

Image2Garment: Simulation-ready Garments from a Single Image

Anonymous CVPR submission

Paper ID

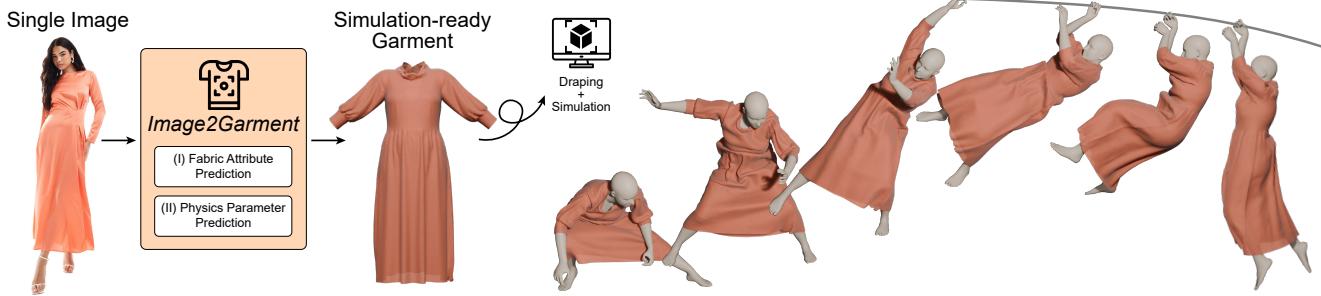


Figure 1. Image2Garment is a feedforward framework that predicts simulation-ready garments from a single image. For this purpose, garment geometry and physical fabric parameters are jointly predicted and used to simulate garment dynamics. Our framework is able to produce complex garments with delicate dynamics, as shown in this example.

Abstract

001 *Estimating physically accurate, simulation-ready garments*
002 *from a single image is challenging due to the absence of*
003 *image-to-physics datasets, and the ill-posed nature of this*
004 *problem. Prior methods either require multi-view capture*
005 *and expensive differentiable simulation or predict only gar-*
006 *ment geometry without the material properties required for*
007 *realistic simulation. We propose a feed-forward frame-*
008 *work that sidesteps these limitations by first fine-tuning a*
009 *vision-language model to infer material composition and*
010 *fabric attributes from real images, and then training a light-*
011 *weight predictor that maps these attributes to the corre-*
012 *sponding physical fabric parameters using a small dataset*
013 *of material-physics measurements. Our approach intro-*
014 *duces two new datasets (FTAG and T2P) and delivers*
015 *simulation-ready garments from a single image without*
016 *iterative optimization. Experiments show that our esti-*
017 *mator achieves superior accuracy in material composition es-*
018 *timation and fabric attribute prediction, and by passing*
019 *them through our physics parameter estimator, we further*
020 *achieve higher fidelity simulations compared to state-of-*
021 *the-art image-to-garment methods.*

022 1. Introduction

023 Generating simulation-ready garments from visual observa-
024 tions is increasingly important for applications in virtual re-

025 ability, gaming, and fashion design. Although recent meth-
026 ods can recover physical fabric properties, they typically
027 rely on multi-view capture setups or manual material mea-
028 surement, both of which are labor-intensive and impractical
029 outside controlled environments. While single-image infer-
030 ence would be an accessible alternative, estimating a gar-
031 ment’s physical properties (e.g., stretch and bend stiffness,
032 density, damping) from an in-the-wild image remains pro-
033 foundly challenging due to limited viewpoints, ambiguous
034 drape cues, and the lack of direct supervision.

035 Recent progress in single-image garment generation,
036 driven by fine-tuning Vision-Language Models (VLMs),
037 has enabled substantial improvements in predicting garment
038 geometry and appearance [12, 45, 51, 56, 81]. However,
039 these methods largely neglect physical parameters nec-
040 essary for faithful simulation (see Fig. 2). Conversely, meth-
041 ods that optimize physical parameters via a differentiable
042 simulation [27, 48, 49, 62, 80] require multi-view inputs,
043 involve slow iterative optimization, and are typically re-
044 stricted to simple garment categories, making them im-
045 practical for real-world deployment. As a result, obtain-
046 ing simulation-ready garments from a single image remains
047 a challenging problem. While optimization-based methods
048 struggle with the limited supervision available from a sin-
049 gle image, directly training a neural network for this task
050 is equally infeasible due to the lack of appropriate data.
051 No large-scale dataset that pairs real-world garment images

052 with the physical parameters required by cloth simulators
 053 exists, and creating such a resource would demand prohibitive manual effort.
 054

055 Our key insight is to reformulate this inverse problem
 056 through a semantically grounded latent decomposition. Al-
 057 though direct image-to-physics supervision is unavailable,
 058 image-to-material information is abundant. Large online
 059 catalogs provide reliable material-composition labels (e.g.,
 060 cotton, silk, polyester blends) as well as complementary
 061 garment attributes such as fabric family, weave structure,
 062 and thickness indicators. Importantly, these descriptors oc-
 063 cupy a structured and relatively low-dimensional space that
 064 has a far more predictable relationship to the physical pa-
 065 rameters used in cloth simulation. As a result, the material-
 066 to-physics mapping is substantially easier to learn and re-
 067 quires only a modest amount of material-physics data that
 068 can realistically be collected from industry-standard sim-
 069 ulators.

070 Building on this observation, we introduce a two-stage
 071 factorization. First, we train a model to infer an inter-
 072 pretable set of material descriptors from the single input im-
 073 age. Second, using a compact set of independently gathered
 074 material-physics annotations, we train a mapper that con-
 075 verts these descriptors into the full set of simulator pa-
 076 rameters. This latent-variable formulation regularizes the learn-
 077 ing task, dramatically reduces data requirements, and re-
 078 solves the ill-posedness that makes direct image-to-physics
 079 prediction impractical. The result is a fully feed-forward,
 080 optimization-free image-to-simulation pipeline.

081 In summary our contributions are:

- Two datasets: (1) The Fabric Attributes from Garment Tags (FTAG) dataset—containing images annotated with material composition, fabric family, and structure type—and (2) the Tag-to-Physics (T2P) Dataset which links fabric attributes with measurable physical parameters directly compatible with an industry standard simulator.
- A feed-forward framework for single-image to simulation-ready garment generation, jointly predicting garment geometry and interpretable fabric descriptors, and mapping them to physically realistic material parameters.
- Extensive experiments demonstrating superior accuracy and speed, showing that our approach predicts accurate garment physics parameters that can be directly plugged in to image-to-simulation workflows.

097 2. Related Work

098 **Garment reconstruction and generation.** Many works
 099 have explored how to reconstruct garments from multi-view
 100 images and videos. Works such as [2, 3, 11, 23, 34, 36, 78]
 101 take videos, multi-view images, or point clouds as input
 102 to recover the garment geometry. However, they do not
 103 model the material composition of the garments, and there-



099 **Figure 2. Impact of garment fabric parameters on simulation.**
 100 We visualize the final frame of a jumping animation for four differ-
 101 ent fabrics (wool, cork–cotton, polyester, and a random material)
 102 each starting from the exact same initial condition. The choice
 103 of garment physics parameters changes the dynamics of the ani-
 104 mation drastically. In turn, this makes it critical to estimate these
 105 parameters accurately in image-to-garment generation settings to
 106 faithfully predict shape, appearance, and dynamics of a garment
 107 from visual observations.

108 fore cannot be directly used for physics simulation. Us-
 109 ing inverse rendering and inverse physics, other approaches
 110 [27, 41, 58, 62, 80] optimize both the garment geometry
 111 and its physical parameters. However, they require studio-
 112 quality videos as input, making the algorithm impractical
 113 for casual users.

114 More recently, generative models have enabled single-
 115 image-to-garment generation. Works such as [4, 29, 31, 33,
 116 39, 63, 71, 72] generate 3D meshes of garments fused with
 117 humans, requiring post-processing to separate the garment
 118 from the human mesh. Others [18, 43, 44, 55, 61] model
 119 the human mesh separately from the garment, but still lack
 120 the necessary physics parameters for simulation. Another
 121 line of works [12, 28, 32, 46, 50, 51, 56, 74, 81] directly
 122 generates the sewing pattern – a CAD representation of gar-
 123 ments – from a single image by training on sewing pattern
 124 datasets [37, 38]. While this representation is more sim-
 125 ulation ready as it automatically includes garment texture
 126 maps, users still need to specify physical parameters man-
 127 ually for each fabric panel.

128 The work closest to ours is Dress 1-to-3 [45], which es-
 129 timates both sewing pattern and physical parameter from
 130 a single image, leveraging inverse optimization. However,
 131 their method requires hours of optimization per garment and
 132 only outputs specific physical parameters that are not gen-
 133 eralizable for other types of simulators. Moreover, code is
 134 not available for this method.

135 **Optimizing physics parameters from images and videos.**
 136 Estimating physics from visual observations has become an
 137 increasingly important topic in computer vision [6, 7, 14,
 138 20, 25, 26, 35, 40, 42, 67–69, 75, 77]. Several works have
 139 also attempted to obtain physics parameters for garments

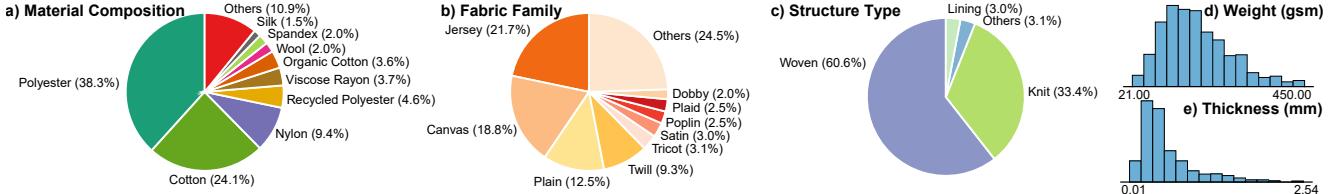


Figure 3. Garment physics parameters dataset. Our dataset contains 1,254 samples linking interpretable fabric attributes to measured fabric mechanical properties used for garment simulation. The fabric attributes include: material composition capturing fiber-level makeup that governs mechanical behavior; fabric family describing microstructure, drape, and surface appearance ; structure type identifying fabric construction that restricts the physically plausible range of simulation parameters; and density and thickness, which directly influence mass. Subfigures show distributions for (a) main fiber compositions, (b) fabric families, (c) structure types, (d) thickness, and (e) weight.

from visual data [8, 10, 27, 47–49, 54, 59, 62, 64, 70, 73, 80]. However, all of these methods require strong visual supervision with multi-view or video setups and most are not compatible with established cloth simulators. This clearly shows the difficulty of predicting physics parameters usable in general physics simulators from a single image.

High-quality estimation of garment physical parameters has been studied primarily in controlled settings. Wang et al. [65] fit elastic models from measurements of real fabric samples, and Bhat et al. [9] estimate parameters by matching simulated dynamics to video observations. The latter work refines these measurement-based approaches, including woven-fabric models [16], friction estimation [60], and differentiable-physics pipelines using multi-view or depth data [76, 79]. These methods achieve high fidelity but depend on specialized capture setups. In contrast, we target parameter estimation from a single RGB image.

3. Method

Overview. Our goal is to recover a *simulation-ready garment* from a single RGB image I . We define a simulation-ready garment as $G = (S, \theta)$, where S is the 3D garment shape and θ are the physical parameters required by a cloth simulator (e.g., stretch, bend, damping, friction). We therefore seek to estimate

$$\hat{G} = \arg \max_G p(G \mid I) = \arg \max_{S, \theta} p(S, \theta \mid I). \quad (1)$$

While single-view shape estimation $p(S \mid I)$ has seen rapid progress, estimating θ directly from I is underexplored: paired image-to-physics data do not exist and the mapping is ambiguous. Our key insight is to introduce a latent variable, *material* M (and auxiliary garment attributes Z ; e.g., weave, thickness, finish), and decompose the posterior as¹

$$p(S, \theta \mid I) = p(S \mid I) \sum_{m, z} p(\theta \mid m, z) p(m, z \mid I). \quad (2)$$

¹We make the assumption that the garment geometry S and physical parameters θ are conditionally independent given the input image I .

Therefore, to obtain a simulation-ready garment G from I , we use a three-stage pipeline: (i) predict a garment geometry G from I ; (ii) infer $p(M, Z \mid I)$ using abundant web-scale image-to-material supervision; and (iii) learn a material-to-physics mapper $f_\phi : (M, Z) \mapsto \theta$ from measured material datasets. In practice, we implement

$$\hat{S} = \arg \max_S p(S \mid I), \quad (3)$$

$$(\hat{M}, \hat{Z}) = \arg \max_{m, z} p(m, z \mid I), \quad (4)$$

$$\hat{\theta} = f_\phi(\hat{M}, \hat{Z}), \quad (5)$$

so that $\hat{G} = (\hat{S}, \hat{\theta})$. This split concentrates uncertainty into the first stage $p(M, Z \mid I)$ (where rich supervision exists) and turns physics estimation into a well-posed supervised mapping $p(\theta \mid M, Z)$.

