

# MagicClay: Sculpting Meshes With Generative Neural Fields

AMIR BARDA, Tel Aviv University, Israel  
 VLADIMIR G. KIM, Adobe Research, USA  
 NOAM AIGERMAN, Université de Montréal, Canada  
 AMIT H. BERMANO, Tel Aviv University, Israel  
 THIBAUT GROUEIX, Adobe Research, USA

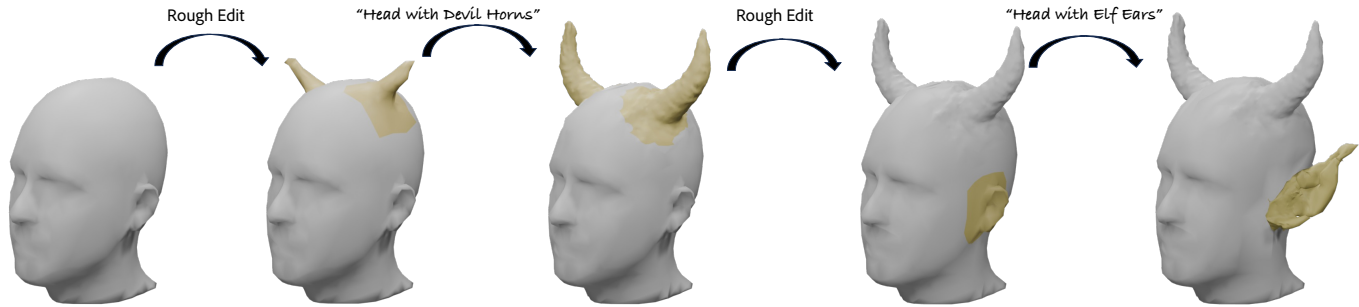


Fig. 1. MagicClay introduces a novel sculpting tool, employing a hybrid mesh-SDF representation. The user selects a region (middle) of an input mesh (left), inputs a text-prompt (bottom), and the region automatically grows to match the prompt (right), while the rest of the shape is unchanged, and the mesh remains topologically valid. This enables operations such as sequential semantic mesh editing.

The recent developments in neural fields have brought phenomenal capabilities to the field of shape generation, but they lack crucial properties, such as incremental control – a fundamental requirement for artistic work. Triangular meshes, on the other hand, are the representation of choice for most geometry related tasks, offering efficiency and intuitive control, but do not lend themselves to neural optimization. To support downstream tasks, previous art typically proposes a two-step approach, where first a shape is generated using neural fields, and then a mesh is extracted for further processing. Instead, in this paper we introduce a hybrid approach that maintains both a mesh and a Signed Distance Field (SDF) representations consistently. Using this representation, we introduce MagicClay – an artist friendly tool for sculpting regions of a mesh according to textual prompts while keeping other regions untouched. Our framework carefully and efficiently balances consistency between the representations and regularizations in every step of the shape optimization; Relying on the mesh representation, we show how to render the SDF at higher resolutions and faster. In addition, we employ recent work in differentiable mesh reconstruction to adaptively allocate triangles in the mesh where required, as indicated by the SDF. Using an implemented prototype, we demonstrate superior generated geometry compared to the state-of-the-art, and novel consistent control, allowing sequential prompt-based edits to the same mesh for the first time.

## 1 INTRODUCTION

The field of 3D shape generation has always been heavily dependent on the representations it uses for the shapes. Recent Neural Field-based representations (i.e., NeRFs [Mildenhall et al. 2021] or SDFs [Park et al. 2019; Chen and Zhang 2019]), have shown remarkable progress to the task [Poole et al. 2022; Wang et al. 2023a] in a very short time. These representations are robust to noisy losses,

and are naturally well-suited for neural frameworks, yielding impressive results and avoiding local minima. On the other hand, these representations are expensive to evaluate (limited by volumetric rendering resolutions), and lack acutely in control (such localized edits or even shape smoothing).

In contrast, triangular meshes are the dominant representation for most 3D applications in industry. This is because meshes are inexpensive, consistent, and intuitive. For instance, when evolving a mesh-based shape, an artist performs edits on geometry, textures, or even topology, and expects any additional updates to retain these attributes. This is not possible currently using Neural Field-based representations. Unfortunately, while the adaptive, or non-uniform, nature of meshes is perhaps their greatest advantage, it is also the reason they are not preferred by current generative frameworks. The sparse gradients induced by meshes tend to limit the ability of optimizations to achieve large deformations in a stable manner.

For this reason, many works turn to implicit functions for coarse generation as a first step, and to meshes in a second step, for the purpose of finer details, or downstream editability. However, as we demonstrate, two-stage pipelines are prone to local minima, and additionally cannot be extended to edit an existing mesh with pre-computed UVs for example.

In this paper, we present MagicClay - a shape evolution and editing framework, based on a hybrid implicit-explicit representation, benefiting from the best of both worlds. MagicClay optimizes a mesh and an SDF jointly in every step throughout the generation process, and introduces a novel sculpting tool for 3D generative modeling. Sculpting is a common approach for 3D modeling, common to commercial 3D modeling software [Blender 2024; ZBrush 2024; SubstanceModeler 2024]. While sculpting currently requires a lot of time and expertise, our new tool allows artists to select a region on

Authors' addresses: Amir Barda, amirbarda@mail.tau.ac.il, Tel Aviv University, Israel; Vladimir G. Kim, Adobe Research, USA; Noam Aigerman, Université de Montréal, Canada; Amit H. Bermano, Tel Aviv University, Israel; Thibault Groueix, Adobe Research, USA.

a mesh to be modified, provide a textual prompt, and hallucinate an updated region (See Figure 1). As we demonstrate, in addition to the contribution to control, our hybrid representation also benefits computational efforts and overall geometric quality, as various priors can be placed more intuitively on the two representations.

The key technical challenge of the hybrid approach is keeping the two representations synced efficiently. To achieve this, we differentially render both representations from various angles, and require consistency in RGB renders, opacity and normal maps. Furthermore, we rely on the in-sync-mesh representation to render the SDF at higher resolutions and faster; Instead of the hundreds of samples per rays, we localize the SDF sampling around the mesh surface, and use as little as three samples. Critically, evolving a mesh consistently and stably is an additional challenge. In terms of resolution, a coarse mesh would not be expressive enough for novel details, and a fine mesh is expensive and unstable. Hence, an adaptive tessellation is required, that evolves along with the shape where required. We rely on recent developments in differentiable mesh reconstruction [Barda et al. 2023] to achieve dynamic mesh topology updates, including face splitting, edge collapse, and edge flips. To texture the mesh despite changes in its topology, we contribute a new strategy based on triangle supersampling. Importantly, using this layer, we can maintain mesh properties throughout the optimization.

