Theseus: A Library for Differentiable Nonlinear Optimization

Anonymous Author(s) Affiliation Address email

Abstract

We present Theseus, an efficient application-agnostic open source library for differ-1 entiable nonlinear least squares (DNLS) optimization built on PyTorch, providing 2 a common framework for end-to-end structured learning in robotics and vision. З Existing DNLS implementations are application specific and do not always incor-4 porate many ingredients important for efficiency. Theseus is application-agnostic, 5 as we illustrate with several example applications that are built using the same 6 underlying differentiable components, such as second-order optimizers, standard 7 costs functions, and Lie groups. For efficiency, Theseus incorporates support for 8 sparse solvers, automatic vectorization, batching, GPU acceleration, and gradient 9 10 computation with implicit differentiation and direct loss minimization. We do extensive performance evaluation in a set of applications, demonstrating significant 11 efficiency gains and better scalability when these features are incorporated. 12

13 **1 Introduction**

Reconciling traditional approaches with deep learning to leverage their complementary strengths is a common thread in a large body of recent work in robotics. In particular, an emerging trend is to differentiate through nonlinear least squares [1] which is a second-order optimization formulation at the heart of many problems in robotics [2–7] and vision [8–13]. Optimization layers as inductive priors in neural models have been explored in machine learning with convex optimization [14, 15] and in meta learning with gradient descent [16, 17] based first-order optimization.

Differentiable nonlinear least squares provides a general scheme to encode inductive priors, as the objective function can be partly parameterized by neural models and partly with engineered domain-specific differentiable models. Here, input tensors define a sum of weighted squares objective function and output tensors are minima of that objective. In contrast, typical neural layers take input

tensors through a linear transformation and some element-wise nonlinear activation function.

The ability to compute gradients end-to-end is retained by differentiating through the optimizer which 25 allows neural models to train on the final task loss, while also taking advantage of priors captured 26 by the optimizer. The flexibility of such a scheme has led to promising state-of-the-art results in a 27 wide range of applications such as structure from motion [18], motion planning [19], SLAM [20, 21], 28 bundle adjustment [22], state estimation [23], image alignment [24] with many other applications like 29 manipulation and tactile sensing [25, 26], control [27], human pose tracking [28, 29] to be explored. 30 However, existing implementations from above are application specific, common underlying tools 31 like optimizers get reimplemented, and features like sparse solvers, batching, and GPU support that 32 impact efficiency are not always included. This has led to a fragmented literature where it is difficult 33 to start working on a new idea or to build on the progress of prior work. 34 To address this gap, we present Theseus, an open source library for differentiable nonlinear least 35

36 squares optimization built on PyTorch. Theseus provides an efficient application-agnostic interface

that consolidates recent efforts and catalyzes future progress in the domain of structured end-to-end

38 learning for robotics and vision. Our contributions are summarized below.

Application agnostic interface. Our implementation provides an easy to use interface to build 39 custom optimization layers and plug them into any neural architecture. (i) The layer can be constructed 40 from a set of available second-order optimizers like Gauss-Newton and Levenberg-Marguardt and 41 a nonlinear least squares objective. (ii) The objective can be constructed with learnable or hand-42 specified cost functions, either by applying one of many common costs already provided in the library, 43 or by building custom costs in-place with support for automatic differentiation through PyTorch [30]. 44 (iii) We also provide differentiable Lie groups for representing 2D/3D positions and rotations [31], 45 and differentiable kinematics wrapping over an existing library [32] for representing robot models. 46 More details are described in Sec. 3. 47 Efficiency based design. Efficiency is a central design consideration and we make several advance-48

ments in improving computation times and memory consumption. (i) As common in prior work, an 49 optimizer implementation using PyTorch's native linear solver would use a dense representation for 50 solving the linear system within the nonlinear optimization. In practice, these optimization problems 51 often have a considerable amount of sparsity that can be exploited [33–36]. In Theseus, we imple-52 ment sparse linear solvers that are differentiable end-to-end and make them efficient with custom 53 CPU and CUDA backends to support batching and GPU acceleration. (ii) Beyond sparse solvers, we 54 extend batching and GPU support to all features in the library and add automatic vectorization of cost 55 functions and other operations to significantly boost efficiency. (iii) Finally, we introduce implicit 56 differentiation [37] and direct loss minimization [38, 39], which have been previously applied to only 57 first order optimizers like gradient descent and convex optimization, to a new class of second-order 58 optimizers. This goes beyond prior work with nonlinear least squares that currently only support 59 differentiation with standard unrolling, which is known to have challenges with compute, memory, 60 and vanishing gradients. More details are described in Sec. 4. 61

Highlights of results. Together, the application-agnostic features let users easily set up a variety of 62 problems like pose graph optimization, tactile state estimation, bundle adjustment, motion planning, 63 and homography estimation, all of which are included as examples in the open source code and 64 described in Sec. 3.1. In evaluations, we find that on a standard GPU, Theseus with a sparse solver is 65 much faster and requires significantly less memory than a dense solver, and when solving a batch of 66 large problems the forward pass of Theseus is up to 4x faster than state-of-the-art C++ based solver 67 Ceres that has limited GPU support and does not support batching and end-to-end learning. We also 68 compare all backward modes to find that with increasing number of optimization iterations, compute 69 and memory increases linearly for unrolling and stays constant for implicit differentiation, while the 70 latter also provides better gradients. More details are described in Sec. 5. 71

72 **2 Background and related work**

73 Nonlinear least squares (NLS) is an optimization problem [1] that finds optimization variables θ

$$\theta^{\star} = \arg\min_{\theta} S(\theta), \qquad S(\theta) = \frac{1}{2} \sum_{i} ||r_{i}(\theta^{i})||^{2} = \frac{1}{2} \sum_{i} ||w_{i}c_{i}(\theta^{i})||^{2}$$
(1)

where the objective $S(\theta)$ is a sum of squared vector-valued residual terms r_i , each a function of $\theta^i \in \theta$, 74 a subset of the optimization variables $\theta = \{\theta_i\}$. Any variable θ_i is a manifold object, for example 75 a Euclidean vector or a matrix Lie group. For flexibility, we represent a residual $r_i(\theta^i) = w_i c_i(\theta^i)$ 76 as a product of a matrix weight w_i and vector cost c_i . Robotics and vision have used this general 77 optimization formulation to tackle many applications [3, 4]. For example, costs capture sensor 78 measurement errors and physical constraints to optimize camera, robot, object, or human poses in 79 estimation and tracking problems like simultaneous localization and mapping (SLAM) [40], structure 80 from motion [13], bundle adjustment [8], visual inertial odometry [2], articulated tracking [12], 81 contact odometry in legged locomotion [25], 3D pose and shape reconstruction of humans [28, 29] or 82 objects [10]. Similarly, costs can also capture constraints and desired future goals to find robot states 83 or actions in motion planning [5], dynamics [6], and control [27] problems. 84

Solving NLS. Problems represented by Eq. (1) are solved by iteratively linearizing the nonlinear 85 objective around the current variables to get the linear system $(\sum_i J_i^{\top} J_i)\delta\theta = (\sum_i J_i^{\top} r_i)$, then 86 solving the linear system to find the update $\delta\theta$, and finally updating the variables $\theta \leftarrow \theta - \delta\theta$, until 87 convergence. Note that in the update the minus operation is more generally a retraction mapping for 88 non-Euclidean variables. In the linear system, $J_i = [\partial r_i / \partial \theta^i]$ are the Jacobians of residuals with 89 respect to the variables and the iterative method above, called Gauss-Newton (GN), is a nonlinear 90 optimizer that is second-order, since $H = (\sum_i J_i^{\top} J_i)$ represents the approximate Hessian. To 91 improve robustness and convergence, variations like Levenberg-Marquardt (LM) damp the linear 92 system, while others use a trust region and adjust step size for the update with line search (e.g., 93

Dogleg). Please refer to [1, 41] for an in depth exploration. In most applications discussed above the
objective structure gives rise to a sparse Hessian since not all costs depend on all variables. Several
general purpose frameworks [33–36] have been built that leverage this sparsity property to efficiently
solve the sparse linear system in every iteration of the nonlinear optimization. While these framework
were not built for deep learning, they are highly efficient and performant on CPU. **NLS with learning.** Data driven learning has been explored to address challenges in hand crafting
costs or features for costs, finding weights to balance different costs, or to find initializations that lead
to better convergence. Some examples include, learning object shape code [42] or environment depth

to better convergence. Some examples include, learning object shape code [42] or environment depth
code [43] for SLAM [44], learning motion priors for planning to manipulate articulated objects [45],
learning relative pose from tactile images to estimate object state during pushing [26], and semantic
2D segmentation fused in 3D mesh for semantic SLAM [7]. These approaches only train features on
a surrogate or intermediate loss and then apply optimization at inference where the true downstream
task loss is available but not utilized. To take full advantage of end-to-end learning, latest approaches
thus are redesigning the optimization to be differentiable.

