

# $\nabla$ SLAM: Automagically differentiable SLAM

Krishna Murthy Jatavallabhula<sup>1</sup>, Ganesh Iyer<sup>\*3</sup>, Soroush Saryazdi<sup>\*4</sup>, and Liam Paull<sup>†1,2</sup>

<sup>1</sup>Université de Montréal, Mila, Robotics and Embodied AI Lab (REAL), <sup>2</sup>Candian CIFAR AI Chair, <sup>3</sup>Carnegie Mellon University, <sup>4</sup>Concordia University

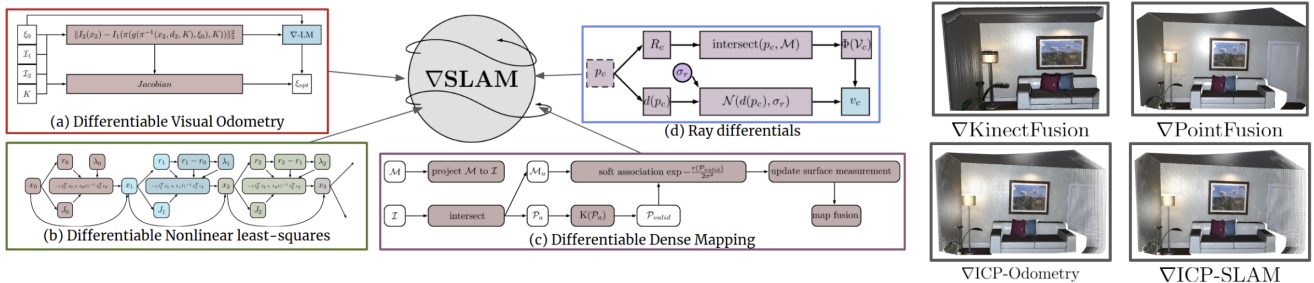


Figure 1.  $\nabla$ SLAM (gradSLAM) is a *fully differentiable* dense simultaneous localization and mapping (SLAM) system. The central idea of  $\nabla$ SLAM is to construct a computational graph representing every operation in a dense SLAM system. We propose differentiable alternatives to several non-differentiable components of traditional dense SLAM systems, such as optimization, odometry estimation, raycasting, and map fusion. This creates a pathway for gradient-flow from 3D map elements to sensor observations (e.g., *pixels*). We implement differentiable variants of three dense SLAM systems that operate on voxels, surfels, and pointclouds.  $\nabla$ SLAM thus is a novel paradigm to integrate representation learning approaches with classical SLAM.

## Abstract

The question of “representation” is central in the context of dense simultaneous localization and mapping (SLAM). While learning-based approaches have the potential to leverage downstream tasks to learn representations, blending such approaches with “classical” SLAM systems has remained an open question. In this work, we propose  $\nabla$ SLAM (gradSLAM), a methodology for posing SLAM systems as differentiable computational graphs, which unifies gradient-based learning and SLAM. We propose differentiable trust-region optimizers, surface measurement and fusion schemes, and raycasting, without sacrificing accuracy. This amalgamation of dense SLAM with computational graphs enables us to backprop all the way from 3D maps to 2D pixels, opening up new possibilities in gradient-based learning for SLAM. A short video explaining the paper and showcasing the results can be found [here](#).

## 1. Introduction

For decades, simultaneous localization and mapping (SLAM) has been a central element of robot perception and

state estimation. SLAM allows robots to operate in previously unseen environments, a core capability that robots must possess for real-world deployment. A large portion of the visual SLAM literature has focused either directly or indirectly on the question of *map representation* [8, 10, 11]. This fundamental choice of representation dramatically impacts the design of processing blocks in the SLAM pipeline, as well as all other downstream tasks that depend on the output of the SLAM system. In particular, for *dense 3D maps* generated from RGB-D cameras, there has been a lack of consensus on the *right* map representation.

Learning representations, as has been done in several other domains is thus appealing for SLAM. However, it is not straightforward because SLAM systems are composed of several subsystems (tracking, mapping, global optimization, etc.) many of which are not inherently differentiable. In this paper, we argue for a declarative approach [7] and propose  $\nabla$ SLAM (gradSLAM): a fully differentiable dense SLAM system that harnesses the power of computational graphs and automatic differentiation to enable learning of dense geometric representations.

