# Metric Selection in Douglas-Rachford Splitting and ADMM

Pontus Giselsson* and Stephen Boyd†

*Abstract*—Recently, several convergence rate results for Douglas-Rachford splitting and the alternating direction method of multipliers (ADMM) have been presented in the literature. In this paper, we show linear convergence rate bounds for Douglas-Rachford splitting under strong convexity and smoothness assumptions. We show that these bounds generalize and/or improve on similar bounds in the literature and that the bounds are tight for the class of problems under consideration. For finite dimensional Euclidean problems, we show how the rate bounds depend on the metric that is used in the algorithm. We also show how to select this metric to optimize the bound. Many real-world problems do not satisfy both the smoothness and strongly convexity assumptions. Therefore, we also propose heuristic metric and parameter selection methods to improve the performance of a wider class of problems. The efficiency of the proposed heuristics is confirmed in a numerical example on a model predictive control problem, where improvements of more than one order of magnitude are observed.

## I. INTRODUCTION

Optimization problems of the form

$$\begin{array}{ll} \text{minimize} & f(x) + g(y) \\ \text{subject to} & \mathcal{A}x = y \end{array} \qquad (1)$$

where $x$ is the variable, $f$ and $g$ are convex, and $\mathcal{A}$ is a bounded linear operator, arise in numerous applications ranging from compressed sensing [8] and statistical estimation [24] to model predictive control [38] and medical imaging [31]. There exist a variety of algorithms for solving convex problems of the form (1), many of which are treated in the monograph [33]. The methods include primal and dual forward-backward splitting methods [10] and their accelerated variants [4], the Arrow-Hurwicz method [1], Douglas-Rachford splitting [15] and Peaceman-Rachford splitting [35], the alternating direction method of multipliers (ADMM) [23], [18], [7] (which is Douglas-Rachford splitting applied to the dual problem [17], [16]), and linearized ADMM [9].

In this paper, we focus on generalized Douglas-Rachford splitting, which includes as special cases Douglas-Rachford splitting and Peaceman-Rachford splitting when applied to the primal problem, and under- and over-relaxed ADMM when applied to the dual problem. These methods have long been known to converge under very mild assumptions, [18], [30], [16]. However, the rate of convergence in the general case has just recently been shown to be $O(1/k)$, [25], [13], [11]. For a restricted class of problems, Lions and Mercier showed in [30] that the Douglas-Rachford algorithm enjoys a linear convergence rate. To the authors' knowledge, this was the sole linear convergence rate results for a long period of

* Department of Automatic Control, Lund University. Email: `pontusg@control.lth.se`.
† Electrical Engineering Department, Stanford University. Email: `boyd@stanford.edu`.

time for these methods. Recently, however, many works have shown linear convergence rates for Douglas-Rachford splitting, Peaceman-Rachford splitting and ADMM in different settings [26], [36], [13], [12], [14], [19], [34], [27], [28], [6], [40], [3], [37], [32]. The works in [26], [13], [6], [36] concern local linear convergence under different assumptions. The works in [27], [28], [40] consider distributed formulations, while the works in [12], [14], [19], [34], [30], [3], [37], [32] consider global convergence. The work in [3] shows linear convergence for the Douglas-Rachford splitting when solving a subspace intersection problem. The work in [37] (which appeared online during the submission procedure of this paper) shows linear convergence for equality constrained problems with upper and lower bounds on variables. The remaining works, [12], [14], [19], [34], [30], [32], show linear convergence under strong convexity and smoothness assumptions. In this paper, we generalize and/or improve on these convergence rate estimates. A detailed description on the improvements and generalizations made is found in Section III-B.

Besides improving on existing results, we also provide an example that shows that the convergence rate bounds in this paper are tight for the classes of problems under consideration. We also provide explicit upper bounds on the over-relaxation parameter $\alpha$ in generalized Douglas-Rachford splitting. We show that this can be greater than one in the linearly convergent case (in the general case, the relaxation factor is limited to $\alpha \in (0, 1)$). A similar finding was reported in [32] in the ADMM case. We supplement this finding in [32] by stating explicit upper bounds on $\alpha$.

When solving problems of the form (1) in finite dimensional Euclidean settings, we can choose a Hilbert space with inner product $\langle x, y \rangle_M = x^T M y$ and induced norm on which to apply the generalized Douglas-Rachford algorithm. The algorithm behaves differently for different choices of $M$ and an appropriate choice can significantly speed up the algorithm, both in theory and in practice. Another contribution of this paper is that we show how to select a metric $M$ to optimize the linear convergence rate factor for problems where $f$ is smooth and strongly convex, $g$ is any proper, closed, and convex function, and $\mathcal{A}$ is surjective, i.e., has full row rank. These results are applied to both the primal and dual problems, and therefore apply both to Douglas-Rachford splitting and ADMM (which is Douglas-Rachford splitting applied to the dual problem). This generalizes, in several directions, the work in [19] in which corresponding results for ADMM applied to solve quadratic programs with linear inequality constraints are provided, see Section III-B for a detailed comparison between the results.

Real-world problems rarely have the properties needed to ensure linear convergence of the generalized Douglas-Rachford algorithm or ADMM. Therefore, we provide heuris-

tic metric and parameter selection methods for cases when some of these assumptions are not met. The heuristics cover most optimization problems that have a quadratic part which is not necessarily strongly convex. Such problems arise, e.g., in model predictive control [38], statistical estimation [24] using, e.g., lasso [41], and compressed sensing [8] which can be used, e.g., in medical imaging [31]. A numerical example on a model predictive control problem is provided that shows the efficiency of the proposed metric selection heuristic. For the considered problem, the execution time is decreased with about one order of magnitude compared to when applying the algorithm on the Euclidean space with the standard inner product and induced norm.

This paper extends and generalizes our conference papers [21], [20].

### A. Notation

We denote by $\mathbf{R}$ the set of real numbers, $\mathbf{R}^n$ the set of real column-vectors of length $n$, and $\mathbf{R}^{m \times n}$ the set of real matrices with $m$ rows and $n$ columns. Further $\overline{\mathbf{R}} := \mathbf{R} \cup \{\infty\}$ denotes the extended real line. Throughout this paper $\mathcal{H}$ denotes a real Hilbert space. Its inner product is denoted by $\langle \cdot, \cdot \rangle$, the induced norm by $\| \cdot \|$, and the identity operator by Id. We specifically consider finite-dimensional Hilbert-spaces $\mathbb{H}_H$ with inner product $\langle x, y \rangle = x^T H y$ and induced norm $\|x\| = \sqrt{x^T H x}$. Sometimes we denote these by $\langle \cdot, \cdot \rangle_H$ and $\| \cdot \|_H$. We also sometimes denote the Euclidean inner-product by $\langle x, y \rangle_2 = x^T y$ and the induced norm by $\| \cdot \|_2$ for clarity. The conjugate function is denoted and defined by $f^*(y) \triangleq \sup_x \{\langle y, x \rangle - f(x)\}$. The power set of a set $\mathcal{X}$, i.e., the set of all subsets of $\mathcal{X}$, is denoted by $2^{\mathcal{X}}$. The graph of an (set-valued) operator $A : \mathcal{X} \to 2^{\mathcal{Y}}$ is defined and denoted by $\mathrm{gph} A = \{(x, y) \in \mathcal{X} \times \mathcal{Y} \mid y \in Ax\}$. The inverse operator $A^{-1}$ is defined through its graph by $\mathrm{gph} A^{-1} = \{(y, x) \in \mathcal{Y} \times \mathcal{X} \mid y \in Ax\}$. Finally, the class of closed, proper, and convex functions $f : \mathcal{H} \to \overline{\mathbf{R}}$ is denoted by $\Gamma_0(\mathcal{H})$.

## II. Background

In this section, we introduce some standard definitions that can be found, e.g. in [2], [39].

*Definition 1 (Strong monotonicity):* An operator $A : \mathcal{H} \to 2^{\mathcal{H}}$ is *$\sigma$-strongly monotone* with $\sigma > 0$ if

$$\langle u - v, x - y \rangle \geq \sigma \|x - y\|^2$$

for all $(x, u) \in \mathrm{gph}(A)$ and $(y, v) \in \mathrm{gph}(A)$.
The definition of *monotonicity* is obtained by setting $\sigma = 0$ in the above definition. In the following definitions, we suppose that $\mathcal{D} \subseteq \mathcal{H}$ is a nonempty subset of $\mathcal{H}$.

*Definition 2 (Lipschitz continuity):* A mapping $A : \mathcal{D} \to \mathcal{H}$ is *$\beta$-Lipschitz continuous* with $\beta \geq 0$ if

$$\|Ax - Ay\| \leq \beta \|x - y\|$$

holds for all $x, y \in \mathcal{D}$. If $\beta = 1$ then $A$ is *nonexpansive* and if $\beta \in (0, 1)$ then $A$ is *$\beta$-contractive*.

*Definition 3 (Averaged mappings):* A mapping $A : \mathcal{D} \to \mathcal{H}$ is *$\alpha$-averaged* if there exists a nonexpansive mapping $B : \mathcal{D} \to \mathcal{H}$ and $\alpha \in (0, 1)$ such that $A = (1 - \alpha)\mathrm{Id} + \alpha B$.

*Definition 4 (Cocoercivity):* A mapping $A : \mathcal{D} \to \mathcal{H}$ is *$\beta$-cocoercive* with $\beta > 0$ if

$$\langle Ax - Ay, x - y \rangle \geq \beta \|Ax - Ay\|^2$$

holds for all $x, y \in \mathcal{D}$.
Mappings that are 1-cocoercive (or equivalently $\frac{1}{2}$-averaged) are also *firmly nonexpansive*.

*Definition 5 (Strong convexity):* A function $f \in \Gamma_0(\mathcal{H})$ is *$\sigma$-strongly convex* with $\sigma > 0$ if $f - \frac{\sigma}{2}\| \cdot \|^2$ is convex.

A strongly convex function has a minimum curvature that is decided by $\sigma$. If $f$ is differentiable, strong convexity can equivalently be defined as that

$$f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle + \frac{\sigma}{2}\|x - y\|^2 \qquad (2)$$

holds for all $x, y \in \mathcal{H}$. Functions with a maximal curvature are called smooth. Next, we present a smoothness definition for convex functions.

*Definition 6 (Smoothness for convex functions):* A function $f \in \Gamma_0(\mathcal{H})$ is *$\beta$-smooth* with $\beta \geq 0$ if it is differentiable and $\frac{\beta}{2}\| \cdot \|^2 - f$ is convex, or equivalently that

$$f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{\beta}{2}\|x - y\|^2 \qquad (3)$$

holds for all $x, y \in \mathcal{H}$.
*Remark 1:* It can be seen from (2) and (3) that for a function that is $\sigma$-strongly convex and $\beta$-smooth, we always have $\beta \geq \sigma$.

## III. Generalized Douglas-Rachford splitting

The generalized Douglas-Rachford algorithm can be applied to solve composite convex optimization problems of the form

$$\text{minimize} \quad f(x) + g(x) \qquad (4)$$

where $f, g \in \Gamma_0(\mathcal{H})$. The solution to (4) is characterized by the following optimality conditions, [2, Proposition 25.1]

$$z = R_{\gamma f} R_{\gamma g} z, \qquad x = \mathrm{prox}_{\gamma f}(z)$$

where $\gamma > 0$, and the prox operator $\mathrm{prox}_{\gamma f}$ and the reflected proximal operator $R_{\gamma f}$ are defined as

$$\mathrm{prox}_{\gamma f}(z) = \arg\min_x \left\{ \gamma f(x) + \tfrac{1}{2}\|x - z\|^2 \right\}$$
$$R_{\gamma f} = 2\mathrm{prox}_{\gamma f} - \mathrm{Id}$$

respectively. In other words, the solution to (4) is found by applying the proximal operator on $z$, where $z$ is a fixed-point to $R_{\gamma f} R_{\gamma g}$. One approach to find a fixed-point to $R_{\gamma f} R_{\gamma g}$ is to iterate the composition

$$z^{k+1} = R_{\gamma f} R_{\gamma g} z^k.$$

This algorithm is known as Peaceman-Rachford splitting, [35]. However, since $R_{\gamma f}$ and $R_{\gamma g}$ are nonexpansive in the general case, so is their composition, and convergence of this algorithm cannot be guaranteed in the general case. The generalized Douglas-Rachford splitting algorithm is obtained by iterating the averaged map of the nonexpansive Peaceman-Rachford operator $R_{\gamma f} R_{\gamma g}$. That is, it is given by the iteration

$$z^{k+1} = ((1 - \alpha)\mathrm{Id} + \alpha R_{\gamma f} R_{\gamma g}) z^k \qquad (5)$$

where $\alpha \in (0,1)$ to guarantee convergence in the general case. (We will, however, see that when additional regularity assumptions are introduced to the problem, $\alpha = 1$, i.e. Peaceman-Rachford splitting, and even some $\alpha > 1$ can be used and convergence can still be guaranteed.) The algorithm known as Douglas-Rachford splitting is obtained by letting $\alpha = \frac{1}{2}$ in (5), but we will use the term Douglas-Rachford splitting for all values of $\alpha$.

