

Efficient Structure-Guided 3D Physical Property Reasoning

Anonymous OpenSUN3D submission

Paper ID

Abstract

001 *Inferring an object’s physical properties such as material*
 002 *type and surface hardness from visual observations is es-*
 003 *sential for augmented reality, robotic perception, and em-*
 004 *odied intelligence. However, existing solutions to phys-*
 005 *ical property reasoning like NeRF2Physics are computa-*
 006 *tionally expensive and error-prone because they interpo-*
 007 *late sparse, noisy CLIP features across dense 3D scenes.*
 008 *This creates a fundamental conflict between the pursuit*
 009 *of high semantic resolution and high reasoning efficiency*
 010 *while making the system sensitive to oblique or low-quality*
 011 *viewpoints. We introduce a lightweight, structure-guided*
 012 *framework that achieves fine-grained semantic consistency*
 013 *for physical property reasoning with orders-of-magnitude*
 014 *lower computational cost. Our key insight is that the 3D*
 015 *structural priors offer a stronger cue for the object’s se-*
 016 *mantic organization, which allows us to avoid the dense in-*
 017 *terpolation for physical property reasoning. We project 2D*
 018 *DINO embeddings into 3D for coarse component segmen-*
 019 *tation, perform adaptive sparse sampling of representative*
 020 *CLIP source points, and apply a view-quality-aware patch*
 021 *selection with probability-weighted aggregation. These de-*
 022 *signs successfully eliminate dense interpolation, suppress*
 023 *noisy viewpoints, and drastically cut the number of CLIP*
 024 *queries. Extensive experiments on ABO dataset demon-*
 025 *strate our method reduces end-to-end runtime from hun-*
 026 *dreds of seconds to mere seconds per scene while improv-*
 027 *ing semantic accuracy, spatial coherence, and downstream*
 028 *physical-property inference.*

029 1. Introduction

030 Reasoning an object’s physical properties from its visual
 031 observations, such as material type, weights, surface rough-
 032 ness, and structural attributes, is a fundamental capability
 033 for visual intelligence and embodied perception [6, 24].
 034 Recent advances that ground 2D Vision-Language Mod-
 035 els (VLMs) in 3D representations have begun to unlock
 036 open-world physical reasoning by connecting an object’s
 037 appearance to its underlying material semantics [13, 19, 23,

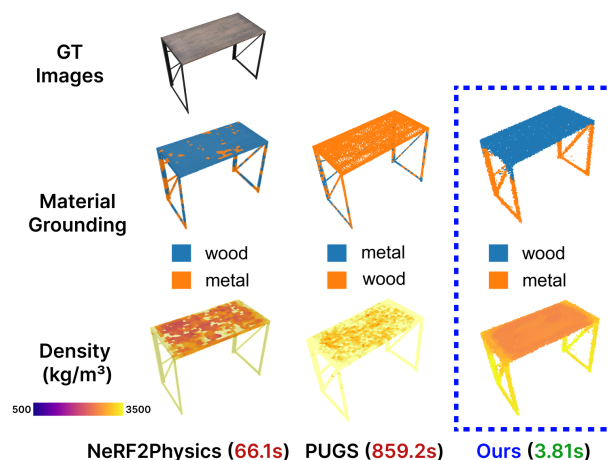


Figure 1. **Material grounding and physical property estimation comparison.** Compared to NeRF2Physics [29] (66.1s) and PUGS [23] (859.2s), our method produces more accurate material grounding and physical property estimation in only 3.81s, achieving over 17× speedup over NeRF2Physics and 225× over PUGS. Our method yields coherent, component-level material assignments (e.g., correctly distinguishing the wooden tabletop from the metal frame) and more uniform density distributions within homogeneous regions. In contrast, prior methods exhibit fragmented material boundaries and inconsistent physical property estimation.

27, 30]. This progress opens new avenues for augmented
 reality interaction [1], more accurate physics-based simu-
 lation [10] and improved robotic manipulation and plan-
 ning [5, 21].

While the detailed techniques may differ, existing vision-
 based physical property reasoning solutions [19, 23, 30] all
 follow a common pipeline: begin with a reconstructed 3D
 point cloud and, for each point, retrieve the corresponding
 multi-view image patches. Each patch is fed into the vision
 encoder [8, 20] to obtain sparse semantic features, which
 are then interpolated across the full 3D point cloud to pre-
 dict a material label for every point. Finally, these per-point
 material estimates are integrated to infer physical properties
 such as hardness or density.

However, such approaches suffer from two limitations.

- Excessively long processing time from exploding CLIP-

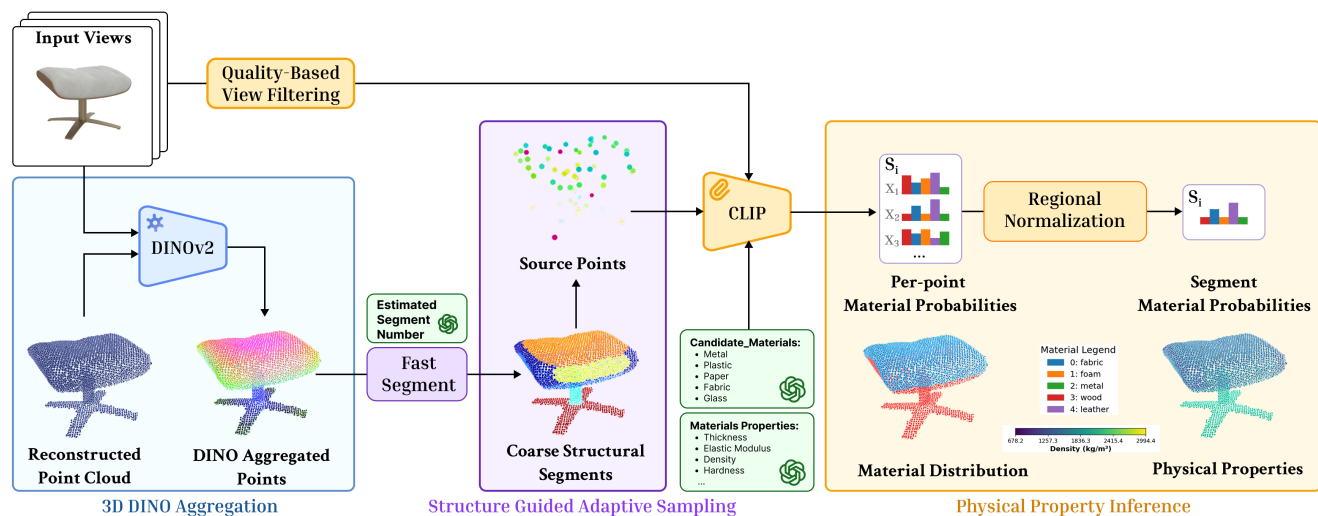


Figure 2. **Overview of S3-PHYS.** Given posed multi-view images, S3-PHYS lifts dense DINO features into 3D to form a structure-aligned feature field, segments the object into coarse components, and samples only a few representative points for CLIP encoding. After quality-based view filtering, we fuse CLIP features to estimate per-component material probabilities and predict downstream physical properties. Our design removes dense interpolation and enables fast, consistent 3D physical-property reasoning.