Differentiable NLS (DNLS) solves the optimization in Eq. (1) and also provides gradients of the solution θ^* with respect to any upstream neural model parameters ϕ that parameterize the objective $S(\theta; \phi)$ and in turn any costs $c_i(\theta^i; \phi)$, weights $w_i(\phi)$, or initialization for variables $\theta_{init}(\phi)$. The goal is to learn these parameters ϕ end-to-end with a downstream learning objective L defined as a function of θ^* . This results in a bilevel optimization setup

inner loop:
$$\theta^{\star}(\phi) = \arg \min_{\theta} S(\theta; \phi),$$
 outer loop: $\phi = \arg \min_{\phi} L(\theta^{\star}(\phi))$ (2)

where the inner loop is second-order DNLS and the outer loop is first-order gradient descent class of optimization that is standard in deep learning. The outer loop performs update $\phi \leftarrow \phi + \delta \phi$ by computing $\delta \phi$ using gradients $\partial \theta^* / \partial \phi$ through inner loop DNLS. Note that more generally the learning objective i.e. outer loss *L* can also depend on other quantities like neural model parameters downstream of θ^* , but we omit them here for clarity.

Recent works with DNLS have outperformed optimization only or learning only methods by com-118 bining the strengths of classical methods with deep learning. For example, learning features for 119 costs to represent depth in bundle adjustment [22] and monocular stereo [46] where an initialization 120 network also learns to predict depth and pose, learning cost weights like motion model weights in 121 video to depth estimation [18], obstacle avoidance weights in 2D motion planning from occupancy im-122 ages [19], learning robust loss weights in image alignment [24] and state-of-the-art dense SLAM [21], 123 and confidence weights for feature matching to optimize camera pose [47]. Other works like, [20] 124 backpropagate reconstruction error to sensor model in a SLAM system, [48] solves large scale bundle 125 adjustment on a GPU, and learn sensor and dynamics models for 2D visual object tracking and visual 126 odometry [23]. These implementations however, are application specific which has led repeated 127 work in building DNLS where features like learnable costs and weights, Lie groups, and kinematics 128 are not always present. Additionally, features that have a significant impact on performance, like 129 sparsity and vectorization of costs are only considered by some [23, 48, 49] or in the case of implicit 130 differentiation for second-order optimization, have not yet been explored. 131

3 Application agnostic interface

Given the lack of a common and efficient framework for DNLS an important goal of Theseus is to provide an application-agnostic interface. In this section, we describe how we enable this with an easy-to-use core API, standard cost functions, and features like Lie groups and kinematics, and illustrate several examples using this interface. We discuss design for efficiency in the next section.

The core API lets users focus on describing the DNLS problem and their interaction with the outer loss L and parameters ϕ within any broader PyTorch model, while the solution and differentiation are seamlessly taken care of under-the-hood. The basic components of the core API are described below with the help of a simple example in Listing 1 (see App. A for more details on the example):

- Variable: refers to either *optimization variables*, θ , or *auxiliary variables* (those constant with
- respect to S, e.g., parameters ϕ or data tensors), which are named wrappers of torch batched tensors stored in Variable.data (lines 3-5).
- CostFunction: defines costs c_i (lines 12-14) and are also responsible for declaring which of its variables are optimization and which are auxiliary (lines 8-9),
- CostWeight: defines weights w_i associated with cost c_i (line 14).
- Objective: defines $S(\theta; \phi)$, and thus the structure of an optimization problem (lines 11, 15) by
- holding all cost functions and weights, and their associated variables. These are implicitly obtained

```
x_true, y_true, v_true = read_data() # shapes (1, N), (1, N), (1, 1)
1
2
3
     x = th.Variable(torch.randn_like(x_true), name="x")
     y = th.Variable(y_true, name="y")
4
     v = th.Vector(1, name="v") # a manifold subclass of Variable for optim_vars
5
6
     def error_fn(optim_vars, aux_vars): # returns y - v * exp(x)
7
         x, y = aux_vars
8
         return y.data - optim_vars[0].data * torch.exp(x.data)
9
10
11
     objective = th.Objective()
     cost_function = th.AutoDiffCostFunction(
12
         [v], error_fn, y_true.shape[1], aux_vars=[x, y],
13
         cost_weight=th.ScaleCostWeight(1.0))
14
     objective.add(cost_function)
15
     layer = th.TheseusLayer(th.GaussNewton(objective, max_iterations=10))
16
17
     phi = torch.nn.Parameter(x_true + 0.1 * torch.ones_like(x_true))
18
     outer_optimizer = torch.optim.RMSprop([phi], lr=1e-3)
19
20
     for epoch in range(10):
         solution, info = layer.forward(
21
             input_data={"x": phi.clone(), "v": torch.ones(1, 1)},
22
             optimizer_kwargs={"backward_mode": th.BackwardMode.IMPLICIT})
23
         outer_loss = torch.nn.functional.mse_loss(solution["v"], v_true)
24
         outer_loss.backward()
25
         outer_optimizer.step()
26
```

Listing 1: Simple DNLS example with Theseus, see App. A for details.

when a CostFunction is added to the Objective, and Variable names are used to infer whichare shared by one or more CostFunction.

• Optimizer: is the inner loop optimization algorithm (e.g. Gauss-Newton) that finds the solution θ^* given objective S (line 16).

TheseusLayer: encapsulates an optimizer and objective, and serves as the interface between the
 DNLS block and other torch modules upstream or downstream (line 16).

The interface between the inner loop optimization and the outer loop's parameters and loss occurs via 155 TheseusLayer. forward (lines 21-23). This receives as input a dictionary mapping variable names 156 to torch tensors, which Theseus then uses to populate the corresponding Variable with the tensor 157 mapped to its name. With the input dictionary users can provide initial values for the optimization 158 variables, data tensors, or current values for parameters ϕ before running the inner loop optimization. 159 The output of forward is another dictionary that maps variable names to tensors with their optimal 160 values found in the inner loop (lines 21, 24); auxiliary variables are not modified during the forward 161 pass. The output tensors can then be combined with other torch modules downstream to compute L. 162 while maintaining the full differentiable computation graph (lines 24-26). 163

We currently provide Gauss-Newton and Levenberg-Marquardt as Optimizer for the inner 164 loop, with the ability to easily add support for more optimizers in the future. Listing 1 uses 165 AutoDiffCostFunction to construct an in-place CostFunction (line 12) which allows automati-166 cally calculating Jacobians J_i with PyTorch. Beyond this, in the library we include standard cost 167 functions with analytical Jacobians broadly used in many applications, like Gaussian measurements, 168 reprojection error, relative pose measurement, motion models, and collision costs. We also include 169 a variety of robust loss functions, useful for example in handling outliers [50], which can be easily 170 integrated with CostFunction. Next we describe support for Lie groups and kinematics. 171

Differentiable Lie groups Lie groups are widely used in robotics and vision to represent 2D/3D 172 positions and rotations [31]. Due to their non-Euclidean geometry, it is difficult to apply them to deep 173 learning, which primarily operates with Euclidean tensors, but recently there is growing interest in 174 making them compatible [23, 51–55]. LieTorch [52] generalizes automatic differentiation on the Lie 175 group tangent space through local parameterization around the identity, but the implementation is 176 complex since every operation requires a custom kernel. In contrast, Theseus computes common Lie 177 group operators, e.g., the exponential and logarithm map, inverse, composition, etc., in closed form, 178 and provides their corresponding analytical derivatives on the tangent space. Following [56], we also 179 implement a projection operator that allows us to project gradients computed by PyTorch's autodiff to 180

the tangent space and use them to easily compute Jacobians and update Lie group variables correctly; a similar strategy has also been implemented in [57]. Additionally, our Lie group implementation includes a heuristic extension that allows using any of PyTorch's first-order optimizers on non-Euclidean manifolds with minimal code changes. All of these make it easy and straightforward to

run optimization and train neural networks with Lie groups variables. More details in App. B.