### A. Linear convergence

Under some regularity assumptions, the convergence of the Douglas-Rachford algorithm is linear. We will analyze the convergence under the following set of assumptions:

*Assumption 1:* Suppose that:

(i) $f$ and $g$ are proper, closed, and convex.

(ii) $f$ is $\sigma$-strongly convex and $\beta$-smooth.

To show linear convergence rates of the Douglas-Rachford algorithm under these regularity assumptions on $f$, we need to characterize some properties of the proximal and reflected proximal operators to $f$. Specifically, we will show that the reflected proximal operator of $f$ is contractive (as opposed to nonexpansive in the general case) and we will also provide a tight contraction factor. The key to arriving at this contraction factor of the reflected proximal operator is the following, to the authors' knowledge, novel (but straightforward) interpretation of the proximal operator.

*Proposition 1:* Assume that $f \in \Gamma_0(\mathcal{H})$ and define $f_\gamma$ as

$$f_\gamma := \gamma f + \tfrac{1}{2}\|\cdot\|^2. \tag{6}$$

Then $\operatorname{prox}_{\gamma f}(y) = \nabla f_\gamma^*(y)$.

*Proof.* Since the proximal operator is the resolvent of $\gamma \partial f$, see [2, Example 23.3], we have $\operatorname{prox}_{\gamma f}(y) = (\mathrm{Id} + \gamma \partial f)^{-1} y = (\partial f_\gamma)^{-1} y$. Since $f \in \Gamma_0(\mathcal{H})$ also $f_\gamma \in \Gamma_0(\mathcal{H})$. Therefore [2, Corollary 16.24] implies that $\operatorname{prox}_{\gamma f}(y) = (\partial f_\gamma)^{-1} y = \nabla f_\gamma^*(y)$, where differentiability of $f_\gamma^*$ follows from [2, Theorem 18.15], since $f_\gamma$ is 1-strongly convex, and since $f = (f^*)^*$ for proper, closed, and convex functions, see [2, Theorem 13.32]. This concludes the proof. $\square$

Using this interpretation of the proximal operator, we can show the following proposition which is proven in Appendix A.

*Proposition 2:* Assume that $f \in \Gamma_0(\mathcal{H})$ is $\sigma$-strongly convex and $\beta$-smooth. Then $\operatorname{prox}_{\gamma f} - \frac{1}{1+\gamma\beta}\mathrm{Id}$ is $\frac{1}{\frac{1}{1+\gamma\sigma} - \frac{1}{1+\gamma\beta}}$-cocoercive if $\beta > \sigma$ and 0-Lipschitz if $\beta = \sigma$.

This result can be used to show the following contraction properties of the reflected proximal operator. A proof to this result, which is one of the main results of the paper, is found in Appendix B.

*Theorem 1:* Suppose that $f \in \Gamma_0(\mathcal{H})$ is $\sigma$-strongly convex and $\beta$-smooth. Then $R_{\gamma f}$ is $\max(\frac{\gamma\beta-1}{\gamma\beta+1}, \frac{1-\gamma\sigma}{\gamma\sigma+1})$-Lipschitz continuous.

This result lays the foundation for the linear convergence rate results in the following theorem, which is proven in Appendix C.

*Theorem 2:* Suppose that Assumption 1 holds. Then the generalized Douglas Rachford algorithm (5) converges linearly

towards a $\bar{z} \in \operatorname{fix}(R_{\gamma f} R_{\gamma g})$ at least with rate $|1 - \alpha| + \alpha \max\left(\frac{\gamma\beta-1}{\gamma\beta+1}, \frac{1-\gamma\sigma}{1+\gamma\sigma}\right)$, i.e.

$$\|z^{k+1} - \bar{z}\| \le \left(|1-\alpha| + \alpha \max\left(\tfrac{\gamma\beta-1}{\gamma\beta+1}, \tfrac{1-\gamma\sigma}{1+\gamma\sigma}\right)\right)^k \|z^0 - \bar{z}\|.$$

*Remark 2:* One interesting consequence of this results is that $\alpha > 1$ can be chosen in the Douglas-Rachford algorithm in (5), when solving problems that satisfy Assumption 1. To get a linear convergence, the rate factor in Theorem 2 should be less than 1, i.e.

$$|1-\alpha| + \alpha \max\left(\tfrac{\gamma\beta-1}{\gamma\beta+1}, \tfrac{1-\gamma\sigma}{1+\gamma\sigma}\right) < 1$$
$$\Rightarrow \qquad \alpha < \frac{2}{1 + \max\left(\frac{\gamma\beta-1}{\gamma\beta+1}, \frac{1-\gamma\sigma}{1+\gamma\sigma}\right)}. \tag{7}$$

This is an explicit upper bound for $\alpha$ which is greater than 1 since the max-expression is strictly less than 1. A similar finding is reported in [32], but no explicit expression for $\alpha$ is provided. To the authors' knowledge, our result is the first *explicit* bound on the relaxation factor $\alpha$ that allows it to be greater than 1.

We can choose the algorithm parameters $\gamma$ and $\alpha$ to optimize the bound on the convergence rate in Theorem 2. This is done in the following proposition.

*Proposition 3:* Suppose that Assumption 1 holds. Then the optimal parameters for the generalized Douglas-Rachford algorithm in (5) are given by $\alpha = 1$ and $\gamma = \frac{1}{\sqrt{\sigma\beta}}$. Further, the optimal rate is given by $\frac{\sqrt{\beta/\sigma}-1}{\sqrt{\beta/\sigma}+1}$.

*Proof.* It is straightforward to verify that $|1 - \alpha| + \alpha \max\left(\frac{\gamma\beta-1}{\gamma\beta+1}, \frac{1-\gamma\sigma}{1+\gamma\sigma}\right)$ is a decreasing function of $\alpha$ for $\alpha \le 1$ and increasing for $\alpha \ge 1$. Therefore the rate factor is optimized by $\alpha = 1$. The $\gamma$ parameter should be chosen to minimize the max-expression $\max\left(\frac{\gamma\beta-1}{\gamma\beta+1}, \frac{1-\gamma\sigma}{1+\gamma\sigma}\right)$. This is done by setting the arguments equal, which gives $\gamma = 1/\sqrt{\sigma\beta}$. Inserting these values into the rate factor expression gives $\frac{\sqrt{\beta/\sigma}-1}{\sqrt{\beta/\sigma}+1}$. $\square$

*Remark 3:* Note that in Proposition 3, $\alpha = 1$ is optimal. That is, the Peaceman-Rachford algorithm gives the best bound on the convergence rate under Assumption 1, even though the Peaceman-Rachford algorithm is not guaranteed to converge in the general case. The reason why we get convergence under the additional assumptions in Assumption 1 is that the one of the reflected proximal operators becomes contractive.

### B. Comparison to other methods

In this section, we discuss in what ways our result in Proposition 3 generalizes and/or improves on the previously known linear convergence rate results in [12], [14], [19], [34], [30] and the linear convergence rate [32] that appeared online during the submission procedure of this paper. Since Douglas-Rachford splitting and ADMM are equivalent in the case where $\mathcal{A} = \mathrm{Id}$ (that is, Douglas-Rachford is self-dual in the sense that it gives equivalent algorithms if applied to the primal and the dual when $\mathcal{A} = \mathrm{Id}$) we can compare Douglas-Rachford
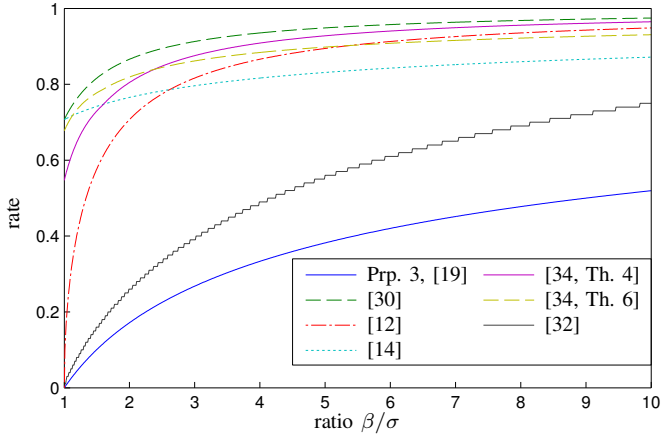
Fig. 1. Comparison between bounds on the linear convergence rate for Douglas-Rachford splitting provided in [12], [14], [19], [34], [30] and in Proposition 3 for different conditioning of the data. Proposition 3 provides the tightest bound on the rate for all ratios $\beta/\sigma$.

convergence rate results with ADMM convergence rate results by letting $\mathcal{A} = \mathrm{Id}$. The convergence rate result in [30] is provided in the case of solving monotone inclusion problems using Douglas-Rachford splitting, while the other [12], [14], [19], [34], [32] are for the convex optimization case treated here (which is a subclass of the monotone inclusion problem class). We will compare to the results in [12], [14], [19], [34], [32] that hold under Assumption 1. In the monotone inclusion problem case considered in [30], Assumption 1 corresponds to that one of the two maximal monotone operators is $\sigma$-strongly monotone and $\beta$-Lipschitz continuous.

In [30, Proposition 4, Remark 10], the linear convergence rate for Douglas-Rachford splitting when solving monotone inclusion problems with one operator being $\sigma$-strongly monotone and $\beta$-Lipschitz continuous is shown to be $\sqrt{1 - \frac{2\gamma\sigma}{(1+\gamma\beta)^2}}$. This rate is optimized by $\gamma = 1/\beta$, which gives a rate of $\sqrt{1 - \frac{\sigma}{2\beta}}$. This was generalized, in the setting of the operators being subdifferentials to proper, closed, and convex functions, to any $\alpha$ (not only $\alpha = \frac{1}{2}$ as in [30]) in [12, Theorem 6]. The rate in [12, Theorem 6] is $\sqrt{1 - \frac{4\alpha\gamma\sigma}{(1+\gamma\beta)^2}}$. The optimal parameters are $\alpha = 1$ and $\gamma = 1/\beta$, and the optimal rate bound becomes $\sqrt{1 - \frac{\sigma}{\beta}}$. Compared to Proposition 3, we see in Figure 1 that Proposition 3 gives a better bound than [12, Theorem 6] (and consequently also compared to [30]) for the plotted ratios $\beta/\sigma$.

The optimal convergence rate bound in [14, Corollary 3.6] is given by $\sqrt{1/(1 + 1/\sqrt{\beta/\sigma})}$, and the optimal parameter coincides with our choice in Proposition 3. Figure 1 shows that the convergence rate bound in Proposition 3 is better than the one provided in [14, Corollary 3.6] for all values of the ratio $\beta/\sigma$. Proposition 3 also generalizes [14, Corollary 3.6] since [14, Corollary 3.6] is stated in the Euclidean setting.

In [34], a new interpretation to Douglas-Rachford splitting is presented. The authors show that if $\gamma$ is small enough ($\gamma < 1/\beta$) and if $f$ is a quadratic function, then Douglas-Rachford splitting is equivalent to a gradient method applied to a function named the Douglas-Rachford envelope. The smoothness and strong convexity parameters for the envelope function can be computed from the corresponding values of the original function $f$. Convergence rate estimates follow from the convergence results for the gradient method. Since the Douglas-Rachford envelope is smooth, they also propose an accelerated Douglas-Rachford algorithm (under the same assumptions, i.e., $f$ quadratic and $\gamma < 1/\beta$). The convergence rate estimates of this also follow from the convergence rate estimates of fast gradient methods. In Figure 1, we plot the convergence rate estimates for the standard Douglas-Rachford algorithm [34, Theorem 4] and the fast Douglas-Rachford algorithm in [34, Theorem 6], both with parameter $\gamma = (\sqrt{2} - 1)/\beta$ (which is proposed in [34]). We see that Proposition 3 gives better rate bounds for all plotted values of the ratio $\beta/\sigma$. Proposition 3 is also more general in applicability.

