

# A Controlled Sim2Sim Study of Online Material Belief for Granular Excavation

Taehwan Yun and HyunJun Jo

**Abstract**—Robotic excavation in granular media depends strongly on material properties that are not directly visible from the initial surface. A controller that optimizes only a target height-map may therefore select motions that match the geometry while violating force limits or failing to move material effectively. This paper studies whether a short raw-RGB material-probe sequence can be converted into an online material belief that improves subsequent excavation decisions. We infer a posterior over density, friction angle, tool friction, and cohesion, and use its mean to condition a finite-budget MPM digging selector. The visual input is raw RGB frames rather than hand-engineered height, flow, or pile descriptors. Across 24 MPM conditions spanning four granular materials, two trench targets, and three initial-bed seeds, the estimated-posterior controller reduces target loss to 2.010, compared with 2.198 for no posterior and 2.629 for an overconfident wrong posterior; the GT-property reference reaches 1.972. Paired tests show that the estimated posterior improves target loss over no posterior ( $p = 0.0025$ ) and wrong posterior ( $p = 0.031$ ), while remaining statistically close to the GT-property reference. An expanded 64-condition audit preserves the target-loss ordering, with 2.046 for estimated posterior versus 2.302/3.234 for no/wrong posterior. Additional checks are treated only as supporting evidence. These results support a bounded Sim2Sim claim: raw RGB material belief can be useful as a controller-facing variable for finite-budget excavation selection, while real-camera closed-loop excavation remains future work.

## I. INTRODUCTION

Granular manipulation is difficult because the visible state of the material does not fully specify how it will react to contact. Two sand beds with similar height maps may differ in compaction, friction, cohesion, and tool interaction. A digging motion that is effective in one bed can stall, overload the robot, or spill material in another. This coupling between geometry and latent material properties is central to robotic excavation, soil sampling, agricultural manipulation, and construction-scale earth moving.

Classical earthmoving theory and soil-tool models explain why resistive force depends on blade geometry, soil strength, and failure state [1], [2]. Robotic excavation methods have then used feedback, impedance, learning, or adaptive policies to cope with uncertain terrain [3], [4]. For smaller-scale granular manipulation, learned dynamics models can predict scooping and dumping outcomes [5], while recent work has made substantial progress on granular parameter inference and differentiable digging. Matl et al. infer granular material properties by matching

observed formations to simulations [6]. Hynninen et al. show that real force/torque traces can identify 11 granular materials [7]. DDBot demonstrates differentiable system identification and high-precision digging optimization for unknown granular materials [8]. These systems motivate a practical question for control: after the robot touches the material, does the resulting belief actually change the next excavation action in a useful way?

We address this question in a controlled Sim2Sim benchmark. The robot first observes raw RGB frames from a short material-probe interaction. An estimator updates a posterior over material properties, and a downstream controller uses this posterior to choose a finite-budget excavation action. The same initial bed, target, and action budget are used for all controller variants, isolating the contribution of the material belief.

The central contribution is a controlled Sim2Sim study of an auditable raw-pixel simulator-library posterior tested through its effect on a finite material-conditioned excavation selector. This is not a full real-robot excavation system; instead, the paper makes a narrower claim that can be audited from the current evidence:

- 1) We formulate posterior-conditioned excavation as a downstream test of material belief, separating property prediction from control utility.
- 2) We introduce a multi-condition MPM ablation comparing no posterior, wrong posterior, estimated posterior, and GT-property control across four materials, two targets, and three initial-bed seeds.
- 3) We connect the main Sim2Sim result to four supporting checks: a RAW-RGB posterior sanity check, a real-camera soil-image bridge, a vision-ambiguous raw-RGB/raw-force modality ablation, and a public real force/torque material-identification dataset. We also add robustness audits for target/seed diversity, exact material anchors, held-out continuous materials, shuffled posteriors, strict trench feasibility, and ranked finite-candidate control.

## II. RELATED WORK

**Soil-tool interaction and excavation control.** Classical earthmoving models estimate cutting and resistive forces from soil-tool geometry and soil failure assumptions [1]. Luengo et al. model and identify soil-tool interaction forces for automated excavation [2]. Later robotic excavation systems use reaction forces to adapt motion across