Differentiable kinematics Many problems such as motion planning or state estimation on high degree of freedom robots like arms or mobile manipulators, involve computation of robot kinematics for collision avoidance or computing distance of end effector to goal. Theseus provides a differentiable implementation of forward kinematics by wrapping over Differentiable Robot Model [32], which builds a differentiable kinematics function from a standard robot model file. Gradients are computed through autodiff, while we also provide a more efficient, analytical manipulator Jacobian. This module can be used within any CostFunction in Theseus.

193 3.1 Example applications

To illustrate the versatility of Theseus, we include a number of example DNLS applications below with more details in App. C. Crucially, to implement these with Theseus, most of the effort is only in defining application-specific components such as data management, neural models, or custom CostFunction. With these defined, putting the full DNLS block together is a few lines of code to setup a TheseusLayer and an outer loop similar to the simple example in Listing 1.

Pose graph optimization (PGO) estimates poses from their noisy relative measurements [58]. With
 DNLS we learn the radius of a Welsh robust cost function for outlier rejection, using the difference
 between estimated and ground truth poses as the outer loss on a synthetic dataset.

Tactile state estimation follows [26], which estimates 2D poses of an object pushed by a robot hand with an image-based tactile sensor [59]. A neural network that predicts relative pose between hand and object from tactile images is learned end-to-end through the TheseusLayer.

Bundle adjustment is the problem of optimizing a 3D reconstruction formed by a set of camera images and a set of landmarks observed and matched across the images [60]. We learn the radius of a soft-kernel that penalizes outlier observations, using the average frame pose error as outer loss.

Motion planning considers a differentiable version of the GPMP2 planning algorithm, inspired by [19], where the outer loss tries to match expert demonstrations. Here we learn a model for initializing optimization variables, and we include the inner loop objective as a term in the outer loss.

Homography estimation. Homography is a linear transformation between corresponding points in two images and can be solved by minimising a dense photometric loss. Robustness to lighting and viewpoint change can be improved with a feature-metric loss based on CNN features [61–66]. In our outer loop, we train a CNN to produce robust features for image alignment.

215 4 Efficiency based design

Theseus enables several different applications with a general interface. Compute and memory efficiency are central to making its usage practical. Next, we explain design considerations to support batching and vectorization, sparsity, and backward modes for differentiation which we demonstrate boost performance in evaluations section.

220 4.1 Batching and vectorization

Parallel processing is important to improve the 221 computational efficiency in machine learning 222 and optimization. In Theseus, we enable two 223 levels of parallelization. First, Theseus natively 224 supports solving a batch of DNLS in parallel. 225 This fits seamlessly in the PyTorch framework, 226 where training and inferences by batches is the 227 default mode of operation, helping reduce the 228 overhead coming from individual operations be-229





ing implemented in Python. Second, inspired by DeepLM [48], and noting that lots of the operations such as costs, gradients and Jacobian computation, and variable updates only differ from each other in terms of the input data, we make use of the single-instruction-multiple-data (SIMD) protocol by automatically detecting and vectorizing operations of the same type, significantly reducing the overhead for forward and backward passes. Using the PGO example, Fig. 1 shows that Theseus achieves significant speedup with automatic vectorization. Note that there is a small application dependent trade-off; here the memory use increases by $\sim 35\%$ for forward and $\sim 10\%$ for backward.

4.2 Handling sparsity with linear solvers beyond PyTorch

Solving NLS requires solving a sequence of linear systems to obtain descent directions. As discussed 238 in Sec. 2, these systems are generally sparse and can be solved much more efficiently if not treated as 239 dense. These includes differentiable sparse solvers that take advantage of the sparsity, comple-240 menting PyTorch's native dense solvers. Importantly, Theseus seamlessly takes care of assembling 241 the cost functions and variables in the objective into sparse data structures that our linear solvers 242 can consume, without any extra burden on the user. Currently, we provide two sparse solvers, one 243 that uses the CPU-based solver CHOLMOD [67], and another, cudaLU, that is based on the cuSolverRF 244 package, itself part of the cuSolver library provided by Nvidia with CUDA. As a bonus feature, we 245 provide access to these solvers as standalone PyTorch functions, so they can, in fact, be used to solve 246 sparse matrices arising outside of NLS or DNLS optimization. 247

CHOLMOD-based solver. CHOLMOD [67] achieves state-of-the-art performance on computation of the 248 Cholesky decomposition of sparse matrices. It exploits parallelism by grouping sparse entries to 249 take advantage of high-performance multi-threaded dense matrix operations in BLAS/LAPACK libraries. 250 CHOLMOD has some limited support for GPU for some of its operations, but the algorithm is strongly 251 CPU-based and the user is expected to provide matrix data on the CPU. One convenient feature is 252 computing the symbolic analysis of a sparse matrix pattern as a separate step and creating a symbolic 253 decomposition object that can be used for all subsequent factorizations. We also take advantage of 254 builtin functionality for sparse multiplication and only provide the Jacobian matrix J to solve for the 255 Hessian matrix $H = J^{\top}J$. Two limitations of the library with respect to Theseus are, first, the lack 256 of proper GPU support, which forces us to provide matrix data on the CPU, and, second, the lack of 257 batching, which requires us to loop to solve every problem in the batch independently. On the other 258 hand, since it runs on CPU it has less memory restrictions than GPU-based solvers (see Sec. 5.1). 259

cudaLU solver. cuSolverRF is designed to accelerate the solution of sets of linear systems by fast LU 260 refactorization when given new coefficients for the same sparsity pattern. To take advantage of this, 261 we implemented custom CUDA kernels for batched sparse matrix-matrix and matrix-vector products, 262 and for solving a batch of sparse linear systems using LU factorization from cuSolverRF. Although 263 this solver leads to a substantial performance boost over PyTorch's dense solver (see Sec. 5.1), the 264 closed-source nature of cuSOLVER results in some challenges and limitations: (i) cuSolverRF does 265 not support separate symbolic decomposition and numeric contexts, so it's not possible to use the 266 same symbolic decomposition to hold in memory separate factors. Since this is necessary in Theseus 267 for unrolling of the inner loop, we work around this limitation by creating a pool of contexts, and we 268 use the least recently used context for factorization. As a consequence, the number of contexts must 269 be set according to the number of iterations that need to be unrolled; (ii) The batch size is fixed once a 270 context is created. Since recreating the contexts is an expensive operation, it means that the batch size 271 has to be constant over the course of outer loop optimization; (iii) It relies on LU factorization, which 272 for symmetric matrices (the case of Theseus) is less efficient than using Cholesky decomposition. 273