The convergence rate estimates provided in [19] coincide with the results provided in Proposition 3. However, the generality of our analysis makes Proposition 3 applicable to a much wider class of problems than the results in [19]. Specifically, [19] considers ADMM applied to Euclidean quadratic problems with linear inequality constraints. We generalize these results to arbitrary real Hilbert spaces (even infinite dimensional), to both Douglas-Rachford splitting and ADMM, to general smooth and strongly convex functions $f$, and, perhaps most importantly, to any proper, closed, and convex function $g$.

Finally, we compare our rate bound to the rate bound in [32]. Figure 1 shows that our bound is indeed tighter for all plotted values of the ratio $\beta/\sigma$. As opposed to all the other rate bounds in this comparison, the rate bound in [32] is not explicit. Rather, a sweep over different rate bound factors is needed. For each guess, a small semi-definite program is solved to assess whether the algorithm is guaranteed to converge with that rate. The quantization level of this sweep is the cause of the steps in the plot in Figure 1.

We have shown that our results generalize the applicability and/or improve on the the linear convergence rate factor compared to existing results in the literature. In the next section, we show that the bounds provided in Theorem 2 and Proposition 3 are tight for the class of problems under consideration.

### C. Tightness of bounds

To show tightness of the linear convergence rate bounds in Theorem 2, we consider a problem of the form (4) with

$$f(x) = \sum_{i=1}^{K} \frac{\lambda_i}{2} \langle x, \phi_i \rangle^2, \tag{8}$$

$$g(x) \equiv 0, \tag{9}$$

where $\{\phi_i\}_{i=1}^{|\mathcal{H}|}$ is an orthonormal basis for $\mathcal{H}$, $|\mathcal{H}|$ is the dimension of the space $\mathcal{H}$ (possibly infinite), and $\lambda_i$ is either $\sigma > 0$ or $\beta > 0$, where $\beta \geq \sigma$. We denote the set of indices $i$ with $\lambda_i = \sigma$ by $\mathcal{I}_\sigma$ and the set of indices $i$ with $\lambda_i = \beta$ by $\mathcal{I}_\beta$ and require that $\mathcal{I}_\sigma \neq \emptyset$ and $\mathcal{I}_\beta \neq \emptyset$.

First, we show that $f$ in (8) is finite for all $x \in \mathcal{H}$. Obviously $f(x) \geq 0$ for all $x \in \mathcal{H}$. We also have for arbitrary $x \in \mathcal{H}$ that

$$f(x) = \sum_{i=1}^{|\mathcal{H}|} \frac{\lambda_i}{2}\langle x, \phi_i\rangle^2 \leq \frac{\beta}{2}\sum_{i=1}^{|\mathcal{H}|}\langle x, \phi_i\rangle^2 = \frac{\beta}{2}\|x\|^2 < \infty$$

where the last equality follows from Parseval's identity. Therefore $f$ and $g$ in (8) and (9) respectively have full domains. That $f$ is proper, closed, and convex holds trivially since $\lambda_i > 0$ for all $i$, and since $f$ is finite everywhere and differentiable. Next, we show that $f \in \Gamma_0(\mathcal{H})$ satisfies Assumption 1(ii), i.e., that $f$ is $\beta$-smooth and $\sigma$-strongly convex.

*Proposition 4:* The function $f$, as defined in (8) with $\lambda_i = \sigma$ for $i \in \mathcal{I}_\sigma$ and $\lambda_i = \beta$ for $i \in \mathcal{I}_\beta$, is $\sigma$-strongly convex and $\beta$-smooth.

*Proof.* We have that

$$\frac{\beta}{2}\|x\|^2 - f(x) = \sum_{i=1}^{|\mathcal{H}|}\frac{\beta - \lambda_i}{2}\langle x, \phi\rangle^2 = \sum_{i\in\mathcal{I}_\sigma}\frac{\beta-\sigma}{2}\langle x, \phi\rangle^2$$

which is convex since $\beta \geq \sigma$. Therefore $f$ is $\beta$-smooth according to Definition 6. We also have

$$f(x) - \frac{\sigma}{2}\|x\|^2 - = \sum_{i=1}^{|\mathcal{H}|}\frac{\lambda_i-\sigma}{2}\langle x, \phi\rangle^2 = \sum_{i\in\mathcal{I}_\beta}\frac{\beta-\sigma}{2}\langle x, \phi\rangle^2$$

which is convex since $\beta \geq \sigma$. Therefore $f$ is $\sigma$-strongly convex according to Definition 5. $\square$

To show that the provided example converges exactly with the rate given in Theorem 2, we need expressions for the proximal operators and reflected proximal operators of $f$ and $g$ in (8) and (9) respectively.

*Proposition 5:* The proximal operator of $f$ in (8) is

$$\text{prox}_{\gamma f}(y) = \sum_{i=1}^{|\mathcal{H}|}\frac{1}{1+\gamma\lambda_i}\langle y, \phi_i\rangle\phi_i \qquad (10)$$

and the reflected proximal operator is

$$R_{\gamma f}(y) = \sum_{i=1}^{|\mathcal{H}|}\frac{1-\gamma\lambda_i}{1+\gamma\lambda_1}\langle y, \phi_i\rangle\phi_i. \qquad (11)$$

*Proof.* We decompose $x = \sum_{i=1}^{|\mathcal{H}|}a_i\phi_i$ where $a_i = \langle x, \phi_i\rangle$ and $y = \sum_{i=1}^{|\mathcal{H}|}b_i\phi_i$ where $b_i = \langle y, \phi_i\rangle$. Then, for $\gamma > 0$, the proximal operator of $f$ is given by:

$$\text{prox}_{\gamma f}(y) = \arg\min_x\left\{\gamma\left(\sum_{i=1}^{|\mathcal{H}|}\frac{\lambda_i}{2}\langle\phi_i, x\rangle^2\right) + \frac{1}{2}\|x-y\|^2\right\}$$

$$= \arg\min_{x=\sum_{i=i}^{|\mathcal{H}|}a_i\phi_i}\left\{\left(\sum_{i=1}^{|\mathcal{H}|}\frac{\gamma\lambda_i}{2}a_i^2\right) + \frac{1}{2}\left\|\sum_{i=1}^{|\mathcal{H}|}(a_i-b_i)\phi_i\right\|^2\right\}$$

$$= \arg\min_{x=\sum_{i=i}^{|\mathcal{H}|}a_i\phi_i}\left\{\frac{1}{2}\sum_{i=1}^{|\mathcal{H}|}\left(\gamma\lambda_i a_i^2 + (a_i-b_i)^2\right)\right\}$$

$$= \sum_{i=1}^{|\mathcal{H}|}\arg\min_{a_i}\frac{1}{2}\left\{\gamma\lambda_i a_i^2 + (a_i-b_i)^2\right\}\phi_i$$

$$= \sum_{i=1}^{|\mathcal{H}|}\frac{1}{1+\gamma\lambda_i}b_i\phi_i = \sum_{i=1}^{|\mathcal{H}|}\frac{1}{1+\gamma\lambda_i}\langle y, \phi_i\rangle\phi_i.$$

The reflected resolvent for $\gamma > 0$ is given by:

$$R_{\gamma f}(y) = 2\text{prox}_{\gamma f}(y) - y$$

$$= 2\sum_{i=1}^{|\mathcal{H}|}\frac{1}{1+\gamma\lambda_i}b_i\phi_i - \sum_{i=1}^{|\mathcal{H}|}b_i\phi_i$$

$$= \sum_{i=1}^{|\mathcal{H}|}\frac{1-\gamma\lambda_i}{1+\gamma\lambda_i}b_i\phi_i = \sum_{i=1}^{|\mathcal{H}|}\frac{1-\gamma\lambda_i}{1+\gamma\lambda_i}\langle y, \phi_i\rangle\phi_i.$$

$\square$

The proximal and reflected proximal operators of $g \equiv 0$ are trivially given by $\text{prox}_{\gamma g} = R_{\gamma g} = \text{Id}$.

Next, these results are used to show that the convergence rate estimates in Theorem 2 are tight for the class of problems under consideration for many choices of algorithm parameters $\alpha$ and $\gamma$. Before we state this result, we need a help lemma.

*Lemma 1:* For $x > 0$, the function $\psi(x) := \frac{1-x}{1+x}$ satisfies $\psi(x) \leq -\psi(y)$ if and only if $y \geq 1/x$.

*Proof.* We have

$$\psi(x) = (1-x)/(1+x) \leq (y-1)(1+y) = -\psi(y)$$
$$\Leftrightarrow \qquad (1-x)(1+y) \leq (y-1)(1+x)$$
$$\Leftrightarrow \qquad 2 \leq 2xy.$$

$\square$

*Theorem 3:* The generalized Douglas-Rachford splitting algorithm (5) when applied to solve (4) with $f$ and $g$ in (8) and (9) respectively, converges exactly with the theoretical upper bound rate

$$|1-\alpha| + \alpha\max\left(\frac{1-\gamma\sigma}{1+\gamma\sigma}, \frac{\gamma\beta-1}{1+\gamma\beta}\right) \qquad (12)$$

in the following cases: (i) $\alpha \in (0,1]$ and $\gamma \in (0, \frac{1}{\sqrt{\beta\sigma}}]$, (ii) $\alpha \in [1, \frac{2}{1+\max\left(\frac{1-\gamma\sigma}{1+\gamma\sigma}, \frac{\gamma\beta-1}{1-\gamma\beta}\right)})$ and $\gamma \in [\frac{1}{\sqrt{\sigma\beta}}, \infty)$ for some algorithm initial condition $z^0$.

*Proof.* For algorithm initial condition $z^0 = \phi_i$ the Douglas-Rachford algorithm evolves according to

$$z^k = \left(1-\alpha + \alpha\frac{1-\gamma\lambda_i}{1+\gamma\lambda_i}\right)^k\phi_i$$

where $\lambda_i$ is either $\sigma$ or $\beta$ depending on if $i \in \mathcal{I}_\sigma$ or $i \in \mathcal{I}_\beta$. This follows immediately from the algorithm in (5), the expression of $R_{\gamma f}$ in Proposition 5, and since $R_{\gamma g} = \mathrm{Id}$. The convergence factor is exactly

$$\left| 1 - \alpha + \alpha \frac{1-\gamma\lambda_i}{1+\gamma\lambda_i} \right|. \tag{13}$$

We need to show that (13) is equal to (12) for the cases (i) and (ii). This holds if $1 - \alpha$ and $\alpha \frac{1-\gamma\lambda_i}{1+\gamma\lambda_i}$ have the same sign and if $\frac{1-\gamma\lambda_i}{1+\gamma\lambda_i} = \max(\frac{1-\gamma\sigma}{1+\gamma\sigma}, \frac{\gamma\beta-1}{1+\gamma\beta})$. First note that Lemma 1 implies that

$$\max(\tfrac{1-\gamma\sigma}{1+\gamma\sigma}, \tfrac{\gamma\beta-1}{1+\gamma\beta}) = \max(\psi(\gamma\sigma), -\psi(\gamma\beta))$$
$$= \begin{cases} \psi(\gamma\sigma) & \text{if } \gamma \leq \frac{1}{\sqrt{\beta\sigma}} \\ -\psi(\gamma\beta) & \text{if } \gamma \geq \frac{1}{\sqrt{\beta\sigma}} \end{cases}$$
$$= \begin{cases} \frac{1-\gamma\sigma}{1+\gamma\sigma} & \text{if } \gamma \leq \frac{1}{\sqrt{\beta\sigma}} \\ \frac{\gamma\beta-1}{1+\gamma\beta} & \text{if } \gamma \geq \frac{1}{\sqrt{\beta\sigma}} \end{cases} \tag{14}$$

where $\psi$ is defined in Lemma 1. This implies that

$$\max(\tfrac{1-\gamma\sigma}{1+\gamma\sigma}, \tfrac{\gamma\beta-1}{1+\gamma\beta}) \geq 0 \tag{15}$$

since $\frac{1-\gamma\sigma}{1+\gamma\sigma} \geq 0$ when $\gamma \leq \frac{1}{\sqrt{\beta\sigma}}$ and $\frac{\gamma\beta-1}{1+\gamma\beta} \geq 0$ when $\gamma \geq \frac{1}{\sqrt{\beta\sigma}}$ (since $\beta \geq \sigma$). Next, we use these observations to show the results for the two cases.