As we demonstrate, our hybrid approach allows localized and sequential mesh editing operations, preserving mesh topology and information on one hand, while allowing radical and semantic evolution on the other. In addition, we show overall higher generated geometric quality, thanks to the priors the two representations impose on each other. The hybrid approach and sculpting tool demonstrate how the merits of both leading representations can be combined. In essence, our framework brings the recent and future breakthroughs in neural shape generation closer to artistic workflows, where work is in incremental steps, giving the artist precision and control over the end result, but with unprecedented expressiveness.

## 2 RELATED WORK

**3D generative models.** In their seminal work, DreamFusion, pool et al. [Poole et al. 2022] show that Text-to-Image diffusion models can be used to provide gradients to optimize a Neural Radiance Field (NeRF) via Score Distillation Sampling (SDS). Magic3D [Lin et al. 2023] achieves better quality by using a two stage approach: the first stage is similar to DreamFusion, and they note that the quality of the generated object is limited by the high cost of performing volumetric rendering for high resolutions images. The second stage uses a differentiable mesh representation to further refine the generated object, as differentially rendering meshes in high resolution is significantly cheaper in both time and memory. Magic123 [Qian et al. 2023] further improves upon Magic3d by using both 3d and 2d diffusion priors. ProlificDreamer [Wang et al. 2023a] proposes an improvement over SDS, the VSD loss, to drive 3D generation from 2d diffusion priors. Fantasia3d [Chen et al. 2023] and TextMesh [Tsalicoglou et al. 2023] decouple the appearance from the geometry by replacing the NeRF with an SDF, and optimizing a color network separately. TextDeformer [Gao et al. 2023] uses CLIP

as prior together with a novel gradient smoothing technique based on mesh Jacobians to deform meshes according to a text prompt.

The choice of using two stages in Magic3D [Lin et al. 2023] highlights the tradeoffs involved in choosing the right 3D representation for 3D generative models. While implicit functions are well suited for coarse generation because they allow topology updates, meshes can be rasterized very efficiently at a high resolution to get fine details in a second step. However, two-stage pipelines are prone to local minima and crucially, cannot be extended to edit an existing mesh with pre-computed UVs. In contrast, we jointly optimise a hybrid SDF and mesh representation, that can be initialized from an existing mesh and maintain all of its properties during optimization, benefiting from the best of both worlds in a *1-step pipeline*: topology updates from the SDF part and fine details from the mesh part.

*Focus on local editing:* Vox-E [Sella et al. 2023] edit a voxel grid via SDS and use the attention layers to encourage localized edits. However, they can not guarantee where the edit will happen because the localization mechanism is based on soft attention.

*Concurrent work:* DreamCraft3D [Sun et al. 2023] builds on the recent SDS works by fine-tuning the diffusion model during the generative process using DreamBooth [Ruiz et al. 2022]. Of note, Instant3D [Li et al. 2023] generates a 3D shape in a single forward pass, without require any costly optimization. While it does not allow for local artistic controls similar to the sculpting application we present, we are excited to leverage ideas from this research direction to accelerate our method in the future.

**Hybrid Representations.** There is no ubiquitous representation in 3D, as there exist in 2D for images, thus several representations exist and have been combined for diverse 3D tasks. The plethora of representations and combination show that there is no one-fit-for-all solution. In this work, we introduce a hybrid representation specialized for generative modeling and focus the related work on hybrids most relevant to this paper. Poursaeed et al. [Poursaeed et al. 2020] uses a coupling of implicit and explicit surface representation for generative 3D modeling, kept in sync by 3D losses. NerfMeshing [Rakotosaona et al. 2023] proposes a improved meshing pipeline for NeRFs. In contrast to [Poursaeed et al. 2020], the coupling is not achieved via coupled regularization losses, but explicitly enforced by projection layers from the SDF to the mesh. Finally, DmTeT [Shen et al. 2021] proposes deep marching Tetrahedra as a hybrid representation for high-resolution 3D Shape synthesis, notably used in the concurrent work Magic3D [Lin et al. 2023]. Our method uses both a set of regularization losses, as well as a dynamic projection layer based on ROAR [Barda et al. 2023] to keep the SDF and mesh part in sync.

**Traditional approaches for sculpting meshes.** Many commercial tools employ digital sculpting metaphor for 3D modeling, such as Zbrush [ZBrush 2024], Mudbox [Autodesk 2024], or Substance-Modeler [SubstanceModeler 2024]. Motivated by these workflows, geometry processing research focused on improving interactive techniques such as mesh deformation [Jacobson et al. 2014], mesh blending for cut-and-paste [Biermann et al. 2002], local parameterization for adding surface details [Schmidt et al. 2006], symmetry-guided autocompletion [Peng et al. 2018], and version control for collaborative editing [Salvati et al. 2015]. Despite these advances, 3D

