# Splat-Nav: Safe Real-Time Robot Navigation in Gaussian Splatting Maps

Timothy Chen[1⋆], Ola Shorinwa[1⋆], Joseph Bruno[3], Aiden Swann[1], Javier Yu[1],
Weijia Zeng[2], Keiko Nagami[1], Philip Dames[3], Mac Schwager[1]

*Abstract*—We present Splat-Nav, a real-time robot navigation pipeline for Gaussian Splatting (GSplat) scenes, a powerful new 3D scene representation. Splat-Nav consists of two components: 1) Splat-Plan, a safe planning module, and 2) Splat-Loc, a robust vision-based pose estimation module. Splat-Plan builds a safe-by-construction polytope corridor through the map based on mathematically rigorous collision constraints and then constructs a Bézier curve trajectory through this corridor. Splat-Loc provides real-time recursive state estimates given only an RGB feed from an on-board camera, leveraging the point-cloud representation inherent in GSplat scenes. Working together, these modules give robots the ability to recursively re-plan smooth and safe trajectories to goal locations. Goals can be specified with position coordinates, or with language commands by using a semantic GSplat. We demonstrate improved safety compared to point cloud-based methods in extensive simulation experiments. In a total of 126 hardware flights, we demonstrate equivalent safety and speed compared to motion capture and visual odometry, but without a manual frame alignment required by those methods. We show online re-planning at more than 2 Hz and pose estimation at about 25 Hz, an order of magnitude faster than Neural Radiance Field (NeRF)-based navigation methods, thereby enabling real-time navigation. We provide experiment videos on our project page at https://chengine.github.io/splatnav/. Our codebase and ROS nodes can be found at https://github.com/chengine/splatnav.

*Index Terms*—Vision-Based Navigation, Collision Avoidance, Localization.

## I. INTRODUCTION

Autonomous robotic operation requires robots to localize themselves within an envrionment, plan safe paths to reach a desired goal location, and have closed-loop trajectory-tracking. Traditionally, the fundamental problems of planning and localization have been performed in maps represented as occupancy grids [1], triangular meshes [2], point clouds [3], and Signed Distance Fields (SDFs) [4], all of which provide well-defined geometry.

However, these explicit scene representations are generally constructed at limited resolutions (to enable real-time opera-

tion), leaving out potentially-important scene details that could be valuable in planning and localization problems.

Neural Radiance Fields (NeRFs) [5] have recently been used to implicitly represent 3D scenes. NeRFs consist of a volumetric density field and a view-dependent color field parameterized by multilayer perceptrons (MLPs). NeRFs generate photo-realistic scene reconstructions, addressing the fundamental limitations of explicit representations; however, NeRFs require running inference on a deep neural network to render the scene, making them impractical for real-time use in robotic path planning. More recently, Gaussian Splatting (GSplat) [6] has emerged as a viable scene representation compared to NeRFs, representing the environment with Gaussian (ellipsoidal) primitives. Compared to NeRFs, GSplats generate higher-fidelity maps at faster rendering rates, with shorter or comparable training times. More importantly for robotics, GSplats, unlike NeRFs, offer a geometrically consistent collision geometry, enabling us to use level sets of these Gaussians to generate an ellipsoidal representation of the scene. These interpretable geometric primitives facilitate the development of rigorous motion planning algorithms that are safe, robust, and real-time.

In this paper, we introduce *Splat-Nav*, a pipeline for drone navigation in GSplat maps with a *monocular* camera. Splat-Nav comprises a lightweight pose estimation module, Splat-Loc, coupled with a planning module, Splat-Plan, to enable safe navigation from RGB-only (monocular) camera observations, as illustrated in Figure 1. Given an incoming RGB frame, Splat-Loc performs Perspective-n-Point (PnP)-based localization, leveraging the GSplat map to estimate the RGB and depth values rendered at candidate poses, which are then used to estimate the drone's pose. Next, Splat-Plan ingests the estimated pose computed by Splat-Loc to generate an initial trajectory, which is subsequently optimized to lie within safe flight corridors constructed from the ellipses that make up the GSplat map. The trajectory is parametrized by smooth, continuously safe Bézier splines that route the robot to a specified position or to a open-vocabulary language-conditioned goal location (i.e., "go to the microwave"). This feature enables the execution of Splat-Nav in a wide array of deployment conditions, such as in search missions where the precise location of targets is not known.

Additionally, the proposed system enables both open-loop trajectory generation and closed-loop re-planning. The latter is important in long trajectories, where existing onboard localization may drift or be subject to noise, impacting the overall safety of the executed trajectory of the robot. In these scenarios, Splat-Loc estimates can either be fused with that

[1] Stanford University, Stanford, CA 94305, USA {chengine, shorinwa, swann, javieryu, knagami, schwager}@stanford.edu

[2] University of California San Diego, San Diego, CA 92093, USA, wez195@ucsd.edu

[3] Temple University, Philadelphia, PA 19122, USA, {brunoj6, pdames}@temple.edu
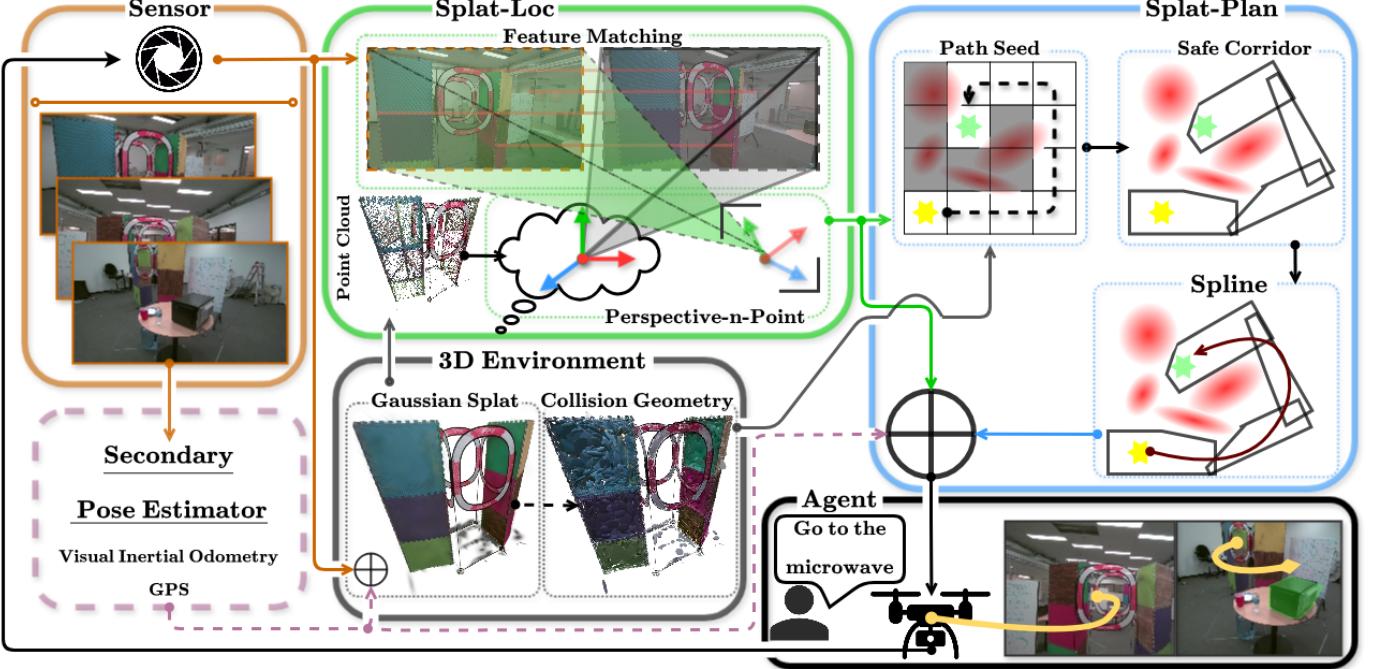
Fig. 1: Splat-Nav, consists of a safe planning module, Splat-Plan, and robust localization module, Splat-Loc, both operating on a Gaussian Splatting environment representation. In Splat-Plan we develop a fast, new ellipsoid-ellipsoid collision test to find a safe flight corridor through the GSplat, and plan a spline through the corridor. In Splat-Loc we localize the robot using only RGB images through a PnP algorithm, using the GSplat to render a point cloud. We use a language-embedded GSplat to enable open-vocabulary specification of goal locations like "go to the microwave."

of the existing localization module or used as a correction mechanism to steer the current motion toward a safer one. Finally, closed-loop re-planning additionally enables changes in goal locations during execution, leading to more dynamic plans.

In extensive simulations we compare Splat-Plan and Splat-Loc with baseline alternatives for planning and localization, respectively. We show Splat-Plan is always safe with respect to the full collision geometry, while four variants of a point-cloud based planner sometimes lead to collisions, or fail to find trajectories. Splat plan achieves similar or better solutions in terms of path length compared to point cloud-based planner in all cases, with similar computation time. Splat-Plan runs at no less than 2 Hz; comparable to point cloud-based solutions for the same scenes, but faster than gradient-based NeRF planners [7] and sampling-based planners (greater than 1 Hz) for similar solution quality. Similarly, we find that Splat-Loc is more accurate, faster, and fails less often compared to baselines. We demonstrate online pose estimation at about 25 Hz on a desktop computer, enabling real-time navigation.

Finally, in an experimental campaign with 124 hardware flights, we show that Splat-Nav (Splat-Plan and Splat-Loc running together) perform as well as motion capture or on board VIO, without the manual frame alignment required for those methods to align the MoCap or VIO frame with the GSplat (since both Splat-Plan and Splat-Loc operate natively in the same GSplat map).

The key contributions of this paper are as follows:

- We develop a fast polytope corridor generation algorithm to enable provably safe planning for drone navigation in GSplat maps.
- We develop a fast camera localization module based on GSplat maps that does not require manual alignment of the pose estimation frame to the planning frame, improving the synergy between planning and pose estimation.
- We demonstrate safe closed-loop re-planning with open-vocabulary goal specification, across a series of 124 hardware experimental trials.

## II. RELATED WORK

We review the related literature in robot planning and localization with different map representations, categorized into three groups: traditional representations (e.g., occupancy grids, meshes, point clouds, and SDFs), NeRFs, and GSplat.

**Planning.** We refer readers to [8] for an excellent explanation of planning algorithms in robotics. Most relevant for this work are the graph-based planners (e.g., A⋆), which compute a path over a grid representation of the environment; sampling-based planners (e.g., PRM [9], RRT and RRT⋆ [10]), which generate a path by sampling candidate states within the configuration space of the robot; and trajectory optimization-based planners (e.g., CHOMP [11] and Traj-Opt [12]), which take an optimization-based approach to planning. Prior work in [13] utilizes an optimization-based approach in path planning, taking a point cloud representation of the environment and converting this into a set of safe polytopes. The resulting safe polytopes are

utilized in computing a safe trajectory, parameterized as a spline from a quadratic program (QP), which can be efficiently solved.

There is also extensive literature on planning based on onboard sensing. Typically, these works present reactive control schemes [14, 15], using the sensed depth to perform collision checking in real time. These methods typically are myopic, reasoning only locally about the scene. Consequently, such methods often converge to local optima, preventing the robot from reaching its goal, especially in cluttered, complex environments. An alternate approach is to construct a map of the environment using depth measurements from Lidar or RGB-D sensors. Often, a Signed Distance Field (SDF) or its truncated variant (TSDF) is constructed from depth data [16, 17], which is encoded within a voxel representation. Such a representation is typical in dynamic robotic motion planning, providing fast collision checking and gradients in planning; however, voxel-based scene models do no provide visually rich or geometrically detailed scene representations compared to NeRFs or GSplats. Point cloud and voxel-based representations require a significant number of points or voxels for high-fidelity scene reconstruction, increasing the computational burden. Prior work [18, 19] introduced a Gaussian Mixture Model (GMM) as a more effective scene representation, which preserves the accuracy of point cloud-based map representations without the additional computational overhead, enabling fast exploration of unknown environments by multi-robot teams. Nevertheless, the aforementioned methods do not achieve photorealistic scene rendering.

More recent research has developed planning methods for highly expressive neural representations, such as NeRFs, which represent the environment as a spatial density field (with color) [5]. Using a NeRF map, NeRF-Nav [7] plans trajectories that minimize the total collision cost for differentially flat robots, e.g., quadrotors. Further, CATNIPS [20] converts the NeRF into a probabilistic voxel grid and then uses this to generate trajectories parameterized as Bézier curves. The work in [21] uses the predicted depth map at sampled poses to enforce step-wise safety using a control barrier function. The above works are complementary, with [7, 20] serving as high-level planners that encourage non-myopic behavior, while [21] can be used as a safety filter for a myopic low-level controller. GSplats are faster to train and provide higher fidelity visual and geometric detail compared to NeRFs [6], making them a strong candidate for scene representations for robot planning. To the best of our knowledge, our work is the first to propose a planning algorithm suitable for GSplat scene representations.

**Localization.** Prior work in robot localization typically utilizes filtering schemes, such as Extended Kalman Filters (EKFs) [22, 23], Particle Filters (PFs) [24, 25], and other related filters [26, 27], to solve the pose localization problem. These methods generally estimate the pose of the robot from low-dimensional observations (measurements), extracted from the high-dimensional observations collected by the robot's onboard sensors, such as cameras. This approach often fails to leverage the entire information available in the raw, high-dimensional measurements. Learning-based filtering methods [28, 29] seek to address this limitation using deep learning to develop end-to-end frameworks for localization, computing a pose estimate directly from raw camera images. Although learning-based approaches can be quite effective given sufficient training data, these methods are often limited to a single robot platform (dynamics model) and thus require separate filters for each robot or environment.

There is some existing work on tracking the pose of a robot equipped with an on-board camera and IMU through a pre-trained NeRF map. For pose localization, these methods compute a pose that minimizes the photometric loss, given an initial guess of the camera's pose. iNeRF [30] does this for single images, and NeRF-Nav [7] and Loc-NeRF [31] both track a trajectory using a sequence of images. Other works consider simultaneous localization and mapping (SLAM) using a NeRF map representation. Existing methods such as [30, 32] all simultaneously optimize the NeRF weights and the robot/camera poses. NeRF-SLAM [33] proposes a combination of an existing visual odometry pipeline for camera trajectory estimation together with online NeRF training for the 3D scene. Although applicable to localization in Gaussian Splatting, photometric loss-based localization methods generally have a small region of convergence and require multiple passes through the scene representation for gradient computation, leading to increased computation times. In this work, we introduce a localization algorithm based on the perspective-n-point problem, which addresses these challenges.

There are a few recent works on SLAM using a GSplat representation of the environment [34, 35, 36]. These SLAM methods use the photometric loss to optimize the camera's pose, suffering from the aforementioned challenges, which we address with our proposed method. Moreover, these SLAM methods do not consider safe trajectory planning and control, as is the focus of this paper.

## III. 3D GAUSSIAN SPLATTING

**Background.** We present a brief introduction to 3D Gaussian Splatting [6], a radiance field method for deriving volumetric scene representations from a set of monocular images. Gaussian Splatting represents non-empty space in a scene using 3D Gaussian primitives, each of which is parameterized by a mean $\mu \in \mathbb{R}^3$ (defining its position), covariance matrix $\Sigma \in \mathbb{S}_{++}$ (related to its spatial extent and orientation), opacity $\alpha \in [0, 1]$, and spherical harmonics (SH) coefficients (defining view-dependent colors). The scene is typically initialized using a sparse point cloud computed via structure-from-motion [37]. To render an image from a given camera pose, the $3D$ Gaussians are projected onto the image plane using an affine approximation of the projective transformation, given by $\Sigma_{2D} = JW\Sigma W^T J^T$, with Jacobian $J$ and viewing transformation $W$. The number of primitives, along with the coefficients for each primitive, is then learned via stochastic gradient descent with a loss function comprising of the photometric loss between the rendered and ground-truth images and the structural similarity (SSIM) index loss (the same as NeRF methods).