We use an existing image-to-garment geometry model to estimate G . To obtain (m, z) and θ , we develop a custom estimation framework from in-the-wild images, leveraging our novel datasets. In the following, we first describe how we obtain these datasets and then how we train the corresponding models.

3.1. Dataset Curation

Fabric attributes from garment tags (FTAG) dataset. We curate a garment material dataset for *fabric composition understanding* in apparel, containing **16,026** images of people wearing garments paired with the corresponding material composition, fabric family, and fabric structure type annotations. The dataset was curated from publicly available online clothing stores and subsequently processed to filter out multi-layered garments and invalid data. Each entry corresponds to a single retail garment with vendor-provided material composition including exact fiber percentages, and a text description of the garment, from which we classify the fabric family and the fabric structure type. Our dataset reflects global fiber production trends [21, 24], with a strong predominance of polyester and cotton, followed by viscose and nylon. We provide further details in the supplementary material.

205 Tag-to-Physics (T2P) Dataset. We curate a dataset to map
206 intrinsic fabric properties to measurable physics parameters
207 for simulation, containing **2,254** fabrics paired with
208 their corresponding fabric attributes and physical measurements.
209 The dataset was compiled from publicly available online sources.
210 Each entry corresponds to a single digitized fabric characterized by (1) *fabric attributes*, including multi-material composition $C = \{(m_j, p_j)\}_{j=1}^N$, where N is the number of constituent materials, m_j denotes the j -th fiber type, and p_j its corresponding percentage, fabric family f , structure type s , areal density ρ (g/m^2), and thickness t (mm); and (2) *physics parameters*, which quantify the mechanical response of fabrics under deformation. We denote the full set of physics parameters as $\mathbf{P} \in \{\rho, \text{friction}, \text{internal damping}, \text{buckling stiffness } (\text{g}\cdot\text{mm}^2/\text{s}^2), \text{buckling ratio}, \text{bending stiffness } (\text{g}\cdot\text{mm}^2/\text{s}^2), \text{shear stiffness } (\text{g}/\text{s}^2), \text{stretch stiffness } (\text{g}/\text{s}^2)\}$. For parameters that are anisotropic, there are distinct values for the weft, warp, and bias directions. Each parameter describes a distinct aspect of a fabric's mechanical behavior—governing resistance to compression, bending, in-plane shear, and tension—and are used in standard commercial garment design software such as Blender [17], MarvelousDesigner [22], and Browzwear [53]. These values are derived from standardized mechanical testing of textiles, providing a supervised mapping from interpretable fabric attributes to their measurable physical response. Fig. 3 visualizes the resulting distributions of fabric families, weave structures, material compositions, thickness, and areal density.

234 3.2. Simulation-ready Garment Pipeline

235 Our single-image to simulation-ready garment generation
 236 pipeline consists of three stages. First, we reconstruct the
 237 garment geometry by estimating its sewing pattern using
 238 ChatGarment [12]. Because ChatGarment is trained on a
 239 diverse set of body poses to estimate sewing patterns draped
 240 on a standard A-pose human body, we can recover the
 241 garment geometry by draping the pattern with a standard cloth
 242 simulator [37]. Next, we predict the garment's material
 243 composition and fabric attributes, including fabric family,
 244 structure type, areal density, and thickness. Finally, we map
 245 these descriptors to the physical parameters required by the
 246 simulator using a learned material-to-physics model. An
 247 overview of the full pipeline is shown in Fig. 4. In the
 248 remainder of this section, we focus on our main technical
 249 contribution: predicting fabric attributes and simulator-
 250 compatible physical parameters from a single image.

251 3.3. Fabric Attribute Prediction

252 In this step, we estimate the material composition M and
 253 garment attributes Z from a single image, where M represents
 254 the fiber-level material mixture and Z includes attributes such as fabric family and structure type. A key ob-

servations is that these attributes differ significantly in how identifiable they are from visual cues alone. In particular, the conditional distributions $p(m | I)$ can be highly multi-modal: many distinct material compositions produce nearly indistinguishable appearances in photographs due to similar textures, colors, and weave patterns. By contrast, higher-level descriptors such as fabric family or structure type tend to exhibit clearer visual signatures and therefore admit a more peaked distribution. Attributes like areal density or thickness, however, often lack direct visual indicators and are therefore the most ambiguous. We use this insight to motivate the model design in the following.

268 Material estimation model. These observations lead us
 269 to estimate the material composition with a generative
 270 model which not only already possesses general knowledge
 271 regarding fabric materials and garments but also supports
 272 multi-modal inputs that are essential for material prediction.
 273 In particular, we formulate the prediction of the fabric at-
 274 tributes as an image captioning task. For this, we prompt
 275 an existing VLM with the unaltered image and a custom
 276 text prompt and ask it to output a structured JSON string
 277 containing per-fiber percentages and other attributes. In the
 278 JSON string, we ask the VLM to include the following fab-
 279 ric properties for each input garment image that are essential
 280 for downstream physical parameter estimation:

- Material composition $\hat{C} = \{(\hat{m}_j, \hat{p}_j)\}_{j=1}^{\hat{N}}$ denotes that the garment contains \hat{p}_j percent of fabric $\hat{m}_j \in \mathcal{M}$, a fixed set of apparel fibers². This representation captures the fiber-level makeup of the fabric, which directly governs its mechanical behavior.
- Fabric family $\hat{f} \in \mathcal{F}$, where \mathcal{F} is a set of common textile families³ such as denim, chiffon, or jersey, describes the fabric's microstructure, drape, and surface appearance.
- Structure type $\hat{s} \in \{\text{knit, woven, others}\}$ identifies the fabric construction. This structural prior restricts the physically plausible range of density, thickness, and stiffness parameters used in downstream simulation.

We use Qwen-2.5VL [5] as our base VLM model and fine-tune it with LoRA [30]. Because there is an imbalance in the frequency of materials in the FTAG dataset, we counteract this by specifying weights w_{t_i} as the inverse of the class frequency of the corresponding material during cross-entropy loss computation. Specifically, the model is fine-tuned using a weighted token-level cross-entropy loss:

$$\mathcal{L}_{\text{VLM}} = - \sum_i w_{t_i} \log p_\theta(t_i | \mathcal{I}, \mathbf{t}_{<i}), \quad 300$$

where \mathbf{t} is the serialized target sequence. 301

²E.g., cotton, polyester, elastane; full list in the Supplementary Material.

³See the Supplementary Material for a full list.

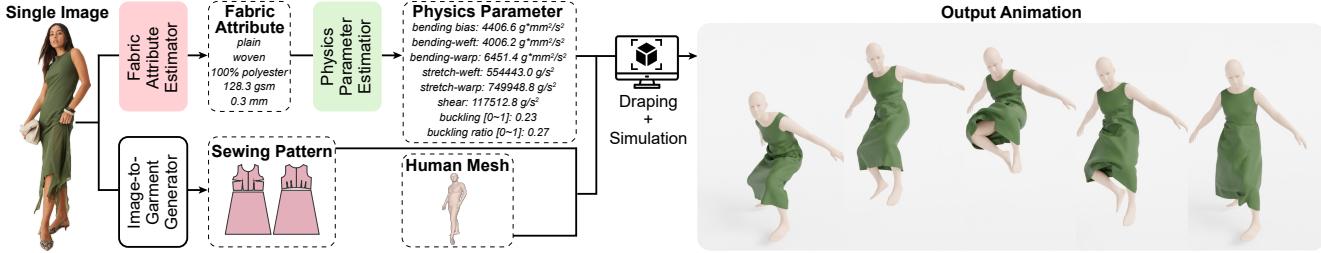


Figure 4. Overview of the Image2Garment pipeline. From a single image, we first generate the garment sewing pattern using Chat-Garment [12]. Then we predict fabric attributes such as material composition, fabric family, structure type, weight and thickness aligned with standardized commercial garment tags. Finally, fabric physics parameters are estimated from the predicted attributes, following [19], yielding mechanically interpretable quantities that describe fabric deformation. The garment geometry and physical parameters are then used to produce simulation-ready garments and physically accurate draping animations for any given motion sequence or body poses such as SMPL [52].

302 **Density–thickness estimator.** Given predicted fabric at-
 303 tributes—composition $\hat{C} = \{\hat{m}_j, \hat{p}_j\}_{j=1}^{\hat{N}}$, fabric family
 304 \hat{f} , and structure type \hat{s} —we estimate areal density $\hat{\rho}$ and
 305 thickness \hat{t} by sampling from our fabric dataset

$$306 \quad \mathcal{D} = \{(C_i, f_i, s_i, \rho_i, t_i)\}_{i=1}^M, \quad C_i = \{(m_{ij}, p_{ij})\}_{j=1}^{N_i}.$$

307 We perform a hierarchical search among fabrics sharing the
 308 same (f_i, s_i) : (1) exact composition match, (2) exact mate-
 309 rial type match, and (3) primary material match. If multiple
 310 candidates exist, $\hat{\rho}$ and \hat{t} are jointly sampled from the same
 311 fabric entry within the most specific non-empty set using
 312 one of three modes—*mean*, *median*, or *random*. We train
 313 using a 70 : 15 : 15 split (train / validation / test) of our
 314 fabric dataset and use 5-fold stratified cross-validation to
 315 select the sampling mode with the lowest validation MAE.
 316 This hierarchical retrieval preserves the empirical corre-
 317 lation between density and thickness while ensuring robust-
 318 ness when exact composition matches are unavailable.

319 3.4. Physics Parameter Prediction

320 We predict physics parameters \mathbf{P} from fabric attributes
 321 (C, f, s, ρ, t) using independent Random Forest Regressors
 322 (RFRs) [13] for bending, shear, stretch, buckling stiffness,
 323 and buckling ratio. Each Random Forest regressor is trained
 324 on ground-truth fabric attributes and corresponding physics
 325 parameters, following [19]. We train using a 70 : 15 : 15
 326 split (train / validation / test) of our fabric dataset and tune
 327 hyper-parameters via a 50-iteration randomized search with
 328 5-fold stratified cross-validation, selecting the configura-
 329 tion that minimizes validation MAE. During inference, we
 330 use predicted attributes $(\hat{C}, \hat{f}, \hat{s}, \hat{\rho}, \hat{t})$ as input to the same
 331 models. Friction μ and internal damping d are fixed con-
 332 stants, while ρ and t are directly taken from predicted fab-
 333 ric attributes. The final estimated physics parameter set is
 334 $\hat{\mathbf{P}} = \{\hat{B}, \hat{S}, \hat{E}, \hat{b}, \hat{r}, \mu, d, \hat{\rho}, \hat{t}\}$.

4. Experiments

335 4.1. Training Configuration

336 We fine-tune Qwen-2.5VL [5] on the FTAG dataset
 337 (Sec. 3.2) using its standard train/val/test split of
 338 9,843/1,231/1,231 samples. To address class imbalance
 339 in material composition, we apply inverse-frequency token
 340 weights. The model is trained for 5,200 iterations with a
 341 global batch size of 9, using LoRA [30] (rank $R = 64$) applied
 342 to all modules except [lm_head, embed_tokens]. We opti-
 343 mize with AdamW (learning rate 1×10^{-5} , weight
 344 decay 0.1, $(\beta_1, \beta_2) = (0.9, 0.999)$) and a cosine schedule.
 345 Images are dynamically resized between 1500×1500 and
 346 1700×1700 pixels using the Qwen2.5-VL image processor.
 347 Training runs on 9 Quadro RTX 6000/8000 GPUs and takes
 348 approximately 80 hours.

350 4.2. Simulation Setup

351 All garment simulations are performed in the Marvelous
 352 Designer (CLO3D) [15] engine, using its standard cloth
 353 solver and parameter definitions. We import the draped
 354 shapes together with our predicted physics parameters and
 355 simulate garments on the SMPL body for each sequence
 356 at 24 fps. Cloth resolution is set to 20 mm particle dis-
 357 tance, with timestep 0.042 s, gravity -9800 cm/s^2 , and
 358 default friction and damping settings (air damping 1.0).
 359 Self-collision and body–cloth collision are enabled using
 360 iteration-based collision detection (50 CG iterations). All
 361 methods are evaluated under identical simulator settings
 362 to ensure that differences in dynamic behavior arise solely
 363 from the predicted geometry and physics parameters.

364 4.3. Evaluation Datasets.

