

# GLOBAL INJECTIVITY IN SECOND-GRADIENT NONLINEAR ELASTICITY AND ITS APPROXIMATION WITH PENALTY TERMS

STEFAN KRÖMER AND JAN VALDMAN

ABSTRACT. We present a new penalty term approximating the Ciarlet-Nečas condition (global invertibility of deformations) as a soft constraint for hyperelastic materials. For non-simple materials including a suitable higher order term in the elastic energy, we prove that the penalized functionals converge to the original functional subject to the Ciarlet-Nečas condition. Moreover, the penalization can be chosen in such a way that for all low energy deformations, self-interpenetration is completely avoided already at all sufficiently small finite values of the penalization parameter. We also present numerical experiments in 2d illustrating our theoretical results.

Nonlinear elasticity, local injectivity, global injectivity, nonsimple materials, Ciarlet-Necas-condition, approximation

## 1. INTRODUCTION

Nonlinear elasticity models the behavior of a solid body subject to relatively strong external forces causing large deformations, albeit not large enough to cause irreversible damage. For a general introduction to the topic, we refer to [25, 9, 3]. It is clear that the deformation of such a body, the map  $y : \Omega \rightarrow \mathbb{R}^d$  mapping the “reference configuration”  $\Omega \subset \mathbb{R}^d$  to a deformed state. Here, typically  $d = 3$  and  $\Omega$  is the domain occupied by the elastic body in its stress-free state before external forces are applied. Clearly, any realistic deformation should always be injective, i.e., the body should not interpenetrate itself in any way. Here, we focus on the framework of hyperelasticity, that is, models where the body does not dissipate energy when deformed, instead simply storing in an internal elastic energy whose mathematical description as a functional depending on  $y$  fully determines its response to external forces. The existence of energy minimizers in this framework was pioneered by John Ball [4]. This theory allows us to enforce a weak form of local injectivity of the deformation, i.e.,  $\det \nabla y > 0$  a.e. in  $\Omega$  for all deformation with finite internal energy, by using an elastic energy density which becomes infinite as  $\det \nabla y \rightarrow 0^+$ . However, having

---

*Date:* January 10, 2019.

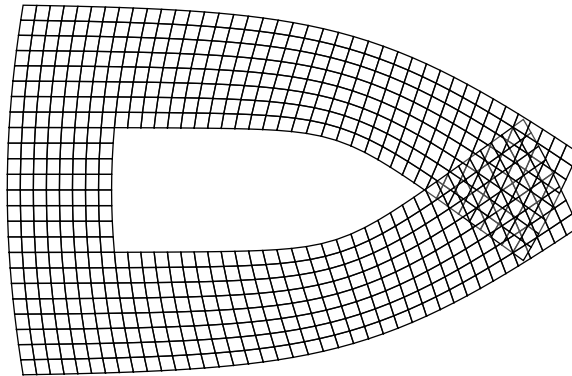


FIGURE 1. Example of deformed domain with overlapping ends.

$\det \nabla y > 0$  a.e. is still not enough to ensure local injectivity everywhere, and globally, a loss of injectivity by, say, two different ends of the body overlapping each other (see Figure 1), is not automatically prevented [2]. For instance, in addition to the positive determinant, to derive local (and global) injectivity, the result of [2] requires global injectivity of  $y$  on  $\partial\Omega$  as a prerequisite, a feature hard to prove and difficult to handle as a constraint (unless it is already given in form of a Dirichlet condition). On the other hand, the positive determinant implies local injectivity in a neighborhood of almost every point in  $\Omega$  given quite natural assumptions on the regularity of the deformation [13]. A direct approach to local injectivity *everywhere* was provided in [17], where a uniform positive lower bound on  $\det \nabla u$  is derived, but this only works in the framework of so-called non-simple materials, where the internal elastic energy is assumed to contain a suitable term involving second order derivatives  $\nabla^2 u$  (or higher order, and the coercivity conditions in [17] automatically entail  $y \in C^1$  by embedding), as opposed to classical hyperelasticity for which the energy only depends on  $\nabla y$ . For further information on the topic of locally injective deformations, we also refer to [5, 14]. In addition, there is literature available discussing settings that allow for cavitation (see, e.g., [18] and the references therein), but we will rule that out by assumption in this article (coercivity in  $W^{1,p}$  with  $p > d$ ). The standard constraint nowadays used to ensure global injectivity of the deformation is the Ciarlet-Nečas condition [10]:

$$\int_{\Omega} \det(\nabla y) \, dx = |y(\Omega)|. \quad (1.1)$$

The numerical treatment of this condition is not well understood, however. In particular, to our knowledge there is so far no example where (1.1) is enforced as a hard constraint on the discrete level in a numerical scheme with provable convergence.

Artificially including terms with derivatives of second or higher order in the energy, possibly multiplied with a small parameter, can also serve as regularization, and among other things, the results of [17] can then be used to obtain local injectivity. However, with such a regularization, there is a risk of a Lavrentiev phenomenon occurring, that is, the minimal energy of the regularized problem might remain strictly above the minimal energy for the original problem without the higher order term, even if the regularization parameter converges to zero. In particular, it is known that the infimum of the internal energy in spaces with higher integrability of  $\nabla y$  may be too large [15]. Similar issues might occur when discretizing, as the typical finite element spaces are all subsets of  $W^{1,\infty}$ . For simple materials, i.e., models without higher order terms, a way out was shown in [21], using an artificially introduced auxiliary field, and similarly in [20]<sup>1</sup>. For the variant with a higher order term also treated in [20], a Lavrentiev phenomenon cannot be ruled out.

In this article, we further elaborate on the approach of [20] in presence of a higher order term

$$\sigma \int_{\Omega} |D^2 y|^s dx \quad (\sigma > 0 \text{ is a fixed parameter}) \quad (1.2)$$

using soft constraints, i.e., everywhere finite terms in the energy depending on control parameters  $\varepsilon_i$ , converging to the singular determinant term and the Ciarlet-Nečas condition (1.1), respectively, as  $\varepsilon_i \rightarrow 0$ ,  $i = 1, 2$ . In particular, the penalization term for the latter in [20] is defined as

$$\frac{1}{\varepsilon_2} \left( \int_{\Omega} \det(\nabla y) dx - |y(\Omega)| \right). \quad (1.3)$$

After a discretizing with mesh size  $h$ , it is shown in [20] that the energy minima converge for these approximations (the energies even  $\Gamma$ -converge) as  $(\varepsilon_1, \varepsilon_2, h) \rightarrow (0, 0, 0)$  in the scaling regime  $\frac{h}{\varepsilon_1} \rightarrow 0$ . They also show that there is also convergence even for  $\sigma = 0$ , but then the suitable scaling regime, while shown to exist, is not explicitly known and therefore effectively impossible to exploit in practice.

How to actually compute (1.3) is not explained in [20], however, and since this term is nonlocal and, in particular, not a standard integral functional, this can in fact be quite problematic to implement. This is further aggravated by the presence of a higher order term like (1.2) which rules out the (direct) use of piecewise affine finite elements. An alternative, more accessible penalization was recently proposed and studied in [6], but only for beams (effectively 1d). Here, after precisely outlining the model we work with in Section 2, we show that instead

---

<sup>1</sup>also for more complicated models including plasticity

of (1.3), alternatives in the form of double integrals can be used (Section 3), for example

$$\frac{1}{\varepsilon_2} \int_{(\Omega \times \Omega)} \frac{1}{\varepsilon_2^d} [|\tilde{x} - x| - \varepsilon_2^{-1} |y(\tilde{x}) - y(x)|]^+ d(x, \tilde{x}). \quad (1.4)$$

A more general class of such terms is defined in (3.3). We show that such terms also lead to a limiting constraint equivalent to the standard Ciarlet-Nečas condition (Theorem 3.3) while the energy minima converge as before, at least for the regularized problem with fixed  $\sigma > 0$  (Theorem 4.6). This new kind of soft constraint makes discrete computations easier and can even be handled by standard packages (although inefficiently). Moreover, as we will see below, we have enough freedom to choose something going along well with particular features of the finite elements used, or to enforce other physically desirable properties like global injectivity even on the discrete level, and not just in the limit (Corollary 3.8). Our theoretical results are complemented by numerical computations carried out for two example problems in 2d, presented in Section 5.

## 2. THE MODEL AND STRUCTURAL ASSUMPTIONS

**2.1. The elastic energy.** As it is standard, we assume that for a deformation  $y \in W^{1,p}(\Omega; \mathbb{R}^d)$ ,  $p > d$ , the elastic energy has the form

$$E^{el}(y) := \int_{\Omega} W(x, \nabla y) dx, \quad y \text{ satisfies (1.1),}$$

where

$$W : \Omega \times \mathbb{R}^{d \times d} \rightarrow \mathbb{R} \cup \{+\infty\} \text{ is a Carathéodory function,} \quad (2.1)$$

i.e.,  $W(x, F)$  is measurable in  $x$  and continuous in  $F$ . Moreover, for all  $F \in \mathbb{R}^{d \times d}$  and all  $x \in \Omega$ ,

$$\begin{aligned} W(x, F) &= +\infty && \text{if } \det F \leq 0, \\ W(x, F) &\geq c_1 (|F|^p + (\det F)^{-q}) - c_2 && \text{if } \det F > 0, \end{aligned} \quad (2.2)$$

with constants  $q > d$  (which is necessary for (2.5) below),  $c_1 > 0$  and  $c_2 \geq 0$ . In addition, we assume that  $f$  is polyconvex, i.e.,

$$\begin{aligned} W(x, F) &= h(x, m(F)), \text{ with a function } h \text{ such that} \\ &h(x, \cdot) \text{ is convex for each } x, \end{aligned} \quad (2.3)$$

where  $m(F) \in \mathbb{R}^{n(d)}$ ,  $n(d) = \sum_{k=1}^d \binom{d}{k}^2$ , denotes the collection of all minors of  $F$ , i.e., all  $k \times k$  sub-determinants with  $1 \leq k \leq d$ . For instance,  $m(F) = (F, \det F) \in \mathbb{R}^5$  for  $d = 2$  and  $m(F) = (F, \text{cof } F, \det F) \in \mathbb{R}^{19}$  for  $d = 3$ . Here, for any  $d$ ,  $\text{cof } F \in \mathbb{R}^{d \times d}$  denotes cofactor matrix so that  $F^{-1} = (\text{cof } F)^T (\det F)^{-1}$  whenever  $F$  is invertible.

*Remark 2.1.* Due to [4, 10] (for a proof of a related lower semicontinuity property of the terms in the Ciarlet-Nečas condition also see [20]), with the assumptions (2.1)–(2.3),  $E^{el}$  always has a minimizer  $y^*$  in  $W^{1,p}(\Omega; \mathbb{R}^d)$ . Like all states with finite energy, it must satisfy

$$\det \nabla y^* > 0 \quad \text{a.e. in } \Omega.$$

**2.2. Approximation including penalization and higher order terms.** Our regularized approximation of  $E^{el}$  is defined as follows:

$$E_{\varepsilon,\sigma}(y) := E_{\varepsilon_1}^{el}(y) + E_{\varepsilon_2}^{CN}(y) + E_{\sigma}^{reg}(y), \quad \varepsilon = (\varepsilon_1, \varepsilon_2).$$

Here, the elastic energy reads

$$E_{\varepsilon_1}^{el}(y) := \int_{\Omega} W_{\varepsilon_1}(x, \nabla y) dx$$

where  $W_{\varepsilon_1} : \Omega \times \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$  can be any everywhere finite approximation of  $W$  such that

$W_{\varepsilon_1}$  is a Carathéodory function and polyconvex;

$$\begin{aligned} c_3(|F|^p + \max\{\varepsilon_1, \det F\}^{-q}) - c_4 \\ \leq W_{\varepsilon_1}(x, F) \leq W(x, F) + \varepsilon_1 \quad \text{for all } F \in \mathbb{R}^{d \times d}; \end{aligned} \quad (2.4)$$

$$W_{\varepsilon_1}(x, F) \geq W(x, F) - \varepsilon_1 \quad \text{if } |F| \leq \frac{1}{\varepsilon_1} \text{ and } \det(F) \geq \varepsilon_1.$$

with constants  $c_3 > 0$ ,  $c_4 \geq 0$ . The term  $E_{\varepsilon_2}^{CN}(y)$  represents a penalization term for the Ciarlet-Nečas condition which we will discuss in detail in the next section. As the final piece of the energy, we added the higher order term

$$E_{\sigma}^{reg}(y) := \sigma \int_{\Omega} |D^2 y|^s dx$$

with some  $s > 1$ . Primarily, we intend to study the limit  $\varepsilon \rightarrow 0$  here, with  $\sigma > 0$  fixed. The limit  $\sigma \rightarrow 0$  would be interesting, too, but seems out of reach at the moment.

We will always work with  $q, s$  admissible for the results of [17], which further restricts these exponents. Altogether, our assumptions on the exponents can be summarized as

$$p > d, \quad s > d, \quad q > \frac{sd}{s-d}. \quad (2.5)$$

*Remark 2.2.* One possible choice for  $W_{\varepsilon_1}$  can always be obtained by replacing  $h(x, \cdot)$  in (2.3) for each  $x$  by the Yosida-type approximation

$$h_{\varepsilon_1}(x, A) := \min \left\{ h(x, \tilde{A}) + \frac{1}{\zeta(\varepsilon_1)} |A - \tilde{A}|^p \mid \tilde{A} \in \mathbb{R}^{n(d)} \right\}$$

with  $\zeta(\varepsilon_1) > 0$  chosen small enough to obtain (2.4);  $\zeta(\varepsilon_1) \rightarrow 0$  as  $\varepsilon_1 \rightarrow 0$ . Of course, in many special cases of  $h$ , fully explicit approximations are possible, too.

*Remark 2.3.* Above, we omitted force terms in the energy, although only to keep the notation short. Similarly, boundary conditions are missing so far. These issues are discussed in greater detail in Remark 4.2 and Remark 4.3 below.

*Remark 2.4.* More general forms of  $E_\sigma^{reg}$  can be used as well if this is desired for modelling purposes. The only features we actually exploit are that with  $\sigma > 0$  and  $s$  as above,

- (i)  $E_\sigma^{reg}(y) \geq \sigma \int_\Omega |D^2 y|^s dx$ ,
- (ii)  $E_\sigma^{reg}(y)$  is uniformly continuous on bounded subsets of  $W^{2,s}$  (cf. Proposition 4.14), and
- (iii)  $y \mapsto E_\sigma^{reg}(y)$ ,  $W^{2,s}(\Omega; \mathbb{R}^d) \rightarrow \mathbb{R} \cup \{+\infty\}$  is sequentially lower semicontinuous with respect to weak convergence in  $W^{2,s}$ .

Moreover, (i) can be further weakened to  $E_\sigma^{reg}(y) \geq \sigma \int_\Omega |D^2 y|^s dx - C$  with some constant  $C$  which can easily be absorbed by other terms. For (iii), it suffices to have  $E_\sigma^{reg}(y)$  as an integral functional depending on  $D^2 y$  with a convex, polyconvex or gradient polyconvex energy density. The notion of gradient polyconvexity and related results can be found in [7].

### 3. VARIANTS OF THE CIARLET-NEČAS CONDITION AND NEW PENALIZATION TERMS

Our starting point for obtaining a new kind of penalization terms is the observation that there are many equivalent ways of stating the Ciarlet-Nečas condition (1.1). For instance, if  $y$  is regular enough such that the coarea formula holds, in particular for  $y \in W^{1,p}$  with  $p > d$  [19], (1.1) is equivalent to

$$\int_{y(\Omega)} (N_y(z) - 1) dz = 0, \quad (3.1)$$

where

$$N_y(z) := \# \{ \tilde{x} \in \Omega \mid y(\tilde{x}) = z \}$$

counts the number of times  $y$  (its continuous representative) reaches the point  $z$  in the deformed configuration. There is self-contact at  $z$  if and only if  $N_y(z) > 1$ . Yet another equivalent way of expressing (1.1) is

$$\int_{\Omega \cap \{x \mid N_y(y(x)) > 1\}} h(x) dx = 0, \quad (3.2)$$

where  $h$  can be any measurable function with  $h > 0$  a.e. in  $\{N_y \circ y > 1\}$ . Notice that the choice of such a function  $h$  does not matter, since (3.2) effectively just states that  $\{N_y \circ y > 1\} \subset \Omega$  is a set of measure zero, and by choosing  $h(x) = (N_y(y(x)))^{-1}(N_y(y(x)) - 1)$ , (3.2) reduces to (3.1) by the coarea formula.