First, for case (i) with $\alpha \in (0, 1]$ and $\gamma \in (0, \frac{1}{\sqrt{\beta\sigma}}]$ we choose $\phi_i$ with $i \in \mathcal{I}_\sigma$ to get that the rate (13) in the example reduces to

$$\left| 1 - \alpha + \alpha\tfrac{1-\gamma\lambda_i}{1+\gamma\lambda_i} \right| = \left| 1 - \alpha + \alpha\tfrac{1-\gamma\sigma}{1+\gamma\sigma} \right|$$
$$= 1 - \alpha + \alpha\max(\tfrac{1-\gamma\sigma}{1+\gamma\sigma}, \tfrac{1-\gamma\beta}{1+\gamma\beta})$$
$$= |1 - \alpha| + \alpha\max(\tfrac{1-\gamma\sigma}{1+\gamma\sigma}, \tfrac{1-\gamma\beta}{1+\gamma\beta})$$

where the second equality follows from (14) since $\gamma \in (0, \frac{1}{\sqrt{\beta\sigma}}]$, from (15), and since $\alpha \in (0, 1]$. That is, (13) coincides with (12). For the second case with $\alpha \in [1, \frac{2}{1+\max\left(\frac{1-\gamma\sigma}{1+\gamma\sigma}, \frac{\gamma\beta-1}{1+\gamma\beta}\right)})$ and $\gamma \in [\frac{1}{\sqrt{\sigma\beta}}, \infty)$ we choose $\phi_i$ with $i \in \mathcal{I}_\beta$ to get that the rate (13) in the example reduces to

$$\left| 1 - \alpha + \alpha\tfrac{1-\gamma\lambda_i}{1+\gamma\lambda_i} \right| = \left| 1 - \alpha + \alpha\tfrac{1-\gamma\beta}{1+\gamma\beta} \right|$$
$$= \alpha - 1 + \alpha\max(\tfrac{1-\gamma\sigma}{1+\gamma\sigma}, \tfrac{1-\gamma\beta}{1+\gamma\beta})$$
$$= |1 - \alpha| + \alpha\max(\tfrac{1-\gamma\sigma}{1+\gamma\sigma}, \tfrac{1-\gamma\beta}{1+\gamma\beta})$$

where the second equality follows from (14) since $\gamma \in [\frac{1}{\sqrt{\beta\sigma}}, \infty)$, from (15), and since $\alpha \geq 1$. That is, (13) coincides with (12) also in this second case. This concludes the proof. $\square$

The convergence rate for the example given by $f$ and $g$ in (8) and (9) respectively coincides with the upper bound on the convergence rate in Theorem 2. The bound in Theorem 2 is therefore tight for the class of problems under consideration and for the combination of algorithm parameters specified in Theorem 3. Especially, the convergence rate bound for the optimal parameters given by $\alpha = 1$ and $\gamma = \frac{1}{\sqrt{\beta\sigma}}$ is tight.

Besides generalizing and/or improving on existing results from the literature, the results in Proposition 3 can guide us in choosing a space on which to perform the Douglas-Rachford algorithm when solving finite-dimensional problems. By selecting the space appropriately, this can significantly improve the convergence properties of the algorithm, both in theory and in practice. This is the topic of the following section.

### D. Metric selection

In this section, we consider finite-dimensional composite convex optimization problems of the form (4), where $f$ and $g$ satisfy:

*Assumption 2:*

(i) The function $f \in \Gamma_0(\mathbb{H}_M)$ is 1-strongly convex if defined on $\mathbb{H}_H$ and 1-smooth if defined on $\mathbb{H}_L$.

(ii) The function $g \in \Gamma_0(\mathbb{H}_M)$.

Examples of functions $f$ that satisfy Assumption 2(i) are piece-wise quadratic functions with Hessians $Q_i$ that are differentiable on the boundary between the regions. The matrix $H$ satisfies $0 \prec H \preceq Q_i$ for all $i$ and $L$ satisfies $L \succeq Q_i$ for all $i$. Obviously, in the general case we have $L \succeq H$ and for a standard quadratic function with Hessian $H$, we have $L = H$. Depending on which space $\mathbb{H}_M$ the functions $f$ and $g$ are defined, the algorithm will have different convergence properties. Proposition 3 suggests that to optimize the convergence rate bound, we should select a space $\mathbb{H}_M$ on which the ratio $\beta/\sigma$ is as small as possible, i.e., on which the conditioning of the function $f$ is as good as possible. Next, we present a result that shows how $\beta$ and $\sigma$ vary with $M$ in $\mathcal{H}_M$.

*Proposition 6:* Suppose that $f \in \Gamma_0(\mathbb{H}_M)$ satisfies Assumption 2(i) and that $M = (D^T D)^{-1}$. Then the strong convexity modulus $\sigma_M(f)$ and the smoothness parameter $\beta_M(f)$ are given by

$$\beta_M(f) = \lambda_{\max}(DLD^T)$$
$$\sigma_M(f) = \lambda_{\min}(DHD^T).$$

*Proof.* Denote by $\nabla_M f$ the gradient of $f$ when defined on $\mathbb{H}_M$ and $\nabla_2 f$ the gradient of $f$ when defined on $\mathbf{R}^n$. Then $\nabla_M f = M^{-1}\nabla_2 f$ since

$$f(x) \geq f(y) + \langle \nabla_M f(x), x - y \rangle_M$$
$$\Leftrightarrow \quad f(x) \geq f(y) + \langle M\nabla_M f(x), x - y \rangle_2$$

where $M\nabla_M f = \nabla_2 f$. Therefore

$$\langle \nabla f(y), x - y \rangle_{M_1} = \langle \nabla f(y), x - y \rangle_{M_2} \tag{16}$$

for any $M_1, M_2 \succ 0$. Further, by letting $M_2 = (D_2^T D_2)^{-1}$, we have

$$\|x\|_{M_1}^2 \geq \lambda_{\min}(D_2 M_1 D_2^T)\|x\|_{M_2}^2 \tag{17}$$
$$\|x\|_{M_1}^2 \leq \lambda_{\max}(D_2 M_1 D_2^T)\|x\|_{M_2}^2. \tag{18}$$

The first inequality holds since

$$\|x\|_{M_1}^2 \geq \lambda_{\min}(D_2 M_1 D_2^T)\|x\|_{M_2}^2$$
$$\Leftrightarrow \quad \|D_2^T x\|_{M_1}^2 \geq \lambda_{\min}(D_2 M_1 D_2^T)\|D_2^T x\|_{M_2}^2$$
$$\Leftrightarrow \quad \|x\|_{D_2 M_1 D_2^T}^2 \geq \lambda_{\min}(D_2 M_1 D_2^T)\|x\|_2^2$$
$$\Leftrightarrow \quad D_2 M_1 D_2^T \succeq \lambda_{\min}(D_2 M_1 D_2^T)I.$$

The inequality (18) is proven similarly. Since $f$ is 1-strongly convex if defined on $\mathbb{H}_H$, the definition of strong convexity for differentiable functions (2) gives

$$
\begin{aligned}
f(x) &\geq f(y) + \langle \nabla f(y), x - y \rangle_H + \tfrac{1}{2}\|x-y\|_H^2 \\
&\geq f(y) + \langle \nabla f(y), x - y \rangle_M + \tfrac{\lambda_{\min}(DHD^T)}{2}\|x-y\|_M^2
\end{aligned}
$$

where (16) and (17) are used in the second inequality. Therefore $f$, when defined on $\mathbb{H}_M$ with $M = (D^TD)^{-1}$, is $\lambda_{\min}(DHD^T)$-strongly convex. Similarly, since $f$ is 1-smooth if defined on $\mathbb{H}_L$, the definition of smoothness (3) gives

$$
\begin{aligned}
f(x) &\leq f(y) + \langle \nabla f(y), x - y \rangle_H + \tfrac{1}{2}\|x-y\|_H^2 \\
&\leq f(y) + \langle \nabla f(y), x - y \rangle_M + \tfrac{\lambda_{\max}(DLD^T)}{2}\|x-y\|_M^2
\end{aligned}
$$

where (16) and (18) are used in the second inequality. Therefore $f$, when defined on $\mathbb{H}_M$ with $M = (D^TD)^{-1}$, is $\lambda_{\max}(DLD^T)$-smooth. This concludes the proof. $\square$

This result indicates that, to optimize the rate in Proposition 3, we should select a metric $M = (D^TD)^{-1}$ that solves

$$
\text{minimize}\ \frac{\beta_M(f)}{\sigma_M(f)} \quad = \quad \text{minimize}\ \frac{\lambda_{\max}(DLD^T)}{\lambda_{\min}(DHD^T)}. \quad (19)
$$

In accordance with Proposition 3, the algorithm parameter $\gamma$ should be chosen as $\gamma = \frac{1}{\sqrt{\lambda_{\max}(DLD^T)\lambda_{\min}(DHD^T)}}$. To select a metric according to (19) can significantly improve the convergence rate bound compared to applying the algorithm in the standard Euclidean space. This is suggested by the following example. Suppose that we minimize a problem with a quadratic function $f$ with Hessian $H$, and that the generalized Douglas-Rachford algorithm is run on the Euclidean space with $M = D = I$. Then Proposition 3 guarantees the rate $\frac{\sqrt{\lambda_{\max}(H)/\lambda_{\min}(H)}-1}{\sqrt{\lambda_{\max}(H)/\lambda_{\min}(H)}+1}$. If we instead apply the generalized Douglas-Rachford algorithm on $\mathbb{H}_M$ with $M = H$ and $D = H^{-1/2}$ (which optimizes (19)), we get rate $\frac{\sqrt{\lambda_{\max}(H^{-1/2}HH^{-1/2})/\lambda_{\min}(H^{-1/2}HH^{-1/2})}-1}{\sqrt{\lambda_{\max}(H^{-1/2}HH^{-1/2})/\lambda_{\min}(H^{-1/2}HH^{-1/2})}+1} = 0$. That is, the algorithm converges in one iteration. The more ill-conditioned the original problem is, the more we improve the rate bound by selecting a better metric for the problem. However, often the functions $f$ and/or $g$ are separable down to the component. In such cases, choosing a non-diagonal $M$ would significantly increase the computational complexity associated with evaluating the prox-operator. Therefore, to get an efficient algorithm both in terms of convergence rate and in terms of complexity within each iteration, the metric matrix $M = (D^TD)^{-1}$ should be chosen to minimize (19), subject to $M$ being diagonal. In [22, Section 6] methods to minimize (19) exactly (it is shown that (19) can be formulated as a convex semi-definite program) as well as computationally cheap methods to reduce the ratio in (19) are presented.

*Remark 4:* Note that the metric selection does not change the problem to be minimized. It only changes the way distances are measured in the algorithm. So the same optimal point (up to numerical accuracy) is returned from the algorithm independent of on which space the algorithm is run.

*Remark 5:* It can be shown that to apply the Douglas-Rachford algorithm on the space $\mathbb{H}_M$ is equivalent to apply Douglas-Rachford splitting on the Euclidean space $\mathbf{R}^n$ to the preconditioned problem

$$
\text{minimize} \quad f_D(x) + g_D(x) \quad (20)
$$

where $M = (D^TD)^{-1}$ and

$$
\begin{aligned}
f_D(x) &:= f(D^Tx) \\
g_D(x) &:= g(D^Tx)
\end{aligned}
$$

and $f, f_D, g, g_D \in \Gamma_0(\mathbf{R}^n)$. Showing this equivalence is omitted for space considerations, but it follows readily by comparing the results form the respective proximal operators.

## IV. ADMM

In this section, besides $\mathcal{H}$ being a real Hilbert space, also $\mathcal{K}$ denotes a real Hilbert space. Here, we consider solving problems of the form

$$
\text{minimize} \quad f(x) + g(\mathcal{A}x) \quad (21)
$$

that satisfy the following assumptions:

*Assumption 3:*
 (i) The function $f \in \Gamma_0(\mathcal{H})$ is $\beta$-smooth and $\sigma$-strongly convex.
 (ii) The function $g \in \Gamma_0(\mathcal{K})$.
 (iii) The bounded linear operator $\mathcal{A} : \mathcal{H} \to \mathcal{K}$ is surjective.