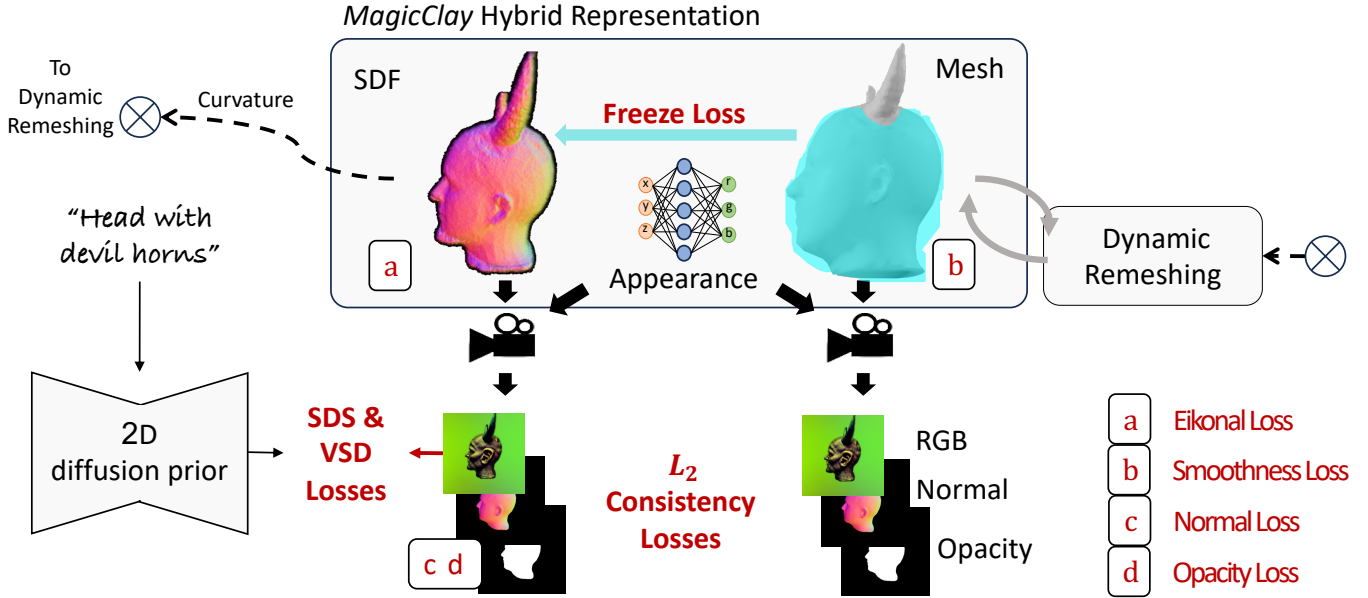


Fig. 2. Overview of the hybrid optimization. We jointly optimize a mesh, an SDF and a shared appearance MLP according to an input prompt. We can either optimize the full geometry, or only a user-selected portion of the mesh for an iterative 3D modeling workflow. We start by differentially rendering both representations, and enforce their consistency. As they are kept in sync, we use the mesh to efficiently sample volumetric rays to render hi-res maps from the SDF in a memory-efficient manner. In addition to the standard SDS loss, the high-res renderings allows using high quality VSD losses [Wang et al. 2023a] to evolve the SDF. We keep the mesh surface and the SDF in sync via multi-view consistency constraints on the RGB pixels, the image opacity and the surface normals. The mesh local topology is updated according to the SDF using ROAR [Barda et al. 2023], splitting triangles where geometry is created, and collapsing edges where needed. Additionally, we leverage representation-specific losses to regularize the optimization: an Eikonal loss on the SDF and a smoothness loss on the mesh.

modeling remains to be only accessible to experts. As an alternative, example-based approaches have been proposed to democratize 3D modeling tools by using existing geometry from a database of stock 3D models to assemble new shapes from parts [Funkhouser et al. 2004]. Subsequent methods have built statistical models over part assemblies [Kalogerakis et al. 2012], and allow high-level semantic control for deformations [Yumer et al. 2015]. Despite their accessibility, these tools are often restricted in their domain, and often rely on heavy annotation of stock assets, and thus have received limited use by professional modelers. In this paper, we utilize pretrained 2D generative data prior to enable semantic controls for local and iterative modeling workflow without the need of preannotated 3D stock data.

### 3 METHOD

Given a mesh, a user-highlighted surface region, and a text prompt that describes the desired target, MagicClay optimizes the shape of the selected region so that the resulting mesh matches the target. To drive the shape optimization, we follow current literature and use the Score Distillation Sampling (SDS) technique [Poole et al. 2022] with differentiable rendering to leverage on text-conditioned 2D diffusion and guide the shape optimization. This approach, however, does not perform well when operated on meshes. Meshes are driven by sparse and irregular samples (vertices), and their connectivity mandates a stable and smooth deformation, avoiding

self-intersections and flip-overs. For this reason, we employ a neural Signed Distance Field (SDF) to drive the mesh shape optimization and topology updates. We thus propose a hybrid representation that captures both a Signed Distance Field (SDF) and the surface, gaining from the advantages of both worlds. While the SDF allows guiding the shape towards larger-scale complex changes, the mesh allows to capture fine details, and localize control to the user-highlighted surface region.

In this section we provide details on the hybrid SDF/Mesh representations (Sec. 3.1), how it can be efficiently optimized with SDS guidance (Sec. 3.2), how to effectively use surface and volumetric priors (Sec. 3.3), and how to update the mesh topology during optimization (Sec. 3.4). Figure 2 overviews the full pipeline.

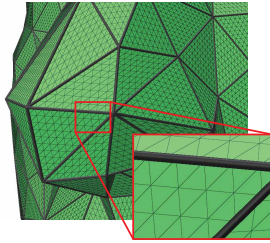
#### 3.1 Hybrid Representation

Our hybrid representation consists of a surface (a mesh), a volume (an SDF), and a shared appearance network encoding RGB colors for an input 3D coordinate. Both the surface and the volumetric representations can be differentially rendered, leveraging the shared appearance network to output images with color, normals, and opacity channels. We now detail the three elements of our hybrid shape representation.

*Surface Representation.* We represent the surface of the shape as a 2-manifold triangular mesh. Mesh topology, or sampling resolution, is locally adapted according to the SDF (see Sec. 3.4 for details). We

encode colors for the mesh using an auxiliary appearance network, derived from the SDF itself (see below). We found this approach simpler and more natural than traditional mesh coloring techniques; Using per-vertex colors is sensitive to triangulation, and would require a large number of vertices to match the resolution of the SDF. Using a texture image requires a complex UV parameterization, usually done a priori on a fixed shape. In addition, our surface is continuously optimized and undergoes through topological changes, re-tessellation, and large-scale deformations, which makes it computationally infeasible to apply traditional UV parameterization techniques during this optimization.

Instead, our hybrid approach offers a simpler approach to shape coloring. To apply the colors from the appearance network to the mesh, we propose to adaptively subdivide each face of the base mesh according to triangle area. Since we only use these subdivided triangles to represent colors they do not have to form a connected mesh unlike traditional subdivision techniques. Thus, we employ MeshColors scheme that was originally proposed for UV-less texturing [Yuksel et al. 2010], and has an efficient GPU implementation. In the inset we illustrate the example subdivision, note how sub-triangles on two adjacent faces do not share the vertices along the edge. During rendering we assign a color to each sub-triangle, by using the coordinate of the three associated supersampled vertices to sample the appearance network.