Backward for custom linear solvers. Obtaining derivatives of the linear system solve with respect 274 to the parameters is crucial operation for DNLS. In particular, we consider optimizing the parameters 275 A and b of a linear system $y = A^{-1}b$ to minimize a downstream function f(y). The derivatives of the 276 loss with respect to the parameters of the linear system can be obtained with implicit differentiation, $\frac{\partial f}{\partial b} = A^{-1} \frac{\partial f}{\partial y}$ and $\frac{\partial f}{\partial A} = -A^{-1} \frac{\partial f}{\partial y} y^{\top}$, as done in Barron and Poole [68]. In Theseus, we implement 277 278 this by connecting the Python interface of our sparse solvers with PyTorch's autograd. Function 279 classes that implement the gradients above in their backward methods. This connects the computation 280 281 graph between the downstream function and any upstream parameters that modify the system via auxiliary variables or values for optimization variables. Furthermore, since the gradients require 282 solving linear systems that use the same matrix as the forward pass, our backward pass can cache 283 factorizations, resulting in much faster backward times compared to dense solvers (see Fig. 2). 284

285 4.3 Backward modes for DNLS

The parameters ϕ upstream of DNLS can be learned end-to-end through the solution $\theta^*(\phi)$ by using the *adjoint derivatives* $\partial \theta^*(\phi) / \partial \phi$. We include four methods for computing them in Theseus.

Unrolling is the standard way in which past work in DNLS has computed the adjoint derivatives. This is often referred to as backpropagation through time or unrolled optimization and is explored in [16, 19, 69–78]. In practice, often only a few steps of unrolling are performed due to challenges with

compute, memory, and vanishing gradients.

Truncated differentiation. Aside from unrolling a few steps, another way of approximating the derivatives is to use truncated backpropagation through time (TBPTT) [79, 80]. Truncation unfortu-

nately results in biased derivatives and many works [81–85] seek to further theoretically understand the properties of TBPTT, including the bias of the estimator and how to unbias it.

Implicit differentiation. If θ^* can be computed exactly, then the implicit function theorem provides a 296 way of computing the adjoint derivatives as done in related work in convex optimization and first-order 297 gradient descent methods [14, 15, 86–91]. We apply the implicit function theorem from Dontchev 298 and Rockafellar [37, Theorem 1B.1] (see App. E) to Eq. (2) to perform implicit differentiation on a 299 new class of second-order NLS optimization. This first requires that we transform Eq. (2) into an 300 implicit function that finds the roots. We do this via the first-order optimality condition, resulting in 301 $q(\theta;\phi) := \nabla_{\theta} S(\theta;\phi)$. Finding $\Theta^{\star}(\phi) := \{\theta \mid q(\theta;\phi) = 0\}$ corresponds to solving Eq. (2). Under 302 mild assumptions, the theorem above gives the adjoint derivative at ϕ 303

$$\mathcal{D}_{\phi}\theta^{\star}(\bar{\phi}) = -\mathcal{D}_{\theta}^{-1}g(\theta^{\star}(\bar{\phi});\bar{\phi})\mathcal{D}_{\phi}g(\theta^{\star}(\bar{\phi});\bar{\phi}).$$
(3)

As Theseus internally uses a (Gauss-)Newton solver, the following proposition provided in App. E

shows how we can use compute Eq. (3) by differentiating a single Newton step at an optimal solution. **Proposition 1.** Differentiating a single Newton iteration h at an optimal and unique θ^* results in the exact computation of the implicit derivative in Eq. (3).

Direct loss minimization. Suppose we have a outer loss as in Eq. (2). The direct loss minimization (DLM) approach uses this loss to augment the inner-loop optimization problem in order to define a finite difference scheme that approaches the true gradient $\nabla_{\phi}L = \lim_{\varepsilon \to 0} g_{\text{DLM}}^{\varepsilon}$, where $g_{\text{DLM}}^{\varepsilon} \triangleq \frac{1}{\varepsilon} \{ \frac{\partial}{\partial \phi} S(\theta^*; \phi) - \frac{\partial}{\partial \phi} S(\theta_{\text{direct}}; \phi) \}$. This was used in prior works that solve optimization problems on structured discrete domains [38, 39, 92, 93], but has so far not seen much use in structured continuous settings. We modify the original DLM formulation to better suit its implementation within Theseus

$$\theta_{\text{direct}} = \arg\min_{\hat{\theta}} S(\hat{\theta}; \phi) + \left\| \varepsilon \hat{\theta} - \frac{1}{2} \nabla_{\theta} L(\theta^*) \right\|^2, \qquad \theta^* = \arg\min_{\hat{\theta}} S(\hat{\theta}; \phi) \tag{4}$$

This is different from the original formulation in two ways: (i) we only assume access to the gradient vector $\nabla_{\theta} L(\theta^*)$ which helps formulate DLM as an algorithm for computing vector-Jacobian products, and (ii) we add a small regularization term to ensure the modified objective for θ_{direct} is a sum-of-squares without affecting the limit as $\varepsilon \to 0$. See App. E for more details.

318 5 Evaluation

We evaluate the performance of Theseus under different settings with PGO and tactile state estimation applications from Sec. 3.1. PGO allows us to easily control the problem scales for performance evaluation; in Sec. 5.1 we profile time and memory consumption of Theseus in an end-to-end setup with forward pass of 10 iterations, backward with implicit differentiation and 20 epochs, and in Sec. 5.2 we evaluate timings of Theseus as a stand-alone NLS optimizer and compare with state-of-the-art Ceres [35]. The tactile state estimation application involves a more complex outer loop model that is useful for comparing all different backward modes which we present in Sec. 5.3.

326 5.1 Profiling forward and backward pass of Theseus for DNLS

We study the performance of Theseus for DNLS on the PGO problem with the synthetic Cube dataset, as described in App. C. We use implicit differentiation to compute gradients of the inner loop, which is run with 10 inner loop iterations and 20 outer loop epochs. In all cases we used an Nvidia V100 GPU with 32GBs of memory for all Python computation, and Intel Xeon 2.2GHz CPU with 20 threads for the CPU-based CH0LM0D linear solver. We evaluate performance using our sparse solvers in Theseus and using PyTorch's Cholesky dense solver.

Fig. 2 shows the average time of a full forward and backward pass for a given batch size, taken by 333 Theseus with different solvers (cudaLU, CHOLMOD and dense) for different problem scales (number 334 of poses and batch size). The two left plots show time as a function of number of poses for two batch 335 sizes (small, 16; large, 128), while the two right plots show time as a function of batch size for two 336 pose settings (small, 256 poses; large, 2048 poses). Note that dense does not scale well with poses or 337 batch size. For a batch size of 16, the largest problem that it can solve before running out of GPU 338 memory has 512 poses; for a batch size of 128, the largest is 256 poses (left two plots). With 2048 339 poses, dense is unable to solve the problem regardless of batch size (right two plots). On the other 340 hand, our cudaLU solver, is able to scale up to 4096 poses with a batch size of 128. Note that CHOLMOD 341 can solve problems even larger, since the linear system is solved on CPU and we have successfully 342 tested up to 8192 poses and batch size 256, for a total of 22GBs of GPU usage. 343

In addition to being more memory efficient, running times of our sparse solvers are also smaller for large enough number of poses/batch size, especially for the backward pass. For the smallest problem



Figure 2: Forward/backward times of Theseus with sparse and dense solvers on different PGO problem scales.

considered (64 poses, batch size 16), the total sum of average forward and backward times are 346 similar, namely, 1.80, 1.88, and 2.10 seconds per batch for cudaLU, dense, and CHOLMOD, respectively. 347 348 Increasing to 128 poses makes cudaLU noticeably faster than dense (3.52 vs 3.87 seconds), while CHOLMOD starts outperforming dense with 256 poses (5.68 vs. 5.91 seconds, for a batch size of 8). As 349 the problem scale increases, the gap between the sparse and dense solvers widens: for the largest 350 problem solvable with dense (512 poses and 16 batch size), we have average total times of 9.84, 351 11.26, and 25.52 seconds for cudaLU, CHOLMOD, and dense, respectively. See App. D for more results 352 and details. 353

5.2 Profiling Theseus as stand-alone NLS optimizer

DNLS typically involves solving numerous optimization problems each epoch where a fast NLS optimizer is essential. We compare Theseus as a stand-alone NLS optimizer with the state-of-the-art Ceres [35] library for solving a batch of PGO problems without any learning involved. We compare all solvers in terms of the total time required to perform 10 iterations on a set of 256 PGO problems. CPU/GPU configurations are same as before. For CH0LM0D, we also include a configuration that runs everything on CPU, including Jacobians and residual computation (labelled CH0LM0D-allcpu).