This also allows us to solve the *inverse mapping* problem (i.e., answer the question: “How much does a specific pixel-measurement contribute to the resulting 3D map?”) something that is not possible with other popular dense visual SLAM systems [8, 11, 15]. The computational graph frame-

\* Authors contributed equally

† No  $\nabla$ students were harmed in the production of this manuscript.

work allows us to obtain a function ( $\mathcal{S}$ ) that relates a pixel in an RGB-D image (or in general, any sensor measurement  $s$ ) to a 3D geometric map  $\mathcal{M}$  of the environment:  $\mathcal{M} = \mathcal{S}(s)$ . As a result, the *gradient*  $\nabla_s \mathcal{S}$  tells us that perturbing the sensor measurement  $s$  by an infinitesimal  $\delta s$  causes the map  $\mathcal{M}$  to change by  $\nabla_s \mathcal{S}(s) \delta s$ .

If an entire SLAM system can be composed from elementary operations, all of which are differentiable, the system allows end-to-end gradient propagation by construction. However, modern *dense* SLAM systems are quite sophisticated, with several non-differentiable subsystems (optimizers, raycasting, surface mapping), that make such a construction challenging. We show how *all* non-differentiable functions in SLAM can be realised as smooth mappings. In particular, we present differentiable versions of nonlinear least squares optimization, raycasting, and rasterization, while maintaining commensurate accuracy with the non-differentiable counterparts.

While there has been progress in deep learning based SLAM systems [3], and declarative sub-components for SLAM [2, 6], to the best of our knowledge, there is no *single* approach that models the entire SLAM pipeline as a differentiable model, and this is the motivation that underlies  $\nabla$ SLAM.

We demonstrate  $\nabla$ SLAM through 3 instantiations, where our differentiable SLAM building blocks are used to realize the following dense SLAM systems: implicit-surface mapping (KinectFusion [11]), surfel-based mapping (PointFusion [8]), and iterative closest point (ICP) mapping (ICP-SLAM). We also demonstrate examples of back-propagating error signals through the entire SLAM system that enable exciting avenues in representation learning for SLAM<sup>1</sup>.

## 2. $\nabla$ SLAM

### 2.1. Overview of $\nabla$ SLAM

The objective of  $\nabla$ SLAM is to make every computation in SLAM exactly realised as a composition of differentiable functions. Wherever exact differentiable realizations are not possible, we desire *as-exact-as-possible* differentiable realizations. Specifically, we propose differentiable alternatives to trust-region optimization (Sec. 2.2), dense mapping (Sec. 2.3), measurement fusion and raycasting (Sec. 2.4).

### 2.2. $\nabla$ LM: A Differentiable Nonlinear Least Squares Solver

Most state-of-the-art SLAM solutions optimize (minimize) nonlinear least squares objectives to obtain local/globally consistent estimates of the robot state and the map. Such objectives are of the form  $\frac{1}{2} \sum \mathbf{r}(\mathbf{x})^2$ , where  $\mathbf{r}(\mathbf{x})$  is a nonlinear function of residuals.

*Trust-region* methods such as Levenberg-Marquardt (LM) are commonly used to optimize these objectives. These methods are not differentiable as at each optimization step, they involve recalibration of optimizer parameters, based on a *lookahead* operation over subsequent iterates [9]. Specifically, after a new iterate is computed, LM solvers need to make a *discrete* decision between damping or undamping the linear system. Furthermore, when undamping, the iterate must be restored to its previous value. This discrete *switching* behavior of LM does not allow for smooth gradients with respect to the optimizer hyperparameters.

We propose a computationally efficient *soft reparametrization* of the damping mechanism to enable differentiability in LM solvers. Our key insight is that, if  $\mathbf{r}_0 = \mathbf{r}(\mathbf{x}_0)^T \mathbf{r}(\mathbf{x}_0)$  is the (squared) norm of the error at the current iterate, and  $\mathbf{r}_1 = \mathbf{r}(\mathbf{x}_1)^T \mathbf{r}(\mathbf{x}_1)$  is the norm of the error at the *lookahead* iterate, the value of  $\mathbf{r}_1 - \mathbf{r}_0$  determines whether to damp or to undamp. And, only when we choose to undamp, we revert to the current iterate. We define two smooth *gating* functions  $Q_x$  and  $Q_\lambda$  based on the generalized logistic function [12] to update the iterate and determine the next damping coefficient.

