

Optimality of the relaxed polar factors by a characterization of the set of real square roots of real symmetric matrices

Lev Borisov,¹ Andreas Fischle,² and Patrizio Neff³

June 30, 2016

Abstract

We consider the problem to determine the optimal rotations $R \in \text{SO}(n)$ which minimize

$$W : \text{SO}(n) \rightarrow \mathbb{R}_0^+, \quad W(R; D) := \|\text{sym}(RD - \mathbf{1})\|^2$$

for a given diagonal matrix $D := \text{diag}(d_1, \dots, d_n) \in \mathbb{R}^{n \times n}$. The function W subject to minimization is the reduced form of the Cosserat shear-stretch energy, which, in its general form, is a contribution in any geometrically nonlinear, isotropic and quadratic Cosserat micropolar (extended) continuum model. We characterize the critical points of the energy $W(R; D)$, determine the global minimizers and the global minimum. This proves the correctness of previously obtained formulae for the optimal Cosserat rotations in dimensions two and three. The key to the proof is a characterization of the entire set of (possibly non-symmetric) real matrix square roots of (possibly non-positive definite) real symmetric matrices which does not seem to be known in the literature.

Keywords: Cosserat theory, micropolar media, Grioli's theorem, rotations, special orthogonal group, (non-symmetric) matrix square root, symmetric square, polar decomposition, relaxed-polar decomposition.

AMS 2010 subject classification: 15A24, 22E30, 74A30, 74A35, 74B20, 74G05, 74G65, 74N15.

Contents

1	Introduction	2
2	Representation of real matrix square roots of symmetric matrices	6
3	Critical points of the Cosserat shear-stretch energy	11
4	Analysis of the decoupled subproblems	14
5	Global minimization of the Cosserat shear-stretch energy	17
6	Discussion	22
	References	22

¹Lev Borisov, Department of Mathematics, Rutgers University, 240 Hill Center, Newark, NJ 07102, United States, email: borisov@math.rutgers.edu

²Corresponding author: Andreas Fischle, Institut für Numerische Mathematik, TU Dresden, Zellescher Weg 12-14, 01069 Dresden, Germany, email: andreas.fischle@tu-dresden.de

³Patrizio Neff, Head of Lehrstuhl für Nichtlineare Analysis und Modellierung, Fakultät für Mathematik, Universität Duisburg-Essen, Thea-Leymann Str. 9, 45127 Essen, Germany, email: patrizio.neff@uni-due.de

1 Introduction

1.1 The problem

We consider the optimality problem for rotations $R \in \text{SO}(n)$ parametrized by a diagonal matrix $D := \text{diag}(d_1, \dots, d_n) \in \text{Diag}(n)$ as stated in the following

Problem 1.1. *Let*

$$W : \text{SO}(n) \times \text{Diag}(n) \rightarrow \mathbb{R}_0^+, \quad W(R; D) := \|\text{sym}(RD - \mathbb{1})\|^2. \quad (1.1)$$

Compute the set of energy-minimizing rotations

$$\text{rpolar}(D) := \arg \min_{R \in \text{SO}(n)} W(R; D) = \arg \min_{R \in \text{SO}(n)} \|\text{sym}(RD - \mathbb{1})\|^2 \subseteq \text{SO}(n). \quad (1.2)$$

We use the notation $\text{sym}(X) := \frac{1}{2}(X + X^T)$, $\text{skew}(X) := \frac{1}{2}(X - X^T)$, $\text{dev}(X) := X - \frac{1}{n} \text{tr}[X] \cdot \mathbb{1}$, $\langle X, Y \rangle := \text{tr}[X^T Y]$ and we denote the induced Frobenius matrix norm by $\|X\|^2 := \langle X, X \rangle = \sum_{1 \leq i, j \leq n} X_{ij}^2$. We call a rotation $R \in \text{SO}(n)$ optimal for given $D \in \text{Diag}(n)$ if it is a global minimizer for the energy $W(R; D)$ defined in (1.1). Technically, the decisive point in the solution of this minimization is the characterization of the set of rotations $R \in \text{SO}(n)$ satisfying the particular symmetric square condition

$$(RD - \mathbb{1})^2 \in \text{Sym}(n)$$

which is equivalent to the Euler-Lagrange equations of (1.1).

1.2 Motivation and previous results

The optimality problem which we consider here is a distinguished special case to which a more general optimality problem arising in the context of Cosserat theory in solid mechanics can be reduced, see [5–8, 28] for preliminary work. Problem 1.1 is the key step which determines the optimal Cosserat rotations which minimize the Cosserat shear-stretch energy in the general case

$$W_{\mu, \mu_c}(\bar{R}; F) := \mu \left\| \text{sym}(\bar{R}^T F - \mathbb{1}) \right\|^2 + \mu_c \left\| \text{skew}(\bar{R}^T F - \mathbb{1}) \right\|^2. \quad (1.3)$$

The two arguments for the Cosserat shear-stretch energy $W_{\mu, \mu_c} : \text{SO}(n) \times \text{GL}^+(n) \rightarrow \mathbb{R}_0^+$ are the deformation gradient field $F := \nabla \varphi : \Omega \rightarrow \text{GL}^+(n)$ and the Cosserat microrotation field $\bar{R} : \Omega \rightarrow \text{SO}(n)$ evaluated at a given point of a body Ω which is subjected to an admissible deformation mapping $\varphi : \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$. The two weights $\mu > 0$ and $\mu_c \geq 0$ can be interpreted as material parameters; the Lamé shear modulus $\mu > 0$ from linear elasticity and the so-called Cosserat couple modulus $\mu_c \geq 0$, see [25] for a discussion.

Cosserat theory is a model class in nonlinear solid mechanics which explicitly introduces an additional field of rotations, an approach which is also commonly referred to as a micropolar continuum theory; see [3] for an introduction including extensive references. This type of models dates back to the original work of the Cosserat brothers [2]. In a hyperelastic approach, the Cosserat shear-stretch energy density $W_{\mu, \mu_c}(R; F)$ is a contribution to the total elastically stored energy in the variational formulation for any geometrically nonlinear, isotropic and quadratic Cosserat-micropolar continuum model, see [2, 4] and [22]. Historically, the Cosserat brothers themselves were laying foundations regarding the physically necessary invariance requirements for a micropolar continuum theory. For example, they proved that the energy density W in such a theory must be a function of the first Cosserat deformation tensor $\bar{U} := R^T F$. They never proposed a specific expression for the local energy density $W = W(\bar{U})$ in order to model specific materials. The chosen quadratic ansatz for $W_{\mu, \mu_c}(\bar{U})$ which we are interested in, is motivated by a direct extension of the quadratic energy in the linear theory of Cosserat models, see, e.g. [19, 29, 30].

Let us introduce the polar factor $R_p(F) \in \text{SO}(n)$ which is obtained from the right polar decomposition $F = R_p(F) U(F)$ of the deformation gradient $F \in \text{GL}^+(n)$. Here, $U(F) := \sqrt{F^T F} \in \text{PSym}(n)$ denotes the positive definite symmetric right Biot-stretch tensor. Furthermore, we recall that the singular values ν_i , $i = 1, \dots, n$, of the deformation gradient $F \in \text{GL}^+(n)$ are defined as the eigenvalues of $U \in \text{PSym}(n)$.

Our original motivation to characterize the energy-minimizing rotations in Problem 1.1 was a study of strain energy densities related to certain distance problems in nonlinear continuum mechanics. For example, one may consider the following euclidean distance function

$$\text{dist}_{\text{euclid}}^2(F, \text{SO}(3)) := \min_{R \in \text{SO}(3)} \|F - R\|^2 . \quad (1.4)$$

Conceptually, this distance function locally measures the distance of a diffeomorphism $\varphi : \Omega \rightarrow \varphi(\Omega)$ to the subgroup of isometric embeddings of the body Ω into \mathbb{R}^3 . The required invariance properties for isotropy are automatically satisfied. Furthermore, this is consistent with the requirement that a global isometry of a body $\Omega \subset \mathbb{R}^3$ (i.e., a rigid body motion) does not produce any deformation energy, because in that case $F = \nabla(Rx + b) = R \in \text{SO}(3)$ which implies

$$\int_{\Omega} \text{dist}_{\text{euclid}}^2(F, \text{SO}(3)) \, dV = 0 .$$

Variations on this general theme lead to the study of corresponding minimization problems on $\text{SO}(n)$ which have been the subject of multiple contributions, see, e.g., [5–7, 20, 27, 32]. Note that in classical nonlinear continuum models, the local rotation of the specimen at a point is not explicitly accounted for in the strain energy, due to the requirement of frame-indifference. Thus, in a classical theory, the local rotation of the specimen induced by a deformation mapping φ is always given by the continuum rotation $R_p(\nabla\varphi)$.

In strong contrast, in Cosserat theory and other generalized continuum theories (so-called complex materials) with rotational degrees of freedom, the local rotation $\bar{R} : \Omega \rightarrow \text{SO}(3)$ of the material appears explicitly. Accordingly, in such a theory, the computation of locally energy-minimizing rotations provides geometrical insight into the qualitative mechanical behavior of a particular constitutive model.

The first result in this area apparently dates back to 1940 when Grioli [13] proved the following remarkable variational characterization of the orthogonal factor $R_p(F)$:

$$\arg \min_{\bar{R} \in \text{SO}(3)} \left\| \bar{R}^T F - \mathbf{1} \right\|^2 = \{R_p(F)\} , \quad \text{and} \quad (1.5)$$

$$\min_{\bar{R} \in \text{SO}(3)} \left\| \bar{R}^T F - \mathbf{1} \right\|^2 = W_{\text{Biot}}(F) := \|U(F) - \mathbf{1}\|^2 . \quad (1.6)$$

Grioli's result implies that $R_p(F)$ is optimal for the Cosserat strain energy minimized in (1.5). Hence, this strain energy can be expected to produce a microrotation field approximating the continuum rotation $\bar{R} \approx R_p(\nabla\varphi)$ and realizing the Biot-energy. In other words, the corresponding Cosserat boundary value problem can be expected to behave essentially in the same way as a classical nonlinear Biot-model and the actual impact of the additional field of microrotations \bar{R} seems rather limited. In order to introduce non-classical effects, one has to look further.

Following [7], one can simplify (1.5) by exploiting *isotropy* of the energy. As it turns out, it is sufficient to consider rotations *relative* to given $R_p(F) \in \text{SO}(3)$. In this relative picture the deformation gradient F is represented by a diagonal matrix $D := \text{diag}(d_1, d_2, d_3)$ where the entries $d_i = \nu_i > 0$ coincide with the singular values of $F \in \text{GL}^+(3)$. Grioli's theorem then takes the following form

$$\arg \min_{\bar{R} \in \text{SO}(3)} \|\bar{R}D - \mathbf{1}\|^2 = \{\mathbf{1}\} , \quad \text{and} \quad (1.7)$$

$$\min_{\bar{R} \in \text{SO}(3)} \|\bar{R}D - \mathbf{1}\|^2 = \|D - \mathbf{1}\|^2 . \quad (1.8)$$

This result generalizes to arbitrary dimension n and we refer to it as *Grioli's theorem* for absolute and relative optimal rotations, respectively.

One possible extension of the quadratic strain energy density (1.7) is to weight the euclidean distance with respect to the orthogonal symmetric and skew-symmetric contributions which gives

$$\mu \|\text{sym}(RD - \mathbf{1})\|^2 + \mu_c \|\text{skew}(RD - \mathbf{1})\|^2 . \quad (1.9)$$

This contribution naturally appears in any quadratic geometrically nonlinear Cosserat strain energy; cf., e.g., [23]. A non-trivial parameter reduction described in [6] shows that the corresponding

energy-minimizing rotations can be determined by the solution of our Problem 1.1. As an aside, one may also consider the expression where the quadratic volumetric contribution has been singled out and weighted independently by the bulk modulus κ , i.e.,

$$\mu \|\text{dev sym}(RD - \mathbb{1})\|^2 + \mu_c \|\text{skew}(RD - \mathbb{1})\|^2 + \frac{\kappa}{2} (\text{tr}[(RD - \mathbb{1})])^2. \quad (1.10)$$

Since a quadratic volume contribution seems not very attractive in comparison with an exact volume contribution of the form $W_{\text{vol}}(\det[F])$, we have abstained from investigating the formulation (1.10).