054 *query complexity.* Achieving fine-grained material resolution
 055 requires querying CLIP for a large number of source
 056 points, and each point must be processed from multiple
 057 visible views. Consequently, the number of CLIP calls
 058 scales multiplicatively with the number of points and views,
 059 quickly dominating runtime and making high-resolution inference
 060 computationally prohibitive.

061 • *Noisy vision features leading to unstable physical reason-*
 062 *ing.* CLIP embeddings captured far away or from oblique
 063 viewpoints are frequently noisy (see Appendix), yet exist-
 064 ing methods treat all visible views uniformly, allowing low-
 065 quality patches to pass through without filtering. These
 066 noisy features propagate through the dense interpolation
 067 stage, degrading material prediction accuracy and produc-
 068 ing unstable estimates of physical properties. They further
 069 increase latency, as thousands of unreliable patches must
 070 still be processed to recover fine-grained detail.

071 In this paper, we propose S3-PHYS, an efficient yet more
 072 accurate framework for physical property reasoning that re-
 073 tains the semantic richness of CLIP-based methods while
 074 addressing the two core bottlenecks discussed above, as
 075 shown in Fig. 2. Our approach introduces three key ideas
 076 that together strike an effective balance between semantic
 077 fidelity and computational efficiency, as elaborated below.

078 **(1) Structure-guided fast source-point selection.** Instead
 079 of querying CLIP for every point-view pair, we use DINO
 080 features [4, 16] as efficient structural priors. DINO pro-
 081 duces dense pixel-level features in a single forward pass
 082 per view, which we lift into 3D to obtain a dense feature
 083 field aligned with geometric and textural boundaries. This
 084 field provides a reliable signal for segmenting the scene into
 085 coarse, semantically coherent components. From each com-

ponent, we then sample only a small set of representative
 086 points as CLIP source points. This significantly reduces
 087 the number of CLIP queries and improves computational
 088 efficiency while preserving the semantic detail needed for
 089 accurate material and physical property reasoning. 090

(2) Regional material normalization and property inference.
 091 Within each structural region, instead of relying
 092 on dense and potentially noisy interpolation, we normalize
 093 the per-point material probabilities of sparse source points
 094 to obtain a single, spatially consistent material probability
 095 distribution for the entire region. This regional normal-
 096 ization suppresses point-level semantic noise and enforces
 097 component-level coherence. The resulting distribution is
 098 then used to estimate the region-level physical property via
 099 kernel regression following [30] (Eq. 6), yielding smoother
 100 and more reliable material and physical property estimates. 101

**(3) View-quality-aware patch filtering for robust CLIP
 102 embeddings.** To further mitigate CLIP noise, we intro-
 103 duce a view-filtering module that evaluates each candidate
 104 view via geometric and photometric heuristics such as sur-
 105 face normal alignment, view angle, and illumination. Only
 106 high-quality patches are retained for CLIP encoding, which
 107 reduces noisy embeddings and decreases the number of
 108 CLIP calls by an order of magnitude. When combined with
 109 DINO-guided segmentation and adaptive sampling, this en-
 110 ables fast, stable, and high-resolution physical reasoning,
 111 lowering end-to-end runtime from hundreds of seconds to
 112 only a few seconds while improving both semantic coher-
 113 ence and physical accuracy. 114

We summarize the contributions of our work as follows: 115

- We provide, to our knowledge, the first analysis show-
 116 ing why existing CLIP-based 3D physical reasoning
 117 pipelines are slow and unstable, revealing two core bot-
 118

119 tlenecks: (i) excessive CLIP queries caused by dense
 120 sampling, and (ii) noisy multi-view features that under-
 121 mine material and physics inference.

- 122 • We introduce a structure-guided sampling strategy that
 123 projects 2D DINO features into 3D to obtain structural
 124 priors. This enables coarse component segmentation and
 125 reduces CLIP source-point queries by orders of magni-
 126 tude while retaining semantic detail.
- 127 • We develop an efficient semantic fusion pipeline combin-
 128 ing view-quality-aware patch filtering with probability-
 129 weighted component voting. This suppresses noisy view-
 130 points, stabilizes material estimation, and removes the
 131 need for dense point-wise interpolation.
- 132 • Extensive experiments on the ABO dataset show that our
 133 method improves semantic accuracy and spatial coher-
 134 ence, achieving **17.4 - 225.4× speedup** over prior CLIP-
 135 based baselines and enabling practical real-time 3D phys-
 136 ical property reasoning.

137 2. Related Work

138 We divide related works into two categories.

139 **Vision-based object’s physical property reasoning.** Infer-
 140 ring physical properties from visual observations has long
 141 been a core challenge in computer vision and cognitive sci-
 142 ence [5, 9, 26]. Prior work has explored dynamic reason-
 143 ing by observing object motion [14, 26] or physical inter-
 144 actions in simulators [18, 28], and static reasoning from
 145 single images [2, 3, 22, 24]. While these methods can es-
 146 timate quantities such as mass, friction, or elasticity, they
 147 rely on task-specific supervision or constrained experimen-
 148 tal setups. Recent neural approaches extend visual physics
 149 reasoning into 3D. *NeRF2Physics* [30] embeds visual cues
 150 within NeRF fields but requires time-consuming volumet-
 151 ric optimization for each scene. 3D Gaussian-based phys-
 152 ical reasoning algorithms [23, 27] accelerate this process via
 153 3DGS reconstruction, yet they still suffer from instability
 154 under noisy multi-view inputs and high computational cost
 155 from redundant CLIP queries. Our method addresses these
 156 issues by introducing a DINO-guided 3D structural prior
 157 for robust component segmentation and adaptive sampling.
 158 This design not only suppresses noisy-view interference but
 159 also reduces CLIP invocations by orders of magnitude, en-
 160 abling fast, consistent, and high-resolution physical reason-
 161 ing within seconds rather than minutes.