We now introduce a new class of penalization terms  $E_{\varepsilon_2}^{CN}(y)$  that – as we will see later – lead to a condition of the form of (3.2) in the limit as  $\varepsilon_2 \rightarrow 0$ . It is defined as follows:

$$E_{\varepsilon_2}^{CN}(y) := \frac{1}{\varepsilon_2^\beta} \int_{(\Omega \times \Omega)} \frac{1}{\varepsilon_2^d} \left[ g(|\tilde{x} - x|) - g\left(\frac{1}{\varepsilon_2} |y(\tilde{x}) - y(x)|\right) \right]^+ d(x, \tilde{x}), \quad (3.3)$$

where  $[a]^+ := \max\{0, a\}$  denotes the positive part,  $\beta > 0$  is a constant and

$$g : [0, \infty) \rightarrow [0, \infty) \text{ is a continuous, strictly increasing function with } g(0) = 0. \quad (3.4)$$

The choice of  $\beta$  and  $g$  is meant to give us some freedom to optimize the behavior of numerical schemes, with prototypical examples for  $g$  being  $g(t) := t$  or  $g(t) = t^2$ .

*Remark 3.1.* The way  $E_{\varepsilon_2}^{CN}$  is defined, its integrand only contributes in  $O(\varepsilon_2)$ -neighborhoods of the self-contact (or self-penetration) set. More precisely, this “aura” never goes beyond a distance of  $\text{diam}(\Omega) \text{Lip}(y^{-1})\varepsilon_2$  away from the self-contact set. Here,  $\text{Lip}(y^{-1})$  is a Lipschitz constant of the local inverse  $y^{-1}$  (which is globally uniform, cf. Lemma 3.8 below). When computing the integral by numerical integration, this also means that a mesh size  $h$  of this order is needed, at least near self-penetration. Otherwise, huge errors are likely.

In the example illustrated in Figure 2, the radius of the aura beyond the self-contact is roughly  $\varepsilon_2$ . In that particular case,  $\text{Lip}(y^{-1})$  is still quite close to 1, and instead of  $\text{diam}(\Omega)$  in the estimate mentioned above, we may actually also use a number close to 1, namely, the distance of the two disjoint subsets of the reference configuration (undeformed domain) where the self-overlap happens (for which  $\text{diam}(\Omega)$  is of course an upper bound).

*Remark 3.2.* By default, all finite dimensional norms  $|\cdot|$  appearing in this article are assumed to be Euclidean. However, that choice does not really matter. For instance, using a different norm inside of  $g$  in (3.3) is possible. The proofs below are only affected insofar as all balls or annuli in  $\mathbb{R}^d$  or their intersections with  $\Omega$  have to be interpreted as balls (or annuli) with respect to that norm. Additional constants will then appear in Cauchy-Schwarz type inequalities, but that only changes the constants appearing in the results, not their general structure. In particular, for discretizations with finite elements defined on cubes, it can be quite convenient to use  $|x|_\infty = \max |x_i|$ ,  $x = (x_1, \dots, x_d)^T \in \mathbb{R}^d$ , instead of the Euclidean norm.

**3.1. Illustration example: the penalization term for a prescribed deformation.** All pictures of this example are displayed in

Figure 2. We assume a “pincers” domain  $\Omega \in \mathbb{R}^2$  covered in the rectangle  $(-3, 2) \times (-1.5, 1.5)$ . Using rotated polar coordinates

$$r = \sqrt{x_1^2 + x_2^2}, \quad t = \arctan(-x_2, -x_1),$$

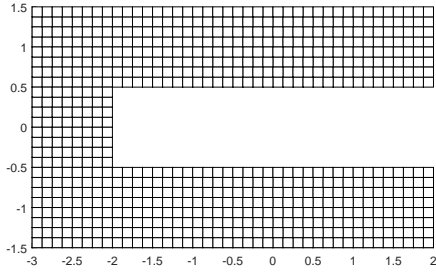
where  $x = (x_1, x_2) \in \Omega$  we define the deformation in the form

$$y(r, t) = -r(\cos(at), \sin(at))$$

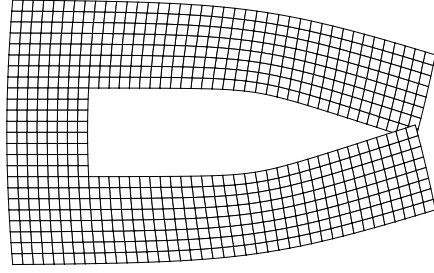
for some parameter  $a > 1$ . For a sufficiently high value of  $a$  (here we choose  $a = 1.1$ ) both pincers parts interpenetrate and the marginal density

$$d_{\varepsilon_2, y}^{CN}(x) := \frac{1}{\varepsilon_2^\beta} \int_{\Omega} \frac{1}{\varepsilon_2^d} \left[ g(|\tilde{x} - x|) - g\left(\frac{1}{\varepsilon_2} |y(\tilde{x}) - y(x)|\right) \right]^+ d\tilde{x}$$

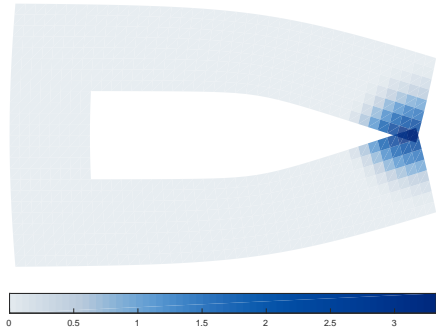
of  $E_{\varepsilon_2}^{CN}$  is evaluated and visualized for  $\varepsilon_2 = 1/2$  and  $\varepsilon_2 = 1/4$ . In both cases we consider  $\beta = 1/2$  and  $g(t) := t$ . Marginal densities are evaluated on rectangular elements by the method of finite elements. Details on implementation are provided in Section 5.



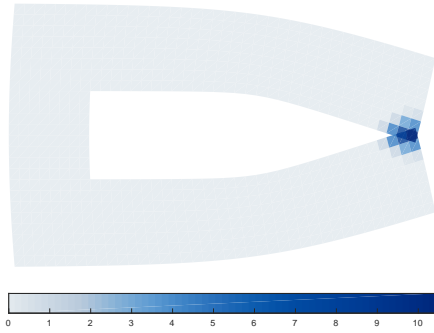
(A) Undeformed domain.



(C) Deformed domain.



(B) Density  $d_{\varepsilon_2, y}^{CN}(x)$  for  $\varepsilon_2 = 1/2$ .



(D) Density  $d_{\varepsilon_2, y}^{CN}(x)$  for  $\varepsilon_2 = 1/4$ .

FIGURE 2. Pincers domain under given deformation.

**3.2. Analytic investigation of the penalty term.** We now analyze the behavior of  $E_{\varepsilon_2}^{CN}$  as  $\varepsilon_2 \rightarrow 0$ .



**Theorem 3.3** (convergence for penalty terms of type (3.3)). *Let  $\Omega \subset \mathbb{R}^d$  be a bounded Lipschitz domain, let  $E_{\varepsilon_2}^{CN}$  be the functional defined in (3.3) with  $\beta > 0$  and  $g$  satisfying (3.4), let  $0 < \alpha \leq 1$ ,  $0 < \delta, M_1, M_2$ , and let  $R := \text{diam}(\Omega) = \sup_{x_1, x_2 \in \Omega} |x_1 - x_2|$ . Then there exist constants  $r, a, A, \bar{\varepsilon} > 0$  only depending on  $d, \Omega, \delta, \alpha, M_1, M_2$  and  $g$  such that for every  $y \in C^{1,\alpha}(\Omega; \mathbb{R}^d)$  with*

$$\det \nabla y \geq \delta > 0 \text{ and } |\nabla y| \leq M_1 \text{ on } \Omega \text{ and } \|\nabla y\|_{C^\alpha(\Omega)} \leq M_2 \quad (3.5)$$

and every  $0 < \varepsilon_2 \leq \bar{\varepsilon}$ ,

$$\varepsilon_2^\beta E_{\varepsilon_2}^{CN}(y) \geq a |P_y(r\varepsilon_2)| \quad (3.6)$$

$$\varepsilon_2^\beta E_{\varepsilon_2}^{CN}(y) \leq A |P_y(R\varepsilon_2)| \quad (3.7)$$

Here, for  $s \geq 0$  (with  $s = r\varepsilon_2$  or  $s = R\varepsilon_2$  above),

$$P_y(s) := \left\{ x \in \Omega \mid \exists \tilde{x} \in \Omega : |y(x) - y(\tilde{x})| \leq s \text{ and } |x - \tilde{x}| > \frac{\varrho}{2} \right\},$$

where  $\varrho = \varrho(\Omega, \delta, M_1, M_2, \alpha) > 0$  denotes the radius of guaranteed local injectivity from Lemma 3.6 below.

*Remark 3.4.* As a consequence of Lemma 3.6,  $P_y(0) = \{N_y \circ y > 1\}$ , i.e., it is precisely the subset in the reference configuration where global injectivity of  $y$  fails. For  $s > 0$ ,  $P_y(s)$  is the set where  $y$  almost fails to be injective up to an error of  $s$ . Moreover,  $P_y(0)$  is the limit  $P_y(s)$  as  $s \searrow 0$  (i.e., their intersection for all  $s > 0$ ;  $P_y(s_1) \subset P_y(s_2)$  if  $s_1 \leq s_2$ ). Thus, both the upper and the lower bound for  $\varepsilon_2^\beta E_{\varepsilon_2}^{CN}$  in (3.6) and (3.7), respectively, converge as  $\varepsilon_2 \rightarrow 0$  by monotone convergence:

$$a |P_y(0)| = \lim_{\varepsilon_2 \rightarrow 0} a |P_y(r\varepsilon_2)| \leq \lim_{\varepsilon_2 \rightarrow 0} A |P_y(R\varepsilon_2)| = A |P_y(0)|$$

Up to the constants, these limits coincide and are functionals of the form of (3.2) with a constant integrand.

*Remark 3.5.* The assumption (3.5) holds in sets with bounded energy, see Proposition 4.13.

For the proof of the theorem, we need the following version of the Inverse Mapping Theorem with additional control, also near  $\partial\Omega$ .

**Lemma 3.6.** *Let  $\Omega \subset \mathbb{R}^N$  be a bounded Lipschitz domain with local Lipschitz constants bounded by a fixed  $L > 0$ , and let  $y \in C^{1,\alpha}(\Omega; \mathbb{R}^d)$  such that (3.5) holds. Then there exists an  $\varrho > 0$  which only depends on  $\delta, M_1, M_2, \alpha$  and  $\Omega$  such that for every  $\bar{x} \in \bar{\Omega}$ ,  $y$  is injective on  $\Omega(\bar{x}, \varrho) := B_\varrho(\bar{x}) \cap \Omega$ . Moreover,  $y$  is bi-Lipschitz with explicitly known constants:*

$$\frac{1}{2} \frac{\delta}{M_1^{d-1}} |x_1 - x_2| \leq |y(x_1) - y(x_2)| \leq M_1 \sqrt{1 + L^2} |x_1 - x_2|, \quad (3.8)$$

for all  $x_1, x_2 \in \Omega(\bar{x}, \varrho)$ . In addition, the inverse  $y^{-1}$  of the restriction of  $y$  to  $\Omega(\bar{x}, \varrho)$  is of class  $C^{1,\alpha}$ .

**Proof.** Since  $\Omega$  is a Lipschitz domain, there exists an  $R_0 = R_0(\Omega) > 0$  such that for all  $x_0 \in \partial\Omega$ ,

$$\begin{aligned} & \text{in a cuboid containing } B_{R_0}(x_0), \partial\Omega \text{ is the graph} \\ & \text{of a Lipschitz map with constant at most } L = L(\Omega). \end{aligned} \quad (3.9)$$

For

$$r_0 := \frac{1}{1 + 2\sqrt{1 + L^2}} R_0,$$

an explicit path connecting  $x_1, x_2 \in \Omega(x_0, r_0)$  in the larger set  $\Omega(x_0, R_0)$  is given by the “V”-shaped piecewise  $C^1$ -path with slope  $L$  a.e., in  $\Omega(x_0, R_0)$  below the graph representing  $\partial\Omega$ . The length of such a path is at most  $\sqrt{1 + L^2} |x_2 - x_1|$ . In particular, if  $\varrho \leq r_0$  and  $x_1, x_2 \in \Omega(x_0, \varrho)$ , such a path never leaves the set  $\Omega(x_0, \frac{R_0}{r_0}\varrho)$ . Hence, (3.9) implies that

$$1 \leq d_{x_0, \varrho} \leq \sqrt{1 + L^2} \quad \text{for all } \varrho \leq r_0, \quad (3.10)$$

where

$$d_{x_0, \varrho} := \sup_{x_1, x_2 \in \Omega(x_0, \varrho)} \inf_p \left\{ \frac{\text{length}(p)}{|x_1 - x_2|} \left| \begin{array}{l} p : [0, 1] \rightarrow \Omega(x_0, \frac{R_0}{r_0}\varrho) \\ \text{is piecewise } C^1, \\ p(0) = x_1, p(1) = x_2 \end{array} \right. \right\}.$$

Notice that  $d_{x_0, \varrho}$  is the worst possible ratio of intrinsic path distance and Euclidean distance in  $\Omega(x_0, \frac{R_0}{r_0}\varrho)$  for pairs of points in the smaller set  $\Omega(x_0, \varrho)$ .

As a first consequence of (3.10),  $y$  is globally Lipschitz on  $\Omega(x_0, \varrho)$  with a Lipschitz constant of at most  $\|Dy\|_{L^\infty(\Omega(x_0, \varrho))} d_{x_0, \varrho} \leq M_1 \sqrt{1 + L^2}$ , which gives the second inequality in (3.8).