Similar optimality questions can also be formulated for a logarithmic non-symmetric microstretch tensor $\log(RD)$ for which we consider

$$W_{\log}(R; D) := \mu \|\text{dev sym} \log(RD)\|^2 + \mu_c \|\text{skew} \log(RD)\|^2 + \frac{\kappa}{2} (\text{tr}[\log(RD)])^2. \quad (1.11)$$

Technicalities aside, one can show that the corresponding minimization problem satisfies

$$\arg \min_{R \in \text{SO}(n)} W_{\log}(R; D) = \{\mathbb{1}\}, \quad \text{and} \quad (1.12)$$

$$\min_{R \in \text{SO}(n)} W_{\log}(R; D) = \mu \|\text{dev} \log D\|^2 + \frac{\kappa}{2} (\text{tr}[\log D])^2, \quad (1.13)$$

see [1, 20, 32] for some essential details. This result has been shown to be closely related to a family of geodesic distances from $F \in \text{GL}^+(3)$ to the subgroup $\text{SO}(3)$ which induce Hencky-type strain energies [26]. We summarize that the minimization problem for the euclidean distance measure (1.7) and the minimization problem for the logarithmic energy (1.12) share a remarkable property: the identity $\mathbb{1} \in \text{SO}(n)$ is always uniquely optimal for any diagonal positive definite $D > 0$. Equivalently, the polar decomposition $R_p(F)$ is always the optimal absolute rotation. In strong contrast, the quadratic energy density (1.9) also admits *non-classical optimal solutions*, see [6, 7, 28] (and [8] for a visualization in the context of an idealized nanoindentation). More precisely, one can choose the parameters μ and μ_c such that

$$\arg \min_{R \in \text{SO}(n)} \left(\mu \|\text{sym}(RD - \mathbb{1})\|^2 + \mu_c \|\text{skew}(RD - \mathbb{1})\|^2 \right) \neq \{\mathbb{1}\}. \quad (1.14)$$

In the absolute picture this implies that the corresponding optimal absolute rotations $\text{rpolar}_{\mu, \mu_c}(F)$ strictly deviate from the continuum rotation $R_p(\nabla \varphi)$. Note that the quadratic stored energy density (1.10) can be obtained from the logarithmic one (1.11) by linearization of $\log(RD) = RD - \mathbb{1} + \text{h.o.t.}$ with respect to $\overline{U} := RD$.

1.3 Technical approach and results

Let us now consider the quadratic formulation

$$\mu \|\text{sym}(RD - \mathbb{1})\|^2 + \mu_c \|\text{skew}(RD - \mathbb{1})\|^2 \quad (1.15)$$

in more detail. The choice of values for the weights (material parameters) $\mu > 0$ and $\mu_c \geq 0$ can be seen to be of crucial importance for the corresponding minimization problem. In fact, one observes two qualitatively different scenarios connected by a bifurcation criterion [6, 7]. In the *classical parameter range* $\mu_c \geq \mu > 0$, we recover a variational characterization of the polar factor $R_p(F)$ as the unique rotation minimizing (1.3) for arbitrary $n \geq 2$. We refer to this characterization as the *generalized Grioli's theorem*, see [13, 31], or [6, Cor. 2.4, p. 5]. Due to the optimality of $R_p(F)$ stated in Grioli's theorem, we have motivated the notation $\text{rpolar}_{\mu, \mu_c}(F)$ for the optimal Cosserat rotations and refer to this set of rotations as the relaxed polar factors of F with weights μ and μ_c .

Let us now restrict our attention to the *non-classical parameter range* $\mu > \mu_c > 0$. Surprisingly, this entire range can be reduced to a single *non-classical limit case* $(\mu, \mu_c) = (1, 0)$, see [6]. The Cosserat shear-stretch energy $W_{1,0}(R; F)$ can then be rewritten in terms of a relative rotation $R \in \text{SO}(n)$ which acts relative to the polar factor $R_p(F)$. In this second reduction step, the parameter $F \in \text{GL}^+(n)$ is replaced by a diagonal matrix $D = \text{diag}(\nu_1, \dots, \nu_n)$, where $\nu_i > 0$, $i = 1, \dots, n$, denote the singular values of the deformation gradient $F \in \text{GL}^+(n)$. Carrying out the beforementioned simplifications, we arrive at Problem 1.1.

Explicit formulae for the critical points and the global minimizers $\text{rpolar}_{\mu, \mu_c}^{\pm}(F)$ for the general form of the quadratic Cosserat shear-stretch energy $W_{\mu, \mu_c}(R; F)$ in dimension $n = 2$ have been presented in [6]. The corresponding minimal energy levels were also provided. In dimension $n = 3$, the following explicit formulae for the solutions to Problem 1.1 were obtained using computer algebra [7, Corollary 2.7]:

Corollary 1.2 (Energy-minimizing relative rotations for $(\mu, \mu_c) = (1, 0)$). *Let $D = \text{diag}(d_1, d_2, d_3)$ such that $d_1 > d_2 > d_3 > 0$. Then the solutions to Problem 1.1 are given by the energy-minimizing relative rotations*

$$\text{rpolar}(D) = \left\{ \begin{pmatrix} \cos \alpha & -\sin \alpha & 0 \\ \sin \alpha & \cos \alpha & 0 \\ 0 & 0 & 1 \end{pmatrix} \right\}, \quad (1.16)$$

where $\alpha \in [-\pi, \pi]$ is an optimal rotation angle satisfying

$$\alpha = \begin{cases} 0 & , \quad \text{if } d_1 + d_2 \leq 2, \\ \pm \arccos(\frac{2}{d_1 + d_2}) & , \quad \text{if } d_1 + d_2 \geq 2. \end{cases} \quad (1.17)$$

In particular, for $d_1 + d_2 \leq 2$, we have $\text{rpolar}(D) = \{\mathbb{1}\}$.

The validation of the formulae (1.16) and (1.17) in [7] was based on brute force stochastic minimization, since a proof of optimality was out of reach. With the present contribution, we close this gap in $n = 3$ and generalize the previously obtained formulae $\text{rpolar}_{1,0}^{\pm}(F)$ from [7, 8] to arbitrary dimension n . Note that the parameter transformation proved in [6] allows to recover the general solution in the non-classical parameter range $\text{rpolar}_{\mu, \mu_c}^{\pm}(F)$ from $\text{rpolar}_{1,0}^{\pm}(F)$ by a rescaling of the deformation gradient, but we shall not detail this here.

Let us now turn towards the techniques which lie at the heart of these new results. The Euler-Lagrange equations for $W(R; D)$ have been derived in [7] and previously in [28]. They characterize the critical points as the solutions of a quadratic matrix equation on the manifold of rotations $\text{SO}(n)$. The key insight for the present development is a new approach to the analysis of the particular condition

$$(RD - \mathbb{1})^2 \in \text{Sym}(n). \quad (1.18)$$

Realizing that this is a *symmetric square condition*

$$(X(R))^2 = S \in \text{Sym}(n), \quad \text{where} \quad X(R) := RD - \mathbb{1} \in \mathbb{R}^{n \times n}, \quad (1.19)$$

one might suspect that the critical points of $W(R; D)$ are connected to real matrix square roots of real symmetric matrices. And indeed, the structure of the set of critical points of $W(R, D)$ can be revealed quite elegantly by a specific characterization of the set of real matrix square roots of real symmetric matrices. Note that this characterization, which is similar in spirit to the standard representation theorem for orthogonal matrices $\text{O}(n)$ as block matrices, seems not to be known in the literature. In Theorem 2.13 we show that the square roots of interest always admit a block-diagonal representation. This allows to reduce the problem from arbitrary dimension $n > 2$ into decoupled one- and two-dimensional subproblems which can then be solved independently. For example, we shall see that in $n = 3$, for a non-classical minimizer, we have to solve a one-dimensional and a two-dimensional subproblem. The one-dimensional problem determines the rotation axis of the optimal rotations, while the two-dimensional subproblem determines the optimal rotation angles.

There is a large body of work on matrix square roots. Mostly the literature focusses on the unique symmetric positive definite, the so-called principal matrix square root, of a symmetric positive definite matrix, see, e.g., the monographs [12, 16–18], or [10] for a compact introduction. Due to its numerous applications, the numerical approximation of the principal matrix square root is also an important theme, see, e.g., [14, 15]. For some recent developments and a geometric approach towards the numerical approximation of square roots, see [33] and references therein. Given the large body of work on the classical subject of matrix square roots, it may be somewhat surprising that the characterization stated in Theorem 2.13 seems not to be known. However, the characterization of the set of matrix square roots of symmetric matrices which we present in Section 2, was originally inspired by the theory of principal angles between linear subspaces, see, e.g., [9] and was motivated by our specific application.

The case of optimal rotations for recurring parameter values d_i , $i = 1, 2, 3$, in the diagonal parameter matrix $D \in \text{Diag}(n)$ has not been treated previously in [7], but is also accessible with the present approach. Note that this case corresponds to the special case of two or more equal principal stretches ν_i which is an important highly symmetric corner case in mechanics.

This paper is structured as follows: after this introduction in Section 1, we present a characterization of the full set of (possibly non-symmetric) real matrix square roots of (possibly non-positive definite) real symmetric matrices. More precisely, in Section 2, we construct an orthogonal change of basis which renders a matrix square root of this type block diagonal with blocks of size one or two. This block structure allows us to characterize the critical points in Section 3 for arbitrary dimension n . This leads to a sequence of decoupled one- and two-dimensional subproblems posed, however, on $O(1)$ and $O(2)$ and we continue with the solution of these subproblems in Section 4. In Section 5 we extract the globally energy-minimizing optimal Cosserat rotations from the set of critical points by a comparison of the realized energy levels. It turns out that the optimal rotations and energy levels are entirely consistent with previous results for $n = 2, 3$. We end with a short discussion of the present results in Section 6.

2 Representation of real matrix square roots of symmetric matrices

In the following sections, we profit from a characterization theorem for real matrix square roots of real symmetric matrices.

Definition 2.1. We say that $X \in \mathbb{R}^{n \times n}$ (not necessarily symmetric) is a real square root of a real symmetric matrix $S \in \text{Sym}(n)$, if it solves the quadratic matrix equation

$$X^2 = S \in \text{Sym}(n) .$$

Remark 2.2. It is insufficient for our purposes to restrict our attention to the unique symmetric positive definite principal matrix square root of S and we do not require its existence. Rather, we are interested in the set of all possible real matrix square roots of S . Moreover, we do not assume that S is positive semi-definite. The results of this section will be eventually applied to the union of the sets of square roots X as $S = X^2 \in \text{Sym}(n)$ varies, i.e., we are interested in characterizing matrices X which square to an unspecified symmetric matrix.

Example 2.3. The identity matrix $\mathbf{1}_2 \in \text{Sym}(2)$ has infinitely many real roots which are simply size two involution matrices. They fall into three distinct classes according to their trace.

$$X = \mathbf{1}, \quad X = -\mathbf{1}, \quad X \in \left\{ \begin{pmatrix} a & b \\ c & -a \end{pmatrix}, a^2 + bc = 1 \right\} . \quad (2.1)$$

Example 2.4. The roots of the negative identity matrix $-\mathbf{1}_2 \in \text{Sym}(2)$ are given by

$$X \in \left\{ \begin{pmatrix} a & b \\ c & -a \end{pmatrix}, a^2 + bc = -1 \right\} . \quad (2.2)$$

Example 2.5. A negative identity matrix of odd size $-\mathbf{1} \in \text{Sym}(2k-1)$ does not have real matrix square roots, since its determinant is negative.

More generally, let us record a criterion for a 2×2 real matrix $X \in \mathbb{R}^{2 \times 2}$ to be a square root of some real symmetric matrix $S \in \text{Sym}(2)$ which will be used in Section 4.

Lemma 2.6. A matrix $X \in \mathbb{R}^{2 \times 2}$ is a real matrix square root of a real symmetric matrix $S = X^2 \in \text{Sym}(2)$ if and only if $X \in \text{Sym}(2)$ or $\text{tr}[X] = 0$.