162 **3D language fields.** A complementary line of work builds
 163 semantic 3D representations using vision–language mod-
 164 els (VLMs). CLIP [20], trained on large-scale image–text
 165 pairs, has proven effective for zero-shot segmentation, lo-
 166 calization, and open-world reasoning. Methods such as Dis-
 167 tilled Feature Fields, LERF [13], and OpenNeRF [7] inject
 168 CLIP features into NeRFs for semantic 3D understanding,
 169 while approaches like Openscene [17, 31] map CLIP fea-

tures onto point clouds, and feature splatting [19] prop-
 agates them into 3DGS representations. However, these
 pipelines typically rely on dense feature interpolation over
 millions of 3D points and exhaustive multi-view fusion,
 which greatly increases computation and exacerbates CLIP
 noise. In contrast, our method leverages CLIP more judi-
 ciously: we extract features only from sparse, high-quality
 patches selected via DINO-guided structural segmentation
 and view-quality heuristics, producing faster and more reli-
 able 3D semantic fields.

3. Methods

In this section, we first present the system overview (§3.1).
 We then detail each design component, including structure-
 guided adaptive sampling (§3.2), efficient semantic fusion
 (§3.3), material proposal (§3.4), and regional-aware phys-
 ical property inference (§3.5).

3.1. Overview

S3-PHYS takes as input a set of posed multi-view im-
 ages and a reconstructed 3D space (e.g., point cloud from
 VGGT [25], NeRF [30], 3DGS [12, 23], or mesh re-
 construction) and produces a spatially consistent physical-
 property map for any query point X . As illustrated in Fig. 2,
 the framework consists of three core stages:

Step One: Structure-guided adaptive sampling. We ex-
 tract DINO features [16] from each input view and project
 them onto the 3D point cloud to obtain a dense per-point
 feature field that encodes both geometry and appearance.
 This lifted feature field serves as an efficient structural prior:
 it can be computed in a single multi-view pass and adds neg-
 ligible overhead in practice (§5.4). Using these features, we
 segment the object into coarse, semantically coherent com-
 ponents and sample only a small number of representative
 points from each component as CLIP source points, subst-
 antially reducing the number of required CLIP queries.

Step Two: Efficient semantic fusion. For each represen-
 tative point, we evaluate candidate views using geometric
 and photometric criteria and retain only the high-quality
 ones. We perform multi-scale CLIP feature fusion on these
 patches to obtain stable and reliable semantic embeddings.

Step Three: Material proposal and physical reasoning.
 In parallel, a vision-language model (VLM) analyzes a set
 of diverse input views to propose candidate materials and
 their physical properties. We fuse these proposals with
 the 3D semantic embeddings through *probability-weighted
 component voting*, yielding component-level material as-
 signments and object-level physical property distributions.

These stages yield fast, noise-resistant, and semanti-
 cally coherent 3D physical-property maps without relying
 on dense interpolation, reducing the end-to-end running la-
 tency from hundreds of seconds to 3.81s (§5).

220 3.2. Structure-Guided Adaptive Sampling

221 We consider the task of inferring per-point material semantics and downstream physical properties from a recon-
 222 structed 3D object. The input includes a point cloud $\{x_i\}$,
 223 and a set of posed multi-view images $\{I_v\}_{v \in V}$. For each
 224 3D point x_i , the goal is to infer a material probability distri-
 225 bution and convert it into physical quantities. Most exist-
 226 ing pipelines, such as NeRF2Physics, all rely on query-
 227 ing numerous source points from the object’s point cloud to
 228 obtain sufficient semantic resolution. This design however
 229 creates a fundamental trade-off: sparse sampling triggers
 230 less CLIP encoding but produces inconsistent semantics,
 231 whereas dense sampling dramatically improves coverage at
 232 the cost of substantial runtime and increased noise.
 233

234 To overcome this challenge, we first use DINO features
 235 as a geometric-semantic prior for adaptive sampling. For
 236 each input view $v \in V$, DINO produces dense pixel-level
 237 features $F_v(\cdot)$ in a single forward pass. We then lift these
 238 2D features into 3D by projecting every point x_i onto each
 239 view through the camera projection function $\pi(x_i, v) \in \mathbb{R}^2$.
 240 Only views in which the point is actually visible contribute
 241 to its aggregated feature.

242 **Aggregated DINO feature field.** Let $M_{\text{vis}}(x_i, v)$ denote
 243 the visibility of point x_i from view v . The aggregated DINO
 244 feature for point x_i is defined as:

$$245 \Psi_{\text{DINO}}(x_i) = \frac{\sum_{v \in V} M_{\text{vis}}(x_i, v) \cdot F_v(\pi(x_i, v))}{\sum_{v_j \in V} M_{\text{vis}}(x_i, v_j)}. \quad (1)$$

246 **Visibility function.** A point is considered visible in a view
 247 if its projected depth is consistent with the depth map:

$$248 M_{\text{vis}}(x_i, v_j) = \begin{cases} 1, & \text{if } \text{dis}(x_i, v) \leq \text{Depth}(\pi(x_i, v)) \\ 0, & \text{otherwise.} \end{cases} \quad (2)$$

249 where $\text{dis}(x_i, v)$ is the depth of x_i along the camera ray of
 250 view v , and $\text{Depth}_v(\cdot)$ is the depth map rendered from that
 251 view.