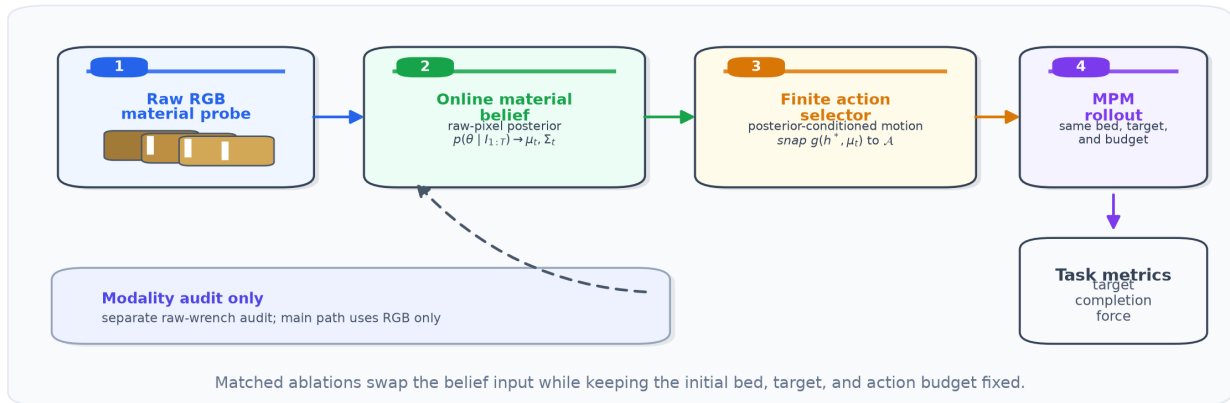


Fig. 1. Overview. The main closed-loop ablation uses raw RGB frames from a short material-probe interaction to update an online material posterior, then uses its mean to choose a finite-budget excavation motion. The ablation replaces the posterior input with no-posterior, wrong-posterior, estimated-posterior, and GT-property beliefs while keeping the initial bed, target, and action budget fixed; a separate modality audit adds raw wrench traces.

successive passes [3] or train policies that adapt to soil variation while respecting machine limits [4]. These methods motivate force-aware control, but they generally do not maintain an explicit posterior over granular properties for downstream target digging.

**Granular dynamics and parameter inference.** BayesSim casts simulator calibration as probabilistic inference over parameters [9]. Matl et al. use simulated granular formations to infer material properties for robotic tasks [6]. Our work follows the probabilistic view, but uses the inferred belief as a controller input and evaluates it by the resulting excavation behavior.

**Data-driven granular manipulation.** Schenck et al. learn predictive models for robotic scooping and dumping of granular media and show that explicitly modeling physical mechanics improves shape-control performance [5]. Material-adaptive graph dynamics [10] and visuo-haptic granular property estimation [11] further suggest that learned models can capture material-dependent deformation. Our method differs in that the learned material belief is evaluated through a posterior-conditioned excavation controller.

**Force-based material identification.** Force and torque measurements are informative for granular media because resistance emerges directly during contact. Hynninen et al. provide a public dataset and show that interactive force measurements can classify granular materials [7]. Related haptic systems recognize granular media through force or dip interactions [12], [13]. We use the public dataset as a real-data sanity check for the sensing premise, while the closed-loop excavation experiments remain simulator-based.

**Differentiable granular digging.** DDBot combines an MPM-based differentiable simulator, system identification, and gradient-based digging skill optimization [8].

That work is closest to ours in task domain. The comparison in this paper is scoped carefully: we implement a DDBot-core target-height planner for a shared benchmark, but do not claim to reproduce the full official DDBot runtime.

**MPM and differentiable simulation.** The material point method was introduced for history-dependent materials that are difficult for purely Lagrangian or Eulerian discretizations [14]. MLS-MPM provides efficient material point simulation with rigid body coupling [15], and DiffTaichi demonstrates differentiable programming for physical simulation [16]. These simulation tools make it possible to evaluate fine-grained granular interaction, but our main question is not simulator fidelity alone. We ask whether a material posterior changes the control decision under a fixed simulator and task budget.

### III. PROBLEM FORMULATION

We consider a robot executing a digging task in a granular bed. Let  $h_0$  be the initial height field and  $h^*$  the target trench or target height field. The robot has a finite budget of candidate tool motions  $a \in \mathcal{A}$ , where each candidate specifies insertion depth, push length, duration, and tool orientation. The environment is an MPM simulator with latent material parameters

$$\theta = [\rho, \phi, \mu_{\text{tool}}, c], \quad (1)$$

where  $\rho$  is density,  $\phi$  is friction angle,  $\mu_{\text{tool}}$  is tool-interface friction, and  $c$  is cohesion.

During probing or early interaction, the robot observes raw visual evidence

$$e_t = f(I_{0:t}^{\text{rgb}}), \quad (2)$$

where  $I^{\text{rgb}}$  denotes rendered camera frames. In the main closed-loop experiments, the visual branch receives the

RGB tensor directly; no pile-height, optical-flow, surface-slope, or mask-area summary is supplied. Raw wrench traces are added only in the modality ablation. The estimator maintains a Gaussian material belief

$$p_t(\theta) = \mathcal{N}(\mu_t, \Sigma_t). \quad (3)$$

In implementation, this belief is a nonparametric raw-pixel posterior. We render a fixed library of 260 short RGB probe videos with sampled material parameters, compare the observation to each library video by mean squared pixel distance over frames, height, width, and color channels, and apply a fixed softmax temperature to obtain library weights. The reported mean and covariance are the weighted first and second moments of the material parameters. Thus the main estimator is intentionally simple and auditable: it is not trained on hand-engineered visual descriptors, and its exact-prototype dependence is checked separately by the anchor-free and held-out continuous-material audits.

The main controller uses the posterior mean  $\mu_t$  as the material input for a finite-budget action selector. The posterior covariance is reported for calibration and coverage audits, but the main MPM control ablation does not claim covariance-aware optimal control. The selector maps the material belief to a desired insertion depth, reach, and push duration, then snaps that desired motion to the nearest member of  $\mathcal{A}$  under a fixed normalized action distance,  $a^* = \arg \min_{a \in \mathcal{A}} d(a, g(h^*, \mu_t))$ . Since  $\mathcal{A}$  is finite, the GT-property variant is a reference for this controller class, not a global optimum. The mapping  $g$  is fixed before evaluation. It first computes a normalized material strength  $s = 0.18\rho_n + 0.30\phi_n + 0.23\mu_{\text{tool},n} + 0.29c_n$  and a force-limited index  $q = 0.65c_n + 0.25\mu_{\text{tool},n} + 0.10\phi_n$ . Mid-strength material is assigned a deeper and longer desired pass, while high-cohesion or high-interface-friction material is assigned a shallower, slower pass. This keeps the controller auditable: the posterior changes the material input to a fixed finite selector rather than training or tuning a new planner per trial.