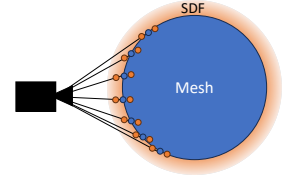
*Signed Distance Functions.* Our volumetric shape representation is chosen off-the-shelf, and conceptually serves as a regularization guiding the mesh evolution using existing state-of-the-art text2shape tools. We use a neural SDF, or an implicit continuous scalar field that can be sampled anywhere in  $\mathbb{R}^3$ , returning a signed shortest distance to the surface (negative on the inside, positive on the outside). We encode the SDF using a multiresolution hash encoding of features defined over a grid which are then mapped to distance value by a small MLP, following instant-NGP [Müller et al. 2022]. As in the mesh case, the shared appearance network is sampled to obtain colors during rendering.

*Appearance Network.* The shared appearance network encodes colors implicitly as a map over  $\mathbb{R}^3$ . It shares the same hash grid as the SDF, but has a smaller MLP head, with a single hidden layer that take hash grid features as input and outputs RGB values.

### 3.2 Hybrid Shape Guidance

In essence, our shape optimization is based on Score-Distillation Sampling (SDS) to distill gradients from a text prompt. The primary motivation to maintain an SDF representation in addition to the mesh is because SDFs are more robust noisy guidance, which is an inherent property of the multi-view SDS approach (see Figure 5). We thus choose to inject the text guidance only to the auxiliary SDF representation, and propagate the changes to the mesh via the consistency losses (Sec. 3.3) and the topology updates (Sec. 3.4).

To apply the text guidance and the consistency losses, we need to render both representation differentially. We use Nvdiffrast [Laine et al. 2020] to render meshes and VolSDF [Yariv et al. 2021] for volumetric rendering of our SDF. Clearly, as mesh rasterization is much cheaper than volumetric rendering, the process is bottlenecked by the resolution at which we can render the SDF, both in terms of speed and memory. Our hybrid representation uniquely enables a strategy to render SDF faster and cheaper, at a higher resolution of 512x512. This is achieved thanks to the consistency between the mesh and SDF representations throughout the optimization. We can significantly reduce the typical 512 samples per ray necessary for rendering the SDF by using the intersection of the ray with the mesh representation (efficiently calculated by the differentiable mesh renderer). Using the intersection as the center of a small spread of samples (typically 3), this allows for high resolution renders of the SDF (i.e. 512x512 and larger), which are otherwise memory prohibitive. The idea of leveraging the surface to reduce the number of network queries per ray emerged in concurrent works, namely Adaptive Shell [Wang et al. 2023b] and HybridNerf [Turki et al. 2023], which shows its generality and success in other settings than ours.



Using this strategy, we render the SDF in 512x512 and apply the VSD loss of those high-res renderings. We also apply regular SDS on lower-res 128x128 renderings by regular VolSDF as we find that this improves the results slightly.

### 3.3 Representation Priors

In addition to the text guidance, we apply representation-specific regularizations and consistency losses that keep both representations in sync.

*Consistency Loss.* The SDF and the mesh are consistent if their images are in 1 to 1 correspondences from any camera angle. We thus supervise the L2 difference between their RGB renderings, normal maps and opacity maps. If the renderings are made at different resolution, we downsize to the lower resolution before computing the L2 loss.

*Enforcing Localization and Freeze Loss.* To localize changes to the user-selected area we first fix the mesh vertices in all non-selected regions during optimization by zeroing out gradients outside of user selection. While localization is harder to achieve for SDF, we add a sampling-based freeze loss, which favors regions around fixed vertices to remain unchanged:

$$s(v_{\text{sampled}}) = 0 \quad (1)$$

where  $v_{\text{sampled}}$  are vertices sampled uniformly over the faces which are not part of the optimization region selected by the user.

*Laplacian (Smoothness) Loss.* While it is harder to regularize the surface of an implicit function to be smooth, the explicit representation of the mesh allows to easily define a smoothness term using the Laplacian of the mesh, defined for each vertex:



$$\delta(x_i) = x_i - \frac{\sum_j x_j}{N}, \quad (2)$$

where  $x_j$  are neighbors of  $x_i$ . The Laplacian vector encodes local geometry changes, with a smooth mesh is defined by low Laplacian vectors, and we use a global smoothness loss:

$$L_{\text{smooth}} = \sum_i \|\delta_i\|. \quad (3)$$

*SDF Eikonal Loss.* To encourage the implicit function to learn a valid SDF representation we use the Eikonal term as a loss. The SDF  $s$  is valid if and only if the loss in Eqn 4 is 0:

$$L_{\text{Eik}} = \sum_x (|\|\nabla s(x)\| - 1|)^2 \quad (4)$$

*SDF opacity and normal Loss.* Inspired by TextMesh [Tsalicoglou et al. 2023], we also binarize the SDF opacity and apply a Binary Cross Entropy loss to encourage discrete 0 or 1 values. To penalize badly oriented normals of the implicit surface, we apply an L2 penalty to the dot product between the normal and the camera direction if it is negative.

### 3.4 Updating the mesh topology

To maintain consistency between the mesh and SDF, it is necessary to perform local topology edits on the mesh in that increase or decrease mesh resolution where required. Continuous Remeshing [Palfinger 2022] pioneered such a local topology update approach, by using the Adam optimizer state as a signal. While this approach works well in a multi-view reconstruction scenario, where the images are sharp, and the camera parameters known, the noise involved in SDS makes the gradients, and by extension the Adam state, very noisy and an unstable signal to trigger those operations. We turn to another work, ROAR [Barda et al. 2023], particularly well-tailored to our hybrid representation. Within this framework, we use the SDF as the signal to trigger mesh triangle splits.

In a nutshell, for each triangle on the mesh, ROAR starts by supersampling the triangle into  $K$  sub-faces, and projects each sub-vertices on the 0-level set of the SDF  $s$  using a projection operator:

$$P(x) = -s(x) \cdot \nabla s(x) \quad (5)$$

This projection results in a piece-wise linear surface that approximate the implicit surface closest to the initial triangle. The decision to split this triangle is based on the curvature score of this piece of projected surface. If the surface is very curved, then the triangle is split using  $\sqrt{3}$ -subdivision [Kobbelt 2000]. Similarly, each edge is assigned a score based on the quadratic distance of its vertices to all the planes in the 1-ring of the edge, which intuitively represents how important is the edge to the geometry. If the score is low, then the edge can be collapsed with Qslim [Garland and Heckbert 2023].