365 **4D-Dress Dataset** We evaluate garment reconstruction
 366 and dynamic draping quality on the 4D-Dress dataset [66].
 367 The dataset contains 50 multi-view video sequences of
 368 clothed subjects from which we select 5 sequences for eval-

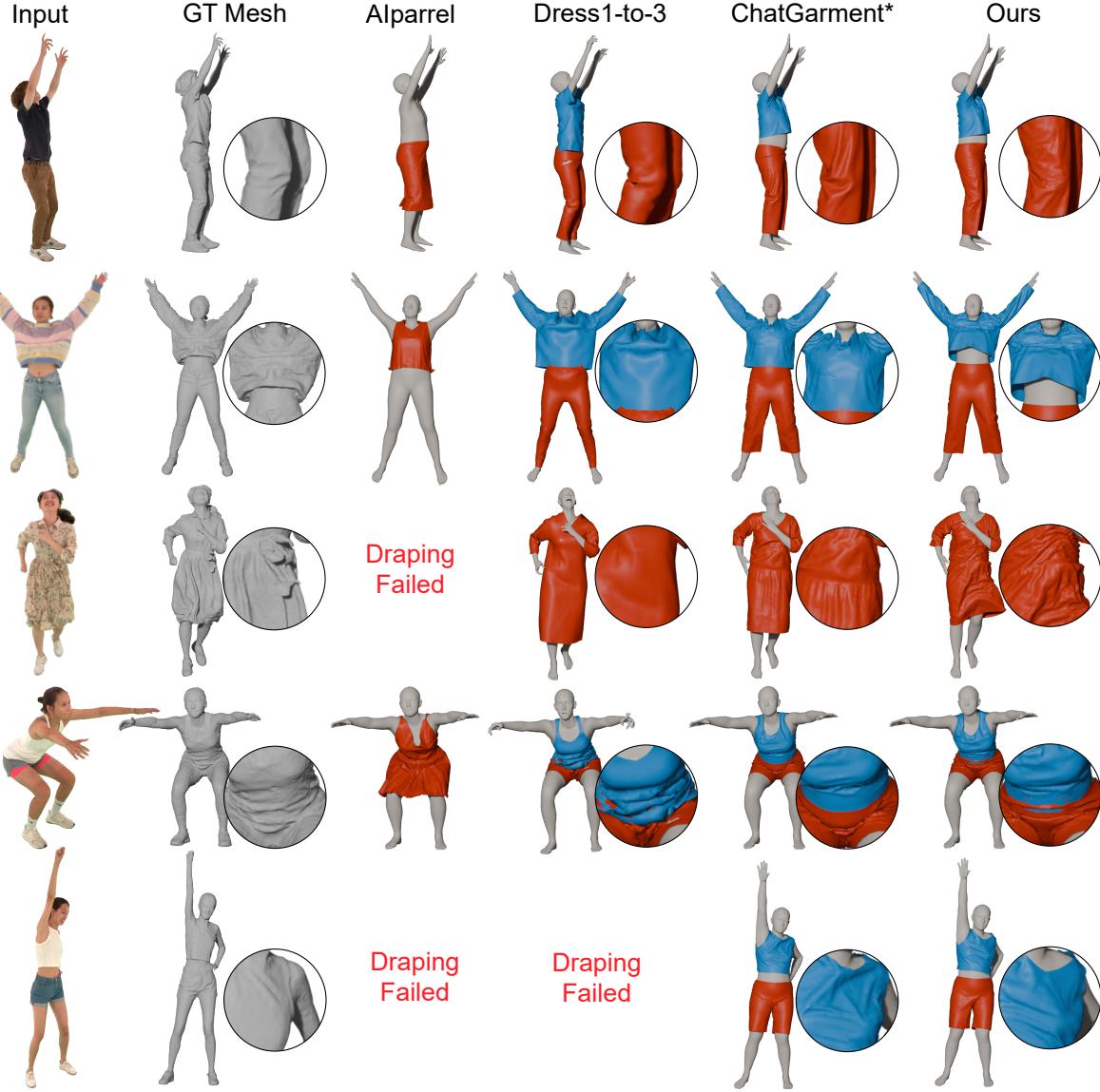


Figure 5. **Qualitative comparison of clothing reconstruction on 4D-Dress [66].** We visualize reconstructed garment meshes from different approaches. The zoomed-in image next to the garment meshes highlights the differences resulting from running the simulation. Compared to baselines, our model produces garments with more realistic shapes and wrinkle patterns that are better aligned with the ground truth geometry. *ChatGarment** uses the same garment geometry as our method but randomly sampled physical parameters, leading to incorrect garment deformation and dynamics.

369 uation. For each sequence, we use the first camera view as
 370 the single-image input and use the provided cloth meshes
 371 for evaluation.

372 **FTAG Dataset** For quantitative evaluation of material
 373 and fabric attribute prediction, we use the test split of our
 374 FTAG dataset, which includes 1,231 garments spanning 70
 375 distinct material compositions (669 single-fiber and 562
 376 multi-fiber blends), 28 fabric families, and 3 structure types,
 377 featuring challenging textures and diverse shapes.

4.4. Baselines

Single-image garment reconstruction. We compare against several state-of-the-art single-image garment reconstruction methods: ChatGarment [12], Aliparel [56], and Dress-1-to-3 [48]. We use the released checkpoints for ChatGarment and Aliparel. For Dress 1-to-3, we obtained the evaluating samples from the authors.

Material and fabric attribute prediction. We use OpenAI’s GPT-5 [57], as a baseline as it is often credited with

378

379

380

381

382

383

384

| Method | Lower | | Upper | |
|-------------------|-------------|-------------|-------------|-------------|
| | CD ↓ | IoU ↑ | CD ↓ | IoU ↑ |
| Dress 1-to-3 [48] | 34.1 | 44.8 | 40.5 | 40.0 |
| ChatGarment [12] | 28.9 | 48.4 | 28.1 | 46.6 |
| Alpparel [56] | 380 | 28.0 | 380 | 28.0 |
| Ours | 27.8 | 48.7 | 27.4 | 47.4 |

Table 1. **Clothing reconstruction benchmark.** We report the average Chamfer Distance (CD ↓) and Intersection over Union (IoU ↑) between the ground-truth garment meshes and the reconstructed clothing averaged over time for the lower and upper garment categories in 4D-Dress. Our method achieves the best score in both metrics, showcasing the importance of accurate physical parameters for garment dynamics. Reported CD is multiplied by 1e4.

387 having the strongest visual recognition and reasoning capabilities.
 388 We test it in both a zero-shot and few-shot manner,
 389 providing it with the same information about plausible materials
 390 that we make available to our own model.

4.5. Evaluation Metrics

392 **Garment prediction.** For the simulation-ready garment
 393 prediction, we focus on the garment geometry across the
 394 entire simulation. We compare it to the ground truth using
 395 Chamfer Distance between the point clouds sampled from
 396 the predicted and ground-truth mesh, as well as the Inter-
 397 section over Union (IoU) of the garment geometry. See the
 398 supplementary for details.

399 **Fabric attribute prediction.** For structure type and fabric
 400 family, we report macro-averaged accuracy and F1-score,
 401 computed by averaging the per-class metrics across all cat-
 402 egories. For material composition, we report per-example
 403 accuracy and F1-score averaged across all test examples.
 404 For material composition accuracy, we exclude true nega-
 405 tives from the calculation, as the large vocabulary of valid
 406 materials combined with the sparsity of material usage (gar-
 407 ments typically contain 1-3 materials) would inflate accu-
 408 racy scores artificially. For numerical values (material per-
 409 centages), we compare predictions with ground truth us-
 410 ing both the Mean Absolute Error (MAE) and Normalized
 411 Mean Absolute Error (NMAE).

4.6. Results

413 **Simulation-ready garment estimation.** Tab. 1 reports
 414 quantitative results on single-image garment reconstruction
 415 and dynamic draping. Across all metrics, our approach out-
 416 performs existing state-of-the-art methods. In particular,
 417 our simulated garments achieve the lowest Chamfer Dis-
 418 tance and highest IoU relative to ground-truth cloth geom-
 419 etry, demonstrating that both our reconstructed sewing pat-

terns and our predicted physics parameters are functionally accurate across the entire simulation.

In addition to accuracy, our method is substantially more efficient than optimization-based approaches such as Dress-1-to-3, which require iterative fitting. In contrast, our feed-forward pipeline produces reconstruction and physics estimates in a single pass, enabling consistent performance on in-the-wild inputs without any human intervention.

Fig. 5 provides a visual comparison on the 4D-Dress dataset. While all methods recover plausible static garment shapes, differences become pronounced under motion. Competing approaches frequently produce garments that over-stiffen, collapse, or drift away from the ground-truth dynamics. Our method, by contrast, generates garment motion and silhouettes that closely follow the ground-truth cloth over entire sequences. This highlights the importance of accurate fabric attribute and physics prediction: our inferred bending, shear, and stretch parameters yield dynamic deformations that better capture the true material behavior.

Fabric attribute estimation. Table 2 reports the performance of our fabric attribute estimation model compared to strong vision-language baselines. Across all fields, our fine-tuned Qwen2.5-VL model substantially outperforms both zero-shot and few-shot ChatGPT. For categorical attributes such as fabric family, structure type, and material type, our method yields markedly higher accuracy and F1-scores, with gains of up to 30% absolute accuracy over zero-shot prompting. These improvements reflect the benefits of domain-specific finetuning and structured JSON supervision, which enable the model to resolve subtle appearance cues that general-purpose VLMs fail to capture. For continuous-valued fields, including material percentage, areal density, and thickness, our approach achieves the lowest MAE and NMAE across all metrics. Notably, thickness prediction error is reduced by nearly half compared to zero-shot ChatGPT, and weight estimation shows similar improvements. These results demonstrate that accurate fabric attribute prediction requires specialized training rather than generic prompting, and they establish our model as a reliable foundation for downstream physics parameter estimation and garment simulation.

Ablations. We conduct an ablation study to assess the importance of major components in our framework, as summarized in Table 3. Fine-tuning proves essential across all three tasks. Without fine-tuning, the zero-shot model’s performance drops by 23, 23, and 35 percentage points on material composition, structure type, and fabric family estimation, respectively, compared to our fully fine-tuned model. This substantial performance gap demonstrates that task-specific adaptation is critical for reliable fabric property estimation from visual data alone. For material composi-

| Attribute Field | ChatGPT (zero-shot) | ChatGPT (few-shot) | Ours |
|---|-----------------------------|---------------------------|----------------------|
| <i>Categorical Fields (Accuracy % / F1-score) ↑</i> | | | |
| Fabric Family | 0.58 / 0.42 | <u>0.61</u> / <u>0.43</u> | 0.75 / 0.72 |
| Structure Type | 0.74 / 0.68 | <u>0.75</u> / <u>0.69</u> | 0.86 / 0.85 |
| Material Type | <u>0.65</u> / <u>0.70</u> | <u>0.66</u> / <u>0.70</u> | 0.71 / 0.75 |
| <i>Continuous Fields (MAE % / NMAE) ↓</i> | | | |
| Material Percentage | 23.3 / 0.45 | <u>22.4</u> / <u>0.43</u> | 19.3 / 0.40 |
| Density (g/m ²) | <u>64.38</u> / <u>0.121</u> | 75.62 / 0.143 | 51.28 / 0.097 |
| Thickness (mm) | <u>0.376</u> / <u>0.053</u> | 0.378 / 0.054 | 0.227 / 0.032 |

Table 2. **Performance of fabric attribute estimation baselines.** ChatGPT (zero-shot) and ChatGPT (few-shot) results are compared with our finetuned Qwen2.5-VL model. Categorical fields are evaluated using Accuracy and F1 Score, while continuous-valued fields are evaluated using MAE and NMAE (%). Weight and Thickness are estimated on the fabric test set using the thickness–density sampler (5-fold Cross Validation), and the remaining fields correspond to the stratified test split of our material composition estimation benchmark.

Table 3. **Ablation study.** We evaluate the impact of fine-tuning, image cropping, majority-only material prediction, and joint prediction with fabric attributes on material composition, structure type, and fabric family estimation tasks.

| Method Variant | Accuracy ↑ | F1-score ↑ |
|-----------------------------|------------|-------------|
| <i>Material Composition</i> | | |
| Full Model | 71% | 0.75 |
| w/ Primary Fiber | 62% | 0.68 |
| w/o Fabric Attributes | <u>69%</u> | <u>0.73</u> |
| w/o Fine-tuning | 48% | 0.56 |
| w Cropped Images | 63% | 0.67 |
| <i>Structure Type</i> | | |
| Full Model | 86% | 0.85 |
| w/o Fine-tuning | 63% | 0.33 |
| w Cropped Images | <u>75%</u> | <u>0.73</u> |
| <i>Fabric Family</i> | | |
| Full Model | 75% | 0.72 |
| w/o Fine-tuning | 40% | 0.31 |
| w Cropped Images | <u>58%</u> | <u>0.48</u> |

tion estimation, predicting only the primary material type (w/ Primary Fiber) leads to a noticeable drop in both accuracy ($71\% \rightarrow 62\%$) and F1-score ($0.75 \rightarrow 0.68$), confirming that capturing the full material composition is critical. Removing jointly predicted fabric attributes causes additional degradation ($71\% \rightarrow 69\%$ accuracy), indicating that these complementary cues provide useful supporting information for material composition estimation. Following [19], we also evaluate using images of the fabrics only. To isolate the fabric region, we apply a 512×512 center crop to each image. Center-cropping consistently harms performance across all tasks, with severe drops in material composition ($71\% \rightarrow 63\%$), structure type ($86\% \rightarrow 75\%$), and especially fabric family estimation ($75\% \rightarrow 58\%$ ac-

curacy). This demonstrates that high-resolution detail and complete garment context—including silhouette, draping patterns, and peripheral regions—provide essential visual information for reliably inferring fabric properties. Overall, each component contributes meaningfully to final estimation quality, with fine-tuning and full-image context being the most critical factors.