252 **Coarse structural segmentation.** With the aggregated
 253 DINO field $\Psi_{\text{DINO}}(x_i)$, our next step is to partition the
 254 object into coarse, semantically coherent components that
 255 serve as regions for adaptive sampling. DINO features natu-
 256 rally align with geometric and textural boundaries, allow-
 257 ing structurally similar regions to be clustered together at
 258 negligible cost. However, the raw DINO embeddings are
 259 very high-dimensional, and K-means clustering is known to
 260 degrade in both accuracy and efficiency in such spaces. To
 261 address this, we apply PCA to obtain a compact represen-
 262 tation of each aggregated feature (8D in our case), greatly
 263 reducing computational cost while preserving the relevant
 264 structure. Then we concatenate each PCA-reduced DINO
 265 feature with its 3D coordinate x_i , enabling the clustering to
 266 jointly consider appearance similarity and spatial proximity.
 267 Formally, the segment results are represented as follows:

$$268 \mathcal{S} = \text{KMeans}([\mathcal{f}_{\text{PCA}}(\Psi_{\text{DINO}}(x_i)), x_i], N_S), \quad (3)$$

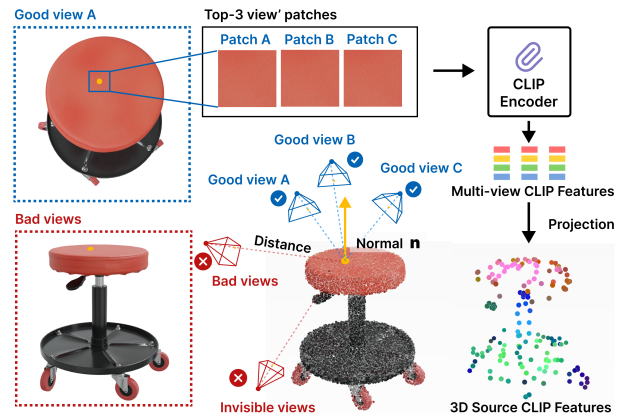


Figure 3. **Quality-based View Filtering** S3-PHYS excludes those bad views that contribute less to semantic fusion, thereby reducing the runtime latency.

269 where the number of clusters N_S is provided by the VLM
 270 during the material-proposal stage (§3.4). This effectively
 271 groups geometrically and visually coherent regions, estab-
 272 lishing a structural prior for sampling. To ensure coverage
 273 of all structures regardless of size, we sample a fixed num-
 274 ber of $k=10$ points from each segment S_i . This strategy
 275 drastically reduces computational overhead without sacri-
 276 ficing coverage. As detailed in Table 1, our approach re-
 277 duces the average CLIP source points from 3,505 to 82
 278 ($42\times$ reduction) while maintaining semantic diversity.

279 3.3. Efficient Semantic Fusion

280 Our semantic fusion module contains two components:
 281 quality-based view filtering and multi-scale feature fusion.

282 3.3.1. Quality-based View Filtering

283 After selecting representative 3D source points, we fuse
 284 CLIP features from multiple views to obtain stable seman-
 285 tic embeddings. However, multi-view images vary widely
 286 in quality: some views provide sharp, front-facing obser-
 287 vations, while others suffer from blur, occlusion, extreme
 288 angles, or specular reflection, as shown in Fig. 3. Directly
 289 fusing CLIP features from all visible views amplifies noise
 290 and increases computation, a key weakness of prior CLIP-
 291 based pipelines. To mitigate this issue, we assess the quality
 292 of each visible view before extracting and fusing its CLIP
 293 features using two geometric factors:

- 294 • **Distance-induced resolution loss.** Faraway views con-
 295 tain fewer texture details and produce weak CLIP fea-
 296 tures, so they should be excluded.
- 297 • **View-surface angle mismatch.** Highly oblique viewing
 298 angles introduce foreshortening and reflection artifacts,
 299 which degrade semantic consistency.

300 We compute an importance score that combines distance
 301 and normal alignment for every point-view pair. Each sam-
 302 pled point has an estimated surface normal. For a point at

303 depth z from a view, the distance score is defined as follows:

$$304 \quad S_{\text{dist}}(z) = \frac{1}{1+z},$$

305 The angle score is defined below:

$$306 \quad S_{\text{angle}}(\mathbf{v}, \mathbf{n}) = \max(0, \mathbf{v} \cdot (-\mathbf{n})),$$

307 where \mathbf{v} is the unit view direction and \mathbf{n} is the outward normal. The importance score of an view is computed as:

$$309 \quad S_{\text{total}} = S_{\text{dist}} \cdot S_{\text{angle}}.$$

310 We rank views by S_{total} and retain only the top- k (where
311 $k=3$) for CLIP feature extraction. According to our ex-
312 periments, this filtering reduces the number of embedded
313 patches by approximately $2.7\times$ while improving robustness
314 by discarding poorly aligned or low-quality observations, as
315 shown in Table 1.

316 3.3.2. Multi-Scale Feature Fusion

317 After selecting the top- k views, we extract CLIP features
318 from patches centered at the projected location of each sam-
319 pled point. Using a single patch size often provides lim-
320 ited context: small patches fail to capture larger structures,
321 while large patches smooth out fine texture. Following
322 LeRF [13], we therefore aggregate features across multiple
323 spatial scales to improve robustness. For each retained view,
324 patches of different scales are encoded by the CLIP image
325 encoder and averaged to form a multi-scale representation.
326 For source point x , features from the top- k views are then
327 averaged to obtain the final semantic embedding:

$$328 \quad \mathbf{f}(x) = \frac{1}{k} \sum_{v=1}^k \left(\frac{1}{|\mathcal{L}|} \sum_{l \in \mathcal{L}} \phi(\mathcal{P}_v^{(l)}(x)) \right), \quad (4)$$

329 where $\phi(\cdot)$ is the CLIP encoder, $\mathcal{P}_v^{(l)}(s)$ is a patch of scale
330 l extracted from view v , and \mathcal{L} is the set of scales. This
331 strategy preserves both fine-grained detail and coarse con-
332 text while keeping computation manageable. As shown in
333 Table 1, even with multi-scale processing, the total number
334 of embedded patches remains $40\times$ smaller than the base-
335 line.

336 **Parallelize multi-scale feature fusion.** To accelerate se-
337 mantic feature fusion, we project all sampled points onto the
338 selected views in one vectorized operation. Angle scores
339 of all selected views are then batch-evaluated using pre-
340 estimated normals, and all valid image patches are consol-
341 idated into a single global batch for CLIP encoding. This
342 transforms the standard nested loop-based computation into
343 a fully parallelized pass, substantially reducing latency.

344 With these optimizations, the entire semantic fusion
345 stage, including CLIP feature extraction and fusion, com-
346 pletes within 1.2 s (Tab. 6).

