# DESCENT PROPERTIES OF AN ANDERSON ACCELERATED GRADIENT METHOD WITH RESTARTING*

WENQING OUYANG†, YANG LIU‡, AND ANDRE MILZAREK†

**Abstract.** Anderson Acceleration (AA) is a popular acceleration technique to enhance the convergence of fixed-point schemes. The analysis of AA approaches often focuses on the convergence behavior of a corresponding fixed-point residual, while the behavior of the underlying objective function values along the accelerated iterates is currently not well understood. In this paper, we investigate local properties of AA with restarting applied to a basic gradient scheme (AA-R) in terms of function values. Specifically, we show that AA-R is a local descent method and that it can decrease the objective function at a rate no slower than the gradient method up to higher-order error terms. These new results theoretically support the good numerical performance of AA(-R) when heuristic descent conditions are used for globalization and they provide a novel perspective on the convergence analysis of AA-R that is more amenable to nonconvex optimization problems. Numerical experiments are conducted to illustrate our theoretical findings.

**Key words.** Anderson Acceleration, Descent Properties, Restarting

**AMS subject classifications.** 90C30, 65K05, 90C06, 90C53

**1. Introduction.** In this work, we consider the smooth optimization problem

$$\min_{x \in \mathbb{R}^n} \ f(x), \tag{1.1}$$

where $f : \mathbb{R}^n \to \mathbb{R}$ is a continuously differentiable function. If the gradient mapping $\nabla f$ is additionally Lipschitz continuous with modulus $L$, then the basic gradient descent method with fixed step size,

$$x^{k+1} = x^k - \frac{1}{L} \nabla f(x^k) =: g(x^k), \tag{1.2}$$

can be utilized to solve problem (1.1). Here, $g : \mathbb{R}^n \to \mathbb{R}^n$ represents the associated gradient step mapping with step size $\frac{1}{L}$. The gradient descent step (1.2) can be viewed as a fixed-point iteration and the fixed-points of $g$ are exactly the stationary points of the objective function $f$.

Anderson Acceleration (AA) applies to fixed-point iterations of the form (1.2) and is a popular technique to accelerate the convergence of such iterative fixed-point schemes. For instance, AA-based algorithms have been applied successfully in computer graphics [31, 52, 27], reinforcement learning [15, 43], machine learning [49], and numerical methods for PDEs [34]. In iteration $k$ and based on the past $m$ iterations $\{x^{k-m}, \ldots, x^k\}$, AA first computes the mixing coefficients $\alpha^k = (\alpha_1^k, \ldots, \alpha_m^k)^\top \in \mathbb{R}^m$ as solution of the following optimization problem:

$$\min_{\alpha \in \mathbb{R}^m} \ \left\| h(x^{k-m}) + \sum_{i=1}^m \alpha_i (h(x^{k-m+i}) - h(x^{k-m})) \right\|^2, \tag{1.3}$$

†School of Data Science (SDS), Shenzhen Research Institute of Big Data (SRIBD), The Chinese University of Hong Kong, Shenzhen, China (wenqingouyang1@link.cuhk.edu.cn and andremilzarek@cuhk.edu.cn).

‡Mathematical Institute, University of Oxford, UK (yang.liu@maths.ox.ac.uk).

where $h(x) := g(x) - x$ denotes the residual map and $m$ is a corresponding memory parameter. AA then performs the accelerated iteration:

$$(1.4) \qquad x^{k+1} = g(x^{k-m}) + \sum_{i=1}^{m} \alpha_i^k (g(x^{k-m+i}) - g(x^{k-m})).$$

The parameter $m$ is usually chosen to be fixed or it is allowed to increase each iteration until a given threshold is reached after which $m$ is reinitialized. We will refer to such a restarted version of the Anderson accelerated gradient scheme (1.4) as AA-R (cf. Algorithm 2.1 in section 2). The goal of this paper is to analyze and establish novel descent properties of AA-R. In particular, we show that AA-R not only decreases the norm of the residual $\|h(x)\|$, but it can also decrease the underlying objective function $f$. Therefore, we answer the following question affirmatively:

>  *Can Anderson accelerated schemes achieve descent on the underlying objective function values?*

**1.1. Related Work and Literature.** Originally proposed by Anderson [1] for solving partial differential equations, AA has gained steadily growing attention during the last decade [31, 30, 19, 23, 27]. Though widely used in various fields and applications, the theoretical analysis and properties of AA are still somewhat limited. AA is known to belong to the class of multi-secant quasi-Newton methods [11, 12, 38]. When applied to linear problems, AA is equivalent to the generalized minimal residual method (GMRES) [47, 35]. For nonlinear problems, AA is also closely related to the nonlinear generalized minimal residual method (NGMRES) [48]. The convergence analysis in [46] shows that AA converges locally r-linearly under a smoothness condition on the map $g$ and uniform boundedness of the coefficients $\{\alpha^k\}_k$, but the obtained linear rate is slower than the rate of the original scheme. Later, in [10], the authors prove that AA can achieve an improved linear rate with additional quadratic error terms which overall yields r-linear convergence. This result is further improved in [33] by assuming sufficient linear independence on the set of difference vectors of the residuals $h(x^k)$ and q-linear convergence of AA is established with a rate faster than the Picard iteration (1.2). Moreover, if the coefficients $\{\alpha^k\}_k$ are assumed to be constant in each iteration, an asymptotic rate is given in [48]. The convergence behavior of AA applied to nonsmooth algorithmic schemes is also considered in [23, 6].

Since AA is known to only converge locally [46, 23], globalization mechanisms are required to use it in practice. A simple and heuristic choice is to check whether AA decreases the objective function value $f$ and to perform a fixed-point iteration if the decrease of the AA step is not sufficient. Such a strategy is utilized in [31, 27, 41, 17]. However, to the best of our knowledge, no consistent global-local convergence proofs are known in this case. Alternatively, one can check whether AA decreases the residual $\|h(x)\|$ and to reject the step if no sufficient decrease is observed. This strategy is more common and has been used in [51, 28, 14, 9]. Transition to local fast convergence of such a globalized AA approach is provided in [28]. Unfortunately, the convergence analyses in [51, 28, 14, 9] require global nonexpansiveness of $g$[1], which often necessitates convexity of $f$. Notably, it is also possible to combine function value- and residual-based globalization techniques, see, e.g., [17, 45].

Restarting strategies are part of many numerical algorithms. For instance, restarting is used in the conjugate gradient method (CG) [36], the generalized minimal residual method (GMRES) [40], and in quasi-Newton methods [24]. Restarting strategies have

---

[1]Though residual-based globalizations of AA without nonexpansiveness are possible, it is not fully clear which type of convergence guarantees can be achieved.

also been widely applied in the context of AA. In [4, 12], AA is restarted whenever the ratio of the square of the current residual to the sum of the squares of the previous residuals exceeds a predetermined constant. Similar ideas are discussed in [51, Section 3.2] and [32]. Restarting strategies are sometimes also used in tandem with regularization techniques to enhance the numerical stability of AA, see, e.g., [19, 43, 41]. A comprehensive comparison between AA with restarting and the original AA scheme on linear problems with parallel implementation can be found in [22]. The results in [22] suggest that the performance of AA with restarting is generally comparable to the performance of the original AA method. Further supporting observations for the effectiveness of restarted AA are provided in [37].