References

- [1] Adobe Creative Cloud. Mixamo, 2021. 5
- [2] Thiendo Alldieck, Marcus Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. Video based reconstruction of 3d people models. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8387–8397, 2018. 2
- [3] Thiendo Alldieck, Marcus Magnor, Bharat Lal Bhatnagar, Christian Theobalt, and Gerard Pons-Moll. Learning to reconstruct people in clothing from a single RGB camera. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [4] Thiendo Alldieck, Gerard Pons-Moll, Christian Theobalt, and Marcus Magnor. Tex2shape: Detailed full human body geometry from a single image. In *IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2019. 2
- [5] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhao-hai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. 4, 5
- [6] Anton Bakhtin, Laurens van der Maaten Wu, Justin Johnson, Quoc Le, and Ross Girshick. Phyre: A new benchmark for physical reasoning. *arXiv preprint arXiv:1908.05656*, 2019. 2
- [7] Peter W Battaglia, Razvan Pascanu, Matthew Lai, Danilo Rezende, and Koray Kavukcuoglu. Interaction networks for

- 521 learning about objects, relations and physics. In *Advances in*
522 *Neural Information Processing Systems (NeurIPS)*, 2016. 2
- 523 [8] Hugo Bertiche, Meysam Madadi, and Sergio Escalera. Neur-
524 al cloth simulation. *Computer Graphics Forum (Eurograph-*
525 *ics)*, 2022. 3
- 526 [9] Kiran S. Bhat, Christopher D. Twigg, and Jessica K. Hod-
527 gins. Estimating cloth simulation parameters from video.
528 In *Proceedings of the 2003 ACM SIGGRAPH/Eurographics*
529 *Symposium on Computer Animation (SCA)*, pages 37–51,
530 2003. 3
- 531 [10] Kiran S Bhat, Christopher D Twigg, and Jessica K Hodgins.
532 Estimating cloth simulation parameters from video. In *ACM*
533 *SIGGRAPH/Eurographics Symposium on Computer Anima-*
534 *tion (SCA)*, 2003. 3
- 535 [11] Bharat Lal Bhatnagar, Garvita Tiwari, Christian Theobalt,
536 and Gerard Pons-Moll. Multi-garment net: Learning to dress
537 3d people from images. In *IEEE International Conference on*
538 *Computer Vision (ICCV)*. IEEE, 2019. 2
- 539 [12] Siyuan Bian, Chenghao Xu, Yuliang Xiu, Artur Grigorev,
540 Zhen Liu, Cewu Lu, Michael J. Black, and Yao Feng. Chat-
541 garment: Garment estimation, generation and editing via
542 large language models, 2025. 1, 2, 4, 5, 6, 7
- 543 [13] Leo Breiman. Random forests. *Machine Learning*, 45(1):
544 5–32, 2001. 5
- 545 [14] Michael B Chang, Tomer Ullman, Antonio Torralba, and
546 Joshua B Tenenbaum. A compositional object-based ap-
547 proach to learning physical dynamics. *arXiv preprint*
548 *arXiv:1612.00341*, 2016. 2
- 549 [15] CLO Virtual Fashion. Marvelous designer, 2024. 3D gar-
550 ment design cloth simulation software. 5, 2
- 551 [16] Donald Clyde, Rasmus Tamstorf, and Joseph Teran. Model-
552 ing and data-driven parameter estimation for woven fabrics.
553 In *Proceedings of the ACM SIGGRAPH/Eurographics Sym-*
554 *posium on Computer Animation (SCA)*, pages 1–10, 2017.
555 3
- 556 [17] Blender Online Community. *Blender - a 3D modelling and*
557 *rendering package*. Blender Foundation, Stichting Blender
558 Foundation, Amsterdam, 2018. 4, 5
- 559 [18] Enric Corona, Albert Pumarola, Guillem Alenyà, Ger-
560 ard Pons-Moll, and Francesc Moreno-Noguer. Smplicit:
561 Topology-aware generative model for clothed people. In
562 *CVPR*, 2021. 2
- 563 [19] Henar Dominguez-Elvira, Alicia Nicás, Gabriel Cirio, Ale-
564 jandro Rodríguez, and Elena Garces. Practical Method to Es-
565 timate Fabric Mechanics from Metadata. *Computer Graph-*
566 *ics Forum*, 2024. 5, 8, 2, 4
- 567 [20] Sebastian Ehrhardt, Aron Monszpart, Niloy Mitra, and An-
568 drea Vedaldi. Unsupervised intuitive physics from visual
569 observations. In *Asian Conference on Computer Vision*
570 (*ACCV*), 2018. 2
- 571 [21] Textile Exchange. Preferred fiber & materials market report
572 2021. *Textile Exchange: Lamesa, TX, USA*, 4, 2021. 3
- 573 [22] CLO Virtual Fashion. Clo 3d – 3d fashion design software,
574 2009. Accessed: 2025-10-07. 4, 5
- 575 [23] Yao Feng, Jinlong Yang, Marc Pollefeys, Michael J. Black,
576 and Timo Bolkart. Capturing and animation of body and
577 clothing from monocular video. In *SIGGRAPH Asia 2022*
578 *Conference Papers*, 2022. 2
- [24] L Fernandez. Distribution of textile fibers production world-
wide in 2020, by type, 2021. 3
- [25] Pablo Garcia, Yue Sun, Taesung Kim, and Animesh Garg.
Learning physics from video: Unsupervised physical param-
eter estimation for continuous dynamical systems. In *Pro-
ceedings of the IEEE/CVF Conference on Computer Vision*
and *Pattern Recognition (CVPR)*, 2025. 2
- [26] Oliver Groth, Kerstin Wimmer, Andrea Vedaldi, and Chris-
tian Rupprecht. Shapestacks: Learning vision-based phys-
ical intuition for generalised object stacking. In *European*
Conference on Computer Vision (ECCV), 2018. 2
- [27] Michelle Guo, Matt Jen-Yuan Chiang, Igor Santesteban,
Nikolaos Sarafianos, Hsiao-yu Chen, Oshri Halimi, Aljaž
Božič, Shunsuke Saito, Jiajun Wu, C Karen Liu, et al. Pgc:
Physics-based gaussian cloth from a single pose. In *Pro-
ceedings of the Computer Vision and Pattern Recognition Confer-
ence*, pages 21215–21225, 2025. 1, 2, 3
- [28] Kai He, Kaixin Yao, Qixuan Zhang, Jingyi Yu, Lingjie Liu,
and Lan Xu. Dresscode: Autoregressively sewing and gen-
erating garments from text guidance, 2024. 2
- [29] Tong He, Yuanlu Xu, Shunsuke Saito, Stefano Soatto, and
Tony Tung. Arch++: Animation-ready clothed human re-
construction revisited. In *Proceedings of the IEEE/CVF In-
ternational Conference on Computer Vision (ICCV)*, pages
11046–11056, 2021. 2
- [30] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-
Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al.
Lora: Low-rank adaptation of large language models. *ICLR*,
1(2):3, 2022. 4, 5
- [31] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and
Tony Tung. ARCH: Animatable Reconstruction of Clothed
Humans . In *2020 IEEE/CVF Conference on Computer*
Vision and Pattern Recognition (CVPR), pages 3090–3099,
Los Alamitos, CA, USA, 2020. IEEE Computer Society. 2
- [32] Moon-Hwan Jeong, Dong-Hoon Han, and Hyeong-Seok Ko.
Garment capture from a photograph. *Comput. Animat. Vir-
tual Worlds*, 26(3–4):291–300, 2015. 2
- [33] Boyi Jiang, Juyong Zhang, Yang Hong, Jinhao Luo, Ligang
Liu, and Hujun Bao. Bcnets: Learning body and cloth shape
from a single image. In *Computer Vision – ECCV 2020: 16th*
European Conference, Glasgow, UK, August 23–28, 2020,
Proceedings, Part XX, page 18–35, Berlin, Heidelberg, 2020.
Springer-Verlag. 2
- [34] Boyi Jiang, Yang Hong, Hujun Bao, and Juyong Zhang. Sel-
frecon: Self reconstruction your digital avatar from monocular
video. In *IEEE/CVF Conference on Computer Vision and*
Pattern Recognition (CVPR), 2022. 2
- [35] Venkatesh Kandukuri, Korbinian Riedhammer, and Vladlen
Koltun. Physical representation learning and parameter iden-
tification from video via differentiable physics. *International*
Journal of Computer Vision (IJCV), 2021. 2
- [36] Maria Korosteleva and Sung-Hee Lee. Neuraltailor: Recons-
tructing sewing pattern structures from 3d point clouds of
garments. *ACM Trans. Graph.*, 41(4), 2022. 2
- [37] Maria Korosteleva and Olga Sorkine-Hornung. Garment-
Code: Programming parametric sewing patterns. *ACM*
Transaction on Graphics, 42(6), 2023. SIGGRAPH ASIA
2023 issue. 2, 4