$$\begin{aligned} \lambda_1 &= Q_\lambda(r_0, r_1) = \lambda_{min} + \frac{\lambda_{max} - \lambda_{min}}{1 + D e^{-\sigma(r_1 - r_0)}} \\ Q_x(r_0, r_1) &= x_0 + \frac{\delta x_0}{1 + e^{-(r_1 - r_0)}} \end{aligned} \quad (1)$$

where  $D$  and  $\sigma$  are tunable hyperparameters that control the slope of the falloff function [12].  $[\lambda_{min}, \lambda_{max}]$  is the range of values the damping function can assume (usually,  $\lambda_{min} = \frac{1}{2}$ ,  $\lambda_{max} = 2$ , when using multiplicative damping with a damping coefficient of 2). This smooth parameterization of the LM update allows the optimizer to be expressed as a fully differentiable computational graph (Fig. 1(b)).

$\nabla$ LM allows differentiable realizations of several SLAM components, such as dense visual odometry estimation [13] (Fig. 1(a)), and Iterative Closest Point (ICP) alignment [1].

### 2.3. Differentiable Mapping

The *mapping* process in SLAM involves a number of differentiable yet non-smooth operations (clipping, indexing, thresholding, new/old decision, active/inactive decision, etc.). We enforce differentiability in this process by the following corrective measures.

1. The surface measurement made at each valid pixel  $p$  in the live frame is a function of  $p$  and also active neighbours of  $p$ ,  $nb(p)$  via a *kernel*  $K(p, nb(p))$ .
2. When a surface measurement is transformed to the global frame, we use *soft* (one-many) rather than *hard* (one-one) associations.
3. Every surface measurement is, by default, assumed to represent a new map element, which is passed to a *dif*

<sup>1</sup>Project page: <http://montrealrobotics.ca/gradSLAM/>



Figure 2. **Qualitative results** on sequences from the ScanNet [4] dataset. We use each of the differentiable SLAM systems ( $\nabla$ KinectFusion,  $\nabla$ PointFusion, and  $\nabla$ ICP-SLAM) to reconstruct parts of the scene. We also show outputs from BundleFusion [5].

*ferentiable fusion step (c.f. Sec 2.4).*

The kernel  $K(p, nbd(p))$  can be a discrete approximation (e.g., constant within a pixel) or can vary at the subpixel level. For faster computation and coarse gradients, we use a bilinear interpolation kernel. While bilinear interpolation is a sensible approximation for image pixels, this is often a poor choice for use in 3D *soft* associations. For forming 3D associations, we leverage characteristics of RGB-D sensors in defining the soft falloff functions. Specifically, we compute, for each point  $P$  in the live surface measurement, a set of closest candidate points in a region  $exp\left(-\frac{r(P)^2}{2\sigma^2}\right)$ , where  $r(P)$  is the radial depth of the point from the camera ray, and  $\sigma$  affects the falloff region. The computational graph for this differentiable mapping process is illustrated in Fig. 1(c).

## 2.4. Other differentiable modules

We also repurpose existing differentiable modules for measurement fusion [8, 11] and raycasting [14].

## 3. Case Studies: KinectFusion, PointFusion, and ICP-SLAM

As concrete demonstrations of  $\nabla$ SLAM, we leverage the aforementioned differentiable SLAM subsystems and compose them to realise three practical SLAM solutions. We implement differentiable versions of the *KinectFusion* [11] algorithm that constructs TSDF-based volumetric maps, the *PointFusion* [8] algorithm that constructs surfel maps, and a pointcloud-based SLAM framework that we call *ICP-SLAM*. Fig. 2 shows qualitative reconstruction results on sequences from ScanNet [4], and Table 1 shows that the accuracy of these differentiable SLAM systems is similar to their non-differentiable counterparts.

Method	ATE	RPE
ICP-Odometry (non-differentiable)	0.029	0.0318
$\nabla$ ICP-Odometry	<b>0.01664</b>	<b>0.0237</b>
ICP-SLAM (non-differentiable)	0.0282	0.0294
$\nabla$ ICP-SLAM	<b>0.01660</b>	<b>0.0204</b>
PointFusion (non-differentiable)	<b>0.0071</b>	<b>0.0099</b>
$\nabla$ PointFusion	0.0072	0.0101
KinectFusion (non-differentiable)	<b>0.013</b>	<b>0.019</b>
$\nabla$ KinectFusion	0.016	0.021

Table 1. **Performance of  $\nabla$ SLAM** compared to non-differentiable counterparts (ATE: Absolute Trajectory Error, RPE: Relative Pose Error).