For better numerical optimization, the anisotropic 3D covariance of each Gaussian is written as: $\Sigma = RSS^T R^T$, where

$R \in \mathrm{SO}(3)$ is a rotation matrix (parameterized by a quaternion) and $S$ is a diagonal scaling matrix (parameterized by a 3D vector). This anisotropic covariance along with adaptive density control (i.e., splitting and merging Gaussians) enable the computation of compact high-quality representations, even in complex scenes, unlike many state-of-the-art point-based rendering methods. Further, 3D Gaussian Splatting obviates the need for volumetric ray-marching required in NeRF methods, enabling high-quality real-time rendering, even from novel views.

**GSplats versus NeRFs.** Gaussian Splatting typically requires less training time than state-of-the-art NeRF methods, while achieving about the same or better photometric quality. The biggest difference is in the rendering speed, where Gaussian Splatting achieves real-time performance [6]. Moreover, 3D Gaussian Splatting enables relatively fast extraction of a mesh representation (Remark 1) of the scene from the Gaussian primitives, and instantaneous extraction of the primitives themselves. In contrast, slower meshing techniques [38] are needed for NeRFs, and the extraction of a point cloud requires slow volumetric rendering of many training viewpoints. In Fig. 2, we visualize the ground-truth mesh, the GSplat mesh, and the associated point cloud extracted from a NeRF of a simulated Stonehenge scene to showcase the collision geometry quality of GSplats over NeRFs. Quantitatively, the GSplat mesh has a smaller Chamfer distance (0.031 with 3M vertices) compared to the NeRF point cloud (0.081 with 4M points) despite having fewer points. We note that the NeRF does not necessarily yield a view-consistent geometry due to volumetric rendering, especially when the point cloud is not post-processed to remove outliers, leading to relatively poor collision geometry despite having good photometric quality.

**Remark 1.** *The original work [6] only projects 3D Gaussians whose 99% confidence interval intersects the view frustum of a camera, effectively restricting the scene representation to the 99% confidence ellipsoid associated with each Gaussian. Consequently, the union of the 99% confidence ellipsoids represents the entirety of the geometry of the scene learned during the training procedure. We find that this cutoff is too conservative, due to the fact that the color of the Gaussians toward the tails of the distribution are close to transparent. Instead, we find that renderings of the $1\sigma$ collision geometry closely matches that of the GSplat depth channel, so we elect to use $1\sigma$-ellipsoid as the collision geometry for the remainder of this work. Future work will seek to explore the calibration of this cutoff.*

**Semantic Gaussian Splatting.** To enable goal specifications for the navigation task in natural-language, we leverage semantic Gaussian Splatting [39, 40, 41], which distills 2D language semantics from vision-language models, e.g., CLIP [42], into 3D GSplat models. In general, these methods assign learnable semantic codes to each Gaussian, supervised by the robust semantic features extracted by 2D foundation models. The semantic GSplats are trained in the same way as non-semantic GSplat via gradient descent. Semantic Gaussian Splatting has been utilized in prior work to enable open-vocabulary robotics
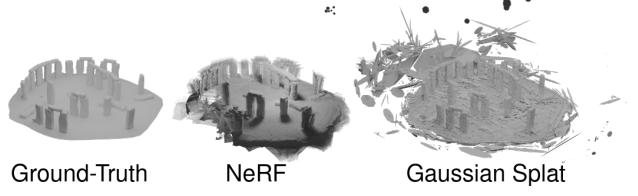


Fig. 2: Visualization of a point cloud from a NeRF and a mesh from a Gaussian Splat in the synthetic scene Stonehenge. The Chamfer Distance between the NeRF and ground-truth is 0.081 (with 4M points). The Chamfer distance between the GSplat and ground-truth is 0.031 (with 3M vertices). The collision geometry (especially the foreground) of the GSplat is better and can be extracted instantaneously from the model parameters compared to the costly rendering procedure from many viewpoints to create a point cloud from the NeRF.

tasks, e.g., robotic manipulation [43, 44].

In the subsequent sections, we present the core contributions of our work in deriving an efficient navigation pipeline for robots, describing how we leverage 3D Gaussian Splatting as the underlying scene representation. Specifically, the quick extraction of simple convex primitives (whose union closely approximates the ground-truth scene geometry) promotes the development of guarantees on safety and solution quality of Splat-Plan and facilitates real-time deployment with low sim-to-real gap while navigating in Gaussian Splatting environments. Similarly, the fast and high-quality color and depth rendering from arbitrary viewpoints of the GSplat enables robust, fast camera localization in Splat-Loc.

## IV. PLANNING WITH SAFE POLYTOPES

Now, we present Splat-Plan, our planner for GSplat maps. Splat-Plan generates safe polytopic corridors (inspired by [13]) that represent the free space of a GSplat map between an initial configuration to a goal configuration. These corridors, and the resulting trajectories through them, are rigorously built on theory derived from tests for intersection between ellipsoids. The method is fast enough to provide real-time operation, provides safety guarantees extending to any scene with a pre-trained GSplat representation, and is not overly-conservative.

We stress that, as with any safety guarantee on a map, our ultimate safety rests on the completeness of the map. If the map does not reflect the presence of an obstacle, our method may collide with the obstacle—we cannot avoid what we cannot see. In practice, we observe that GSplat maps provide fast and efficient representations of the underlying ground-truth geometry, as validated in our hardware experiments.

We would also like to motivate the use of the full collision geometry of GSplats for planning compared to conventional representations like point clouds in an RGB setting. It is common to extract the means of the GSplat to form a point cloud. However, in feature-less regions, we observe that the point cloud can be quite sparse. Meanwhile, the full collision geometry spanned by the ellipsoids covers the full surface. This phenomenon can be observed in Fig. 1, where the render of the ellipsoidal representation of the collision geometry closely

mimics the RGB render from the GSplat. However, the point cloud extracted from the means is very sparse. While usable for localization, such a sparse representation leaves a large sim-to-real gap when planning safe trajectories close to those areas. Another option is to sample the surface of these primitives for a point cloud, but even with this modification, point cloud-based planners are not as robust as Splat-Plan (Section VI).

Before presenting the planning problem, we make the following assumptions on the representations of the robot $\mathcal{R}$ and the map $\mathcal{G}$ considered in this work. We assume that the robot is represented by a union of the ellipsoids in the non-empty set $\{\mathcal{E}_{\mathcal{R},i}\}_{i=1}^{d}$, where $d$ denotes the cardinality of the set, i.e., $\mathcal{R} \subseteq \cup_{i=1}^{d}\mathcal{E}_{\mathcal{R},i}$. For simplicity, we consider a singleton set $\mathcal{E}_{\mathcal{R}}$, noting that the subsequent discussion applies directly to the non-singleton case by running the collision check for all robot ellipsoids. One can also convert a mesh or point cloud of a robot to an ellipsoid by finding the minimal bounding ellipsoid (or sphere).

We represent non-empty space in the environment with $\gamma\%$ confidence ellipsoids obtained from the GSplat map, as discussed in Remark 1, given by:

$$\mathcal{E}_j = \{x \in \mathbb{R}^3 \mid (x - \mu_j)^T \Sigma_j^{-1}(x - \mu_j) \leq \chi_3^2(\gamma)\}, \quad (1)$$

where $\mu_j \in \mathbb{R}^3$ denotes the mean of Gaussian $j$, $\Sigma_j \in \mathbb{S}_{++}$ denotes its covariance matrix, and $\chi_3^2(\gamma)$ denotes the $\gamma$th percentile of the chi-square distribution with three degrees of freedom. The union of these ellipsoids, given by $\mathcal{G} = \{\mathcal{E}_j\}_{j=1}^{N}$, defines the occupied space in the environment. To simplify notation, we express the ellipsoid in (1) in standard form: $\mathcal{E}_j = \{x \in \mathbb{R}^3 \mid (x - \mu_j)^T \Sigma_j^{-1}(x - \mu_j) \leq 1\}$, where we overload notation with $\Sigma_j := \chi_3^2(\gamma)\Sigma_j$. Based on Remark 1, we set $\gamma = 0.2$ to be safe with respect to the entirety of the supervised scene.

**Remark 2** (Online Gaussian Splatting). *Our planning algorithm requires a GSplat map. While this map can be trained online using real-time SLAM methods for radiance fields [36, 35, 34], which is a very new and active area of research, we limit the scope of this work to only plan in pre-trained maps.*

**Remark 3** (Handling Uncertainty of the Scene Representation). *We can vary the value of $\gamma$ (from that used during the training procedure) based on the quality of the GSplat map and uncertainty in different regions of the GSplat map. In general, larger values of $\gamma$ inflate the volume of the confidence ellipsoids associated with each Gaussian, resulting in greater safety margins and more conservative planning. The converse holds if smaller values of $\gamma$ are selected. Moreover, for simplicity, we utilized a uniform value of $\gamma$. However, the value of $\gamma$ can vary among the ellipsoids, allowing the planner to account for varying levels of uncertainty in different regions of the GSplat map. Likewise, the volume of the ellipsoid representing the robot can be increased/decreased to account for uncertainty in the pose of the robot.*

**Remark 4** (Dynamic Scenes). *We limit our discussion to planning in static scenes. However, we note that our method readily applies to planning in dynamic scenes, under the assumption that a dynamic Gaussian Splatting scene representation can be constructed. We discuss more about planning in dynamic scenes in Section VIII.*

**Problem Statement.** Given a bounding ellipsoid $\mathcal{E}_{\mathcal{R}}$ for the robot and a GSplat map $\mathcal{G}$, we seek to find a smooth, feasible path $x(t)$ for a robot to navigate from an initial configuration $x(0) = x_0$ to a specified goal configuration $x(T) = x_f$, such that there are no collisions in the continuum, i.e., $\mathcal{E}_{\mathcal{R}}(x(t)) \cap \mathcal{E}_j = \emptyset, \forall \mathcal{E}_j \in \mathcal{G}, \forall t \in [0, T]$.

**Collision Detection.** We leverage the ellipsoidal representations of the robot and the environment to derive an efficient collision-checking algorithm, based on [45], where we take advantage of GPU parallelization for faster computation. We build upon [45] rather than on other existing ellipsoid-to-ellipsoid intersection tests, because of its amenability to significant GPU parallelization. We do not utilize the GJK algorithm [46], since we do not require knowledge of the distance between the two ellipsoids. For completeness, we restate the collision-checking method from [45, Proposition 2].

**Theorem 1.** *Given two ellipsoids $\mathcal{E}_a, \mathcal{E}_b$ (with means $\mu_a, \mu_b$ and covariances $\Sigma_a, \Sigma_b$) and the concave function $K : (0, 1) \to \mathbb{R}$,*

$$K(s) = (\mu_b - \mu_a)^T \left[\frac{1}{1-s}\Sigma_a + \frac{1}{s}\Sigma_b\right]^{-1}(\mu_b - \mu_a),$$

*$\mathcal{E}_a \cap \mathcal{E}_b = \emptyset$ if and only if there exists $s \in (0, 1)$ such that $K(s) > 1$.*

We will let $\Sigma_{a,b}(s) = \frac{1}{1-s}\Sigma_a + \frac{1}{s}\Sigma_b$ for compactness throughout the remainder of the work. Using this, we can rewrite $K(s) = (\mu_b - \mu_a)^T \Sigma_{a,b}^{-1}(s)(\mu_b - \mu_a)$.

Note that $K(s)$ is concave in $s$ and convex with respect to the means and variances. Theorem 1 is a complete test that will always indicate whether two ellipsoids are in collision or not. We note, however, that solving the feasibility problem in Theorem 1 can be challenging, particularly in large-scale problems, where the feasibility problem has to solved for many pairs of ellipsoids with an associated matrix inversion procedure in each problem. In general, GSplat environments consists of hundreds of thousands to millions of Gaussians [6]. Consequently, we eliminate the matrix inversion by operating in a shared basis for both $\Sigma_A$ and $\Sigma_B$, for faster collision-checking, detailed in the following Proposition.

**Proposition 1.** *By solving the generalized eigenvalue problem for $\Sigma_a$ and $\Sigma_b$, we obtain generalized eigenvalues $\lambda_i$ and the corresponding matrix of generalized eigenvectors $\phi$. Then*

$$\Sigma_{a,b}^{-1}(s) = \phi \, \mathbf{diag}\left(\frac{s(1-s)}{1 + s(\lambda_i - 1)}\right)\phi^T.$$

*Proof.* We present the proof in Appendix A. $\square$

**Corollary 1.** *Given a GSplat representation with $\Sigma_j = RSS^T R^T$, let $SS^T = \mathbf{diag}\left(\lambda_i^2\right)$. If we choose to parameterize our robot body as a sphere with covariance $\Sigma_{\mathcal{R}} = \kappa^2 \mathcal{I}$, then*

$$\Sigma_{a,b}^{-1}(s) = R \, \mathbf{diag}\left(\frac{s(1-s)}{\kappa^2 + s(\lambda_i^2 - \kappa^2)}\right)R^T.$$

From hereafter, we will treat Corollary 1 as the general formula expression of $K(s)$ for both sphere and ellipsoid to

ellipsoid, substituting the appropriate variables for $R, \kappa$, and $\lambda$. For readers interested in visualizing the behavior of the aforementioned proposition and corollary, we direct readers to [47, Figure 2] to understand how the shape of $K(s)$ changes as ellipsoids move through space.

**Extension to Linear Motion.** Proposition 1 (and Corollary 1) can be extended to account for linear motion of ellipsoidal bodies. Consider a line segment starting at point $x_0$ and ending at point $x_1$, and let $\delta_x = x_1 - x_0$. Then the line segment can be parameterized as $\ell(t) = x_0 + t\delta_x$ for $t \in [0, 1]$. In our case, we consider a moving ellipsoid $\mathcal{E}_a$, which starts at $x_0 = \mu_a$. Let

$$K(s, t) = (\mu_a + t\delta_x - \mu_b)^T \Sigma_{a,b}^{-1}(s)(\mu_a + t\delta_x - \mu_b). \quad (2)$$

The ellipsoid $\mathcal{E}_a$ must satisfy Proposition 1 at all points along the line $\ell(t)$, so the safety test[1] is

$$\min_{t \in [0,1]} \max_{s \in (0,1)} K(s, t) > 1. \quad (3)$$

We seek to solve a portion of Eq. (3) using Corollary 2.