For the first inequality in (3.8) let  $\bar{x} \in \bar{\Omega}$ . If either  $B_{\frac{1}{2}r_0}(\bar{x}) \subset \Omega$  or  $\bar{x} \in \partial\Omega$ , we may take  $x_0 := \bar{x}$  and (3.10) holds for all  $\varrho \leq \frac{1}{2}r_0$ . In case we are given  $\bar{x} \in \bar{\Omega}$  “in between” with  $B_{\frac{1}{2}r_0}(\bar{x}) \cap \partial\Omega \neq \emptyset$ , there always exists  $x_0 = x_0(\bar{x}) \in \partial\Omega$  such that  $\Omega(\bar{x}, \frac{1}{2}r_0) \subset \Omega(x_0, r_0)$  and we again have (3.10). In addition, for any  $\bar{x} \in \bar{\Omega}$ , all  $\varrho \leq r_0$  and any pair  $x_1, x_2 \in \Omega(\bar{x}, \frac{1}{2}\varrho) \subset \Omega(x_0, \varrho)$  connected with a  $C^1$ -path  $p : [0, 1] \rightarrow \Omega(x_0, \frac{R_0}{r_0}\varrho)$ ,  $p(0) = x_1$ ,  $p(1) = x_2$ , we have that

$$\begin{aligned} |y(x_1) - y(x_2)| &= \left| \int_0^1 Dy(p(t)) \dot{p}(t) dt \right| \\ &\geq \left| \int_0^1 Dy(x_0) \dot{p}(t) dt \right| - \left| \int_0^1 |Dy(p(t)) - Dy(x_0)| |\dot{p}(t)| dt \right| \\ &\geq \left| \int_0^1 Dy(x_0) \dot{p}(t) dt \right| - M_2 \|p - x_0\|_{L^\infty(0,1)}^\alpha \int_0^1 |\dot{p}(t)| dt \\ &\geq |Dy(x_0)(x_2 - x_1)| - M_2 \left( \frac{R_0}{r_0} \varrho \right)^\alpha \text{length}(p). \end{aligned}$$

Since this is true for all such paths  $p$  and  $d_{x_0, \varrho} \leq \sqrt{1 + L^2}$  by (3.10), we infer that

$$\begin{aligned} & |y(x_1) - y(x_2)| \\ & \geq |Dy(x_0)(x_2 - x_1)| - M_2 \left( \frac{R_0}{r_0} \varrho \right)^\alpha \sqrt{1 + L^2} |x_1 - x_2| \end{aligned} \quad (3.11)$$

As  $Dy(x_0)$  is invertible with  $|Dy(x_0)^{-1}| \leq \frac{M_1^{N-1}}{\delta}$ , we also have that

$$|Dy(x_0)(x_2 - x_1)| \geq \frac{\delta}{M_1^{N-1}} |x_1 - x_2|. \quad (3.12)$$

We now choose  $\varrho$  small enough so that  $M_2 \left( \frac{R_0}{r_0} \varrho \right)^\alpha \sqrt{1 + L^2} \leq \frac{1}{2} \frac{\delta}{M_1^{N-1}}$ , and for that choice, (3.11), (3.12) yield that

$$|y(x_1) - y(x_2)| \geq \frac{1}{2} \frac{\delta}{M_1^{N-1}} |x_1 - x_2|$$

for all  $x_1, x_2 \in \Omega(\bar{x}, \frac{\varrho}{2})$ , proving the first inequality in (3.8). This in turn implies that  $y$  is locally injective. Finally, to see the asserted  $C^{1, \alpha}$ -regularity of  $y^{-1}$ , observe that  $Dy^{-1}(z) = Dy(y^{-1}(z))^{-1}$ . Therefore, due to the Lipschitz regularity of  $y^{-1}$  provided by (3.8),  $Dy^{-1} \in C^\alpha$  just like  $Dy$ .  $\square$

**Proof of Theorem 3.3.** As before, we use the following shorthand notation for  $x \in \Omega$  and  $s > 0$ :

$$\Omega(x, s) := B_s(x) \cap \Omega,$$

Next, we introduce and study a few auxiliary sets related to the definition  $P_y(s)$  that will be needed in the rest of the proof.

**(i) Auxiliary sets related to  $P_y(s)$ :  $Q_y(s, x)$  and  $X_y(s, x)$ .**

For  $s \geq 0$  and  $x \in P_y(s)$  let  $Q_y(s, x)$  denote the set of all admissible choices of  $\tilde{x}$  in the definition of  $P_y(s)$ , additionally including those that are close to  $x$ :

$$Q_y(s, x) := \{\tilde{x} \in \Omega \mid |y(x) - y(\tilde{x})| \leq s\}.$$

We claim that for  $s$  small enough,  $Q_y(s, x)$  is separated into subsets of small balls that are pairwise far apart: For every

$$s < K \varrho \quad \text{where } K := \frac{\delta}{4M_1^{d-1}},$$

we have that

$$Q_y(s, x) \subset (\Omega \setminus \Omega(x_0, \varrho)) \cup \Omega(x_0, Ks) \quad \text{for all } x_0 \in Q_y(s, x). \quad (3.13)$$

For the proof of (3.13), take any  $z \in \Omega$ ,  $z \notin (\Omega \setminus \Omega(x_0, \varrho)) \cup \Omega(x_0, Ks)$ . Then

$$\frac{4M_1^{d-1}}{\delta} s < |z - x_0| \leq \varrho,$$

and the uniform bi-Lipschitz property (3.8) entails that

$$2s < |y(z) - y(x_0)|.$$

Since  $x_0 \in Q_y(s, x)$  and thus  $|y(x) - y(x_0)| \leq s$ , we infer that  $s < |y(z) - y(x)|$ , i.e.,  $z \notin Q_y(s, x)$ .

As a consequence of (3.13), for  $s$  small enough as above, we can select a finite set  $X_y(s, x)$  such that

$$\begin{aligned} X_y(s, x) &\subset Q_y(s, x), \quad x \in X_y(s, x), \\ |x_1 - x_2| &\geq \varrho \quad \text{for all } x_1, x_2 \in X_y(s, x), x_1 \neq x_2, \end{aligned} \quad (3.14)$$

and  $Q_y(s, x)$  is contained in a disjoint union of small balls centered at points in  $X_y(s, x)$ :

$$Q_y(s, x) \subset \bigcup_{x_0 \in X_y(s, x)} \Omega(x_0, Ks). \quad (3.15)$$

Finally, notice that by the definition of  $P_y(s)$ ,  $Q_y(s, x) \setminus \Omega(x, \varrho) \neq \emptyset$ . Therefore,  $X(s, x)$  always contains at least one more point besides  $x$ :

$$\#X(s, x) \geq 2 \quad \text{for all } x \in P_y(s). \quad (3.16)$$

**(ii) Splitting  $E_{\varepsilon_2}^{CN}$ .**

For any  $s \geq 0$ , splitting  $\Omega = P_y(s) \cup (\Omega \setminus P_y(s))$  gives that

$$\varepsilon_2^\beta E_{\varepsilon_2}^{CN}(y) = \int_{P_y(s)} \frac{1}{\varepsilon_2^d} J_{y, \varepsilon_2}(x) dx + \int_{\Omega \setminus P_y(s)} \frac{1}{\varepsilon_2^d} J_{y, \varepsilon_2}(x) dx, \quad (3.17)$$

where

$$J_{y, \varepsilon_2}(x) := \int_{\Omega} \left[ g(|\tilde{x} - x|) - g\left(\frac{|y(\tilde{x}) - y(x)|}{\varepsilon_2}\right) \right]^+ d\tilde{x}.$$

Below, we estimate the two terms on the right hand side of (3.17) separately, for suitable choices of  $s$  depending on  $\varepsilon_2$ .

**(iii) Proof of (3.6).**

We use  $s := r\varepsilon_2$  in (3.17), with some constant  $0 < r \leq 1$  to be determined later. Since  $J_{y, \varepsilon_2} \geq 0$ , we get that

$$\varepsilon_2^\beta E_{\varepsilon_2}^{CN}(y) \geq \int_{P_y(r\varepsilon_2)} \frac{1}{\varepsilon_2^d} J_{y, \varepsilon_2}(x) dx. \quad (3.18)$$

The integrand of  $J_{y, \varepsilon_2}$  is also non-negative, and therefore, using (3.15), for each  $x \in P_y(r\varepsilon_2)$  we also have that

$$J_{y, \varepsilon_2}(x) \geq \sum_{x_0 \in X(r\varepsilon_2, x) \setminus \{x\}} I_{y, x, x_0}(\varepsilon_2) \quad (3.19)$$

where

$$I_{y, x, x_0}(\varepsilon_2) := \int_{\Omega(x_0, r\varepsilon_2)} \left[ g(|\tilde{x} - x|) - g\left(\frac{1}{\varepsilon_2} |y(\tilde{x}) - y(x)|\right) \right]^+ d\tilde{x}. \quad (3.20)$$

We will now proceed to estimate  $I_{y,x,x_0}(\varepsilon_2)$  for all  $x \in P_y(r\varepsilon_2)$  and  $x_0 \in X(r\varepsilon_2, x) \setminus \{x\}$ , which by (3.14) in particular implies that  $|x - x_0| \geq \varrho$ . For  $\tilde{x} \in \Omega(x_0, r\varepsilon_2)$ , the latter yields that

$$|\tilde{x} - x| \geq |x - x_0| - |x_0 - \tilde{x}| \geq \frac{\varrho}{2}$$

as long as  $2r\varepsilon_2 \leq \varrho$ . Later, it will be convenient to also have that  $r\varepsilon_2 \leq \frac{1}{2}R = \frac{1}{2} \text{diam}(\Omega)$ . As long as  $r \leq 1$ , it altogether suffices if

$$\varepsilon_2 \leq \bar{\varepsilon} := \min \left\{ \frac{\varrho}{2}, \frac{R}{2} \right\}$$

Moreover, recall that  $x_0 \in X(r\varepsilon_2, x) \subset Q(r\varepsilon_2, s)$ . By the definition of  $Q(r\varepsilon_2, s)$  in step (i), this entails that  $|y(x_0) - y(x)| \leq r\varepsilon_2$ , and consequently,

$$\begin{aligned} |y(\tilde{x}) - y(x)| &\leq |y(x_0) - y(x)| + |y(\tilde{x}) - y(x_0)| \\ &\leq r\varepsilon_2 + |y(\tilde{x}) - y(x_0)| \leq r\varepsilon_2 + M_1\sqrt{L^2 + 1}|\tilde{x} - x_0|, \end{aligned}$$

where we also used the local Lipschitz continuity of  $y$  given by (3.8). With these observation and the monotonicity of  $g$ , both expressions in  $g$  can be estimated and in this way, (3.20) implies that

$$\begin{aligned} I_{y,x,x_0}(\varepsilon_2) &\geq \int_{\Omega(x_0, r\varepsilon_2)} \left[ g\left(\frac{\varrho}{2}\right) - g\left(r + \frac{1}{\varepsilon_2}M_1\sqrt{L^2 + 1}|\tilde{x} - x_0|\right) \right]^+ d\tilde{x} \\ &\geq \int_{\Omega(x_0, r\varepsilon_2)} \left[ g\left(\frac{\varrho}{2}\right) - g\left(r + rM_1\sqrt{L^2 + 1}\right) \right]^+ d\tilde{x}. \end{aligned} \tag{3.21}$$

Here,  $r + rM_1\sqrt{L^2 + 1} \leq \frac{\varrho}{4}$  for

$$r := \min \left\{ \frac{\varrho}{4}(1 + M_1\sqrt{L^2 + 1})^{-1}, 1 \right\}.$$

Substituting  $t := r\varepsilon_2 \leq \bar{\varepsilon} \leq \frac{1}{2}R$ , we conclude that

$$I_{y,x,x_0}(\varepsilon_2) \geq a\varepsilon_2^d \tag{3.22}$$

with the constant

$$0 < a := r^d \left( g\left(\frac{\varrho}{2}\right) - g\left(\frac{\varrho}{4}\right) \right) \inf_{t,x_0} \left\{ \frac{|\Omega(x_0, t)|}{t^d} \mid \begin{array}{l} 0 < t \leq \frac{1}{2}R, \\ x_0 \in \Omega \end{array} \right\}.$$

Notice that the infimum above is a geometric constant which only depends on  $\Omega$ . It is determined by the smallest possible the volume fractions  $|\Omega \cap B_t(x_0)| / |B_t(x_0)|$ . Such fractions are bounded away from zero because  $\Omega$ , being a Lipschitz domain, satisfies an interior cone condition.

Combined with (3.18), (3.19) and (3.16), (3.22) yields (3.6).

**(iv) Proof of (3.7).**

This time, we use (3.17) with  $s = R\varepsilon_2$ ,  $R = \text{diam} \Omega$ , and distinguish the cases  $x \in P_y(R\varepsilon_2)$  and  $x \in \Omega \setminus P_y(R\varepsilon_2)$ .

**Case 1:  $x \in \Omega \setminus P_y(R\varepsilon_2)$ .** We claim that for such  $x$ ,  $J_{y,\varepsilon_2}(x) = 0$  for sufficiently small  $\varepsilon_2$ . If  $x \in \Omega \setminus P_y(R\varepsilon_2)$  then for all  $\tilde{x} \in \Omega$ ,

$$|y(\tilde{x}) - y(x)| > R\varepsilon_2 \quad \text{or} \quad |\tilde{x} - x| < \frac{\varrho}{2}.$$

In the former case, the integrand in  $J_{y,\varepsilon_2}(x)$  vanishes since  $|x - \tilde{x}| \leq \text{diam } \Omega = R$ . In the latter case, Lemma 3.6 can be applied, and due to the monotonicity of  $g$  and the lower bound in (3.8), the integrand in  $J_{y,\varepsilon_2}(x)$  vanishes again, at least if

$$\varepsilon_2 \leq \tilde{\varepsilon}, \quad \tilde{\varepsilon} := \frac{\delta}{2M_1^{d-1}}.$$

Hence,

$$J_{y,\varepsilon_2}(x) = 0 \quad \text{if } x \in \Omega \setminus P_y(R\varepsilon_2) \text{ and } \varepsilon_2 \leq \tilde{\varepsilon}. \quad (3.23)$$

**Case 2:  $x \in P_y(R\varepsilon_2)$ .** Let

$$\varepsilon_2 \leq \bar{\varepsilon} := \frac{\varrho}{R}K, \quad \text{with } K = \frac{\delta}{4M_1^{d-1}} \text{ as in (3.13)}$$

( $\bar{\varepsilon}$  here differs from its old namesake). Since  $|y(\tilde{x}) - y(x)| < R\varepsilon_2$  if and only if  $\tilde{x} \in Q_y(R\varepsilon_2, x)$ , the integrand in  $J_{y,\varepsilon_2}(x)$  vanishes for all other  $\tilde{x}$ :

$$\left[ g(|\tilde{x} - x|) - g\left(\frac{|y(\tilde{x}) - y(x)|}{\varepsilon_2}\right) \right]^+ \leq \left[ g(R) - g\left(\frac{|y(\tilde{x}) - y(x)|}{\varepsilon_2}\right) \right]^+ = 0$$

if  $|y(\tilde{x}) - y(x)| \geq R\varepsilon_2$ , since  $g$  is increasing. For  $\varepsilon_2 \leq \bar{\varepsilon}$  and  $x \in P_y(R\varepsilon_2)$ , we can therefore use (3.15) to estimate  $J_{y,\varepsilon_2}(x)$  as follows:

$$\begin{aligned} & J_{y,\varepsilon_2}(x) \\ & \leq \int_{Q_y(R\varepsilon_2, x)} \left[ g(R) - g\left(\frac{|y(\tilde{x}) - y(x)|}{\varepsilon_2}\right) \right]^+ d\tilde{x}, \\ & \leq \sum_{x_0 \in X_y(x, R\varepsilon_2)} \int_{\Omega(x_0, KR\varepsilon_2)} \left[ g(R) - g\left(\frac{|y(\tilde{x}) - y(x)|}{\varepsilon_2}\right) \right]^+ d\tilde{x}, \quad (3.24) \\ & \leq \sum_{x_0 \in X_y(x, R\varepsilon_2)} |\Omega(x_0, KR\varepsilon_2)| g(R). \\ & \leq (\sharp X_y(x, R\varepsilon_2))(KR)^d |B_1(0)| g(R). \end{aligned}$$

This is bounded by a suitable constant  $A$  because  $\sharp X_y(x, R\varepsilon_2)$ , the number of elements of  $X_y(x, R\varepsilon_2)$ , is bounded by a constant only depending on  $\varrho$  and  $R = \text{diam } \Omega$ , as a consequence of (3.14).  $\square$

Theorem 3.3 provides additional insights on the behavior of  $E_{\varepsilon_2}^{CN}$ :

**Corollary 3.7.** *In the situation of Theorem 3.3, let  $\varepsilon_2 \leq \bar{\varepsilon}$  and suppose in addition that  $y$  is more than a distance of  $R\varepsilon_2$  away from any self-contact, i.e.,*

$$|y(x_1) - y(x_2)| > R\varepsilon_2 \quad \text{for all } |x_1 - x_2| > \frac{\rho}{2}, \quad (3.25)$$

with  $R = \text{diam } \Omega$  as before. Then  $E_{\varepsilon_2}^{CN}(y) = 0$ .