**1.2. Contributions.** The convergence analyses of AA [46, 10, 33, 28] focus on the decrease of the residual $\|h(x)\|$, which is natural since AA aims to minimize this norm in the AA subproblem (1.3). However, as mentioned, such globalization strategies usually require global nonexpansiveness of $g$ in order to obtain global convergence results, see, e.g., [51, 28, 14, 9]. Hence, the heuristic idea to base the acceptance of an AA step on the decrease of the objective function value seems attractive, since the Picard iteration (1.2) can decrease the function value even if $f$ is nonconvex. So far, there seems to be no theoretical backing ensuring that AA(-R) can achieve descent on the underlying objective function — even if strong convexity is assumed and an appropriate initial point is selected. Our aim is to investigate this gap and to show that AA-R can decrease the objective function value locally. This result provides theoretical guarantees for algorithms that utilize descent conditions for $f$ as globalization mechanism without hindering the local fast convergence of AA-R. We now summarize our main contributions:

- To the best of our knowledge, we establish the first descent properties of AA-R iterations for the gradient descent scheme (1.2). On the one hand, this illuminates the success of algorithms that have used heuristic descent-type conditions to globalize AA(-R) [31, 27, 41]. On the other hand, our findings can be utilized in the design of novel globalization techniques for AA-R methods.
- We verify that the iterates generated by AA-R in one restarting cycle are equivalent to the iterates generated by GMRES when being run on a perturbed linear system model (for the same amount of iterations). This model, without perturbation, is exactly the quadratic expansion of $f$ and the Hessian of this model is symmetric and positive definite if we assume local strong convexity. Hence, AA-R is close to running the conjugate residual method (CR) [44] on such a quadratic model of $f$. Motivated by these observations, we analyze and specify the error between GMRES and CR under small perturbations of the system matrix which will allow us to link AA-R, CR, and CG. Based on classical results for CG [20], we then show that AA-R not only decreases the objective function value, but the overall achieved descent is actually no smaller than the one obtained by performing a gradient descent step with step size $\frac{1}{L}$ up to higher-order error terms. Some byproducts of our results indicate that CR itself decreases the objective function value no slower than the gradient descent method.
- We design a practical function value-based globalization mechanism for AA-R. Unlike residual-based globalizations, this allows to apply AA-R directly to nonconvex problems without requiring any adjustments of the underlying AA-R scheme. We illustrate the numerical performance of our simple globalization and numerically confirm the derived theoretical descent guarantees of AA-R

on several nonconvex large-scale problems.

**1.3. Organization.** This work is organized as follows. In section 2, we introduce the algorithmic details of AA-R and list the standing assumptions. In section 3, we derive the core descent properties of AA-R. This is done step by step. In subsection 3.1, we first establish q-linear convergence. The mentioned equivalence between AA-R and GMRES is shown in subsection 3.3. Next, in subsection 3.4, we analyze the error between the sequences generated by CR and GMRES which allows to connect GMRES and CR. The detailed connection between CR and CG is investigated in subsection 3.5. An objective function value-based globalization of AA-R is presented in section 4. Finally, in section 5, we verify our theoretical results and test the proposed globalized AA-R algorithm on several examples.

**1.4. Notation.** Throughout this work, we consider the fixed-point mapping $g(x) = x - \nabla f(x)/L$ and the corresponding residual $h(x) := g(x) - x = -\nabla f(x)/L$. For a given sequence of iterates $\{x^k\}_k$, we define the terms:

$$M_h^k := \max_{k-\hat{m} \leq i \leq k} \|h(x^i) - h(x^{k-\hat{m}})\|, \quad M_x^k := \max_{k-\hat{m} \leq i \leq k} \|x^i - x^{k-\hat{m}}\|,$$

where $\hat{m} = \mathrm{mod}(k, m+1)$. We further introduce the following matrices and notations:

$$X_k := [x^{k-\hat{m}+1} - x^{k-\hat{m}}, \ldots, x^k - x^{k-\hat{m}}] \in \mathbb{R}^{n \times \hat{m}},$$
$$H_k := [h(x^{k-\hat{m}+1}) - h(x^{k-\hat{m}}), \ldots, h(x^k) - h(x^{k-\hat{m}})] \in \mathbb{R}^{n \times \hat{m}},$$
$$G_k := [g(x^{k-\hat{m}+1}) - g(x^{k-\hat{m}}), \ldots, g(x^k) - g(x^{k-\hat{m}})] \in \mathbb{R}^{n \times \hat{m}},$$

$\hat{x}^0 := x^0$, $\hat{x}^k := x^{k-\hat{m}} + X_k \alpha^k$, $g^k := g(x^k)$, $\hat{g}^k := g^{k-\hat{m}} + G_k \alpha^k$, $h^k := h(x^k)$, and $\hat{h}^k := \hat{g}^k - \hat{x}^k$. The definition of $\hat{x}^k$ follows [10, Equation (2.4)]. We further note that the AA subproblem (1.3) can be viewed as finding the minimal value of the linearized residual of $\hat{x}^k$ which is $\hat{h}^k$. Based on these notations, we can express the solution to (1.3) explicitly by $\alpha^k = -(H_k^\top H_k)^{-1} H_k^\top h^{k-\hat{m}}$ provided that $H_k^\top H_k$ is invertible.

For given $n \in \mathbb{N}$, we set $[n] := \{1, \ldots, n\}$. For a matrix $A$, $\sigma_{\max}(A)$ ($\lambda_{\max}(A)$) denotes the largest singular value (eigenvalue) of $A$ and $\sigma_{\min}(A)$ ($\lambda_{\min}(A)$) is the smallest singular value (eigenvalue) of $A$. The condition number of $A$ is given by $\kappa(A) := \sigma_{\max}(A)/\sigma_{\min}(A)$. Unless specified otherwise, the norm $\|\cdot\|$ refers to the standard Euclidean norm for vectors and the spectral norm for matrices. We will use $\|\cdot\|_F$ to denote the Frobenius norm of a matrix. The space $\mathcal{K}^t(A, b) := \mathrm{span}\{b, Ab, \ldots, A^{t-1}b\}$ is used to denote the $t$-th Krylov space generated by $A$ and $b$.

**2. Anderson Acceleration with Restarting.** Throughout this paper, we assume that:

ASSUMPTION 2.1. *There is some $r > 0$ and a stationary point $x^\star \in \mathbb{R}^n$ of $f$ (i.e., $\nabla f(x^\star) = 0$) such that:*
  (A.1) *The function $f$ is $L$-smooth on $\mathbb{R}^n$.*
  (A.2) *The function $f$ is $\mu$-strongly convex on $\mathbb{B}_r(x^\star) := \{x : \|x - x^\star\| < r\}$.*
  (A.3) *The Hessian $\nabla^2 f$ is Lipschitz continuous with modulus $L_H$ on $\mathbb{B}_r(x^\star)$.*

We note that the conditions formulated in Assumption 2.1 are common in the convergence analysis of Anderson accelerated gradient methods, see, e.g., [46, 10] for comparison. Let $\kappa_r := \frac{L}{\mu}$ denote the condition number of the Hessian $\nabla^2 f$ on $\mathbb{B}_r(x^\star)$. Then, under Assumption 2.1, it follows:

$$(2.1) \quad \|g(x) - g(y)\| \leq \left\| I - \frac{1}{L}\bar{F} \right\| \|x - y\| \leq \left(1 - \frac{1}{\kappa_r}\right) \|x - y\|, \quad \forall\, x, y \in \mathbb{B}_r(x^\star),$$

where $\bar{F} := \int_0^1 \nabla^2 f(y + t(x - y)) \, dt$. Hence, $g$ is contractive on $\mathbb{B}_r(x^\star)$ with Lipschitz constant $1 - \frac{1}{\kappa_r}$.

We study local properties of AA with restarting applied to the gradient mapping $g(x) = x - \frac{1}{L}\nabla f(x)$. In particular, given some initial point $x^0$ which is sufficiently close to $x^\star$, we apply AA on $\{x^0, \ldots, x^{k-1}\}$ to obtain the new iterate $x^k$. After $m$ iterations ($m$ is the fixed memory parameter), this procedure is stopped and AA is restarted with $x^{m+1}$ as new initial point of the next cycle. The full algorithm of the restarted AA scheme for problem (1.1) – AA-R – is shown below in Algorithm 2.1.