Metric	LM	$\nabla$ LM
Convergence iters	7.6 $\pm$ 3.4	7.4 $\pm$ 3.8
$\ a_{opt} - a_{init}\ _1$	<b>0.011</b> $\pm$ 0.011	0.399 $\pm$ 0.162
$\ b_{opt} - b_{init}\ _1$	<b>0.331</b> $\pm$ 0.174	0.574 $\pm$ 0.137
$\ c_{opt} - c_{init}\ _1$	0.114 $\pm$ 0.021	<b>0.084</b> $\pm$ 0.032
Total error	<b>0.184</b>	0.207

Table 2.  $\nabla$ LM performs quite similarly to its non-differentiable counterpart. (Convergence tolerance  $10^{-6}$ )

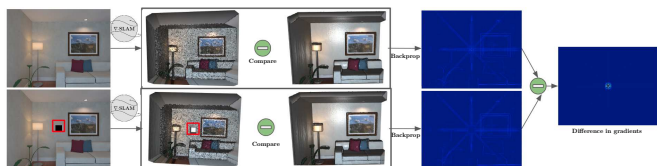


Figure 3. **Analysis of gradients:**  $\nabla$ SLAM enables gradients to flow through to the input images. *Top:* An RGB-D image pair (depth not shown) is passed through  $\nabla$ KinectFusion. The resulting map is compared with a precise (ground-truth) map. The comparison error is backpropagated through the SLAM system, to the depth map (blue colormap). *Bottom:* An occluder is added to the center of the RGB-D pair. This occluder results in a gaping hole. But, using the backpropagated gradients, one can identify the set of image/depthmap pixels that result in the reconstruction error.

## 4. Experiments and results

### 4.1. Differentiable Optimization

We design a test suite of nonlinear curve fitting problems to measure the performance of  $\nabla$ LM (Sec 2.2) to its non-differentiable counterpart. We uniformly sample the parameters  $p = a, b, c$ , with initial guess  $a_0, b_0, c_0$  from the exponential family:  $p \sim y = aexp\left(-\frac{(x-b)^2}{2c^2}\right)$ . For 1000 sampled problem sets, we optimize using both LM and  $\nabla$ LM, and measure the following quantities: iterations to converge, quality of the solution (i.e., discrepancy between estimated and true parameters). Notice from Table 2 how  $\nabla$ LM performs similarly to LM (a slight performance drop is noticeable, due to smoothing).

### 4.2. Analysis of Gradients

The computational graph approach of  $\nabla$ SLAM allows us to recover meaningful gradients of 2D (or 2.5D) measurements with respect to a 3D surface reconstruction. We provide an analysis of what these multi-view gradients correlate to in the input image and depth space. In Fig. 3, the top row shows an RGB-D image differentially transformed—



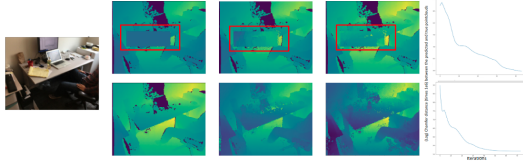


Figure 4. **End-to-end gradient propagation:** (Top): A chunk of a depth map is *chopped*. The resultant sequence is reconstructed using  $\nabla$ PointFusion and the pointcloud is compared to a *clean* one reconstructed using the unmodified depth map. The Chamfer distance between these two pointclouds is used to define a reconstruction error between the two pointclouds, which is backpropagated through to the input depth map and updated by gradient descent. (Bottom): Similar to the Fig. 3, we show that  $\nabla$ SLAM can *fill-in* holes in the depthmap by leveraging multi-view gradient information.

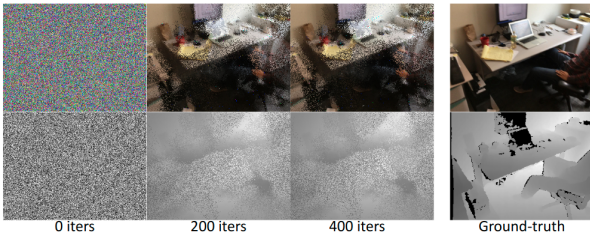


Figure 5. **RGB-D completion using end-to-end gradient propagation:** Three RGB-D images and a *noise image* are passed through  $\nabla$ PointFusion, and compared to a clean reconstruction obtained from four RGB-D images. The reconstruction loss is used to optimize the *noise image* by gradient descent. We can recover most of the artifacts from the raw RGB and depth images. Note that finer features are hard to recover from a random initialization, as the overall *SLAM function* is only locally differentiable.