Proof. In dimension $n = 2$, the Cayley-Hamilton theorem implies that

$$S = X^2 = \text{tr}[X] X - (\det[X]) \mathbf{1} .$$

As the square $S = X^2$ is symmetric, the skew part of the right hand side vanishes

$$0 = \text{skew}(S) = \text{skew}(X^2) = \text{skew}(\text{tr}[X] X - (\det[X]) \mathbf{1}) = \text{tr}[X] \text{skew}(X) .$$

This finishes the argument. ■

Let $S \in \text{Sym}(n)$ be a real symmetric matrix. Then S has real eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_m$, $m \leq n$. Furthermore, there exists a decomposition of \mathbb{R}^n into mutually orthogonal eigenspaces $E_{\lambda_i}, i = 1, \dots, m$, of S which allows us to write

$$\mathbb{R}^n = E_{\lambda_1} \oplus E_{\lambda_2} \oplus \dots \oplus E_{\lambda_m} . \quad (2.3)$$

Note that the eigenspaces are preserved by S , i.e., $SE_{\lambda} = E_{\lambda}$, unless $\lambda = 0$ in which case $E_0 = \ker S$ and hence $SE_0 = \mathbf{0} \subseteq E_0$. This implies in particular that S does not mix its eigenspaces, i.e.,

$$\forall 1 \leq i \leq m : \quad SE_{\lambda_i} \subseteq E_{\lambda_i}$$

and, equivalently,

$$\forall 1 \leq i \neq j \leq m : \quad (SE_{\lambda_i}) \cap E_{\lambda_j} = \mathbf{0} .$$

This implies the existence of a basis of \mathbb{R}^n with transition matrix $T \in \text{GL}(n)$ in which S takes block diagonal form

$$T^{-1}ST = \tilde{S} = \text{diag}(\tilde{S}_{\lambda_1}, \tilde{S}_{\lambda_2}, \dots, \tilde{S}_{\lambda_m}) , \quad \tilde{S}_{\lambda_i} \in \mathbb{R}^{\dim E_{\lambda_i} \times \dim E_{\lambda_i}}, \quad i = 1, \dots, m .$$

It is a standard result from linear algebra, that in this particular case, each of the blocks \tilde{S}_{λ_i} , $i = 1, \dots, m$, is a multiple of a suitable identity matrix $\tilde{S}_{\lambda_i} = \lambda_i \mathbb{1}$ and so \tilde{S} is diagonal. Furthermore, if we choose an orthogonal basis for each eigenspace E_{λ_i} , individually, the change of basis matrix $T \in \text{O}(n)$ is orthogonal and we have

$$T^{-1}ST = \tilde{S} = \text{diag}(\underbrace{\lambda_1, \dots, \lambda_1}_{\dim E_{\lambda_1}}, \underbrace{\lambda_2, \dots, \lambda_2}_{\dim E_{\lambda_2}}, \dots, \underbrace{\lambda_m, \dots, \lambda_m}_{\dim E_{\lambda_m}}) .$$

Lemma 2.7 (Eigenspaces of $S = X^2 \in \text{Sym}(n)$ are not mixed by X). *Let $X \in \mathbb{R}^{n \times n}$ be a matrix square root of a symmetric matrix $S = X^2 \in \text{Sym}(n)$ and let $E_{\lambda_i}, i = 1, \dots, m$, denote the eigenspaces of S . Then X preserves the eigenspaces of S , i.e.,*

$$XE_{\lambda_i} \subseteq E_{\lambda_i} .$$

Proof. Let $v \in E_{\lambda}$, then

$$S(Xv) = X^3v = XSv = X\lambda v = \lambda(Xv) . \quad (2.4)$$

Hence, $Xv \in E_{\lambda}$ and since $v \in E_{\lambda}$ was arbitrary, we have $XE_{\lambda} \subseteq E_{\lambda}$. ■

Corollary 2.8. *Let $S \in \text{Sym}(n)$ and $T \in \text{O}(n)$ such that $\tilde{S} = T^{-1}ST$ is diagonal. Then any real matrix square root X of $S = X^2 \in \text{Sym}(n)$ is a block matrix of the form*

$$T^{-1}XT = \tilde{X} = \text{diag}(\tilde{X}_{\lambda_1}, \tilde{X}_{\lambda_2}, \dots, \tilde{X}_{\lambda_m}) , \quad \tilde{X}_{\lambda_i} \in \mathbb{R}^{\dim E_{\lambda_i} \times \dim E_{\lambda_i}}, \quad i = 1, \dots, m . \quad (2.5)$$

In particular $\tilde{S} = \tilde{X}^2 \in \text{Sym}(n)$ and

$$\tilde{X}_i^2 = \tilde{S}_i = \lambda_i \mathbb{1} .$$

Remark 2.9. *The preceding Corollary 2.8 reduces the subsequent characterization of real matrix square roots of symmetric matrices formidably, because it shows that it suffices to consider each of the X -invariant eigenspaces $E_{\lambda_i}, i = 1, \dots, m$, of S individually.*

We shortly recall the definition of the orthogonal complement V^{\perp} of a linear subspace $V \subseteq \mathbb{R}^n$,

$$V^{\perp} := \{w \in \mathbb{R}^n \mid w \perp V\} = \{w \in \mathbb{R}^n \mid \forall v \in V : \langle v, w \rangle = 0\} ,$$

which induces an orthogonal decomposition of $\mathbb{R}^n = V \oplus V^{\perp}$. In what follows, we exploit the well-known fact that for $Y \in \mathbb{R}^{n \times n}$,

$$YV^{\perp} \subseteq V^{\perp} \iff Y^TV \subseteq V . \quad (2.6)$$

Indeed, let $w \in V^{\perp}$, then $0 = \langle Yw, v \rangle = \langle w, Y^Tv \rangle$. Since the choice of $w \in V^{\perp}$ was arbitrary, we have that $Y^Tv \perp V^{\perp}$ which shows $Y^Tv \in V$, because $\mathbb{R}^n = V \oplus V^{\perp}$. The reverse implication is completely analogous.

Lemma 2.10 (Block lemma). *Let $S = \lambda \mathbb{1} \in \text{Sym}(n)$ be a multiple of an identity matrix of size $n \geq 1$. Then any real square root $Y \in \mathbb{R}^{n \times n}$ of $S = Y^2 = \lambda \mathbb{1} \in \text{Sym}(n)$ admits an orthogonal change of coordinates $T \in O(n)$ which renders it block-diagonal*

$$\tilde{Y} := T^{-1}YT = \text{diag}(\tilde{Y}_1, \tilde{Y}_2, \dots, \tilde{Y}_r) = \begin{pmatrix} \tilde{Y}_1 & 0 & \dots & 0 \\ 0 & \tilde{Y}_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \tilde{Y}_r \end{pmatrix}. \quad (2.7)$$

The square blocks $\tilde{Y}_i, i = 1, \dots, r$, are either of dimension 1 or 2 and satisfy $\tilde{Y}_i^2 = \lambda \mathbb{1}$.

Proof. The proof proceeds by induction on n . The base case of induction $n \in \{1, 2\}$ holds, since Y is already block-diagonal with blocks of size one or two. For the induction step let us assume that the statement holds for matrices of size $n - 1$ and $n - 2$.

Our strategy is to prove the existence of a one- or two-dimensional subspace V of \mathbb{R}^n such that both V and its orthogonal complement V^\perp are left invariant by Y , i.e.,

$$\dim V \in \{1, 2\}, \quad YV \subseteq V \quad \text{and} \quad YV^\perp \subseteq V^\perp. \quad (2.8)$$

Thus, if we pick an orthonormal basis of V and V^\perp , this is equivalent to the statement that orthogonal conjugates of Y and Y^T are block matrices of the form

$$Q^{-1}YQ = \left(\begin{array}{c|c} \tilde{Y}_1 & 0 \\ \hline 0 & Z \end{array} \right). \quad (2.9)$$

Since $(Q^{-1}YQ)^2 = Q^{-1}Y^2Q = \lambda \mathbb{1}$, we get $Z^2 = \lambda \mathbb{1}$, so by the induction assumption there exists T_0 such that

$$T_0^{-1}ZT_0 = \begin{pmatrix} \tilde{Z}_1 & 0 & \dots & 0 \\ 0 & \tilde{Z}_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \tilde{Z}_s \end{pmatrix}. \quad (2.10)$$

Then the orthogonal matrix

$$T = Q \begin{pmatrix} \mathbb{1} & 0 \\ 0 & T_0 \end{pmatrix} \in O(n) \quad (2.11)$$

satisfies

$$T^{-1}YT = \begin{pmatrix} \tilde{Y}_1 & 0 & \dots & 0 \\ 0 & \tilde{Z}_1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \tilde{Z}_s \end{pmatrix}, \quad (2.12)$$

which completes the induction step.

To finish the argument, we have to construct the invariant subspace V of \mathbb{R}^n of dimension 1 or 2.

Since $(Y^T)^2 = S^T = S = \lambda \mathbb{1}$, the symmetric matrices YY^T and Y^TY commute

$$(YY^T)(Y^TY) = Y(Y^T)^2Y = Y(\lambda \mathbb{1})Y = \lambda S = \lambda^2 \mathbb{1} = (Y^TY)(YY^T). \quad (2.13)$$

Therefore, the operators YY^T and Y^TY are simultaneously diagonalizable and we can find a common eigenvector w of both. Let us normalize w so that $\|w\| = 1$ and note that there exist values $\alpha, \beta \in [0, \infty)$ satisfying

$$Y^TYw = \alpha w \quad \text{and} \quad YY^Tw = \beta w. \quad (2.14)$$

Our next step is to choose the invariant subspace V . We have to distinguish several cases.

Case 1: $Yw \in \text{span}(\{w\})$, $Y^T w \in \text{span}(\{w\})$, in other words, w is an eigenvector of Y and Y^T . We select $V = \text{span}(\{w\})$ and construct an orthogonal matrix with first column given by $q_1 = w$, i.e.,

$$Q = (w|q_2|\dots|q_n) \in O(n). \quad (2.15)$$

An associated change of basis for Y and Y^T introduces the following zero patterns

$$Q^{-1}YQ = \left(\begin{array}{c|c} * & * \\ \hline 0 & \\ \vdots & \\ \vdots & * \\ \hline 0 & \end{array} \right) \quad \text{and} \quad Q^{-1}Y^TQ = \left(\begin{array}{c|c} * & * \\ \hline 0 & \\ \vdots & \\ \vdots & * \\ \hline 0 & \end{array} \right). \quad (2.16)$$

Since $Q^{-1}Y^TQ = (Q^{-1}YQ)^T$ these matrices are transposes of each other which implies that we obtain a block matrix of the form

$$Q^{-1}YQ = \left(\begin{array}{c|cccc} * & 0 & \cdot & \cdot & 0 \\ \hline 0 & & & & \\ \vdots & & & & \\ \vdots & & * & & \\ \hline 0 & & & & \end{array} \right), \quad (2.17)$$

which is of the form described in (2.9).

Case 2: $Yw \in \text{span}(\{w\})$, $Y^T w \notin \text{span}(\{w\})$, in other words w is an eigenvector of Y but not of Y^T . Consider the subspace $V = \text{span}(\{w, Y^T w\})$. Then the image of V under Y satisfies

$$YV = \text{span}(\{Yw, YY^T w\}) \subseteq \text{span}(\{w, w\}) \subseteq V \quad (2.18)$$

$$Y^T V = \text{span}(\{Y^T w, (Y^T)^2 w\}) \subseteq \text{span}(\{Y^T w, \lambda w\}) \subseteq V. \quad (2.19)$$

We now pick an orthonormal basis w_1, w_2 of $V = \text{span}(\{w_1, w_2\}) = \text{span}(\{w, Y^T w\})$ and extend it to an orthogonal matrix

$$Q = (w_1|w_2|q_3|\dots|q_n) \in O(n). \quad (2.20)$$

Then, similar to Case 1, an associated change of basis for Y and Y^T introduces a zero pattern

$$Q^{-1}YQ = \left(\begin{array}{cc|c} * & * & * \\ * & * & * \\ \hline 0 & 0 & \\ \vdots & \vdots & \\ \vdots & \vdots & * \\ \hline 0 & 0 & \end{array} \right) \quad \text{and} \quad Q^{-1}Y^TQ = \left(\begin{array}{cc|c} * & * & * \\ * & * & * \\ \hline 0 & 0 & \\ \vdots & \vdots & \\ \vdots & \vdots & * \\ \hline 0 & 0 & \end{array} \right). \quad (2.21)$$

As before, since $Q^{-1}Y^TQ = (Q^{-1}YQ)^T$ the two matrices are transposes of each other which creates a 2-block in the upper left corner

$$Q^{-1}YQ = \left(\begin{array}{cc|cccc} * & * & 0 & \dots & 0 \\ * & * & 0 & \dots & 0 \\ \hline 0 & 0 & & & \\ \vdots & \vdots & & & \\ \vdots & \vdots & & * & \\ \hline 0 & 0 & & & \end{array} \right), \quad (2.22)$$

which is of the form described in (2.9).