---

**Algorithm 2.1** AA with Restarting (AA-R)

---

**Require:** Choose the initial point $x^0$ and the memory parameter $m \in \mathbb{N}$. Set the current memory parameter as $\hat{m} = 0$.
1: **for** $k = 0, 1, \ldots$ **do**
2:     $\hat{m} = \mathrm{mod}(k, m + 1)$.
3:     **if** $\hat{m} = 0$ **then**
4:         Set $x^{k+1} = g(x^k)$.
5:     **else**
6:         Calculate the coefficient $\alpha^k$ based on the sequence $\{h(x^k), \ldots, h(x^{k-\hat{m}})\}$ via solving (1.3) and set $x^{k+1} = g^{k-\hat{m}} + G_k \alpha^k$.
7:     **end if**
8: **end for**

---

**3. Convergence and Descent Properties of AA-R.** Most classical convergence analyses of AA are based on the same idea — linearization [46, 10, 23, 28]. It is well-known that AA(-R) is equivalent to GMRES if the mapping $g$ is affine [35] and therefore, AA(-R) can be viewed as GMRES being applied to a linear approximation of $g$ along with some linearization error. When specialized to the gradient mapping, these analyses ignore the structural information that $g$ has a symmetric Hessian. Taking this information into account, we can deduce that the system matrix of the GMRES procedure is essentially close to a symmetric positive definite matrix, which means that GMRES is close to CR in this case. This observation motivates us to utilize classical tools for CR to show that AA-R locally performs descent steps for $f$. The main technical difficulty lies in the fact that the iterates generated by AA-R only coincide with the iterates generated by GMRES after one additional gradient step. We resolve this complication by connecting GMRES and CG (via CR) and by analyzing the relevant properties via a CG-based perspective.

**3.1. Q-Linear Convergence of AA-R.** We first present an additional assumption and several basic properties of AA-R (and AA) that allow to establish q-linear convergence. This will serve as a foundation for our later results.

ASSUMPTION 3.1.    (A.4) *The condition number of $X_k^\top X_k$ is bounded by $M^2$ for every $k \in \mathbb{N}$.*

Let us note that the analogous assumption on $H_k$ is more common, since $H_k$ appears directly in the computation of the coefficient $\alpha_k$ in (1.3). We address this issue in the following proposition and show that condition (A.4) is actually equivalent to assuming that the condition number of $H_k^\top H_k$ is bounded locally. Hence, in practice, (A.4) can be ensured by monitoring the condition number of $H_k^\top H_k$. For instance, we can restart the current cycle whenever the condition number of $H_k^\top H_k$ exceeds a given tolerance. Alternative strategies are further discussed in [33, Section 5.1.3]. It

is also possible to mitigate condition (A.4) and related boundedness assumptions on the mixing coefficients $\{\alpha^k\}_k$, [46, 10, 23, 6] via algorithmic independence checks or adaptive depth mechanisms, see, e.g., [7]. However, such adjustments naturally affect the achievable convergence and acceleration results.

PROPOSITION 3.2. *Suppose that the conditions* (A.1)–(A.3) *are satisfied. Then, the following statements hold:*
  (i) *For every $M_X > 0$ there is a neighborhood $U_1$ of $x^\star$ such that if $x^{k-\hat{m}} \ldots, x^k \in U_1$ and $\kappa(X_k^\top X_k) \leq M_X^2$, then it holds that $\kappa(H_k^\top H_k) \leq 4\kappa_r^2 \kappa(X_k^\top X_k)$.*
  (ii) *For every $M_H > 0$ there is a neighborhood $U_2$ of $x^\star$ such that if $x^{k-\hat{m}}, \ldots, x^k \in U_2$ and $\kappa(H_k^\top H_k) \leq M_H^2$, then we have $\kappa(X_k^\top X_k) \leq 4\kappa_r^2 \kappa(H_k^\top H_k)$.*

*Proof.* Without loss of generality and in order to simplify the notation, we assume $\hat{m} = m$ and $k = m$. Let us first define $U_1 := \mathbb{B}_{\delta_1}(x^\star)$ where $\delta_1 := \min\{r, (1 - \frac{1}{\sqrt{2}})\frac{\mu}{\sqrt{m}M_X L_H}\}$. We further set $b_i := \nabla f(x^i) - \nabla f(x^0) - \nabla^2 f(x^0)(x^i - x^0)$ and $B_k := [b_1, \ldots, b_k]$. Utilizing [25, Lemma 4.1.1], it follows $\|b_i\| \leq \frac{L_H}{2}\|x^i - x^0\|^2$ for all $i \in [k]$ and we obtain

$$\|B_k\| \leq \|B_k\|_F \leq \sqrt{m} \max_{1 \leq i \leq k} \|b_i\| \leq \frac{\sqrt{m}L_H}{2} \max_{1 \leq i \leq k} \|x^i - x^0\|^2 = \frac{\sqrt{m}L_H(M_x^k)^2}{2}.$$

Defining $E_k := B_k(X_k^\top X_k)^{-1}X_k^\top$ ($E_k$ is well-defined since $X_k^\top X_k$ is invertible) and $A_k = \nabla^2 f(x^0) + E_k$, a direct calculation yields

$$(3.1) \qquad A_k(x^i - x^0) = \nabla f(x^i) - \nabla f(x^0) \quad \forall\, i \in [k].$$

In other words, we have $A_k X_k = -LH_k$. Moreover, the norm of $E_k$ can be estimated as follows

$$(3.2) \qquad \begin{aligned} \|E_k\| &\leq \|B_k\|\|(X_k^\top X_k)^{-1}X_k^\top\| \\ &= \|B_k\|\sqrt{\|(X_k^\top X_k)^{-1}\|} \leq \frac{\|B_k\|_F M_X}{\|X_k\|} \leq \frac{\sqrt{m}L_H M_X M_x^k}{2}, \end{aligned}$$

where we used $\|X_k\| \geq \max_{1 \leq i \leq k} \|X_k e_i\| = M_x^k$ with $e_i \in \mathbb{R}^k$ being the $i$-th unit vector. Due to $M_x^k \leq 2\delta_1$, we can further infer $\|E_k\| \leq \sqrt{m}L_H M_X \delta_1 \leq (1 - \frac{1}{\sqrt{2}})\mu$. Therefore, it holds that $\sigma_{\max}(A_k) \leq \lambda_{\max}(\nabla^2 f(x^0)) + \sigma_{\max}(E_k) \leq \sqrt{2}L$ and $\sigma_{\min}(A_k) \geq \lambda_{\min}(\nabla^2 f(x^0)) - \sigma_{\max}(E_k) \geq \frac{1}{\sqrt{2}}\mu$. Consequently, we have $\kappa(A_k) \leq 2\kappa_r$. This allows to bound the condition number of $H_k$:

$$\kappa(H_k^\top H_k) = \kappa(X_k^\top A_k^\top A_k X_k) \leq \kappa(X_k^\top X_k)\kappa(A_k^\top A_k) \leq 4\kappa_r^2 M_X^2$$

and proves part (i). We now turn to the proof of the second statement. We define $U_2 := \mathbb{B}_{\delta_2}(x^\star)$, $\delta_2 := \min\{r, (1 - \frac{1}{\sqrt{2}})\frac{\mu^2}{\sqrt{m}LL_H M_H}\}$, $\tilde{b}_i = \frac{1}{L}[(\nabla^2 f(x^0))^{-1}(\nabla f(x^i) - \nabla f(x^0)) - (x^i - x^0)]$, and $\tilde{B}_k = [\tilde{b}_1, \ldots, \tilde{b}_k]$. As before, we obtain

$$\|\tilde{B}_k\| = \frac{1}{L}\|\nabla f(x^0)^{-1}B_k\| \leq \frac{1}{L\mu}\|B_k\| \leq \frac{\sqrt{m}L_H(M_x^k)^2}{2L\mu}.$$

Let $j \in [k]$ be given with $M_x^k = \max_{1 \leq i \leq k}\|x^i - x^0\| = \|x^j - x^0\|$. Using (A.2), it then holds that

$$(3.3) \qquad \|H_k\| \geq M_h^k \geq \|h^j - h^0\| = \frac{1}{L}\|\nabla f(x^j) - \nabla f(x^0)\| \geq \frac{1}{\kappa_r}\|x^j - x^0\| = \frac{M_x^k}{\kappa_r}.$$