We refer the interested reader to the ROAR paper [Barda et al. 2023] for more details, but the important point to note that ROAR offers a principled way to perform edge collapses and face splits in the sense that each iteration of ROAR strictly decrease an energy - the difference between the highest face score and the lowest edge score. It thus exhibits a convergence behavior after enough iterations. We also note that manifoldness is guaranteed to be preserved throughout the iterations.

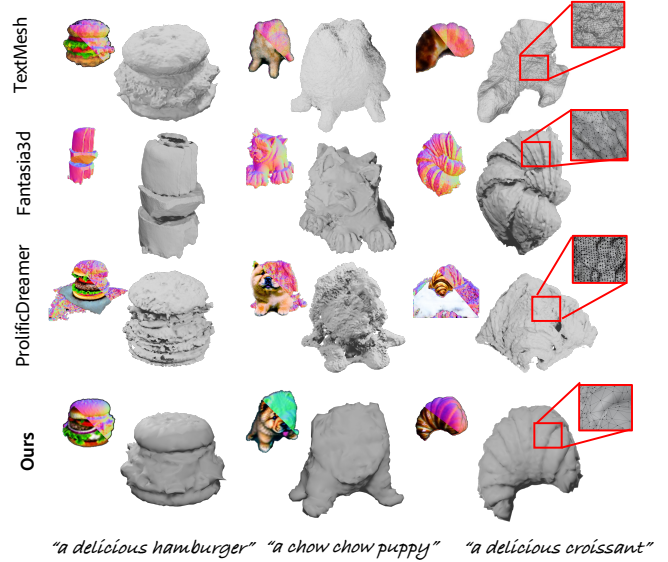


Fig. 3. **Comparison on text-to-3D from scratch.** We compare the quality of the triangular meshes extracted from various state-of-the-art generative methods: Fantasia3d [Chen et al. 2023], ProlificDreamer [Wang et al. 2023a] and TextMesh [Tsalicoglou et al. 2023]. While all methods produce realistic RGB renderings, only our hybrid representation generates smooth geometry.

## 4 EXPERIMENTS

We implement our pipeline in Threestudio [Guo et al. 2023], and use the implementations of other methods provided in the framework for comparisons with Stable Diffusion v1.5 as the backbone diffusion model. All experiments presented in this work were executed on a single A100-40GB GPU.

In the rest of this section we compare our representation to prior work on text-conditioned 3D generation (Sec. 4.1), demonstrate its utility in a mesh sculpting application (Sec. 4.2), and compare to a text-driven mesh deformation baseline (Sec. 4.3). We then provide a simple illustrative experiment to motivate the hybrid representation when using SDS guidance with noisy gradients (Sec 4.4), and finally ablate our method (Sec 4.5).

### 4.1 Comparison with Generative Methods

Since MagicClay is a modeling tool, we are primarily interested in evaluating the quality of the geometry and thus focus on mesh renderings without texture. Note that existing 3D generative techniques are not designed to edit a part of an existing mesh, and thus we compare performance of our hybrid approach on the task of text-to-3D generation. We compare against three recent approaches: Fantasia3d [Chen et al. 2023], ProlificDreamer [Wang et al. 2023a] and TextMesh [Tsalicoglou et al. 2023]. The representation used in TextMesh [Tsalicoglou et al. 2023] is an SDF, so we run marching cube on their output to evaluate the quality of the mesh. ProlificDreamer uses a two stage approach: first a NeRF-base generation, then a refinement using an explicit representation. Fantasia3d uses DMtet [Shen et al. 2021] to extract a mesh from an SDF.

We present our results in Figure 3. Even though each approach generates a representation that can produce realistic RGB renders,

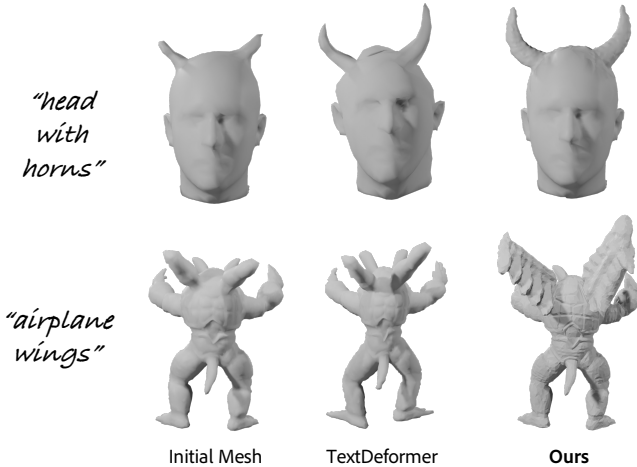


Fig. 4. **Comparisons to TextDeformer** [Gao et al. 2023]. Given the user-modified input mesh, we either run TextDeformer or our method to modify it towards the target prompt.

extracted meshes often exhibit significant surface artifacts, which make them hardly recognizable without texture (see “Chow Chow puppy” by ProlificDreamer or “Croissant” by TextMesh). By comparison, our geometries are recognizable and smooth thanks to the fact that our hybrid representation enables an explicit regularization of the surface. This validates that MagicClay successfully bridges the generative capabilities of implicit radiance fields with the surface-level controls of meshes.

## 4.2 Mesh Sculpting

We further demonstrate that our method can be used to enable novel iterative 3D sculpting workflows. Starting with an initial mesh, an artist can select region of interest (and optionally sculpt a coarse adjustment to that region) along with a textual prompt describing the desired updated shape. MagicClay generates a modified mesh, which could be iteratively refined with new elements. Note that hybrid representation is essential to this application. First, selecting and adjusting the region of interest is easily accomplished using standard mesh editing tools [Blender 2024]. Second, the generated result tends to be more expressive with additional SDF representation guidance. Third, using the mesh allows us to keep non-selected surface regions intact by zeroing out their deformation gradients, which guarantees that the change will only affect the user-selected region. Refer to Figures 1, 9 and supplemental video for example results. MagicClay generates a high quality edits, that match the rough local edit and adhere to the user’s text prompt.