The assumption of $\mathcal{A}$ being a surjective bounded linear operator reduces to $\mathcal{A}$ being a real matrix with full row rank in the Euclidean case. Problems of the form (21) cannot be directly efficiently solved using generalized Douglas-Rachford splitting. Therefore, we instead solve the (negative) Fenchel dual problem, which is given by (see [2, Definition 15.19])

$$
\text{minimize} \quad d(\mu) + g^*(\mu) \quad (22)
$$

where $g^* \in \Gamma_0(\mathcal{K})$ and $d \in \Gamma_0(\mathcal{K})$ is

$$
d(\mu) := f^*(-\mathcal{A}^*\mu) \quad (23)
$$

where $\mathcal{A}^* : \mathcal{K} \to \mathcal{H}$ is the adjoint operator of $\mathcal{A}$, defined as the unique operator that satisfies $\langle \mathcal{A}x, \mu \rangle = \langle x, \mathcal{A}^*\mu \rangle$ for all $x \in \mathcal{H}$ and $\mu \in \mathcal{K}$. Applying Douglas-Rachford splitting (i.e. generalized Douglas-Rachford splitting with $\alpha = \frac{1}{2}$) to the dual is well known to be equivalent to applying ADMM to the primal, see [17], [16]. To apply generalized Douglas-Rachford splitting to the dual for other choices of $\alpha$ is known as ADMM with over-relaxation for $\alpha \in (\frac{1}{2}, 1]$ and ADMM with under-relaxation for $\alpha \in (0, \frac{1}{2})$ (here we show that also $\alpha > 1$ is possible under Assumption 3). Therefore, the results we obtain in this section applies to relaxed ADMM.

### A. Linear convergence

To optimize the bound on the linear convergence rate in Proposition 3 when applied to solve the dual problem (22), we need to quantify the strong convexity and smoothness parameters for $d$. This is done in the following proposition.

*Proposition 7:* Suppose that Assumption 3 holds. Then $d \in \Gamma_0(\mathcal{K})$ is $\frac{\|\mathcal{A}^*\|^2}{\sigma}$-smooth and $\frac{\theta^2}{\beta}$-strongly convex, where $\theta > 0$ always exists and satisfies $\|\mathcal{A}^*\mu\| \geq \theta\|\mu\|$ for all $\mu \in \mathcal{K}$.

*Proof.* Since $f$ is $\sigma$-strongly convex, [2, Theorem 18.15] gives that $f^*$ is $\frac{1}{\sigma}$-smooth and that $\nabla f^*$ is $\frac{1}{\sigma}$-Lipschitz continuous. Therefore, $\nabla d$ satisfies

$$\begin{aligned}
\|\nabla d(\mu) - \nabla d(\nu)\| &= \|\mathcal{A}\nabla f^*(-\mathcal{A}^*\mu) - \mathcal{A}\nabla f^*(-\mathcal{A}^*\nu)\| \\
&\leq \frac{\|\mathcal{A}\|}{\sigma}\|\mathcal{A}^*(\mu - \nu)\| \\
&\leq \frac{\|\mathcal{A}^*\|^2}{\sigma}\|\mu - \nu\|
\end{aligned}$$

since $\|\mathcal{A}\| = \|\mathcal{A}^*\|$. This is equivalent to that $d$ is $\frac{\|\mathcal{A}^*\|^2}{\sigma}$-smooth, see [2, Theorem 18.15].

Next, we show the strong convexity result. The property that $f$ is $\beta$-smooth implies through [2, Theorem 18.15] that $f^*$ is $\frac{1}{\beta}$-strongly convex and that $\nabla f^*$ is $\frac{1}{\beta}$-strongly monotone. This implies that $\nabla d$ satisfies

$$\begin{aligned}
&\langle \nabla d(\mu) - \nabla d(\nu), \mu - \nu \rangle \\
&= \langle -\mathcal{A}(\nabla f^*(-\mathcal{A}^*\mu) - \nabla f^*(-\mathcal{A}^*\nu)), \mu - \nu \rangle \\
&= \langle \nabla f^*(-\mathcal{A}^*\mu) - \nabla f^*(-\mathcal{A}^*\nu), -\mathcal{A}^*\mu + \mathcal{A}^*\nu \rangle \\
&\geq \tfrac{1}{\beta}\|\mathcal{A}^*(\mu - \nu)\|^2 \geq \tfrac{\theta^2}{\beta}\|\mu - \nu\|^2.
\end{aligned}$$

This, by [2, Theorem 18.15], is equivalent to $d$ being $\frac{\theta^2}{\beta}$-strongly convex. That $\theta > 0$ follows from [2, Fact 2.18 and Fact 2.19]. Specifically, [2, Fact 2.18] says that $\ker\mathcal{A}^* = (\mathrm{ran}\mathcal{A})^\perp = \emptyset$, since $A$ is surjective. Since $\mathrm{ran}\mathcal{A} = \mathcal{K}$ (again by surjectivity), it is closed. Then [2, Fact 2.19] states that there exists $\theta > 0$ such that $\|\mathcal{A}^*\mu\| \geq \theta\|\mu\|$ for all $\mu \in (\ker\mathcal{A}^*)^\perp = (\emptyset)^\perp = \mathcal{K}$. This concludes the proof. $\square$

This result gives us the following immediate corollary.

*Corollary 1:* Suppose that Assumption 3 holds and that generalized Douglas-Rachford is applied to solve the dual problem (22) (or equivalently ADMM is applied to solve the primal (21)). Then the algorithm converges at least with the rate $|1 - \alpha| + \alpha \max\left(\frac{\gamma\hat{\beta}-1}{1+\gamma\hat{\beta}}, \frac{1-\gamma\hat{\sigma}}{1+\gamma\hat{\sigma}}\right)$ where $\hat{\beta} = \frac{\|\mathcal{A}^*\|^2}{\sigma}$ and $\hat{\sigma} = \frac{\theta^2}{\beta}$. Further, the algorithm parameters $\gamma$ and $\alpha$ that optimize the rate bound are $\alpha = 1$ and $\gamma = \frac{1}{\sqrt{\hat{\beta}\hat{\sigma}}} = \frac{\sqrt{\beta\sigma}}{\sqrt{\|\mathcal{A}^*\|^2\theta^2}}$. The optimized linear convergence rate bound factor is $\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}$, where $\kappa = \frac{\hat{\beta}}{\hat{\sigma}} = \frac{\|\mathcal{A}^*\|^2\beta}{\theta^2\sigma}$.

*Proof.* This follows directly from Propositions 7 and 3 and Theorem 2. $\square$

*Remark 6:* Also in the dual case, the $\alpha$-parameter can be chosen greater than one. The upper bound on $\alpha$ is given in (7) with $\beta$ and $\sigma$ replaced by $\hat{\beta} = \frac{\|\mathcal{A}^*\|^2}{\sigma}$ and $\hat{\sigma} = \frac{\theta^2}{\beta}$ respectively.

The convergence rate bounds in Corollary 1 depend both on the conditioning of the function $f$ and the conditioning of the linear operator $\mathcal{A}^*$. The better the conditioning, the faster the rate. However, some of the parameters might be hard to compute or estimate, especially $\theta$. In the following section, we show how to compute this in finite-dimensional Hilbert spaces $\mathbb{H}_K$. We also show how to select the space $\mathbb{H}_K$ (i.e., select matrix $K$) to optimize the bound on the convergence rate.

### B. Metric selection

In this section, we still consider problems of the form (21) and we suppose that the following assumptions hold:

*Assumption 4:*
(i) The function $f \in \Gamma_0(\mathbb{H}_M)$ is 1-strongly convex if defined on $\mathbb{H}_H$ and 1-smooth if defined on $\mathbb{H}_L$.
(ii) The function $g \in \Gamma_0(\mathbb{H}_K)$.
(iii) The bounded linear operator $\mathcal{A} : \mathbb{H}_H \to \mathbb{H}_K$ is surjective.

Items (i) and (ii) are the same as in Assumption 2, and the assumption on the bounded linear operator is added due to the more general problem formulation treated here.

Also here, we solve (21) by applying Douglas-Rachford splitting on the dual problem (22) (or equivalently by applying ADMM directly on the primal (21)). In this case, we can select the space $\mathbb{H}_K$ on which to define the dual problem and apply the algorithm. To aid in the selection of such a space, we show in the following proposition how the strong convexity modulus and smoothness constant of $d \in \Gamma_0(\mathbb{H}_K)$ depend on the space on which $d$ is defined.

*Proposition 8:* Suppose that Assumption 4 holds, that $A \in \mathbf{R}^{m \times n}$ satisfies $Ax = \mathcal{A}x$ for all $x$, and that $K = E^T E$. Then $d \in \mathbb{H}_K$ is $\|EAH^{-1}A^T E^T\|$-smooth and $\lambda_{\min}(EAL^{-1}A^T E^T)$-strongly convex, where $\lambda_{\min}(EAL^{-1}A^T E^T) > 0$.

*Proof.* First, we relate $\mathcal{A}^* : \mathbb{H}_K \to \mathbb{H}_M$ to $A$, $M$, and $K$. We have

$$\begin{aligned}
\langle \mathcal{A}x, \mu \rangle_K &= \langle Ax, K\mu \rangle_2 = \langle x, A^T K\mu \rangle_2 = \langle x, M^{-1}A^T K\mu \rangle_M \\
&= \langle M^{-1}A^T K\mu, x \rangle_M = \langle \mathcal{A}^*\mu, x \rangle_M.
\end{aligned}$$

Thus, $\mathcal{A}^*\mu = M^{-1}A^T K\mu$ for all $\mu \in \mathbb{H}_K$.

Next, we show that the space $\mathbb{H}_M$ on which $f$ and $f^*$ are defined does not influence the shape of $d$. We denote by $f_H$, $f_L$, and $f_e$ the function $f$ defined on $\mathbb{H}_H$, $\mathbb{H}_L$ and $\mathbf{R}^n$ respectively and by $\mathcal{A}_H^* : \mathbb{H}_K \to \mathbb{H}_H$, $\mathcal{A}_L^* : \mathbb{H}_K \to \mathbb{H}_L$, and $A^T : \mathbf{R}^m \to \mathbf{R}^n$ the operator $\mathcal{A}^*$ defined on different spaces. Further, let $d_H := f_H^* \circ -\mathcal{A}_H^*$, $d_L := f_L^* \circ -\mathcal{A}_L^*$, and $d_e := f_e^* \circ -A^T$. By these definitions both $d_L$ and $d_H$ are defined on $\mathbb{H}_K$, while $d_e$ is defined on $\mathbf{R}^m$. Next we show that $d_L$ and $d_H$ are identical for any $\mu$:

$$\begin{aligned}
d_H(\mu) &= f_H^*(-\mathcal{A}_H^*\mu) = \sup_x \left\{ \langle -\mathcal{A}_H^*\mu, x \rangle_H - f_H(x) \right\} \\
&= \sup_x \left\{ \langle -HH^{-1}A^T K\mu, x \rangle_2 - f_e(x) \right\} \\
&= \sup_x \left\{ \langle -LL^{-1}A^T K\mu, x \rangle_2 - f_e(x) \right\} \\
&= \sup_x \left\{ \langle -\mathcal{A}_L^*\mu, x \rangle_L - f_L(x) \right\} = d_L(\mu)
\end{aligned}$$

where $\mathcal{A}_H^*\mu = H^{-1}A^T K\mu$ is used. This implies that we can show properties of $d \in \mathbb{H}_K$ by defining $f$ on any space $\mathbb{H}_M$. Thus, Proposition 7 gives that 1-strong convexity of $f$ when defined on $\mathbb{H}_H$ implies $\|\mathcal{A}^*\|^2$-smoothness of $d$, where

$$\begin{aligned}
\|\mathcal{A}^*\| &= \sup_\mu \left\{ \|\mathcal{A}^*\mu\| \mid \|\mu\| \leq 1 \right\} \\
&= \sup_\mu \left\{ \|H^{-1}A^T K\mu\|_H \mid \|\mu\|_K \leq 1 \right\} \\
&= \sup_\mu \left\{ \|H^{-1/2}A^T E^T E\mu\|_2 \mid \|E\mu\|_2 \leq 1 \right\} \\
&= \sup_\nu \left\{ \|H^{-1/2}A^T E^T\nu\|_2 \mid \|\nu\|_2 \leq 1 \right\} \\
&= \|H^{-1/2}A^T E^T\|_2.
\end{aligned}$$