**Corollary 2.** *Note that $K(s, t)$ is a convex scalar function in $t$ because $\delta_x^T \Sigma_{a,b}^{-1}(s)\delta_x > 0$ for all $\delta_x \neq 0$ since covariance matrices are symmetric and positive definite. The $t$ that minimizes $K(s, t)$ is*

$$\hat{t}(s) = -\frac{(\mu_a - \mu_b)^T \Sigma_{a,b}^{-1}(s)\delta_x}{\delta_x^T \Sigma_{a,b}^{-1}(s)\delta_x},$$

*so the optimal value will be $t^*(s) = clamp(\hat{t}(s), 0, 1)$. Then we can write the safety check as*

$$K(s) = K(s, t^*(s)) = K(s, 0) + (\mu_a - \mu_b)^T \Sigma_{a,b}^{-1}(s)\delta_x t^*(s)$$

*Proof.* Sion's minimax theorem states that switching the order of the minimum and maximum yields identical solutions when $K(s, t)$ is concave in $s$ and convex in $t$. Additionally, $s$ must lie in a convex set and $t$ in a compact, convex set, which Eq. (3) admits. Consequently, the max-min problem results is an inner minimization problem of a quadratic, which is solved in closed-form. The point-wise minimum of concave functions $K(s, t^*(s))$ is concave, hence the outer maximization is still over a concave function. $\square$

As a byproduct of concavity of $K(s)$, we have the following corollary:

**Corollary 3.** *By concavity of $K(s)$, any approximate solution $\hat{s}$ in Theorem 1, Proposition 1, Corollary 1, or Corollary 2 results in $K(s^*) \geq K(\hat{s})$. Hence, no approximate solution will yield false negatives (i.e., miss a collision).*

**Optimization.** While Corollary 3 is a nice blanket certificate, we can craft approximate solutions that exponentially converge to the optimal $s^*$ such that false *positives* tend to 0 (i.e., a perfect approximation). Bisection searches (especially in 1D) are efficient, simple ways to guarantee exponential convergence for bounded variables in smooth convex/concave functions,

---

[1] Note that the sliding of the ellipsoid along a line forms capsules, making Corollary 2 also a necessary and sufficient collision test between this type of geometry with an ellipsoid.

i.e., $||\hat{s}_i - s^*|| \leq \epsilon$ for any desired $\epsilon$. We propose to solve $\max_{s \in [0,1]} K(s)$ using Algorithm 1.

---

**Algorithm 1:** $K(s)$ Bisection Search

**Input:** number of iterations $k$;
**Output:** maximal estimator $\hat{s}$;
```
// Initialize lower and upper bounds
```
$s_l \leftarrow 0, s_h \leftarrow 1$;
**for** $i \leftarrow 0$ **to** $k$ **do**
      `// Test midpoint`
      $\hat{s}_i \leftarrow \frac{s_l + s_h}{2}$;
      `// Find minimal t`
      $t_i^* \leftarrow \text{clamp}\left(-\frac{\delta_x^T \Sigma_{a,b}^{-1}(\hat{s}_i)(x_0-\mu)}{\delta_x^T \Sigma_{a,b}^{-1}(\hat{s}_i)\delta_x}, 0, 1\right)$;
      `// Find which way to move s`
      **if** $\nabla_s K(\hat{s}_i, t_i^*) \geq 0$ **then**
          $s_l \leftarrow \hat{s}_i$;
      **else**
          $s_h \leftarrow \hat{s}_i$;
**end**
$\hat{s} \leftarrow \hat{s}_k$;

---

**Corollary 4.** *The distance of $s$ to $s^*$ converges at a rate of $\epsilon = \frac{1}{2^{k+1}}$ through $k$ iterations using Algorithm 1.*

*Proof.* The bisection method guarantees convergence to a root of a continuous function $f(s)$ in the interval $[a, b]$ if $f(a)$ and $f(b)$ have different signs. The method achieves a rate of

$$||s_k - s^*|| \leq \frac{||b - a||}{2^{k+1}}. \quad (4)$$

Note that $K(s)$ evaluates to 0 at both $s = \{0, 1\}$, and $K(s)$ is concave and non-linear. Therefore, we know a unique global maxima occurs between 0 and 1 and that the gradient $f(s) = \nabla_s K(s_k, t_k^*)$ is positive at $s = 0$ and negative at $s = 1$. $\square$

Additionally, note that Algorithm 1 is batchable across many queries to different ellipsoids and is more efficient than performing uniform sampling for the same tolerance. For all tests, we use $k = 10$.

**Computing Safe Polytopes.** We like to again emphasize that having convex primitives (ellipsoids) as an environment representation facilitates the development of interpretable algorithms for planning. This is especially true in the construction of safe trajectories within convex safe polytopes, which define obstacle-free regions of the robot's configuration space. We build upon prior work on convex decomposition of configuration spaces such as [48, 13]. In this work, we leverage the ellipsoidal primitives to create polytopes that define the safe regions of space through the use of supporting hyperplanes. The ellipsoidal representation of the environment obtained from GSplat enables the direct computation of these convex obstacle-free regions without the need for a convex optimization procedure. Furthermore, our method is fast enough to run in real time. In the following proposition, we describe the generation of safe polytopes for a given robot.

**Proposition 2.** *Given a seed point $x^*$ for a candidate robot position and a collision set $\mathcal{G}^* = \{\mathcal{E}_j\}$, a supporting hyperplane for the ellipsoid robot can be derived from Proposition 1 or Corollary 1 for any desired buffer $\epsilon > 0$:*

$$\underbrace{\Delta_j^T \Sigma_{x^*,j}^{-1}(s^*)}_{a_j} x \geq \underbrace{(1+\epsilon)k_j^* + \Delta_j^T \Sigma_{x^*,j}^{-1}(s^*)\mu_j}_{b_j},$$

*where $\Delta_j = x^* - \mu_j$, $\Sigma_{x^*,j}^{-1}(s^*)$ uses the robot shape and $\Sigma_j$, and $(k_j^*)^2 = K(s^*) = \Delta_j^T \Sigma_{x^*,j}^{-1}(s^*)\Delta_j > 0$, for $s^* \in (0,1)$. By stacking the hyperplane constraints $(a_j, b_j)$, we arrive at a polytope $Ax \geq b$ that is guaranteed to be safe.*

*Proof.* We provide the proof in Appendix B. $\square$

**Corollary 5.** *Proposition 2 can be extended for the $K(s)$ in Corollary 2 by substituting $x^* = x_0 + t^*\delta_x$, where $t^* \in [0,1]$.*

*Proof.* We provide the proof in Appendix C. $\square$

Proposition 2 and Corollary 5 guarantee manageability, coined by [49], which refers to the encapsulation of the seed object by the free-space partition. This property is important to guarantee connected-ness of each part of the safe flight corridor, which in turn admits a feasible trajectory that resides solely within the corridor.

**Generating Safe Paths.** We present Splat-Plan, similar in spirit to the safe flight corridor method from [13]. There are four primary components: (1) feasible path seeding through graph-based search, (2) construction of a collision set around each part of the path, (3) generation of hyperplane constraints, and (4) smooth path planning posed as a spline optimization. For Splat-Plan, we leverage Corollary 2, Algorithm 1, and Corollary 5 to generate safe polytopes along the seed path. Within the polytopes, we plan Bézier curves, which can be formulated as a quadratic program.

*1) Seed Path:* There are two primary flavors of graph-based paths that are popular in the literature: those that use random trees (e.g. RRT) and those on uniform grids. We will detail how both can be used as an initialization.

Methods like RRT primarily rely on a module for collision detection at test points as well as a module to test for collision along a line. The use of Corollary 2 serves both functions. Unfortunately, the probabilistic completeness of these algorithms make them undesirable for real-time execution.

The use of a uniform grid to run algorithms like Dijkstra Search are optimal and typically faster than those of random trees **if** there exists a cheap subroutine that converts the scene representation into a uniform grid. Specifically, we would like to avoid expensive collision checking between each disjoint sub-region of $3D$ space with the environment. Conversion from point clouds to binary voxel grids circumvents this issue by binning every point and assigning it an $(i, j, k)$ index.

While there are many ways one could convert the ellipsoidal representation into a conservative occupancy grid, we propose the following method that is parallelizable and efficient, and show in Section VI that it is not too conservative. Without loss of generality, we assume that the robot is a sphere, which can be done by applying the necessary rotation and stretching

for all ellipsoids such that the robot ellipsoid is a sphere. For every ellipsoid, we calculate its axis-aligned bounding box.

The Minkowski sum of the ellipsoid with a sphere does not present an ellipsoid. However, at the extremal points which represent intersections of the bounding box with the ellipsoid, the normal of the ellipsoid is in the principal directions. The bounding box is defined as the following

$$\mathcal{B}_j = [\mu_j^i - \sqrt{\Sigma_{ii}}, \ \mu_j^i + \sqrt{\Sigma_{ii}}]_{i=1}^3, \tag{5}$$

where $\mu_j^i$ is the $i$-th element in the mean for ellipsoid $j$, and $\Sigma_{ii}$ represents the $i$-th diagonal term of the associated covariance.

For those $\mathcal{B}_j$ whose side lengths are not within the resolution of each grid cell $v_{x,y,z}$, we subdivide them by their largest side length relative to $v_{x,y,z}$. We iterate on this process until all subdivisions of $\mathcal{B}_j$ are smaller than $v_{x,y,z}$. At this point, we can calculate all 8 vertices for every subdivision of $\mathcal{B}_j$ and bin them similar to the point cloud case. This procedure can leverage batch operations on GPU and ensures that we construct an over-approximation of the collision geometry.

To account for the extent of the robot body, we convert the robot sphere into a kernel (similar to [20]) and perform a MaxPool3D operation. The resultant grid represents where the robot can be centered and be safe or unsafe. More sophisticated subdivision routines may be used to reduce the conservativism of the grid. Once the final grid is constructed, we run Dijkstra to find the seed path represented as an ordered set of connected line segments $\mathcal{L} = \{\ell_i\}_{i=1}^L$.

*2) Collision Set:* Along the seed path, rather than checking collisions between the robot and *every* ellipsoid in the scene, we would like to quickly find a subset of these primitives in the local vicinity of the robot to check against for efficiency reasons. In fact, Proposition 1 or Corollary 1 can directly be used to define a ball or ellipsoid collision set centered around the seed path, but may contain unnecessary information at the cost of additional compute.

Instead, following the paradigm of [13], we can rapidly define a bounding box oriented along $\ell_i$ and pinpoint ellipsoids that live within it without incurring the additional cost of reasoning about the linear motion of the robot body. We define a radius $r_s = \frac{v_{max}^2}{2a_{max}}$, which is the maximum stopping distance (dependent on the maximum velocity and acceleration), such that the facets of the box are no less than $r_s + \kappa$ away from $\ell_i$. This bounding box will be denoted as $A_i^{bb} x \leq B_i^{bb}$.

To check all ellipsoids that are at least partially contained within the box, we check for the minimum signed distance between each hyperplane $\min_{x \in \mathcal{E}_j} a_i^{bb} x \leq b_i^{bb}$ with every ellipsoid $\mathcal{E}_j$ in the scene. Ellipsoids that have negative signed distance for every hyperplane in the box will be at least partially contained. To perform this check, the plane and ellipsoid undergo an affine transformation to produce a new plane and an origin-centered sphere. The signed distance of the new plane from the origin must be less than 1, namely

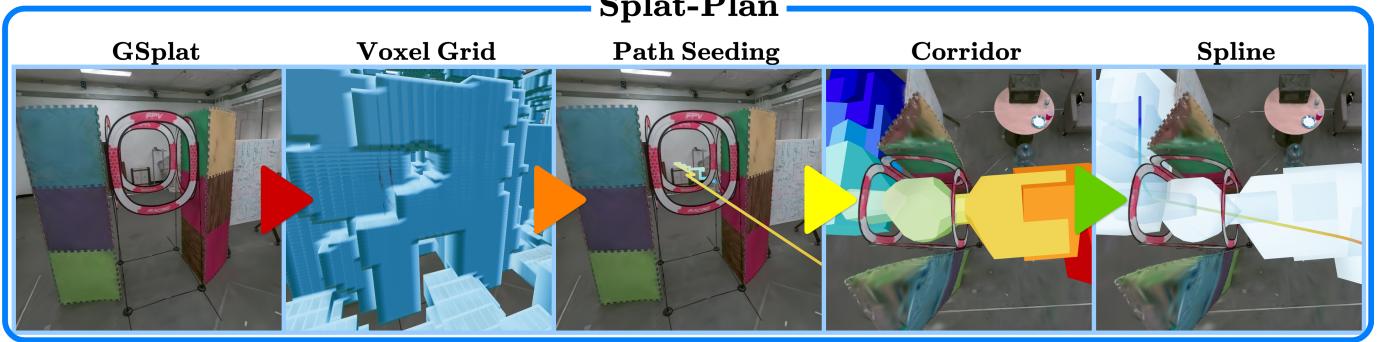$$\frac{a_i^{bb}\mu_j - b_i^{bb}}{||a_i^{bb} R_j S_j||_2} \leq 1. \tag{6}$$

Fig. 3: Splat-Plan, as described by Algorithm 2. Given a GSplat and its corresponding ellipsoidal collision geometry, Splat-Plan generates a binary occupancy grid representing the collision-less free space. Next, a seed path is created using graph-based search. Corollary 5 synthesizes a set of connected polytopes forming a corridor around the feasible path. Finally, a quadratic program is solved for the control points of a sequence of Bézier curves that lives entirely within the corridor and hence is safe.

*3) Polytope Generation:* The creation of polytopes around the line segment $\ell_i$ can be done through Corollary 5 and appending these constraints to the bounding box constraints $a_i^{bb}x \leq b_i^{bb} - \kappa\|a_i^{bb}\|_2$. Note that if we were to create a halfspace for every ellipsoid in the constraint set, we would overly constrain the free space, leading to a smaller-than-necessary polytope. This phenomenon arises from the fact that, given an existing set of halfspaces, ellipsoids that are outside of the set can still contribute non-redundant halfspaces to the existing set. Moreover, having more halfspaces than necessary in the polytope representation can significantly slow down the proceeding spline optimization.

Therefore, we adopt a greedy algorithm like [13]. Every time we form a new halfspace, we use Eq. (6) to eliminate from our collision set all ellipsoids that violate this halfspace. Of the remaining ellipsoids, we create a new halfspace for the one that had the smallest $K(s^*)$. We iterate this process until no ellipsoids remain in the collision set.

Due to manageability, we can further reduce the complexity of our corridor representation by retrieving a smaller number $P$ of polytopes than line segments $L$. For the current part of the seed path, we construct the minimal collision set and the polytope $(A_p, b_p)$. Then, we check subsequent line segments, represented as the endpoints, with the current polytope. The first instance where the line segment is not fully contained in $(A_p, b_p)$, we construct a new minimal collision set and polytope and repeat the process until the end of the seed path. Keeping more polytopes enables smoother paths (e.g. less opportunity for pinch points) at the expense of higher computation in the spline optimization phase.

*4) Spline Optimization:* Given the safe flight corridor represented as $P$ polytopes and initial and final configurations $(x_0, x_f)$, we compute a set of $P$ Bézier curves (parametrized by $M + 1$ control points $c_p^m$ and Bernstein basis $\beta^m(t)^2$ with progress $t \in [0, 1]$) representing the trajectory of the robot

---

²For notational simplicitiy, we refer to the variable as both the conventional basis and its time derivatives up to some specified order $D$.

using the path-length minimization problem:

$$\min_{c_p^m, x_p(t)} \sum_{p=1}^{P} \sum_{m=0}^{M-1} \|c_{m+1}^p - c_m^p\|_2^2 \tag{7a}$$

subject to

Safety: $A_p c_p^m \leq B_p,\ p=1,...,P; m=0,...,M$ (7b)

Configuration: $x_0(0) = x_0$ (7c)

$$x_P(1) = x_f$$

Continuity: $x_p(1) = x_{p+1}(0),\ p = 1,...,P-1$ (7d)

Bézier curve: $x_p(t) = \sum_{m=0}^{M} \beta^m(t)c_p^m, p=1,...,P-1$ (7e)

Dynamics: $x_p(t+1) = f(x_p(t), u).$ (7f)

Without the dynamics constraints (7f), the optimization problem reduces to a quadratic program that can be solved in real-time, producing a trajectory that can be tracked by differentially-flat robots. The quadratic program is solved natively using Clarabel [50].