## 4.3 Comparison with TextDeformer

Our method is the first to tackle the problem of interactive, localized mesh sculpting via text-based prompts. A naive alternative would be to simply use the existing text-driven mesh deformation technique [Gao et al. 2023] on the user-modified mesh, still aiming to achieve desired changes in the geometry. In Figure 4 we compare our method to this baseline. Note how our method is able to add geometrically complex large-scale details due to guidance from SDF and topological updates. Our method’s changes are also restricted only

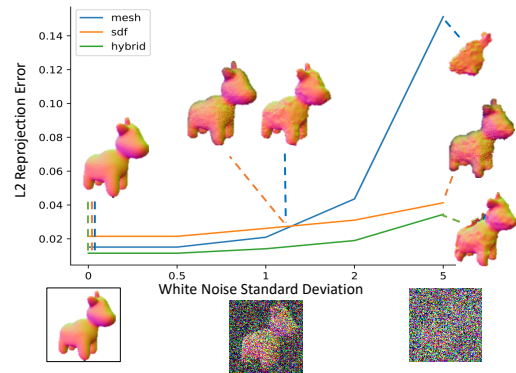


Fig. 5. **Mesh and SDF robustness to noisy gradients.** We optimize a mesh, an SDF or our hybrid representation with multi-view reconstruction losses after applying various noise levels to the ground truth renderings. We report the L2 reprojection error against the ground truth renders. The SDF exhibit more robustness than the mesh to high noise regime, and our hybrid outperforms both.

to user-modified region, and do not lead to large-scale deformations in the other parts of the input.

## 4.4 Analysis of Mesh and SDF robustness to noise

We now illustrate the motivation for our hybrid representation by a simple controlled experiment, where we aim to reconstruct a fixed 3D target with different levels of noise in the guidance. We formulate it as a reconstruction problem to have a clear ground truth and eliminate ambiguity arising due to text-based objectives. Even though we use synthetic noise, we expect these findings to apply in an SDS setting, where gradients are also noisy due to random noising step performed at each SDS iteration [Poole et al. 2022].

Given multi-view renderings of a fixed (true) 3D model, we add uncorrelated per-pixel Gaussian noise to each image, and compute L2 pixel-wise loss to guide our shape representation towards the target. As we increase the noise level (by increasing standard deviation) we find that different representations are more prone to errors in reconstructing the target. We use L2 re-projection error with respect to the ground truth shape as our evaluation metric, and compare vanilla Mesh, SDF representations to our hybrid approach (using our mesh rendering), and show results in Figure 5. Note that as the amount of noise increases, the quality of the mesh reconstruction degrades sharply. At the highest noise regimes (standard deviation of 5), the mesh reconstruction degenerates to a blob, while the SDF reconstruction is still recognizable despite surface irregularities. Importantly, the hybrid representation performs better than both individual representation at all levels of noise, and the benefits are the strongest at higher noise level.

## 4.5 Ablations

We perform several ablations to justify the design of our system.

*No face supersampling for mesh colors.* Figure 6 illustrates the need for the mesh-driven super-sampling of the appearance network based on Mesh Colors schema (Sec 3.1). We optimize our

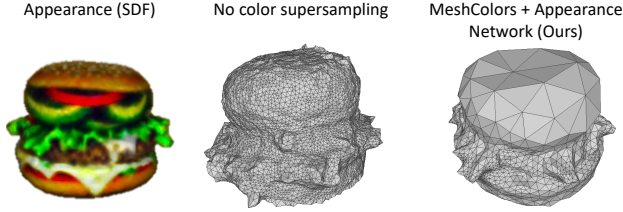


Fig. 6. **Ablation: no color supersampling.** Using mesh-guided supersampling in conjunction with the appearance network allows to decouple geometry and appearance. Using this approach (top right), large faces are used for the mesh, even though the rendering still presents high frequency colors (bottom). When using a color-per-vertex scheme (top middle), significantly larger mesh resolution is required to achieve similar appearance.

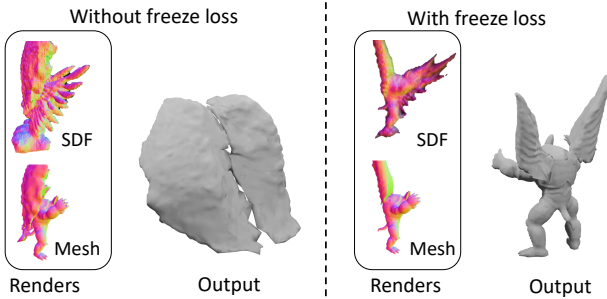


Fig. 7. **Ablation: not enforcing localization.** Without localization and freeze losses we observe that shape changes can propagate beyond the user-highlighted area, potentially destroying the original content. Here armadillo was erased by “angel wings.”

representation with respect to a target image (first column), either sampling colors per-vertex (second column), or using our adaptive sampling using MeshColors (third column). Note that without MeshColors sampling a much higher resolution is needed, which leads to poorer reconstruction and an over-tesellated mesh.

*Not enforcing localization, no freeze loss.* We remove the mechanism for enforcing localization via fixing non-selected surface regions and nearby SDF values as discussed in Sec. 3.3. In Figure 7 we show that without this feature, the shapes undergo unintended global changes, potentially erasing the original shape.

*No topology updates.* The topological updates (Sec. 3.4) during the optimization significantly improve the results as they allow to gradually add resolution. Optimizing a fixed-resolution mesh would either result in a shape that only marginally differs from input if the initial resolution is too high (Fig 8, left), or lacks fine details if the initial resolution is too coarse (Fig 8, right).

## 5 CONCLUSION

We presented MagicClay, a generative sculpting tool, backed by our new hybrid SDF and mesh representation. We demonstrated the importance of the hybrid representation through careful analysis and baselines. Key to the success of the generative process is our new rendering strategy that leverages the mesh part of the hybrid representation to localize ray sampled in the volumetric rendering of the SDF. We believe MagicClay is an important step towards

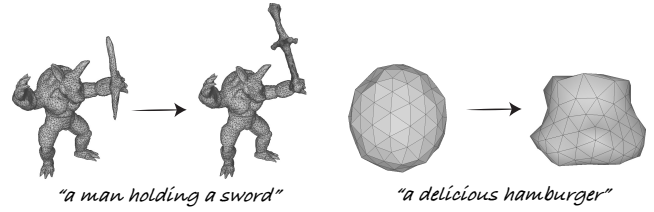


Fig. 8. **Ablation: no topology updates.** Optimizing the mesh without topology update results in the final generated object being limited by the initial resolution. **Left** When starting with a fine mesh the optimization will often get stuck since each vertex has tiny effect on the objective. **Right** When starting from a coarse mesh, no fine details can be created.

turning the recent advancements in text-to-image-from-scratch into an actual modeling tool usable by artists in an iterative workflow.