- 637 [38] Maria Korosteleva, Timur Levent Kesdogan, Fabian Kem-
638 per, Stephan Wenninger, Jasmin Koller, Yuhan Zhang, Mario
639 Botsch, and Olga Sorkine-Hornung. GarmentCodeData: A
640 dataset of 3D made-to-measure garments with sewing pat-
641 terns. In *Computer Vision – ECCV 2024*, 2024. 2, 5
- 642 [39] Verica Lazova, Eldar Insafutdinov, and Gerard Pons-Moll.
643 360-Degree Textures of People in Clothing from a Single Im-
644 age . In *2019 International Conference on 3D Vision (3DV)*,
645 pages 643–653, Los Alamitos, CA, USA, 2019. IEEE Com-
646 puter Society. 2
- 647 [40] Long Le, Ryan Lucas, Chen Wang, Chuhao Chen, Dinesh Ja-
648 yaraman, Eric Eaton, and Lingjie Liu. Pixie: Fast and gen-
649 eralizable supervised learning of 3d physics from pixels. *arXiv*
650 preprint arXiv:2508.17437, 2025. 2
- 651 [41] Changmin Lee, Jihyun Lee, and Tae-Kyun Kim. Mpma-
652 vatar: Learning 3d gaussian avatars with accurate and robust
653 physics-based dynamics, 2025. 2
- 654 [42] Adam Lerer, Sam Gross, and Rob Fergus. Learning phys-
655 ical intuition of block towers by example. In *Proceedings*
656 of the 33rd International Conference on Machine Learning
657 (ICML), 2016. 2
- 658 [43] Ren Li, Corentin Dumery, Benoit Guillard, and Pascal Fua.
659 Garment Recovery with Shape and Deformation Priors. In
660 *Proceedings of the IEEE/CVF Conference on Computer Vi-
661 sion and Pattern Recognition*, 2024. 2
- 662 [44] Ren Li, Cong Cao, Corentin Dumery, Yingxuan You, Hao Li,
663 and Pascal Fua. Single view garment reconstruction using
664 diffusion mapping via pattern coordinates. In *Proceedings*
665 of the Special Interest Group on Computer Graphics and In-
666 teractive Techniques Conference Conference Papers, pages
667 1–11, 2025. 2
- 668 [45] Ren Li, Cong Cao, Corentin Dumery, Yingxuan You, Hao Li,
669 and Pascal Fua. Single view garment reconstruction using
670 diffusion mapping via pattern coordinates, 2025. 1, 2
- 671 [46] Siran Li, Chen Liu, Ruiyang Liu, Zhendong Wang, Gaofeng
672 He, Yong-Lu Li, Xiaogang Jin, and Huamin Wang. Garma-
673 genet: A multimodal generative framework for sewing pat-
674 tern design and generic garment modeling. *arXiv* preprint
675 arXiv:2504.01483, 2025. 2
- 676 [47] Xiaoyu Li, Xinyu Chen, Joshua B. Tenenbaum, and Jiajun
677 Wu. Neural garment dynamics via manifold-aware trans-
678 formers. In *Eurographics*, 2024. 3
- 679 [48] Xuan Li, Chang Yu, Wenxin Du, Ying Jiang, Tianyi Xie,
680 Yunuo Chen, Yin Yang, and Chenfanfu Jiang. Dress-1-to-
681 3: Single image to simulation-ready 3d outfit with diffu-
682 sion prior and differentiable physics. *ACM Transactions on*
683 *Graphics (TOG)*, 44(4):1–16, 2025. 1, 6, 7
- 684 [49] Yunzhu Li, Xinyang Liu, Guandao Yang, Gordon Wet-
685 zstein, Joshua B Tenenbaum, and Jiajun Wu. Diffavatar:
686 Simulation-ready 3d avatars with differentiable physics. In
687 *Proceedings of the IEEE/CVF Conference on Computer Vi-
688 sion and Pattern Recognition (CVPR)*, 2024. 1, 3
- 689 [50] Lijuan Liu, Xiangyu Xu, Zhijie Lin, Jiabin Liang, and
690 Shuicheng Yan. Towards garment sewing pattern recon-
691 struction from a single image. *ACM Transactions on Graphics*
692 (*SIGGRAPH Asia*), 2023. 2
- 693 [51] Shengqi Liu, Yuhao Cheng, Zhuo Chen, Xingyu Ren, Wen-
694 han Zhu, Lincheng Li, Mengxiao Bi, Xiaokang Yang, and
695 Yichao Yan. Multimodal latent diffusion model for com-
696 plex sewing pattern generation. *International Conference on*
697 *Computer Vision (ICCV)*, 2025. 1, 2
- 698 [52] Matthew Loper, Naureen Mahmood, Javier Romero, Ger-
699 ard Pons-Moll, and Michael J. Black. SMPL: A skinned
700 multi-person linear model. *ACM Trans. Graphics (Proc.
SIGGRAPH Asia)*, 34(6):248:1–248:16, 2015. 5
- 701 [53] Browzwear Solutions Pte Ltd. Browzwear — digital ap-
702 parel design and development platform, 2025. Accessed:
703 2025-11-12. 4, 2
- 704 [54] Weijia Mao, Chengcheng Tang, Robert Tong, Zhaoqi Wang,
705 and Ming C Li. Video classification of cloth simulations us-
706 ing deep networks. In *Pacific Graphics (PG)*, 2020. 3
- 707 [55] Gyeongsik Moon, Hyeongjin Nam, Takaaki Shiratori, and
708 Kyoung Mu Lee. 3d clothed human reconstruction in the
709 wild. In *European Conference on Computer Vision (ECCV)*,
710 2022. 2
- 711 [56] Kiyohiro Nakayama, Jan Ackermann, Timur Levent Kes-
712 dogan, Yang Zheng, Maria Korosteleva, Olga Sorkine-
713 Hornung, Leonidas J Guibas, Guandao Yang, and Gordon
714 Wetzstein. Aipparel: A multimodal foundation model for
715 digital garments. In *Proceedings of the Computer Vision and*
716 *Pattern Recognition Conference*, pages 8138–8149, 2025. 1,
717 2, 6, 7
- 718 [57] OpenAI. Gpt-5, 2025. Large language model accessed via
719 the OpenAI API. 6
- 720 [58] Bo Peng, Yunfan Tao, Haoyu Zhan, Yudong Guo, and
721 Juyong Zhang. Pica: Physics-integrated clothed avatar.
722 *IEEE Transactions on Visualization and Computer Graph-
723 ics*, pages 1–15, 2025. 2
- 724 [59] Gerard Pons-Moll, Sergi Pujades, Shihao Hu, and Michael J.
725 Black. Clothcap: Seamless 4d clothing capture and retarget-
726 ing. In *ACM Transactions on Graphics (SIGGRAPH)*, 2017.
727 3
- 728 [60] Ahmed Rasheed, Omar AlSibai, Luca Bergamini, Stelian
729 Coros, and Bernhard Thomaszewski. Learning to measure
730 the static friction coefficient in cloth contact. In *Proceed-
731 ings of the IEEE/CVF Conference on Computer Vision and*
732 *Pattern Recognition (CVPR)*, pages 4880–4889, 2020. 3
- 733 [61] Li Ren, Benoit Guillard, Edoardo Remelli, and Pascal Fua.
734 DIG: Draping Implicit Garment over the Human Body. In
735 *Asian Conference on Computer Vision*, 2022. 2
- 736 [62] Boxiang Rong, Artur Grigorev, Wenbo Wang, Michael J.
737 Black, Bernhard Thomaszewski, Christina Tsalicoglou,
738 and Otmar Hilliges. Gaussian Garments: Reconstruct-
739 ing simulation-ready clothing with photorealistic appearance
740 from multi-view video. In *International Conference on 3D*
741 *Vision 2025*, 2025. 1, 2, 3
- 742 [63] Shunsuke Saito, , Zeng Huang, Ryota Natsume, Shigeo Mor-
743 ishimura, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned
744 implicit function for high-resolution clothed human digitiza-
745 tion. *arXiv* preprint arXiv:1905.05172, 2019. 2
- 746 [64] Igor Santesteban, Elena Garces, and Miguel A. Otaduy.
747 Snug: Self-supervised neural dynamic garments. In *Pro-
748 ceedings of the IEEE/CVF Conference on Computer Vision*
749 *and Pattern Recognition (CVPR)*, 2022. 3

- 751 [65] Huamin Wang, Ravi Ramamoorthi, and James F. O'Brien.
752 Data-driven elastic models for cloth: Modeling and measure-
753 ment. *ACM Transactions on Graphics (TOG)*, 30(4):109,
754 2011. 3
- 755 [66] Wenbo Wang, Hsuan-I Ho, Chen Guo, Boxiang Rong, Artur
756 Grigorev, Jie Song, Juan Jose Zarate, and Otmar Hilliges.
757 4d-dress: A 4d dataset of real-world human clothing with se-
758 mantic annotations. In *Proceedings of the IEEE Conference*
759 *on Computer Vision and Pattern Recognition (CVPR)*, 2024.
760 5, 6
- 761 [67] Nicholas Watters, Andrea Tacchetti, Theophane Weber, Raz-
762 van Pascanu, Peter Battaglia, and Daniel Zoran. Visual
763 interaction networks: Learning a physics simulator from
764 video. In *Advances in Neural Information Processing Sys-
765 tems (NeurIPS)*, 2017. 2
- 766 [68] Jiajun Wu, Ilker Yildirim, Joseph J Lim, William T Free-
767 man, and Joshua B Tenenbaum. Galileo: Perceiving physical
768 object properties by integrating a physics engine with deep
769 learning. In *Advances in Neural Information Processing Sys-
770 tems (NeurIPS)*, 2015.
- 771 [69] Jiajun Wu, Erika Lu, Pushmeet Kohli, William T Freeman,
772 and Joshua B Tenenbaum. Learning to see physics via visual
773 de-animation. In *Advances in Neural Information Processing
774 Systems (NeurIPS)*, 2017. 2
- 775 [70] Donglai Xiang, Guandao Yang, Yang Zheng, Yunzhu Li, and
776 Gordon Wetzstein. Inverse garment and pattern modeling
777 with a differentiable simulator. In *CVPR Workshops*, 2024.
778 3
- 779 [71] Yuliang Xiu, Jinlong Yang, Dimitrios Tzionas, and
780 Michael J. Black. ICON: Implicit Clothed humans Obtained
781 from Normals. In *Proceedings of the IEEE/CVF Conference
782 on Computer Vision and Pattern Recognition (CVPR)*, pages
783 13296–13306, 2022. 2
- 784 [72] Yuliang Xiu, Jinlong Yang, Xu Cao, Dimitrios Tzionas,
785 and Michael J. Black. ECON: Explicit Clothed humans
786 Optimized via Normal integration. In *Proceedings of the
787 IEEE/CVF Conference on Computer Vision and Pattern
788 Recognition (CVPR)*, 2023. 2
- 789 [73] Yonghao Xu, Min Tang, and Robert Tong. Learning-based
790 bending stiffness parameter estimation for cloth simulation.
791 In *Computer Graphics Forum (Eurographics)*, 2021. 3
- 792 [74] Shan Yang, Zherong Pan, Tanya Amer, Ke Wang, Licheng
793 Yu, Tamara Berg, and Ming C. Lin. Physics-inspired gar-
794 ment recovery from a single-view image. *ACM Trans.
795 Graph.*, 37(5), 2018. 2
- 796 [75] Kexin Yi, Chuang Gan, Yunzhu Li, Jiajun Wu, Antonio Tor-
797 ralba, and Joshua B Tenenbaum. Cleverer: Collision events
798 for video representation and reasoning. In *International Con-
799 ference on Learning Representations (ICLR)*, 2020. 2
- 800 [76] Dohyeong Yoon et al. Bayesian differentiable physics for
801 cloth digitalization. *Journal of Computational Design and
802 Engineering*, 12(8):29–44, 2025. 3
- 803 [77] Tianyuan Zhang, Hong-Xing Yu, Rundi Wu, Brandon Y.
804 Feng, Changxi Zheng, Noah Snavely, Jiajun Wu, and
805 William T. Freeman. Physdreamer: Physics-based interac-
806 tion with 3d objects via video generation. *arXiv:2404.13026*,
807 2024. 2
- 808 [78] Yiqun Zhao, Chenming Wu, Binbin Huang, Yihao Zhi, Chen
809 Zhao, Jingdong Wang, and Shenghua Gao. Surfel-based
810 gaussian inverse rendering for fast and relightable dynamic
811 human reconstruction from monocular videos. *IEEE Trans-
812 actions on Pattern Analysis and Machine Intelligence*, pages
813 1–17, 2025. 2
- 814 [79] Qiuyu Zheng, Robert Smith, and Sanjiban Choudhury. Dif-
815 ferentiable cloth parameter identification and state estima-
816 tion in manipulation. *arXiv preprint arXiv:2311.05141*,
817 2023. 3
- 818 [80] Yang Zheng, Qingqing Zhao, Guandao Yang, Yifan Wang,
819 Donglai Xiang, Florian Dubost, Dmitry Lagun, Thabo
820 Beeler, Federico Tombari, Leonidas Guibas, and Gordon
821 Wetzstein. Physavatari: Learning the physics of dressed 3d
822 avatars from visual observations. In *European Conference
823 on Computer Vision (ECCV)*, 2024. 1, 2, 3
- 824 [81] Feng Zhou, Ruiyang Liu, Chen Liu, Gaofeng He, Yong-Lu
825 Li, Xiaogang Jin, and Huamin Wang. Design2garmentcode:
826 Turning design concepts to tangible garments through pro-
827 gram synthesis. In *Proceedings of the Computer Vision and
828 Pattern Recognition Conference*, pages 23712–23722, 2025.
829 1, 2

Image2Garment: Simulation-ready Garments from a Single Image

Supplementary Material

830 5. Supplementary Website

831 Please open the attached supplementary website to view our
 832 results in video, including dynamic draping simulations and
 833 qualitative comparisons.

834 6. Overview

835 This supplementary material provides additional details on
 836 our datasets, model hyperparameters, and evaluation proto-
 837 cols. We also include further quantitative and qualitative re-
 838 sults to complement the main paper. Specifically, we cover:

- 839 • Dataset curation and normalization for the Fabric At-
 tributes from Garment Tags (FTAG) and Tag-to-Physics
 (T2P) datasets;
- 840 • Implementation details of the density–thickness estimator,
 and the physics parameter regressor;
- 841 • Detailed evaluation metrics and protocols for both fabric
 attribute prediction and dynamic draping;
- 842 • Additional quantitative comparisons and qualitative ex-
 amples, including typical failure cases.

848 7. Additional Dataset Details

849 7.1. FTAG: Fabric Attributes from Garment Tags

850 Material Composition Estimation Dataset

851 We curate a garment material dataset for *fabric attribute*
 852 *prediction*, containing **16,026** images of models wearing re-
 853 tail garments paired with three key labels: *material com-
 854 position*, *fabric structure type*, and *fabric family*. Images
 855 were collected from publicly available online sources and
 856 processed through extensive cleaning, de-duplication, and
 857 normalization.

858 Each entry corresponds to a single retail garment with
 859 vendor-provided **material composition**

$$860 C = \{(m_j, p_j)\}_{j=1}^N,$$

861 where each fiber type m_j belongs to a controlled vocabu-
 862 lary, including:

863 We further classify the garment’s **fabric structure type**:

864 *woven*, *knit*, *others*

865 where woven fabrics correspond to interlaced yarns (e.g.,
 866 denim, poplin, twill), knit fabrics correspond to interlooped
 867 yarns (e.g., jersey, rib, fleece), and others include non-
 868 woven categories such as mesh, lace, and sequin fabrics.

869 Each garment also receives a **fabric family** label:

870 *chiffon*, *crepe*, *denim*, *satin*, *taffeta*, *poplin*, *twill*, *can-*
 871 *vas*, *velvet*, *corduroy*, *flannel*, *organza*, *voile*, *lace*,

tulle, *mesh*, *crochet*, *broderie*, *jacquard*, *sequin*, *metallic*,
leather, *fleece*, *jersey*, *rib knit*, *french terry*, *interlock*,
pique, *milano knit*, *ponte*, *scuba*, *tricot*, *plissé*,
gingham, *chambray*, *plain*, *cheesecloth*, *terry*, *unknown*

selected from a controlled taxonomy based on visual tex-
 876 ture, weave pattern, drape, and surface appearance.