347 3.4. Material Proposal

348 To ensure coverage of all visible materials, we select 2-4
349 views I_m with the largest pose differences from the input set
350 and concatenate them into a single composite image. This
351 composite image, together with a task-specific text prompt,
352 is provided to the VLM to generate a detailed semantic de-
353 scription and a list of material candidates. For each mate-
354 rial, the model predicts its approximate thickness within the
355 object as well as corresponding physical property values.
356 Formally, this process produces a material–property dictio-
357 nary:

$$358 \quad \mathcal{M} = \{(k_i, y_i, \theta_i)\}_{i=1}^K,$$

359 where k_i is the material name, y_i is the associated set of
360 physical properties or value ranges, and θ_i is the estimated
361 thickness ratio. In addition, the VLM provides an estimated
362 segment number N_S , which is used as the target cluster
363 count in §3.2. The complete prompt and query design are
364 provided in the Appendix.

365 3.5. Regional-aware Physical Inference

366 Given the spatially aggregated semantic features (§3.3.2)
367 and the VLM-generated material proposals (§3.4), we next
368 infer a consistent material distribution and its correspond-
369 ing physical properties for each structural segment and sub-
370 sequently estimate its corresponding physical properties to
371 obtain the final object-level attributes. Because DINO fea-
372 tures offer much finer granularity than CLIP embeddings,
373 we assume that each DINO-derived segment S_i is domi-
374 nated by a single material class.

375 **Regional material normalization.** To get robust material
376 estimation, we first compute the material affinity score for
377 each point and then apply regional material normalization.
378 Specifically, for any point x and material candidate k , we
379 define: $\omega_x(k) = \phi_{\text{CLIP}}(\mathbf{f}(x), \mathbf{t}(k))$, where $\phi_{\text{CLIP}}(\cdot)$ denotes
380 cosine similarity, $\mathbf{f}(x)$ is the fused semantic feature at point
381 x , $\mathbf{t}(k)$ is the CLIP text embedding of the k -th material la-
382 bel. Thus, $\omega_x(k)$ measures how well point x semantically
383 matches material k in CLIP space.

384 Rather than relying on noisy point-wise predictions, we
385 aggregate affinities over each structural segment to improve
386 robustness. We compute the spatially averaged material
387 score $\bar{\omega}_{S_i}(k)$ for the segment S_i using the equation:

$$388 \quad \bar{\omega}_{S_i}(k) = \frac{1}{|S_i|} \sum_{x \in S_i} \omega_x(k). \quad (5)$$

389 We then assign the material with the maximum aggregated
390 score to segment S_i .

391 **Physical property inference.** Once each segment S_i has
392 been assigned a dominant material, we estimate its physical
393 properties using the material–property dictionary provided
394 by the VLM as shown in Fig. 2. Let $\{y_k\}_{k=1}^K$ denote the in-
395 trinsic physical values (e.g., density) associated with the K
396 material candidates. To make the prediction robust, rather

Table 1. Workload reduction from each module in semantic fusion. SAS: Structure-Guided Adaptive Sampling; QVF: Quality-based View Filtering; MFF: Multi-scale Feature Fusion. Results from 5 scenes. Our full pipeline reduces processed patches from 28,040 to 702 (40 \times).

	Baseline	SAS (§3.2)	SAS + QVF (§3.3.1)	SAS + QVF + MFF (§3.3.2)
Avg. embedded source points	3,505	78	78	78
Avg. patches per point	8	8	3	9
Total embedded patches	28,040	624	234	702

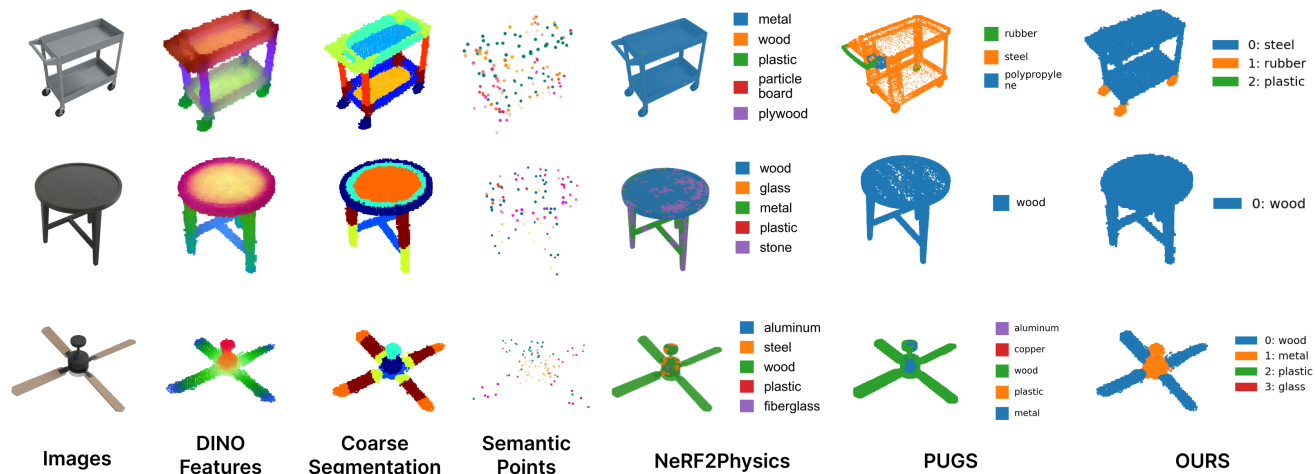


Figure 4. **Qualitative comparison on representative ABO objects.** For each object, we show the input RGB view, DINO features, coarse segmentations, semantic points and material grounding. Compared with NeRF2Physics and PUGS, ours system produces cleaner, more coherent component boundaries and more accurate material assignments, leading to more reliable physical-property predictions.