Therefore, setting $\tilde{E}_k := \tilde{B}_k (H_k^\top H_k)^{-1} H_k^\top$, it follows

$$\|\tilde{E}_k\| \leq \|\tilde{B}_k\| \sqrt{\|(H_k^\top H_k)^{-1}\|} \leq \frac{\sqrt{m} L_H (M_x^k)^2 M_H}{2 L \mu \|H_k\|} \leq \frac{\sqrt{m} L_H M_H M_x^k}{2\mu^2}$$

and by the definition of $\delta_2$, we have $M_x^k \leq 2\delta_2$ and $\|\tilde{E}_k\| \leq \frac{\sqrt{m} L_H M_H \delta_2}{\mu^2} \leq (1 - \frac{1}{\sqrt{2}}) \frac{1}{L}$. Next, defining $\tilde{A}_k := (\nabla^2 f(x^0))^{-1} + \tilde{E}_k$, we can again infer

$$\tilde{A}_k (\nabla f(x^i) - \nabla f(x^0)) = x^i - x^0, \quad i \in [k] \quad \Longrightarrow \quad L \tilde{A}_k H_k = -X_k,$$

$\sigma_{\max}(\tilde{A}_k) \leq \lambda_{\max}(\nabla^2 f(x^0)^{-1}) + \sigma_{\max}(\tilde{E}_k) \leq \frac{\sqrt{2}}{\mu}$, and $\sigma_{\min}(\tilde{A}_k) \geq \lambda_{\min}(\nabla^2 f(x^0)^{-1}) - \sigma_{\max}(\tilde{E}_k) \geq 1/(\sqrt{2}L)$. This yields $\kappa(\tilde{A}_k) \leq 2\kappa_r$ and $\kappa(X_k^\top X_k) \leq \kappa(H_k^\top H_k) \kappa(\tilde{A}_k^\top \tilde{A}_k) \leq 4\kappa_r^2 M_H^2$. $\qquad\square$

We note that the matrix $A_k = \nabla^2 f(x^0) + E_k$ defined in the proof of Proposition 3.2 is a key technical ingredient in this paper and will be used in the subsequent sections. The matrix $A_k$ consists of the symmetric Hessian $\nabla^2 f(x^0)$ and the perturbation matrix $E_k$. Our goal in the next subsections is to suitably control the norm of $E_k$ promoting a link between AA-R and CR.

Based on condition (A.4), we now verify q-linear convergence of the sequence $\{\|h^k\|\}_k$. Let us remark that assumption (A.4) (or its equivalent formulation for the matrices $H_k^\top H_k$) is generally stronger than the condition appearing in [33]. Namely, in [33, Theorem 5.1], q-linear convergence of AA(-R) is shown under a sufficient linear independence condition on each of the columns of (a permutation of) $H_k$. Here, we will work with the slightly stronger assumption (A.4) as it allows us to study the behavior of $(H_k^\top H_k)^{-1} H_k^\top$ under perturbations which is required to link AA-R and CR. More details can be found in the proof of Theorem 3.12. In addition, in [33], contraction and Lipschitz differentiability of $\nabla f$ is assumed on the whole space $\mathbb{R}^n$, while we consider the local case in a neighborhood of $x^\star$. Following the derivation in [33], we first analyze the behavior of the residuals $\{\|h^k\|\}_k$ in one cycle of AA-R.

PROPOSITION 3.3. *Let the conditions* (A.1)–(A.4) *be satisfied and let us further assume* $g(x^i), x^i \in U_1$ *with* $i = k - \hat{m}, \ldots, k$ *and* $\hat{x}^k, \hat{g}^k \in U_1$ *(where* $U_1$ *is introduced in* Proposition 3.2 *for* $M_X = M$*). Then, it holds that:*

$$(3.4) \qquad \|h^{k+1}\| \leq \left(1 - \frac{1}{\kappa_r}\right) \|h^k\| + \mathcal{O}\left(\|h^k\| \left(\sum\nolimits_{i=k-\hat{m}}^k \|h^i\|\right)\right).$$

*Proof.* This result basically follows from [33, Theorem 5.1]. A comprehensive proof is presented in Appendix A. $\qquad\square$

In order to transfer the statement in Proposition 3.3 to the full sequence $\{\|h^k\|\}_k$, we need to show that AA-R stays indeed local.

LEMMA 3.4. *Let the conditions* (A.1)–(A.3) *hold and assume* $x^{k-\hat{m}}, \ldots, x^k \in \mathbb{B}_r(x^\star)$. *If* $\kappa(H_k^\top H_k) \leq M_H^2$, *then we have:*

$$\|\hat{x}^k - x^{k-\hat{m}}\| \leq \sqrt{m} M_H \kappa_r \|h^{k-\hat{m}}\| \quad and \quad \|\hat{g}^k - x^{k-\hat{m}}\| \leq (1 + \sqrt{m} M_H \kappa_r) \|h^{k-\hat{m}}\|.$$

*Proof.* We start with bounding the coefficient $\alpha^k$. As mentioned, the closed-form expression for $\alpha^k$ is given by $\alpha^k = -(H_k^\top H_k)^{-1} H_k^\top h^{k-\hat{m}}$. Therefore, it holds that

$$\|\alpha^k\| \leq \sqrt{\|(H_k^\top H_k)^{-1}\|} \|h^{k-\hat{m}}\| \leq \frac{M_H}{\|H_k\|} \|h^{k-\hat{m}}\| \leq \frac{M_H \kappa_r}{M_x^k} \|h^{k-\hat{m}}\|,$$

where we used (3.3) in the last inequality. Due to $\|X_k\|_F \leq \sqrt{\hat{m}} M_x^k$, the definition of $\hat{h}^k$, and (1.3), this implies

$$\|\hat{x}^k - x^{k-\hat{m}}\| = \|X_k \alpha^k\| \leq \|X_k\|\|\alpha^k\| \leq \|X_k\|_F \|\alpha^k\| \leq \sqrt{m} M_H \kappa_r \|h^{k-\hat{m}}\|$$

and $\|\hat{g}^k - x^{k-\hat{m}}\| \leq \|\hat{h}^k\| + \|\hat{x}^k - x^{k-\hat{m}}\| \leq (1 + \sqrt{m} M_H \kappa_r)\|h^{k-\hat{m}}\|$.                $\square$

PROPOSITION 3.5. *Let* (A.1)–(A.4) *be satisfied and let* $\{x^k\}_k$ *be generated by* Algorithm 2.1. *Then there exists a neighborhood* $U$ *of* $x^\star$ *such that if* $x^0 \in U$, *it follows* $\{x^k\}_k \subset U$ *and* $\|h(x^{k+1})\| \leq (1 - \frac{1}{2\kappa_r})\|h(x^k)\|$ *for all* $k \in \mathbb{N}$.