877 Processing and Normalization Details

All garment entries come with vendor-provided material composition and fabric descriptions that accompany each product image, allowing us to access standardized fiber information and free-form fabric terminology. The **material composition** is already reported in a consistent structured format (e.g., “95% Polyester, 5% Elastane”), allowing us to extract fiber types directly and retain only materials belonging to our validated vocabulary. During this step, we also map synonymous fiber names to a canonical form (e.g., *Spandex* and *Lycra Elastane* → *Elastane*), ensuring that all materials follow a unified naming scheme. Samples containing unrecognized fibers or malformed percentage formats are discarded.

For the **fabric family** and **fabric structure type**, the vendor descriptions include a wide variety of naming conventions (e.g., “satin-style”, “ribbed jersey”, “denim-like weave”). To ensure consistency, we perform an exhaustive search over these free-form descriptors and map each variant to its canonical category in our controlled taxonomy. Examples include normalizing “satin-style” and “sateen” to *satin*, “ribbed knit” to *rib knit*, and “mesh lace” to *mesh*.

The **structure type** is also provided in the vendor meta-
 900 data and we standardize it into the categories *woven*, *knit*, or
 901 *others*. We remove entries where the structure type contra-
 902 dictions the fabric family (e.g., “woven jersey”) or where the
 903 information is incomplete or ambiguous.

Our dataset also contains garments for which vendors report material composition under multiple headers (e.g., *Main*, *Shell*, *Lining*, *Upper*), each corresponding to a different visible or internal layer of the garment. Because our goal is to predict a single material composition from a single image, we retain only entries whose composition is provided under a primary header (*Main* or other non-
 906 secondary headers) and discard garments where the compo-
 907 sition is split across multiple layers. This removes cases
 908 such as jackets worn over shirts or dresses with separate
 909 lining materials, where a single image cannot be associ-
 910 ated with a unique fiber composition. The resulting FTAG
 911 dataset statistics and distributions are shown in Figure 6.

We additionally discard any sample for which one or
 912 more of the three attributes —material composition, struc-
 913 ture type, and fabric family—do not have a valid value.

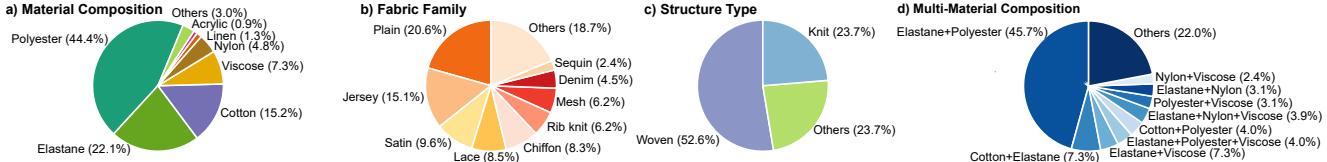


Figure 6. Fabric attribute dataset. Our Fabric Attributes from Garment Tags (FTAG) dataset contains 16,026 garment images paired with interpretable, vendor-provided fabric metadata. Each sample includes (i) *material composition*, specifying the fiber-level makeup that governs mechanical behavior; (ii) *fabric family*, characterizing microstructure, drape, and surface appearance; and (iii) *structure type*, indicating how yarns are interlaced or knitted to form the fabric. Subfigures show the distributions of (a) main fiber compositions, (b) fabric families, (c) structure types, and (d) multi-material fiber combinations.

ture type, or fabric family—was missing, inconsistent, or ambiguous. After applying these filtering steps, we obtain a final cleaned dataset consisting of **12,305** samples. We split this dataset in a stratified manner based on the joint label of material composition, structure type, and fabric family, resulting in **9,843** training, **1,231** validation, and **1,231** test samples for our fabric attribute estimation benchmark.

7.2. Tag-to-Physics (T2P) Dataset

We curate a dataset to convert intrinsic fabric properties to simulation physics parameters, containing **1,277** fabrics paired with corresponding fabric attributes and CLO 3D physics parameters. The dataset is collected from publicly available online sources and reflects real-world materials.

Each entry corresponds to a single digitized fabric with (1) fabric attributes, which consist of:

- *Material composition* including exact fiber percentages, $C \in \{\text{Acrylic, Alpaca, Angora, Cashmere, Cotton, Cupro, Hemp, Jute, Linen, Lyocell, Metallic, Modal, Nylon, Organic Cotton, PE (Polyethylene), PP (Polypropylene), PU (Polyurethane), PVC (Polyvinyl chloride), Polyester, Raccoon, Recycled Nylon, Recycled Polyester, Silk, Spandex, Supima Cotton, TENCEL™, TPE (Thermoplastic Elastomer), TPU (Thermoplastic Poly Urethane), Triacetate, Viscose Rayon, Viscose from Bamboo, Wool, vinyl acetate copolymer}\}$;
- *Fabric family* $f \in \{\text{Canvas, Challis, Chambray / Oxford, Chiffon, Circular Knit Spacer, Clip Jacquard, Corduroy, Crepe / CDC, Denim, Dobby, Dobby / Jacquard, Double Knit / Interlock, Double Weave, Flannel, Fleece, French Terry, Gauze / Double Gauze, Georgette, Interlining, Jacquard / Brocade, Jersey, Lace, Loop Terry, Memory, Mesh / Tulle, Organza, PVC, Pique, Plaid, Plain, Pointelle, Polar Fleece, Ponte, Poplin, Raschel, Rib, Ripstop, Satin, Seersucker, Sequin, Stretch Lining, Taffeta, Taffeta Lining, Tricot, Tweed, Twill, Vegan Fur, Vegan Leather, Vegan Suede, Velour, Velvet, Voile, Waffle}\}$;
- *Structure type* $s \in \{\text{Knit, Woven, Lining, Others}\}$;
- Areal density $\rho \in \mathbb{R} (\text{g/m}^2)$ and thickness $t \in \mathbb{R} (\text{mm})$.

(2) CLO 3D physics parameters quantify the mechanical response of fabrics under deformation:

$\in \{\rho, \text{friction, internal damping, buckling stiffness, } \{b_{bias-left}, b_{bias-right}, b_{warp}, b_{weft}\}, \text{buckling ratio, } \{r_{bias-left}, r_{bias-right}, r_{warp}, r_{weft}\}, \text{bending stiffness, } \{B_{bias-left}, B_{bias-right}, B_{warp}, B_{weft}\} (\text{g}\cdot\text{mm}^2/\text{s}^2), \text{shear stiffness, } \{S_{left}, S_{right}\} (\text{g}/\text{s}^2), \text{and stretch stiffness, } \{E_{warp}, E_{weft}\} (\text{g}/\text{s}^2)\}$.

Each parameter describes a distinct aspect of the fabric's mechanical behavior, governing resistance to compression, bending, in-plane shear, and tension, and is directly compatible with commercial garment design software such as Marvelous Designer [15], Browzwear [53], CLO 3D, and proprietary ones such as Seddi Author [19]. Each fabric sample is measured using the CLO Fabric Kit; we use these measurements without additional scale factors, enabling a supervised mapping from interpretable fabric attributes to measurable mechanical response. Figure 3 in the main paper visualizes the distributions of fabric families, structure types, material compositions, thickness, and areal density.

7.3. Material Prediction Metrics

We report two primary metrics for evaluating material prediction performance:

Average Accuracy:

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^N \frac{\text{TP}_i}{\text{TP}_i + \text{FP}_i + \text{FN}_i} \quad (6)$$

Average accuracy across all examples, where for each example i , we compute the ratio of correctly predicted materials (true positives) to the total number of unique materials in either the ground truth or predictions. N is the total number of examples.

Average F1 Score:

$$\text{F1} = \frac{1}{N} \sum_{i=1}^N \frac{2 \times \text{Precision}_i \times \text{Recall}_i}{\text{Precision}_i + \text{Recall}_i} \quad (7)$$

where $\text{Precision}_i = \frac{\text{TP}_i}{\text{TP}_i + \text{FP}_i}$ and $\text{Recall}_i = \frac{\text{TP}_i}{\text{TP}_i + \text{FN}_i}$.

Average F1 score across all examples, computed as the harmonic mean of precision and recall for each example.

993 For each example i , TP_i represents correctly identified ma-
994 terials, FP_i represents incorrectly predicted materials, and
995 FN_i represents ground truth materials that were missed.

996 7.4. Material Percentage Prediction Metrics

997 We report two metrics for evaluating the accuracy of pre-
998 dicted material percentages:

999 **Average Mean Absolute Error (MAE):**

$$1000 \quad \text{MAE} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|M_i|} \sum_{m \in M_i} |p_{i,m}^{\text{gt}} - p_{i,m}^{\text{pred}}| \quad (8)$$

1001 Average mean absolute error across all examples, where
1002 for each example i , we compute the absolute difference be-
1003 tween ground truth percentage $p_{i,m}^{\text{gt}}$ and predicted percent-
1004 age $p_{i,m}^{\text{pred}}$ for each material m in the set M_i (all unique ma-
1005 terials in either ground truth or predictions for example i).
1006 $|M_i|$ is the total number of materials for example i , and N
1007 is the total number of examples.

1008 **Average Normalized Mean Absolute Error (NMAE):**

$$1009 \quad \text{NMAE} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|M_i|} \sum_{m \in M_i} \frac{|p_{i,m}^{\text{gt}} - p_{i,m}^{\text{pred}}|}{\max(p_{i,m}^{\text{gt}}, p_{i,m}^{\text{pred}})} \quad (9)$$

1010 Average normalized mean absolute error across all ex-
1011 amples, where the absolute error for each material is nor-
1012 malized by the maximum of the ground truth and predicted
1013 percentages to account for varying material proportions.
1014 This provides a scale-invariant measure of percentage pre-
1015 diction accuracy.

1016 7.5. Structure Type Prediction Metrics

1017 We report two metrics for evaluating fabric structure type
1018 classification:

1019 **Overall Accuracy:**

$$1020 \quad \text{Accuracy} = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(s_i^{\text{gt}} = s_i^{\text{pred}}) \quad (10)$$

1021 Overall accuracy across all examples, where s_i^{gt} is the
1022 ground truth structure type for example i , s_i^{pred} is the pre-
1023 dicted structure type, and $\mathbb{1}(\cdot)$ is the indicator function that
1024 equals 1 when the condition is true and 0 otherwise. N
1025 is the total number of examples.

1026 **Macro-Averaged F1 Score:**

$$1027 \quad \text{F1}_{\text{macro}} = \frac{1}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} \text{F1}_s \quad (11)$$

1028 where

$$1029 \quad \text{F1}_s = \frac{2 \times \text{Precision}_s \times \text{Recall}_s}{\text{Precision}_s + \text{Recall}_s} \quad (12)$$

and $\text{Precision}_s = \frac{\text{TP}_s}{\text{TP}_s + \text{FP}_s}$, $\text{Recall}_s = \frac{\text{TP}_s}{\text{TP}_s + \text{FN}_s}$. 1030

Macro-averaged F1 score computed by first calculating
the F1 score for each structure type class $s \in \mathcal{S}$ (where \mathcal{S}
is the set of all structure types: {knit, woven, others}), then
averaging across all classes. TP_s , FP_s , and FN_s represent
true positives, false positives, and false negatives for class
 s , respectively. 1031

1032 7.6. Fabric Family Prediction Metrics

We report two metrics for evaluating fabric family classifi-
cation: 1033

Overall Accuracy: 1034

$$1035 \quad \text{Accuracy} = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(f_i^{\text{gt}} = f_i^{\text{pred}}) \quad (13) \quad 1036$$

Overall accuracy across all examples, where f_i^{gt} is the
ground truth fabric family for example i , f_i^{pred} is the pre-
dicted fabric family, and $\mathbb{1}(\cdot)$ is the indicator function that
equals 1 when the condition is true and 0 otherwise. N is
the total number of examples. 1037

Macro-Averaged F1 Score: 1038

$$1039 \quad \text{F1}_{\text{macro}} = \frac{1}{|\mathcal{F}|} \sum_{f \in \mathcal{F}} \text{F1}_f \quad (14) \quad 1040$$

where 1041

$$1042 \quad \text{F1}_f = \frac{2 \times \text{Precision}_f \times \text{Recall}_f}{\text{Precision}_f + \text{Recall}_f} \quad (15) \quad 1043$$

and $\text{Precision}_f = \frac{\text{TP}_f}{\text{TP}_f + \text{FP}_f}$, $\text{Recall}_f = \frac{\text{TP}_f}{\text{TP}_f + \text{FN}_f}$. 1044

Macro-averaged F1 score computed by first calculating
the F1 score for each fabric family class $f \in \mathcal{F}$ (where \mathcal{F}
is the set of all fabric families, including Poplin, Jersey, Lace,
Satin, Chiffon, etc.), then averaging across all classes. TP_f ,
 FP_f , and FN_f represent true positives, false positives, and
false negatives for class f , respectively. 1045

1046 7.7. Density and Thickness Estimation Metrics

We use our T2P dataset to estimate density and thickness
values from ground truth fabric family, structure type, and
material composition. We also prompt GPT to perform the
same estimation for comparison (see Table 2 in the main
paper). 1047

We report two metrics for evaluating continuous parameter
estimation: 1048

Mean Absolute Error (MAE): 1049

$$1050 \quad \text{MAE} = \frac{1}{N} \sum_{i=1}^N |y_i^{\text{gt}} - y_i^{\text{pred}}| \quad (16) \quad 1051$$

Mean absolute error between ground truth values y_i^{gt} and
predicted values y_i^{pred} across all N examples. 1052

1070 **Normalized Mean Absolute Error (NMAE):**

$$1071 \quad \text{NMAE} = \frac{\text{MAE}}{y_{\max} - y_{\min}} \quad (17)$$

1072 Normalized mean absolute error, where MAE is divided
1073 by the range of ground truth values ($y_{\max} - y_{\min}$) to provide
1074 a scale-invariant measure of prediction accuracy.