397 than taking a hard material assignment, we compute a soft,
 398 probability-weighted estimate by applying a temperature-
 399 controlled softmax over the aggregated material affinities
 400 $\bar{\omega}_{S_i}(k)$. Here T is a temperature parameter that controls
 401 confidence sharpness (*lower T makes the distribution more*
 402 *peaked*, while higher T encourages smoother mixing across
 403 materials). Formally, the predicted physical property for
 404 segment S_i is represented as:

$$405 \quad \rho(S_i) = \frac{\sum_{k=1}^K \exp(\bar{\omega}_{S_i}(k)/T) y_k}{\sum_{k=1}^K \exp(\bar{\omega}_{S_i}(k)/T)}. \quad (6)$$

406 This yields a single, stable physical-property estimate for
 407 the segment S_i , which is assigned uniformly to all points in
 408 S_i . To derive global object-level physical quantities such as
 409 total mass $\hat{\zeta}$, we accumulate material contributions from all
 410 segments as follows:

$$411 \quad \hat{\zeta} = \sum_i \rho(S_i) \theta_i b^2 \lambda, \quad (7)$$

412 where θ_i is the VLM-predicted thickness for segment S_i , b
 413 is the adaptive voxel size, and λ is the geometric correction
 414 factor from [30].

415 4. Implementation Details

416 We run all experiments on a workstation equipped with an
 417 NVIDIA A6000Ada GPU and an AMD Ryzen Threadrip-
 418 per PRO 5975WX CPU Processor, using PyTorch. We

use DINOv2-B/14 [16] to extract 2D features, reduced to 8
 dimensions by PCA before projection into 3D, and adopt
 OpenCLIP ViT-B/16 [11] as the image encoder for multi-
 scale patch embedding. The number of clusters N_S is esti-
 mated by the VLM during the material proposal stage, typi-
 cally ranging from 6–14 depending on object complexity.
 For each segment, we sample $k_{sample} = 10$ representative
 source points and retain the top-3 views based on the view-
 quality score. Feature fusion and CLIP encoding are fully
 vectorized and executed in parallel, allowing the semantic
 fusion stage to complete within 1.2 s. The material propo-
 sal stage, based on GPT-4o [15] with 2 composite views,
 runs asynchronously. The full pipeline, including feature
 extraction, sampling, fusion, and voting, processes a com-
 plete scene in approximately 3.81 s on average, achieving
 17.4 - 225.4 \times speedup over previous CLIP-based methods
 such as NeRF2Physics and PUGS.

5. Experiments

5.1. Metrics

We employ four complementary metrics to evaluate both
 physical accuracy and computational efficiency. (1) **Mass
 estimation error.** Following NeRF2Physics [30] and
 PUGS [23], we evaluate four quantitative measures between
 the predicted and measured mass on the ABO dataset (500
 objects): (i) Absolute Deviation Error (ADE), (ii) Absolute

Table 2. Accuracy results on ABO dataset (500 objects). Lower is better for ADE/ALDE/APE; higher is better for MnRE.

Method	ADE↓	ALDE↓	APE(%)↓	MnRE(%)↑
NeRF2Physics [30]	12.725	0.736	1.040	0.564
PUGS [23]	9.461	0.661	0.767	0.576
Ours	8.485	0.657	0.751	0.571

444 Log-Deviation Error (ALDE), (iii) Absolute Percentage Error (APE), and (iv) Mean Relative Error (MnRE). (2) **Material segmentation IoU**, assessing spatial consistency across
 445 100 ABO objects. (3) **Efficiency**, measured as the average
 446 scene processing time excluding point cloud reconstruction,
 447 for fair comparison with PUGS and NeRF2Physics.
 448

449 For the segmentation metric, we unify all materials
 450 into 10 categories and manually annotate material masks
 451 on selected reference views for both datasets. The mean
 452 Intersection-over-Union (mIoU) is then computed between
 453 our predicted material distributions and the annotated
 454 ground truth.
 455

456 5.2. Evaluation

457 We compare S3-PHYS against two recent representa-
 458 tive CLIP-based baselines—NeRF2Physics [30] and
 459 PUGS [23]—which also perform material-aware physical
 460 property reasoning from visual inputs. Table 2 and Table 3
 461 summarize the quantitative results for mass estimation and
 462 material segmentation.

463 Across all four mass-related metrics (ADE, ALDE, APE,
 464 and MnRE), our method consistently achieves competitive
 465 results, reducing ADE by 33% relative to NeRF2Physics
 466 and 10% relative to PUGS. The improvement mainly arises
 467 from our component-aware sampling and Regional aware
 468 physical inference, which enhance the semantic consistency
 469 of CLIP embeddings and eliminate interpolation artifacts
 470 commonly observed in previous methods.

471 Per-class results reveal complementary strengths: our
 472 method achieves clear gains on high-coverage, struc-
 473 turally coherent materials including **Plastic** (0.132 vs.
 474 0.059/0.033), **Fabric** (0.752 vs. 0.720/0.634), and **Wood**
 475 (0.736 vs. 0.640/0.638), and remains competitive on **Ce-
 476 ramic** (0.850 vs. 0.865 for PUGS), while NeRF2Physics
 477 leads on **Metal** (0.297 vs. 0.202 ours) and PUGS on **Rub-
 478 ber** (0.367 vs. 0.135 ours). **Glass** yields near-zero IoU
 479 across all methods due to sparse pixel presence, and the
 480 **Other** category shows large inter-method variance, with
 481 NeRF2Physics scoring 0.362 while both PUGS and ours
 482 yield zero, likely reflecting its heterogeneous composi-
 483 tion. Overall, the mIoU gains reflect stronger dataset-
 484 wide robustness on semantically coherent materials, with
 485 challenges on visually ambiguous or rare materials shared
 486 across methods.

487 In addition, our framework produces smoother and more
 488 spatially coherent material predictions. Qualitative com-



Figure 5. **Qualitative results in different physical properties.** Our method predicts four distinct physical properties—thermal conductivity, density, elastic modulus, and Shore hardness—across diverse object categories. The predicted distributions reflect material-level consistency: e.g., the wooden chair exhibits uniformly low thermal conductivity and moderate density, while the metal components of the mechanic stool show distinctly higher hardness and elastic modulus values compared to its rubber wheels.

489 comparisons in Fig. 4 show that NeRF2Physics and PUGS often exhibit fragmented or noisy material boundaries, while S3-PHYS yields uniform and physically plausible segmentation across diverse object categories. This demonstrates that leveraging DINO features as a structural prior not only improves quantitative performance but also enhances cross-component semantic coherence—an essential property for reliable physical reasoning in downstream applications. Fig. 5 exhibits S3-PHYS broad adaptability to multiple downstream physical properties.