*Proof.* We define $S_\epsilon = \{x \in \mathbb{R}^n : \|h(x)\| \leq \epsilon\} \cap \mathbb{B}_r(x^\star)$. Due to (A.2), we obtain

$$\|h(x)\| = \frac{1}{L}\|\nabla f(x)\| \geq \frac{\mu}{L}\|x - x^\star\| = \frac{1}{\kappa_r}\|x - x^\star\| \quad \forall\, x \in \mathbb{B}_r(x^\star).$$

Let $U_1$ be defined as in Proposition 3.2 for $M_X = M$. The previous inequality implies that there is some $\epsilon_1 > 0$ such that $S_{\epsilon_1} \subset U_1$. Moreover, by Proposition 3.3, there exists another neighborhood $S_{\epsilon_2}$ such that if $x^k, \dots, x^{k-\hat{m}}, \hat{x}^k, \hat{g}^k \in S_{\epsilon_2}$, then we have

$$(3.5) \qquad\qquad \|h(\hat{g}^k)\| \leq (1 - (2\kappa_r)^{-1})\|h(x^k)\|.$$

We now take $\bar{\epsilon} = \min\{\epsilon_1, \frac{\epsilon_2}{2 + 2\sqrt{m}\kappa_r^2 M}, \frac{r}{1 + (2\sqrt{m}\kappa_r M + 1)\kappa_r}\}$ and set $U = S_{\bar{\epsilon}}$. Let us further suppose $x^0 \in U$. We use an induction to show $\|h^k\| \leq (1 - (2\kappa_r)^{-1})\|h^{k-1}\|$ and $x^k \in U$ for all $k$. It is clear that we only need to prove this conclusion for $k = 1, \dots, m+1$, since the analysis is identical for the next cycle of the restarted AA method. We start with $k = 1$. By definition, we have $x^1 = g(x^0)$ and according to (2.1), it follows:

$$\|x^1 - x^\star\| = \|g(x^0) - g(x^\star)\| \leq (1 - \kappa_r^{-1})\|x^0 - x^\star\| \leq r.$$

This proves $x^1 \in \mathbb{B}_r(x^\star)$. Moreover, it holds that:

$$\|h(x^1)\| = \|g(x^1) - x^1\| = \|g(x^1) - g(x^0)\| \leq (1 - \kappa_r^{-1})\|h(x^0)\| \leq \bar{\epsilon},$$

which shows $\|h(x^1)\| \leq (1 - (2\kappa_r)^{-1})\|h(x^0)\|$ and $x^1 \in S_{\bar{\epsilon}}$. Next, let us assume that the induction hypothesis is true for $i = 1, \dots, k$ and let us prove the conclusion for $i = k + 1$. Applying Proposition 3.2 and (A.4), we obtain $\kappa(H_k^\top H_k) \leq 4\kappa_r^2 M^2$ and using Lemma 3.4, we can further infer:

$$\|x^{k+1} - x^\star\| \leq \|\hat{g}^k - x^0\| + \|x^0 - x^\star\| \leq (1 + 2\sqrt{m}M\kappa_r^2)\|h^0\| + \kappa_r\|h^0\| \leq r.$$

This establishes $x^{k+1} \in \mathbb{B}_r(x^\star)$. Similarly, we also have $\hat{x}^k \in \mathbb{B}_r(x^\star)$ and it holds that

$$\|h(\hat{x}^k)\| \leq \|h(\hat{x}^k) - h(x^0)\| + \|h^0\| \leq \|\hat{x}^k - x^0\| + \bar{\epsilon} \leq (1 + 2\sqrt{m}\kappa_r^2 M)\bar{\epsilon} \leq \epsilon_2$$

and $\|h(\hat{g}^k)\| \leq \epsilon_2$. Hence, by (3.5), we can conclude $\|h^{k+1}\| = \|h(\hat{g}^k)\| \leq (1 - \frac{1}{2\kappa_r})\|h^k\|$ and $x^{k+1} \in S_{\bar{\epsilon}}$.                $\square$

**3.2. Local Descent Properties of AA-R.** We now formulate and present one of the main theoretical results of this paper.

THEOREM 3.6. *Let the conditions* (A.1)–(A.4) *hold and let* $\{x^k\}_k$ *be generated by* Algorithm 2.1. *There is a neighborhood* $U$ *of* $x^\star$ *such that if* $x^0 \in U$, *then we have:*

$$(3.6) \qquad\qquad f(x^{k+1}) \leq f(g(x^k)) + \mathcal{O}(\|\nabla f(x^{k-\hat{m}})\|^3)$$

*for all* $k \in \mathbb{N}$, *where* $\hat{m} = \mathrm{mod}(k, m+1)$.

As already outlined, the proof of Theorem 3.6 relies on subtle connections between AA-R, GMRES, CR, and CG. We will establish and discuss these connections step-by-step in the subsequent subsections.

Before proceeding with further details, let us briefly discuss GMRES, CR, and CG, cf. [39]. All of these algorithms are designed to solve linear systems of form $Ax = b$ via Krylov subspace techniques. For GMRES and CR, the $k$-th iterate is the point in the $k$-th Krylov subspace $\mathcal{K}^k(A, b)$ with minimal residual norm $\|Ax - b\|$. Here, CR typically requires the matrix $A$ to be symmetric positive semidefinite, which can be exploited in (faster) implementations. No additional assumptions (on $A$) need to be made when applying GMRES. CG is connected to CR and requires $A$ to be symmetric, positive (semi)definite. Instead of finding elements with minimal residual norm, CG aims at minimizing the quadratic form $\frac{1}{2}x^\top Ax - b^\top x$ within the subspace $\mathcal{K}^k(A, b)$.

We now summarize the core components of our proof. In subsection 3.3, we show that the AA-R iterate $x^{k+1}$ coincides with an iterate $\bar{x}_G^k$ that can be generated via performing an additional gradient step on the GMRES iterate $x_G^k$. Here, GMRES is applied to a non-symmetric perturbed linear system $A(x - x^0) = b$ that is connected to AA-R. In addition, the gradient step $g(x^k)$ can be viewed as applying two gradient steps on the previous GMRES iterate $x_G^{k-1}$ resulting in $\tilde{x}_G^{k-1}$. Since the system matrix $A$ can be interpreted as a perturbed version of the symmetric Hessian $\nabla^2 f(x^0)$, our idea is to run CR on the linear system $B(x - x^0) = b$ with $B = \nabla^2 f(x^0)$ and $b = -\nabla f(x^0)$ and to bound the differences between the CR and GMRES iterates. In subsection 3.4, we verify that this error has order $\mathcal{O}(\|b\|^2)$. Finally, in subsection 3.5, we connect CR and CG and use the rich computational properties of CG, [20], to show that the CR iterates achieve the desired descent on a local quadratic model of $f$. In the last subsection, we combine these different components to prove that the AA-R iterate $x^{k+1}$ itself decreases the objective function value up to certain higher-order error terms.

**3.3. Connecting AA-R and GMRES.** We now establish equivalence of GMRES and AA-R when running one restarting cycle. Let us introduce the matrix $A := A_m = \nabla^2 f(x^0) + E_m$, where $A_m$ and $E_m$ have been defined in the proof of Proposition 3.2, i.e., it holds that

$$(3.7) \qquad A = \nabla^2 f(x^0) + E_m, \quad E_m = B_m(X_m^\top X_m)^{-1} X_m^\top, \quad B_m = [b_1, \ldots, b_m],$$

where $b_i = \nabla f(x^i) - \nabla f(x^0) - \nabla^2 f(x^0)(x^i - x^0)$, $i \in [m]$.

The matrix $A$ can be utilized to construct a new perturbed gradient mapping $\bar{g}(x) := x - \frac{1}{L}(A(x - x^0) + \nabla f(x^0))$. Recalling (3.1), it can be shown that the function $\bar{g}$ is exact at $x^k$ for all $k = 0, \ldots, m$, i.e., we have $\bar{g}(x^k) = g(x^k)$. We now study GMRES applied to the linear system $A(x - x^0) = -\nabla f(x^0)$. While there are various implementations of GMRES [40, 5], each of these variants will yield the same iteration sequence $\{x_G^k\}_k$. The following proposition (which holds for general input data $A \in \mathbb{R}^{n \times n}$ and $b, x^0 \in \mathbb{R}^n$) is taken from [35, Equation (4)] and characterizes the iterates $x_G^k$ more explicitly.

PROPOSITION 3.7. *Let $A \in \mathbb{R}^{n \times n}$ be nonsingular and let $b, x^0 \in \mathbb{R}^n$ be given. Suppose further that $\{x_G^k\}_k$ is generated by GMRES to solve the system $A(x - x^0) = b$ with $x_G^0 = x^0$. Then, we have:*

$$x_G^k = \mathrm{argmin}_{x \in x^0 + \mathcal{K}^k(A,b)} \|A(x - x^0) - b\|^2.$$

Next, we verify the equivalence of AA-R and GMRES (in one restarting cycle) in the general nonlinear setting. As we will see, the matrix $A$ in (3.7) and the perturbed

gradient mapping $\bar{g}$ will play an important role when connecting AA-R and GMRES. We further note that the linear case has been already covered in [35, Proposition 2].