#### IV. METHOD

Figure 1 summarizes the system. The main estimator receives raw RGB frames and updates a posterior over material properties. The controller uses the posterior mean to select a material-conditioned desired motion from the same finite action budget. The reported posterior uncertainty is used as a sensing-quality diagnostic rather than as a separate risk term in the main controller.

The main ablation uses four belief inputs:

- **No posterior:** a nominal prior, without contact-derived material update.
- **Wrong posterior:** a deliberately incorrect belief, used to test whether extra belief inputs can hurt control.
- **Estimated posterior:** the online material estimate.
- **GT property:** true simulator material parameters, passed through the same finite candidate controller.

The wrong-posterior condition is essential. If the estimated posterior only acts as an arbitrary conditioning vector, then a wrong posterior should not produce a systematic failure. If the controller truly depends on material semantics, then a wrong posterior can lead to an unsafe or ineffective action.

#### V. EXPERIMENTAL DESIGN

The experiments are organized around three questions.

**Q1: Does the posterior improve excavation control?** We compare the four belief conditions over 24 MPM trials: four material settings, two trench targets, and three initial-bed seeds. Within each trial, all variants use the same initial bed, target corridor, and candidate action budget. The primary metric is a target loss combining transport error, depth error, target-mass error, force-limit violation, and lateral spillage. We also report trench completion and force violation separately. The wrong-posterior condition is an overconfident hard-dry belief that intentionally selects a deeper, longer pass than is appropriate for force-limited materials. For trial  $i$ , the reported target loss is

$$L_i = 0.90E_{\text{transport}} + 0.85E_{\text{depth}} + 0.35E_{\text{mass}} + 4.0 \times 10^{-4}V_{\text{force}} + 3.5 \times 10^{-5}M_{\text{spill}}, \quad (4)$$

where  $E_{\text{transport}}$  is normalized transport error,  $E_{\text{depth}} = |d/d^* - 1|$ ,  $E_{\text{mass}}$  is normalized target-zone mass error,  $V_{\text{force}} = \max(0, F_{\text{max}} - F_{\text{lim}})$ , and  $M_{\text{spill}}$  is lateral spillage mass in simulator units. Lower values are better. These weights are fixed before the ablation and are not tuned per material or per controller; the main table also reports completion and force violation separately so that the scalar loss does not hide the tradeoff. Completion is the achieved-depth ratio and is not capped at one; values above one therefore indicate over-digging rather than extra success.

**Q2: Is the estimator robust beyond one nominal material?** We evaluate whether the posterior can be inferred from raw RGB frames rather than hand-engineered vision summaries. A held-out procedural RGB benchmark reports normalized property MAE and two-sigma posterior coverage. This is a sanity check for the visual input contract, not a real-camera result. To partially bridge this gap, we also evaluate a raw-pixel model on the public real-camera soil-image and particle-size-distribution dataset of [17], [18]. This bridge uses resized RGB tensors directly and splits by soil sample id to avoid image-level leakage. We further stress-test camera appearance shift by applying global gain, bias, gamma, and noise perturbations directly to the raw RGB probe videos. To distinguish static appearance from interaction evidence, we also compare first-frame, full-sequence, and late-frame posterior inference.

We also run a modality ablation under deliberately ambiguous visual evidence. The visual branch still receives raw RGB frame tensors, but the frames are color-

TABLE I

MAIN RAW-RGB POSTERIOR CONTROL ABLATION OVER 24 MPM TRIALS. STRICT IS THE COUNT SATISFYING DEPTH COMPLETION IN  $[0.9, 1.1]$  WITH NO FORCE VIOLATION; GT ACT. IS EXACT AGREEMENT WITH THE GT-PROPERTY FINITE ACTION.

Belief	Loss ↓	Comp.	Viol. N ↓	Strict	GT act.
No post.	$2.198 \pm 0.587$	$1.19 \pm 0.22$	$446 \pm 928$	2/24	0/24
Wrong post.	$2.629 \pm 1.706$	$0.59 \pm 0.51$	$2852 \pm 4476$	0/24	0/24
Est. post.	$2.010 \pm 0.508$	$1.03 \pm 0.41$	$420 \pm 882$	5/24	12/24
GT prop.	$1.972 \pm 0.552$	$0.98 \pm 0.43$	$338 \pm 864$	2/24	24/24

TABLE II

PAIRED TARGET-LOSS TESTS OVER THE SAME 24 MPM CONDITIONS. POSITIVE  $\Delta L$  MEANS THE ESTIMATED POSTERIOR HAS LOWER LOSS THAN THE COMPARISON. CI IS A BOOTSTRAP 95% CONFIDENCE INTERVAL FOR THE PAIRED MEAN DELTA.

Comparison	$\Delta L$	95% CI	$p_t$	$p_W$	Wins
No post. – Est.	0.188	[0.082, 0.295]	0.0025	0.0024	17/24
Wrong post. – Est.	0.620	[0.145, 1.176]	0.0310	0.0170	17/24
GT prop. – Est.	-0.038	[-0.133, 0.052]	0.442	0.528	13/24

normalized so material appearance cues are weak. We then compare RGB-only, raw-wrench-only, and RGB+wrench posterior inputs on the same 24 MPM material-target-seed conditions.