499 5.3. Ablation Studies

500 We conduct ablations on a 100-object subset of ABO to
 501 evaluate the contribution of each module in our pipeline.
 502 Starting from a naive CLIP-only baseline with kNN-based
 503 physical inference [30], we progressively add Structure-
 504 guided Adaptive Sampling (SAS), Regional-aware Physical
 505 Inference (RPI), Quality-based View Filtering (QVF), and
 506 Multi-Scale Feature Fusion (MFF). For a fair comparison,
 507 the naive baseline uses voxel downsampling to obtain ap-
 508 proximately 80 source points, which matches the average
 509 number of points selected by SAS.

510 Table 4 summarizes the results. The naive baseline is the
 511 fastest, as it does not require DINO feature extraction or fu-
 512 sion. However, its accuracy is the lowest: with only a small
 513 number of source points, kNN-based inference struggles to
 514 preserve spatial granularity. Adding SAS alone does not re-
 515 solve this issue. Although the sampling becomes more struc-
 516 tured, the underlying kNN inference still collapses when the
 517 source set is sparse.

518 Introducing the RPI module leads to a substantial im-
 519 provement: MnRE increases by over 12% without any run-

Table 3. **Per-class IoU and scene average mIoU comparison across methods.** Higher is better. The *Stone* category does not appear in the 100-object annotated subset and is therefore marked as N/A.

Method	Plastic	Rubber	Fabric	Metal	Wood	Ceramic	Glass	Stone	Other	Overall mIoU
NeRF2Physics	0.059	0.222	0.720	0.297	0.640	0.475	0.015	–	0.362	0.304
PUGS	0.033	0.367	0.634	0.235	0.638	0.865	0.000	–	0.000	0.428
Ours	0.132	0.135	0.752	0.202	0.736	0.850	0.000	–	0.000	0.461

Table 4. **Ablation on 100-object subset of ABO.** We progressively add SAS (Structure-guided adaptive Sampling), RPI (Regional-aware Physical Inference), QVF (Quality-based View Filtering), and MFF (Multi-Scale Feature Fusion) modules. Each delta percentage for MnRE relates to naive pipeline, and each delta runtime relates to its previous row. Naive denotes the baseline using sparse voxel sampling without further optimization.

Ablation Settings	MnRE \uparrow	Time(s) \downarrow
Naive	0.550	1.73
+ SAS	0.551 (+0.18%)	4.01 (+2.28s)
+ SAS + RPI	0.595 (+8.18%)	3.97 (-0.04s)
+ SAS + RPI + QVF	0.598 (+8.73%)	3.26 (-0.71s)
+ SAS + RPI + QVF + MFF	0.602 (+9.45%)	3.44 (+0.18s)

time overhead. This demonstrates that region-level physical inference is significantly more robust than point-level voting when dealing with sparsely sampled observations.

Adding QVF further reduces runtime by 0.84s by discarding low-quality views before aggregation, while slightly improving MnRE. Finally, integrating MFF yields the best overall accuracy with only a modest increase in runtime, still lower than the unfiltered variants. Overall, each module contributes complementary benefits, jointly improving both accuracy and efficiency.

5.4. Efficiency Analysis

We further evaluate the runtime efficiency of our framework on the ABO dataset, measuring the average latency of each processing stage. As summarized in Table 5, our method achieves an average end-to-end processing time of only 3.81 s per scene (excluding reconstruction), corresponding to a 17.4 \times speedup over NeRF2Physics and a 225.5 \times speedup over PUGS. The improvement mainly stems from our structure-guided adaptive sampling and efficient semantic fusion, which substantially reduce the number of CLIP embeddings and patch queries required for material reasoning. The runtime scales linearly with the number of 3D points, demonstrating the method’s suitability for large-scale deployment and real-time perception scenarios.

A detailed stage-wise breakdown is presented in Table 6. Among all modules, DINO feature extraction and semantic fusion account for the majority of processing time, contributing 42.2% and 30.7% of total runtime, respectively. DINO feature aggregation and adaptive sampling are relatively lightweight, accounting for 22.6% and 3.9%, while

Table 5. **Efficiency results on ABO dataset (500 objects).** Lower is better for runtime; higher is better for speedup.

Method	Time(s) \downarrow	Speedup \uparrow
NeRF2Physics [30]	66.1	x1.0
PUGS [23]	859.2	x0.08
Ours	3.81	x17.4

Table 6. Stage-wise latency breakdown per scene on ABO (excluding reconstruction).

Stage	Time(s) \downarrow	Share(%)
DINO Feature Extraction	1.61	42.2
DINO Feature Aggregation	0.86	22.6
Structure Guided Adaptive Sampling	0.15	3.9
Efficient Semantic Fusion	1.17	30.7
Property inference	0.02	0.5
Total (Ours)	3.81	100

property inference is negligible (0.5%). This balanced runtime distribution indicates that the computational bottleneck has been largely mitigated, with no single component dominating the overall latency.

6. Conclusion

We have presented the design and evaluation of S3-PHYS, a lightweight, structure-guided framework for fast and reliable visual physical property reasoning. By combining DINO-derived structural priors with structure-guided adaptive sparse sampling, quality-based view filtering, and regional-aware physical inference, S3-PHYS achieves accurate and spatially coherent material understanding while reducing inference time from hundreds of seconds to only a few seconds per scene (achieving up to 225 \times speedup), which is an order-of-magnitude improvement over prior pipelines. Experiments on ABO dataset confirm strong gains in material grounding quality and physical property estimation. Current limitations include sensitivity to homogeneous surfaces in DINO-based segmentation and CLIP’s difficulty distinguishing visually similar materials; future work will explore physically grounded cues and end-to-end differentiable formulations to address these challenges. Looking ahead, our approach establishes a practical foundation for real-time material perception in AR/VR, robotics, and multimodal digital twins.