PROPOSITION 3.8. *Let $\{x^k\}_{k=0,\ldots,m+1}$ be generated by Algorithm 2.1 and suppose that (A.4) is satisfied. Let the sequence $\{x_G^k\}_{k=0,\ldots,m}$ be generated by GMRES applied to $A(x - x^0) = -\nabla f(x^0)$ with initial point $x_G^0 = x^0$, where $A$ is defined as in (3.7). Suppose that the matrix $A$ is nonsingular. Setting $\hat{x}^0 := x^0$, it then holds that*

$$\bar{x}_G^k := \bar{g}(x_G^k) = x^{k+1} \quad and \quad x_G^k = \hat{x}^k \quad \forall\ k = 0, \ldots, m.$$

*In particular, we have $\kappa((\bar{X}_G^k)^\top \bar{X}_G^k) \leq M^2$ for each $k = 1, \ldots, m$, where $\bar{X}_G^k := [\bar{x}_G^0 - x^0, \ldots, \bar{x}_G^{k-1} - x^0]$.*

*Proof.* We prove Proposition 3.8 by induction. The base case $k = 0$ is obviously satisfied. Next, let us suppose that the induction hypothesis is true for any $i \leq k - 1$. By definition, we have $\hat{x}^k = x^0 + \sum_{i=1}^k \alpha_i^k (x^i - x^0)$ where $\alpha^k$ is the solution to

$$\min_\alpha\ \left\| h(x^0) + \sum_{i=1}^k \alpha_i (h(x^i) - h(x^0)) \right\|^2.$$

Based on the definition of the matrix $A$ (see (3.7)) and by (3.1), we further obtain:

$$(3.8) \qquad A(x^i - x^0) = \nabla f(x^i) - \nabla f(x^0) = -L(h(x^i) - h(x^0)) \quad \forall\ i = 0, \ldots, m.$$

Hence, $\alpha^k$ is also the solution to the problem $\min_\alpha \|A[\sum_{i=1}^k \alpha_i(x^i - x^0)] + \nabla f(x^0)\|^2$ and it holds that $\hat{x}^k = \operatorname{argmin}_{x \in x^0 + \operatorname{span}\{x^1 - x^0, \ldots, x^k - x^0\}} \|A(x - x^0) + \nabla f(x^0)\|^2$. In addition, using assumption (A.4), we can infer that the vectors $\{x^1 - x^0, \ldots, x^k - x^0\}$ are linearly independent. Applying Proposition 3.7, we have $x_G^k - x^0 \in \mathcal{K}^k(A, -\nabla f(x^0))$ and thus, it follows $\bar{x}_G^k - x^0 = x_G^k - x^0 - \frac{1}{L}(A(x_G^k - x^0) + \nabla f(x^0)) \in \mathcal{K}^{k+1}(A, -\nabla f(x^0))$ (for all $k$). Combining these observations and using the induction hypothesis, this yields $\operatorname{span}\{x^1 - x^0, \ldots, x^k - x^0\} = \operatorname{span}\{\bar{x}_G^1 - x^0, \ldots, \bar{x}_G^{k-1} - x^0\} = \mathcal{K}^k(A, -\nabla f(x^0))$, where the last equality follows from $\dim(\operatorname{span}\{x^1 - x^0, \ldots, x^k - x^0\}) = k$. Thus, by Proposition 3.7, we can deduce:

$$x_G^k = \operatorname{argmin}_{x \in x^0 + \mathcal{K}^k(A, -\nabla f(x^0))} \|A(x - x^0) + \nabla f(x^0)\|^2$$
$$= \operatorname{argmin}_{x \in x^0 + \operatorname{span}\{x^1 - x^0, \ldots, x^k - x^0\}} \|A(x - x^0) + \nabla f(x^0)\|^2 = \hat{x}^k$$

and thanks to (3.8), we obtain

$$\bar{x}_G^k = \bar{g}(x_G^k) = \bar{g}(\hat{x}^k) = \hat{x}^k - L^{-1} \cdot \sum_{i=1}^k \alpha_i^k A(x^i - x^0) + h(x^0) = \hat{g}^k = x^{k+1}.$$

The last assertion in Proposition 3.8 now follows from $\bar{x}_G^k = x^{k+1}$ and (A.4). □

**3.4. Connecting GMRES and CR.** In this section, we study GMRES and CR applied to the general linear systems

$$A(x - x^0) = b \quad and \quad B(x - x^0) = b, \quad A, B \in \mathbb{R}^{n \times n}, \quad b, x^0 \in \mathbb{R}^n.$$

Let $x_G^k$ denote the $k$-th iteration of GMRES applied to $A(x - x^0) = b$ and let $x_R^k$ denote the $k$-th iteration of CR applied to the linear system $B(x - x^0) = b$ (in our case, we will have $B = \nabla^2 f(x^0)$ and $b = -\nabla f(x^0)$). Here, we want to investigate and bound the distance between the iterates $x_G^k$ and $x_R^k$. In our analysis, we further assume $x_G^0 = x_R^0 = x^0$ and $b \neq 0$. We will largely utilize the following simple fact:

PROPOSITION 3.9. *Let $a_1, b_1, a_2, b_2$ be given scalars, vectors, or matrices with appropriate dimensions such that $a_1 b_1$, $a_2 b_2$, $a_1 - a_2$, $b_1 - b_2$, and $a_1 b_1 - a_2 b_2$ are well-defined. Then, we have $\|a_1 b_1 - a_2 b_2\| \leq \|a_1 - a_2\|\|b_1\| + \|b_1 - b_2\|\|a_2\|$.*

As usual, the concrete implementation of CR is not of our concern and we only require the following property of CR:

PROPOSITION 3.10. *[13, Section 2.2] Suppose $B \in \mathbb{R}^{n \times n}$ is symmetric and positive definite and let $b, x^0 \in \mathbb{R}^n$ be given. Let $\{x_R^k\}_k$ be generated by CR applied to the linear system $B(x - x^0) = b$ with $x_R^0 = x^0$. Then, we have:*

$$x_R^k = \operatorname{argmin}_{x \in x^0 + \mathcal{K}^k(B, b)} \|B(x - x^0) - b\|^2.$$

Based on our earlier discussion, we now introduce several additional terms:

$$(3.9) \quad \begin{aligned} \bar{x}_G^k &= x_G^k - \frac{1}{L}(A(x_G^k - x^0) - b), & \bar{x}_R^k &= x_R^k - \frac{1}{L}(B(x_R^k - x^0) - b), \\ \tilde{x}_G^k &= \bar{x}_G^k - \frac{1}{L}(A(\bar{x}_G^k - x^0) - b), & \tilde{x}_R^k &= \bar{x}_R^k - \frac{1}{L}(B(\bar{x}_R^k - x^0) - b), \\ \bar{X}_G^k &= [\bar{x}_G^0 - x^0, \ldots, \bar{x}_G^k - x^0], & \bar{X}_R^k &= [\bar{x}_R^0 - x^0, \ldots, \bar{x}_R^k - x^0]. \end{aligned}$$

Our first lemma in this section allows to connect the residuals of $\bar{x}_G^k$ and $x_G^k$.