Due to the convex hull property of Bézier curves, constraining the control points to lie in the polytopes ensures that all points along the curves will lie in the corridor and hence guarantees safety in the continuum. Additionally, Splat-Plan is sound and complete, summarized in the following corollary.

**Corollary 6.** *Given that there exists a smooth, safe path in an arbitrarily complex GSplat environment, Splat-Plan (Algorithm 2) will always return a feasible, safe, and smooth path as the limit of the voxel occupancy grid used for path seeding goes to 0, the number of steps in Algorithm 1 goes to $\infty$, and the number of control points parametrizing the spline also goes to $\infty$.*

*Proof.* As the occupancy grid resolution grows increasingly small, the unoccupied grid converges toward the true collision-free space of the scene. In the limit, Dijkstra will find an initial safe, feasible path toward the goal. Next, the constructed set of polytopes forming the corridor will (1) always be safe, (2) connected, and (3) contain the initial path assuming a

sufficient number of iterations of Algorithm 1 is performed due to Corollary 4 and manageability. Due to the Stone-Weierstrass theorem [51], given an arbitrary smooth curve, a Bézier curve of sufficient degree can exactly recover it. Under mild conditions, a smooth curve exists within the corridor, and given a sufficiently expressive parametrization and enough time, the quadratic program will find a solution.

Certainly, without these conditions, Dijkstra could fail to find a path for finite resolutions. Similarly, Algorithm 1 and Corollary 5 could return a conservative estimate that does not contain the line segment if given an insufficient number of iterations. Finally, Eq. (7) could return infeasibility if the degree of the Bézier curve is not expressive enough. Failure of these three components will lead Splat-Plan to not return a solution. However, in our experimental results, we find that this is not the case in practical settings. $\square$

**Querying Waypoints.** For simplicity of notation, we will refer to the output of Splat-Plan as $X(T)$, which takes in metric time, finds the associated spline $X_p$, and queries the spline for its position and derivatives at the local spline time $\Delta_p T = T - T_{p,start}$. In our hardware demonstrations, we un-normalize the Bézier curves in order to approximately achieve the desired $v_{\max}$. Specifically, by knowing the associated subset of the seed path (and the total length $L_p$) for each polytope, we can enforce in metric time the duration of the individual splines $\Delta T_{p,0} = \frac{L_p}{v_{max}}$. Our local Bézier curve is re-mapped using the formula $X_p(T) = x_p(\frac{\Delta_p T}{\Delta T_{p,0}})$. Eq. (7) can again be used to constrain $X_p$ to lie within the safety corridor and enforce continuity in metric time. The entire Splat-Plan algorithm is visualized in Fig. 3.

## V. MONOCULAR POSE ESTIMATION

In this section, we present our pose estimation module, Splat-Loc, for localizing a robot in a GSplat representation of its environment. This is essential to the overall functionality of the SplatNav pipeline as the safety guarantees of Splat-Plan only hold if the robot is able to consistently and accurately estimate its pose in the GSplat map. Splat-Loc only requires a monocular RGB camera, which enables it to work on a broad range of hardware platforms, including those beyond robots (such as mobile phones). Furthermore, Splat-Loc can be used either as a stand-alone pose estimation system or in conjunction with an independent pose estimation system (onboard VIO, external motion capture, etc).

**Problem Formulation.** Formally, we wish to estimate the pose of a robot at a particular time $\hat{T}_t \in \text{SE}(3)$ given a color image $I_t \in \mathbb{R}^{H \times W \times 3}$. The true camera pose $T_t$ is unknown. A pose in SE(3) is parameterized by a rotation matrix $R \in \text{SO}(3)$ and a translation vector $\tau \in \mathbb{R}^3$

$$T_t = \left( \begin{array}{c|c} R_t & \tau_t \\ \hline \mathbf{0}_{1 \times 3} & 0 \end{array} \right). \tag{8}$$

In the case that the navigating robot has an independent pose estimation system, we would like to use those pose estimates as initializations for Splat-Loc's optimization procedures, and also correct these poses using the estimates from Splat-Loc.

---

**Algorithm 2:** Splat-Plan

**Input:** $x_0$, $x_f$, grid resolution, lower, and upper bounds $(d, u_\ell, u_h)$, $v_{max}$, $a_{max}$, $\kappa$, $\mathcal{G}$, num. iters. $k$;
**Output:** Bézier spline $x(t)$;
```
// Create the voxel grid
```
$V \leftarrow \text{GSplatVoxelGrid}(d, u_\ell, u_h, \kappa, \mathcal{G})$;
```
// Path Seeding: creates safe but
   non-smooth path
```
$\{\ell\}_{i=1}^L \leftarrow V(x_0, x_f)$;
```
// Iterate through the line segments
```
**for** $i \leftarrow 1$ **to** $L$ **do**
    ```// Skip to next line segment if
      contained```
    **if** *IsInPolytope*$(\ell_i, (A_{curr}, B_{curr}))$ **then**
        | continue;
    ```// Create collision set```
    $\mathcal{G}_i, (A_i^{bb}, B_i^{bb}) \leftarrow$
      $\text{GetCollisionSet}(\ell_i, \kappa, \mathcal{G}, v_{max}, a_{max})$;
    ```// Initialize polytope```
    $A_i, B_i \leftarrow A_i^{bb}, B_i^{bb}$;
    ```// Create polytope```
    **while** $|\mathcal{G}_i| > 0$ **do**
        $\{K_j\} \leftarrow \text{Algorithm 1}(\ell_i, \mathcal{G}_i, k)$;
        $K \leftarrow \min(\{K_j\})$;
        ```// Create halfspace```
        $(A, B) \leftarrow \text{Corollary 5}(\ell_i, K)$;
        ```// Add halfspace to polytope```
        $A_i \leftarrow \text{append}(A_i, A), B_i \leftarrow \text{append}(B_i, B)$;
        ```// Reject redundant ellipsoids```
        $\mathcal{G}_i \leftarrow \text{Eq. (6)}(\mathcal{G}_i, (A, B))$;
    **end**
    ```// Set current polytope```
    $A_{curr} \leftarrow A_i, B_{curr} \leftarrow B_i$;
**end**
```
// Optimize Bézier splines
```
$x(t) \leftarrow \text{Optimize}(x_0, x_f, \{(A_p, B_p)\}_{p=1}^P)$

---

We assume knowledge of the camera's calibration including the intrinsic matrix and distortion coefficients for projective geometry. These are easily computable, and are often available from the camera manufacturer.

**Lightweight Monocular Pose Estimator.** At its core, Splat-Loc uses the fast rendering capabilities of GSplats and standard tools from camera tracking to formulate Perspective-n-Point (PnP) problems, which can be reliably solved using off-the-shelf optimizers, and produces accurate estimates of the robot pose. As input for the pose estimation procedure, we have the color image and a coarse initial guess for the pose estimate, $\hat{T}_{t,0}$. This guess can either come from an independent localization module (e.g. VIO) or can simply be the previous time step's estimate. We begin by rendering an RGB image using the GSplat map with the camera pose set to the initial guess and simultaneously generate a local point-cloud within the camera's view, effectively using the GSplat as a monocular depth estimator.

Next, a local feature extractor is used to compute visual

features (keypoints and descriptors) in both the camera image and the rendered image. Each keypoint has an associated pixel coordinate $(u, v) \in \mathbb{R}^2$, and let $m$ and $n$ respectively be the number of keypoints in the camera and rendered images. A feature matcher is used to determine correspondences between the visual features in the camera image and the rendered image. Let $\ell \leq \min\{m, n\}$ be the number of successfully matched features. In our experiments we found that the feature extractor SuperPoint [52] used in conjunction with the transformer-based LightGlue [53] feature matcher had the best performance (see Section VI for more details).

Using the rendered depth image and the camera intrinsics matrix, the keypoints from the rendered color image can be projected into the 3D to produce a point cloud. Let $\hat{p}_j \in \mathbb{R}^3$ be the position of the $j$th projected keypoint where $j \in 1, \ldots, n$. Finally, we seek to minimize the following reprojection error in order to find the relative pose transform that transforms $\hat{T}_{t,0}$ to $\hat{T}_t$

$$\bar{T}_t = \arg\min_{T \in \mathrm{SE}(3)} \sum_{k=1}^{\ell} ||\rho_k - KT\hat{p}_k||_2, , \qquad (9)$$

where $\rho_k = [u_k, v_k, 1]^\top$ and subsequently recover our estimated pose $\hat{T}_t = \bar{T}_t \hat{T}_{t,0}$. Eq. (9) is the Perspective-n-Point problem, and is a nonlinear least-squares optimization problem that we solve using the Levenberg-Marquardt algorithm. In practice, we use Random Sample Consensus (RANSAC) to remove outliers from the set of matched features which results in more robust solutions of Eq. (9). We illustrate this procedure in Figure 1. In Section VI, we highlight the accuracy of incremental estimation in real-world experiments while a drone navigates a cluttered environment.

**Global Initialization.** The above pose estimation procedure requires common overlap between $I_t$ and $\hat{I}_t$, necessitating a reasonably accurate initial estimate of the robot's pose $\hat{T}_{t,0}$, which may not be available in many practical settings. When a good initial guess of the robot's pose is unavailable, we execute a global pose estimation procedure. Note that this only needs to be performed once, and then subsequent pose estimate steps can be performed using the solution from the previous iteration.

One approach requires a monocular depth estimator, e.g., [54, 55], to augment the RGB image obtained by the robot with depth information, which is used to generate a point cloud (in the camera frame). Another is to randomly sample $\mathrm{SE}(3)$ for pose initializations and return the pose estimate from the P$n$P run that has the lowest reprojection error. However, we instead generate a point cloud of the scene from the GSplat means $\{\mu_j\}_{j=1}^N$, enabling the formulation of a point-cloud registration problem:

$$\hat{T}_0 = \arg\min_{T \in \mathrm{SE}(3)} \sum_{(p,q) \in \mathcal{C}} ||p - Tq||_2^2, \qquad (10)$$

where $\mathcal{C}$ denotes the set of correspondences, associating the point $p$ in the map cloud to a point $q$ in the point cloud from the camera. If we are given a known set of correspondences, we can compute the optimal solution of (10) using Umeyama's method [56].

In practice, we do not have prior knowledge of the set of correspondences $\mathcal{C}$ between the two point clouds. To address this challenge, we apply standard techniques in feature-based global point-cloud registration. We begin by computing 33-dimensional Fast Point Feature Histograms (FPFH) descriptors [57] for each point in the point-cloud, encoding the local geometric properties of each point. Prior work has shown that visual attributes can play an important role in improving the convergence speed of point-cloud registration algorithms [58], something that FPFH does not do. To solve this, we augment the FPFH feature descriptor of a given point with its RGB color. We then identify putative sets of correspondences using a nearest-neighbor query based on the augmented FPFH descriptors, before running RANSAC to iteratively identify and remove outliers in $\mathcal{C}$. The RANSAC convergence criterion is based on the distance between the aligned point clouds and the length of a pair of edges defined by the set of correspondences.

**Non-invasive Pose Correction.** While fusing Splat-Loc poses with existing pose estimates like VIO is beyond the scope of this work, we will address challenges that arises when using Splat-Plan to plan high-level plans in a GSplat while using existing pose estimates to stabilize (i.e., for control). Fundamentally, discrepancies between the one in which the GSplat is trained in $\mathcal{T}_{\mathrm{gs}}$ and the running coordinate frame of the existing localization module $\mathcal{T}_{(\mathrm{control},t)}$ can vary with time, either due to noise or drift. Yet, poses from Splat-Loc are inherently tied to the GSplat coordinate frame, leading to potentially more informative state estimates of whether the robot is in collision or not. In turn, these estimates can be passed into Splat-Plan to create safer trajectories if necessary, as depicted in Fig. 1. However, the trajectory that Splat-Plan returns again lives in $\mathcal{T}_{\mathrm{gs}}$ and not necessarily the running coordinate frame of the existing localization, which is crucially used for control. To overcome this mismatch, we necessarily need to transform the outputs of Splat-Plan into the control localization frame. Namely, there exists a transform ${}^{\mathrm{control},t}T_{\mathrm{gs}} : \mathcal{T}_{\mathrm{gs}} \rightarrow \mathcal{T}_{(\mathrm{control},t)}$ that maps poses in the GSplat frame to ones in the control localization frame. Therefore, the waypoints that we send to the robot are ${}^{(\mathrm{control},t)}T_{\mathrm{gs}}(X(T))$, which is depicted in Fig. 1 as the input to the robot.

## VI. EXPERIMENTS

We demonstrate the effectiveness of our navigation pipeline for GSplat maps, examining its performance in real-world scenes on hardware and in simulation. In addition, we perform ablative studies comparing our algorithms against existing methods.

### A. Simulation Results

*1) Test Environments:* We benchmark Splat-Plan and Splat-Loc independently on four different environments: **Stonehenge**, a fully-synthetic scene, and three real-world scenes **Statues**, **Flightroom**, and **Old Union**. For **Stonehenge**, we captured image-pose pairs by rendering the **Stonehenge** mesh in Blender. For the other scenes, we recorded a video from a mobile phone and processed the image frames through structure-from-motion (COLMAP [37]) to retrieve corresponding camera poses and intrinsics.

*2) Splat-Loc Evaluations:* We compare Splat-Loc to existing pose estimation methods, including a baseline GS-Loc, based on the localization component of existing GSplat SLAM methods [34, 35]. We leverage finite differences to estimate the gradient of the photometric loss function utilized in the pose estimator, which might not be particularly fast or robust, especially for larger errors in the initial pose estimate. While these methods optimize over the re-rendering loss composed of the photometric loss, and in some cases, depth and semantic-related loss terms, in our baseline, we optimize only over the photometric loss, since we assume the robot in these evaluations does not have an RGB-D camera for depth measurements. As a result, our baseline essentially matches the GSplat SLAM method in [36]. In addition, we compare our pose estimator to the Point-to-plane Iterative Closest Point (ICP) [59] and Colored-ICP [58] algorithms, assuming these point-cloud methods have privileged 3D information that the incremental estimation of Splat-Loc does not have. Furthermore, we examine two variants of our pose estimator: Splat-Loc-Glue, which utilizes LightGlue for feature matching; and Splat-Loc-SIFT, which utilizes SIFT for feature matching.

In each scene, we run 10 trials (of 100 frames each) of each pose estimation algorithm. We evaluate the rotation error (R.E.) and translation error (T.E.) with respect to the ground-truth pose, the computation time (C.T.) per frame, and the overall success rate (S.R.). Here, success indicates the generation of a solution regardless of its quality. The performance of pose estimation algorithms often depends on the error associated with the initial estimate of the pose. As such, we test our system across a range of different errors in the initial estimate of the pose. In this study, we assume an initial estimate of the pose is available. We generate the initial estimate by taking the ground truth pose then applying a rotation $\delta_R$ about a random axis and the translation $\delta_t$ in a random direction.

We provide the summary statistics of the error in the pose estimates computed by each algorithm, in addition to the computation time on a trial with 100 frames in the **Statues** scene in Table I. We note that all methods had a perfect success rate in this problem. The GS-Loc algorithm achieves the lowest accuracy and requires the greatest computation time, unlike Colored-ICP, Splat-Loc-SIFT, and Splat-Loc-Glue, which achieve much-higher accuracy with a rotation error less than a degree and a translation error less than 15cm. GS-Loc requires a computation time of about 36.15 s per frame, which is about two orders of magnitude slower than the next-slowest method ICP, which requires a computation time of about 110 ms. Colored-ICP, Splat-Loc-SIFT, and Splat-Loc-Glue require less than 100 ms of computation time. Compared to all methods, Splat-Loc-Glue yields pose estimates with the lowest mean rotation and translation error, less than $0.06°$ and 4 mm, respectively, and achieves the fastest mean computation time, less than 42 ms. The computation time of Splat-Loc may be about a standard deviation greater during the first call, which may be due to the time spent loading the models and initializing the GPU kernels.