575

References

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

- [1] Syed Shah Alam, Samiha Susmit, Chieh-Yu Lin, Mohammad Masukujjaman, and Yi-Hui Ho. Factors affecting augmented reality adoption in the retail industry. *Journal of Open Innovation: Technology, Market, and Complexity*, 7(2):142, 2021. 1
- [2] Sean Bell, Kavita Bala, and Noah Snavely. Intrinsic images in the wild. *ACM TOG*, 33(4):1–12, 2014. 3
- [3] Sean Bell, Paul Upchurch, Noah Snavely, and Kavita Bala. Material recognition in the wild with the materials in context database. In *CVPR*, pages 3479–3487, 2015. 3
- [4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *CVPR*, pages 9650–9660, 2021. 2
- [5] Boyuan Chen, Fei Xia, Brian Ichter, Kanishka Rao, Keerthana Gopalakrishnan, Michael S Ryoo, Austin Stone, and Daniel Kappler. Open-vocabulary queryable scene representations for real world planning. *arXiv preprint arXiv:2209.09874*, 2022. 1, 3
- [6] Wei Chow, Jiageng Mao, Boyi Li, Daniel Seita, Vitor Guizilini, and Yue Wang. Physbench: Benchmarking and enhancing vision-language models for physical world understanding. *arXiv preprint arXiv:2501.16411*, 2025. 1
- [7] Francis Engelmann, Fabian Manhardt, Michael Niemeyer, Keisuke Tateno, Marc Pollefeys, and Federico Tombari. Opennerf: open set 3d neural scene segmentation with pixel-wise features and rendered novel views. *arXiv preprint arXiv:2404.03650*, 2024. 3
- [8] Michael Fischer, Iliyan Georgiev, Thibault Groueix, Vladimir G Kim, Tobias Ritschel, and Valentin Deschaintre. Sama: Material-aware 3d selection and segmentation. *arXiv preprint arXiv:2411.19322*, 2024. 1
- [9] Katerina Fragkiadaki, Pulkit Agrawal, Sergey Levine, and Jitendra Malik. Learning visual predictive models of physics for playing billiards. *arXiv preprint arXiv:1511.07404*, 2015. 3
- [10] Yuanming Hu, Luke Anderson, Tzu-Mao Li, Qi Sun, Nathan Carr, Jonathan Ragan-Kelley, and Frédo Durand. DiffTaichi: Differentiable programming for physical simulation. *arXiv preprint arXiv:1910.00935*, 2019. 1
- [11] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hananeh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. 6
- [12] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM TOG*, 42(4):139–1, 2023. 3
- [13] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language embedded radiance fields. In *ICCV*, pages 19729–19739, 2023. 1, 3, 5
- [14] Xuan Li, Yi-Ling Qiao, Peter Yichen Chen, Krishna Murthy Jatavallabhula, Ming Lin, Chenfanfu Jiang, and Chuang Gan. Pac-nerf: Physics augmented continuum neural radiance fields for geometry-agnostic system identification. *arXiv preprint arXiv:2303.05512*, 2023. 3
- [15] OpenAI. Gpt-4o (“o” for omni) system card. Technical report, OpenAI, 2024. arXiv:2410.21276, https://doi.org/10.48550/arXiv.2410.21276. 6
- [16] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 2, 3, 6
- [17] Songyou Peng, Kyle Genova, Chiyu Jiang, Andrea Tagliasacchi, Marc Pollefeys, Thomas Funkhouser, et al. Openscene: 3d scene understanding with open vocabularies. In *CVPR*, pages 815–824, 2023. 3
- [18] Lerrel Pinto, Dhiraj Gandhi, Yuanfeng Han, Yong-Lae Park, and Abhinav Gupta. The curious robot: Learning visual representations via physical interactions. In *ECCV*, pages 3–18. Springer, 2016. 3
- [19] Ri-Zhao Qiu, Ge Yang, Weijia Zeng, and Xiaolong Wang. Feature Splatting: Language-Driven Physics-Based Scene Synthesis and Editing, Apr. 2024. arXiv:2404.01223 [cs]. 1, 3
- [20] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1, 3
- [21] Ashutosh Saxena, Justin Driemeyer, and Andrew Y Ng. Robotic grasping of novel objects using vision. *The International Journal of Robotics Research*, 27(2):157–173, 2008. 1
- [22] Lavanya Sharan, Ce Liu, Ruth Rosenholtz, and Edward H Adelson. Recognizing materials using perceptually inspired features. *IJCV*, 103(3):348–371, 2013. 3
- [23] Yinghao Shuai, Ran Yu, Yuantao Chen, Zijian Jiang, Xiaowei Song, Nan Wang, Jv Zheng, Jianzhu Ma, Meng Yang, Zhicheng Wang, et al. Pugs: Zero-shot physical understanding with gaussian splatting. *arXiv preprint arXiv:2502.12231*, 2025. 1, 3, 6, 7, 8
- [24] Trevor Standley, Ozan Sener, Dawn Chen, and Silvio Savarese. image2mass: Estimating the mass of an object from its image. In *Conference on Robot Learning*, pages 324–333. PMLR, 2017. 1, 3
- [25] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. VggT: Visual geometry grounded transformer. In *CVPR*, pages 5294–5306, 2025. 3
- [26] Jiajun Wu, Ilker Yildirim, Joseph J Lim, Bill Freeman, and Josh Tenenbaum. Galileo: Perceiving physical object properties by integrating a physics engine with deep learning. *NeurIPS*, 28, 2015. 3
- [27] Xinli Xu, Wenheng Ge, Dicong Qiu, ZhiFei Chen, Dongyu Yan, Zhuoyun Liu, Haoyu Zhao, Hanfeng Zhao, Shunsi Zhang, Junwei Liang, et al. Gaussianproperty: Integrating physical properties to 3d gaussians with lms. In *CVPR*, pages 7231–7240, 2025. 1, 3
- [28] Shaoxiong Yao and Kris Hauser. Estimating tactile models of heterogeneous deformable objects in real time. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 12583–12589. IEEE, 2023. 3
- [29] Albert J. Zhai, Yuan Shen, Emily Y. Chen, Gloria X. Wang, 634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692

- 693 Xinlei Wang, Sheng Wang, Kaiyu Guan, and Shenlong
694 Wang. Physical Property Understanding from Language-
695 Embedded Feature Fields. pages 28296–28305, 2024. 1
696 [30] Albert J Zhai, Yuan Shen, Emily Y Chen, Gloria X Wang,
697 Xinlei Wang, Sheng Wang, Kaiyu Guan, and Shenlong
698 Wang. Physical property understanding from language-
699 embedded feature fields. In *CVPR*, pages 28296–28305,
700 2024. 1, 2, 3, 6, 7, 8
701 [31] Junbo Zhang, Runpei Dong, and Kaisheng Ma. Clip-fo3d:
702 Learning free open-world 3d scene representations from 2d
703 dense clip. In *CVPR*, pages 2048–2059, 2023. 3