*Limitations.* Our method is inherently constrained by the quality of the SDS gradients. This stems from the inability of the current generative model to generate *consistent* multi-view images during score-distillation sampling. Each view tracts the optimization in a different direction which results in a noisy process, rendering the emergence of fine details more difficult. That said, MagicClay employs current text-to-shape methods without adaptation, and can easily integrate and benefit from the unavoidable further development of the field. Second, MagicClay is not interactive, as running MagicClay takes  $\bar{1}$  hr per prompt on an A100 GPU. We know from the NeRF literature that a couple of dozens of multi-view images are sufficient to reconstruct 3D shapes very efficiently. Hence, we believe that as diffusion models become better at generating consistent images, the number of iterations required by MagicClay will drop significantly, making it much faster.

*Future work.* We see several venues for future research. First, we see opportunities to leverage inpainting and depth-conditioned diffusion models. Indeed, this generative process transforms the full object in each rendering, whereas it is clear that some part of the generated image should stay the same as the 3D edit is localized. We think that leveraging this insight would reduce the amount of noise in the 3D optimization helping the optimization reach a better geometry and faster. Second, we would like to explore using image targets with our system. We believe this could help making edits more specific, and also could allow more control, where users can highlight which part of the image should affect the generated shape. Finally, embellishing the surface representation to capture high resolution geometric details (such as normal and bump maps) and connecting it to SDF representation with appropriate consistency losses, can further improve the quality of the resulting shapes.

## REFERENCES

- Autodesk. 2024. Mudbox. <https://www.autodesk.com/products/mudbox>.
- Amir Barda, Yotam Erel, Yoni Kasten, and Amit H. Bermano. 2023. ROAR: Robust Adaptive Reconstruction of Shapes Using Planar Projections. arXiv:2307.00690 [cs.GR]
- Henning Biermann, Ioana Martin, Fausto Bernardini, and Denis Zorin. 2002. Cut-and-Paste Editing of Multiresolution Surfaces.
- Blender. 2024. <http://www.blender.org>.
- Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. 2023. Fantasia3D: Disentangling Geometry and Appearance for High-quality Text-to-3D Content Creation. arXiv:2303.13873 [cs.CV]

- Zhiqin Chen and Hao Zhang. 2019. Learning Implicit Fields for Generative Shape Modeling.
- Thomas Funkhouser, Michael Kazhdan, Philip Shilane, Patrick Min, William Kiefer, Ayellet Tal, Szymon Rusinkiewicz, and David Dobkin. 2004. Modeling by Example. *ACM Transactions on Graphics* (2004).
- William Gao, Noam Aigerman, Thibault Groueix, Vladimir G. Kim, and Rana Hanocka. 2023. TextDeformer: Geometry Manipulation using Text Guidance. arXiv:2304.13348 [cs.CV]
- Michael Garland and Paul S. Heckbert. 2023. Surface Simplification Using Quadratic Error Metrics. , 8 pages. <https://doi.org/10.1145/3596711.3596727>
- Yuan-Chen Guo, Ying-Tian Liu, Ruizhi Shao, Christian Laforte, Vikram Voleti, Guan Luo, Chia-Hao Chen, Zi-Xin Zou, Chen Wang, Yan-Pei Cao, and Song-Hai Zhang. 2023. threestudio: A unified framework for 3D content generation. <https://github.com/threestudio-project/threestudio>.
- Alec Jacobson, Zhigang Deng, Ladislav Kavan, and J.P. Lewis. 2014. Skinning: Real-time Shape Deformation.
- Evangelos Kalogerakis, Siddhartha Chaudhuri, Daphne Koller, and Vladlen Koltun. 2012. A Probabilistic Model of Component-Based Shape Synthesis. *ACM Transactions on Graphics* 31, 4 (2012).
- Leif Kobbelt. 2000. Sqrt(3)-Subdivision. *ACM SIGGRAPH 2000 (05 2000)*.
- Samuli Laine, Janne Hellsten, Tero Karras, Yeongho Seol, Jaakko Lehtinen, and Timo Aila. 2020. Modular Primitives for High-Performance Differentiable Rendering. *ACM Transactions on Graphics* 39, 6 (2020).
- Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan Sunkavalli, Greg Shakhnarovich, and Sai Bi. 2023. Instant3d: Fast text-to-3d with sparse-view generation and large reconstruction model.
- Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiao-hui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. 2023. Magic3D: High-Resolution Text-to-3D Content Creation. arXiv:2211.10440 [cs.CV]
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. , 99–106 pages.
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM Trans. Graph.* 41, 4, Article 102 (July 2022), 15 pages. <https://doi.org/10.1145/3528223.3530127>
- Werner Palfinger. 2022. Continuous remeshing for inverse rendering. *Computer Animation and Virtual Worlds* 33 (07 2022). <https://doi.org/10.1002/cav.2101>
- Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. 2019. DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation.
- Mengqi Peng, Jun Xing, and Li-Yi Wei. 2018. Autocomplete 3D Sculpting.
- Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. 2022. DreamFusion: Text-to-3D using 2D Diffusion. arXiv:2209.14988 [cs.CV]
- Omid Poursaeed, Matthew Fisher, Noam Aigerman, and Vladimir G. Kim. 2020. Coupling Explicit and Implicit Surface Representations for Generative 3D Modeling. In *Computer Vision – ECCV 2020*, Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (Eds.), Springer International Publishing, Cham, 667–683.
- Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skorokhodov, Peter Wonka, Sergey Tulyakov, and Bernard Ghanem. 2023. Magic123: One Image to High-Quality 3D Object Generation Using Both 2D and 3D Diffusion Priors. arXiv:2306.17843 [cs.CV]
- Marie-Julie Rakotosaona, Fabian Manhardt, Diego Martin Arroyo, Michael Niemeyer, Abhijit Kundu, and Federico Tombari. 2023. NeRFMeshing: Distilling Neural Radiance Fields into Geometrically-Accurate 3D Meshes. arXiv:2303.09431 [cs.CV]
- Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2022. DreamBooth: Fine Tuning Text-to-image Diffusion Models for Subject-Driven Generation.
- Gabriele Salvati, Christian Santoni, Valentina Tivaldo, and Fabio Pellacini. 2015. Mesh-Histo: Collaborative Modeling by Sharing and Retargeting Editing Histories. *ACM Trans. Graph.* (2015).
- R Schmidt, C Grimm, and B Wyvill. 2006. Interactive decal compositing with discrete exponential maps.
- Etai Sella, Gal Fiebelman, Peter Hedman, and Hadar Averbuch-Elor. 2023. Vox-E: Text-guided Voxel Editing of 3D Objects. arXiv:2303.12048 [cs.CV]
- Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. 2021. Deep Marching Tetrahedra: a Hybrid Representation for High-Resolution 3D Shape Synthesis.
- SubstanceModeler. 2024. Substance Modeler. <https://www.adobe.com/ie/products/substance3d-modeler.html>.
- Jingxiang Sun, Bo Zhang, Ruizhi Shao, Lizhen Wang, Wen Liu, Zhenda Xie, and Yebin Liu. 2023. DreamCraft3D: Hierarchical 3D Generation with Bootstrapped Diffusion Prior. arXiv:2310.16818 [cs.CV]
- Christina Tsalicoglou, Fabian Manhardt, Alessio Tonioni, Michael Niemeyer, and Federico Tombari. 2023. TextMesh: Generation of Realistic 3D Meshes From Text Prompts. arXiv:2304.12439 [cs.CV]
- Haitem Turki, Vasu Agrawal, Samuel Rota Bulò, Lorenzo Porzi, Peter Kotschieder, Deva Ramanan, Michael Zollhöfer, and Christian Richardt. 2023. HybridNeRF: Efficient Neural Rendering via Adaptive Volumetric Surfaces. arXiv:2312.03160 [cs.CV]
- Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. 2023a. ProlificDreamer: High-Fidelity and Diverse Text-to-3D Generation with Variational Score Distillation. arXiv:2305.16213 [cs.LG]
- Zian Wang, Tianchang Shen, Merlin Nimier-David, Nicholas Sharp, Jun Gao, Alexander Keller, Sanja Fidler, Thomas Müller, and Zan Gojcic. 2023b. Adaptive Shells for Efficient Neural Radiance Field Rendering. , 15 pages. <https://doi.org/10.1145/3618390>
- Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. 2021. Volume rendering of neural implicit surfaces.
- Cem Yuksel, John Keyser, and Donald House. 2010. Mesh Colors. *ACM Trans. Graph.* 29 (03 2010). <https://doi.org/10.1145/1731047.1731053>
- Ersin Yumer, Siddhartha Chaudhuri, Jessica Hodgins, and Levent Burak Kara. 2015. Semantic Shape Editing Using Deformation Handles.
- ZBrush. 2024. <https://www.maxon.net/en/zbrush>.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009