Case 3: $Yw \notin \text{span}(\{w\})$, in other words w is *not* an eigenvector of Y . We consider the subspace $V = \text{span}(\{w, Yw\})$. The inclusion

$$YV = \text{span}(\{Yw, Y^2w\}) = \text{span}(\{Yw, \lambda w\}) \subseteq V \quad (2.23)$$

is immediate. In order to prove the invariance $Y^TV \subseteq V$, we need to consider the following two subcases:

Case 3a: $\lambda \neq 0$. In this case Y and Y^T are invertible and so $Y^TYw = \alpha w$ with $\alpha > 0$. This allows us to express w as follows

$$\left(\frac{1}{\alpha}Y^TY\right)w = \frac{\alpha}{\alpha}w = w. \quad (2.24)$$

We have to compute

$$Y^TV = \text{span}(\{Y^Tw, Y^TYw\}) = \text{span}(\{Y^Tw, \alpha w\}). \quad (2.25)$$

To this end, we expand

$$Y^Tw = Y^T\left(\frac{1}{\alpha}Y^TY\right)w = \frac{1}{\alpha}(Y^2)^TYw = \frac{1}{\alpha}S^TYw = \frac{1}{\alpha}Y^3w = \frac{1}{\alpha}YSw = \frac{\lambda}{\alpha}Yw \in V \quad (2.26)$$

which shows that $Y^TV \subseteq V$.

Case 3b: $\lambda = 0$. Consider the product

$$(Y^TY)(YY^T)w = Y^TSY^Tw = S^2w = \lambda^2w = 0. \quad (2.27)$$

Since we also have, $Y^TYw = \alpha w$ and $YY^Tw = \beta w$, it follows that $(Y^TY)(YY^T)w = \alpha\beta w = 0$. Hence, $\alpha\beta = 0$. If $\beta = 0$, then

$$YY^Tw = 0 \implies \langle w, YY^Tw \rangle = 0 \implies \langle Y^Tw, Y^Tw \rangle = \|Y^Tw\|^2 = 0. \quad (2.28)$$

Since $Y^Tw = 0 \in V$, the subspace V is invariant under both Y and Y^T . The second case $\alpha = 0$ is not possible. To see this, we similarly compute

$$Y^TYw = 0 \implies Yw = 0 \quad (2.29)$$

which shows that w is an eigenvector of Y . This contradicts our assumptions for Case 3 (but note that this situation is handled in Case 1 or 2).

This completes the construction of the invariant subspace V and the proof of the lemma. ■

Remark 2.11. The case $\lambda > 0$ of Lemma 2.10 can be deduced from the theory of principal angles (see, e.g., [9]) for the eigenspaces of Y with eigenvalues $\sqrt{\lambda}$ and $-\sqrt{\lambda}$. We are not aware of a similar connection in the case $\lambda \leq 0$.

Remark 2.12. The condition on Y is also sufficient, i.e., any matrix with the described block structure is a real matrix square root of a symmetric matrix. It is also possible to show that solutions to $Y^2 = \lambda \mathbb{1}$ exist if and only if $\lambda \geq 0$, or n is even.

We are now ready to formulate the main result of this section.

Theorem 2.13. For any real square root $X \in \mathbb{R}^{n \times n}$ of a symmetric matrix $S = X^2 \in \text{Sym}(n)$ there exists an orthogonal change of coordinates $T \in O(n)$ such that the transformed square root is block-diagonal

$$\tilde{X} := T^{-1}XT = \text{diag}(\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_r) = \begin{pmatrix} \tilde{X}_1 & 0 & \dots & 0 \\ 0 & \tilde{X}_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \tilde{X}_r \end{pmatrix} \quad (2.30)$$

with square blocks $\tilde{X}_j, j = 1, \dots, r$, that are either of size 1 or 2. Each block \tilde{X}_j is a real square root of a multiple of an identity matrix $\mathbb{1}$, i.e.,

$$\tilde{X}_j^2 = \mu_j \mathbb{1}, \quad \mu_j \in \mathbb{R}.$$

Proof. This is a straightforward application of the block lemma to each eigenblock of $S = X^2$. ■

Remark 2.14. Each eigenspace E_{λ_i} of S in Lemma 2.7 is possibly decomposed into multiple subspaces by the Lemma 2.10. As a result, the eigenvalues $\mu_j, 1 \leq j \leq r$, of \tilde{X}_j^2 in Theorem 2.13 are equal to the eigenvalues $\lambda_i, 1 \leq i \leq m$, in the notation of Lemma 2.7 with, possibly, different indices. Several μ_j in the statement of Theorem 2.13 may be equal to the same λ_i in the sense of Lemma 2.7. For example, we might have the following

$$(\mu_1, \mu_2, \mu_3, \mu_4, \mu_5) = (\lambda_1, \lambda_1, \lambda_2, \lambda_3, \lambda_3).$$

Remark 2.15. An equivalent reformulation of the theorem is the following. For a matrix X whose square is symmetric, there exists a decomposition of \mathbb{R}^n into an orthogonal direct sum of X -invariant subspaces V_i of dimension one or two such that X^2 is a multiple of the identity matrix on each V_i . The list of columns of the change of basis matrix $T \in O(n)$ in Theorem 2.13 is obtained by concatenation of orthonormal bases of V_i . Note that each V_i is also invariant under X^T .

Remark 2.16. Given X and S the decomposition into invariant subspaces is not unique. In particular, a subspace of dimension two can sometimes be further decomposed into two one-dimensional subspaces.

Remark 2.17. Our description of matrices which square to a symmetric matrix is similar in spirit to the well-known characterization of orthogonal matrices. Every orthogonal matrix is orthogonally conjugated to a block diagonal matrix with blocks of size one and two, see, e.g., [11, Thm. 12.5, p. 354].

3 Critical points of the Cosserat shear-stretch energy

In this section we investigate the critical points $R \in SO(n)$ of the energy subject to minimization in Problem 1.1

$$W(R; D) = \|\text{sym}(RD - \mathbb{1})\|^2 \quad (3.1)$$

for a given diagonal matrix $D = \text{diag}(d_1, \dots, d_n)$. We give a complete description of them under some mild assumptions on the diagonal entries d_i based on a criterion which we derive next.

The Lie algebra $\mathfrak{so}(n)$ of the matrix group of rotations $SO(n)$ is given by the subspace of skew-symmetric matrices, i.e., $\mathfrak{so}(n) = \text{Skew}(n)$. Furthermore, the Frobenius inner product gives rise to the orthogonal decomposition

$$\mathbb{R}^{n \times n} = \text{Sym}(n) \oplus_{\perp} \text{Skew}(n) = \text{Sym}(n) \oplus_{\perp} \mathfrak{so}(n)$$

of real square matrices into the subspaces of symmetric matrices and skew-symmetric matrices.

Lemma 3.1 (Symmetric square condition). *Let $D := \text{diag}(d_1, \dots, d_n) \in \mathbb{R}^{n \times n}$ be a diagonal matrix. A rotation $R \in SO(n)$ is a critical point of the function*

$$W(R; D) := \|\text{sym}(RD - \mathbb{1})\|^2$$

if and only if the matrix $(RD - \mathbb{1})^2$ is symmetric.

Proof. In order to compute critical points in the submanifold $SO(n) \subset \mathbb{R}^{n \times n}$, we have to locate zeroes of the tangent mapping $dW : TSO(n) \rightarrow T\mathbb{R}^{n \times n} \cong \mathbb{R}^{n \times n}$. To this end, we compute the derivatives of the energy $W(R; D)$ along a family of smooth curves

$$c_A : (-\varepsilon, \varepsilon) \rightarrow SO(n), \quad c_A(t) := \exp(tA)R \in SO(n), \quad A \in \mathfrak{so}(n), \quad (3.2)$$

in the manifold of rotations. The right-trivialization of the tangent space at $R \in \text{SO}(n)$ allows to identify $T_R \text{SO}(n) = \mathfrak{so}(n) \cdot R = \text{Skew}(n) \cdot R$ and so we can always express a tangent vector $\xi \in T_R \text{SO}(n)$ in the form $\xi = AR \in \text{Skew}(n) \cdot R$. This family of curves satisfies

$$\forall \xi = AR \in T_R \text{SO}(n) : \quad \left. \frac{d}{dt} \right|_{t=0} c_A(t) = AR = \xi . \quad (3.3)$$

Thus, for every possible tangent direction $\xi = AR \in T_R \text{SO}(n)$, there is precisely one curve of the family which emanates from $R \in \text{SO}(n)$ into this direction ξ .

A rotation R is a critical point of the energy $W(R; D)$ if and only if

$$\forall A \in \mathfrak{so}(n) : \quad \left. \frac{d}{dt} (W \circ c_A)(t) \right|_{t=0} = 0 .$$

It is well-known that the matrix exponential is given by $(\mathbb{1} + tA)$ to first order in t and we write $\exp(tA) \sim (\mathbb{1} + tA)$. Thus, by the chain rule, we also have

$$(W \circ c_A)(t) \sim (W \circ (\mathbb{1} + tA)R)(t) .$$

We expand the expression

$$\begin{aligned} W \circ (\mathbb{1} + tA)R &= \|\text{sym}((\mathbb{1} + tA)RD - \mathbb{1})\|^2 = \|\text{sym}(RD - \mathbb{1}) + t \text{sym}(ARD)\|^2 \\ &= \|\text{sym}(RD - \mathbb{1})\|^2 + 2t \langle \text{sym}(RD - \mathbb{1}), \text{sym}(ARD) \rangle + t^2 \|\text{sym}(ARD)\|^2 \end{aligned}$$

and obtain the expression for the first derivative dW from the term linear in t . In other words

$$\left. \frac{d}{dt} (W \circ c_A)(t) \right|_{t=0} = 2 \langle \text{sym}(RD - \mathbb{1}), \text{sym}(ARD) \rangle . \quad (3.4)$$

Hence, a point R is a critical point for the energy W if and only if it satisfies

$$\forall A \in \mathfrak{so}(n) : \quad \text{sym}(RD - \mathbb{1}) \perp \text{sym}(ARD) .$$

Since $\text{Sym}(n) \perp \text{Skew}(n)$, we may add $\text{skew}(ARD)$ on the right hand side which gives us the equivalent condition

$$\forall A \in \mathfrak{so}(n) : \quad \text{sym}(RD - \mathbb{1}) \perp ARD .$$

Expanding the definition of the Frobenius inner product, we find

$$\begin{aligned} 0 &= \langle \text{sym}(RD - \mathbb{1}), ARD \rangle = \text{tr} [\text{sym}(RD - \mathbb{1})^T ARD] = \text{tr} [RD \text{sym}(RD - \mathbb{1})A] \\ &= \langle \text{sym}(RD - \mathbb{1})(RD)^T, A \rangle . \end{aligned} \quad (3.5)$$

Since this condition must hold for all $A \in \text{Skew}(n)$, it follows that

$$\text{sym}(RD - \mathbb{1})DR^T \in \text{Sym}(n) .$$

We now multiply by a factor of 2 and expand the definition of $\text{sym}(X) := \frac{1}{2}(X + X^T)$ which leads us to

$$\begin{aligned} 2 \text{sym}(RD - \mathbb{1})DR^T &= (RD + DR^T - 2\mathbb{1})DR^T = RD^2R^T + (DR^T)^2 - 2DR^T \\ &= (DR^T - \mathbb{1})^2 + (RD^2R^T - \mathbb{1}) . \end{aligned} \quad (3.6)$$

The second term on the right hand side is always symmetric and the effective condition for a critical point is thus

$$(DR^T - \mathbb{1})^2 \in \text{Sym}(n) . \quad (3.7)$$

Finally, observing that symmetry is invariant under transposition, we conclude that

$$((DR^T - \mathbb{1})^2)^T = (RD - \mathbb{1})^2 \in \text{Sym}(n) \quad (3.8)$$

is a sufficient and necessary condition for a critical point $R \in \text{SO}(n)$ of $W(R; D)$. ■

Remark 3.2. We immediately observe that $R = \mathbb{1}$ solves the condition (3.8) and is always a critical point of the energy $W(R; D) := \|\text{sym}(RD - \mathbb{1})\|^2$. However, in general, it will not be the global minimizer.