Taking the square gives the smoothness claim. To show the strong-convexity claim, we use that 1-smoothness of $f$ when defined on $\mathbb{H}_L$ implies $\theta^2$-strong convexity of $d$ where $\theta > 0$ satisfies $\|\mathcal{A}^*\mu\| \geq \theta\|\mu\|$ for all $\mu \in \mathbb{H}_K$, see Proposition 7. Such a $\theta$ is given by

$$
\begin{aligned}
\|\mathcal{A}^*\mu\|_L^2 = \|L^{-1}A^T K\mu\|_L^2 &= \|L^{-1/2}A^T E^T(E\mu)\|_2^2 \\
&= \|E\mu\|_{EAL^{-1}A^T E^T}^2 \\
&\geq \lambda_{\min}(EAL^{-1}A^T E^T)\|E\mu\|_2^2 \\
&= \lambda_{\min}(EAL^{-1}A^T E^T)\|\mu\|_K^2.
\end{aligned}
$$

The smallest eigenvalue $\lambda_{\min}(EAL^{-1}A^T E^T) > 0$ since $A$ is surjective, i.e. has full row rank, and $E$ and $L$ are positive definite matrices. This concludes the proof. $\qquad\square$

This result shows how the smoothness constant and strong convexity modulus of $d \in \Gamma_0(\mathbb{H}_K)$ change with the space $\mathbb{H}_K$ on which $d$ is defined. Combining this with Proposition 3, we get that the bound on the convergence rate for Douglas-Rachford splitting applied to the dual problem (22), or equivalently ADMM applied to the primal (21), is optimized by choosing $K = E^T E$ where $E$ solves

$$
\text{minimize} \quad \frac{\lambda_{\max}(EAH^{-1}A^T E^T)}{\lambda_{\min}(EAL^{-1}A^T E^T)} \tag{24}
$$

and by choosing $\gamma = \frac{1}{\sqrt{\lambda_{\max}(EAH^{-1}A^T E^T)\lambda_{\min}(EAL^{-1}A^T E^T)}}$. As for Douglas-Rachford splitting applied to the primal problem, using a non-diagonal $K$ usually gives prohibitively expensive prox-evaluations. Therefore, we propose to select a diagonal $K = E^T E$ that minimizes (24). The reader is referred to [22, Section 6] for different methods to achieve this exactly and approximately.

*Remark 7:* Also in this dual case, the Douglas-Rachford algorithm applied on the space $\mathbb{H}_K$ is equivalent to solving a preconditioned problem on the Euclidean space, namely:

$$
\text{minimize} \quad d_E(\nu) + g_E^*(\nu) \tag{25}
$$

where $K = E^T E$,

$$
\begin{aligned}
d_E(\nu) &:= d(E^T \nu) \\
g_E^*(\nu) &:= g^*(E^T \nu),
\end{aligned}
$$

and $d, d_E, g^*, g_E^* \in \Gamma_0(\mathbf{R}^m)$. Note that the matrix that defines the space satisfies $K = E^T E$, where $E$ is the preconditioner matrix, while in the primal case the corresponding relation is $M = (D^T D)^{-1}$. The reason is that in the dual formulation, the shape of $d$ and $g^*$ change depending on which space they are defined. This is not the case in the primal formulation.

Relating this to ADMM, it can also be shown that solving the dual problem on space $\mathbb{H}_K$ using Douglas-Rachford splitting is equivalent to solving the preconditioned problem

$$
\begin{aligned}
\text{minimize} \quad & f(x) + g(y) \\
\text{subject to} \quad & EAx = Ey.
\end{aligned} \tag{26}
$$

using ADMM. Details are omitted for space considerations.

## V. HEURISTIC METRIC SELECTION

In this section, we discuss metric and parameter selection when some of the assumptions needed to have linear convergence are not met. We focus here on quadratic problems of the form

$$
\text{minimize} \quad \underbrace{\tfrac{1}{2}x^T Q x + q^T x + \hat{f}(x)}_{f(x)} + g(Ax) \tag{27}
$$

where $Q \in \mathbf{R}^{n \times n}$ is positive semi-definite, $q \in \mathbf{R}^n$, $\hat{f} \in \Gamma_0(\mathbf{R}^n)$, $g \in \Gamma_0(\mathbf{R}^n)$ and $A \in \mathbf{R}^{m \times n}$. One set of assumptions that guarantee linear convergence for Douglas-Rachford splitting applied to the primal or the dual is that $Q$ is positive definite, $\hat{f} \equiv 0$, and that $A$ has full row rank. Here, we consider situations in which some of these assumptions are not met. Specifically, we consider situations where (some of) the following items violate the linear convergence assumptions:

(i) $Q$ is not positive definite, but positive semi-definite.
(ii) $\hat{f} \not\equiv 0$, but instead the indicator function of a convex constraint set (or some other non-smooth function without curvature).
(iii) $A$ does not have full row rank.

In the first case, we loose strong convexity in the primal formulation and smoothness in the dual formulation. In the second case, we loose smoothness in the primal formulation and strong convexity in the dual formulation. The third case is not applicable to the primal case (since $A = I$ there), but in the dual formulation this results in loss of strong convexity.

We first discuss the primal formulation and assume that the assumptions that give linear convergence are violated using both (i) and (ii). Then, we have quadratic curvature only in the range space of $Q$. In the null space of $Q$, the function $f$ is governed by the function $\hat{f}$ (which is either 0 of $\infty$ if it is the indicator function of a convex constraint set). Therefore, we propose to select a diagonal metric $M = (D^T D)^{-1}$ that optimizes the conditioning on the range space of $Q$, i.e., that solves

$$
\text{minimize} \quad \frac{\lambda_{\max}(D^T Q D)}{\lambda_{\min>0}(D^T Q D)}
$$

where $\lambda_{\min>0}$ denotes the smallest non-zero eigenvalue. Also, we propose to select the $\gamma$-parameter to reflect the curvature on the range space of $Q$, i.e., $\gamma = \frac{1}{\sqrt{\lambda_{\max}(D^T Q D)\lambda_{\min>0}(D^T Q D)}}$.

For the dual case, i.e., the ADMM case, we propose to select the metric as if $\hat{f} \equiv 0$ (which it is if the assumptions to get linear convergence are not violated by point (ii)). To do this, we define the quadratic part of $f$ in (27) to be $f_{\mathrm{pc}}(x) := \tfrac{1}{2}x^T Q x + q^T x$ and introduce the function $d_{\mathrm{pc}} = f_{\mathrm{pc}}^* \circ -A^T$. The heuristic metric selection will be based on this function. The function $f_{\mathrm{pc}}$ is given by

$$
\begin{aligned}
f_{\mathrm{pc}}^*(y) &= \sup_x \{\langle y, x \rangle - f_{\mathrm{pc}}(x)\} \\
&= \begin{cases} \tfrac{1}{2}(y-q)^T Q^\dagger(y-q) & \text{if } (y-q) \in \mathcal{R}(Q) \\ \infty & \text{else} \end{cases}
\end{aligned}
$$

where $Q^\dagger$ is the pseudo-inverse of $Q$ and $\mathcal{R}$ denotes the range space. This gives

$$d_{\mathrm{pc}}(\mu) = \begin{cases} \frac{1}{2}(A^T\mu + q)^T Q^\dagger (A^T\mu - q) & \text{if } (A^T\mu + q) \in \mathcal{R}(Q) \\ \infty & \text{else} \end{cases}$$

The quadratic part of the approximated dual $d_{\mathrm{pc}}$ is given by $AQ^\dagger A^T$, and is defined on a sub-space only (if $Q$ is not positive definite). As in the primal case, we propose to select a diagonal metric $K = E^T E$ such that the quadratic part of the, in some cases approximate, dual function is well conditioned on its domain. That is, we propose to select a metric $K = E^T E$ such that the pseudo condition number of $AQ^\dagger A^T$ is minimized. This is computed by solving

$$\text{minimize} \quad \frac{\lambda_{\max}(EAQ^\dagger A^T E^T)}{\lambda_{\min > 0}(EAQ^\dagger A^T E^T)}.$$

This reduces to the optimal metric choice in the case where linear convergence is achieved, i.e., where none of items (i), (ii), or (iii) are met, and can be used as a heuristic when any of the points (i), (ii), and/or (iii) violate the assumptions needed to get linear convergence. The $\gamma$-parameter is also chosen in accordance with the above reasoning and Corollary 1 as $\gamma = \frac{1}{\sqrt{\lambda_{\max}(EAQ^\dagger A^T E^T)\lambda_{\min > 0}(EAQ^\dagger A^T E^T)}}$.

In the particular case where $\hat{f}$ in (27) is the indicator function for an affine subspace, i.e., when $\hat{f} = I_{Bx=b}$. Then $d$ can be written as

$$d(\mu) = \tfrac{1}{2}\mu^T AP_{11}A^T\mu + \xi^T\mu + \chi$$

where $\xi \in \mathbf{R}^n$, $\chi \in \mathbf{R}$, and

$$\begin{bmatrix} Q & B^T \\ B & 0 \end{bmatrix}^{-1} = \begin{bmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{bmatrix}. \tag{28}$$

Then we can choose metric by minimizing the pseudo condition number of $AP_{11}A^T$, which is the Hessian of $d$, and select $\gamma$ as $\gamma = \frac{1}{\sqrt{\lambda_{\max}(EAP_{11}A^T E^T)\lambda_{\min > 0}(EAP_{11}A^T E^T)}}$.

Minimization of the pseudo condition number $\lambda_{\max}/\lambda_{\min > 0}$ can be posed as a convex optimization problem and be solved exactly, see [22, Section 6]. Also computationally cheap heuristics to select $E$ that reduce the pseudo condition number can be found there.

## VI. NUMERICAL EXAMPLE

In this section, we evaluate the metric and parameter selection method by applying ADMM to the (small-scale) aircraft control problem from [29], [5]. As in [5], the continuous time model from [29] is sampled using zero-order hold every 0.05 s. The system has four states $x = (x_1, x_2, x_3, x_4)$, two outputs $y = (y_1, y_2)$, two inputs $u = (u_1, u_2)$, and obeys the following dynamics

$$x^+ = \begin{bmatrix} 0.999 & -3.008 & -0.113 & -1.608 \\ -0.000 & 0.986 & 0.048 & 0.000 \\ 0.000 & 2.083 & 1.009 & -0.000 \\ 0.000 & 0.053 & 0.050 & 1.000 \end{bmatrix} x + \begin{bmatrix} -0.080 & -0.635 \\ -0.029 & -0.014 \\ -0.868 & -0.092 \\ -0.022 & -0.002 \end{bmatrix} u,$$

$$y = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} x$$

where $x^+$ denotes the state in the next time step. The system is unstable, the magnitude of the largest eigenvalue of the dynamics matrix is 1.313. The outputs are the attack and pitch
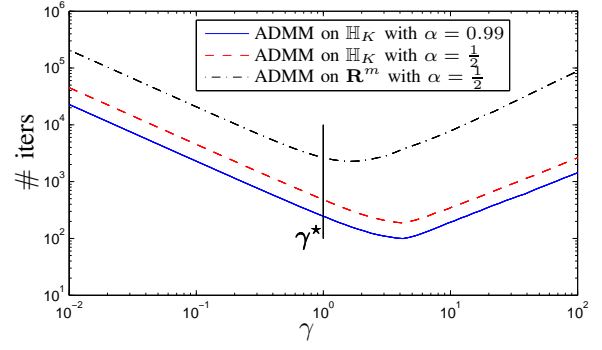


Fig. 2. Average number of iterations for different $\gamma$-values, different metrics, and different relaxations $\alpha$.

angles, while the inputs are the elevator and flaperon angles. The inputs are physically constrained to satisfy $|u_i| \le 25°$, $i = 1, 2$. The outputs are soft constrained and modeled using the piece-wise linear cost function

$$h(y, l, u, s) = \begin{cases} 0 & \text{if } l \le y \le u \\ s(y - u) & \text{if } y \ge u \\ s(l - y) & \text{if } y \le l \end{cases}$$

Especially, the first output is penalized using $h(y_1, -0.5, 0.5, 10^6)$ and the second is penalized using $h(y_2, -100, 100, 10^6)$. The states and inputs are penalized using

$$\ell(x, u, s) = \tfrac{1}{2}\big((x - x_r)^T Q(x - x_r) + u^T Ru\big)$$

where $x_r$ is a reference, $Q = \mathrm{diag}(0, 10^2, 0, 10^2)$, and $R = 10^{-2}I$. Further, the terminal cost is $Q$, and the control and prediction horizons are $N = 10$. The numerical data in Figure 2 is obtained by following a reference trajectory on the output. The objective is to change the pitch angle from $0°$ to $10°$ and then back to $0°$ while the angle of attack satisfies the (soft) output constraints $-0.5° \le y_1 \le 0.5°$. The constraints on the angle of attack limits the rate on how fast the pitch angle can be changed. By stacking vectors and forming appropriate matrices, the full optimization problem can be written on the form

$$\text{minimize} \quad \underbrace{\tfrac{1}{2}z^T Qz + r_t^T z + I_{Bz=bx_t}(z)}_{f(z)} + \underbrace{\sum_{i=1}^m h(z_i', \underline{d}_i, \bar{d}_i, 10^6)}_{g(z')}$$

$$\text{subject to} \quad Cz = z'$$

where $x_t$ and $r_t$ may change from one sampling instant to the next.