**Q3: What happens in a force-dominant target benchmark?** We compare against a DDBot-core target-height planner in a shared benchmark with three hidden material-strength cases and five seeds. The baseline follows the DDBot core abstraction of optimizing a 5D digging skill for final target height, using either GT strength or a vision-nominal strength. Our controller uses the same target and task simulator, but first performs a force probe, forms a strength posterior, and replans for five closed-loop strokes. Safe target reach means final height-map error below 8.0 and peak force below 3300 N. This comparison is not an official DDBot reproduction.

Finally, we test whether the sensing assumption is plausible on real data by evaluating a lightweight classifier on the public 11-material force/torque dataset of [7].

## VI. RESULTS

### A. Posterior-Conditioned Excavation

Figure 2 and Table I show the main result. The estimated-posterior controller reduces mean target loss from 2.198 to 2.010 relative to no posterior and from 2.629 to 2.010 relative to the wrong posterior. Table II reports paired tests over the same 24 matched conditions:  $p_t = 0.0025$ ,  $p_W = 0.0024$ , and bootstrap CI [0.082, 0.295] against no posterior; and  $p_t = 0.031$ ,  $p_W = 0.017$ , and CI [0.145, 1.176] against wrong posterior. The estimated posterior is slightly worse than the GT-property reference, 2.010 versus 1.972, and this difference is not statistically significant ( $p_t = 0.44$ ), with bootstrap CI [-0.133, 0.052] crossing zero. Because the GT row uses the same finite action selector and the same rollout budget, it is an

TABLE III

MATERIAL AND TARGET BREAKDOWN.  $\Delta L$  COLUMNS ARE PAIRED TARGET-LOSS DELTAS; POSITIVE VALUES FAVOR THE ESTIMATED POSTERIOR.

Group	No $\Delta$	Wrg. $\Delta$	No W	Wrg. W
Loose dry	0.133	-0.227	3/6	3/6
Dense sand	0.167	0.010	6/6	3/6
Cohesive wet	0.138	2.253	3/6	6/6
Angular grit	0.316	0.443	5/6	5/6
Shallow target	0.392	0.209	12/12	8/12
Deep target	-0.016	1.030	5/12	9/12

oracle material-input reference rather than a global optimum; occasional estimated-posterior wins are therefore interpreted as rollout and metric effects, not as evidence that the estimator is better than GT material parameters. This is the desired ordering for the claim: the learned belief improves control over missing or incorrect belief and approaches the oracle material input.

The strict-success column in Table I is deliberately a literal trench-band audit, not an oracle-optimality metric. A trial is strict only if depth completion falls inside  $[0.9, 1.1]$  and the force budget is not violated; policies that over-dig, under-dig, or trade depth for lower target loss can therefore fail this binary column even when their mean target loss is better. For this reason the GT-property row should be read jointly with target loss and exact GT-action match: it verifies the material-conditioned selector’s reference action, whereas strict success tests one narrow operational trench specification.

Table III makes the failure modes explicit. Estimated posterior improves over no posterior on all four material means and is strongest on angular grit, while the wrong hard-dry posterior fails sharply on cohesive wet sand and angular grit. The loose-dry case is an exception against the wrong posterior because the aggressive hard-dry pass is not force-unsafe there. The deep target is also harder to separate from no posterior because nominal over-digging can accidentally reach part of the deeper trench. These breakdowns are why we claim average paired improvement under matched conditions, not universal dominance in every material-target pair. The largest no-posterior counterexample occurs for loose dry sand with the deep target and bed seed 7: the estimated posterior snaps to the same finite action as the GT-property reference and scores 1.346, whereas the no-posterior robust nominal pass scores 1.073. This case illustrates that the GT-property row is a material-input reference for the implemented selector, not a global oracle for every scalar-loss instance.

A controller-regret audit gives the same interpretation at the action level. The estimated posterior selects the exact same discrete action as the GT-property controller in 12/24 conditions, compared with 0/24 for both no-posterior and wrong-posterior controllers. Its mean normalized action distance to the GT policy is 0.462, compared with 1.336 for no posterior and 2.088 for wrong