Fig. 3 shows speedup obtained by Theseus with 361 batching, vectorization and sparse solvers, over 362 Ceres as a function of increasing number of 363 poses or batch size. We vary the number of 364 poses for two fixed batch sizes (small, 16; large, 365 128), and vary the batch size for two fixed num-366 367 ber of poses (small, 256; large, 2048). Although Ceres is faster than all of our solvers 368 when the number of poses and batch size are 369 small, as these increase Theseus shows sig-370 nificant speedup by being able to solve larger 371



Figure 3: Speedup of Theseus (forward) over Ceres (black dashed) on different PGO problem sizes.

batches of problems in parallel. Since typical use case of Theseus involves large batches and number
 of variables during end-to-end learning with DNLS, the speedups in this setting against a performant
 NLS solver highlights the significance of our efficiency-based design choices.

375 5.3 Backward modes analysis

We explore the trade-offs between our different backward modes using the tactile state estimation 376 application in Sec. 3.1; the learnable components here include a neural network, and thus more closely 377 follow the type of applications that motivate Theseus. We compare the following backward modes: 378 derivative unrolling (Unroll), implicit differentiation (Implicit), truncated differentiation (Trunc), 379 and direct loss minimization (DLM); for Trunc we include results when truncating 5 and 10 steps. 380 We compare all modes along 3 axis of performance, namely, validation loss after 100 epochs (outer 381 loop), run time during training, and peak GPU memory consumption of TheseusLayer. CPU/GPU 382 configurations are same as before except GPU with 16GB of memory is used. For time and memory 383 we present separate results for forward and backward pass, and all numbers are averaged over 700 (7 384 batches for 100 epochs). Below we discuss our main findings from this analysis, and more results 385 and details can be found in App. E. 386

Fig. 4 shows average run times for all backward modes as a function of the maximum number of iterations in the inner loop optimization. We observe that the time used in the forward pass (Fig. 4, far left) increases roughly linearly for all modes, all having similar times except for Unroll and Trunc-5, which are slower after 30 inner loop iterations. However, we observe stark differences in the backward pass time (Fig. 4, center left), where Unroll is the only method that has a linear dependence



Figure 4: Time and memory consumption of different backward modes in tactile state estimation.

on the number of inner loop iterations. All other methods have a constant footprint for computing derivatives, independent of the number of inner loop iterations. As expected, increasing the number of iterations through which we backprop (5 or 10 for Trunc, all iterations for Unroll) increases the time necessary for a backward pass (Implicit < Trunc-5 < Trunc-10). DLM's backward pass time is larger than Unroll's for a small number of iterations, but performs better when the number of iterations is large (> 20 in this example).

Figure 4 (center right) shows the average peak memory consumption of the backward modes. In this case, the trends observed for the backward pass memory consumption is similar to the trends in time. In particular, Unroll's memory footprint increases linearly with the number of inner loop terations, from ~ 34 MBs to ~ 262 MBs; for all other methods the memory consumption remains constant. The best memory profiles in this example is obtained with Implicit and DLM backward modes, with ~ 33 MBs and ~ 43 MBs, respectively. These trends also hold for the forward pass memory consumption.

Figure 4 also shows the validation losses obtained with all backward modes (far right). The best 405 validation loss, after 100 epochs of training, is obtained using Implicit, followed by the Trunc. 406 Most methods improve up to a certain point with an increasing number of inner loop iterations, but 407 methods that do backward through unrolled inner loop iterations (Unroll and Trunc) start performing 408 worse after 10 inner loop iterations. One exception is DLM, which doesn't improve much with the 409 number of iterations, but is also the best method when only 2 inner loop iterations are performed. 410 411 As a point of caution, we stress that, unlike the timing and memory results, the relative training performance between different backward modes is likely to be application dependent, and is affected 412 by hyperparameters such as the step size used for the inner loop optimizer (0.05 in this example), 413 and the outer optimizer's learning rate. Our experiments suggest that implicit differentiation is a 414 good default to use for differentiable optimization, considering its low time/memory footprint, and 415 potential for better end-to-end performance with proper hyperparameter tuning. 416

417 6 Discussion

Summary. Theseus provides nonlinear least squares as a differentiable layer and enables easily 418 building and training end-to-end architectures for robotics and vision applications. We illustrate 419 several example applications using the same application-agnostic interface and demonstrate significant 420 improvement in performance with our efficiency based design. Following how autodiff and GPU 421 acceleration (among others) have led to the evolution of PyTorch in contrast to numpy, we can 422 similarly view sparsity and implicit differentiation on top of autodiff and GPU acceleration as the key 423 ingredients that power Theseus in contrast to solvers like Ceres that typically only support sparsity. 424 **Limitations.** Theseus currently has a few limitations. The nonlinear solvers we currently support 425 apply constraints in a soft manner (i.e., using weighted costs). Hard constraints can be handled with 426 methods like augmented Lagrangian or sequential quadratic programs [94, 95], and differentiating 427 through them are active research topics. The current implementation of LM does not support damping 428 to be learnable. Some limitations and trade-offs with the sparse linear solvers are also discussed 429 in Sec. 4.2. Online learning applications may require frequently editing the objective and depending 430 on the problem size there may be a nontrivial overhead that is not currently optimized as we explored 431 only non-incremental setting in this work. Additional performance gains can be extracted by moving 432 some of our Python implementation to C++ but we prioritized flexibility in evolving the API in the 433 434 short-term. The automatic Jacobian computation with AutodiffCostFunction is based on torch autograd jacobian, which unnecessarily computes gradients across independent batch entries and thus 435 results in significant memory and compute for highly parameterized cost functions. We do not yet 436 support distributed training beyond what PyTorch natively supports. We will explore these features 437 and optimizations in the future as the library continues to evolve. 438