Lastly, we examine the performance of the pose estimation algorithms in problems with a larger error in the initial estimate of the pose, with $\delta_R = 30°$ and $\delta_t = 0.5$ m in the

TABLE I: Comparison of baseline pose estimation algorithms in simulation in the **Statues** scene with $\delta_R = 20°$ and $\delta_t = 0.1$ m.

| Algorithm | R.E. (deg.) | T.E. (cm) | C.T. (msec.) | S.R. (%) |
|---|---|---|---|---|
| ICP [59] | $73.1 \pm 45.9$ | $129 \pm 75$ | $107 \pm 19.2$ | 100 |
| Colored-ICP [58] | $0.83 \pm 0.37$ | $1.31 \pm 0.60$ | $43.3 \pm 9.70$ | 100 |
| Splat-Loc-SIFT (ours) | $0.09 \pm 0.06$ | $0.56 \pm 0.75$ | $63.3 \pm 2.39$ | 100 |
| Splat-Loc-Glue (ours) | $\mathbf{0.05 \pm 0.03}$ | $\mathbf{0.33 \pm 0.27}$ | $\mathbf{41.2 \pm 73.2}$ | 100 |
| GS-Loc [34, 35, 36] | $122 \pm 33.8$ | $245 \pm 91.2$ | $36200 \pm 5440$ | 100 |

synthetic **Stonehenge** scene. We present the performance of each algorithm on each metric in Table II, where we note that ICP and Colored-ICP do not provide accurate estimates of the robot's pose. Moreover, the pose estimation errors achieved by ICP and Colored-ICP have a significant variance. In contrast, Splat-Loc-SIFT and Splat-Loc-Glue yield pose estimates of high accuracy with average rotation and translation errors less than 0.5 deg. and 5mm, respectively. However, Splat-Loc-SIFT achieves a lower success rate, compared to Splat-Loc-Glue, which achieves a perfect success rate.

TABLE II: Comparison of baseline pose estimation algorithms in the **Stonehenge** scene with $\delta_R = 30°$ and $\delta_t = 0.5$ m.

| Algorithm | R.E. (deg.) | T.E. (cm) | C.T. (msec.) | S.R. (%) |
|---|---|---|---|---|
| ICP [59] | $131 \pm 22.6$ | $370 \pm 554$ | $122 \pm 153$ | 100 |
| Colored-ICP [58] | $94.9 \pm 51.3$ | $57.4 \pm 28.8$ | $488 \pm 104$ | 20 |
| Splat-Loc-SIFT (ours) | $\mathbf{0.217 \pm 0.0369}$ | $0.334 \pm 0.0563$ | $139 \pm 3.15$ | 70 |
| Splat-Loc-Glue (ours) | $0.220 \pm 0.203$ | $\mathbf{0.315 \pm 0.210}$ | $\mathbf{45.1 \pm 0.611}$ | 100 |

*3) Splat-Plan Evaluations:* Splat-Plan is benchmarked against three different methods: a point-cloud planner [13], a sampling-based planner (RRT* using Proposition 1), and a NeRF-based planner (NeRF-Nav [7]). Furthermore, we perform ablations against variations of the point-cloud planner in order to expose flaws when planning against point clouds compared to the full scene geometry. For each simulation scene, we train a dense and sparse GSplat, totaling 8 scenes. In every scene, we run 100 start and goal locations distributed in a circle around the boundary of the scene.

In the simulated tests, we represent the robot using balls of various sizes in order to generate interesting trajectories due to the fact that the simulated scenes are not trained in metric scale.[3] . Additional parameters, such as the number of Gaussians, can be found in Table III.

TABLE III: Parameters across scenes for experiments. Number of Gaussians is reported for both dense and sparse variants of the same scene.

| | Radius | $V_{max}$ | $A_{max}$ | Resolution | N. Gauss (K) |
|---|---|---|---|---|---|
| Stonehenge | 0.01 | 0.1 | 0.1 | $150^3$ | 116/12 |
| Statues | 0.03 | 0.1 | 0.1 | $100^3$ | 201/18 |
| Flight | 0.02 | 0.1 | 0.1 | $100^3$ | 281/4 |
| Old Union | 0.01 | 0.1 | 0.1 | $100^3$ | 525/87 |
| Maze | 0.25 | 0.5 | 1.0 | $80^3$ | 100 |
| Maze (fast) | 0.25 | 1.5 | 1.0 | $80^3$ | 100 |

While point cloud-based planners are ubiquitously used, they can sometimes fall short when the scene geometry is not dense

---

[3]Nerfstudio adopts the NeRF conventions in scaling the scene to fit within the confines of a two-unit-length cube centered at the origin, with the poses of the camera residing within a $[-1, 1]^3$-bounding box. We disable this feature for the hardware **Maze** scene.

or if the scene is very cluttered. To this end, we developed four variants of the Safe Flight Corridor (SFC) [13]. SFC-1 ingests the GSplat means as a point cloud, runs Dijkstra to retrieve a feasible initial path seed, creates collision sets with respect to the point cloud, synthesizes a polytope corridor that marginally intersects with the point cloud, and finally deflates the polytopes by the robot radius. These polytopes are fed to the same spline optimizer (7) that Splat-Plan uses. SFC-2 executes the same pipeline as SFC-1, but the point cloud representation is sampled from the surface of the ellipsoids. We sample 20 points from each ellipsoid in the scene to simulate a typical amount of points a Lidar or depth image would produce (approximately 2-5 million points). SFC-3 uses the Splat-Plan occupancy grid to retrieve a feasible path seed, while the means are still used to create polytopes. Finally, SFC-4 uses the Splat-Plan occupancy grid, synthesizes polytopes using the means, but deflates the polytope by the robot radius and the maximum eigenvalue of the ellipsoid whose mean was used to create a particular halfspace in the polytope. These variants are all potential solutions to apply SFC to GSplat environments. We summarize the tradeoffs of all methods in Table IV.

TABLE IV: Splat-Plan strikes a favorable tradeoff compared to existing methods in terms of safety, non-conservativeness (NC), smoothness, solution feasibility, and real-time execution.

| | Safe | NC | Smooth | Feasible | RT | Env. |
|---|---|---|---|---|---|---|
| NeRF-Nav | ✗ | N/A | ✗ | ✓ | ✗ | NeRF |
| RRT* | ✗ | ✗ | ✗ | ✓ | ✗ | GS |
| SFC-1 | ✗ | ✓ | ✓ | ✗ | ✓ | GS |
| SFC-2 | ✗ | ✓ | ✓ | ✗ | ✓ | GS |
| SFC-3 | ✗ | ✓ | ✓ | ✗ | ✓ | GS |
| SFC-4 | ✗ | ✗ | ✓ | ✗ | ✓ | GS |
| **Splat-Plan** | ✓ | ✓ | ✓ | ✓ | ✓ | GS |

Visually, the paths generated by Splat-Plan are smooth, safe, and non-conservative (Fig. 4). This fact is validated in Fig. 6, where Splat-Plan's trajectories in blue are safe (minimum distances greater than 0 with respect to the GSplat collision geometry). Unfortunately, because many of these scenes were captured in the real-world, no ground-truth mesh exists. Moreover, we inspect the point cloud and mesh created by COLMAP and notice poor overall reconstruction of the collision geometry. Therefore, we elected to use the GSplat ellipsoidal geometry in place of the ground-truth geometry due to its high-quality approximation.

Notice that these trajectories are non-conservative compared to the SFC methods (low path lengths and high polytope volume in Figs. 5 and 6). More importantly, we see that Splat-Plan never fails to return a trajectory, highlighted by the 0 failure rate. All other methods have failures, other than NeRF-Nav by virtue of it being an end-to-end optimization method. Finally, Splat-Plan has comparable execution times to SFC. Note that as SFC does not use GPU, we rewrote the codebase in Pytorch to yield comparable times to Splat-Plan.

Finally, in terms of memory, we observe that in the scene with the most Gaussians (**Old Union**), GPU memory usage hovered around 3.1 GB, with the GSplat itself requiring 1.6 GB and the binary occupancy grid, 1.5GB.

## B. Hardware Results

*1) Test Environment:* We test Splat-Nav in the **Maze** scene using a drone. Images to train **Maze** were captured using the RGB camera onboard the drone. We utilize Nerfstudio [60] to train the Semantic GSplat, using its default parameters (which includes estimating the camera poses for each image frame from structure-from-motion via COLMAP [37]). In Figure 8, we show the true training images captured by the drone, the rendered RGB image from the GSplat at the same camera pose, and the semantic relevancy for the associated language query. First, we note that the rendered image is photorealistic, highlighting the remarkable visual quality of the trained Gaussian Splat. Second, the semantic relevancy spatially agrees with the expected location of the queried object, making the semantic field suitable for open-vocabulary goal querying.

*2) Hardware:* We test our pipeline on the Modal AI development drone platform measuring 29 cm x 20 cm x 10 cm (diagonal length of 36.6 cm). In the hardware tests, we approximate the robot using a sphere with diameter $0.5\,\mathrm{m}$. Readers can find our test parameters in Table III. An OptiTrack motion capture system is solely used for evaluation purposes. Any other markers, such as ArUco tags, in the scene are purely cosmetic.

*3) Implementation:* We run the pose estimator and the planner ROS2 nodes on a desktop computer with an Nvidia RTX 4090 GPU and an Intel i9 13900K CPU, which communicates with the drone via WiFi. We emphasize that both modules are running asynchronously. At a frequency of about $3\,\mathrm{Hz}$, the drone transmits images from its cameras and associated VIO poses to the desktop computer. Splat-Loc ingests the VIO pose $\hat{T}_{t,0}$ and the image $I_t$ to compute the pose estimate $\hat{T}_t$ of the drone body after applying rigid body transforms to transform the camera pose to the body frame. We run the estimator continuously, synchronized with the stream of images published from the drone via ROS2.

The planning module ingests $\hat{T}_t$ as $x_0$ and computes a safe trajectory for the drone to follow toward the language-conditioned goal $x_f$. The re-plan node, which runs Splat-Plan based on $x_0$, updates the spline(s) $X(T)$ as frequently as possible. The waypoint node, which operates asynchronously from the re-plan node, measures the running time since the $X(T)$ was last updated and returns positions, velocities, acceleration, and jerk at this running time. The waypoint node runs at $10\,\mathrm{Hz}$. This architecture allows the drone to continue following a smooth spline even when Splat-Plan is still computing the next plan. However, we find that simply sending position waypoints can cause the drone to jerk when $X(T)$ is updated, as successive splines need not be close to one another. To rectify this issue, we forward integrate the waypoint velocities to get positions. Finally, these positions undergo $^{(\mathrm{control},t)}T_{\mathrm{gs}}(X(T))$ before being sent to the drone. Additionally, we run a Kalman Filter to smooth $^{(\mathrm{control},t)}T_{\mathrm{gs}}$, as both the Splat-Loc and VIO pose estimates can be somewhat noisy.

*4) Goal Specification:* In the **Maze**, we specify the goal locations for the drone via natural language, comprising of the following objects: a keyboard, beachball, phonebook, and microwave (depicted in Figure 8). We query the semantic

**Stonehenge** 🪨      **Statues** 🗿      **Flightroom** 🚁      **Old Union** 🪑
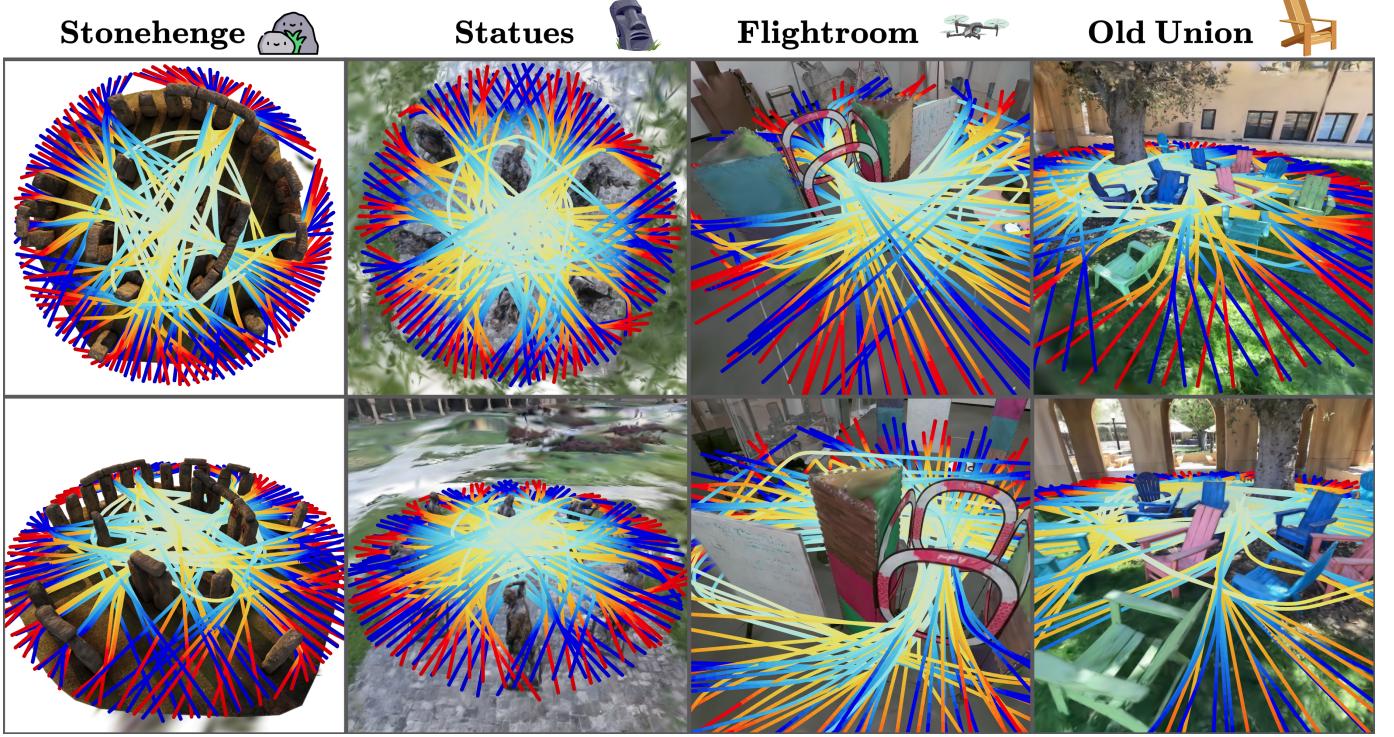


Fig. 4: Qualitative results of 100 safe trajectories using Splat-Plan with start (blue) and goal (red) states spread over a circle. We see that the trajectories are safe, but not conservative.

Gaussian Splat for the location of these objects using the following text prompts: "keyboard," "beachball," "phonebook" and "microwave," corresponding to these objects, without negative prompts. These objects are placed in locations that require dynamic motions, such as hard turns and elevation maneuvers, to reach. Moreover, all tests begin at the same position at hover, and the objects are positioned relative to this position so that they are not immediately visible when the drone first begins flight.