## Posterior-conditioned excavation: evidence stack

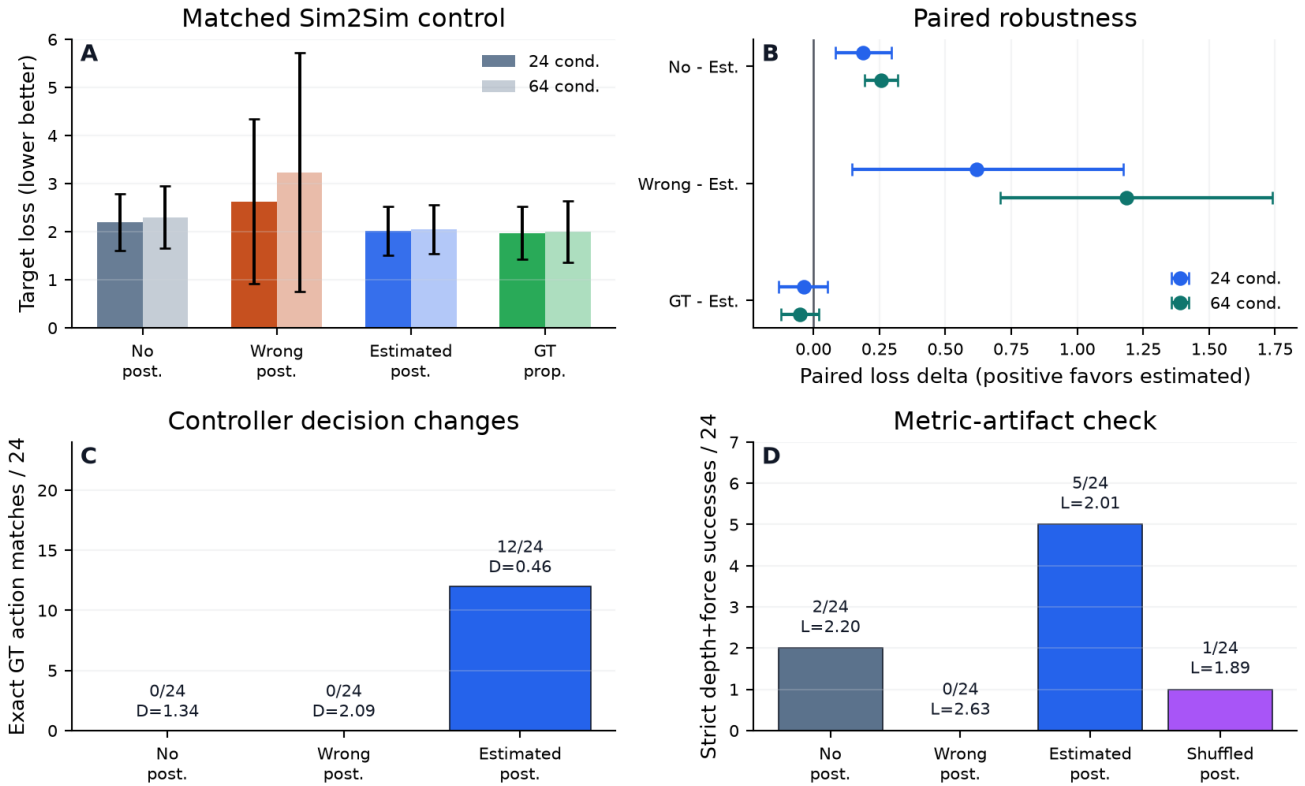


Fig. 2. Evidence stack for the main control claim. (A) Target-loss means over the 24-condition main benchmark and the 64-condition extension. (B) Paired bootstrap intervals preserve the estimated-posterior advantage over no/wrong posterior and overlap the GT-property reference. (C) The estimated posterior changes the selected discrete action toward the GT-property action. (D) A strict depth-and-force audit shows that the shuffled-posterior scalar-loss improvement comes from conservative under-digging rather than literal trench success.

posterior. The corresponding absolute target-loss regret is 0.147, compared with 0.282 and 0.942. Thus the posterior improves not only the reported task metric, but also the controller’s selected action relative to the GT-property reference.

Loss-weight sensitivity gives the same qualitative conclusion against no posterior under geometry-heavy, reported, and force-safety-heavy weights; the wrong-posterior comparison is weaker when geometry dominates ( $p = 0.166$ ) and stronger when force safety is weighted. A 64-condition audit over four targets, four bed seeds, and four materials preserves the ordering: 2.046 versus 2.302/3.234 for no/wrong posterior (54/64 and 50/64 paired wins), close to the GT-property controller at 1.996 with CI [-0.124, 0.021].

Two negative-control audits bound the interpretation. First, a shuffled-posterior audit can lower scalar target loss to 1.889 by under-digging, but it worsens action distance to the GT-property action from 0.462 to 0.904 and reduces exact GT-action matches from 12/24 to 0/24. Second, a full ranked-candidate audit that scores all finite candidates does not support the stronger claim that the current

posterior is a globally reliable candidate optimizer: the estimated-versus-no target-loss CI is [-0.136, 0.213]. It still moves the selected action toward the GT-property action, with 15/24 exact GT matches. We therefore describe the implemented controller as a posterior-conditioned finite selector, not as a globally optimal excavation optimizer.

The held-out continuous-material audit further tests whether the result depends on the four named benchmark prototypes. We introduce four held-out materials and evaluate the same two targets and three bed seeds, giving 24 additional matched MPM conditions. The top posterior match is never one of the named benchmark materials (0/24 conditions), and the normalized property MAE is 0.058. The estimated posterior reduces target loss to 2.189 versus 2.375 for no posterior and 2.913 for wrong posterior, with paired CIs [0.082, 0.300] and [0.178, 1.371]. It remains close to the GT-property finite selector at 2.095, with CI [-0.225, 0.018] crossing zero. At the action level, estimated posterior matches the GT-property action in 12/24 conditions, compared with 6/24 for no posterior and 0/24 for wrong posterior.