LEMMA 3.11. *Let $B \in \mathbb{R}^{n \times n}$ be a symmetric, positive definite matrix with $L \geq \lambda_{\max}(B)$, $\lambda_{\min}(B) \geq \mu > 0$ and suppose that $A \in \mathbb{R}^{n \times n}$ satisfies $\|A - B\| < \mu$. Let $x, b \in \mathbb{R}^n$ be given and set $\bar{x} = x - L^{-1}(Ax - b)$. It holds that:*

$$\|A\bar{x} - b\| \leq \|Ax - b\|.$$

*Proof.* First notice that $\sigma_{\min}(A) \geq \sigma_{\min}(B) - \|A - B\| > 0$, which shows that $A$ is nonsingular. We can then define $x^* := A^{-1}b$ and rewrite $Ax - b = A(x - x^*)$. Furthermore, it holds that

$$\|A\bar{x} - b\| = \|A(\bar{x} - x^*)\| = \|(I - L^{-1}A)A(x - x^*)\| \leq \|I - L^{-1}A\|\|A(x - x^*)\|.$$

Hence, it suffices to verify $\|I - L^{-1}A\| \leq 1$. Indeed, we have $\|I - L^{-1}A\| \leq \|I - L^{-1}B\| + \|L^{-1}(A - B)\| \leq 1 - \frac{\mu}{L} + \frac{\mu}{L} \leq 1$ which finishes the proof. $\square$

Next, we present our main result of this subsection.

THEOREM 3.12. *Let $B \in \mathbb{R}^{n \times n}$ be a symmetric, positive definite matrix with $\lambda_{\max}(B) \leq L$ and $\lambda_{\min}(B) \geq \mu > 0$ and let $A \in \mathbb{R}^{n \times n}$, $b, x^0 \in \mathbb{R}^n$, and $\mathbb{N} \ni m \leq n$ be given. Let the sequences $\{x_G^k\}_k$ and $\{x_R^k\}_k$ be generated by GMRES and CR applied to the linear systems $A(x - x^0) = b$ and $B(x - x^0) = b$ with $x_G^0 = x_R^0 = x^0$, respectively. Suppose further that there are constants $C_1, C_2, C_3 > 0$ such that:*

  *(i) $\|A - B\| \leq C_1 \|b\|$.*
  *(ii) For each $k = 0, \ldots, m$, we have $\|\bar{x}_G^k - x_G^0\| \leq C_2 \|b\|$.*
  *(iii) We have $\kappa((\bar{X}_G^k)^\top \bar{X}_G^k) \leq C_3$ for all $k = 1, \ldots, m$.*
*There exists a constant $\epsilon_\sharp > 0$ such that if $\|b\| \leq \epsilon_\sharp$, then there is $C > 0$ such that:*

$$\|x_G^k - x_R^k\| \leq C\|b\|^2, \quad \|\bar{x}_G^k - \bar{x}_R^k\| \leq C\|b\|^2, \quad \|\tilde{x}_G^k - \tilde{x}_R^k\| \leq C\|b\|^2, \quad \forall\, 0 \leq k \leq m.$$

*Proof.* Without loss of generality, we can assume $b \neq 0$. We define the following quantities recursively: $\zeta_0 = 0$, $c_{k,1} = L^{-1}C_1C_2 + L^{-2}C_1 + \zeta_k$, $c_{k,2} = (k+1)C_1C_2 + L(\sum_{i=0}^k c_{i,1}^2)^{\frac{1}{2}}$, $c_{k,3} = 3c_{k,2}L(k+1)C_2$, $c_4 = \frac{16}{9\mu^2}C_3L^2$, $c_{k,5} = c_4 c_{k,2} + \frac{7}{2}c_4^2 c_{k,3}(k+1)C_2$,

$\zeta_{k+1} = \frac{7}{4}L(k+1)C_2 c_{k,5} + \sqrt{c_4}c_{k,2}$, $c_j = \max_{k=0,\ldots,m} c_{k,j}$, for all $j = 1,2,3,5$ and $\epsilon_\sharp := \frac{1}{2}\min\{(L^2 C_3 \sum_{i=0}^m c_{i,1}^2)^{-\frac{1}{2}}, \frac{\mu}{2C_1}, \frac{L}{c_2}, \frac{1}{c_3 c_4}\}$. Next, let us assume $\|b\| \le \epsilon_\sharp$. Due to $C_1 \epsilon_\sharp \le \frac{1}{4}\mu \le \frac{1}{4}L$, we then immediately obtain:

$$(3.10) \quad \|A\| \le \|B\| + \|A - B\| \le \tfrac{5}{4}L \quad \text{and} \quad \sigma_{\min}(A) \ge \lambda_{\min}(B) - \|A - B\| \ge \tfrac{3}{4}\mu.$$

This shows that $A$ is nonsingular. Our goal is now to establish $\|x_G^k - x_R^k\| \le \zeta_k \|b\|^2$ by induction. The base case $k = 0$ is trivial. Let us suppose that the induction hypothesis is true for all $0 \le i \le k$. By Proposition 3.7 and Proposition 3.10, we have:

$$x_G^{k+1} = \operatorname*{argmin}_{x \in x^0 + \mathcal{K}^{k+1}(A,b)} \|A(x - x^0) - b\|^2, \quad x_R^{k+1} = \operatorname*{argmin}_{x \in x^0 + \mathcal{K}^{k+1}(B,b)} \|B(x - x^0) - b\|^2.$$

Furthermore, mimicking the proof of Proposition 3.8 and using (iii), we can deduce

$$(3.11) \qquad \bar{x}_G^k - x^0 \in \mathcal{K}^{k+1}(A,b) \quad \text{and} \quad \bar{x}_R^k - x^0 \in \mathcal{K}^{k+1}(B,b)$$

for all $k$ and $\operatorname{span}\{\bar{x}_G^0 - x^0, \ldots, \bar{x}_G^k - x^0\} = \mathcal{K}^{k+1}(A,b)$. In addition, applying (i) and (ii), it holds that

$$
\begin{aligned}
\|\bar{x}_G^i - \bar{x}_R^i\| &= \|x_G^i - L^{-1}A(x_G^i - x^0) - x_R^i + L^{-1}B(x_R^i - x^0)\| \\
&= \|L^{-1}(A - B)(x^0 - x_G^i) + (I - L^{-1}B)(x_G^i - x_R^i)\| \\
&\le L^{-1}\|A - B\|(\|\bar{x}_G^i - x^0\| + \|\bar{x}_G^i - x_G^i\|) + \|I - L^{-1}B\|\|x_G^i - x_R^i\| \\
&\le C_1 L^{-1}\|b\|(\|\bar{x}_G^i - x^0\| + \|\bar{x}_G^i - x_G^i\|) + \|x_G^i - x_R^i\| \\
(3.12) \qquad &\le C_1 L^{-1}\|b\|(C_2\|b\| + L^{-1}\|b\|) + \zeta_i \|b\|^2 = c_{i,1}\|b\|^2,
\end{aligned}
$$

where we used $\|I - L^{-1}B\| \le 1 - \frac{\mu}{L} \le 1$ and Proposition 3.7 to show that:

$$\|\bar{x}_G^i - x_G^i\| = L^{-1}\min_{x \in x^0 + \mathcal{K}^i(A,b)}\|A(x - x^0) - b\| \le L^{-1}\|A(x^0 - x^0) - b\| = L^{-1}\|b\|.$$

Therefore, we are able to bound the error between $\bar{X}_G^k$ and $\bar{X}_R^k$:

$$\|\bar{X}_G^k - \bar{X}_R^k\| \le \left(\sum_{i=0}^k \|\bar{x}_G^i - \bar{x}_R^i\|^2\right)^{1/2} \le \left(\sum_{i=0}^k c_{i,1}^2\right)^{1/2}\|b\|^2 \le \frac{1}{2L\sqrt{C_3}}\|b\|,$$

where we applied the definition of $\epsilon_\sharp$. Furthermore, due to (iii), we can infer:

$$(3.13) \qquad \|\bar{X}_G^k\| \ge \|\bar{x}_G^0 - x^0\| = \frac{\|b\|}{L} \implies \sigma_{\min}(\bar{X}_G^k) \ge \frac{\sigma_{\max}(\bar{X}_G^k)}{\sqrt{C_3}} \ge \frac{\|b\|}{(L\sqrt{C_3})}.$$