439 **References**

- [1] Jorge Nocedal and Stephen Wright. *Numerical optimization*. Springer Science & Business
 Media, 2006. 1, 2, 3
- (2) Christian Forster, Luca Carlone, Frank Dellaert, and Davide Scaramuzza. On-manifold
 preintegration for real-time visual-inertial odometry. *IEEE Transactions on Robotics*, 33(1):
 1–21, 2016. 1, 2
- [3] Frank Dellaert, Michael Kaess, et al. Factor graphs for robot perception. *Foundations and Trends*® *in Robotics*, 6(1-2):1–139, 2017. 2
- [4] Timothy D. Barfoot. State estimation for robotics. *State Estimation for Robotics*, (1987):
 1–368, 2017. doi: 10.1017/9781316671528. 2
- [5] Mustafa Mukadam, Jing Dong, Xinyan Yan, Frank Dellaert, and Byron Boots. Continuous time gaussian process motion planning via probabilistic inference. *The International Journal* of Robotics Research, 37(11):1319–1340, 2018. 2, 19, 20
- [6] Mandy Xie and Frank Dellaert. A unified method for solving inverse, forward, and hybrid
 manipulator dynamics using factor graphs. *arXiv preprint arXiv:1911.10065*, 2019. 2
- [7] Antoni Rosinol, Marcus Abate, Yun Chang, and Luca Carlone. Kimera: an open-source
 library for real-time metric-semantic localization and mapping. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1689–1696. IEEE, 2020. 1, 3
- [8] Bill Triggs, Philip F McLauchlan, Richard I Hartley, and Andrew W Fitzgibbon. Bundle
 adjustment—a modern synthesis. In *International workshop on vision algorithms*, pages
 298–372. Springer, 1999. 1, 2
- [9] Hanna Pasula, Stuart Russell, Michael Ostland, and Yaacov Ritov. Tracking many objects with
 many sensors. In *IJCAI*, volume 99, pages 1160–1171. Citeseer, 1999.
- [10] Richard Szeliski and Sing Bing Kang. Recovering 3d shape and motion from image streams
 using nonlinear least squares. *Journal of Visual Communication and Image Representation*, 5
 (1):10–28, 1994. 2
- [11] Aleksandr Segal, Dirk Haehnel, and Sebastian Thrun. Generalized-icp. In *Robotics: science and systems*, volume 2, page 435. Seattle, WA, 2009.
- [12] Tanner Schmidt, Richard A Newcombe, and Dieter Fox. Dart: Dense articulated real-time
 tracking. In *Robotics: Science and Systems*, volume 2, pages 1–9. Berkeley, CA, 2014. 2
- I3 Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceed-ings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. 1, 2
- [14] Brandon Amos and J Zico Kolter. Optnet: Differentiable optimization as a layer in neural
 networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume* 70, pages 136–145. JMLR. org, 2017. 1, 7
- [15] Akshay Agrawal, Brandon Amos, Shane Barratt, Stephen Boyd, Steven Diamond, and J Zico
 Kolter. Differentiable convex optimization layers. In *Advances in Neural Information Process- ing Systems*, pages 9558–9570, 2019. 1, 7
- [16] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast
 adaptation of deep networks. In *International conference on machine learning*, pages 1126–
 1135. PMLR, 2017. 1, 6
- [17] Edward Grefenstette, Brandon Amos, Denis Yarats, Phu Mon Htut, Artem Molchanov,
 Franziska Meier, Douwe Kiela, Kyunghyun Cho, and Soumith Chintala. Generalized in ner loop meta-learning. *arXiv preprint arXiv:1910.01727*, 2019. 1
- [18] Zachary Teed and Jia Deng. Deepv2d: Video to depth with differentiable structure from motion. *arXiv preprint arXiv:1812.04605*, 2018. 1, 3

- [19] Mohak Bhardwaj, Byron Boots, and Mustafa Mukadam. Differentiable gaussian process
 motion planning. In 2020 IEEE International Conference on Robotics and Automation (ICRA),
 pages 10598–10604. IEEE, 2020. 1, 3, 5, 6, 20
- [20] Krishna Murthy Jatavallabhula, Ganesh Iyer, and Liam Paull. ∇ slam: Dense slam meets automatic differentiation. In 2020 IEEE International Conference on Robotics and Automation (ICRA), pages 2130–2137. IEEE, 2020. 1, 3
- [21] Zachary Teed and Jia Deng. Droid-slam: Deep visual slam for monocular, stereo, and rgb-d
 cameras. Advances in Neural Information Processing Systems, 34, 2021. 1, 3
- [22] Chengzhou Tang and Ping Tan. BA-Net: Dense bundle adjustment networks. *7th International Conference on Learning Representations, ICLR 2019*, 2019. 1, 3
- [23] Brent Yi, Michelle A Lee, Alina Kloss, Roberto Martín-Martín, and Jeannette Bohg. Differ entiable factor graph optimization for learning smoothers. In 2021 IEEE/RSJ International
 Conference on Intelligent Robots and Systems (IROS), pages 1339–1345. IEEE, 2021. 1, 3, 4
- [24] Zhaoyang Lv, Frank Dellaert, James M Rehg, and Andreas Geiger. Taking a deeper look at the
 inverse compositional algorithm. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4581–4590, 2019. 1, 3
- [25] Ross Hartley, Maani Ghaffari Jadidi, Lu Gan, Jiunn-Kai Huang, Jessy W Grizzle, and Ryan M
 Eustice. Hybrid contact preintegration for visual-inertial-contact state estimation using factor
 graphs. In 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS),
 pages 3783–3790. IEEE, 2018. 1, 2
- [26] Paloma Sodhi, Michael Kaess, Mustafa Mukadam, and Stuart Anderson. Learning tactile
 models for factor graph-based estimation. In 2021 IEEE International Conference on Robotics
 and Automation (ICRA), pages 13686–13692. IEEE, 2021. 1, 3, 5, 19
- [27] Markus Giftthaler, Michael Neunert, Markus Stäuble, Jonas Buchli, and Moritz Diehl. A
 family of iterative gauss-newton shooting methods for nonlinear optimal control. In 2018
 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 1–9.
 IEEE, 2018. 1, 2
- [28] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and
 Michael J Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a
 single image. In *European conference on computer vision*, pages 561–578. Springer, 2016. 1,
 2
- [29] Taosha Fan, Kalyan Vasudev Alwala, Donglai Xiang, Weipeng Xu, Todd Murphey, and
 Mustafa Mukadam. Revitalizing optimization for 3d human pose and shape estimation: A
 sparse constrained formulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11457–11466, 2021. 1, 2
- [30] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative
 style, high-performance deep learning library. *Advances in neural information processing systems*, 32:8026–8037, 2019. 2
- [31] Joan Solà, Jeremie Deray, and Dinesh Atchuthan. A micro Lie theory for state estimation in robotics. *arXiv*, pages 1–17, 2018. URL http://arxiv.org/abs/1812.01537. 2, 4
- [32] Franziska Meier, Austin Wang, Giovanni Sutanto, Yixin Lin, and Paarth Shah. Differentiable
 and learnable robot models. *arXiv preprint arXiv:2202.11217*, 2022. 2, 5
- [33] Giorgio Grisetti, Rainer Kümmerle, Hauke Strasdat, and Kurt Konolige. g2o: A general frame work for (hyper) graph optimization. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pages 9–13, 2011. 2, 3, 18
- [34] Frank Dellaert. Factor graphs and gtsam: A hands-on introduction. Technical report, Georgia
 Institute of Technology, 2012.