Remark 3.3 (Critical points and real matrix square roots). *Introducing the notation*

$$X(R) := RD - \mathbb{1} ,$$

we see that $R \in \text{SO}(n)$ is a critical point of $W(R; D)$ if and only if

$$S(R) := (X(R))^2 = (RD - \mathbb{1})^2 \in \text{Sym}(n) .$$

In other words, for any critical point R of $W(R; D)$, $X(R) = RD - \mathbb{1}$ is a real square root of a real symmetric matrix. This connects the set of critical points for $W(R; D)$ to our previously derived characterization of the set of real square roots of a real symmetric matrix stated in Theorem 2.13.

Our next step is to apply Theorem 2.13 and Remark 2.15 to the special case $X(R) = RD - \mathbb{1}$. As we shall see, this implies quite restrictive conditions on $R \in \text{SO}(n)$.

Let us make the following assumption on the diagonal matrix D .

Assumption 3.4. *The entries of the diagonal matrix $D = \text{diag}(d_1, \dots, d_n)$, which parametrizes the energy $W(R; D)$, do not vanish and do not cancel each other additively, i.e.,*

$$d_i \neq 0 \quad \text{and} \quad d_i + d_j \neq 0, \quad 1 \leq i, j \leq n .$$

This ensures that $\ker(D) = \mathbf{0}$ and that any D^2 -invariant subspace is also D -invariant. Note that if the entries of $D = \text{diag}(d_1, \dots, d_n)$ are positive, this assumption is satisfied. For the original problem in Cosserat theory which stimulated the present work [6–8], the entries of D are the singular values $\nu_i > 0$, $i = 1, \dots, n$, of the deformation gradient $F \in \text{GL}^+(n)$.

The following insight is a key to our discussion.

Lemma 3.5 (Simultaneous invariance of R and D). *Suppose that the eigenvalues of D satisfy the above assumption. Let V be a subspace invariant under $X(R) = RD - \mathbb{1}$, such that V^\perp is also invariant under $X(R)$. Then both V and V^\perp are invariant under D and R .*

Proof. Recall first that V and V^\perp are both invariant under $X(R)$ if and only if V is invariant under both $X(R)$ and $X(R)^T$; cf. (2.6).

By assumption the subspace V is invariant under both $RD = X + \mathbb{1}$ and $(RD)^T = DR^T = X^T + \mathbb{1}$. Therefore

$$D^2V = (DR^T)(RD)V \subseteq (DR^T)V \subseteq V .$$

From the assumption on D , we have $DV \subseteq V$. Since D has only nonzero eigenvalues D is invertible and so $DV = V$. It follows that

$$RDV \subseteq V \quad \implies \quad RV \subseteq V .$$

Since $R \in \text{SO}(n)$ is invertible, we have $RV = V$. Reversing the roles of V and V^\perp , we can apply the same argument to V^\perp . ■

By Theorem 2.13, as phrased in Remark 2.15, there exists a sequence of pairwise orthogonal vector spaces V_i , $i = 1, \dots, r$, with $1 \leq \dim V_i \leq 2$ which decompose $\mathbb{R}^n = V_1 \oplus_\perp V_2 \oplus_\perp \dots \oplus_\perp V_r$. These correspond to a block-diagonal representation of $X(R) := RD - \mathbb{1}$. The existence of an associated orthogonal change of basis matrix $T \in O(n)$ is also assured by Theorem 2.13. Furthermore, by Lemma 3.5, both R and D are also block-diagonal with respect to this choice of basis. This means, in particular, that *any solution* R satisfying the symmetric square condition $(X(R))^2 = (RD - \mathbb{1})^2 \in \text{Sym}(n)$ admits a block-diagonal representation. Since this condition characterizes

the critical points by Lemma 3.1, *any critical point* of $W(R; D)$ admits a representation in block-diagonal form

$$\tilde{R} = T^{-1}RT = \text{diag}(\tilde{R}_1, \dots, \tilde{R}_r) = \begin{pmatrix} \tilde{R}_1 & 0 & \dots & 0 \\ 0 & \tilde{R}_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \tilde{R}_r \end{pmatrix} \in O(n) \subset \mathbb{R}^{n \times n}, \quad (3.9)$$

where the blocks on the diagonal satisfy $\tilde{R}_i \in O(n_i)$, $i = 1, \dots, r$, with $n_i \in \{1, 2\}$ and $\sum_i n_i = n$. This shows that, in the basis provided by $T \in O(n)$, any critical point $R \in O(n)$ can be constructed from solutions $\tilde{R}_i \in O(n_i)$ of one- and two-dimensional subproblems

$$\left(\tilde{X}(\tilde{R}_i) \right)^2 \in \text{Sym}(n_i). \quad (3.10)$$

Note that these subproblems are now posed on the space of *orthogonal*, rather than *special orthogonal* matrices.

Assumption 3.6. *For the purpose of clarity of exposition, we make an additional, stronger assumption on the diagonal matrix $D = \text{diag}(d_1, \dots, d_n)$, namely*

$$d_1 > d_2 > \dots > d_n > 0.$$

The slightly more general case of possibly non-distinct positive entries d_i can be treated similarly which we will indicate in running commentary.

Remark 3.7 (Implications of D -invariance). *Under the Assumption 3.6, the D -invariance of the subspaces V_i shown in Lemma 3.5 implies a strong restriction: the V_i are necessarily coordinate subspaces in the standard basis of \mathbb{R}^n . Thus, we can index these data by partitions of the index set $\{1, \dots, n\}$ into disjoint subsets of size one or two. Furthermore, by picking a standard coordinate basis for each V_i , we can ensure that the change of basis matrix $T \in O(n)$ is a permutation matrix.*

We summarize that this particular structure allows to reduce the optimization Problem 1.1 to a finite list of decoupled one- and two-dimensional subproblems. However, we have to consider minimization with respect to *orthogonal* matrices $R \in O(n)$ instead of $R \in \text{SO}(n)$. This will be the content of the next section.

4 Analysis of the decoupled subproblems

Let $I \subseteq \{1, \dots, n\}$ be a one-element subset $\{i\}$ or a two-element subset $\{i, j\}$ and let D_I be the associated restriction of D given by

$$\begin{cases} D_I := (d_i), & \text{if } I = \{i\}, \\ D_I := \begin{pmatrix} d_i & 0 \\ 0 & d_j \end{pmatrix}, & \text{if } I = \{i, j\}. \end{cases}$$

In this section we solve for critical points of the function

$$W(R_I; D_I) := \|\text{sym}(R_I D_I - \mathbf{1})\|^2$$

for $R_I \in O(|I|)$ and compute the corresponding critical values. This corresponds to the solution of the decoupled lower-dimensional subproblems as described in the previous section.

Theorem 4.1 (Critical points: size one). *For $I = \{i\}$ we have the submatrix $D_I = (d_i)$ and $R_I = \pm \mathbf{1} = (\pm 1)$. The realized critical energy levels are*

$$W(+\mathbf{1}; D_I) = (d_i - 1)^2 \quad \text{and} \quad W(-\mathbf{1}; D_I) = (d_i + 1)^2. \quad (4.1)$$

Proof. There are only two orthogonal matrices in dimension one and the result is immediate. ■

For the case $|I| = 2$, we consider the two separate cases $\det[R_I] = 1$ and $\det[R_I] = -1$.

Theorem 4.2 (Critical points: size two and positive determinant). *The critical points R_I with $\det[R_I] = 1$ are described as follows. For any values d_i and d_j the matrices $R_I = \pm \mathbb{1}$ are critical points with the critical values $(d_i - 1)^2 + (d_j - 1)^2$ and $(d_i + 1)^2 + (d_j + 1)^2$, respectively. In addition, if $d_i + d_j > 2$, there are two non-diagonal critical points*

$$R_I = \begin{pmatrix} \cos \alpha & -\sin \alpha \\ \sin \alpha & \cos \alpha \end{pmatrix}, \quad \text{with} \quad \cos \alpha = \frac{2}{d_i + d_j} \quad (4.2)$$

which attain the same critical value

$$W(R_I; D_I) = \frac{1}{2}(d_i - d_j)^2. \quad (4.3)$$

Proof. By Lemma 3.1 R_I is a critical point if and only if $(R_I D_I - \mathbb{1})^2$ is symmetric. We may thus apply Lemma 2.6 which implies $R_I D_I - \mathbb{1} \in \text{Sym}(2)$ or $\text{tr}[R_I D_I - \mathbb{1}] = 0$. Using the explicit representation

$$R_I = \begin{pmatrix} \cos \alpha & -\sin \alpha \\ \sin \alpha & \cos \alpha \end{pmatrix},$$

the symmetry condition $R_I D_I - \mathbb{1} \in \text{Sym}(2)$ is equivalent to $(d_i + d_j) \sin \alpha = 0$ which has two solutions $R_I = \pm \mathbb{1}$. The trace condition $\text{tr}[R_I D_I - \mathbb{1}] = 0$ is equivalent to $(d_i + d_j) \cos \alpha = 2$ which can be solved for α if and only if $d_i + d_j \geq 2$. It gives rise to two non-diagonal solutions if and only if $d_i + d_j > 2$.

In the first case $R_I = \pm \mathbb{1}$, the critical values are immediately seen to be $(d_i - 1)^2 + (d_j - 1)^2$ and $(d_i + 1)^2 + (d_j + 1)^2$, respectively.

In the second case, the critical values are calculated as follows. Observing that

$$\text{sym}(R_I D_I - \mathbb{1}) = \begin{pmatrix} d_i \cos \alpha - 1 & \frac{1}{2}(d_j - d_i) \sin \alpha \\ \frac{1}{2}(d_j - d_i) \sin \alpha & d_j \cos \alpha - 1 \end{pmatrix} \quad (4.4)$$

we use $(d_i + d_j) \cos \alpha = 2$ to get

$$\begin{aligned} \|\text{sym}(R_I D_I - \mathbb{1})\|^2 &= (d_i \cos \alpha - 1)^2 + (d_j \cos \alpha - 1)^2 + \frac{1}{2}(d_j - d_i)^2 \sin^2 \alpha \\ &= (d_i^2 + d_j^2) \cos^2 \alpha - 2(d_i + d_j) \cos \alpha + 2 + \frac{1}{2}(d_j - d_i)^2 (1 - \cos^2 \alpha) \\ &= \frac{1}{2}(d_j - d_i)^2 + \frac{1}{2}(d_i + d_j)^2 \cos^2 \alpha - 2(d_i + d_j) \cos \alpha + 2 \\ &= \frac{1}{2}(d_i - d_j)^2 + 2 - 4 + 2 = \frac{1}{2}(d_i - d_j)^2. \end{aligned} \quad (4.5)$$

This shows the claim. ■

Theorem 4.3 (Critical points: size two and negative determinant). *The critical points R_I with $\det[R_I] = -1$ are described as follows. For any values d_i and d_j the diagonal matrices $R_I = \pm \text{diag}(1, -1)$ are critical points with the critical values $(d_i - 1)^2 + (d_j + 1)^2$ and $(d_i + 1)^2 + (d_j - 1)^2$, respectively. In addition, for $|d_i - d_j| > 2$, there are two non-diagonal critical points*

$$R_I = \begin{pmatrix} \cos \alpha & \sin \alpha \\ \sin \alpha & -\cos \alpha \end{pmatrix}, \quad \text{with} \quad \cos \alpha = \frac{2}{|d_i - d_j|}, \quad (4.6)$$

which attain the same critical value

$$W(R_I; D_I) = \frac{1}{2}(d_i + d_j)^2. \quad (4.7)$$

Proof. By Lemma 3.1 R_I is a critical point if and only if $(R_I D_I - \mathbb{1})^2$ is symmetric. We may thus apply Lemma 2.6 which implies $R_I D_I - \mathbb{1} \in \text{Sym}(2)$ or $\text{tr}[R_I D_I - \mathbb{1}] = 0$. Using the explicit representation

$$R_I = \begin{pmatrix} \cos \alpha & \sin \alpha \\ \sin \alpha & -\cos \alpha \end{pmatrix}$$

the symmetry condition $R_I D_I - \mathbb{1} \in \text{Sym}(2)$ is equivalent to

$$(d_i - d_j) \sin \alpha = 0 \quad (4.8)$$

which has two solutions $R_I = \pm \text{diag}(1, -1)$ since $d_i \neq d_j$ due to Assumption 3.6. The trace condition $\text{tr}[R_I D_I - \mathbb{1}] = 0$ is equivalent to $(d_i - d_j) \cos \alpha = 2$ which can be solved for α if and only if $|d_i - d_j| \geq 2$. Thus there are two non-diagonal solutions if and only if $|d_i - d_j| > 2$.

