentiable physics from Google. [Online]. Available: https://arxiv.org/abs/2106.13281

[39] A. Sanchez-Gonzalez, J. Godwin, T. Pfaff, R. Ying, J. Leskovec, and P. Battaglia, "Learning to simulate complex physics with graph networks," in *International Conference on Machine Learning (ICML)*, 2020, pp. 8459–8468, verified: arXiv accessible - Graph Network Simulators from DeepMind. [Online]. Available: https://arxiv.org/abs/2002.09405

[40] M. Turpin, N. Michael, and V. Kumar, "Goal assignment and trajectory planning for large teams of interchangeable robots," vol. 37, no. 4, 2014, pp. 401–415, verified: DOI accessible - Centralized multi-robot trajectory planning. [Online]. Available: https://doi.org/10.1007/s10514-014-9412-1

[41] M. Čáp, P. Novák, A. Kleiner, and M. Selecký, "Prioritized planning algorithms for trajectory coordination of multiple mobile robots," *IEEE Transactions on Automation Science and Engineering*, vol. 12, no. 3, 2015, pp. 835–849, verified: DOI accessible - Priority-based multi-robot planning. [Online]. Available: https://doi.org/10.1109/TASE.2015.2445780

[42] J. van den Berg, M. Lin, and D. Manocha, "Reciprocal n-body collision avoidance," in *Robotics Research*. Springer, 2011, pp. 3–19, verified: DOI accessible - Reciprocal velocity obstacles for distributed collision avoidance. [Online]. Available: https://doi.org/10.1007/978-3-642-19457-3_1

[43] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 6379–6390, 2017, verified: arXiv accessible - MADDPG for multi-agent reinforcement learning. [Online]. Available: https://arxiv.org/abs/1706.02275

[44] Q. Li, F. Gama, A. Ribeiro, and A. Prorok, "Learning decentralized controllers for robot swarms with graph neural networks," pp. 671–682, 2020, verified: arXiv accessible - GNN-based learning for multi-robot coordination. [Online]. Available: https://arxiv.org/abs/1903.10527

[45] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. J. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, H. Michalewski, I. Mordatch, K. Pertsch, K. Rao, K. Reymann, M. Ryoo, G. Salazar, P. Sanketi, P. Sermanet, J. Singh, A. Singh, R. Soricut, H. Tran, V. Vanhoucke, Q. Vuong, A. Wahid, S. Welker, P. Wohlhart, J. Wu, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich, "Rt-1: Robotics transformer for real-world control at scale," *arXiv preprint arXiv:2212.06817*, 2022, verified: arXiv accessible - Large-scale robot learning with transformers from Google. [Online]. Available: https://arxiv.org/abs/2212.06817

[46] O. Nachum, S. Gu, H. Lee, and S. Levine, "Data-efficient hierarchical reinforcement learning," *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 3303–3313, 2018, verified: arXiv accessible

- HIRO for hierarchical reinforcement learning. [Online]. Available: https://arxiv.org/abs/1805.08296

[47] M. C. Machado, M. G. Bellemare, and M. Bowling, "A laplacian framework for option discovery in reinforcement learning," in *International Conference on Machine Learning (ICML)*, 2017, pp. 2295–2304, verified: arXiv accessible - Unsupervised option discovery using graph Laplacian. [Online]. Available: https://arxiv.org/abs/1703.00956

[48] S. Krishnan, A. Garg, S. Patil, C. Lea, G. Hager, P. Abbeel, and K. Goldberg, "Transition state clustering: Unsupervised surgical trajectory segmentation for robot learning," in *International Symposium on Robotics Research (ISRR)*. Springer, 2017, pp. 91–106, verified: DOI accessible - Trajectory segmentation for skill discovery. [Online]. Available: https://doi.org/10.1007/978-3-319-60916-4_6

[49] A. Paraschos, C. Daniel, J. Peters, and G. Neumann, "Probabilistic movement primitives," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2013, pp. 2616–2624, verified: URL accessible - ProMP for task-parameterized skill learning. [Online]. Available: https://papers.nips.cc/paper/5177-probabilistic-movement-primitives