- [35] Sameer Agarwal and Keir Mierle. Ceres solver: Tutorial & reference. *Google Inc*, 2(72):8, 2012. 7, 8, 18
- [36] Jing Dong and Zhaoyang Lv. miniSAM: A flexible factor graph non-linear least squares optimization framework. *CoRR*, abs/1909.00903, 2019. URL http://arxiv.org/abs/1909.
 00903. 2, 3
- [37] Asen L Dontchev and R Tyrrell Rockafellar. *Implicit functions and solution mappings*, volume
 543. Springer, 2009. 2, 7, 21
- [38] Tamir Hazan, Joseph Keshet, and David McAllester. Direct loss minimization for structured prediction. *Advances in neural information processing systems*, 23, 2010. 2, 7
- [39] Yang Song, Alexander Schwing, Raquel Urtasun, et al. Training deep neural networks via
 direct loss minimization. In *International conference on machine learning*, pages 2169–2177.
 PMLR, 2016. 2, 7, 22
- [40] Frank Dellaert and Michael Kaess. Square root sam: Simultaneous localization and mapping
 via square root information smoothing. *The International Journal of Robotics Research*, 25
 (12):1181–1203, 2006. 2
- [41] Carl T Kelley. Iterative methods for optimization. SIAM, 1999. 3
- [42] Edgar Sucar, Kentaro Wada, and Andrew Davison. NodeSLAM: Neural Object Descriptors
 for Multi-View Shape Reconstruction. *Proceedings 2020 International Conference on 3D Vision, 3DV 2020*, pages 949–958, 2020. doi: 10.1109/3DV50981.2020.00105. 3
- [43] Michael Bloesch, Jan Czarnowski, Ronald Clark, Stefan Leutenegger, and Andrew J Davison.
 Codeslam—learning a compact, optimisable representation for dense visual slam. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2560–2568, 2018. 3
- [44] Jan Czarnowski, Tristan Laidlow, Ronald Clark, and Andrew J. Davison. DeepFactors: Real-Time Probabilistic Dense Monocular SLAM. *IEEE Robotics and Automation Letters*, 5(2): 721–728, 2020. doi: 10.1109/LRA.2020.2965415. 3
- [45] Muhammad Asif Rana, Mustafa Mukadam, Seyed Reza Ahmadzadeh, Sonia Chernova, and
 Byron Boots. Towards robust skill generalization: Unifying learning from demonstration and
 motion planning. In *Conference on Robot Learning*, pages 109–118. PMLR, 2017. 3
- [46] Ronald Clark, Michael Bloesch, Jan Czarnowski, Stefan Leutenegger, and Andrew J Davison.
 Ls-net: Learning to solve nonlinear least squares for monocular stereo. *arXiv preprint arXiv:1809.02966*, 2018. 3
- [47] Barbara Roessle and Matthias Nießner. End2end multi-view feature matching using differen tiable pose optimization. 2022. URL https://arxiv.org/abs/2205.01694. 3
- [48] Jingwei Huang, Shan Huang, and Mingwei Sun. Deeplm: Large-scale nonlinear least squares
 on deep learning frameworks using stochastic domain decomposition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10308–10317,
 2021. 3, 5
- [49] Borglab. Swiftfusion. https://github.com/borglab/SwiftFusion, 2020. 3
- [50] Jonathan T Barron. A general and adaptive robust loss function. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4331–4339, 2019.
 4
- [51] Valentin Peretroukhin, Matthew Giamou, W. Nicholas Greene, David Rosen, Jonathan Kelly,
 and Nicholas Roy. A Smooth Representation of Belief over SO(3) for Deep Rotation Learning
 with Uncertainty. (3), 2020. doi: 10.15607/rss.2020.xvi.007. URL http://arxiv.org/abs/
 2006.01031. 4

- [52] Zachary Teed and Jia Deng. Tangent Space Backpropagation for 3D Transformation Groups.
 Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pages 10333–10342, 2021. doi: 10.1109/CVPR46437.2021.01020. URL http: //arxiv.org/abs/2103.12032. 4
- [53] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2019-June:5738–5746, 2019. doi: 10.1109/CVPR.2019.00589.
- [54] Congyue Deng, Or Litany, Yueqi Duan, Adrien Poulenard, Andrea Tagliasacchi, and Leonidas
 Guibas. Vector Neurons: A General Framework for SO(3)-Equivariant Networks. (3), 2021.
 doi: 10.1109/iccv48922.2021.01198. URL http://arxiv.org/abs/2104.12229.
- [55] Philippe Hansen-Estruch, Wenling Shang, Lerrel Pinto, Pieter Abbeel, and Stas Tiomkin.
 GEM: Group Enhanced Model for Learning Dynamical Control Systems. 2021. URL
 http://arxiv.org/abs/2104.02844. 4
- [56] P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2009. 4
- 596 [57] Brent Yi. jaxlie. https://github.com/brentyi/jaxlie, 2021. 5
- [58] D.M. Rosen, L. Carlone, A.S. Bandeira, and J.J. Leonard. SE-Sync: A certifiably correct algorithm for synchronization over the special Euclidean group. *Intl. J. of Robotics Research*, 38(2–3):95–125, March 2019. 5, 18
- [59] Mike Lambeta, Po-Wei Chou, Stephen Tian, Brian Yang, Benjamin Maloon, Victoria Rose
 Most, Dave Stroud, Raymond Santos, Ahmad Byagowi, Gregg Kammerer, et al. Digit: A
 novel design for a low-cost compact high-resolution tactile sensor with application to in-hand
 manipulation. *IEEE Robotics and Automation Letters*, 5(3):3838–3845, 2020. 5, 19
- [60] Sameer Agarwal, Noah Snavely, Steven M Seitz, and Richard Szeliski. Bundle adjustment in the large. In *European conference on computer vision*, pages 29–42. Springer, 2010. 5, 19
- [61] Jan Czarnowski, Stefan Leutenegger, and Andrew J. Davison. Semantic texture for robust
 dense tracking. In *Proceedings of the IEEE International Conference on Computer Vision* (*ICCV*) Workshops, Oct 2017. 5, 20
- [62] Jing Dong, Byron Boots, Frank Dellaert, Ranveer Chandra, and Sudipta N. Sinha. Learning to align images using weak geometric supervision. 2018.
- [63] Chengzhou Tang and Ping Tan. BA-Net: Dense Bundle Adjustment Network. In *ICLR*, 2019.
- [64] Paul-Edouard Sarlin, Ajaykumar Unagar, Måns Larsson, Hugo Germain, Carl Toft, Viktor
 Larsson, Marc Pollefeys, Vincent Lepetit, Lars Hammarstrand, Fredrik Kahl, and Torsten
 Sattler. Back to the Feature: Learning robust camera localization from pixels to pose. In
 CVPR, 2021. URL https://arxiv.org/abs/2103.09213.
- [65] Zhaoyang Lv, Frank Dellaert, James Rehg, and Andreas Geiger. Taking a deeper look at the
 inverse compositional algorithm. In *CVPR*, 2019.
- [66] Chaoyang Wang, Hamed Kiani Galoogahi, Chen-Hsuan Lin, and Simon Lucey. Deep-lk
 for efficient adaptive object tracking. 2018 IEEE International Conference on Robotics and
 Automation (ICRA), pages 627–634, 2018. 5
- [67] Yanqing Chen, Timothy A Davis, William W Hager, and Sivasankaran Rajamanickam. Algorithm 887: Cholmod, supernodal sparse cholesky factorization and update/downdate. ACM
 Transactions on Mathematical Software (TOMS), 35(3):1–14, 2008. 6
- [68] Jonathan T Barron and Ben Poole. The fast bilateral solver. In *European conference on computer vision*, pages 617–632. Springer, 2016. 6

- [69] Barak A Pearlmutter and Jeffrey Mark Siskind. Reverse-mode ad in a functional framework:
 Lambda the ultimate backpropagator. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 30(2):1–36, 2008. 6
- [70] Chongjie Zhang and Victor Lesser. Multi-agent learning with policy prediction. In *Twentyfourth AAAI conference on artificial intelligence*, 2010.
- [71] Dougal Maclaurin, David Duvenaud, and Ryan Adams. Gradient-based hyperparameter
 optimization through reversible learning. In *International conference on machine learning*,
 pages 2113–2122. PMLR, 2015.
- [72] David Belanger and Andrew McCallum. Structured prediction energy networks. In *Interna- tional Conference on Machine Learning*, pages 983–992, 2016.
- [73] Luke Metz, Ben Poole, David Pfau, and Jascha Sohl-Dickstein. Unrolled generative adversarial
 networks. *CoRR*, abs/1611.02163, 2016. URL http://arxiv.org/abs/1611.02163.
- [74] Tian Han, Yang Lu, Song-Chun Zhu, and Ying Nian Wu. Alternating back-propagation
 for generator network. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
 volume 31, 2017.
- [75] David Belanger, Bishan Yang, and Andrew McCallum. End-to-end learning for structured
 prediction energy networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 429–439. JMLR. org, 2017.
- [76] David Belanger. *Deep energy-based models for structured prediction*. PhD thesis, University
 of Massachusetts Amherst, 2017.
- [77] Jakob N Foerster, Richard Y Chen, Maruan Al-Shedivat, Shimon Whiteson, Pieter Abbeel, and
 Igor Mordatch. Learning with opponent-learning awareness. *arXiv preprint arXiv:1709.04326*,
 2017.
- [78] Vishal Monga, Yuelong Li, and Yonina C Eldar. Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing. *IEEE Signal Processing Magazine*, 38(2):
 18–44, 2021. 6
- [79] Paul J Werbos. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10):1550–1560, 1990. 6
- [80] Herbert Jaeger. *Tutorial on training recurrent neural networks, covering BPPT, RTRL, EKF and the" echo state network" approach*, volume 5. GMD-Forschungszentrum Informationstechnik
 Bonn, 2002. 6
- [81] Corentin Tallec and Yann Ollivier. Unbiasing truncated backpropagation through time. *arXiv preprint arXiv:1705.08209*, 2017. 7
- [82] Yuhuai Wu, Mengye Ren, Renjie Liao, and Roger Grosse. Understanding short-horizon bias
 in stochastic meta-optimization. *arXiv preprint arXiv:1803.02021*, 2018.
- [83] Renjie Liao, Yuwen Xiong, Ethan Fetaya, Lisa Zhang, KiJung Yoon, Xaq Pitkow, Raquel Urta sun, and Richard Zemel. Reviving and improving recurrent back-propagation. In *International Conference on Machine Learning*, pages 3082–3091. PMLR, 2018.
- [84] Amirreza Shaban, Ching-An Cheng, Nathan Hatch, and Byron Boots. Truncated back propagation for bilevel optimization. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1723–1732. PMLR, 2019.
- [85] Paul Vicol, Luke Metz, and Jascha Sohl-Dickstein. Unbiased gradient estimation in unrolled
 computation graphs with persistent evolution strategies. In *International Conference on Machine Learning*, pages 10553–10563. PMLR, 2021. 7
- [86] Brandon Amos. *Differentiable Optimization-Based Modeling for Machine Learning*. PhD
 thesis, Carnegie Mellon University, May 2019. 7