In the first case $R_I = \pm \text{diag}(1, -1)$, the critical values are immediately seen to be $(d_i - 1)^2 + (d_j + 1)^2$ and $(d_i + 1)^2 + (d_j - 1)^2$, respectively.

In the second case, the critical values are calculated as follows. Observing that

$$\text{sym}(R_I D_I - \mathbb{1}) = \begin{pmatrix} d_i \cos \alpha - 1 & \frac{1}{2}(d_i + d_j) \sin \alpha \\ \frac{1}{2}(d_i + d_j) \sin \alpha & -d_j \cos \alpha - 1 \end{pmatrix} \quad (4.9)$$

we use $|d_i - d_j| \cos \alpha = 2$ to get

$$\begin{aligned} \|\text{sym}(R_I D_I - \mathbb{1})\|^2 &= (d_i \cos \alpha - 1)^2 + (d_j \cos \alpha + 1)^2 + \frac{1}{2}(d_i + d_j)^2 \sin^2 \alpha \\ &= (d_i^2 + d_j^2) \cos^2 \alpha - 2(d_i - d_j) \cos \alpha + 2 + \frac{1}{2}(d_i + d_j)^2 (1 - \cos^2 \alpha) \\ &= \frac{1}{2}(d_i + d_j)^2 + \frac{1}{2}(d_i - d_j)^2 \cos^2 \alpha - 2(d_i - d_j) \cos \alpha + 2 \\ &= \frac{1}{2}(d_i + d_j)^2 + 2 - 4 + 2 = \frac{1}{2}(d_i + d_j)^2. \end{aligned} \quad (4.10)$$

This shows the claim. ■

Remark 4.4 (The positive choice $\det[R_I] = +1$ minimizes energy). *A direct comparison of the energy levels realized by the different choices for the determinant of R_I is instructive. Summarizing our preceding results, we have for $|I| = 1$, i.e., for a block of size one*

$$\det[R_I] = +1 \quad \mapsto \quad (d_i - 1)^2, \quad (4.11)$$

$$\det[R_I] = -1 \quad \mapsto \quad (d_i + 1)^2 \geq (d_i - 1)^2. \quad (4.12)$$

Similarly, for $|I| = 2$, i.e., for a block of size two, we obtain

$$\det[R_I] = +1 \quad \mapsto \quad \frac{1}{2}(d_i - d_j)^2, \quad (4.13)$$

$$\det[R_I] = -1 \quad \mapsto \quad \frac{1}{2}(d_i + d_j)^2 \geq \frac{1}{2}(d_i - d_j)^2. \quad (4.14)$$

The estimates follow from our Assumption 3.6 on the entries $d_i > 0$ of the diagonal matrix $D > 0$.

Remark 4.5. *The diagonal critical points $R_I = \pm \mathbb{1}$ and $R_I = \pm \text{diag}(1, -1)$ reduce to size one blocks (or index subsets $|I| = 1$) in the block decomposition (3.9).*

Remark 4.6 (On non-distinct entries of D). *If we relax the Assumption 3.6 and allow for*

$$d_1 \geq d_2 \geq \dots \geq d_n > 0$$

then there are degenerate critical points with $\det[R_I] = -1$ if and only if $d_i = d_j$. The corresponding critical value is the same as that realized by the diagonal matrices $\pm \text{diag}(1, -1)$.

5 Global minimization of the Cosserat shear-stretch energy

Combining the results of the two preceding sections, we can now describe the critical values of the Cosserat shear-stretch energy $W(R; D)$ which are attained at the critical points. The main result of this section is a procedure (algorithm) which traverses the set of critical points in a way that reduces the energy at every step of the procedure and finally terminates in the subset of global minimizers.

Technically, we label the critical points by certain partitions of the index set $\{1, \dots, n\}$ containing only subsets I with one or two elements. In the last section, we have seen that the subsets I and a choice of sign for $\det[R_I]$ uniquely characterize a critical point $R \in \text{SO}(n)$.

Let us give an outline of the energy-decreasing traversal strategy starting from a given labeling partition (i.e., critical point):

1. Choose the positive sign $\det[R_I] = +1$ for each subset of the partition (cf. Remark 4.4 and Remark 5.3).
2. Disentangle all overlapping blocks for $n > 3$ (cf. Lemma 5.8).
3. Successively shift all 2×2 -blocks to the lowest possible index, i.e., collect the blocks of size two as close to the upper left corner of the matrix R as possible (cf. Lemma 5.4).
4. Introduce as many additional 2×2 -blocks by joining adjacent blocks of size 1 as the constraint $d_i + d_j > 2$ allows (cf. Lemma 5.4).

At the end of this section, we provide an Example 5.12.

The next theorem expresses the value of $W(R; D)$ realized by a critical point in terms of the labeling partition and choice of determinants $\det[R_I]$ which characterize it.

Theorem 5.1 (Characterization of critical points and values). *Under the Assumption 3.6 on the entries $d_1 > d_2 > \dots > d_n > 0$ of $D \in \text{Diag}(n)$, the critical points $R \in \text{SO}(n)$ can be classified according to partitions of the index set $\{1, \dots, n\}$ into subsets of size one or two and choices of signs for the determinant $\det[R_I]$ for each subset I . The subsets of size two $I = \{i, j\}$ satisfy*

$$\begin{cases} d_i + d_j > 2, & \det[R_I] = +1, \quad \text{and} \\ |d_i - d_j| > 2, & \det[R_I] = -1. \end{cases}$$

The critical values are given by

$$W(R; D) = \sum_{\substack{I=\{i\} \\ \det[R_I]=1}} (d_i - 1)^2 + \sum_{\substack{I=\{i\} \\ \det[R_I]=-1}} (d_i + 1)^2 + \sum_{\substack{I=\{i,j\} \\ \det[R_I]=1}} \frac{1}{2}(d_i - d_j)^2 + \sum_{\substack{I=\{i,j\} \\ \det[R_I]=-1}} \frac{1}{2}(d_i + d_j)^2.$$

Proof. A suitable partition of the index set $\{1, \dots, n\}$ can be constructed as detailed in Section 3. The contributions of the subsets I of size one and two are given by the theorems of Section 4. It suffices to consider the non-diagonal critical points for the subproblems of size two, because the diagonal cases can be accounted for by splitting the subset $I = \{i, j\}$ into two subsets $\{i\}$ and $\{j\}$ of size one, see Remark 4.5. ■

Remark 5.2 (On non-distinct entries of D). *If we relax the Assumption 3.6 and allow for*

$$d_1 \geq d_2 \geq \dots \geq d_n > 0$$

then the D - and R -invariant subspaces V_i are not necessarily coordinate subspaces. This produces non-isolated critical points but does not change the formula for the critical values.

In order to compute the global minimizers $R \in \text{SO}(n)$ for the Cosserat shear-stretch energy $W(R; D)$, we have to compare all the critical values which correspond to the different partitions and choices of the signs of the determinants in the statement of Theorem 5.1. In what follows, we prove the various reduction steps.

Remark 5.3. Notice that $|d_i - d_j| > 2$ implies that $d_i + d_j > 2$. Therefore, it is always possible to replace negative determinant choices by positive ones. In the process the value of $W(R; D)$ is reduced. Therefore, if R is a critical point which is a global minimizer of $\|\text{sym}(RD - \mathbb{1})\|^2$, it only contains R_I with determinant $\det[R_I] = 1$.

This allows us to assume that $\det[R_I] = 1$ for all subsets I without any loss of generality.

The following lemma shows that blocks of size two are always favored *whenever they exist*.

Lemma 5.4 (Comparison lemma). *If $d_i + d_j > 2$ then the difference between the critical values of $W(R; D)$ corresponding to the choice of a size two subset $I = \{i, j\}$ as compared to the choice of two size one subsets $\{i\}, \{j\}$ is given by*

$$-\frac{1}{2}(d_i + d_j - 2)^2.$$

Proof. We subtract the corresponding contributions of the subsets and simplify

$$\frac{1}{2}(d_i - d_j)^2 - (d_i - 1)^2 - (d_j - 1)^2 = -\frac{1}{2}(d_i + d_j - 2)^2.$$

This proves the claim. ■

Let us rewrite $W(R; D)$ in a slightly different form in order to distill the contributions of the size two blocks in the partition.

Corollary 5.5. *For the choices of $\det[R_I] = 1$ there holds*

$$W(R; D) = \|\text{sym}(RD - \mathbb{1})\|^2 = \sum_{i=1}^n (d_i - 1)^2 - \frac{1}{2} \sum_{I=\{i,j\}} (d_i + d_j - 2)^2.$$

Proof. The first term in the formula is the value realized by $W(R; D)$ for the trivial partition into n subsets of size one. By virtue of the Comparison Lemma 5.4 each block of size two reduces the critical value by the amount $\frac{1}{2}(d_i + d_j - 2)^2$. ■

Let us now consider the case of dimension $n = 3$ explicitly in order to prepare the exposition of the higher dimensional case.

Theorem 5.6. *Let $d_1 > d_2 > d_3 > 0$. If $d_1 + d_2 \leq 2$ then the global minimum of $W(R; D)$ occurs at $R = \mathbb{1}$ and is given by*

$$W(R; D) = (d_1 - 1)^2 + (d_2 - 1)^2 + (d_3 - 1)^2.$$

If $d_1 + d_2 > 2$ then the global minimum is realized by either of two critical points of the form

$$R = \begin{pmatrix} \cos \alpha & -\sin \alpha & 0 \\ \sin \alpha & \cos \alpha & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{with} \quad (d_1 + d_2) \cos \alpha = 2.$$

In this case the global minimum is

$$W(R; D) = (d_1 - 1)^2 + (d_2 - 1)^2 + (d_3 - 1)^2 - \frac{1}{2}(d_1 + d_2 - 2)^2 = \frac{1}{2}(d_1 - d_2)^2 + (d_3 - 1)^2.$$

Proof. If $d_1 + d_2 \leq 2$ then $d_i + d_j \leq 2$ for all index pairs (i, j) and there are no blocks of size two at the global minimum. If $d_1 + d_2 > 2$ then the choice of partition $\{1, 2\} \sqcup \{3\}$ is admissible. Corollary 5.5 shows that this is always favorable compared to the partition into three size one subsets $\{1\} \sqcup \{2\} \sqcup \{3\}$. Whether or not other size two subsets are admissible according to the inequalities $d_i + d_j > 2$, the partition $\{1, 2\} \sqcup \{3\}$ is always optimal. This follows from the ordering $d_1 > d_2 > d_3 > 0$ which implies that the partition-dependent term $\frac{1}{2}(d_i + d_j - 2)^2$ in Corollary 5.5 is maximized for $I = \{i, j\} = \{1, 2\}$. ■

In mechanics one often makes assumptions on the symmetries of the deformation gradient $F \in \text{GL}^+(n)$ and it may then have non-distinct singular values $\nu_i = \nu_j$, $i \neq j$. It is thus of interest to investigate the three-dimensional case with $d_1 \geq d_2 \geq d_3 \geq 0$.