**Proof.** This is a direct consequence of (3.7) and the definition of  $P_y(R\varepsilon_2)$ : (3.25) implies that  $P_y(R\varepsilon_2) = \emptyset$ .  $\square$

Another interesting consequence of Theorem 3.3 is

**Corollary 3.8** (Global invertibility for finite  $\varepsilon_2$ ). *Suppose that the assumptions of Theorem 3.3 hold and let  $C > 0$ . If  $\beta > d$  in (3.3) then there exists a constant  $0 < \tilde{\varepsilon} \leq \bar{\varepsilon}$  which only depends on  $\beta, C, d, \Omega, \delta, \alpha, M_1, M_2$  and  $g$ , such that for all  $\varepsilon_2 < \tilde{\varepsilon}$  and all  $y \in C^{1,\alpha}(\Omega; \mathbb{R}^d)$  satisfying (3.5),*

$$E_{\varepsilon_2}^{CN}(y) \leq C \quad \text{implies that } y \text{ is globally injective.} \quad (3.26)$$

**Proof.** We will prove (3.26) indirectly. Suppose that  $y$  is not globally injective, i.e.,  $y(x_1) = y(x_2)$  for a pair of points  $x_1, x_2 \in \Omega, x_1 \neq x_2$ . In view of (3.6), it suffices to show that then

$$\varepsilon_2^{-\beta} a |P_y(r\varepsilon_2)| > C \quad \text{for all } \varepsilon_2 < \tilde{\varepsilon} \quad (3.27)$$

with a suitable choice of  $\tilde{\varepsilon} > 0$ . We claim that

$$|P_y(r\varepsilon_2)| \geq c\varepsilon_2^d \quad \text{for all } \varepsilon_2 \leq \hat{\varepsilon} \quad (3.28)$$

with constants  $\hat{\varepsilon} > 0, c > 0$  yet to be determined. From (3.28), we immediately get (3.27) with  $\tilde{\varepsilon} := \min \left\{ \hat{\varepsilon}, \left( \frac{ca}{C} \right)^{\frac{1}{\beta-d}} \right\} > 0$ .

To prove (3.28), first notice that as a Lipschitz domain,  $\Omega$  satisfies an interior cone condition, i.e., there is a (cut off) cone of the form

$$V = B_\mu(0) \cap \{z \in \mathbb{R}^d \mid z \cdot e > \nu |z|\}$$

(with a fixed unit vector  $e \in \mathbb{R}^d$  and constants  $\nu < 1, \mu > 0$ ) which only depends on  $\Omega$  such that for each  $x \in \Omega, x + QV \subset \Omega$  with a suitable rotation  $Q = Q(x) \in SO(d)$ . In particular, there is  $Q_1, Q_2 \in SO(d)$  such that  $x_1 + Q_1V \subset \Omega$  and  $x_2 + Q_2V \subset \Omega$ . By Lemma 3.6, we know that  $|x_1 - x_2| \geq \rho$ , and therefore  $x_1, x_2 \in P_y(0)$ . By the local Lipschitz continuity (3.8) of  $y$  with constant  $M_1\sqrt{1+L^2}$  and the definition of  $P_y(r\varepsilon_2)$  in Theorem 3.3, we see that as a consequence, for  $j = 1, 2$ ,

$$(x_j + Q_jV) \cap B_{\lambda\varepsilon_2}(x_j) \subset P_y(r\varepsilon_2), \quad \text{with } \lambda := \frac{r}{M_1\sqrt{1+L^2}}, \quad (3.29)$$

provided that  $\varepsilon_2 \leq \bar{\varepsilon}$  and  $\lambda\varepsilon_2 \leq \rho$ . If  $\lambda\varepsilon_2 \leq \frac{\rho}{2}$ , we also know that  $B_{\lambda\varepsilon_2}(x_1)$  and  $B_{\lambda\varepsilon_2}(x_2)$  are disjoint. Since  $|(x_j + Q_jV) \cap B_{\lambda\varepsilon_2}(x_j)| = \frac{|V|}{|B_\mu(0)|} (\lambda r)^d$  as long as  $\lambda\varepsilon_2 \leq \mu$ , (3.29) entails (3.28) with  $c := 2 \frac{|V|}{|B_\mu(0)|} \lambda^d$ , for all  $\varepsilon_2 < \hat{\varepsilon} := \min \left\{ \bar{\varepsilon}, \frac{\rho}{2\lambda}, \frac{\mu}{\lambda} \right\}$ .  $\square$

*Remark 3.9.* The proof of Corollary 3.8 also works if  $x_1, x_2 \in \partial\Omega$ , and we take  $y(x_1)$  and  $y(x_2)$  as the uniquely determined values of the continuous extension of  $y \in C^{1,\alpha}$  to  $\bar{\Omega}$ . Hence, self-contact on the surface is also prevented for all  $\varepsilon_2$  small enough. In fact, one can see with similar arguments that whenever  $\beta > d$ , a universal bound on the penalty term  $E_{\varepsilon_2}^{CN}(y)$  as in (3.26) even enforces a positive minimal distance between different pieces of the body's surface (different in the sense that they are not closer than the radius  $\varrho$  of local invertibility in the reference configuration). This minimal distance converges to zero as  $\varepsilon_2 \rightarrow 0$ .

In the final piece of this section, we discuss the stability of  $E_{\varepsilon_2}^{CN}(y)$  with respect to perturbations in  $y$ . For fixed  $\varepsilon_2$ ,  $E_{\varepsilon_2}^{CN}$  is obviously continuous in  $L^\infty$ , but that continuity is not uniform in the limit  $\varepsilon_2 \rightarrow 0$ . From Theorem 3.3 and the definition of the sets  $P_y(r\varepsilon_2)$ ,  $P_y(R\varepsilon_2)$  we can infer that  $E_{\varepsilon_2}^{CN}(y)$  does not change too much if  $y$  is replaced by some perturbed deformation  $z$  with  $\|y - z\|_{L^\infty} \leq \frac{r}{3}\varepsilon_2$ , because then  $P_y(\frac{r}{3}\varepsilon_2) \subset P_z(r\varepsilon_2) \subset P_y(\frac{5}{3}r\varepsilon_2)$  (and similar inclusions also hold with  $R$  instead of  $r$ ). Here,  $z$  of course may depend on  $\varepsilon_2$ . However, it is important to be able to handle also perturbations that are small but not controlled by  $\varepsilon_2$ :

**Proposition 3.10.** *In the situation of Theorem 3.3, suppose that  $y, z \in C^{1,\alpha}(\Omega; \mathbb{R}^d)$  both satisfy (3.5). Then for every  $0 < \gamma < \frac{\varrho}{2}$  (with  $\varrho$  from Lemma 3.6), there exists a constant  $\lambda > 0$  such that*

$$\tilde{P}_y^{(\gamma)}(0) \subset P_z(0) \quad \text{if } \|y - z\|_{L^\infty} \leq \lambda, \quad (3.30)$$

where

$$\tilde{P}_y^{(\gamma)}(0) := \left\{ x_1 \in \Omega \mid \begin{array}{l} \exists x_2 \in \Omega \text{ with } \text{dist}(x_2; \partial\Omega) > \gamma, \\ y(x_1) = y(x_2) \text{ and } |x_1 - x_2| > \frac{\varrho}{2} + \gamma \end{array} \right\}.$$

Here,  $\lambda$  may depend on  $\gamma$  and the constants appearing in Theorem 3.3 but not on  $y, z$  or  $\varepsilon_2$ .

*Remark 3.11.* Since  $P_z(0) \subset P_z(r\varepsilon_2)$ , (3.30) and (3.6) in particular imply that

$$E_{\varepsilon_2}^{CN}(z) \geq a\varepsilon_2^{-\beta} |\tilde{P}_y^{(\gamma)}(0)| \quad \text{if } \|y - z\|_{L^\infty} \leq \lambda. \quad (3.31)$$

Moreover,  $\tilde{P}_y^{(\gamma)}(0)$  is always an open set (because if  $y(x_1) = y(x_2)$ , then  $y$  also self-intersects on whole neighborhoods of  $x_1$  and  $x_2$  since  $y$  is locally bi-Lipschitz due to Lemma 3.6). Therefore, whenever  $P_y(0) \neq \emptyset$  we can find  $\gamma = \gamma(y) > 0$  such that  $\tilde{P}_y^{(\gamma)}(0) \neq \emptyset$  and thus  $|\tilde{P}_y^{(\gamma)}(0)| > 0$ , and the right hand side of the inequality in (3.31) then blows up as  $\varepsilon_2 \rightarrow 0$ . Hence, only deformations  $y$  with  $P_y(0) = \emptyset$  (i.e.,  $y$  is globally invertible) can be reached in the limit along a sequence for which  $E_{\varepsilon_2}^{CN}$  remains bounded.



**Proof of Proposition 3.10.** Let  $x_1 \in \Omega$  with  $x_1 \in \tilde{P}_y^{(\gamma)}(0)$ . By definition of  $\tilde{P}_y^{(\gamma)}(0)$ , there exists  $x_2 \in \Omega$  with  $\text{dist}(x_2; \partial\Omega) > \gamma$ ,  $y(x_1) = y(x_2)$  and  $|x_1 - x_2| > \frac{\rho}{2} + \gamma$ . Both  $y$  and  $z$  are locally bi-Lipschitz due to Lemma 3.6, and the constants explicitly given in (3.8) do not depend on  $y$  or  $z$ . Hence, a whole neighborhood of  $y(x_1) = y(x_2)$  is contained in  $y(\Omega)$ . More precisely,

$$B_\tau(y(x_2)) \subset y(B_\gamma(x_2)) \subset y(\Omega) \quad \text{where } \tau := \frac{\gamma}{L_{y^{-1}}}.$$

Here,  $L_{y^{-1}} \geq 1$  can be any Lipschitz constant of the local inverse  $y^{-1}$  of  $y$  near  $x_2$ , for instance  $L_{y^{-1}} := \max\{1, \frac{2M_1^{d-1}}{\delta}\}$  is admissible by (3.8), and this particular choice is also independent of  $x_2$  and  $y$ . Analogously,

$$B_\tau(z(x_2)) \subset z(B_\gamma(x_2)).$$

Therefore, for every  $z$  with  $|y(x_i) - z(x_i)| < \gamma := \frac{1}{2}\tau$ ,  $i = 1, 2$ , we obtain that

$$z(x_1) \in B_\gamma(y(x_2)) \subset B_\tau(z(x_2)) \subset z(B_\gamma(x_2)).$$

This implies that  $x_1 \in P_z(0)$ : There exists  $\tilde{x}_2 \in B_\gamma(x_2) \subset \Omega$  such that  $z(x_1) = z(\tilde{x}_2)$  and  $|x_1 - \tilde{x}_2| \geq |x_1 - x_2| - \gamma > \frac{\rho}{2}$ .  $\square$

#### 4. CONVERGENCE OF ENERGIES

In this section, for  $y \in W^{1,p}(\Omega; \mathbb{R}^d)$ , we prove that in the limit as  $\varepsilon = (\varepsilon_1, \varepsilon_2) \rightarrow 0$ , the penalized energy

$$E_{\varepsilon, \sigma}(y) = \begin{cases} E_{\varepsilon_1}^{el}(y) + E_\sigma^{reg}(y) + E_{\varepsilon_2}^{CN}(y) & \text{if } y \in W^{2,s}, \\ +\infty & \text{else,} \end{cases}$$

with

$$E_{\varepsilon_1}^{el}(y) := \int_{\Omega} W_{\varepsilon_1}(x, \nabla y) dx,$$

converges to the original energy

$$E_\sigma(y) = \begin{cases} E^{el}(y) + E_\sigma^{reg}(y) & \text{if } y \in W^{2,s} \text{ and (1.1) holds,} \\ +\infty & \text{else,} \end{cases}$$

which includes the Ciarlet-Nečas condition (1.1) as a built-in constraint. Here, recall that

$$E_{\varepsilon_1}^{el}(y) = \int_{\Omega} W_{\varepsilon_1}(x, \nabla y) dx, \quad E^{el}(y) = \int_{\Omega} W(x, \nabla y) dx.$$

In addition, we also consider the convergence of discrete Galerkin approximations. For that, let  $h > 0$  (typically a mesh size) and let  $Y_h$  denote associated finite dimensional subspaces of  $(W^{2,s} \cap W^{1,p})(\Omega; \mathbb{R}^d)$

(typically  $Y_h \subset W^{2,\infty}$ ) such that the maximal approximation error  $\mathcal{E}(h)$  satisfies

$$\mathcal{E}(h) := \sup_{y \in W^{2,s}} \inf_{y_h \in Y_h} (\|y - y_h\|_{W^{2,s} \cap W^{1,p}}) \xrightarrow{h \rightarrow 0} 0. \quad (4.1)$$

The corresponding finite dimensional approximations of  $E_{\varepsilon,\sigma}$  are

$$E_{\varepsilon,\sigma}^h(y) := \begin{cases} E_{\varepsilon_1}^{el}(y) + E_{\sigma}^{reg}(y) + E_{\varepsilon_2}^{CN}(y) & \text{if } y \in Y_h, \\ +\infty & \text{else,} \end{cases}$$

As defined,  $E_{\varepsilon,\sigma}^h$  is assumed to be exact on  $Y_h$ . In this context, we will not discuss the question of how to calculate the integrals in  $E_{\varepsilon,\sigma}^h$  in practice. The easiest possible approach is of course based on additional approximations using standard methods in numerical integration. For our analysis, additional error terms that might appear at this stage do not matter as long as they still converge to zero as  $(h, \varepsilon) \rightarrow 0$ . However, it is useful to optimize the evaluation of the double integral in  $E_{\varepsilon_2}^{CN}$  for performance reasons, since only small neighborhoods of the self-contact set (or any almost self-contact) actually contribute.

*Remark 4.1.* Artificially assigning the value  $+\infty$  in the definitions of the functionals is just a way of encoding a restricted class of admissible functions. Be warned that there are still other “inadmissible” deformations with infinite energy in case of  $E_{\sigma}$ , namely any  $y \in W^{2,s} \cap W^{1,p}$  for which  $\int_{\Omega} W(x, \nabla y) dx = +\infty$  because  $\det \nabla y$  is too close to zero or even non-positive on a non-negligible set.

*Remark 4.2* (additional force terms). As already briefly mentioned, we did not add any terms corresponding to exterior forces, but only to keep the notation short. Since we actually prove  $\Gamma$ -convergence, our results are stable with respect to the addition of any term that is continuous with respect to the topology used for the states in the  $\Gamma$ -limit (see [12], e.g.). For us, that is the weak topology of  $W^{2,s}$  (or the weak topology of  $W^{1,p}$ , which is a weaker topology but still leads to the same result for fixed  $\sigma > 0$ ). Continuous perturbations in the weak topology of  $W^{1,p}$  in particular include linear body force terms like

$$\int_{\Omega} y \cdot g_{\text{body}} dx, \quad \text{with a } g_{\text{body}} \in L^1(\Omega; \mathbb{R}^d). \quad (4.2)$$

Similarly, one can add linear boundary force terms like

$$\int_{\partial\Omega} y \cdot g_{\text{surface}} d\mathcal{H}^{d-1}(x), \quad \text{with a } g_{\text{surface}} \in L^1(\partial\Omega; \mathbb{R}^d), \quad (4.3)$$

where the space  $L^1$  on  $\partial\Omega$  is understood with respect to the Hausdorff measure (surface measure)  $\mathcal{H}^{d-1}$ . Moreover, since  $p > d$  and  $\Omega \subset \mathbb{R}^d$  is Lipschitz,  $W^{1,p}(\Omega)$  is compactly embedded into  $C(\bar{\Omega})$ . Due to this

compact embedding, any nonlinear force terms that are continuous on  $C(\bar{\Omega}; \mathbb{R}^d)$  or  $C(\partial\Omega; \mathbb{R}^d)$  are allowed as well, like

$$\int_{\Omega} G_{\text{body}}(x, y) dx, \quad \text{with a } G_{\text{body}} \in C(\bar{\Omega} \times \mathbb{R}^d), \text{ or}$$

$$\int_{\partial\Omega} G_{\text{surface}}(x, y) d\mathcal{H}^{d-1}(x), \quad \text{with a } G_{\text{surface}} \in C(\partial\Omega \times \mathbb{R}^d).$$

Finally, we could exploit the added regularity in form of terms that are weakly continuous in  $W^{2,s}$ , which allows even bulk and boundary terms involving  $\nabla y$ .