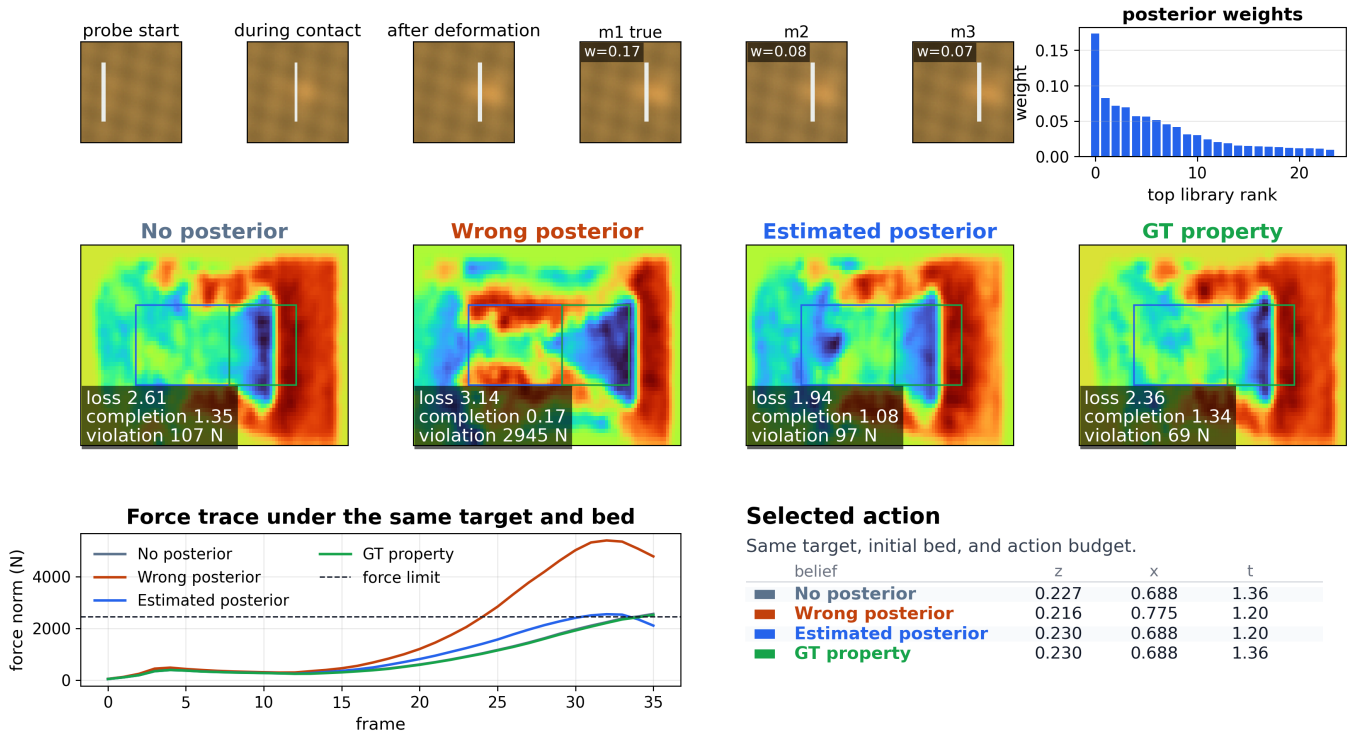


Fig. 3. Qualitative matched rollout for cohesive wet sand under the same shallow-trench target, initial bed seed, and action budget. The top row shows the raw RGB probe sequence, top posterior library matches, and posterior weights. The middle row shows final height maps for no posterior, wrong posterior, estimated posterior, and GT property control; blue and green boxes mark cut and deposit regions. The bottom row shows force traces and selected action parameters. The wrong posterior selects an unsafe material assumption and violates the force limit, while the estimated posterior changes the action toward the GT-property rollout. This condition illustrates the paired statistics in Table I rather than replacing the aggregate results.

TABLE IV  
RAW-RGB POSTERIOR SANITY CHECK OVER 160 HELD-OUT SAMPLES.

Property	MAE	Mean $\sigma$
$\rho$	15.9	44.1
$\phi$	3.49	4.95
$\delta$	2.07	4.04
$c$	0.78	1.89

### B. RAW-RGB Posterior Sanity Check

The raw-RGB posterior sanity check gives a mean normalized MAE of 0.077 and a mean two-sigma coverage of 0.981 over 160 held-out samples, verifying the raw pixel input path rather than real-camera video generalization. An exact-material anchor audit removes the four named benchmark prototypes from the posterior library: property MAE rises from 0.044 to 0.079, yet closed-loop target loss is 2.015, essentially matching the anchored posterior at 2.010 and still improving over no/wrong posterior with CIs [0.084, 0.281] and [0.133, 1.170]. This weakens an exact lookup explanation within the same procedural rendering family.

A 256-sample calibration audit gives conservative 1/2/3-sigma coverages of 0.846/0.993/1.000, so sigma is reported as an uncertainty diagnostic rather than as a covariance-

optimal control input.

The appearance-shift audit shows that direct raw-pixel matching is sensitive to global camera shifts: normalized MAE rises from 0.077 to 0.241 under gain, bias, gamma, and noise, while affine-normalized raw RGB reduces the shifted MAE to 0.103. This does not solve real-camera transfer, but it reduces the fixed render-palette shortcut risk.

The temporal audit separates static appearance from temporal evidence. The first rendered frame is already informative (0.082 versus 0.077 for the full sequence), but under color-normalized RGB the full sequence and late frames reduce MAE from 0.196 to 0.172 and 0.169. Thus deformation helps most when appearance cues are weak.