This is the optimization problem formulation discussed in Section V where item (ii) violates the assumptions that guarantee linear convergence. In Figure 2, the performance of the ADMM algorithm for different values of $\gamma$ and for different metrics is presented. Since the numerical example treated here is a model predictive control application, we can spend much computational effort offline to compute a metric that will be used in all samples in the online controller. We compute a diagonal metric $K = E^T E$ that minimizes the pseudo condition number of $ECP_{11}C^T E^T$, where $P_{11}$ is

implicitly defined in (28). This matrix $K$ defines the space $\mathbb{H}_K$ on which the algorithm is applied. In Figure 2, the performance of ADMM when applied on $\mathbb{H}_K$ with relaxations $\alpha = \frac{1}{2}$ and $\alpha = 0.99$, and ADMM applied on $\mathbf{R}^m$ with $\alpha = \frac{1}{2}$ is shown. In this particular example, improvements of about one order of magnitude are achieved when applied on $\mathbb{H}_K$ compared to when applied on $\mathbf{R}^m$. Figure 2 also shows that ADMM with over-relaxation performs better than standard ADMM. The proposed $\gamma$-parameter selection is denoted by $\gamma^\star$ in Figure 2 ($E$ or $C$ is scaled to get $\gamma^\star = 1$ for all examples). Figure 2 shows that $\gamma^\star$ does not coincide with the empirically found best $\gamma$, but still gives gives a reasonable choice of $\gamma$ in all cases.

## VII. CONCLUSIONS

We have shown tight linear convergence rate bounds for Douglas-Rachford splitting and ADMM. Based on these results, we have presented methods to select metric and algorithm parameters for these methods. We have also provided a numerical example to evaluate the proposed metric and parameter selection methods for ADMM. Performance improvements of about one order of magnitude, compared to when ADMM is applied on the Euclidean space, are reported.

## REFERENCES

[1] K. J. Arrow, L. Hurwicz, and H. Uzawa. *Studies in Linear and Nonlinear Programming*. Stanford University Press, 1958.

[2] H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer, 2011.

[3] H. H. Bauschke, J. Y. B. Cruz, T. T. A. Nghia, H. M. Phan, and X. Wang. The rate of linear convergence of the Douglas-Rachford algorithm for subspaces is the cosine of the Friedrichs angle. *Journal of Approximation Theory*, 185(0):63–79, 2014.

[4] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sciences*, 2(1):183–202, 2009.

[5] A. Bemporad, A. Casavola, and E. Mosca. Nonlinear control of constrained linear systems via predictive reference management. *IEEE Transactions on Automatic Control*, 42(3):340–349, 1997.

[6] D. Boley. Local linear convergence of the alternating direction method of multipliers on quadratic or linear programs. *SIAM Journal on Optimization*, 23(4):2183–2207, 2013.

[7] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3(1):1–122, 2011.

[8] E. J. Candes, J. Romberg, and T. Tao. Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 52(2):489–509, Feb 2006.

[9] A. Chambolle and T. Pock. A first-order primal-dual algorithm for convex problems withapplications to imaging. *Journal of Mathematical Imaging and Vision*, 40(1):120–145, 2011.

[10] P. L. Combettes and V. R. Wajs. Signal recovery by proximal forward-backward splitting. *SIAM journal on Multiscale Modeling and Simulation*, 4(4):1168–1200, 2005.

[11] D. Davis and W. Yin. Convergence rate analysis of several splitting schemes. Available http://arxiv.org/abs/1406.4834, August 2014.

[12] D. Davis and W. Yin. Faster convergence rates of relaxed Peaceman-Rachford and ADMM under regularity assumptions. Available: http://arxiv.org/abs/1407.5210, July 2014.

[13] L. Demanet and X. Zhang. Eventual linear convergence of the Douglas-Rachford iteration for basis pursuit. Available: http://arxiv.org/abs/1301.0542, May 2013.

[14] W. Deng and W. Yin. On the global and linear convergence of the generalized alternating direction method of multipliers. Technical Report CAAM 12-14, Rice University, 2012.

[15] J. Douglas and H. H. Rachford. On the numerical solution of heat conduction problems in two and three space variables. *Trans. Amer. Math. Soc.*, 82:421–439, 1956.

[16] J. Eckstein. *Splitting methods for monotone operators with applications to parallel optimization*. PhD thesis, MIT, 1989.

[17] D. Gabay. Applications of the method of multipliers to variational inequalities. In M. Fortin and R. Glowinski, editors, *Augmented Lagrangian Methods: Applications to the Solution of Boundary-Value Problems*. North-Holland: Amsterdam, 1983.

[18] D. Gabay and B. Mercier. A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Computers and Mathematics with Applications*, 2(1):17–40, 1976.

[19] E. Ghadimi, A. Teixeira, I. Shames, and M. Johansson. Optimal parameter selection for the alternating direction method of multipliers (ADMM): Quadratic problems. *IEEE Transactions on Automatic Control*, 60(3):644–658, March 2015.

[20] P. Giselsson. Tight linear convergence rate bounds for Douglas-Rachford splitting and ADMM. Available: http://arxiv.org/abs/1503.00887.

[21] P. Giselsson and S. Boyd. Diagonal scaling in Douglas-Rachford splitting and ADMM. In *53rd IEEE Conference on Decision and Control*, pages 5033–5039, Los Angeles, CA, December 2014.

[22] P. Giselsson and S. Boyd. Metric selection in fast dual gradient methods. 2014. Submitted.

[23] R. Glowinski and A. Marroco. Sur l'approximation, par éléments finis d'ordre un, et la résolution, par pénalisation-dualité d'une classe de problémes de dirichlet non linéaires. *ESAIM: Mathematical Modelling and Numerical Analysis - Modélisation Mathématique et Analyse Numérique*, 9:41–76, 1975.

[24] T. Hastie, R. Tibshirani, and J. Friedman. *The elements of statistical learning: data mining, inference and prediction*. Springer, 2nd edition, 2009.

[25] B. He and X. Yuan. On the $o(1/n)$ convergence rate of the Douglas-Rachford alternating direction method. *SIAM Journal on Numerical Analysis*, 50(2):700–709, 2012.

[26] R. Hesse, D. R. Luke, and P. Neumann. Alternating projections and Douglas-Rachford for sparse affine feasibility. *IEEE Transactions on Signal Processing*, 62(18):4868–4881, September 2014.

[27] M. Hong and Z.-Q. Luo. On the linear convergence of the alternating direction method of multipliers. Available: http://arxiv.org/abs/1208.3922, March 2013.

[28] F. Iutzeler, P. Bianchi, P. Ciblat, and W. Hachem. Explicit convergence rate of a distributed alternating direction method of multipliers. Available: http://arxiv.org/abs/1312.1085, December 2013.

[29] P. Kapasouris, M. Athans, and G. Stein. Design of feedback control systems for unstable plants with saturating actuators. In *Proceedings of the IFAC Symposium on Nonlinear Control System Design*, pages 302–307. Pergamon Press, 1990.

[30] P. L. Lions and B. Mercier. Splitting algorithms for the sum of two nonlinear operators. *SIAM Journal on Numerical Analysis*, 16(6):964–979, 1979.

[31] M. Lustig, D. Donoho, and J. M. Pauly. Sparse MRI: The application of compressed sensing for rapid MR imaging. *Magnetic Resonance in Medicine*, 58(6):1182–1195, 2007.

[32] R. Nishihara, L. Lessard, B. Recht, A. Packard, and M. Jordan. A general analysis of the convergence of ADMM. February 2015. Available: http://arxiv.org/abs/1502.02009.

[33] N. Parikh and S. Boyd. Proximal algorithms. *Foundations and Trends in Optimization*, 1(3):123–231, 2014.

[34] P. Patrinos, L. Stella, and A. Bemporad. Douglas-Rachford splitting: Complexity estimates and accelerated variants. In *Proceedings of the 53rd IEEE Conference on Decision and Control*, Los Angeles, CA, December 2014.

[35] D. W. Peaceman and H. H. Rachford. The numerical solution of parabolic and elliptic differential equations. *Journal of the Society for Industrial and Applied Mathematics*, 3(1):28–41, 1955.

[36] H. M. Phan. Linear convergence of the Douglas-Rachford method for two closed sets. Available: http://arxiv.org/abs/1401.6509, October 2014.

[37] A. Raghunathan and S. Di Cairano. ADMM for convex quadratic programs: Linear convergence and infeasibility detection. November 2014. Available: http://arxiv.org/abs/1411.7288.

[38] J. B. Rawlings and D. Q. Mayne. *Model Predictive Control: Theory and Design*. Nob Hill Publishing, Madison, WI, 2009.

[39] R. T. Rockafellar and R. J-B. Wets. *Variational Analysis*. Springer, Berlin, 1998.

[40] W. Shi, Q. Ling, K. Yuan, G. Wu, and W. Yin. On the linear convergence of the ADMM in decentralized consensus optimization. *IEEE Transactions on Signal Processing*, 62(7):1750–1761, April 2014.

[41] R Tibshirani. Regression shrinkage and selection via the lasso. *J. Royal. Statist. Soc B.*, 58(1):267–288, 1996.

## APPENDIX A
### PROOF TO PROPOSITION 2

*Proof.* Since $f$ is $\sigma$-strongly convex and $\beta$-smooth, $f_\gamma$ is $(1 + \gamma\sigma)$-strongly convex and $(1 + \gamma\beta)$-smooth. Therefore [2, Theorem 18.15] and [2, Theorem 13.32] directly imply that $f_\gamma^*$ is $\frac{1}{1+\gamma\sigma}$-smooth and $\frac{1}{1+\gamma\beta}$-strongly convex. From the smoothness definition in Definition 6, we get that

$$\frac{1}{2(1+\gamma\sigma)}\| \cdot \|^2 - f_\gamma^* \tag{29}$$

$$= \left(\frac{1}{2(1+\gamma\sigma)} - \frac{1}{2(1+\gamma\beta)}\right)\| \cdot \|^2 - (f_\gamma^* - \frac{1}{2(1+\gamma\beta)}\| \cdot \|^2)$$

is convex. Further, Definition 5 implies that $f_\gamma^* - \frac{1}{2(1+\gamma\beta)}\| \cdot \|^2$ is convex, and therefore (29) is the definition of $\frac{1}{2(1+\gamma\sigma)} - \frac{1}{2(1+\gamma\beta)}$-smoothness of $f_\gamma^* - \frac{1}{2(1+\gamma\beta)}\| \cdot \|^2$. Let $\beta = \sigma$, and we get that $f_\gamma^* - \frac{1}{2(1+\gamma\beta)}\| \cdot \|^2$ is 0-smooth, or equivalently by applying [2, Theorem 18.15], that $\nabla f_\gamma^* - \frac{1}{1+\gamma\beta}\mathrm{Id} = \mathrm{prox}_{\gamma f} - \frac{1}{1+\gamma\beta}\mathrm{Id}$ is 0-Lipschitz. When $\beta > \sigma$, [2, Theorem 18.15] implies that $\frac{1}{2(1+\gamma\sigma)} - \frac{1}{2(1+\gamma\beta)}$-smoothness of $f_\gamma^* - \frac{1}{2(1+\gamma\beta)}\| \cdot \|^2$ is equivalent to $\frac{1}{\frac{1}{1+\gamma\sigma} - \frac{1}{1+\gamma\beta}}$-cocoercivity of $\nabla f_\gamma^* - \frac{1}{1+\gamma\beta}\mathrm{Id} = \mathrm{prox}_{\gamma f} - \frac{1}{1+\gamma\beta}\mathrm{Id}$. This concludes the proof. □

## APPENDIX B
### PROOF TO THEOREM 1

*Proof.* To show this result, We first need the following lemmas.