1075 7.8. Density–Thickness Estimator

1076 Given predicted discrete attributes $(\hat{C}, \hat{f}, \hat{s})$ from the VLM,
1077 we estimate areal density $\hat{\rho}$ and thickness \hat{t} using the T2P
1078 dataset $\mathcal{D} = \{(C_i, f_i, s_i, \rho_i, t_i)\}_{i=1}^M$. We perform hierarchical
1079 retrieval:

- 1080 1. **Exact attribute match:** we first search for fabrics with
1081 $(f_i, s_i) = (\hat{f}, \hat{s})$ and a composition C_i matching \hat{C} ex-
1082 actly up to a small tolerance in percentages.
- 1083 2. **Material-set match:** if no exact match exists, we relax
1084 to fabrics whose *set* of fibers matches that of \hat{C} , ignoring
1085 small percentage differences.
- 1086 3. **Primary fiber match:** if the previous sets are empty, we
1087 match fabrics whose primary fiber (highest percentage in
1088 C_i) equals the primary fiber in \hat{C} .

1089 From the most specific non-empty candidate set, we esti-
1090 mate $(\hat{\rho}, \hat{t})$ either by taking the mean, median, or by ran-
1091 domly sampling a single fabric entry. We select the aggrega-
1092 tion mode by performing 5-fold stratified cross-validation
1093 on the T2P dataset (train:val:test = 70:15:15) and choosing
1094 the mode with lowest validation MAE. The best performing
1095 mode was **mean aggregation**. This retrieval-based strategy
1096 preserves the empirical coupling between areal density and
1097 thickness while remaining robust when exact composition
1098 matches are not available.

1099 7.9. Physics Parameter Estimator

1100 We predict physics parameters P from fabric attributes
1101 (C, f, s, ρ, t) using a collection of independent Random
1102 Forest Regressors (RFRs), following the practical recipe of
1103 Dominguez-Elvira *et al.* [19]. We form a feature vector by
1104 concatenating:

- 1105 • normalized fiber percentages over the vocabulary,
- 1106 • one-hot encodings of fabric family f and structure type s ,
- 1107 • scalar features ρ and t (optionally log-transformed).

1108 We train five separate Random Forest Regressors, one
1109 for each physics parameter group. We use a 70/15/15
1110 train/val/test split and perform a 50-iteration random-
1111 ized hyperparameter search with 5-fold stratified cross-
1112 validation to select the number of trees and features, the
1113 max depth, and the minimum samples per split and leaf.

1114 For training all models, we use Mean Absolute Error
1115 as the loss function. The best hyperparameters selected
1116 via cross-validation are: for bending, stretch, shear, and
1117 buckling stiffness, $n_{\text{estimators}} = 100$, $\text{max_depth} = 20$,
1118 $\text{min_samples_split} = 0.0420$, $\text{min_samples_leaf} = 0.0094$,

max_features = 0.6911; for buckling ratio, $n_{\text{estimators}} = 1119$
200, $\text{max_depth} = 30$, $\text{min_samples_split} = 0.0703$,
1120 $\text{min_samples_leaf} = 0.0016$, max_features = 0.9595.
1121

1122 During inference, we feed predicted attributes
($\hat{C}, \hat{f}, \hat{s}, \hat{\rho}, \hat{t}$) into the same regressors to obtain \hat{P} . Friction
1123 μ and internal damping d are kept fixed across fabrics,
1124 while $\hat{\rho}$ and \hat{t} are directly taken from the density–thickness
1125 estimator.
1126

1127 8. Garment Geometry and Simulation Details

1128 8.1. ChatGarment-based Sewing Pattern Estima- 1129 tion

1130 For garment geometry, we adopt ChatGarment as our
1131 image-to-sewing-pattern backbone. Given a single im-
1132 age, ChatGarment predicts a sewing pattern representation
1133 draped on a canonical SMPL body in A-pose. We export
1134 the predicted pattern and re-simulate it using our cloth sim-
1135 ulator, ensuring that all methods (ours and baselines) are
1136 evaluated under identical numerical simulation settings.

1137 We use the default ChatGarment resolution and panel pa-
1138 rameterization and do not fine-tune ChatGarment on 4D-
1139 Dress to keep the comparison fair and emphasize the ben-
1140 efits of improved physics prediction.

1141 9. Evaluation Protocols and Metrics

1142 9.1. Dynamic Draping Metrics

1143 To evaluate dynamic draping on 4D-Dress, we compare
1144 our simulated garment meshes against ground-truth cloth
1145 meshes over the entire sequence. We use two metrics:

- 1146 • **Chamfer Distance (CD) ↓:** We sample points from both
1147 the predicted and ground-truth garment meshes at each
1148 frame and compute the symmetric Chamfer distance. The
1149 reported CD is averaged over frames and scaled by 10^4
1150 for readability.
- 1151 • **Intersection-over-Union (IoU) ↑:** We voxelize each
1152 mesh using 5cm voxels and compute the intersection-
1153 over-union of the occupied voxels between predicted and
1154 ground-truth garments, then average over frames.

1155 We report separate metrics for upper- and lower-body
1156 garments, as in Table 1 of the main paper. Garments span-
1157 ning upper and lower body, count towards both categories
1158 results.

1159 9.2. Fabric Attribute Prediction Metrics

1160 For fabric attribute estimation on FTAG, we evaluate:

- 1161 • **Structure type and fabric family:** macro-averaged accu-
1162 racy and F1-score, computed by averaging per-class met-
1163 rics across all categories.
- 1164 • **Material type:** multi-label accuracy and F1-score over
1165 the material vocabulary, ignoring true negatives to avoid
1166 inflated scores due to sparsity.

- 1167 • **Material percentages, density, and thickness:** mean ab-
1168 solute error (MAE) and normalized MAE (NMAE) over
1169 all test samples.

1170 These metrics are reported for our fine-tuned Qwen2.5-
1171 VL model and for ChatGPT baselines (zero- and few-shot)
1172 in the main paper.

1173 10. Additional Results

1174 10.1. Synthetic Dataset

1175 **Dataset creation.** To compare our method further with
1176 exiting state-of-the-art methods we create the first synthetic
1177 evaluation dataset designed with full garment prediction in
1178 mind. In particular our dataset contains three outfits sam-
1179 pled from [38] that each have been assigned a material per
1180 garment. We drape this garment over SMPL [52] bodies
1181 and create distinct animations from Mixamo [1]. We then
1182 use CLO 3D’s [22] simulator to simulate the garment un-
1183 der the animation. Finally, we render the garment at each
1184 step in the animation using Blender [17]. From this data,
1185 we use the first rendered frame as an input and the remain-
1186 ing frames and meshes for evaluation. Compared to existing
1187 evaluation benchmarks, this data has the advantage that
1188 it includes more complex garments and animations, and it
1189 does not suffer from measurement inaccuracies.

1190 **Quantitative Evaluation.** In this section, we evaluate our
1191 method on the dataset we created and compare it against
1192 state-of-the-art garment shape estimators paired with ran-
1193 dom material samples. Tab. 4 reports the quantitative re-
1194 sults. Our method achieves the best average performance
1195 across all 3D metrics and most 2D metrics. The only excep-
1196 tion is SSIM, where our score is slightly lower; however,
1197 this can be attributed to the uniformly high SSIM values
1198 achieved by all methods, which reduces the discriminative
1199 power of this metric.

1200 **Qualitative Evaluation.** Here, we present qualitative
1201 comparisons between our method and state-of-the-art gar-
1202 ment generation approaches using random physics par-
1203 ameters. Fig. 7 shows results for the first sequence in our
1204 dataset. While Alpparel and DMap struggle to recover
1205 accurate garment shapes, our method produces a plausi-
1206 ble geometry and estimates physics parameters that remain
1207 consistent with the ground truth throughout the animation.
1208 Compared to ChatGarment, our predictions avoid the exces-
1209 sive deformation observed in its outputs.

1210 Fig. 8 presents results for the second sequence. As be-
1211 fore, our predicted shape and physics parameters yield the
1212 closest visual match to the ground truth across the anima-
1213 tion. In contrast, ChatGarment and Alpparel consistently
1214 generate dresses that appear unnaturally folded.

Finally, Fig. 9 shows the third sequence. While Alpparel again fails to estimate the correct shape, our method remains close to the ground truth. During the rapid, high-energy motion in this sequence, the ground-truth garment slides off the human body due to its loose fit, a behavior that our method successfully reproduces but none of the baselines capture.

Ablations. In addition to the parameter ablations pre-
1215 sented in the main paper, we conduct further ablations at
1216 the 4D shape level using our curated dataset. To isolate the
1217 effect of physics parameter prediction, we use the ground-
1218 truth meshes and vary only the physics parameters. The
1219 quantitative results are reported in Tab. 5. Our predicted
1220 physics parameters yield substantial improvements across
1221 all metrics and sequences. Although both our method and
1222 the baseline exhibit some variance, our approach consis-
1223 tently outperforms the baseline overall, highlighting the im-
1224 portance of accurate physics parameter estimation.

1225 10.2. Additional Image-to-Simulation Examples

1226 **Qualitative Evaluation.** In the supplementary figures, we
1227 provide additional in-the-wild examples produced by our
1228 Image2Garment pipeline. For each example, we show the
1229 original sequence (first image as input), the reconstructed
1230 sewing pattern, the final simulation frame, and close-ups
1231 of characteristic wrinkles. Figures 10, 11, and 12 com-
1232 pare our results against simulations generated using random
1233 fabric materials. Across diverse garment categories, our
1234 method produces more realistic silhouettes, time-consistent
1235 wrinkle patterns, and garment dynamics that closely match
1236 those observed in the source video clips. Please refer to our
1237 project website for the full video results.

1238 10.3. Typical Failure Modes

1239 We observe several characteristic failure cases:

- **Ambiguous material appearance:** for visually ambigu-
1240 ous fabrics (e.g., cotton vs. viscose plain-weave), the
1241 model occasionally mispredicts the dominant fiber, lead-
1242 ing to slightly over- or under-damped dynamics.
- **Highly layered garments:** garments with multiple visi-
1243 ble layers (e.g., jackets worn over dresses) are filtered out
1244 in FTAG, but in-the-wild images may still contain layer-
1245 ing; in such cases, our single-layer physics may not fully
1246 capture the observed dynamics.
- **Extreme accessories and trims:** heavy buttons, zippers,
1247 or appliqués are not explicitly modeled and can lead to
1248 local discrepancies between simulated and real wrinkles
1249 near those regions.

1250 These cases suggest interesting future directions such as
1251 explicit layer reasoning, accessory-aware physics modeling,
1252 and leveraging multi-view or temporal cues when available.