### C. Real-Camera RAW-RGB Bridge

To test whether real camera pixels carry related granular information, we use the public close-range soil-image dataset of [17], [18]. It provides smartphone RGB images paired with particle-size distributions. We resize each image to a  $64 \times 64 \times 3$  raw RGB tensor and train a ridge model directly on pixels; no color moments, texture descriptors, or other visual summaries are used. Evaluation is grouped by soil sample id, so images of the same soil sample never appear in both train and test. On 135 images from 25 image-matched soil samples, the raw-RGB model reduces

TABLE V

DDBOT-CORE STRESS-TEST BENCHMARK OVER 15 TRIALS. VALUES ARE MEAN  $\pm$  STD. SPILLAGE IS REPORTED IN  $10^{-4}$  M<sup>3</sup>; REACH IS THE FRACTION OF TRIALS THAT REACH THE TARGET-ERROR THRESHOLD WITHOUT A FORCE-LIMIT VIOLATION.

Method	HM err.	Viol. N	Spill $10^{-4}$ m <sup>3</sup>	Score	Reach
DDBot-core GT	10.91 $\pm$ 0.54	923 $\pm$ 944	1.10 $\pm$ 0.18	13.68 $\pm$ 3.35	0.00
DDBot-core nom.	10.91 $\pm$ 0.54	920 $\pm$ 941	1.10 $\pm$ 0.18	13.67 $\pm$ 3.34	0.00
Ours force-post.	6.76 $\pm$ 0.81	23 $\pm$ 31	2.57 $\pm$ 0.28	6.82 $\pm$ 0.85	0.60

sample-level log-D50 MAE from 0.896 to 0.787 relative to a train-mean baseline. It reduces the  $> 2$  mm coarse-fraction MAE from 29.5% to 23.3% and improves gravel-dominant classification accuracy from 0.56 to 0.72. The result is only a bridge experiment, not real robotic excavation, but it reduces the concern that the visual branch works only on procedural renders.

#### D. Vision-Ambiguous Modality Ablation

The modality audit isolates the sensing role of force when vision is intentionally ambiguous. Raw RGB alone gives a normalized posterior MAE of 0.195. Adding raw wrench traces reduces the MAE to 0.084; force-only evidence gives 0.095. In control, RGB+force reduces mean force violation from 436 N to 368 N relative to RGB only. Its target-loss improvement over RGB-only is positive but modest:  $\Delta L = 0.060$  over 24 paired trials, with  $p = 0.198$  by paired t-test and  $p = 0.057$  by one-sided Wilcoxon test. Force-only has the lowest target loss among estimated variants, 2.018 versus 2.108 for RGB-only and 2.047 for RGB+force, while RGB+force has the lowest force violation. The GT-property reference reaches 1.978. We use this ablation as evidence that contact force is informative under visual ambiguity, not as proof that multimodal fusion strictly dominates every single-modality controller.

#### E. Supporting DDBot-Core Force-Dominant Stress Test

The force-posterior controller reduces final height-map error from 10.91 to 6.76 and force violation from 920 N to 23 N relative to the target-only nominal DDBot-core planner. The tradeoff is higher spillage, because the closed-loop controller moves more material while staying under the force limit. All rows share the same target-height generator, hidden material-strength cases, seeds, 5D digging-skill abstraction, and simulator metric; the comparison differs in whether force-derived posterior feedback is allowed during replanning. This result should be read as a controlled benchmark of the failure mode where target geometry alone is under-informative and force feedback is needed. It is not a claim of superiority over the full DDBot system or its released runtime.

#### F. Real Force/Torque Sanity Check

On the public dataset, Raw+HFMH reaches 89.5% accuracy across 11 materials. Using only the first 0.8 s of interaction still gives 86.1% accuracy. This does not prove

TABLE VI

REAL FORCE/TORQUE CLASSIFICATION. HFMH IS A FORCE/TORQUE FEATURE SUMMARY FROM THE PUBLIC DATASET BASELINE, NOT A VISUAL SUMMARY.

Input	Acc.	Std.
Raw sequence	0.736	0.039
HFMH summary	0.883	0.020
Raw + HFMH	0.895	0.028
Raw + HFMH, first 0.8 s	0.861	0.020

TABLE VII

EVIDENCE BOUNDARY FOR THE CURRENT SUBMISSION.

Item	Evidence in paper	Not claimed
Control	24/64 matched MPM ablations	Real-robot excavation
Vision	Raw procedural RGB video tensors	Real-camera video transfer
Real pixels	Static soil RGB with PSD labels	Closed-loop robot control
Force	Force/torque data and modality audit	Hidden wrench in main RGB result
DDBot	Scoped DDBot-core stress test	Full official DDBot superiority

real-world excavation, but it supports the sensing premise that early contact force can carry material information before the full motion is complete.

## VII. DISCUSSION AND LIMITATIONS

The experiments support a controller-facing interpretation of material posterior. The posterior is not only a label or parameter estimate; it is a decision variable that changes which excavation action the controller selects. The main ablation is therefore stronger than a standalone property-prediction metric: it tests whether the estimate matters for a downstream manipulation task.

The improvement is not uniform over all materials. The estimated posterior has lower target loss than the no-posterior and wrong-posterior variants in 17 of 24 paired trials, but an incorrect hard-dry posterior can occasionally select a pass that is adequate for loose dry sand. Its failure is clearer in force-limited cohesive or angular materials, where the same overconfident belief produces large force violations. We therefore claim average paired improvement under this benchmark, not universal dominance for every granular condition.