Remark 5.7 (On non-distinct entries of D). Assume $d_1 \geq d_2 \geq d_3 > 0$. Our results imply the following.

If $d_1 + d_2 \leq 2$, then all V_i are of dimension 1. Since the restriction of a given minimizer R to each V_i satisfies $R|_{V_i} = \mathbb{1}$, we see that $R = \mathbb{1}$. The global minimum of the Cosserat shear-stretch energy is given by

$$W(R; D) = (d_1 - 1)^2 + (d_2 - 1)^2 + (d_3 - 1)^2.$$

If $d_1 + d_2 > 2$, then for a global minimizer R there is a one-dimensional R -invariant subspace which is also D -invariant with associated eigenvalue d_3 . Therefore, R is a rotation with axis in the d_3 -eigenspace of D . The rotation angle satisfies the relation $(d_1 + d_2) \cos \alpha = 2$ and the global minimum of the energy is given by

$$W(R; D) = (d_1 - 1)^2 + (d_2 - 1)^2 + (d_3 - 1)^2 - \frac{1}{2}(d_1 + d_2 - 2)^2 = \frac{1}{2}(d_1 - d_2)^2 + (d_3 - 1)^2.$$

This case further splits into several subcases all realizing the same energy level according to the multiplicity of the eigenvalue d_3 :

If $d_1 \geq d_2 > d_3$, i.e., the multiplicity of d_3 is one, then there are two isolated global minimizers which are rotations with rotation angle $\arccos(2/(d_1 + d_2))$ with respect to either of the two half-axes in $\text{span}(\{e_3\})$ (as in the case of distinct entries of D discussed in Theorem 5.6).

If $d_1 > d_2 = d_3$, i.e., the multiplicity of d_3 is two, then the global minimizers R form a one-dimensional family of rotations with rotation angle $\arccos(2/(d_1 + d_2))$ and rotation half-axes in the d_3 -eigenplane $\text{span}(\{e_2, e_3\})$ of D .

If $d_1 = d_2 = d_3$, i.e., the multiplicity of d_3 is three, then there is a two-dimensional family of global minimizers R which are rotations with rotation angle $\arccos(2/(d_1 + d_2))$ about arbitrary half-axes in \mathbb{R}^3 .

It is interesting that the set of global minimizers is *connected* in the last two cases where $d_2 = d_3$. This allows for a continuous transition between minimizers with opposite half-axes which are inverses of each other.

To study the global minimizers for the Cosserat shear-stretch energy in arbitrary dimension $n \geq 4$, we need to investigate the relative location of the size two subsets of the partition.

Lemma 5.8. Let $R \in \text{SO}(n)$ be a global minimizer for $W(R; D)$. Then R cannot contain overlapping size two subsets, i.e., $I = \{i_1, i_4\}$, $J = \{i_2, i_3\}$, with $i_1 < i_2 < i_3 < i_4$.

Proof. We assume that R is a global minimizer corresponding to a partition containing two overlapping subsets as described above and derive a contradiction.

It suffices to consider the case $i_1 = 1, i_2 = 2, i_3 = 3$ and $i_4 = 4$ with the general case being completely analogous. We recall the ordering $d_1 > d_2 > d_3 > d_4 > 0$.

There are two cases to consider:

Case 1: $d_3 + d_4 > 2$. In this case, we can consider another critical point \mathring{R} corresponding to the partition $\{1, 2\} \sqcup \{3, 4\}$ instead of $\{1, 4\} \sqcup \{2, 3\}$. By Corollary 5.5 we have

$$\begin{aligned} W(R; D) - W(\mathring{R}; D) &= \frac{1}{2}(d_1 + d_2 - 2)^2 + \frac{1}{2}(d_3 + d_4 - 2)^2 - \frac{1}{2}(d_1 + d_4 - 2)^2 - \frac{1}{2}(d_2 + d_3 - 2)^2 \\ &= d_1 d_2 + d_3 d_4 - d_1 d_4 - d_2 d_3 = (d_1 - d_3)(d_2 - d_4) > 0. \end{aligned}$$

Thus R is not a global minimum of $W(R; D)$.

Case 2: $d_3 + d_4 \leq 2$. In this case, we can not have the size two subset $\{3, 4\}$. However, it is

possible to decrease the value of $W(R; D)$ by choosing another critical point \mathring{R} corresponding to the partition $\{1, 2\} \sqcup \{3\} \sqcup \{4\}$ instead of $\{1, 4\} \sqcup \{2, 3\}$. By Corollary 5.5 we have

$$\begin{aligned} W(R; D) - W(\mathring{R}; D) &= \frac{1}{2}(d_1 + d_2 - 2)^2 - \frac{1}{2}(d_1 + d_4 - 2)^2 - \frac{1}{2}(d_2 + d_3 - 2)^2 \\ &\geq \frac{1}{2}(d_1 + d_2 - 2)^2 - \frac{1}{2}(d_1 + (2 - d_3) - 2)^2 - \frac{1}{2}(d_2 + d_3 - 2)^2 \\ &= \frac{1}{2}(d_1 + d_2 - 2)^2 - \frac{1}{2}(d_1 - d_3)^2 - \frac{1}{2}(d_2 + d_3 - 2)^2 \\ &= (d_1 - d_3)(d_2 + d_3 - 2) > 0. \end{aligned}$$

In the first inequality we use the fact that for $d_1 + d_4 \geq 2$ the function $(d_1 + d_4 - 2)^2$ is increasing in d_4 and $d_4 \leq 2 - d_3$ by assumption. This shows that R is not a global minimum of $W(R; D)$. We arrive at a contradiction in both cases which proves the statement. ■

We are now ready to state and prove the general n -dimensional case.

Theorem 5.9. *Let $d_1 > d_2 > \dots > d_n > 0$ be the entries of D . Let us fix the maximum k for which $d_{2k-1} + d_{2k} > 2$. Any global minimizer $R \in \text{SO}(n)$ corresponds to the partition of the form*

$$\{1, 2\} \sqcup \{3, 4\} \sqcup \dots \sqcup \{2k-1, 2k\} \sqcup \{2k+1\} \sqcup \dots \sqcup \{n\}$$

and the global minimum of $W(R; D)$ is given by

$$\begin{aligned} W^{\text{red}}(D) &:= \min_{R \in \text{SO}(n)} \|\text{sym}(RD - \mathbb{1})\|^2 = \sum_{i=1}^n (d_i - 1)^2 - \frac{1}{2} \sum_{i=1}^k (d_{2i-1} + d_{2i} - 2)^2 \\ &= \frac{1}{2} \sum_{i=1}^k (d_{2i-1} - d_{2i})^2 + \sum_{i=2k+1}^n (d_i - 1)^2. \end{aligned}$$

Proof. Lemma 5.8 shows that a global minimizer $R \in \text{SO}(n)$ can not have a partition with *overlapping* size two subsets. As in the proof of Theorem 5.6 (the $n = 3$ case) we can decrease the value of $W(R; D)$ by shifting down the indices of all size two subsets as far as possible. Therefore the optimal partition is of the form

$$\{1, 2\} \sqcup \{3, 4\} \sqcup \dots \sqcup \{2l-1, 2l\} \sqcup \{2l+1\} \sqcup \dots \sqcup \{n\}$$

for some $l \leq k$. By Corollary 5.5 the global minimum is realized by the critical points corresponding to the maximal possible choice $l = k$. The value of $W(R; D)$ at a global minimizer is computed by inserting the corresponding optimal partition into Theorem 5.1 and Corollary 5.5. ■

Remark 5.10. *The number of global minimizers in the above theorem is 2^k , where k is the number of blocks of size two in the preceding characterization of a global minimizer as a block diagonal matrix. All global minimizers are block diagonal similar to the $n = 3$ case (Theorem 5.6).*

Remark 5.11 (On non-distinct entries of D). *If we relax the Assumption 3.6 and allow for*

$$d_1 \geq d_2 \geq \dots \geq d_n > 0$$

then the global minimizers may or may not be isolated. The formula for the reduced energy as stated in Theorem 5.9 is, however, not affected.

The following example illustrates the energy-minimizing traversal of critical points which always terminates in a global minimizer.

Example 5.12. *Let $D = \text{diag}(4, 2, 1, \frac{1}{2}, \frac{1}{4})$. Theorem 5.1 shows that the critical points can be characterized by certain partitions¹ of the index set $\{1, 2, 3, 4, 5\}$ and a choice of a sign*

¹More precisely, a labeling partition uniquely characterizes sets of critical points which generate the same critical value. A block of size two, for example, characterizes two different symmetric solutions corresponding to the choice of sign for the rotation angle α . Both choices, however, yield the same value for the energy.

for each subset I of the partition. Thus, we introduce the convenient notation of a pair of a subset and a sign (I, \pm) , where the sign encodes a possible choice for the determinant $\det[R_I] = \pm 1$.

Setup: We consider a critical point $R^{(0)}$ corresponding to the labeling partition

$$\mathcal{P}^{(0)} = \{(\{1\}, +), (\{2, 5\}, -), (\{3\}, -), (\{4\}, -)\}. \quad (5.1)$$

Note that $d_2 + d_5 = 2 + \frac{1}{4} > 2$, i.e., the 2×2 -block corresponding to $I = \{2, 5\}$ exists, as required for a valid partition characterizing a critical point $R^{(0)}$. The corresponding critical value of the summation formula in the statement of Theorem 5.1 is given by

$$\begin{aligned} W^{(0)} = W(R^{(0)}; D) &= \underbrace{(4-1)^2}_{(\{1\}, +)} + \underbrace{(1+1)^2}_{(\{3\}, -)} + \underbrace{\left(1 + \frac{1}{2}\right)^2}_{(\{4\}, -)} + \underbrace{\frac{1}{2} \left(2 + \frac{1}{4}\right)^2}_{(\{2, 5\}, -)} \\ &= \frac{569}{32} \approx 17.78. \end{aligned} \quad (5.2)$$

Step 1 (Choice of positive sign): We consistently choose the positive sign for the determinant in the labeling partition which gives

$$\mathcal{P}^{(1)} = \{(\{1, 5\}, +), (\{2\}, +), (\{3\}, +), (\{4\}, +)\}. \quad (5.3)$$

This updated partition characterizes a different critical point $R^{(1)}$ realizing a lower energy level

$$\begin{aligned} W^{(1)} = W(R^{(1)}; D) &= \underbrace{(4-1)^2}_{(\{1\}, +)} + \underbrace{(1-1)^2}_{(\{3\}, +)} + \underbrace{\left(1 - \frac{1}{2}\right)^2}_{(\{4\}, +)} + \underbrace{\frac{1}{2} \left(2 - \frac{1}{4}\right)^2}_{(\{2, 5\}, +)} \\ &= \frac{345}{32} \approx 10.28. \end{aligned} \quad (5.4)$$

Step 2 (Disentanglement): The next step of the procedure is to remove overlap of 2×2 -blocks. In our example, we only have one such block and there is nothing to do, i.e., $\mathcal{P}^{(2)} = \mathcal{P}^{(1)}$.

Step 3 (Index shift): We now decrement the indices of the 2×2 -blocks as much as possible, i.e., we string them together starting in the upper left corner. Shifting the $\{2, 5\}$ -block to $\{1, 2\}$, we obtain the following new partition

$$\mathcal{P}^{(3)} = \{(\{1, 2\}, +), (\{3\}, +), (\{4\}, +), (\{5\}, +)\}. \quad (5.5)$$

The energy level realized by a corresponding critical point $R^{(3)}$ is

$$\begin{aligned} W^{(3)} = W(R^{(3)}; D) &= \underbrace{(1-1)^2}_{(\{3\}, +)} + \underbrace{\left(1 - \frac{1}{2}\right)^2}_{(\{4\}, +)} + \underbrace{\left(1 - \frac{1}{4}\right)^2}_{(\{5\}, +)} + \underbrace{\frac{1}{2} (4-2)^2}_{(\{1, 2\}, +)} \\ &= \frac{45}{16} \approx 2.81. \end{aligned} \quad (5.6)$$

Step 4 (Exhaustion by 2×2 -blocks): In this step, we try to create as many 2×2 -blocks as possible. We first locate the pair of subsets of size one with minimal indices which is $(\{3\}, \{4\})$. Since $d_3 + d_4 = 1 + \frac{1}{2} \leq 2$, no further 2×2 -block exists. Thus, $\mathcal{P}^{(4)} = \mathcal{P}^{(3)}$.