*Lemma 2:* The function $\psi(x) = \frac{1-x}{1+x}$ is strictly decreasing for $x > -1$.

*Proof.* Let's define $\psi$ on $x > -1$. Then it is differentiable and $\psi'(x) = -\frac{2}{(1+x)^2} < 0$ for all $x > -1$. This concludes the proof. □

*Lemma 3:* Assume that $\beta > 0$. Then $\frac{1}{2\beta}$-cocoercivity of $\beta\mathrm{Id} + A$ is equivalent to $\beta$-Lipschitz continuity of $A$.

*Proof.* From the definition of cocoercivity, Definition 4, it follows directly that $\beta\mathrm{Id} + A$ is $\frac{1}{2\beta}$-cocoercive if and only if $\frac{1}{2\beta}(\beta\mathrm{Id} + A)$ is 1-cocoercive. This, in turn is equivalent to that $2\frac{1}{2\beta}(\beta\mathrm{Id} + A) - \mathrm{Id} = \frac{1}{\beta}A$ is nonexpansive [2, Proposition 4.2 and Definition 4.4]. Finally, from the definition of Lipschitz continuity, Definition 2, it follows directly that $\frac{1}{\beta}A$ is nonexpansive if and only if $A$ is $\beta$-Lipschitz continuous. This concludes the proof. □

*Lemma 4:* Assume that $A$ is $\frac{1}{\alpha}$-cocoercive. Then $A + \beta\mathrm{Id}$ is $\frac{1}{\alpha+\beta}$-cocoercive for any $\beta > 0$.

*Proof.* By applying Lemma 3 to $B = A - \frac{\alpha}{2}\mathrm{Id}$, we get that $\frac{1}{\alpha}$-cocoercivity of $B + \frac{\alpha}{2}\mathrm{Id} = A$ is equivalent to $\frac{\alpha}{2}$-Lipschitz continuity of $A - \frac{\alpha}{2}\mathrm{Id}$. This, in turn implies that $C = A - \frac{\alpha-\beta}{2}\mathrm{Id}$ is $\frac{\alpha+\beta}{2}$-Lipschitz due to the triangle inequality:

$$\|(A - (\tfrac{\alpha}{2} - \tfrac{\beta}{2})\mathrm{Id})x - (A - (\tfrac{\alpha}{2} - \tfrac{\beta}{2})\mathrm{Id})y\|$$

$$= \|(A - \tfrac{\alpha}{2})x - (A - \tfrac{\alpha}{2})y - \tfrac{\beta}{2}(x - y)\|$$

$$\leq \|(A - \tfrac{\alpha}{2})x - (A - \tfrac{\alpha}{2})y\| + \tfrac{\beta}{2}\|x - y\|$$

$$\leq (\tfrac{\alpha}{2} + \tfrac{\beta}{2})\|x - y\| = \tfrac{\alpha+\beta}{2}\|x - y\|$$

where we have used that $\beta > 0$. Applying Lemma 3 to $C = A - \frac{\alpha-\beta}{2}\mathrm{Id}$, we get that $C + \frac{\alpha+\beta}{2}\mathrm{Id} = (A - \frac{\alpha-\beta}{2}\mathrm{Id}) + \frac{\alpha+\beta}{2}\mathrm{Id} = A + \beta\mathrm{Id}$ is $\frac{1}{2\frac{(\alpha+\beta)}{2}} = \frac{1}{(\alpha+\beta)}$-cocoercive. This concludes the proof. □

*Lemma 5:* Suppose that $A + \alpha\mathrm{Id}$ is $\frac{1}{\alpha+\beta}$-cocoercive with $\alpha + \beta > 0$. Then $A$ is $\max(\alpha, \beta)$-Lipschitz continuous.

*Proof.* Let $B := A + \alpha\mathrm{Id}$. Then for $\alpha \geq \beta$, we have $\alpha > 0$ and

$$\langle Bx - By, x - y\rangle \geq \tfrac{1}{\beta+\alpha}\|Bx - By\|^2 \geq \tfrac{1}{2\alpha}\|Bx - By\|^2.$$

Using Lemma 3, this implies that $A$ is $\alpha$-Lipschitz continuous when $\alpha \geq \beta$. When $\alpha \leq \beta$, we have $\beta > 0$. Applying Lemma 4 on $B = A + \alpha\mathrm{Id}$ implies that $A + \beta\mathrm{Id} = B + (\beta - \alpha)\mathrm{Id}$ is $\frac{1}{(\alpha+\beta)+(\beta-\alpha)} = \frac{1}{2\beta}$-cocoercive. Therefore, Lemma 3 implies that $A$ is $\beta$-Lipschitz when $\beta \geq \alpha$. To conclude, $A$ is $\max(\alpha, \beta)$-Lipschitz, and the proof is complete. □

Now, we are ready to show the result. First, we show that $R_{\gamma f} + \frac{\gamma\beta-1}{1+\gamma\beta}\mathrm{Id}$ is $\frac{1}{2\left(\frac{1-\gamma\sigma}{1+\gamma\sigma} + \frac{\gamma\beta-1}{1+\gamma\beta}\right)}$-cocoercive when $\beta > \sigma$:

$$\langle(R_{\gamma f} + \tfrac{\gamma\beta-1}{1+\gamma\beta}\mathrm{Id})x - (R_{\gamma f} + \tfrac{\gamma\beta-1}{1+\gamma\beta}\mathrm{Id})y, x - y\rangle$$

$$= \langle(R_{\gamma f} - \tfrac{1-\gamma\beta}{1+\gamma\beta}\mathrm{Id})x - (R_{\gamma f} - \tfrac{1-\gamma\beta}{1+\gamma\beta}\mathrm{Id})y, x - y\rangle$$

$$= 2\langle(\mathrm{prox}_{\gamma f} - \tfrac{1}{1+\gamma\beta}\mathrm{Id})x - (\mathrm{prox}_{\gamma f} - \tfrac{1}{1+\gamma\beta}\mathrm{Id})y, x - y\rangle$$

$$\geq \tfrac{2}{\frac{1}{1+\gamma\sigma} - \frac{1}{1+\gamma\beta}}\|(\mathrm{prox}_{\gamma f} - \tfrac{1}{1+\gamma\beta}\mathrm{Id})x - (\mathrm{prox}_{\gamma f} - \tfrac{1}{1+\gamma\beta}\mathrm{Id})y\|^2$$

$$= \tfrac{1}{2\left(\frac{1}{1+\gamma\sigma} - \frac{1}{1+\gamma\beta}\right)}\|(2\mathrm{prox}_{\gamma f} - \mathrm{Id} - \tfrac{1-\gamma\beta}{1+\gamma\beta}\mathrm{Id})x$$
$$\qquad\qquad - (2\mathrm{prox}_{\gamma f} - \mathrm{Id} - \tfrac{1-\gamma\beta}{1+\gamma\beta}\mathrm{Id})y\|^2$$

$$= \tfrac{1}{\frac{1-\gamma\sigma}{1+\gamma\sigma} - \frac{1-\gamma\beta}{1+\gamma\beta}}\|(R_{\gamma f} - \tfrac{1-\gamma\beta}{1+\gamma\beta}\mathrm{Id})x - (R_{\gamma f} - \tfrac{1-\gamma\beta}{1+\gamma\beta}\mathrm{Id})y\|^2$$

$$= \tfrac{1}{\frac{1-\gamma\sigma}{1+\gamma\sigma} + \frac{\gamma\beta-1}{1+\gamma\beta}}\|(R_{\gamma f} + \tfrac{\gamma\beta-1}{1+\gamma\beta}\mathrm{Id})x - (R_{\gamma f} + \tfrac{\gamma\beta-1}{1+\gamma\beta}\mathrm{Id})y\|^2$$

where the inequality follows from Proposition 2. Now, since $\gamma\sigma < \gamma\beta$, Lemma 2 implies that $\frac{1-\gamma\sigma}{1+\gamma\sigma} + \frac{\gamma\beta-1}{1+\gamma\beta} = \frac{1-\gamma\sigma}{1+\gamma\sigma} - \frac{1-\gamma\beta}{1+\gamma\beta} > 0$. Therefore Lemma 5 can be applied and it implies that $R_{\gamma f}$ is $\max(\frac{1-\gamma\sigma}{1+\gamma\sigma}, \frac{\gamma\beta-1}{\gamma\beta+1})$-Lipschitz (since $R_{\gamma f} + \frac{\gamma\beta-1}{1+\gamma\beta}\mathrm{Id}$ is $\frac{1}{2\left(\frac{1-\gamma\sigma}{1+\gamma\sigma} + \frac{\gamma\beta-1}{1+\gamma\beta}\right)}$-cocoercive). When $\beta = \sigma$, we get

$$\|R_{\gamma f}x - R_{\gamma f}y\| = \|2\mathrm{prox}_{\gamma f}x - 2\mathrm{prox}_{\gamma f}y - (x - y)\|$$

$$= \|2((\mathrm{prox}_{\gamma f} - \tfrac{1}{1+\gamma\beta}\mathrm{Id})x - (\mathrm{prox}_{\gamma f} - \tfrac{1}{1+\gamma\beta}\mathrm{Id})y)$$
$$\qquad\qquad - (1 - \tfrac{2}{1+\gamma\beta})(x - y)\|$$

$$\leq 2\|(\mathrm{prox}_{\gamma f} - \tfrac{1}{1+\gamma\beta}\mathrm{Id})x - (\mathrm{prox}_{\gamma f} - \tfrac{1}{1+\gamma\beta}\mathrm{Id})y\|$$
$$\qquad\qquad + |(1 - \tfrac{2}{1+\gamma\beta})|\|x - y\|$$

$$= |(1 - \tfrac{2}{1+\gamma\beta})|\|x - y\| = \max(\tfrac{1-\gamma\beta}{1+\gamma\beta}, \tfrac{\gamma\beta-1}{1+\gamma\beta})\|x - y\|$$

where the second last equality follows from Proposition 2. Noting that $\sigma = \beta$ concludes the proof. □

## APPENDIX C
### PROOF TO THEOREM 2

*Proof.* By [2, Corollary 23.10], $R_{\gamma g}$ is nonexpansive and by Theorem 1, $R_{\gamma f}$ is $\delta := \max(\frac{1-\gamma\sigma}{1+\gamma\sigma}, \frac{\gamma\beta-1}{\gamma\beta+1})$-contractive. Therefore the composition $R_{\gamma g}R_{\gamma f}$ is also $\delta$-contractive since

$$\|R_{\gamma f}R_{\gamma g}z_1 - R_{\gamma f}R_{\gamma g}z_2\| \leq \delta\|R_{\gamma g}z_1 - R_{\gamma g}z_2\| \leq \delta\|z_1 - z_2\|. \tag{30}$$

for any $z_1, z_2 \in \mathcal{H}$. Now, let $T = (1 - \alpha)I + \alpha R_{\gamma f}R_{\gamma g}$ be the generalized Douglas-Rachford operator in (5). Since $\bar{z}$ is a fixed-point to $R_{\gamma f}R_{\gamma g}$ it is also a fixed-point to $T$, i.e., $\bar{z} = T\bar{z}$. Thus

$$
\begin{aligned}
\|z^{k+1} - \bar{z}\| &= \|Tz^k - T\bar{z}\|^2 \\
&= \|(1 - \alpha)(z^k - \bar{z}) + \alpha(R_{\gamma f}R_{\gamma g}z^k - R_{\gamma f}R_{\gamma g}\bar{z})\| \\
&\leq |1 - \alpha|\|z^k - \bar{z}\| + \alpha\|R_{\gamma f}R_{\gamma g}z^k - R_{\gamma f}R_{\gamma g}\bar{z}\| \\
&\leq (|1 - \alpha| + \alpha\delta)\|z^k - \bar{z}\| \\
&= \left(|1 - \alpha| + \alpha\max(\tfrac{1-\gamma\sigma}{1+\gamma\sigma}, \tfrac{\gamma\beta-1}{\gamma\beta+1})\right)\|z^k - \bar{z}\|
\end{aligned}
$$

where (30) is used in the second inequality. This concludes the proof. □