using  $\nabla$ SLAM—into a (noisy) TSDF surface measurement, and then compared to a more precise global TSDF map. The bottom row is similarly transformed, with the difference being the presence of a small ( $40 \times 40$  px) occluder. Element-wise comparison of aligned volumes gives us a reconstruction error, whose gradients are backpropagated through to the input depthmap using the computational graph maintained by  $\nabla$ SLAM. Inspecting the gradients with respect to the input indicates the per pixel contribution of the occluding surface to the volumetric error. In Fig. 4, we similarly introduce such occluders (top row) and pixel noise (bottom row) in one of the depth maps of a sequence and reconstruct the scene using  $\nabla$ PointFusion. We then calculate the Chamfer distance between the noisy and true surfel maps and backpropagate the error with respect to each pixel. The minimized loss leads to the targeted recovery of the noisy and occluded regions. We additionally show an RGB-D image completion task (from uniform noise) in Fig. 5.

Thus,  $\nabla$ SLAM provides a rich interpretation of the computed gradients: they denote the contribution of each pixel towards the eventual 3D reconstruction.

## 5. Conclusion

We introduce  $\nabla$ SLAM, a declarative, computational graph approach to SLAM.  $\nabla$ SLAM enables gradient-based learning for localization and mapping based tasks by providing explicit gradients with respect to the input sensor observations. We showcase how the gradients propagate through the tracking, mapping, and fusion stages. Future efforts will focus on learning  $\nabla$ SLAM components to optimize downstream task performance.  $\nabla$ SLAM potentially enables a variety of self-supervised learning applications, by equipping gradient-based learning architectures with *spatial understanding*.

## References

- [1] Paul J Bes, Neil D McKay, et al. A method for registration of 3-d shapes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 14(2):239–256, 1992. 2
- [2] Bert De Brabandere, Wouter Van Gansbeke, Davy Neven, Marc Proesmans, and Luc Van Gool. End-to-end lane detection through differentiable least-squares fitting. *arXiv*, 1902.00293, 2019. 2
- [3] Jan Czarnowski, Tristan Laidlow, Ronald Clark, and Andrew Davison. Deepfactors: Real-time probabilistic dense monocular slam. In *IEEE RAL*, 2020. 2
- [4] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, 2017. 3
- [5] Angela Dai, Matthias Nießner, Michael Zollöfer, Shahram Izadi, and Christian Theobalt. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface re-integration. *ACM Transactions on Graphics*, 2017. 3
- [6] Ravi Garg, Vijay Kumar BG, Gustavo Carneiro, and Ian Reid. Un-supervised cnn for single view depth estimation: Geometry to the rescue. In *ECCV*, 2016. 2
- [7] Stephen Gould, Richard Hartley, and Dylan Campbell. Deep declarative networks: A new hope. *arXiv:1909.04866*, 2019. 1
- [8] Maik Keller, Damien Lefloch, Martin Lambers, Shahram Izadi, Tim Weyrich, and Andreas Kolb. Real-time 3d reconstruction in dynamic scenes using point-based fusion. In *International Conference on 3D Vision*, 2013. 1, 2, 3
- [9] Michael Lampton. Damping–undamping strategies for the levenberg–marquardt nonlinear least-squares method. *Computers in Physics*, 1997. 2
- [10] Raúl Mur-Artal and Juan D. Tardós. ORB-SLAM2: an open-source SLAM system for monocular, stereo and RGB-D cameras. *IEEE TRO*, 33(5):1255–1262, 2017. 1
- [11] Richard A Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J Davison, Pushmeet Kohli, Jamie Shotton, Steve Hodges, and Andrew W Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *ISMAR*, 2011. 1, 2, 3
- [12] FJ Richards. A flexible growth function for empirical use. *Journal of experimental Botany*, 1959. 2
- [13] Frank Steinbrücker, Jürgen Sturm, and Daniel Cremers. Real-time visual odometry from dense rgb-d images. In *ICCV Workshops*, 2011. 2
- [14] Shubham Tulsiani, Tinghui Zhou, Alexei A. Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *CVPR*, 2017. 3
- [15] Thomas Whelan, Michael Kaess, Hordur Johannsson, Maurice Fallon, John J Leonard, and John McDonald. Real-time large-scale dense rgb-d slam with volumetric fusion. *IJRR*, 2015. 1