*Remark 4.3* (boundary conditions). We also did not add any explicit boundary condition so far. Still, a weak form of a natural Neumann type boundary condition with the outer normal  $\nu$  to  $\partial\Omega$  is built in on all pieces of the boundary  $\Lambda_N \subset \partial\Omega$  where  $y$  is not subject to explicit other boundary conditions (if any), e.g.:

$$\frac{\sigma}{s} |D^2 y|^{s-2} D^2 y : (\nu \otimes \nu) + D_F W(x, \nabla y) \cdot \nu + g_{\text{surface}} = 0 \quad \text{on } \Lambda_N.$$

Here, we assumed that exactly one surface term was added to the energy, namely (4.3). Dirichlet conditions on  $\Lambda_D$ , the rest of the boundary, could be added. The limit of  $E_{\varepsilon_2}^{CN}$  is not directly affected by that since the results of Section 3 obviously also hold for any restricted class of states. Still, extra efforts in the proof of Theorem 4.6 (ii) below would be required to make sure that the Dirichlet condition is always respected when we manipulate states. The extra requirements for the boundary data that would be needed then are the following: If we impose

$$y = y_0 \text{ on } \Lambda_D, \text{ with } \Lambda_D \subset \partial\Omega \text{ relatively open,}$$

the given boundary data  $y_0 : \Lambda_D \rightarrow \mathbb{R}^d$  must have an extension to a state  $y_0 \in W^{2,s}(\Omega; \mathbb{R}^d)$  which is far enough from any self-penetration so that

- (i)  $E_{\sigma}(y_0) < +\infty$ ;
- (ii)  $E_{\varepsilon_2}^{CN}(y_0) \rightarrow 0$  as  $\varepsilon_2 \rightarrow 0$ .

In particular,  $y_0$  must satisfy the Ciarlet-Nečas condition (1.1), and if  $\beta > d$  (recall that  $\beta$  is the parameter formally governing the blow-up rate of  $E_{\varepsilon_2}^{CN}$ ),  $y_0$  must not have self-contact on the boundary, cf. Corollary 3.8.

*Remark 4.4.* In their basic form without additional terms,  $E_{\varepsilon, \sigma}$  and  $E_{\sigma}$  are translation invariant, i.e., constant vectors can be added to  $y$  without changing the energy. In particular,  $E_{\varepsilon, \sigma}$  and  $E_{\sigma}$  are only coercive when these constants are removed. This can be easily achieved by working in the quotient space  $W^{2,s}(\Omega; \mathbb{R}^d)/\mathbb{R}^d$  or  $W^{1,p}(\Omega; \mathbb{R}^d)/\mathbb{R}^d$ . Alternatively, if translation invariance is broken by boundary conditions or additional terms in the energy, it suffices if these somehow fix the

constant (e.g., by a Dirichlet condition) or control it (e.g., by a coercive nonlinear force term).

For fixed  $\varepsilon$  and  $h$ , we always have the existence of an energy minimizer:

**Proposition 4.5** (existence of minimizers). *Let  $\sigma > 0$  be fixed and suppose that (2.1)–(2.3), (2.4) and (2.5) hold. Then for every fixed  $h, \varepsilon_1, \varepsilon_2 > 0$ ,  $E_{\varepsilon, \sigma}$  and  $E_{\varepsilon, \sigma}^h$  attain their minima in  $W^{2,s}$  and  $Y_h \subset W^{2,s}$ , respectively.*

*Proof.* All three summands of  $E_{\varepsilon, \sigma}$  (or  $E_{\varepsilon, \sigma}^h$ ) are weakly lower semicontinuous in  $W^{2,s}$ :  $E^{reg}$  is convex and thus weakly lower semicontinuous. The other two terms are even weakly continuous, because  $W^{2,s}$  (or its closed subspace  $Y_h$ ) is compactly embedded in  $W^{1,\infty}$  and  $L^\infty$ ,  $y \mapsto \int_\Omega W_{\varepsilon_1}(x, \nabla y) dx$  is continuous in  $W^{1,\infty}$  and  $y \mapsto E_{\varepsilon_2}^{CN}(y)$  is continuous in  $L^\infty$ . Due to the definition of  $E^{reg}$  and the lower bounds for  $W$  and  $W_\varepsilon$ ,  $E_{\varepsilon, \sigma}$  and  $E_{\varepsilon, \sigma}^h$  are also coercive with respect to the seminorm  $\|y\| := \|D^2 y\|_{L^s} + \|Dy\|_{L^p}$  on  $W^{2,s}$ , which by Poincaré's inequality is a norm on the quotient space  $W^{2,s}/\mathbb{R}^d$  where functions differing only up to an additive constant vector are considered equivalent. (The quotient space is only needed when the translation invariance of the energy is not broken by boundary conditions or force terms.) Hence, we get the existence of minimizers by the direct method of the calculus of variations.  $\square$

Our main results provides convergence of  $E_{\varepsilon, \sigma}^h$  and its minimum as  $(h, \varepsilon) \rightarrow 0$ :

**Theorem 4.6.** *Let  $\sigma > 0$  be fixed and suppose that (2.1)–(2.3), (2.4) and (2.5) hold. Then for every  $(h(k), \varepsilon(k)) = (h(k), \varepsilon_1(k), \varepsilon_2(k)) \in (0, \infty)^3$ ,  $k \in \mathbb{N}$ , with  $h(k) \rightarrow 0$ ,  $\varepsilon_1(k) \rightarrow 0$  and  $\varepsilon_2(k) \rightarrow 0$  as  $k \rightarrow \infty$ , we have the following two properties for all  $y \in W^{2,s}(\Omega; \mathbb{R}^d)$ :*

(i) *For every sequence  $y_k \rightharpoonup y$  in  $W^{2,s}$  (weakly),*

$$\liminf_{k \rightarrow \infty} E_{\varepsilon(k), \sigma}^{h(k)}(y_k) \geq E_\sigma(y);$$

(ii) *there exists a sequence  $y_k \rightarrow y$  in  $W^{2,s}$  (strongly) such that*

$$\lim_{k \rightarrow \infty} E_{\varepsilon(k), \sigma}^{h(k)}(y_k) = E_\sigma(y).$$

*This also remains true for the case  $h = 0$  if we define  $E_{\varepsilon, \sigma}^0 := E_{\varepsilon, \sigma}$ .*

*Remark 4.7* ( $\Gamma$ -convergence). If (ii) is slightly weakened to

(ii)' *there exists a sequence  $y_k \rightharpoonup y$  in  $W^{2,s}$  (weakly) such that*

$$\lim_{k \rightarrow \infty} E_{\varepsilon(k), \sigma}^{h(k)}(y_k) = E_\sigma(y),$$

then (i) and (ii)' are exactly the definition of  $\Gamma(W^{2,s}$ -weak)-convergence of  $E_{\varepsilon, \sigma}^h$  to  $E_\sigma$  as  $(h, \varepsilon) \rightarrow 0$ , i.e.,  $\Gamma$ -convergence with respect to the weak topology in  $W^{2,s}$ .

*Remark 4.8* (convergence of minimizers). For  $\sigma > 0$  fixed, the family of functionals  $(E_{\varepsilon,\sigma}^h)_{(\varepsilon,h)}$  is equi-coercive in  $W^{2,s}/\mathbb{R}^d$  due to (2.4) and the obvious properties of  $E_\sigma^{reg}$ . As a consequence,  $\Gamma$ -convergence automatically implies that up to a subsequence (and subtracting suitable constants if necessary), minimizers of  $E_{\varepsilon,\sigma}^h$  weakly converge in  $W^{2,s}$  to a minimizer of the limit functional  $E_\sigma$ . Moreover,  $\Gamma(W^{2,s}$ -weak)-convergence of  $E_{\varepsilon,\sigma}^h$  to  $E_\sigma$  is equivalent to  $\Gamma(W^{1,p}$ -weak)-convergence.

*Remark 4.9.* Unlike the corresponding result in [20], we do not need to prescribe any “stability criterion” linking  $\varepsilon_1$  and  $h$ . This is essentially a consequence of Lemma 4.11 below that was slightly improved compared to its predecessor in [17], see also Remark 4.12. However, if  $E_{\varepsilon_2}^{CN}$  or other integrals in the energy are not computed exactly as assumed within this section, but only approximated by numerical integration (even in  $Y_h$ ), we typically need  $h$  of the order of  $\varepsilon_2$  or smaller to keep the approximation error on a tolerable level, cf. Remark 3.1.

*Remark 4.10.* In [23] (also see [22] for domains with Lipschitz boundaries), equilibrium equations for hyperelastic minimizers were derived. The Ciarlet-Nečas condition there gives rise to a boundary force term with a force density given in form of a Radon measure concentrated on the self-contact set on the boundary (if any). In a sense, our penalty term in the limit as  $\varepsilon_2 \rightarrow 0$  should represent an associated energy. However, a direct connection on the technical level is not quite obvious.

For the proof of Theorem 4.6, we strongly rely on a result of [17] that yields a uniform positive lower bound  $\det \nabla y \geq \delta > 0$  for all deformations with bounded energy, provided that the energy contains a higher order term that controls the norm of  $J(x) := \det \nabla y(x)$  in a Hölder space. This also uses that as a Lipschitz domain,  $\Omega$  has an interior cone property: For each  $x \in \Omega$ , there exists a rotation  $Q_x \in SO(d)$  such that  $x + Q_x V \subset \Omega$ , where

$$V := B_\mu(0) \cap \{z = (z_1, \dots, z_d) \in \mathbb{R}^d \mid z_1 > \nu |z|\} \subset \mathbb{R}^d,$$

is a fixed (cut-off) cone given by suitable constants  $\mu > 0$ ,  $\nu < 1$  independent of  $x$ . For such domains, we have the following variant of [17, Lemma 4.1]. Here, we also use slightly weaker assumptions and state additional explicit information about the constant  $\delta$ , but essentially, it is still based on the same ideas.

**Lemma 4.11.** *Suppose that  $\Omega \subset \mathbb{R}^d$  is bounded domain with an interior cone property as introduced above, and let  $J \in C^\alpha(\Omega)$ ,  $\alpha \in (0, 1)$ . In addition, suppose that*

$$\int_\Omega \max\{\delta, J(x)\}^{-q} dx \leq C \text{ and } \sup_{\substack{x_1, x_2 \in \Omega \\ |x_1 - x_2| < \mu}} \frac{|J(x_1) - J(x_2)|}{|x_1 - x_2|^\alpha} \leq M \quad (4.4)$$

where  $q \geq d/\alpha$ ,  $M > 0$  are constants and

$$\delta := \kappa^{-1}(C), \quad \kappa(t) := d \frac{|V|}{\mu^d} \int_0^\mu (t + M |r|^\alpha)^{-q} r^{d-1} dr, \quad t > 0.$$

Then  $J(x) > \delta > 0$  for all  $x \in \Omega$ .

*Remark 4.12.* Since  $\delta$  depends on  $C$ , the first condition in (4.4) looks somewhat implicit. Typically, we a priori have it with some other constant instead of  $\delta$ , say,  $\gamma \geq 0$ . One can then compute  $\delta$  (which does not depend on  $\gamma$ ) and check a posteriori whether  $\gamma \leq \delta$ . If this is the case, we automatically get (4.4) with  $\delta$ , too, because then trivially  $\int_\Omega \max\{\delta, J(x)\}^{-q} dx \leq \int_\Omega \max\{\gamma, J(x)\}^{-q} dx \leq C$ . We state (4.4) in this slightly complicated form to point out that the singular function  $J \mapsto J^{-q}$  appearing there can be modified near the origin, removing the singularity, as long as one does not change the value for  $J \geq \delta$ . As we show in detail in Proposition 4.13 below, this is quite useful because it means we can verify (4.4) for  $J = \det \nabla y$  also using energy bounds for our approximate elastic energy functionals involving the modified energy densities  $W_{\varepsilon_1}$ , at least if  $\varepsilon_1$  is small enough.

**Proof of Lemma 4.11.** First notice that  $\kappa$  is a strictly decreasing function with  $\kappa(t) \rightarrow 0$  as  $t \rightarrow +\infty$ , and  $\kappa(t) \rightarrow +\infty$  as  $t \searrow 0$  since  $q \geq d/\alpha$ . Hence,  $\kappa^{-1} : (0, \infty) \rightarrow (0, \infty)$  is well defined and also strictly decreasing. Moreover, while  $\kappa$  was defined using polar coordinates, the integral also can be written in standard coordinates:

$$\kappa(t) = \int_V (t + M |x|^\alpha)^{-q} dx = \int_{QV} (t + M |x|^\alpha)^{-q} dx, \quad (4.5)$$

for all  $Q \in SO(d)$ . Now let  $x_0 \in \Omega$  and choose  $Q = Q(x_0) \in SO(d)$  such that  $x_0 + QV \subset \Omega$ . Depending on  $x_0$ , we define  $K \in C^\alpha(\Omega)$ ,

$$K(x) := \begin{cases} J(x) & \text{if } J(x_0) \geq \delta, \\ J(x) + \delta - J(x_0) & \text{if } J(x_0) < \delta. \end{cases}$$

Since  $K \geq J$  and their difference is a constant, (4.4) implies that

$$\int_\Omega \max\{\delta, K(x)\}^{-q} dx \leq C \quad \text{and} \quad \sup_{\substack{x_1, x_2 \in \Omega \\ |x_1 - x_2| < \mu}} \frac{|K(x_1) - K(x_2)|}{|x_1 - x_2|^\alpha} \leq M \quad (4.6)$$

In addition,  $K(x_0) \geq \delta$ , and from (4.6) and (4.5) we thus get that

$$\begin{aligned} \kappa(K(x_0)) &\leq \int_{\Omega \cap B_\mu(x_0)} (K(x_0) + M |x - x_0|^\alpha)^{-q} dx \\ &= \int_{\Omega \cap B_\mu(x_0)} \max\{\delta, K(x_0) + M |x - x_0|^\alpha\}^{-q} dx \\ &< \int_\Omega \max\{\delta, K(x)\}^{-q} dx \leq C. \end{aligned}$$

Hence,  $K(x_0) > \kappa^{-1}(C) = \delta$ , and therefore  $J(x_0) = K(x_0) > \delta$ .  $\square$

We can now derive additional properties for deformations with bounded energy.

**Proposition 4.13.** *Let  $\Omega \subset \mathbb{R}^d$  be a bounded Lipschitz domain, let  $\sigma > 0$ , suppose that (2.1), (2.2), (2.4) and (2.5) hold, and let  $K > 0$ . Then there is a constant  $\bar{\varepsilon}_1 > 0$  such that for every  $0 < \varepsilon_1 \leq \bar{\varepsilon}_1$  and every*

$$y \in \mathcal{B}_K(\varepsilon_1) := \{y \in W^{2,s}(\Omega; \mathbb{R}^d) \mid E_\sigma^{reg}(y) + E_{\varepsilon_1}^{el}(y) \leq K\},$$

(3.5) holds with  $\alpha := \frac{s-d}{s}$  and suitable constants  $\delta > 0$ ,  $M_1, M_2 \geq 0$  independent of  $y$  and  $\varepsilon_1$ . That is,

$$\det \nabla y \geq \delta > 0 \text{ and } |\nabla y| \leq M_1 \text{ on } \Omega, \text{ and } \|\nabla y\|_{C^\alpha(\Omega)} \leq M_2. \quad (4.7)$$

This also holds for  $\varepsilon_1 = 0$  if we replace  $E_{\varepsilon_1}^{el}$  by  $E^{el}$ .