Table 4. **Quantitative comparison of image-to-garment prediction.** We show the results we get for Geometry (CD, IoU), image-space-reconstruction (PSNR, SSIM, LPIPS), and physics parameter accuracy (NMAE) with respect to the ground truth of our created data. Best results in **bold**. Arrows indicate optimization direction.

| Sequence | Method | #Frames | CD↓ | IoU↑ | PSNR↑ | SSIM↑ | LPIPS↓ |
|---------------|--------------|---------|--------------|-------------|--------------|--------------|--------------|
| Jumping Jack | DMap* | 133 | 2030 | 1.87 | 16.86 | 0.910 | 0.129 |
| | AIparrel* | | 98.1 | 14.6 | 20.62 | 0.942 | 0.064 |
| | ChatGarment* | | 88.9 | 20.0 | <u>21.66</u> | <u>0.951</u> | <u>0.056</u> |
| | Ours | | 64.4 | 21.6 | 22.10 | 0.954 | 0.053 |
| Joyful Jump | DMap* | 91 | — | — | — | — | — |
| | AIparrel* | | 46.8 | 17.4 | 24.64 | <u>0.965</u> | 0.034 |
| | ChatGarment* | | <u>7.73</u> | <u>37.5</u> | <u>27.88</u> | 0.970 | 0.021 |
| | Ours | | 7.49 | 38.6 | 28.05 | 0.970 | 0.021 |
| Northern Spin | DMap* | 125 | — | — | — | — | — |
| | AIparrel* | | <u>257</u> | 14.0 | 17.86 | 0.881 | 0.156 |
| | ChatGarment* | | 525 | <u>13.9</u> | <u>18.32</u> | <u>0.879</u> | <u>0.150</u> |
| | Ours | | 163 | 12.4 | 18.91 | 0.866 | 0.145 |
| Average | DMap* | | 2030 | 1.87 | 16.86 | 0.910 | 0.129 |
| | AIparrel* | | <u>134.0</u> | 15.3 | 21.04 | 0.929 | 0.085 |
| | ChatGarment* | | 207.2 | <u>23.8</u> | <u>22.62</u> | 0.933 | <u>0.076</u> |
| | Ours | | 78.3 | 24.2 | 23.02 | 0.930 | 0.073 |

Table 5. **Ablation of the Physics Parameter Prediction.** Geometry (CD, IoU), image reconstruction (PSNR, SSIM, LPIPS), and physics parameter accuracy (NMAE) with respect to ground truth. Our method estimates physical parameters through a feedforward method, while Random Parameters are sampled within the bounds of the simulator. Best results in **bold**. Arrows indicate optimization direction.

| Sequence | Method | #Frames | CD↓ | IoU↑ | PSNR↑ | SSIM↑ | LPIPS↓ |
|---------------|-------------------|---------|-------------|-------------|--------------|--------------|--------------|
| Jab Cross | Random Parameters | 105 | 5.42 | 46.1 | 29.04 | 0.968 | 0.020 |
| | Ours | | 2.50 | 58.4 | 30.99 | 0.974 | 0.014 |
| Jumping Jack | Random Parameters | 133 | 79.4 | 22.2 | 22.14 | 0.951 | 0.053 |
| | Ours | | 31.2 | 28.3 | 24.63 | 0.960 | 0.035 |
| Northern Spin | Random Parameters | 150 | 68.6 | 25.1 | 23.92 | 0.932 | 0.063 |
| | Ours | | 1.35 | 76.4 | 33.60 | 0.970 | 0.009 |
| Average | Random Parameters | | 51.1 | 31.1 | 25.03 | 0.950 | 0.045 |
| | Ours | | 11.7 | 54.4 | 29.74 | 0.968 | 0.019 |

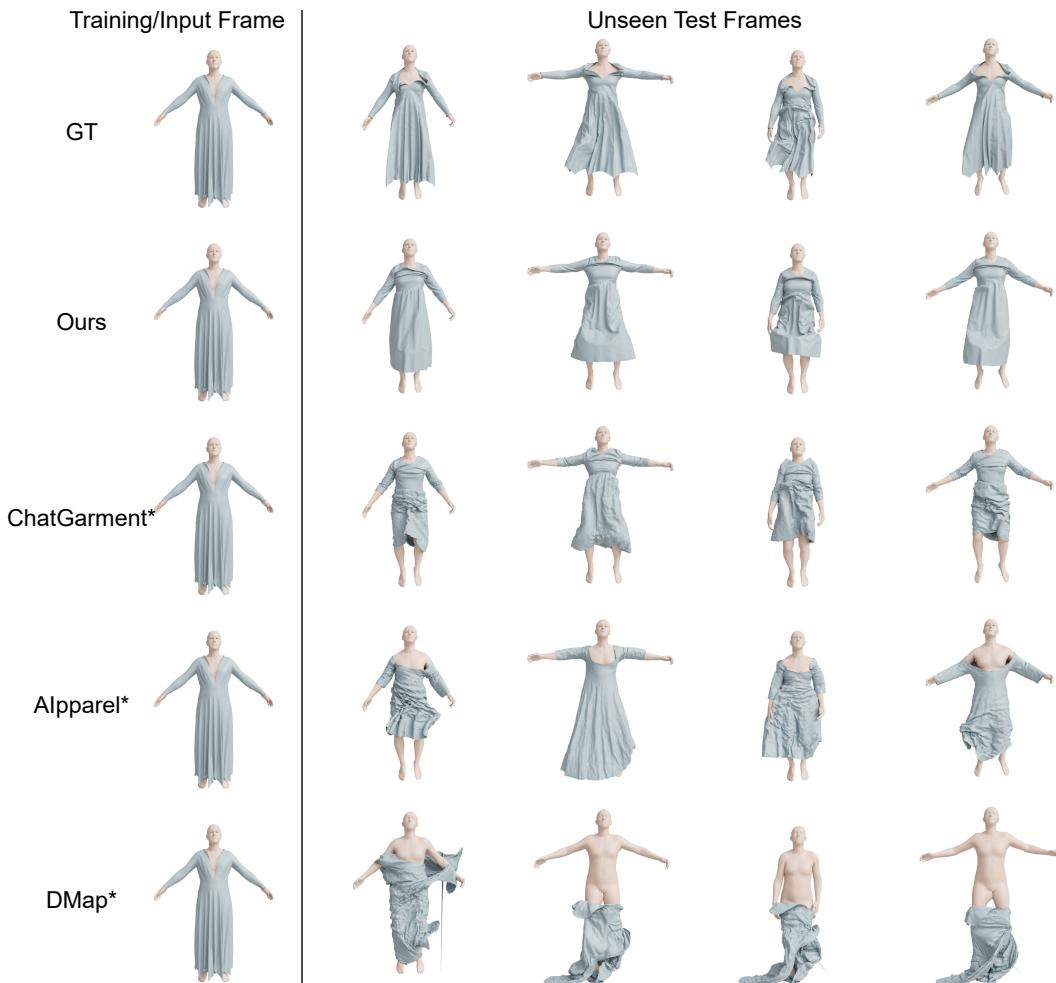


Figure 7. Qualitative Comparison on Jumping Jack. The left most column shows the input frames presented to each method and the columns on the right show the animated and rendered garments at different points in time of the animation. Our method's result most closely resemble the ground truth (top row).

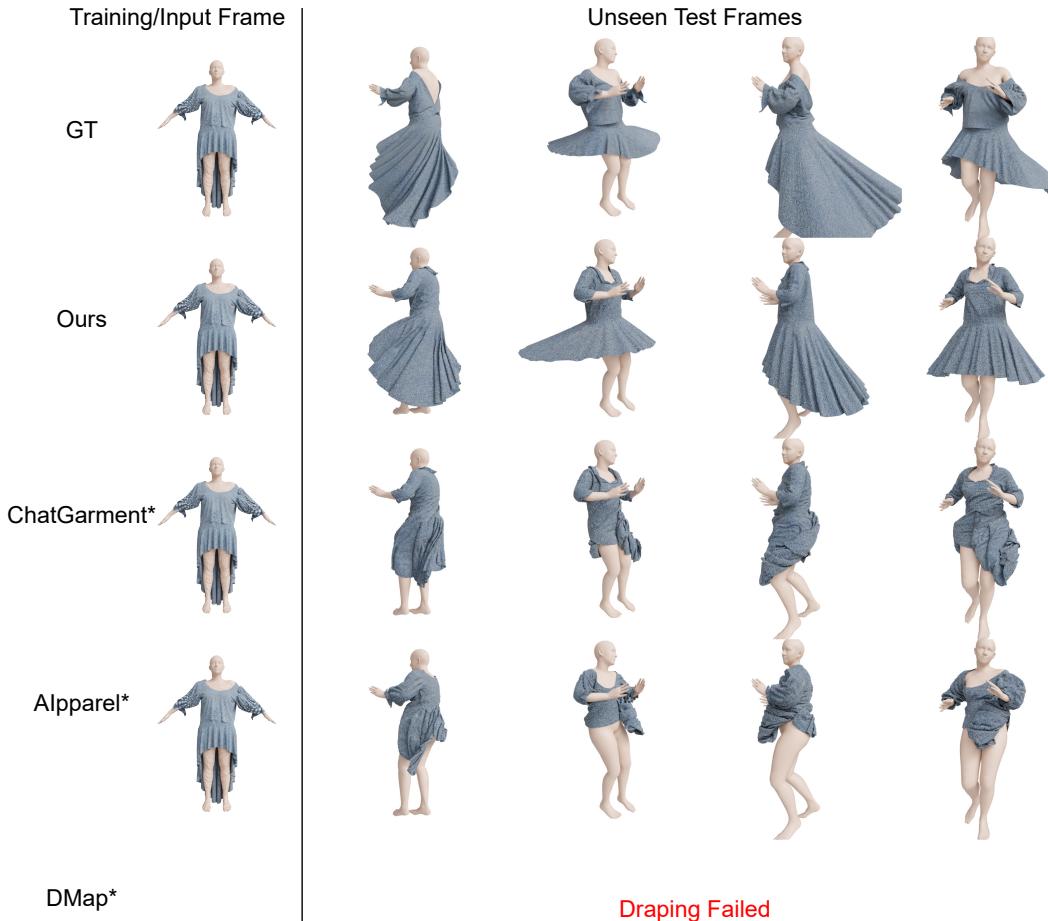


Figure 8. **Qualitative Comparison on Northern Spin.** The left most column shows the input frames presented to each method and the columns on the right show the animated and rendered garments at different points in time of the animation. Our method's result most closely resemble the ground truth (top row).

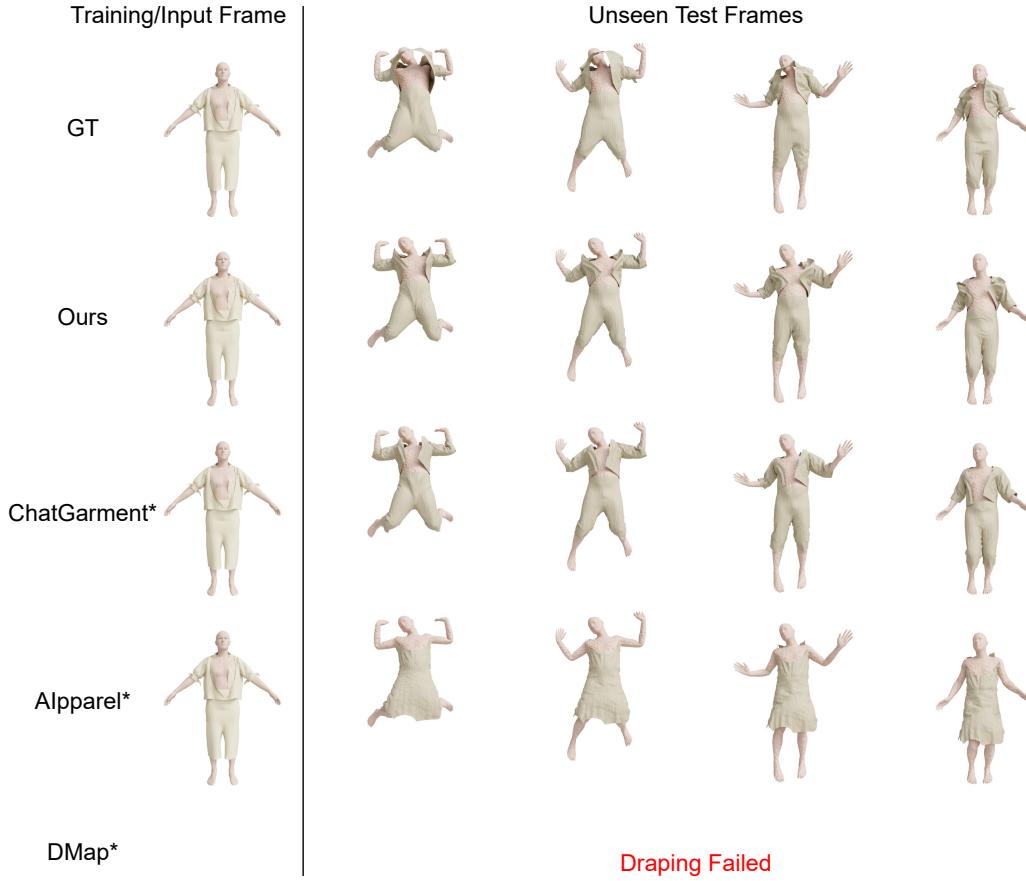


Figure 9. **Qualitative Comparison on Joyful Jumping.** The left most column shows the input frames presented to each method and the columns on the right show the animated and rendered garments at different points in time of the animation. Our method's result most closely resemble the ground truth (top row).



Figure 10. **Qualitative Comparison in-the-wild video.** The top row shows the original sequence. We only use the left-most frame as input. The rows below show renderings of the garments after simulation.



Figure 11. **Qualitative Comparison on in-the-wild video.** The top row shows the original sequence. We only use the left-most frame as input. The rows below show renderings of the garments after simulation.



Figure 12. **Qualitative Comparison on in-the-wild video.** The top row shows the original sequence. We only use the left-most frame as input. The rows below show renderings of the garments after simulation.