*5) Control Schemes:* Our hardware tests consist of three different control schemes, coined Open-loop, Closed-loop VIO, and Splat-Loc. Open-loop tests do not not re-plan, and therefore does not use Splat-Loc estimates. One trajectory is created at the start $T = 0$, and the control node returns the corresponding waypoint at that point in time. No forward integration of the velocities is necessary since only one trajectory is ever created. Closed-loop VIO and Splat-Loc are re-planning control schemes where the re-plan node updates the trajectory $X(T)$ as frequently as possible. Closed-loop VIO uses the VIO estimate as $x_0$ and no additional transform is applied to the waypoint. Conversely, Splat-Loc uses the Splat-Loc pose estimate as $x_0$, and the smoothed $^{(\text{control},t)}T_{\text{gs}}$ is applied to the Splat-Plan trajectories to transform them into the VIO control frame. Our hardware tests consist of all combinations of goal locations and control schemes. In addition, we run these combinations 10 times for statistical significance, yielding a total of 120 flights.

*6) Splat-Loc Evaluations:* We validate the performance of Splat-Loc in hardware experiments in the **Maze** scene, showing that Splat-Loc achieves relatively the same level of accuracy as the onboard VIO in estimating the drone's pose, without

requiring any special calibration or re-initialization procedures for frame alignment, which the onboard VIO requires. In Table V, we provide the rotation and translation errors of the Splat-Loc estimates, with the MOCAP poses as the ground-truth estimates. We note that Splat-Loc achieves rotation errors of about 3 deg and translation errors of about 4 cm, which is comparable to the accuracy of the VIO estimates, shown in Table VI. However, Splat-Loc failed in one of the closed-loop trials with the "keyboard" goal location. As a result, the rotation and translation errors for this goal location is higher compared to the those of the other goal locations. The failure case is visualized in Figure 9, where the drone goes past the keyboard. We note that the failure likely occurred because the drone's camera was pointing towards an area of the scene which was not really covered in the video used in training the GSplat. We discuss strategies for addressing such failure cases in Section VIII. In Figure 8, we show the estimated trajectories of the drone using MOCAP, the onboard VIO, and Splat-Loc, demonstrating the effectiveness of Splat-Loc. Essentially, all the pose estimators achieve comparable estimation accuracy. However, unlike the MOCAP system, Splat-Loc does not require a specialized hardware system and is amenable to any monocular camera. Moreover, Splat-Loc runs at about 25 Hz on average, which is fast-enough for real-time operation. The bulk of the computation time is utilized in computing the feature matches and in solving the PnP problem, which requires about 10 milliseconds.

*7) Splat-Plan Evaluations:* Visualizations of 120 trajectories across four goal locations and three control schemes can be found in Fig. 9. Note that all flights were collision-free with
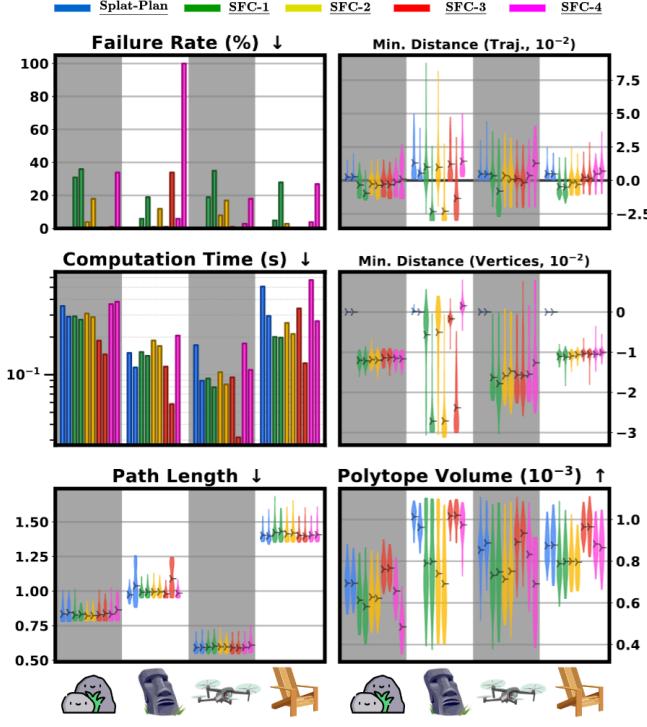
Fig. 5: A comparison of trajectories generated by Splat-Plan and four variants of SFC [13] for the same 100 start and end locations for each scene (Stonehenge, Statues, Flightroom, and Old Union). Bars on the left represent trajectories planned using the dense scene, while the sparse scene is on the right. We evaluate the safety (Distance to GSplat and failure rate), conservativeness (path length and polytope volume), and computation time. We also show the distance of the polytope vertices to the GSplat to illustrate that Splat-Plan is sound. Violin plots showcase the spread of the minimum distance, path length, and polytope volume across all 100 trajectories, with markings indicating the mean. Splat-Plan displays competitive non-conservativeness and computation time, while exhibiting superior safety and success rates.

TABLE V: Rotation Error (R.E.), Translation Error (T.E.), and the Success Rate (S.R.) of Splat-Loc with MOCAP as the ground-truth reference.

| Goal | R.E. (deg.) | T.E. (cm) | S.R.(%) |
|------|-------------|-----------|---------|
| Beachball | $2.82 \pm 0.12$ | $3.93 \pm 0.86$ | 100 |
| Keyboard | $5.61 \pm 3.50$ | $8.14 \pm 0.35$ | 90 |
| Microwave | $2.82 \pm 0.23$ | $3.59 \pm 0.42$ | 100 |
| Phonebook | $2.83 \pm 0.25$ | $3.37 \pm 0.95$ | 100 |

respect to the true scene except for one flight using Splat-Loc to navigate to the keyboard. The drone was oriented toward the edge of the scene, where features were few and the GSplat quality was poor. The poor quality can be attributed to the lack of training images pointing toward the edges of the scene, as we wanted to reconstruct the foreground in the highest quality. These qualitative results indicate that, within the confines of our controlled setting, all control schemes work equally well.
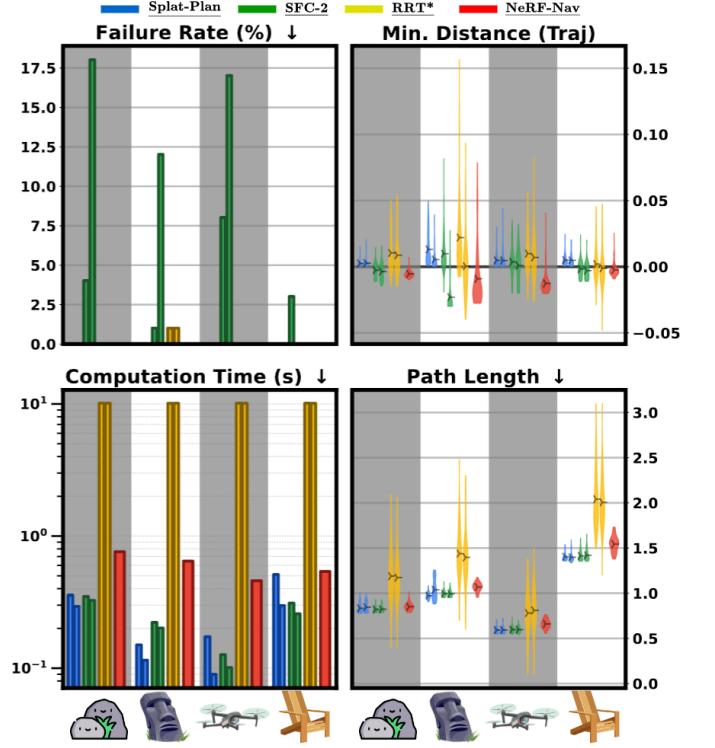


Fig. 6: A comparison of trajectories generated by Splat-Plan, SFC [13], RRT*, and NeRF-Nav [7]) for the same 100 start and end locations for each scene (Stonehenge, Statues, Flightroom, and Old Union). Bars on the left represent trajectories planned using the dense scene, while the sparse scene is on the right. We evaluate the safety (Distance to GSplat and failure rate), conservativeness (path length), and computation time. Violin plots showcase the spread of the minimum distance and path length across all 100 trajectories, with markings indicating the mean. Splat-Plan is safer than competing methods, relatively fast, non-conservative, and never fails across all 8 problem settings.

TABLE VI: Rotation Error (R.E.), Translation Error (T.E.), and the Success Rate (S.R.) of VIO pose estimates with MOCAP as the ground-truth reference.

| Goal | R.E. (deg.) | T.E. (cm) | S.R.(%) |
|------|-------------|-----------|---------|
| Beachball | $2.42 \pm 0.13$ | $3.87 \pm 0.89$ | 100 |
| Keyboard | $2.36 \pm 0.27$ | $3.28 \pm 1.20$ | 100 |
| Microwave | $2.36 \pm 0.14$ | $4.29 \pm 1.30$ | 100 |
| Phonebook | $2.40 \pm 0.23$ | $3.63 \pm 1.23$ | 100 |

These results are promising for Splat-Loc from a convenience point of view. We noticed that the VIO of the drone would drift in subsequent runs, necessitating the reinitialization of the VIO at the start of every run. In addition, as the VIO is not calibrated to be in the GSplat frame, we manually aligned the frames by zero-ing the VIO of the drone at the same position for all flights and for collection of training data. Meanwhile, Splat-Loc needed no such alignment, and was kept running continuously throughout all experiments without zero-ing (even
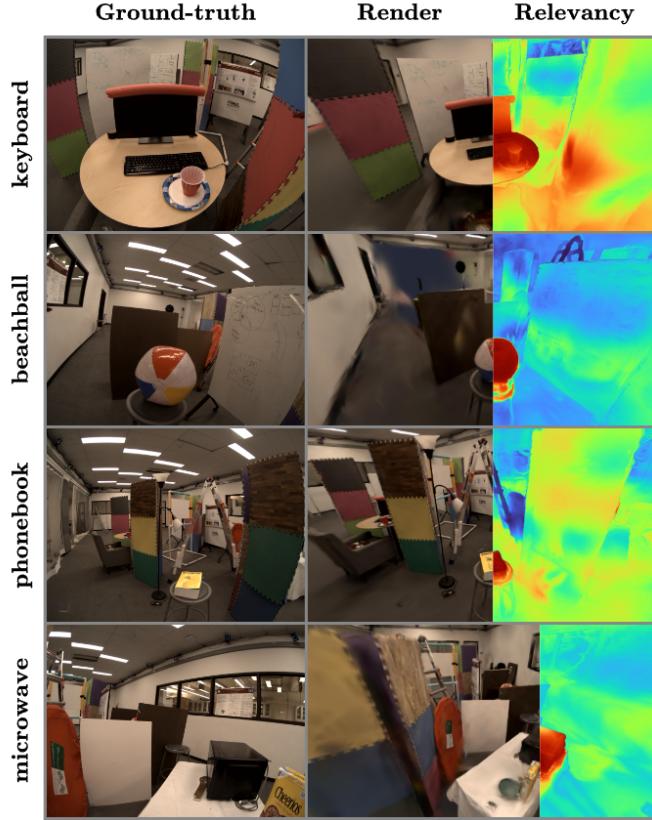
Fig. 7: Natural language-specified goal locations in the **Maze** scene. The rendered RGB image from the GSplat demonstrate good reconstruction of the ground-truth. Additionally, the semantic relevancy spatially agrees with the provided language query.
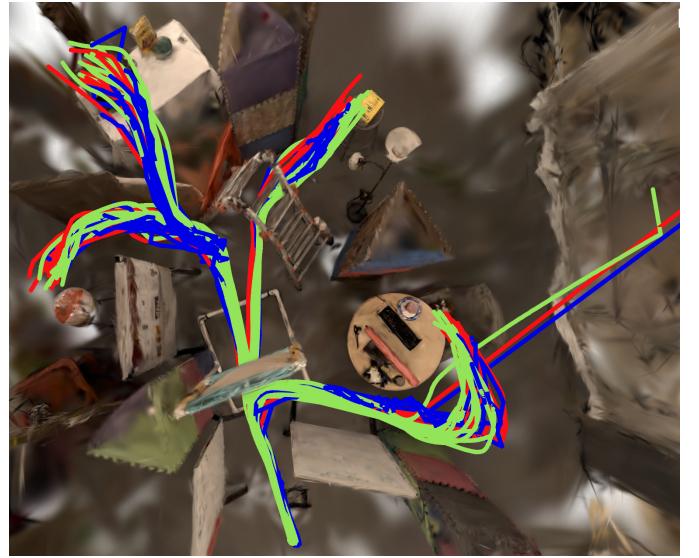


Fig. 8: Pose estimates of Splat-Nav using motion capture (green), onboard VIO (red), and Splat-Loc (blue). Splat-Loc gives comparable performance without requiring a manual frame alignment between the GSplat and a separate localization frame.

in control schemes that do not use Splat-Loc).

Qualitatively, we see similar trends in Fig. 10. All control schemes are unsafe at different times, but in similar amounts. Note that some curves dip below 0, yet are verifiably safe in real-life. This is due to a variety of reasons, the most of prominent of which are: the difference in the set robot radius (0.25 cm) versus the true radius (0.18 cm), errors in aligning the motion capture frame into the frame of the GSplat (because the GSplat was not trained using motion capture), and the tracking capabilities of the drone. Note that the safety violation of all control schemes is relatively small compared to the size of the drone, which allows error in low-level tracking to obfuscate advantages of one method over another, especially in cluttered environments.

*8) Fast Control:* We stress test Splat-Plan by increasing $v_{\max}$ until the onboard VIO could no longer track the desired waypoint with enough accuracy to avoid collision, which was $1.5\,\mathrm{m/s}$. These speeds, coupled with the clutter in the environment, allowed for dynamic flight, which is visualized in the right column of Fig. 9. We point readers toward the associated videos hosted on our website (https://chengine.github.io/splatnav/) to better visualize the trajectories.

*9) Closed-loop Endurance:* Finally, we stress test the Splat-Loc re-planning pipeline through endurance flights. The

pipeline is left to continually execute. Once the drone reaches a goal location, another goal location is set. We demonstrate collision-free flight over the order of minutes, which can again be visualized on our website (https://chengine.github.io/splatnav/).

## VII. CONCLUSION

We introduce an efficient navigation pipeline termed *Splat-Nav* for robots operating in GSplat environments. Splat-Nav consists of a guaranteed-safe planning module *Splat-Plan*, which allows for real-time planning ($> 2\,\mathrm{Hz}$) by leveraging the ellipsoidal representation inherent in GSplats for efficient collision-checking and safe corridor generation, facilitating real-time online replanning. Splat-Plan demonstrates superior performance in terms of conservativeness, safety, success rate and comparable computation times compared to point-cloud and NeRF methods on the same scene. Moreover, our proposed pose estimation module *Splat-Loc* computes high-accuracy pose estimates faster ($25\,\mathrm{Hz}$) and more reliably compared to existing pose estimation algorithms for radiance fields, such as NeRFs. We present extensive hardware and simulation results, highlighting the effectiveness of Splat-Nav.