**Proof.** We only discuss the case  $\varepsilon_1 > 0$ ; the case  $\varepsilon_1 = 0$  is a similar and more straightforward application of Lemma 4.11. Since  $E_\sigma^{reg}(y) = \sigma \int_\Omega |D^2 y|^s dx$  and  $W_{\varepsilon_1}$  satisfies the lower bound stated in (2.4),  $y \in \mathcal{B}_K(\varepsilon_1)$  implies that  $\|D^2 y\|_{L^s}^s + c_3 \|\nabla y\|_{L^p}^p \leq K + c_4$ . In particular,  $\nabla y$  is bounded in  $W^{1,s}(\Omega; \mathbb{R}^{d \times d})$ , which is continuously embedded in  $C^\alpha(\Omega; \mathbb{R}^{d \times d})$  and  $C(\bar{\Omega}; \mathbb{R}^{d \times d})$ . Hence, we immediately get the last two inequalities in (4.7). It remains to show the lower bound for  $\det \nabla y$ . This will be obtained by applying Lemma 4.11, and we therefore have to verify (4.4) for  $J := \det \nabla y$ . Since  $\|\nabla y\|_{L^\infty} \leq M_1$  and  $\|\nabla y\|_{C^\alpha} \leq M_2$ , the Hölder semi-norm of  $\det \nabla y$  (a polynomial of degree  $d$  in  $\nabla y$ ) appearing in (4.4) is bounded by

$$M := dM_1^{d-1}M_2.$$

For a proof of the first inequality in (4.4), let  $\gamma > 0$ . For  $y \in \mathcal{B}_K(\varepsilon_1)$ , we have that  $\int_\Omega W_{\varepsilon_1}(x, \nabla y) dx \leq K$ , and as a consequence of this and (2.4), we obtain the following estimate for all  $\varepsilon_1 \leq \gamma$ :

$$\int_\Omega \max\{\gamma, \det \nabla y(x)\}^{-q} dx \leq \frac{c_4}{c_3} |\Omega| + \frac{1}{c_3} \int_\Omega W_{\varepsilon_1}(x, \det \nabla y) dx \leq C$$

with  $C := \frac{c_4}{c_3} |\Omega| + \frac{1}{c_3} K$ .

Notice that  $C$  does not depend on  $\gamma$  or  $\varepsilon_1$ . Hence, we may use  $\bar{\varepsilon}_1 := \gamma := \delta$ , with the constant  $\delta$  of Lemma 4.11 (with  $C$  and  $M$  as defined above). The lemma then entails that  $J = \det \nabla y \geq \delta > 0$ , i.e., the first inequality in (4.7).  $\square$

The final missing ingredients for the proof of Theorem 4.6 are some uniform continuity properties of the terms in the energy.

**Proposition 4.14.** *Let  $\Omega \subset \mathbb{R}^d$  be a bounded domain,  $\sigma > 0$  and  $s \geq 1$ . Then  $y \mapsto E_\sigma^{reg}(y) = \sigma \int_\Omega |D^2 y(x)|^s dx$  is uniformly (Lipschitz) continuous on all bounded subset of  $W^{2,s}(\Omega; \mathbb{R}^d)$ .*

**Proof.** This is a simple consequence of Hölder's inequality and the elementary  $s$ -Lipschitz continuity of  $G \mapsto |G|^s$ ,  $\mathbb{R}^{d \times d \times d} \rightarrow \mathbb{R}$ :

$$||G_1|^s - |G_2|^s| \leq s(|G_1|^{s-1} + |G_2|^{s-1})|G_1 - G_2|. \quad \square$$

**Proposition 4.15.** *Let  $\Omega \subset \mathbb{R}^d$  be a bounded Lipschitz domain, suppose that (2.1), (2.2), (2.4) and (2.5) hold, and let  $K > 0$ . Then there exists a constant  $\bar{\varepsilon}_1 > 0$  and a modulus of continuity  $\theta$  (i.e.,  $\theta : [0, \infty) \rightarrow [0, \infty)$  continuous and increasing with  $\theta(0) = 0$ ) such that for every  $0 < \varepsilon_1 \leq \bar{\varepsilon}_1$ ,*

$$\begin{aligned} |E_{\varepsilon_1}^{el}(y_1) - E^{el}(y_2) dx| &\leq |\Omega| \varepsilon_1 + \theta(\|y_1 - y_2\|_{W^{1,\infty}}) \\ &\text{for all } y_1, y_2 \in \mathcal{B}_K(\varepsilon_1), \end{aligned} \quad (4.8)$$

where  $\mathcal{B}_K(\varepsilon_1)$  is the set defined in Proposition 4.13. We emphasize that  $\theta$  is independent of  $y_1, y_2$  and  $\varepsilon_1$ .

**Proof.** By (2.4),

$$|W_{\varepsilon_1}(x, F) - W(x, F)| \leq \varepsilon_1 \quad \text{if } |F| \leq \frac{1}{\varepsilon_1} \text{ and } \det F \geq \varepsilon_1.$$

In view of Proposition 4.13, it therefore suffices to show that for a suitable modulus of continuity  $\theta$ ,

$$\begin{aligned} \left| \int_{\Omega} W(x, \nabla y_1(x)) dx - \int_{\Omega} W(x, \nabla y_2(x)) dx \right| \\ \leq \theta(\|\nabla y_1 - \nabla y_2\|_{W^{1,\infty}}) \quad \text{for all } y_1, y_2 \in \tilde{\mathcal{B}}_K, \end{aligned} \quad (4.9)$$

i.e., the uniform continuity of  $y \mapsto \int_{\Omega} W(x, \nabla y) dx$  with respect to the topology of  $W^{1,\infty}$  on the set

$$\tilde{\mathcal{B}}_K := \{y \in W^{2,s}(\Omega; \mathbb{R}^d) \mid |\nabla y| \leq M_1 \text{ and } \det \nabla y \geq \delta \text{ in } \Omega\}.$$

If  $W$  does not depend on  $x$ , (4.9) is obvious as for each  $x$ ,  $W(x, \cdot) : \mathbb{R}^{d \times d} \rightarrow \mathbb{R} \cup \{+\infty\}$  is continuous and therefore uniformly continuous on any compact set where it is finite. For a general Carathéodory function  $W$ , (4.9) still follows from similar reasoning, as a consequence of the Scorza-Dragnoni theorem (continuity of  $W$  on compact sets with complements of arbitrarily small measure in  $\Omega$ , see [11], e.g.).  $\square$

**Proof of Theorem 4.6.** We will only provide a proof for the case involving Galerkin approximations with  $h(n) > 0$ ,  $h(n) \rightarrow 0$ . The case  $h = 0$  is similar and slightly simpler.

(i) **“Lower bound”:** Let  $y_n \rightharpoonup y$  as  $n \rightarrow \infty$ , weakly in  $W^{2,s}$ . By compact embedding, this implies that  $y_n \rightarrow y$  strongly in  $W^{1,\infty}$ . Passing to a suitable subsequence (not relabeled), we may assume that  $e_0 := \liminf E_{\sigma, \varepsilon(n)}^{h(n)}(y_n) = \lim E_{\sigma, \varepsilon(n)}^{h(n)}(y_n)$ . In addition, we may assume that  $e_0 < +\infty$  because otherwise there is nothing to show. With  $K := e_0 + 1$ , we have  $E_{\sigma, \varepsilon(n)}^{h(n)}(y_n) \leq K$  for all  $n$  sufficiently large. Since  $E_{\varepsilon_2(n)}^{CN} \geq 0$ , we infer that  $y_n \in \mathcal{B}_K(\varepsilon_1(n))$  for all such  $n$ , where  $\mathcal{B}_K(\varepsilon_1(n))$



is the set introduced in Proposition 4.13. Due to Proposition 4.15, we therefore have that

$$\lim_{n \rightarrow \infty} \int_{\Omega} W_{\varepsilon_1(n)}(x, \nabla y_n) dx = \int_{\Omega} W(x, \nabla y) dx. \quad (4.10)$$

Moreover, by weak lower semicontinuity of the convex functional  $E_{\sigma}^{reg}$ ,

$$\liminf_{n \rightarrow \infty} E_{\sigma}^{reg}(y_n) \geq E_{\sigma}^{reg}(y). \quad (4.11)$$

Besides (4.10) and (4.11), it also trivially holds that  $E_{\varepsilon_2(n)}^{CN} \geq 0$ . We thus get  $\liminf E_{\sigma, \varepsilon(n)}^{h(n)}(y_n) \geq E_{\sigma}(y)$  as asserted, provided that  $y$  satisfies the Ciarlet-Nečas condition. For the proof of the latter, first observe that due to Proposition 4.13, (3.5) is satisfied for each  $y_n$  (in place of  $y$ ), and Theorem 3.3 is applicable, at least as long as  $n$  is also large enough so that  $\varepsilon_2(n) \leq \bar{\varepsilon}$ . We also know that  $E_{\varepsilon_2(n)}^{CN}(y_n)$  is bounded from above because  $e_0 < +\infty$  and the other terms in  $E_{\sigma, \varepsilon(n)}^{h(n)}(y_n)$  are bounded from below. Since  $y_n \rightarrow y$  in  $L^{\infty}$ , Proposition 3.10 and Remark 3.11 therefore yield that  $\{N_{y_n} \circ y_n\} = P_y(0) = \emptyset$ . In particular,  $y$  is globally invertible and the Ciarlet-Nečas condition holds for  $y$ .

**(ii) Existence of a (strongly converging) “recovery sequence”:** Let  $y \in W^{2,s}(\Omega; \mathbb{R}^d)$ . If  $E_{\sigma}(y) = +\infty$ , any sequence  $(y_n) \subset Y_{h(n)}$  with  $y_n \rightarrow y$  in  $W^{2,s}$  is suitable, in particular  $\lim E_{\sigma, \varepsilon(n)}^{h(n)}(y_n) = +\infty = E_{\sigma}(y)$  as a consequence of (i). We thus may assume that  $E_{\sigma}(y) < +\infty$ . In particular, the Ciarlet-Nečas condition holds for  $y$  and Proposition 4.13 (for the case  $h = 0$ ) is applicable with  $K := E_{\sigma}(y)$ , which yields (4.7). Due to (4.7), Lemma 3.6 can be applied. Hence,  $y$  is also locally bi-Lipschitz, and for such maps, the Ciarlet-Nečas condition is equivalent to global invertibility in the classical sense.

Nevertheless,  $y$  may still exhibit self-contact on the boundary. This is problematic for our construction because we might lose control of  $E_{\varepsilon_2(n)}^{CN}$ . We therefore first artificially create a little gap around the boundary, with the ultimate goal of finding a suitable sequence  $y_n \in Y_{h(n)}$  approximating  $y$  and its energy while  $E_{\varepsilon_2}^{CN}(y_n) = 0$  for all  $n$  (large enough). To create this gap, let  $\Psi_j : \Omega \rightarrow \Omega$  be a sequence of globally invertible maps of class  $C^{\infty}$  which “shrink”  $\Omega$  into a slightly smaller set and converge to the identity, more precisely,

$$\Psi_j(\Omega) \subset \Omega^{(j)} := \left\{ x \in \Omega \mid \text{dist}(x; \partial\Omega) > \frac{1}{j} \right\}, \quad (4.12)$$

$$\|\Psi_j - id\|_{W^{2,\infty}} \xrightarrow{j \rightarrow \infty} 0.$$

Such maps  $\Psi_j$  are easy to define locally in a neighborhood of a boundary point where  $\partial\Omega$  can be represented as the graph of a Lipschitz function. Globally, the local pieces can be glued together using a decomposition

of unity; we omit the details. As a consequence of (4.12),

$$z_j \xrightarrow{j \rightarrow \infty} y \text{ in } W^{2,s} \text{ (and } W^{1,\infty}) \text{ for } z_j := y \circ \Psi_j. \quad (4.13)$$

The function  $z_j$ , like other possible perturbations  $z$  of  $y$  in  $W^{2,s}$ , inherits (4.7) with slightly modified constants: There is a constant  $\kappa > 0$  such that for all  $z \in W^{2,s} \subset C^{1,\alpha} \subset W^{1,\infty}$  with  $\|z - y\|_{W^{2,s}} \leq \kappa$ ,

$$\det \nabla z \geq \tilde{\delta} > 0 \text{ and } |\nabla z| \leq \tilde{M}_1 \text{ on } \Omega, \text{ and } \|\nabla z\|_{C^\alpha(\Omega)} \leq \tilde{M}_2, \quad (4.14)$$

where  $\tilde{\delta} := \frac{1}{2}\delta$ ,  $\tilde{M}_1 := M_1 + 1$  and  $\tilde{M}_2 := M_2 + 1$ . In particular,  $z = z_j$  is admissible if  $j$  is sufficiently large. Due to (4.13), (4.14) and Proposition 4.15,

$$\begin{aligned} |E_{\varepsilon_1(n)}^{el}(z) - E^{el}(y)| &\leq |\Omega| \varepsilon_1(n) + \theta(c\|z - y\|_{W^{2,s}}), \\ &\text{for all } z \in W^{2,s} \text{ with } \|z - y\|_{W^{2,s}} \leq \kappa. \end{aligned} \quad (4.15)$$

Here,  $c > 0$  denotes the embedding constant such that  $c\|z - y\|_{W^{2,s}} \leq \|z - y\|_{W^{1,\infty}}$ , and  $\kappa > 0$  is the constant governing the admissible functions in (4.14).

Next, we claim that essentially due to the gap around the boundary we have for any fixed  $j$ , there exists  $n_0 = n_0(j)$  such that for all  $n \geq n_0$ ,  $E_{\varepsilon_2(n)}^{CN}(z_j) = 0$ , and the same also holds for close enough perturbations  $z$  of  $z_j$ :

$$E_{\varepsilon_2(n)}^{CN}(z) = 0 \text{ for } n \geq n_0(j), z \text{ with } \|z - z_j\|_{W^{2,s}} < 2\mathcal{E}(h(n)), \quad (4.16)$$

where  $\mathcal{E}(h(n))$  is the maximal Galerkin approximation error from (4.1). We choose  $n_0 = n_0(j)$  (w.l.o.g. increasing) such that for all  $n \geq n_0$ ,

$$\frac{1}{j} \geq s(n) := \frac{2M_1^{d-1}}{\delta} t(n), \quad t(n) := R\varepsilon_2(n) + 2\mathcal{E}(h(n)). \quad (4.17)$$

and  $s(n) \leq \varrho$ . Here,  $R := \text{diam}(\Omega)$ ,  $\delta > 0$ ,  $M_1 > 0$  are the constants from (4.7), and  $\varrho$  is the radius of local invertibility from Lemma 3.6. We claim that for any such  $n$ , all  $z$  with  $\|z - z_j\|_{W^{2,s}} < 2\mathcal{E}(h(n))$  and any pair  $x_1, x_2 \in \Omega$  with  $|x_1 - x_2| > \frac{\varrho}{2}$ ,

$$|z(x_1) - z(x_2)| \geq |z_j(x_1) - z_j(x_2)| - 2\mathcal{E}(h(n)) \geq R\varepsilon_2(n). \quad (4.18)$$

Given (4.18), Corollary 3.7 immediately implies (4.16). We prove (4.18) indirectly (the second inequality; the first follows from the triangle inequality). Suppose that

$$|z_j(x_1) - z_j(x_2)| < t(n) = R\varepsilon_2(n) + 2\mathcal{E}(h(n)). \quad (4.19)$$

Since  $\Psi_j(x_1) \in \Omega_j \subset \Omega$  (and  $\Psi_j(x_2)$  likewise), we have that

$$\text{dist}(\Psi_j(x_1); \partial\Omega) > \frac{1}{j} \geq s(n) \quad (4.20)$$

by (4.12) and (4.17). But on the other hand, if (4.20) holds, then  $B_{s(n)}(\Psi_j(x_1)) \subset \Omega$ . Since  $y$  is bi-Lipschitz on  $B_{s(n)}(\Psi_j(x_1))$  ( $s(n) \leq \varrho$ ;

also notice that the constant factor linking  $s(n)$  and  $t(n)$  is exactly the constant from the lower bound in (3.8), we infer from (4.19) that

$$z_j(x_2) \in B_{t(n)}(z_j(x_1)) = B_{t(n)}(y(\Psi_j(x_1))) \subset y(B_{s(n)}(\Psi_j(x_1))).$$

This is impossible because  $z_j(x_2) = y(\Psi_j(x_2))$  and  $y$  and  $\Psi_j$  are injective, which concludes the proof of (4.16).