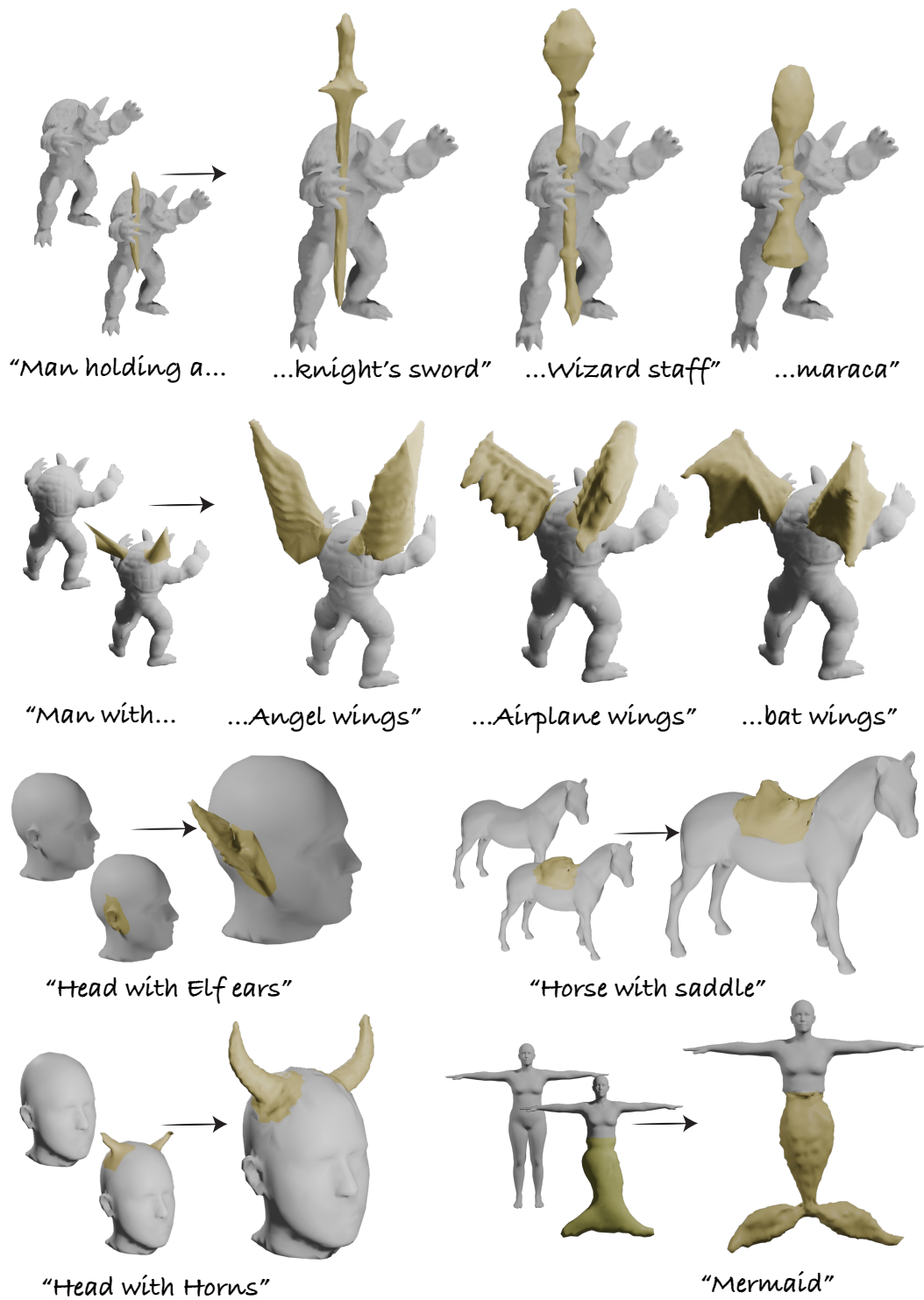


Fig. 9. **Sculpting gallery.** *Left:* from a source mesh, the user performs a rough edit in under two minutes, highlighted in yellow. *Right:* MagicClay refines it to match the provided prompt.