- [87] Justin Domke. Generic methods for optimization-based modeling. In *Artificial Intelligence and Statistics*, pages 318–326. PMLR, 2012.
- [88] Stephen Gould, Basura Fernando, Anoop Cherian, Peter Anderson, Rodrigo Santa Cruz, and
 Edison Guo. On differentiating parameterized argmin and argmax problems with application
 to bi-level optimization. *arXiv preprint arXiv:1607.05447*, 2016.
- [89] Jonathan Lorraine, Paul Vicol, and David Duvenaud. Optimizing millions of hyperparameters
 by implicit differentiation. In *International Conference on Artificial Intelligence and Statistics*,
 pages 1540–1552. PMLR, 2020.
- [90] Mathieu Blondel, Quentin Berthet, Marco Cuturi, Roy Frostig, Stephan Hoyer, Felipe Llinares López, Fabian Pedregosa, and Jean-Philippe Vert. Efficient and modular implicit differentiation.
 arXiv preprint arXiv:2105.15183, 2021.
- [91] Matan Atzmon, Niv Haim, Lior Yariv, Ofer Israelov, Haggai Maron, and Yaron Lipman.
 Controlling neural level sets. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,
 volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper/
 2019/file/b20bb95ab626d93fd976af958fbc61ba-Paper.pdf. 7
- [92] Marin Vlastelica, Anselm Paulus, Vít Musil, Georg Martius, and Michal Rolínek. Differentiation of blackbox combinatorial solvers. *arXiv preprint arXiv:1912.02175*, 2019. 7
- [93] Mathias Niepert, Pasquale Minervini, and Luca Franceschi. Implicit mle: Backpropagating
 through discrete exponential family distributions. *Advances in Neural Information Processing Systems*, 34, 2021. 7
- [94] Paloma Sodhi, Sanjiban Choudhury, Joshua G. Mangelson, and Michael Kaess. ICS: Incre mental Constrained Smoothing for State Estimation, 2020. 9
- [95] André F.T. Martins, Mário A.T. Figueiredo, Pedro M.Q. Aguiar, Noah A. Smith, and Eric P.
 Xing. An augmented Lagrangian approach to constrained MAP inference. *Proceedings of the* 28th International Conference on Machine Learning, ICML 2011, pages 169–176, 2011. 9
- [96] Taosha Fan, Hanlin Wang, Michael Rubenstein, and Todd Murphey. Cpl-slam: Efficient and
 certifiably correct planar graph-based slam using the complex number representation. *IEEE Transactions on Robotics*, 36(6):1719–1737, 2020. 18
- [97] Daniel Martinec and Tomas Pajdla. Robust rotation and translation estimation in multiview
 reconstruction. In 2007 IEEE Conference on Computer Vision and Pattern Recognition, pages
 1–8. IEEE, 2007. 18
- [98] Amit Singer and Yoel Shkolnisky. Three-dimensional structure determination from common
 lines in cryo-em by eigenvectors and semidefinite programming. *SIAM journal on imaging sciences*, 4(2):543–572, 2011. 18
- [99] Mihai Cucuringu, Yaron Lipman, and Amit Singer. Sensor network localization by eigenvector
 synchronization over the euclidean group. ACM Transactions on Sensor Networks (TOSN), 8
 (3):1–42, 2012. 18
- [100] Jiaji Zhou, James A. Bagnell, and Matthew T. Mason. A fast stochastic contact model for planar
 pushing and grasping: Theory and experimental validation. In Nancy M. Amato, Siddhartha S.
 Srinivasa, Nora Ayanian, and Scott Kuindersma, editors, *Robotics: Science and Systems XIII, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA, July 12-16, 2017,* 2017. doi: 10.15607/RSS.2017.XIII.040. URL http://www.roboticsproceedings.org/
 rss13/p40.html. 19
- [101] Paloma Sodhi, Eric Dexheimer, Mustafa Mukadam, Stuart Anderson, and Michael Kaess. Leo:
 Learning energy-based models in factor graph optimization. In *Conference on Robot Learning*,
 pages 234–244. PMLR, 2022. 19

- [102] Bruce D. Lucas and Takeo Kanade. An iterative image registration technique with an applica tion to stereo vision. In *Proceedings of the 7th International Joint Conference on Artificial Intelligence Volume 2*, IJCAI'81, page 674–679, San Francisco, CA, USA, 1981. Morgan
 Kaufmann Publishers Inc. 20
- [103] Simon Baker and Iain A. Matthews. Lucas-kanade 20 years on: A unifying framework.
 International Journal of Computer Vision, 56(3):221–255, 2004. URL https://doi.org/10.
 1023/B:VISI.0000011205.11775.fd. 20
- [104] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Deep image homography
 estimation. RSS 2016 Workshop: Limits and Potentials of Deep Learning in Robotics, 2016.
- 728 URL http://arxiv.org/abs/1606.03798.20
- [105] Ulisse Dini. Analisi infinitesimale. Lithografia Gorani, 1878. 21

730 Checklist

731	1. For all authors
732 733	 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
734	(b) Did you describe the limitations of your work? [Yes]
735 736	(c) Did you discuss any potential negative societal impacts of your work? [No] We do not foresee any beyond what the field of differentiable optimization already have.
737 738	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
739	2. If you are including theoretical results
740 741	(a) Did you state the full set of assumptions of all theoretical results? [N/A](b) Did you include complete proofs of all theoretical results? [N/A]
742	3. If you ran experiments
743 744	(a) Did you include the code, data, and instructions needed to reproduce the main experi- mental results (either in the supplemental material or as a URL)? [Yes]
745 746	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
747 748	(c) Did you report error bars (e.g., with respect to the random seed after running experi- ments multiple times)? [Yes]
749 750	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
751	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
752	(a) If your work uses existing assets, did you cite the creators? [Yes]
753	(b) Did you mention the license of the assets? [Yes]
754	(c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
755	(d) Did you discuss whether and how consent was obtained from people whose data you're
756	using/curating? [N/A]
757 758	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
759	5. If you used crowdsourcing or conducted research with human subjects
760	(a) Did you include the full text of instructions given to participants and screenshots if
761	applicable? [N/A]
762	(b) Did you describe any potential participant risks, with links to Institutional Review
763	Board (IRB) approvals, if applicable? [N/A]
764 765	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]