In particular, we may use  $z = z_j^{(h(n))}$  in (4.16) with a sufficiently close Galerkin approximation  $z_j^{(h(n))} \in Y_{h(n)}$  of  $z_j$ . Now take  $j(n) \rightarrow \infty$  as  $n \rightarrow \infty$ , but slow enough so that  $n_0(j(n)) \leq n$ . With this choice,

$$\|y_n - z_{j(n)}\|_{W^{2,s}} < 2\mathcal{E}(h(n)) \xrightarrow{n \rightarrow \infty} 0 \quad \text{for } y_n := z_{j(n)}^{(h(n))} \in Y_{h(n)}$$

which together with (4.13) implies that  $y_n \rightarrow y$  in  $W^{2,s}$ . By Proposition 4.14, we see that

$$E_\sigma^{reg}(y_n) \rightarrow E_\sigma^{reg}(y) \quad \text{as } n \rightarrow \infty, \quad (4.21)$$

and from (4.15), we infer that

$$E_{\varepsilon_1(n)}^{el}(y_n) \rightarrow E^{el}(y) \quad \text{as } n \rightarrow \infty. \quad (4.22)$$

Finally, (4.16) yields that

$$E_{\varepsilon_2(n)}^{CN}(y_n) = 0 \quad \text{for all } n \text{ large enough} \quad (4.23)$$

(more precisely,  $n$  large enough so that (4.14) holds for  $z = z_{j(n)}$  and  $z = y_n$ ). Altogether,  $E_{\sigma, \varepsilon(n)}^{h(n)}(y_n) \rightarrow E_\sigma(y)$  as asserted.  $\square$

## 5. NUMERICAL EXPERIMENTS

We consider  $d = 2$  and the approximate deformation

$$y = (y_1, y_2) \in W^{2,s}(\Omega; \mathbb{R}^2)$$

is searched for as the (ideally global) minimizer of

$$E_{\varepsilon, \sigma, \mu}(y) = E_{\varepsilon_1}^{el}(y) + \mu E_{\varepsilon_2}^{CN}(y) + E_\sigma^{reg}(y) + E^{body}(y) \quad (5.1)$$

over a finite dimensional space. Both deformation components  $y_1, y_2$  are discretized in the space of the Bogner-Fox-Schmit (BFS) rectangular elements [8] that provide continuous differentiability of approximations. Here,

$$E^{body}(y) := \int_{\Omega} g_{\text{body}}(x) \cdot y(x) \, dx$$

is the energy contribution of a body force of type (4.2), with  $g_{\text{body}}$  as specified below (in Model II;  $g_{\text{body}} = 0$  in Model I). We here fix the constants

$$\varepsilon_1 := \frac{1}{100}, \quad \sigma := 1, \quad s := 4, \quad p := 4 \quad \text{and} \quad q := 6$$

(in particular, (2.5) is satisfied for  $d = 2$ ). The Ciarlet-Nečas penalty term  $E_{\varepsilon_2}^{CN}$  is included with a weight  $\mu$  for which we tested several values

between 0 and 1; notice that  $\mu = 0$  completely switches off the penalty term. As before,  $E_{\varepsilon_2}^{CN}$  is given by (3.3), and for the computational tests, we chose  $g(t) = t$ ,  $\beta = 1.8$  or  $\beta = 2.2$  (two values close to  $d = 2$ , the threshold for Corollary 3.8), and also experiment with different values for  $\varepsilon_2$ , typically adapted to the grid size  $h$ . As an example for the higher order term, we employ

$$E_{\sigma}^{reg}(y) = \sigma \int_{\Omega} |D^2 y(x)|^s dx.$$

Finally, we use the elastic part of the energy

$$E_{\varepsilon_1}^{el}(y) = \int_{\Omega} W_{\varepsilon_1}^{el}(\nabla y(x)) dx$$

with density chosen as follows, for all  $F \in \mathbb{R}^{2 \times 2}$  and  $J := \det F$ :

$$W_{\varepsilon_1}^{el}(F) := |F|^p - d^{\frac{p}{2}} - \frac{p}{q} d^{\frac{p}{2}-1} + \frac{p}{q} d^{\frac{p}{2}-1} \begin{cases} J^{-q} & \text{if } J \geq \varepsilon_1, \\ -q\varepsilon_1^{-q-1}(J - \varepsilon_1) + \varepsilon_1^{-q} & \text{if } J < \varepsilon_1. \end{cases}$$

Here, recall that  $d = 2$ , although the example could also be used for higher dimensions. Above,

$$\nabla y(x) \in \mathbb{R}^{2 \times 2}, \quad D^2 y(x) \in \mathbb{R}^{2 \times 2 \times 2}$$

and denote the gradient and the Hessian of  $y : \Omega \rightarrow \mathbb{R}^2$ , respectively, and the norms  $|\cdot|$  are euclidean (Frobenius):  $|F| := (\sum_{i,j} F_{ij}^2)^{\frac{1}{2}}$  for  $F = (F_{ij}) \in \mathbb{R}^{2 \times 2}$ ,  $|G| := (\sum_{i,j,k} G_{ijk}^2)^{\frac{1}{2}}$  for  $G = (G_{ijk}) \in \mathbb{R}^{2 \times 2 \times 2}$ .

*Remark 5.1.*  $W_{\varepsilon_1}^{el}$  is polyconvex and frame indifferent. Moreover, if  $0 < \varepsilon_1 < 1$ ,  $W_{\varepsilon_1}^{el}(F) \geq 0$  with equality if and only if  $F$  is a rotation matrix. The second part of  $W_{\varepsilon_1}^{el}$  is a  $C^1$ -function in  $J = \det F$ , and the two cases define a truncated version of  $J \mapsto J^{-q}$  using an affine extension for  $J < \varepsilon_1$ . For the shapes, body forces and boundary conditions used in our examples, there is no incentive for the material to create spots with high local compression. As a consequence, the results are independent of the choice of  $\varepsilon_1 \ll 1$ , as the computed deformations stay far away from the regime  $\det \nabla y < \varepsilon_1$  anyway (as the optimal deformations are expected to), for any reasonably small choice of  $\varepsilon_1$ .

*Remark 5.2.* For the actual computations, we have replaced the non-differentiable functions  $[\cdot]^+$  and  $g(|\cdot|)$  appearing in the definition of  $E_{\varepsilon_2}^{CN}$  by  $C^1$ -approximations  $h([\cdot]^+)$  and  $h(g(|\cdot|))$ , where

$$h(x) := \begin{cases} 0, & \text{for } x \leq 0, \\ x^2/(2a) & \text{for } 0 \leq x \leq a, \\ x - a/2, & \text{for } a \leq x \end{cases}$$

for some small parameter  $a > 0$  (we set  $a = 1/10$  in all computations). This smoothing is introduced in order to avoid the risk that the actual minimizer sits at a point where the functional does not have a well defined derivative. The latter could cause serious problems for the solver. While changing  $g$  is fully covered by the theory developed above and can at most affect constants appearing in the theoretical results, changing  $[\cdot]^+$  even slightly around zero can potentially effect the scaling of  $E_{\varepsilon_2}^{CN}$  as  $\varepsilon_2 \rightarrow 0$ . But in practice the asymptotics as  $\varepsilon_2 \rightarrow 0$  cannot easily be observed numerically anyway.

5.1. **Model I.** We consider a reference configuration  $\Omega = \Omega_1 \cup \Omega_2 \subset \mathbb{R}^2$  which consists of two rectangular boxes  $\Omega_1 = (0, 2) \times (0.5, 1.5)$  and  $\Omega_2 = (0, 2) \times (-1.5, 0.5)$ . See Figure 3 for illustration. We impose non-

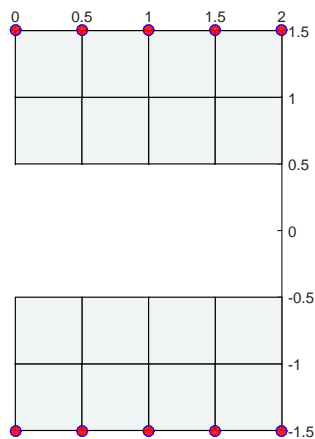


FIGURE 3. Model I : Coarse rectangular mesh with Dirichlet boundary nodes indicated by full dots.

homogeneous Dirichlet boundary conditions on a part of the boundary:

$$\begin{aligned} y_1 &= x_1 + m_1, y_2 = x_2 - m_2 & \text{for } (x_1, x_2) \in \Lambda_{D,1}, \\ y_1 &= x_1, \quad y_2 = x_2 + m_2 & \text{for } (x_1, x_2) \in \Lambda_{D,2}, \end{aligned}$$

where  $\Lambda_{D,1} := (0, 2) \times \{1.5\}$  and  $\Lambda_{D,2} := (0, 2) \times \{-1.5\}$  and  $m_1, m_2 \in \mathbb{R}$  are parameters. There is no linear body force term considered in this model. We consider a sequence of minimization problems with parameters

$$m_2 \in \{0.4, 0.5, 0.6, 0.7\}$$

and the same parameters  $m_1 = 0.2$  and  $\mu = 1$  in two different cases

$$\varepsilon_2 = 1/4, \quad \varepsilon_2 = 1/8.$$

Figure 5 displays how much the total energy  $E_{\varepsilon, \sigma, \mu}(y)$  and the scaled penetration energy  $E_{\varepsilon_2}^{CN}(y)$  depend on the parameter  $m_2$ . Unsurprisingly, both the energy and the influence of  $E_{\varepsilon_2}^{CN}(y)$  grow the larger  $m_2$  gets, i.e., the more the two pieces are pushed against each other by the

boundary conditions. Since there is no linear body force term considered in this model, the difference energy  $E_{\varepsilon,\sigma,\mu}(y) - E_{\varepsilon_2}^{CN}(y)$  converts to  $E_{\varepsilon_1}^{el}(y)$  and  $E_{\sigma}^{reg}(y)$ . Figure 4 shows a few resulting optimized domains. Choosing  $\varepsilon_2 = 1/8$  is more or less the smallest reasonable choice for  $\varepsilon_2$  given the grid size (cf. Remark 3.1). In this case, one can already see effects of errors due to the numerical integration in the marginal density of  $E_{\varepsilon_2}^{CN}$ , which is much more uniform and intuitive for  $\varepsilon_2 = 1/4$ .

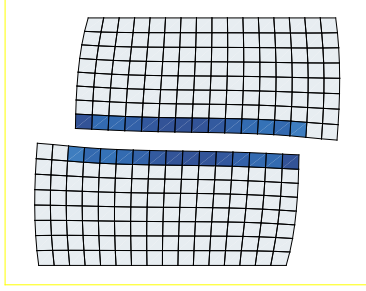
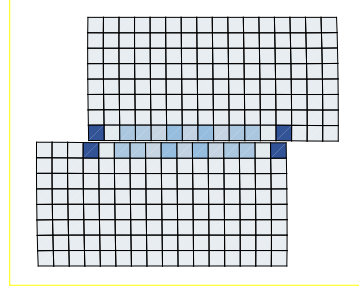
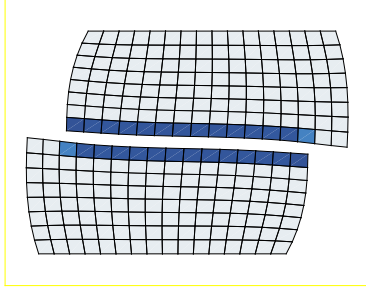
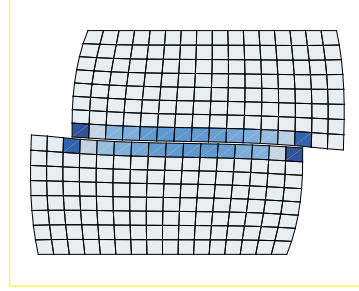
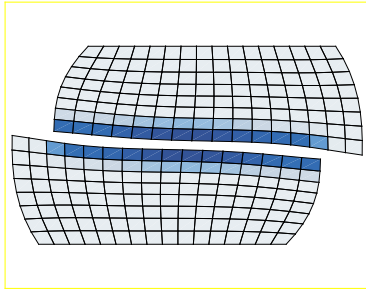
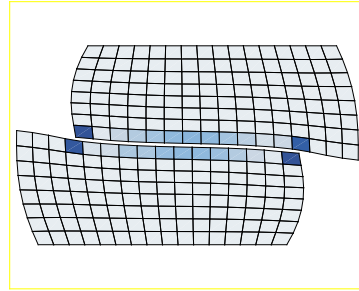
(A)  $m_2 = 0.5, \varepsilon_2 = 1/4$ .(D)  $m_2 = 0.5, \varepsilon_2 = 1/8$ .(B)  $m_2 = 0.6, \varepsilon_2 = 1/4$ .(E)  $m_2 = 0.6, \varepsilon_2 = 1/8$ .(C)  $m_2 = 0.7, \varepsilon_2 = 1/4$ .(F)  $m_2 = 0.7, \varepsilon_2 = 1/8$ .

FIGURE 4. Model I: Optimized deformed domains with underlying marginal density of  $E_{\varepsilon_2}^{CN}(y)$ .

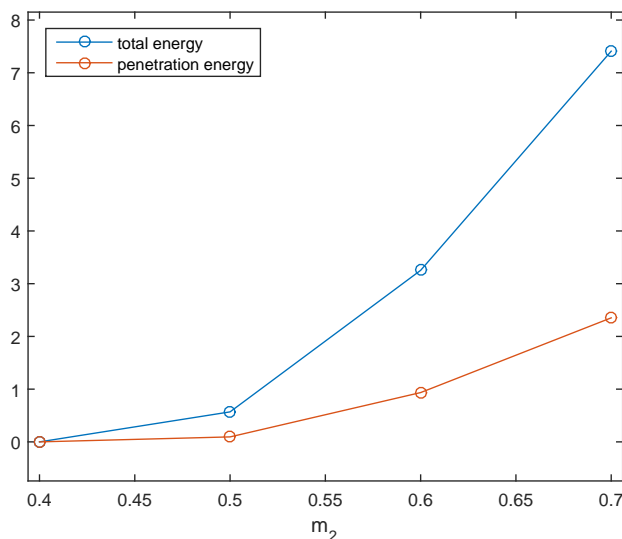


FIGURE 5. Model I: Total energy  $E_{\varepsilon, \sigma, \mu}(y)$  and penetration energy  $E_{\varepsilon_2}^{CN}$  displayed versus the parameter  $m_2$ . We consider the case  $\varepsilon_2 = 1/4$ .

5.2. **Model II.** We consider the same “pincers” domain  $\Omega \subset \mathbb{R}^2$  as in example of Subsection 3.1 and subject  $\Omega$  to the linear body force density

$$g_{\text{body}}(x_1, x_2) = \nu(0, -H(x_1) \text{sign}(x_2)) \quad \text{on } (x_1, x_2) \in \Omega,$$

where  $H$  denotes the Heaviside step function. In addition, we impose Dirichlet boundary conditions on a part of the boundary:

$$y_1 = x_1, \quad y_2 = x_2 \quad \text{for } (x_1, x_2) \in \Lambda_D,$$

where  $\Lambda_D := \{0\} \times (-1/2, 1/2)$ . See Figure 6 for illustration.

In this model, we measure the response of an elastic continuum to various scaling of the Ciarlet-Nečas penalty term  $\mu E_{\varepsilon_2}^{CN}$ . We consider a sequence of minimization problems with various multipliers

$$\mu \in \{10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$$

in two different cases

$$\varepsilon_2 = 1/2, \quad \varepsilon_2 = 1/4.$$

and the same loading given by  $\nu = 0.2$ . Figure 7 displays how much the total energy  $E_{\varepsilon, \sigma, \mu}(y)$  and the scaled penetration energy  $\mu E_{\varepsilon_2}^{CN}(y)$  depend on the multiplier  $\mu$ .

For lower values of  $\mu$  the scaled penetration term allows for a penetration of both pincers parts. For higher values of  $\mu$  the scaled penetration term one can usually prevent penetration altogether. Here, recall that

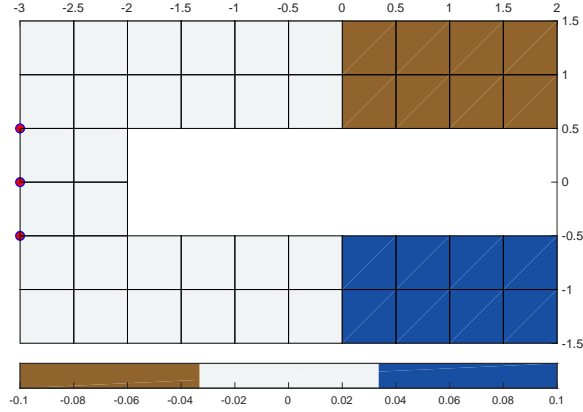


FIGURE 6. Model II : body force density in  $x_2$  direction forcing both pincers part to move against each other and Dirichlet boundary nodes indicated by full dots.