The current evidence remains bounded: closed-loop excavation is simulated; closed-loop RGB is procedural; the real-camera bridge is static soil imagery; the DDBot comparison is DDBot-core; and the GT-property row is finite-budget, not globally optimal. A shuffled-posterior audit can also lower scalar loss via conservative under-digging, so we pair target loss with action-regret, breakdowns, and strict depth-and-force feasibility: estimated succeeds in 5/24 rollouts versus 2/24, 0/24, and 1/24 for no, wrong, and shuffled posterior. The material-to-action mapping is also heuristic: loss-weight and ranked-candidate audits test two important axes, but a complete deployment study should perturb the selector mapping itself and include real calibration shift. A real-camera excavation test must therefore collect synchronized raw RGB, force/torque, robot state, and pre/post height observations while replaying the

same no/wrong/estimated/GT controller ablation under matched initial beds and targets.

### VIII. CONCLUSION

This paper asks whether online material belief improves granular excavation control. Across 24 MPM material-target-seed conditions, an estimated posterior reduces target loss relative to no-posterior and wrong-posterior controllers and remains close to the GT-property reference. A RAW-RGB sanity check verifies that the visual posterior can be formed from pixel inputs rather than hand-engineered visual summaries. A scoped DDBot-core stress test, real-camera soil images, and public force/torque traces provide limited support outside procedural rendering. These results motivate real-to-sim posterior estimation as the next stage toward material-aware robotic excavation.

### REFERENCES

- [1] A. R. Reece, "Paper 2: The fundamental equation of earth-moving mechanics," *Proceedings of the Institution of Mechanical Engineers, Conference Proceedings*, vol. 179, no. 6, 1964.
- [2] O. Luengo, S. Singh, and H. Cannon, "Modeling and identification of soil-tool interaction in automated excavation," in *Proceedings of the 1998 IEEE/RSJ International Conference on Intelligent Robots and Systems*, vol. 3, 1998, pp. 1900–1906.
- [3] G. J. Maeda, D. C. Rye, and S. P. N. Singh, "Iterative autonomous excavation," in *Field and Service Robotics*, ser. Springer Tracts in Advanced Robotics. Springer, 2014, vol. 92, pp. 369–382.
- [4] P. Egli, D. Gaschen, S. Kerscher, D. Jud, and M. Hutter, "Soil-adaptive excavation using reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 9778–9785, 2022.
- [5] C. Schenck, J. Tompson, S. Levine, and D. Fox, "Learning robotic manipulation of granular media," in *Proceedings of the 1st Annual Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, vol. 78. PMLR, 2017, pp. 239–248.
- [6] C. Matl, Y. Narang, R. Bajcsy, F. Ramos, and D. Fox, "Inferring the material properties of granular media for robotic tasks," in *IEEE International Conference on Robotics and Automation*, 2020.
- [7] S. Hynninen, T. N. Le, and V. Kyrki, "Interactive identification of granular materials using force measurements," *arXiv preprint arXiv:2403.17606*, 2024.
- [8] X. Yang, M. Wei, Y.-K. Lai, and Z. Ji, "Ddbot: Differentiable physics-based digging robot for unknown granular materials," *arXiv preprint arXiv:2510.17335*, 2025.
- [9] F. Ramos, R. C. Possas, and D. Fox, "Bayessim: Adaptive domain randomization via probabilistic inference for robotics simulators," in *Robotics: Science and Systems*, 2019.
- [10] K. Zhang, B. Li, K. Hauser, and Y. Li, "Adaptigraph: Material-adaptive graph-based neural dynamics for robotic manipulation," *arXiv preprint arXiv:2407.07889*, 2024.
- [11] Z. Zhang, G. Zheng, X. Ji, G. Chen, R. Jia, W. Chen, G. Chen, L. Zhang, and J. Pan, "Understanding particles from video: Property estimation of granular materials via visuo-haptic learning," *IEEE Robotics and Automation Letters*, 2025.
- [12] Z. Zhang, G. Chen, W. Chen, R. Jia, L. Zhang, and J. Pan, "A joint learning of force feedback of robotic manipulation and textual cues for granular materials classification," *IEEE Robotics and Automation Letters*, vol. 10, pp. 7166–7173, 2025.
- [13] X. Wang, S. Zhang, Z. Zhao, L. Zhu, and A. Song, "Dipme: Haptic recognition of granular media for tangible interactive applications," *arXiv preprint arXiv:2411.08641*, 2024.
- [14] D. Sulsky, Z. Chen, and H. L. Schreyer, "A particle method for history-dependent materials," *Computer Methods in Applied Mechanics and Engineering*, vol. 118, no. 1–2, pp. 179–196, 1994.
- [15] Y. Hu, Y. Fang, Z. Ge, Z. Qu, Y. Zhu, A. Pradhana, and C. Jiang, "A moving least squares material point method with displacement discontinuity and two-way rigid body coupling," *ACM Transactions on Graphics*, vol. 37, no. 4, 2018.
- [16] Y. Hu, L. Anderson, T.-M. Li, Q. Sun, N. Carr, J. Ragan-Kelley, and F. Durand, "DiffTaichi: Differentiable programming for physical simulation," in *International Conference on Learning Representations*, 2020.
- [17] E. Soranzo and A. D'Souza, "Close range images of soils and their particle size distributions," 2025. [Online]. Available: <https://zenodo.org/records/14725633>
- [18] E. Soranzo, "Dataset of close-range soil images and corresponding particle size distributions," *Data in Brief*, vol. 60, p. 111631, 2025.