## VIII. LIMITATIONS AND FUTURE WORK

We only tested Splat-Nav in pre-constructed scenes. Existing GSplat SLAM algorithms do not run in real-time [36, 35], limiting the application of our method in online mapping. In future work, we seek to examine the derivation of real-time GSplat mapping methods, integrated with the planning and pose estimation algorithms proposed in this work. Additionally, the results from Section VI-A2 suggest that we can incorporate Splat-Loc as a localization module within online GSplat SLAM
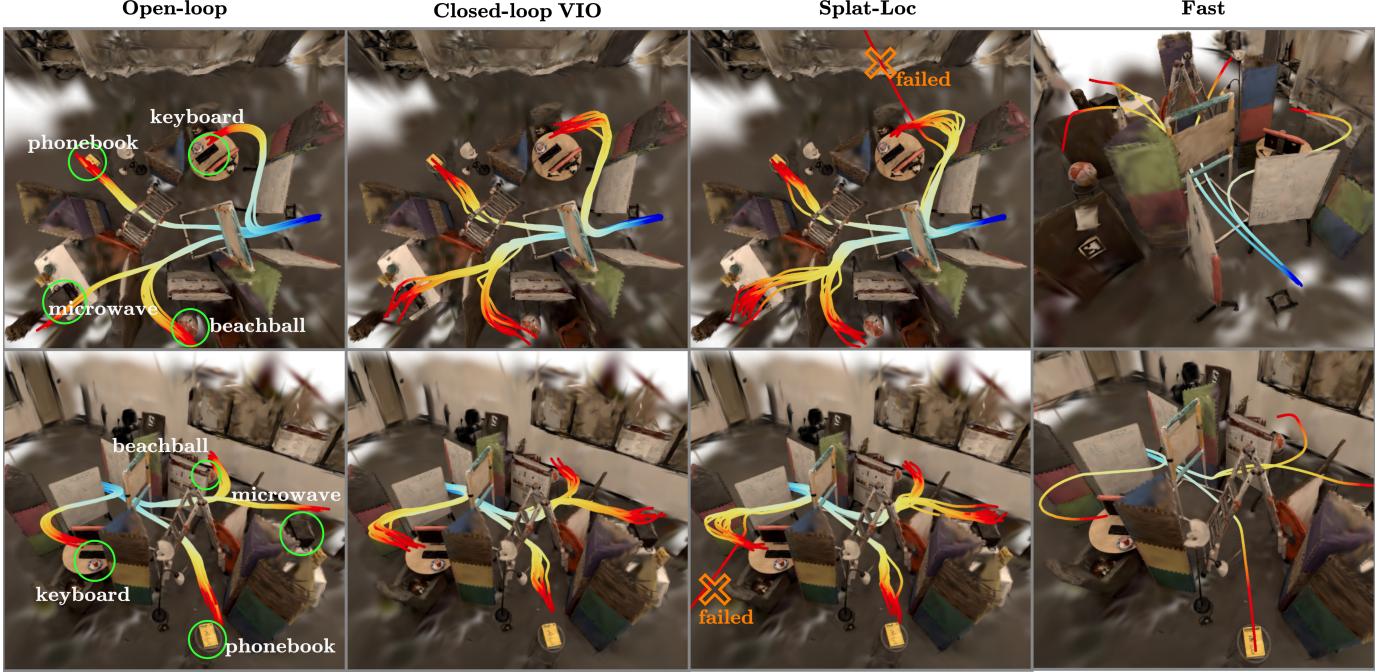
Fig. 9: Ground-truth trajectories of the drone navigating, projected onto the **Maze** GSplat. The drone is subjected to different goal locations and control schemes. The recorded flight trajectories are best viewed on our website https://chengine.github.io/splatnav/.

algorithms to improve localization accuracy and the resulting map quality.

We assumed that the pre-constructed scenes were correct. Safety of the planned trajectories depends on the quality of the underlying GSplat map. As noted in Remark 3, we can use different confidence levels of the ellipsoids to account for uncertainty in an object in the GSplat map. Splat-Plan cannot do anything if an obstacle is completely missing from the scene, which is a fundamental limitation of the GSplat map representation. Likewise, Splat-Plan could fail if the initialization graph-search procedure which utilizes Dijkstra fails to find a path to the goal, which could occur in maps with a coarse resolution. Future work will examine uncertainty quantification of different regions within a GSplat scene to aid the design of active-planning algorithms that enable a robot to collect additional observations in low-quality regions, such as areas with missing/non-existent geometry, while updating the GSplat scene representation via online mapping.

We only tested Splat-Nav in a static scene. This could be a limitation in many practical problems. Using NeRFs and GSplat for dynamic environments remains an open area of research, especially in scenes without prerecorded motion [61, 62, 63]. Splat-Plan is fast enough to be extended easily to problems with dynamic scenes so long as the underlying dynamic GSplat representation is available.

The performance of Splat-Loc depends on the presence of informative features in the scene. We can address this in two ways: through planning and by incorporating additional sensor data. Future work will explore the design of planning algorithms that bias the path towards feature-rich regions, improving localization accuracy during path execution. Future work will also incorporate IMU data to improve the robustness of the pose estimator, particularly in featureless regions of the scene where the PnP-RANSAC procedure might fail.

Splat-Plan and Splat-Nav require loading the GSplat model onto the GPU, which takes up about $10\,\text{GB}$ of GPU memory. Many drone platforms do not have the onboard compute resources to load the GSplat model, hindering onboard computation. Future work will seek to reduce the memory-usage demands of GSplat models, e.g., using sparse GSplat models.

## APPENDIX A
### PROOF OF PROPOSITION 1

This proposition was stated by [64] without proof. We provide the proof here. The goal is to diagonalize both $\Sigma_a$ and $\Sigma_b$ in such a way that they share the same eigenbasis. In this way, the matrix inversion amounts to a reciprocal of the eigenvalues. We can find such an eigenbasis through the generalized eigenvalue problem, which solves the following system of equations $\Sigma_a \phi_i = \lambda_i \Sigma_b \phi_i$. These generalized eigenvalues and eigenvectors satisfy the identity $\Sigma_a \phi = \Sigma_b \phi \Lambda$.

To solve the generalized eigenvalue problem, we utilize the fact that $\Sigma_a$ and $\Sigma_b$ are symmetric, positive-definite, so the Cholesky decomposition of $\Sigma_b = LL^T$, where $L$ is lower triangular. This also means that we can represent $\Sigma_a$ by construction as: $\Sigma_a = LCL^T = LV\Lambda V^T L^T$, where $C$ is also symmetric, positive-definite, so its eigen-decomposition is $C = V\Lambda V^T$, where the eigenvectors $V$ are orthonormal $V^T V = I$. To retrieve the original eigenbasis $\phi$, we solve the triangular system $L^T \phi = V$. A consequence of this solution is that: $V^T V = \phi^T LL^T \phi = \phi^T \Sigma_b \phi = I$. Moreover, this means that $\Sigma_b = \phi^{-T} \phi^{-1}$. We can show that $\Lambda$ and $\phi$ indeed solve
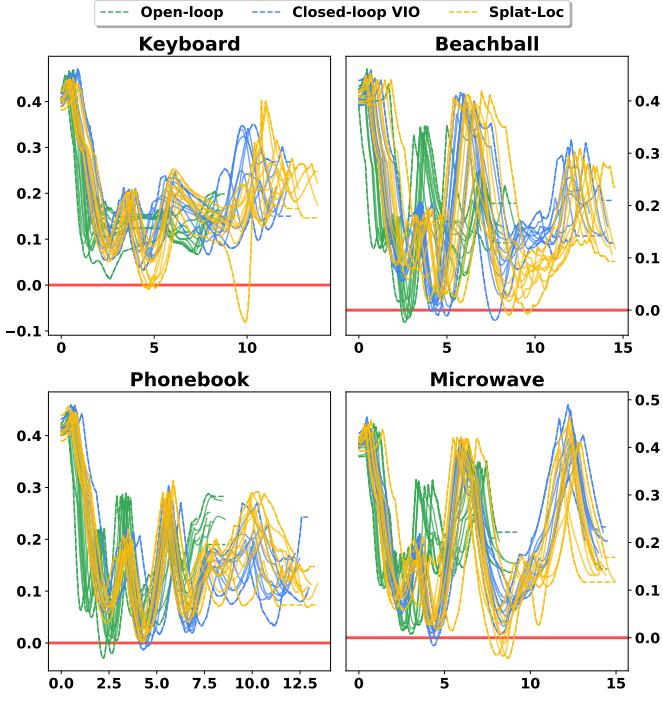
Fig. 10: Distances (meters) of goal-conditioned trajectories to the Gaussian Splat as measured by the motion capture poses across three different control schemes. The abscissa represents time in seconds. The spread and the individual trajectories are visualized.

the generalized eigenvalue problem through the following substitution:

$$\Sigma_a \phi = LV\Lambda V^T L^T \phi = \underbrace{LL^T}_{\Sigma_b} \phi \Lambda \phi^T \underbrace{LL^T}_{\Sigma_b} \phi = \Sigma_b \phi \Lambda \underbrace{\phi^T \Sigma_b \phi}_{I},$$

where $\Lambda = \mathbf{diag}(\lambda_i)$ is the diagonal matrix of eigenvalues.

As a result, we can write the matrix inversion in Theorem 1 as the following:

$$\left[\frac{1}{1-s}\Sigma_a + \frac{1}{s}\Sigma_b\right]^{-1} = \left[\frac{1}{1-s}\Sigma_b\phi\Lambda\phi^T\Sigma_b + \frac{1}{s}\Sigma_b\right]^{-1}$$

$$= \left[\phi^{-T}\left(\frac{1}{1-s}\Lambda + \frac{1}{s}I\right)\phi^{-1}\right]^{-1}$$

$$= \phi\,\mathbf{diag}\left(\frac{s(1-s)}{1+s(\lambda_i-1)}\right)\phi^T.$$

## APPENDIX B
### PROOF OF PROPOSITION 2

From Proposition 2, we know that we must find an $x$ such that $K(s) > 1$ for all ellipsoids in the test set $\mathcal{G}^*$. We will approximate this safe set as a polytope. Given a test point $x^*$ and the $j$th ellipsoid in the collision test set $\mathcal{G}^*$, we can use our collision test (Corollary 2) to derive these polytopes. Notice that $\mathcal{E}_a$ represents the position of the robot (i.e., the mean is test point $\mu_a = x^*$). For any value of $\bar{s} \in (0,1)$, the safety test is a quadratic constraint of the test point $x^*$, namely $\Delta_j^T\Sigma_{x^*,j}^{-1}\Delta_j > 1$, where we omit the dependence of $\Sigma$ on $s$ for brevity.

To derive the supporting hyperplane for the safety check, we will first find the point on the ellipsoid and then linearize our test about this point. Let $f_j(x) = (x-\mu_j)^T\Sigma_{x^*,j}^{-1}(x-\mu_j)$ for any arbitrary $x$. Note that $f_j(x^*) = \Delta_j^T\Sigma_{x^*,j}^{-1}\Delta_j = k_j^2$. We can immediately see that point $x_0 = \mu_j + \frac{1+\epsilon}{k_j}\Delta_j$ will be outside of the ellipsoid for any $\epsilon > 0$ since:

$$f_j(x_0) = \frac{(1+\epsilon)^2}{k_j^2}\Delta_j^T\Sigma_{x^*,j}^{-1}\Delta_j = (1+\epsilon)^2 > 1. \quad (11)$$

Note that this point $x_0$ lies on the ray starting from the center of the ellipsoid $\mu_j$ and passing through the center of the robot at $x^*$. We then linearize the constraint $f_j(x)$ about $x_0$. Taking the derivative yields $\frac{df_j}{dx}(x) = 2(x-\mu_j)^T\Sigma_{x^*,j}^{-1}$. The linear approximation of the constraint is then:

$$f_j(x) \approx f_j(x_0) + \left[\frac{df_j}{dx}(x_0)\right](x-x_0) \geq 1. \quad (12)$$

Plugging in $x_0 = \mu_j + \frac{1+\epsilon}{k_j}\Delta_j$ and simplifying yields the given expression $\Delta_j^T\Sigma_{x^*,j}^{-1}x \geq (1+\epsilon)k_j + \Delta_j^T\Sigma_{x^*,j}^{-1}\mu_j$.

We need to also prove that the above constraint is a supporting hyperplane by showing that all feasible points when the constraint is equality is necessarily outside of the ellipsoid parametrized by $\Sigma_{x^*,j}^{-1}$. Recall that $x_0 = \mu_j + \frac{\Delta_j}{k_j}$ is a point both on the surface of the ellipsoid and on the hyperplane. Therefore, all feasible points on the hyperplane can be expressed as $x = x_0 + \delta$, where $\Delta_j^T\Sigma_{x^*,j}^{-1}\delta = 0$. If the hyperplane supports the ellipsoid parametrized by $\Sigma_{x^*,j}^{-1}$, then necessarily $f(x) = (x-\mu_j)^T\Sigma_{x^*,j}^{-1}(x-\mu_j) \geq 1$. For all points on the hyperplane, the ellipsoid constraint evaluates to

$$f_j(x_0+\delta) = \left(\frac{1+\epsilon}{k_j}\Delta_j + \delta\right)^T\Sigma_{x^*,j}^{-1}\left(\frac{1+\epsilon}{k_j}\Delta_j + \delta\right)$$

$$= \underbrace{\frac{(1+\epsilon)^2}{k_j^2}\Delta_j^T\Sigma_{x^*,j}^{-1}\Delta_j}_{=(1+\epsilon)^2} + \underbrace{\delta^T\Sigma_{x^*,j}^{-1}\delta}_{\geq 0}$$

$$+ \underbrace{\frac{2(1+\epsilon)}{k_j}\Delta_j^T\Sigma_{x^*,j}^{-1}\delta}_{=0} > 1.$$

Hence, the constraint in Proposition 2 is a supporting hyperplane.

## APPENDIX C
### PROOF OF COROLLARY 5

Our objective in this section is to show that the hyperplane created by Corollary 5 always contains the line segment $\ell(t) = x_0 + t\delta_x$ for $t \in [0,1]$. When the minimum unconstrained $t^*$ occurs outside $t \in [0,1]$, then necessarily the constrained minimum occurs at either of the two endpoints due to convexity of $K$ with respect to $t$. Therefore, we consider the case where $t^* \in (0,1)$ and the case where $t^* \in \{0,1\}$. Let $x^* = x_0 + t^*\delta_x$.

#### A. Optimum in the Interior

Note that $t^*$ occurs when the normal of this ellipsoid $\Sigma_{x^*,b}^{-1}(x_0+t^*\delta_x-\mu)$ is perpendicular to the direction of motion

$\delta_x$. Given that $K(s^*) \geq 1$ (Corollary 2), Section B states that the point $x^*$ is safe and satisfies the hyperplane constraint

$$(x^* - \mu_j)^T \Sigma_{x^*,b}^{-1}(s^*)(x - \mu_j) \geq k_j^*. \tag{13}$$

All points $x$ along the line can be parameterized as $x^* + \tilde{t}\delta_x$ for $\tilde{t} = t - t^*$, which due to orthogonality of $\delta_x$ to the ellipsoid normal, yields

$$(x^* - \mu_j)^T \Sigma_{x_0,b}^{-1}(s^*)(x^* + \tilde{t}\delta_x - \mu_j) =$$
$$(x^* - \mu_j)^T \Sigma_{x_0,b}^{-1}(s^*)(x^* - \mu_j) + \tilde{t}\underbrace{(x^* - \mu_j)^T \Sigma_{x_0,b}^{-1}(s^*)\delta_x}_{0} =$$
$$(x^* - \mu_j)^T \Sigma_{x_0,b}^{-1}(s^*)(x^* - \mu_j) \geq k_j^*. \tag{14}$$

Therefore, the entire line segment is contained in the hyperplane and is not in collision.