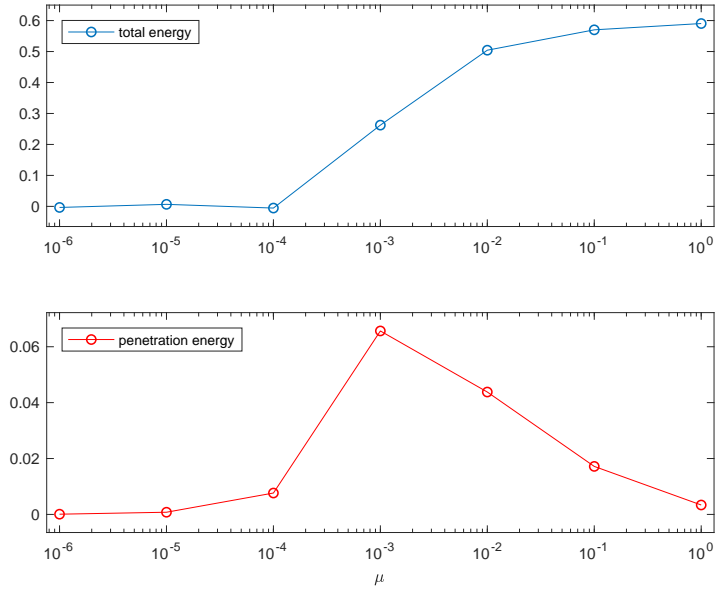
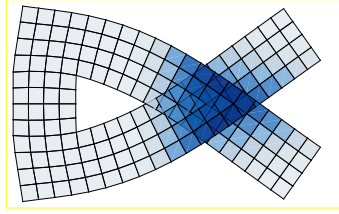


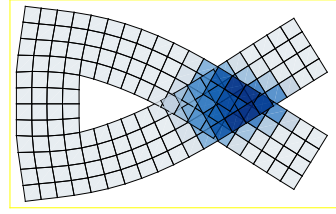
FIGURE 7. Model II: Total energy  $E_{\varepsilon,\sigma,\mu}(y)$  and penetration energy  $\mu E_{\varepsilon_2}^{CN}$  displayed versus the parameter  $\mu$ . We consider the case  $\varepsilon_2 = 1/2$ .

from the point of view of the theory, we only know for sure that reducing  $\varepsilon_2$  will eventually prevent penetration if the scaling exponent  $\beta$  in  $E_{\varepsilon_2}^{CN}$  is big enough (Corollary 3.8). Of course, at finite scales increasing  $\mu$  has a similar effect, and as long as the grid size  $h$  is fixed, we cannot arbitrarily reduce  $\varepsilon_2$ .

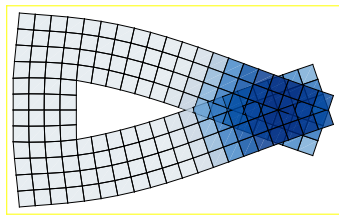




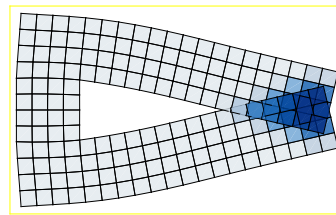
(A)  $\mu_2 = 10^{-4}, \epsilon_2 = 1/2$ .



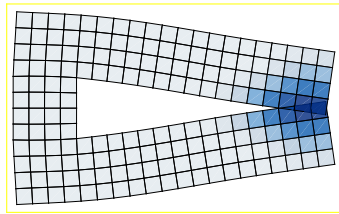
(E)  $\mu_2 = 10^{-4}, \epsilon_2 = 1/4$ .



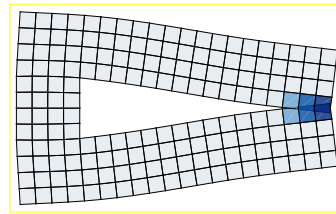
(B)  $\mu_2 = 10^{-3}, \epsilon_2 = 1/2$ .



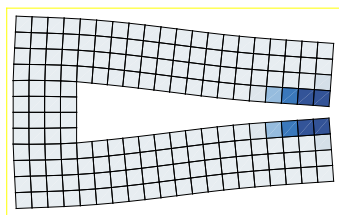
(F)  $\mu_2 = 10^{-3}, \epsilon_2 = 1/4$ .



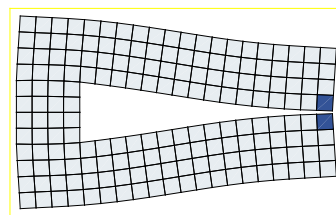
(C)  $\mu_2 = 10^{-2}, \epsilon_2 = 1/2$ .



(G)  $\mu_2 = 10^{-2}, \epsilon_2 = 1/4$ .



(D)  $\mu_2 = 10^{-1}, \epsilon_2 = 1/2$ .



(H)  $\mu_2 = 10^{-1}, \epsilon_2 = 1/4$ .

FIGURE 8. Model II: Optimized deformed domains with underlying marginal density of  $\mu E_{\epsilon_2}^{CN}(y)$ .

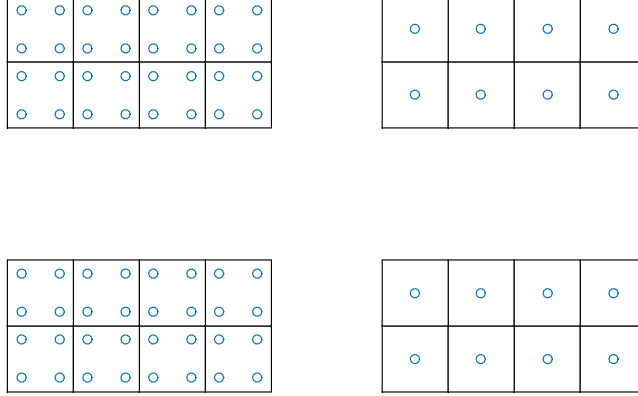


FIGURE 9. Model I : Four Gauss integration points (left) used for evaluation of all energy parts except the penalty term, midpoints (right) used in the evaluation of the penetration penalty term.

**5.3. Remarks on implementation.** Our Matlab code is based on former codes related to [1, 16, 24] that allow for a vectorized assembly of finite element matrices. The code available for download and includes an own implementation of the Bogner-Fox-Schmit (BFS) rectangular elements for a uniformly refined rectangular mesh, where all rectangular elements are for simplicity of the same size  $hx_1 \times hx_2$ . Both Model I and Model II rectangular meshes are of this kind. The basis functions on each rectangle are based on bicubic polynomials, i.e. tensor products of 4 cubic (Hermite) polynomials. They have 16 degrees of freedom with 4 degrees in each of its 4 corner nodes approximating: a function value, its gradient and the second mixed derivative. Therefore, a given scalar function  $u \in C^1(\Omega)$  is represented by a matrix of the size  $nn \times 4$  in the form

$$\underline{\mathbf{u}} = \begin{pmatrix} u(x_1^1, x_2^1) & \frac{\partial u}{\partial x_1}(x_1^1, x_2^1) & \frac{\partial u}{\partial x_2}(x_1^1, x_2^1) & \frac{\partial^2 u}{\partial x_1 \partial x_2}(x_1^1, x_2^1) \\ \dots & \dots & \dots & \dots \\ u(x_1^{nn}, x_2^{nn}) & \frac{\partial u}{\partial x_1}(x_1^{nn}, x_2^{nn}) & \frac{\partial u}{\partial x_2}(x_1^{nn}, x_2^{nn}) & \frac{\partial^2 u}{\partial x_1 \partial x_2}(x_1^{nn}, x_2^{nn}) \end{pmatrix},$$

where  $nn$  denotes the total number of mesh nodes and  $(x_1^i, x_2^i)$  for  $i = 1, \dots, nn$  their corresponding coordinates. The construction of BFS elements additionally guarantees  $\frac{\partial^2 u}{\partial x_1 \partial x_2} \in C(\Omega)$  but the remaining second-order derivatives are generally discontinuous. Based on a global numbering of nodes, the matrix  $\underline{\mathbf{u}}$  is further reformatted as a column vector  $\mathbf{u}$  with  $4 \cdot nn$  entries. For our 2d nonlinear elasticity computations, we approximate both components  $y_1, y_2$  by the BFS elements and resulting vector variable  $y = (y_1, y_2) \in C^1(\Omega; \mathbb{R}^2)$  has  $8 \cdot nn$  entries.

Since energy parts of  $E_{\varepsilon, \sigma}$  are generally non-quadratic functionals, all two-dimensional integrals are evaluated using the Gaussian quadrature,

where integration points are tensor products of one-dimensional Gauss integration points. We deploy four Gauss quadrature points in each rectangle as illustrated in Figure 9 (left).

The penetration penalty term  $E_{\varepsilon_2}^{CN}$  is a nonlocal functional. We make no additional assumptions on the location of the penetration and evaluate all pairwise euclidean distances

$$|(x_1^i, x_2^i) - (x_1^j, x_2^j)|, \quad |(y_1^i, y_2^i) - (y_1^j, y_2^j)|$$

in a double loop over  $i, j = 1, \dots, ne$ . Here, vectors above denote coordinates of rectangles midpoints (cf. the right part of Figure 9) and their corresponding deformations and  $ne$  the number of mesh rectangles. The x-distances above are precomputed, the y-distances need to be recomputed in every evaluation of the penetration penalty term.

*Remark 5.3* (Possible implementation improvement). Even if no assumptions are made on the location of the penetration, it is possible (but not implemented here) to optimize the evaluation of  $E_{\varepsilon_2}^{CN}$  as follows:

- Instead of a full double loop, first only go through all pairs of elements located at the boundary of the domain. Create a list of those elements that contribute to  $E_{\varepsilon_2}^{CN}$ .
- Then start to search for other contributing elements in the interior by repeatedly checking all elements that are neighbors of those that are already known to give a positive contribution.
- Stop when no new contributing neighbor elements are found.

While the full double loop requires a number of steps of the order of  $h^{-2d}$  for the mesh size  $h$ , the double loop through the elements at the boundary only needs  $h^{-2(d-1)}$ . As long as there is no deep penetration (penetration depth of the order of  $h$  or less) and  $\varepsilon_2 = O(h)$  ( $\varepsilon_2 \geq h$ , but not that much bigger), a subsequent search for contributing neighbor elements does not increase that significantly, either.

#### ACKNOWLEDGEMENTS

Our research was supported by the Czech Science Foundation (GA ĀR), through the grants GA18-03834S (SK and JV) and GA17-04301S (JV).

#### REFERENCES

- [1] Immanuel Anjam and Jan Valdman. Fast MATLAB assembly of FEM matrices in 2D and 3D: edge elements. *Appl. Math. Comput.*, 267:252–263, 2015.
- [2] J.M. Ball. Global invertibility of Sobolev functions and the interpenetration of matter. *Proc. R. Soc. Edinb., Sect. A, Math.*, 88:315–328, 1981.
- [3] John M. Ball. Convexity conditions and existence theorems in nonlinear elasticity. *Arch. Rational Mech. Anal.*, 63(4):337–403, 1977.

- [4] John M. Ball. Convexity conditions and existence theorems in nonlinear elasticity. *Arch. Rational Mech. Anal.*, 63(4):337–403, 1977.
- [5] John M. Ball. Some open problems in elasticity. In *Geometry, mechanics, and dynamics*, pages 3–59. Springer, New York, 2002.
- [6] S. Bartels and P. Reiter. Stability of a simple scheme for the approximation of elastic knots and self-avoiding inextensible curves. Preprint arXiv:1804.02206 (2018).
- [7] B. Benešová, M. Kružík, and A. Schlömerkemmer. A note on locking materials and gradient polyconvexity. Preprint arXiv:1706.04055, 2017. To appear in M3AS.
- [8] F. K. Bogner, R. L. Fox, and L. A. Schmit. The generation of Inter-element Compatible Stiffness and Mass Matrices by the Use of Interpolation Formulas. *Proceedings of the Conference on Matrix Methods in Structural Mechanics*, pages 397–444, 1965.
- [9] Philippe G. Ciarlet. *Mathematical elasticity. Vol. I*, volume 20 of *Studies in Mathematics and its Applications*. North-Holland Publishing Co., Amsterdam, 1988. Three-dimensional elasticity.
- [10] Philippe G. Ciarlet and Jindřich Nečas. Injectivity and self-contact in nonlinear elasticity. *Arch. Ration. Mech. Anal.*, 97:173–188, 1987.
- [11] Bernard Dacorogna. *Direct methods in the calculus of variations*, volume 78 of *Applied Mathematical Sciences*. Springer, New York, second edition, 2008.
- [12] Gianni Dal Maso. *An introduction to  $\Gamma$ -convergence*. Number 8 in Progress in Nonlinear Differential Equations and their Applications. Birkhäuser, Basel, 1993.
- [13] I. Fonseca and W. Gangbo. Local invertibility of Sobolev functions. *SIAM J. Math. Anal.*, 26(2):280–304, 1995.
- [14] Irene Fonseca and Wilfrid Gangbo. *Degree theory in analysis and applications*. Oxford: Clarendon Press, 1995.
- [15] M. Foss, W. J. Hrusa, and V. J. Mizel. The Lavrentiev gap phenomenon in nonlinear elasticity. *Arch. Ration. Mech. Anal.*, 167(4):337–365, 2003.
- [16] P. Harasim and J. Valdman. Verification of functional a posteriori error estimates for an obstacle problem in 2D. *Kybernetika*, 50:978–1002, 2014.
- [17] Timothy J. Healey and Stefan Krömer. Injective weak solutions in second-gradient nonlinear elasticity. *ESAIM, Control Optim. Calc. Var.*, 15(4):863–871, 2009.
- [18] Duvan Henao and Carlos Mora-Corral. Regularity of inverses of sobolev deformations with finite surface energy. *Journal of Functional Analysis*, 268(8):2356–2378, 2015.
- [19] Jan Malý, David Swanson, and William P. Ziemer. The co-area formula for Sobolev mappings. *Trans. Am. Math. Soc.*, 355(2):477–492, 2003.
- [20] Alexander Mielke and Tomáš Roubíček. Rate-independent elastoplasticity at finite strains and its numerical approximation. *Math. Models Methods Appl. Sci.*, 26(12):2203–2236, 2016.
- [21] Pablo V. Negrón Marrero. A numerical method for detecting singular minimizers of multidimensional problems in nonlinear elasticity. *Numer. Math.*, 58(2):135–144, 1990.
- [22] Aaron Z. Palmer. Variations of deformations with self-contact on Lipschitz domains. *Set-Valued Var. Anal (online first)*, pages 1–11, 2018.
- [23] Aaron Z. Palmer and Timothy J. Healey. Injectivity and self-contact in second-gradient nonlinear elasticity. *Calc. Var. Partial Differ. Equ.*, 56(4):11, 2017.
- [24] Talal Rahman and Jan Valdman. Fast MATLAB assembly of FEM matrices in 2D and 3D: nodal elements. *Appl. Math. Comput.*, 219(13):7151–7158, 2013.

- [25] Miroslav Šilhavý. *The mechanics and thermodynamics of continuous media*. Texts and Monographs in Physics. Springer, Berlin, 1997.

STEFAN KRÖMER, THE CZECH ACADEMY OF SCIENCES, INSTITUTE OF INFORMATION THEORY AND AUTOMATION, POD VODÁRENSKOU VĚŽÍ 4, 182 08 PRAHA 8, CZECH REPUBLIC (CORRESPONDING AUTHOR), *E-mail address*: SKROEMER@UTIA.CAS.CZ

JAN VALDMAN, THE CZECH ACADEMY OF SCIENCES, INSTITUTE OF INFORMATION THEORY AND AUTOMATION, POD VODÁRENSKOU VĚŽÍ 4, 182 08 PRAHA 8, CZECH REPUBLIC, *E-mail address*: VALDMAN@UTIA.CAS.CZ