Result: The finally obtained labeling partition

$$\mathcal{P} = \mathcal{P}^{(4)} = \{(\{1, 2\}, +), (\{3\}, +), (\{4\}, +), (\{5\}, +)\} \quad (5.7)$$

characterizes a global minimizer. With the notation of Theorem 5.9 the maximal number of 2×2 -blocks is $k = 1$ and we have $2^k = 2$ global minimizers of the form

$$\text{rpolar}(D) = \begin{pmatrix} \begin{array}{cc|ccc} \cos \alpha_1 & -\sin \alpha_1 & 0 & 0 & 0 \\ \sin \alpha_1 & \cos \alpha_1 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{array} \end{pmatrix}, \quad \text{with} \quad \cos(\alpha_1) = \frac{2}{d_1 + d_2} = \frac{1}{3}. \quad (5.8)$$

Inserting the global minimizers into the energy, we obtain the reduced energy

$$W^{\text{red}}(D) := W(\text{rpolar}(D); D) = \frac{45}{16} \approx 2.81. \quad (5.9)$$

Just to give a comparison, the identity matrix $\mathbb{1} \in \text{SO}(n)$ realizes the energy level

$$\begin{aligned} W(\mathbb{1}; D) &= \underbrace{(4-1)^2}_{(\{1\}, +)} + \underbrace{(2-1)^2}_{(\{2\}, +)} + \underbrace{(1-1)^2}_{(\{3\}, +)} + \underbrace{\left(1 - \frac{1}{2}\right)^2}_{(\{4\}, +)} + \underbrace{\left(1 - \frac{1}{4}\right)^2}_{(\{5\}, +)} \\ &= \frac{173}{16} \approx 10.81. \end{aligned} \quad (5.10)$$

Thus, the identity $\mathbb{1} \in \text{SO}(n)$ is not a global minimizer.

Remark 5.13 (Optimality of $\mathbb{1}$). *Our results imply that the identity matrix $\mathbb{1} \in \text{SO}(n)$ is globally optimal for $W(R; D)$ with $D > 0$, if and only if there exists no 2×2 -block with a positive choice of $\det[R_I]$, i.e.,*

$$\max_{1 \leq i \neq j \leq n} (d_i + d_j) \leq 2.$$

This corresponds to the tension-compression asymmetry described in [6–8] for dimensions $n = 2, 3$.

6 Discussion

For the sake of clarity of exposition, we have restricted our attention to the case of a diagonal and positive definite parameter matrix $D > 0$, i.e., $d_i > 0$. Our technical approach, however, readily carries over to the more general case $d_i \neq 0$ with minor modifications. The construction

$$\left\| \text{sym} \left\{ \left[R \left(\frac{\mathbb{1}}{\left| \begin{smallmatrix} \mathbb{1} \\ -\mathbb{1} \end{smallmatrix} \right|} \right) \right] \left[\left(\frac{\mathbb{1}}{\left| \begin{smallmatrix} \mathbb{1} \\ -\mathbb{1} \end{smallmatrix} \right|} \right) D \right] - \mathbb{1} \right\} \right\|^2 \quad (6.1)$$

allows to reduce such a parameter matrix D to $|D| := \text{diag}(|d_1|, \dots, |d_n|) > 0$ which is positive definite. Note that the minimization must then be carried out in the appropriate connected component of the orthogonal matrices $\text{O}(n)$. We also expect that the degenerate case where some $d_i = 0$ can be handled with our techniques as well.

The matrix group of rotations $\text{SO}(3)$ equipped with its natural bi-invariant Riemannian metric

$$g(\xi, \eta)|_R := g(R^T \xi, R^T \eta)|_{\mathbb{1}} := \langle R^T \xi, R^T \eta \rangle = \langle \xi, \eta \rangle \quad (6.2)$$

is a Riemannian manifold $(\text{SO}(3), g)$. In [24], the dynamics of the following Riemannian gradient flow² was investigated

$$R^T \dot{R} = \text{skew}(R^T D) \iff \dot{R} = -\text{grad} \left(\frac{1}{2} \|RD - \mathbb{1}\|^2 \right). \quad (6.3)$$

The flow (6.3) converges to $R = \mathbb{1}$ for appropriate initial conditions which is consistent with Grioli's theorem; cf. Section 1. Similarly, one can study the gradient flow for the energy $\frac{1}{2} \|\text{sym}(RD - \mathbb{1})\|^2$ given by

$$R^T \dot{R} = -\frac{1}{2} \text{skew}((R^T D - \mathbb{1})^2) \iff \dot{R} = -\text{grad} \left(\frac{1}{2} \|\text{sym}(RD - \mathbb{1})\|^2 \right). \quad (6.4)$$

Our present results on critical points of $W(R; D)$ determines the possible asymptotic solutions for the gradient flow (6.4). A characterization of *local* minimizers is currently missing. For example, it is not clear whether every local minimizer is automatically a global minimizer which holds in dimension $n = 2$. It seems likely, that this holds in $n = 3$ as well. The classification of local extrema of $W(R; D)$ is a completely open question in $n \geq 4$.

Acknowledgments: Lev Borisov was partially supported by NSF grant DMS-1201466. Andreas Fischle was supported by German Research Foundation (DFG) grant SA2130/2-1 and, previously, partially supported by DFG grant NE902/2-1 (also: SCHR570/6-1).

²For an introductory exposition of gradient flows on Riemannian manifolds, see, e.g., [21].

References

- [1] L. Borisov, P. Neff, S. Sra, and C. Thiel. The sum of squared logarithms inequality in arbitrary dimensions. *arXiv preprint arXiv:1508.04039*, 2015. <http://arxiv.org/abs/1508.04039>, to appear in Lin. Alg. Appl.
- [2] E. Cosserat and F. Cosserat. *Théorie des corps déformables*. Librairie Scientifique A. Hermann et Fils (engl. translation by D. Delphenich 2007, available online at https://www.uni-due.de/~hm0014/Cosserat_files/Cosserat09_eng.pdf), reprint 2009 by Hermann Librairie Scientifique, ISBN 978 27056 6920 1, Paris, 1909.
- [3] V. A. Eremeyev, L. P. Lebedev, and H. Altenbach. *Foundations of Micropolar Mechanics*. Springer, 2012.
- [4] A. C. Eringen. *Microcontinuum Field Theories. Vol. I: Foundations and Solids*. Springer, 1999.
- [5] A. Fischle. The planar Cosserat model: minimization of the shear energy on $SO(2)$ and relations to geometric function theory. (diploma thesis). 2007. (available online: http://www.uni-due.de/~hm0014/Supervision_files/dipl_final_online.pdf).
- [6] A. Fischle and P. Neff. The geometrically nonlinear Cosserat micropolar shear–stretch energy. Part I: A general parameter reduction formula and energy-minimizing microrotations in 2D. *arXiv preprint arXiv:1507.05480*, 2015. <http://arxiv.org/abs/1507.05480>, to appear in Z. angew. Math. Mechanik.
- [7] A. Fischle and P. Neff. The geometrically nonlinear Cosserat micropolar shear–stretch energy. Part II: Non-classical energy-minimizing microrotations in 3D and their computational validation. *arXiv preprint arXiv:1509.06236*, 2015. <http://arxiv.org/pdf/1509.06236v1>.
- [8] A. Fischle, P. Neff, and D. Raabe. The relaxed-polar mechanism of locally optimal Cosserat rotations for an idealized nanoindentation and comparison with 3D-EBSD experiments. *arXiv preprint arXiv:1603.06633*, 2016. <http://arxiv.org/abs/1603.06633>.
- [9] A. Galántai. *Projectors and projection methods*, volume 6. Springer Science & Business Media, 2013.
- [10] J. Gallier. Logarithms and square roots of real matrices. *arXiv preprint arXiv:0805.0245*, 2008. <http://arxiv.org/abs/0805.0245>.
- [11] J. Gallier. *Geometric methods and applications: for computer science and engineering*, volume 38. Springer Science & Business Media, 2. edition, 2011.
- [12] F. R. Gantmacher. *The Theory of Matrices*, volume I. AMS Chelsea, 1. edition.
- [13] G. Grioli. Una proprietà di minimo nella cinematica delle deformazioni finite. *Boll. Un. Math. Ital.*, 2:252–255, 1940.
- [14] N. J. Higham. Newton’s method for the matrix square root. *Mathematics of Computation*, 46(174):537–549, 1986.
- [15] N. J. Higham. Computing real square roots of a real matrix. *Lin. Alg. Appl.*, 88:405–430, 1987.
- [16] N. J. Higham. *Functions of Matrices: Theory and Computation*. SIAM, Philadelphia, PA, USA, 2008.
- [17] R.A. Horn and C.R. Johnson. *Matrix Analysis*. Cambridge University Press, New York, 1985.
- [18] R.A. Horn and C.R. Johnson. *Topics in Matrix Analysis*. Cambridge University Press, New York, 1991.
- [19] J. Jeong, H. Ramézani, I. Münch, and P. Neff. A numerical study for linear isotropic Cosserat elasticity with conformally invariant curvature. *Z. Angew. Math. Mech.*, 89(7):552–569, 2009.
- [20] J. Lankeit, P. Neff, and Y. Nakatsukasa. The minimization of matrix logarithms: On a fundamental property of the unitary polar factor. *Lin. Alg. Appl.*, 449:28–42, 2014.
- [21] J. M. Lee. *Introduction to Smooth Manifolds*. Graduate Texts in Mathematics. Springer, 2002.
- [22] G. A. Maugin. On the structure of the theory of polar elasticity. *R. Soc. Lond. Philos. Trans. Ser. A Math. Phys. Eng. Sci.*, 356(1741):1367–1395, 1998.
- [23] P. Neff. Existence of minimizers for a geometrically exact cosserat solid. *PAMM*, 4(1):548–549, 2004.
- [24] P. Neff. Local existence and uniqueness for quasistatic finite plasticity with grain boundary relaxation. *Quart. Appl. Math.*, 63:88–116, 2005.
- [25] P. Neff. The Cosserat couple modulus for continuous solids is zero viz the linearized Cauchy-stress tensor is symmetric. *Z. Angew. Math. Mech.*, 86:892–912, 2006.
- [26] P. Neff, B. Eidel, and R. J. Martin. Geometry of logarithmic strain measures in solid mechanics. *arXiv preprint arXiv:1505.02203*, 2015. <http://arxiv.org/pdf/1505.02203v1> to appear in Arch. Rat. Mech. Analysis.
- [27] P. Neff, B. Eidel, F. Osterbrink, and R. Martin. A Riemannian approach to strain measures in nonlinear elasticity. *C. R. Acad. Sci. Paris (Mecanique)*, 342(4):254–257, 2014.
- [28] P. Neff, A. Fischle, and I. Münch. Symmetric Cauchy-stresses do not imply symmetric Biot-strains in weak formulations of isotropic hyperelasticity with rotational degrees of freedom. *Acta Mech.*, 197:19–30, 2008.
- [29] P. Neff and J. Jeong. A new paradigm: the linear isotropic Cosserat model with conformally invariant curvature energy. *Z. Angew. Math. Mech.*, 89(2):107–122, 2009.
- [30] P. Neff, J. Jeong, and A. Fischle. Stable identification of linear isotropic Cosserat parameters: bounded stiffness in bending and torsion implies conformal invariance of curvature. *Acta Mech.*, 211(3-4):237–249, 2010.
- [31] P. Neff, J. Lankeit, and A. Madeo. On Grioli’s minimum property and its relation to Cauchy’s polar decomposition. *Int. J. Engng. Sci.*, 80:209–217, 2014.
- [32] P. Neff, Y. Nakatsukasa, and A. Fischle. A logarithmic minimization property of the unitary polar factor in the spectral and Frobenius norms. *SIAM J. Matrix Anal. Appl.*, 35(3):1132–1154, 2014.
- [33] S. Sra. On the matrix square root via geometric optimization. *arXiv preprint arXiv:1507.08366*, 2015. <http://arxiv.org/abs/1507.08366>.