### B. Optimum on the Boundary

Boundary optima will not necessarily yield orthogonality of $\delta_x$ with the normal of the ellipsoid $\Sigma_{x^*,j}^{-1}$. Without loss of generality, define $x^* = x_0$, which is the closest point on the line segment to the ellipsoid $\mathcal{E}_j$ since the direction of the line segment can always be flipped to achieve this result. The unconstrained optima will occur at some point $x_0 + \hat{t}\delta_x$, where $\hat{t} < 0$. Therefore,

$$\delta_x^T \Sigma_{x^*,b}^{-1}(x_0 + \hat{t}\delta_x - \mu_j) = 0$$
$$\delta_x^T \Sigma_{x^*,b}^{-1}(x_0 - \mu_j) + \hat{t}\underbrace{\delta_x^T \Sigma_{x^*,b}^{-1}\delta_x}_{\geq 0} = 0 \tag{15}$$
$$\therefore (x^* - \mu_j)^T \Sigma_{x^*,j}^{-1}\delta_x \geq 0.$$

Following the reasoning of Eq. (14), for all points on the line $x = x^* + t\delta_x \ \forall \ t \in [0, 1]$, we have $(x^* - \mu_j)^T \Sigma_{x^*,j}^{-1}(s^*)(x - \mu_j) \geq k_j^*$. Again, in this case, the line segment still remains within the polytope. Hence, we have proven Corollary 5.

### APPENDIX D
### GLOBAL POSE INITIALIZATION OF SPLAT-LOC

We evaluate the global point-cloud alignment initialization procedure from Section V when utilized by Splat-Loc-Glue, in the real-world scene Statues and synthetic scene Stonehenge, with the results in Table VII. We see that the success rate is 80% in both scenes, and it runs at approximately 20-30 Hz. When it does succeed, the pose errors are on the order of 0.3 deg and 1 cm. Although these estimates are relatively good, we note that these are less accurate than the recursive pose estimates discussed above.

TABLE VII: Performance of the pose initialization module in Splat-Loc.

| Scene | R.E. (deg.) | T.E. (cm) | C.T. (msec.) | S.R. (%) |
|---|---|---|---|---|
| Statues | $0.37 \pm 0.72$ | $1.34 \pm 2.33$ | $33.9 \pm 0.57$ | 80 |
| Stonehenge | $0.21 \pm 0.21$ | $0.30 \pm 0.23$ | $45.8 \pm 1.70$ | 80 |

### REFERENCES

[1] A. Elfes, "Using occupancy grids for mobile robot perception and navigation," *Computer*, vol. 22, no. 6, pp. 46–57, Jun. 1989.

[2] H. Edelsbrunner, "Surface Reconstruction by Wrapping Finite Sets in Space," in *Discrete and Computational Geometry: The Goodman-Pollack Festschrift*, ser. Algorithms and Combinatorics, B. Aronov, S. Basu, J. Pach, and M. Sharir, Eds. Berlin, Heidelberg: Springer, 2003, pp. 379–404.

[3] P. Kim, J. Chen, and Y. K. Cho, "Slam-driven robotic mapping and registration of 3d point clouds," *Automation in Construction*, vol. 89, pp. 38–48, 2018.

[4] S. Osher and R. P. Fedkiw, *Level Set Methods and Dynamic Implicit Surfaces*. New York: Springer, 2003.

[5] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," in *European Conference on Computer Vision (ECCV)*, 2020.

[6] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3D Gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics*, vol. 42, no. 4, July 2023. [Online]. Available: https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/

[7] M. Adamkiewicz, T. Chen, A. Caccavale, R. Gardner, P. Culbertson, J. Bohg, and M. Schwager, "Vision-only robot navigation in a neural radiance world," *IEEE Robotics and Automation Letters (RA-L)*, vol. 7, no. 2, pp. 4606–4613, 2022.

[8] S. LaValle, "Planning algorithms," *Cambridge University Press google schola*, vol. 2, pp. 3671–3678, 2006.

[9] L. E. Kavraki, P. Svestka, J.-C. Latombe, and M. H. Overmars, "Probabilistic roadmaps for path planning in high-dimensional configuration spaces," *IEEE transactions on Robotics and Automation*, vol. 12, no. 4, pp. 566–580, 1996.

[10] S. Karaman, M. R. Walter, A. Perez, E. Frazzoli, and S. Teller, "Anytime motion planning using the rrt," in *2011 IEEE international conference on robotics and automation*. IEEE, 2011, pp. 1478–1483.

[11] N. Ratliff, M. Zucker, J. A. Bagnell, and S. Srinivasa, "Chomp: Gradient optimization techniques for efficient motion planning," in *2009 IEEE international conference on robotics and automation*. IEEE, 2009, pp. 489–494.

[12] J. Schulman, Y. Duan, J. Ho, A. Lee, I. Awwal, H. Bradlow, J. Pan, S. Patil, K. Goldberg, and P. Abbeel, "Motion planning with sequential convex optimization and convex collision checking," *The International Journal of Robotics Research*, vol. 33, no. 9, pp. 1251–1270, 2014.

[13] S. Liu, M. Watterson, K. Mohta, K. Sun, S. Bhattacharya, C. J. Taylor, and V. Kumar, "Planning dynamically feasible trajectories for quadrotors using safe flight corridors in 3D complex environments," *IEEE Robotics and Automation Letters (RA-L)*, vol. 2, no. 3, pp. 1688–1695, 2017.

[14] J. Van den Berg, M. Lin, and D. Manocha, "Reciprocal velocity obstacles for real-time multi-agent navigation,"

in *2008 IEEE international conference on robotics and automation.* Ieee, 2008, pp. 1928–1935.

[15] X. Wu, S. Chen, K. Sreenath, and M. W. Mueller, "Perception-aware receding horizon trajectory planning for multicopters with visual-inertial odometry," *IEEE Access*, vol. 10, pp. 87 911–87 922, 2022.

[16] H. Oleynikova, Z. Taylor, M. Fehr, R. Siegwart, and J. Nieto, "Voxblox: Incremental 3D Euclidean signed distance fields for on-board MAV planning," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017.

[17] L. Han, F. Gao, B. Zhou, and S. Shen, "Fiesta: Fast incremental euclidean distance fields for online motion planning of aerial robots," *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 4423–4430, 2019.

[18] K. Goel, W. Tabib, and N. Michael, "Rapid and high-fidelity subsurface exploration with multiple aerial robots," in *Experimental Robotics: The 17th International Symposium.* Springer, 2021, pp. 436–448.

[19] K. Goel and W. Tabib, "Incremental multimodal surface mapping via self-organizing gaussian mixture models," *IEEE Robotics and Automation Letters*, vol. 8, no. 12, pp. 8358–8365, 2023.

[20] T. Chen, P. Culbertson, and M. Schwager, "Catnips: Collision avoidance through neural implicit probabilistic scenes," *arXiv preprint arXiv:2302.12931*, 2023.

[21] M. Tong, C. Dawson, and C. Fan, "Enforcing safety for vision-based controllers via control barrier functions and neural radiance fields," in *IEEE International Conference on Robotics and Automation (ICRA).* IEEE, 2023, pp. 10 511–10 517.

[22] L. J. Guibas, R. Motwani, and P. Raghavan, "The robot localization problem," *SIAM Journal on Computing*, vol. 26, no. 4, pp. 1120–1138, 1997.

[23] A. Eman and H. Ramdane, "Mobile robot localization using extended kalman filter," in *2020 3rd International Conference on Computer Applications & Information Security (ICCAIS).* IEEE, 2020, pp. 1–5.

[24] D. Fox, S. Thrun, W. Burgard, and F. Dellaert, "Particle filters for mobile robot localization," in *Sequential Monte Carlo methods in practice.* Springer, 2001, pp. 401–428.

[25] Q.-b. Zhang, P. Wang, and Z.-h. Chen, "An improved particle filter for mobile robot localization based on particle swarm optimization," *Expert Systems with Applications*, vol. 135, pp. 181–193, 2019.

[26] I. Ullah, Y. Shen, X. Su, C. Esposito, and C. Choi, "A localization based on unscented kalman filter and particle filter localization algorithms," *IEEE Access*, vol. 8, pp. 2233–2246, 2019.

[27] J. Biswas and M. Veloso, "Depth camera based indoor mobile robot localization and navigation," in *2012 IEEE International Conference on Robotics and Automation.* IEEE, 2012, pp. 1697–1702.

[28] P. Karkus, D. Hsu, and W. S. Lee, "Particle filter networks with application to visual localization," in *Conference on robot learning.* PMLR, 2018, pp. 169–178.

[29] R. Jonschkowski, D. Rastogi, and O. Brock, "Differen-

tiable particle filters: End-to-end learning with algorithmic priors," *arXiv preprint arXiv:1805.11122*, 2018.

[30] L. Yen-Chen, P. Florence, J. T. Barron, A. Rodriguez, P. Isola, and T.-Y. Lin, "Inerf: Inverting neural radiance fields for pose estimation," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS).* IEEE, 2021, pp. 1323–1330.

[31] D. Maggio, M. Abate, J. Shi, C. Mario, and L. Carlone, "Loc-nerf: Monte carlo localization using neural radiance fields," in *IEEE International Conference on Robotics and Automation (ICRA).* IEEE, 2023, pp. 4018–4025.

[32] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys, "Nice-SLAM: Neural implicit scalable encoding for SLAM," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022.

[33] A. Rosinol, J. J. Leonard, and L. Carlone, "NeRF-SLAM: Real-Time Dense Monocular SLAM with Neural Radiance Fields," *arXiv preprint arXiv:2210.13641*, 2022.

[34] C. Yan, D. Qu, D. Wang, D. Xu, Z. Wang, B. Zhao, and X. Li, "GS-SLAM: Dense visual slam with 3d gaussian splatting," *arXiv preprint arXiv:2311.11700*, 2023.

[35] V. Yugay, Y. Li, T. Gevers, and M. R. Oswald, "Gaussian-SLAM: Photo-realistic dense slam with gaussian splatting," *arXiv preprint arXiv:2312.10070*, 2023.

[36] H. Matsuki, R. Murai, P. H. Kelly, and A. J. Davison, "Gaussian splatting slam," *arXiv preprint arXiv:2312.06741*, 2023.

[37] J. L. Schönberger and J.-M. Frahm, "Structure-from-motion revisited," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[38] W. E. Lorensen and H. E. Cline, "Marching cubes: A high resolution 3d surface construction algorithm," in *Seminal graphics: pioneering efforts that shaped the field*, 1998, pp. 347–353.

[39] M. Qin, W. Li, J. Zhou, H. Wang, and H. Pfister, "Langsplat: 3d language gaussian splatting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 20 051–20 060.

[40] S. Zhou, H. Chang, S. Jiang, Z. Fan, Z. Zhu, D. Xu, P. Chari, S. You, Z. Wang, and A. Kadambi, "Feature 3dgs: Supercharging 3d gaussian splatting to enable distilled feature fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 21 676–21 685.

[41] O. Shorinwa, J. Sun, and M. Schwager, "Fast-splat: Fast, ambiguity-free semantics transfer in gaussian splatting," *arXiv preprint arXiv:2411.13753*, 2024.

[42] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning (ICML).* PMLR, 2021, pp. 8748–8763.

[43] M. Ji, R.-Z. Qiu, X. Zou, and X. Wang, "Graspsplats: Efficient manipulation with 3d feature splatting," *arXiv preprint arXiv:2409.02084*, 2024.

[44] O. Shorinwa, J. Tucker, A. Smith, A. Swann, T. Chen,

R. Firoozi, M. D. Kennedy, and M. Schwager, "Splat-mover: Multi-stage, open-vocabulary robotic manipulation via editable gaussian splatting," in *8th Annual Conference on Robot Learning*, 2024.

[45] I. Gilitschenski and U. D. Hanebeck, "A robust computational test for overlap of two arbitrary-dimensional ellipsoids in fault-detection of kalman filters," in *2012 15th International Conference on Information Fusion*, 2012, pp. 396–401.

[46] E. G. Gilbert, D. W. Johnson, and S. S. Keerthi, "A fast procedure for computing the distance between complex objects in three-dimensional space," *IEEE Journal on Robotics and Automation*, vol. 4, no. 2, pp. 193–203, 1988.

[47] I. Gilitschenski and U. D. Hanebeck, "A robust computational test for overlap of two arbitrary-dimensional ellipsoids in fault-detection of kalman filters," in *2012 15th International Conference on Information Fusion*. IEEE, 2012, pp. 396–401.

[48] S. Ruan, K. L. Poblete, H. Wu, Q. Ma, and G. S. Chirikjian, "Efficient path planning in narrow passages for robots with ellipsoidal components," *IEEE Transactions on Robotics*, vol. 39, no. 1, pp. 110–127, 2022.

[49] Q. Wang, Z. Wang, M. Wang, J. Ji, Z. Han, T. Wu, R. Jin, Y. Gao, C. Xu, and F. Gao, "Fast iterative region inflation for computing large 2-d/3-d convex regions of obstacle-free space," 2024. [Online]. Available: https://arxiv.org/abs/2403.02977

[50] P. J. Goulart and Y. Chen, "Clarabel: An interior-point solver for conic programs with quadratic objectives," 2024.

[51] L. De Branges, "The stone-weierstrass theorem," *Proceedings of the American Mathematical Society*, vol. 10, no. 5, pp. 822–824, 1959.

[52] D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superpoint: Self-supervised interest point detection and description," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 224–236.

[53] P. Lindenberger, P. Sarlin, and M. Pollefeys, "Lightglue: Local feature matching at light speed. arxiv 2023," *arXiv preprint arXiv:2306.13643*, 2023.

[54] S. F. Bhat, R. Birkl, D. Wofk, P. Wonka, and M. Müller, "Zoedepth: Zero-shot transfer by combining relative and metric depth," *arXiv preprint arXiv:2302.12288*, 2023.

[55] L. Yang, B. Kang, Z. Huang, X. Xu, J. Feng, and H. Zhao, "Depth anything: Unleashing the power of large-scale unlabeled data," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 10 371–10 381.

[56] S. Umeyama, "Least-squares estimation of transformation parameters between two point patterns," *IEEE Transactions on Pattern Analysis & Machine Intelligence*, vol. 13, no. 04, pp. 376–380, 1991.

[57] R. B. Rusu, N. Blodow, and M. Beetz, "Fast point feature histograms (fpfh) for 3d registration," in *2009 IEEE international conference on robotics and automation*. IEEE, 2009, pp. 3212–3217.

[58] J. Park, Q.-Y. Zhou, and V. Koltun, "Colored point cloud registration revisited," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 143–152.

[59] Y. Chen and G. Medioni, "Object modelling by registration of multiple range images," *Image and vision computing*, vol. 10, no. 3, pp. 145–155, 1992.

[60] M. Tancik, E. Weber, R. Li, B. Yi, T. Wang, A. Kristoffersen, J. Austin, K. Salahi, A. Ahuja, D. McAllister, A. Kanazawa, and E. Ng, "Nerfstudio: A framework for neural radiance field development," in *SIGGRAPH*, 2023.

[61] G. Wu, T. Yi, J. Fang, L. Xie, X. Zhang, W. Wei, W. Liu, Q. Tian, and X. Wang, "4d gaussian splatting for real-time dynamic scene rendering," *arXiv preprint arXiv:2310.08528*, 2023.

[62] A. Pumarola, E. Corona, G. Pons-Moll, and F. Moreno-Noguer, "D-nerf: Neural radiance fields for dynamic scenes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10 318–10 327.

[63] C. Gao, A. Saraf, J. Kopf, and J.-B. Huang, "Dynamic view synthesis from dynamic monocular video," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5712–5721.

[64] N. A. (https://math.stackexchange.com/users/3060/nick alger), "Detect if two ellipses intersect," Mathematics Stack Exchange, 2020, uRL:https://math.stackexchange.com/q/3678498 (version: 2021-05-09). [Online]. Available: https://math.stackexchange.com/q/3678498