Consequently, it holds that $\sigma_{\min}(\bar{X}_R^k) \ge \sigma_{\min}(\bar{X}_G^k) - \|\bar{X}_G^k - \bar{X}_R^k\| \ge \frac{1}{2L\sqrt{C_3}}\|b\| > 0$. Thus, the column vectors of $\bar{X}_R^k$ are also linearly independent and by (3.11), it follows $\operatorname{span}\{\bar{x}_R^0 - x^0, \ldots, \bar{x}_R^k - x^0\} = \mathcal{K}^{k+1}(B,b)$. Combining the previous arguments, we can now rewrite $x_G^{k+1}$ and $x_R^{k+1}$ as:

$$
\begin{aligned}
x_G^{k+1} &= \operatorname{argmin}_{x \in x^0 + \operatorname{span}\{\bar{x}_G^0 - x^0, \ldots, \bar{x}_G^k - x^0\}} \|A(x - x^0) - b\|^2, \\
x_R^{k+1} &= \operatorname{argmin}_{x \in x^0 + \operatorname{span}\{\bar{x}_R^0 - x^0, \ldots, \bar{x}_R^k - x^0\}} \|B(x - x^0) - b\|^2.
\end{aligned}
$$

The closed-form expressions of $x_G^{k+1}$ and $x_R^{k+1}$ are therefore given by:

$$x_G^{k+1} = x^0 + Y_G^k((Y_G^k)^\top Y_G^k)^{-1}(Y_G^k)^\top b, \quad x_R^{k+1} = x^0 + Y_R^k((Y_R^k)^\top Y_R^k)^{-1}(Y_R^k)^\top b$$

where $Y_G^k = A\bar{X}_G^k$ and $Y_R^k = B\bar{X}_R^k$. Our first task is to estimate the error between $Y_G^k$ and $Y_R^k$. Using (i) and (ii), we have:

$$\|Y_G^k - Y_R^k\| \leq \|(A - B)\bar{X}_G^k\| + \|B(\bar{X}_G^k - \bar{X}_R^k)\|$$

$$\leq (k+1)C_1\|b\|\|\bar{X}_G^k\|_\infty + L\left(\sum_{i=0}^k c_{i,1}^2\right)^{1/2}\|b\|^2 \leq c_{k,2}\|b\|^2 \leq c_2\|b\|^2.$$

Moreover, applying (3.10), we can infer $\|Y_G^k\| \leq \|A\|\|\bar{X}_G^k\| \leq \frac{5L}{4}\|\bar{X}_G^k\| \leq \frac{5L(k+1)C_2}{4}\|b\|$ and due to $c_2\epsilon_\sharp \leq \frac{L}{2} \leq \frac{L(k+1)}{2}$, we have $c_2\|b\|^2 \leq \frac{L(k+1)C_2}{2}\|b\|$. This allows to establish a bound for the norm of $Y_R^k$:

$$\|Y_R^k\| \leq \|Y_G^k\| + \|Y_G^k - Y_R^k\| \leq \tfrac{7}{4}L(k+1)C_2\|b\|.$$

Therefore, by Proposition 3.9, we can bound the norm of the term $(Y_G^k)^\top Y_G^k - (Y_R^k)^\top Y_R^k$ as follows:

$$\|(Y_G^k)^\top Y_G^k - (Y_R^k)^\top Y_R^k\| \leq \|Y_G^k - Y_R^k\|(\|Y_G^k\| + \|Y_R^k\|)$$

$$\leq c_{k,2}\|b\|^2\left(\tfrac{7}{4}L(k+1)C_2 + \tfrac{5}{4}L(k+1)C_2\right)\|b\| = 3c_{k,2}L(k+1)C_2\|b\|^3 = c_{k,3}\|b\|^3.$$

Our next task is to bound $\|((Y_G^k)^\top Y_G^k)^{-1}\|$. Using $A^\top A \succeq \frac{9}{16}\mu^2 I$ and (3.13), we have:

$$\|((Y_G^k)^\top Y_G^k)^{-1}\| = \|((\bar{X}_G^k)^\top A^\top A\bar{X}_G^k)^{-1}\|$$

$$\leq \frac{16}{9\mu^2}\|((\bar{X}_G^k)^\top \bar{X}_G^k)^{-1}\| \leq \frac{16C_3}{9\mu^2\|\bar{X}_G^k\|^2} \leq \frac{16C_3L^2}{9\mu^2\|b\|^2} = \frac{c_4}{\|b\|^2},$$

Since $\epsilon_\sharp$ is chosen such that $1 - c_{k,3}c_4\|b\| \geq \frac{1}{2} > 0$, we can now apply Banach's perturbation lemma, see, e.g., [16, Theorem 2.3.4], which implies:

$$\|((Y_G^k)^\top Y_G^k)^{-1} - ((Y_R^k)^\top Y_R^k)^{-1}\| \leq \frac{\|((Y_G^k)^\top Y_G^k)^{-1}\|^2\|(Y_G^k)^\top Y_G^k - (Y_R^k)^\top Y_R^k\|}{1 - \|((Y_G^k)^\top Y_G^k)^{-1}\|\|(Y_G^k)^\top Y_G^k - (Y_R^k)^\top Y_R^k\|}$$

$$\leq \frac{c_4^2 c_{k,3}}{\|b\| - c_4 c_{k,3}\|b\|^2} \leq \frac{2c_4^2 c_{k,3}}{\|b\|}.$$

Consequently, applying Proposition 3.9, it follows:

$$\|Y_G^k((Y_G^k)^\top Y_G^k)^{-1} - Y_R^k((Y_R^k)^\top Y_R^k)^{-1}\|$$

$$\leq \|((Y_G^k)^\top Y_G^k)^{-1}\|\|Y_G^k - Y_R^k\| + \|((Y_G^k)^\top Y_G^k)^{-1} - ((Y_R^k)^\top Y_R^k)^{-1}\|\|Y_R^k\|$$

$$\leq c_4 c_{k,2} + \tfrac{7}{2}c_4^2 c_{k,3}(k+1)C_2 = c_{k,5}$$

and $\|Y_G^k((Y_G^k)^\top Y_G^k)^{-1}\| = \sqrt{\|((Y_G^k)^\top Y_G^k)^{-1}\|} \leq \frac{\sqrt{c_4}}{\|b\|}$. Altogether, this yields:

$$\|x_G^{k+1} - x_R^{k+1}\| = \|Y_G^k((Y_G^k)^\top Y_G^k)^{-1}(Y_G^k)^\top b - Y_R^k((Y_R^k)^\top Y_R^k)^{-1}(Y_R^k)^\top b\|$$

$$\leq \|Y_G^k((Y_G^k)^\top Y_G^k)^{-1} - Y_R^k((Y_R^k)^\top Y_R^k)^{-1}\|\|Y_R^k\|\|b\| + \|Y_G^k((Y_G^k)^\top Y_G^k)^{-1}\|\|Y_R^k - Y_G^k\|\|b\|$$

$$\leq \tfrac{7}{4}L(k+1)C_2 c_{k,5}\|b\|^2 + \sqrt{c_4}c_{k,2}\|b\|^2 = \zeta_{k+1}\|b\|^2.$$

This shows $\|x_G^k - x_R^k\| \leq \zeta_k\|b\|^2$ by induction. Mimicking (3.12), we now obtain:

$$\|\tilde{x}_G^k - \tilde{x}_R^k\| \leq \|L^{-1}(A - B)(\bar{x}_G^k - x^0)\| + \|(I - L^{-1}B)(\bar{x}_G^k - \bar{x}_R^k)\|$$

$$\leq L^{-1}C_1C_2\|b\|^2 + \|\bar{x}_G^k - \bar{x}_R^k\| \leq (L^{-1}C_1C_2 + c_1)\|b\|^2.$$

Therefore, it suffices to choose $C := \max\{\max_{k=0,\ldots,m}\zeta_k, L^{-1}C_1C_2 + c_1\}$. $\qquad\square$

**3.5. Connecting CR and CG.** In this subsection, we assume that the matrix $B$ is symmetric and positive definite. Suppose we apply CR to the linear system $B(x - x^0) = b$ starting at $x^0$. Then, by Proposition 3.10, we have:

$$x_R^k = \operatorname{argmin}_{x \in x^0 + \mathcal{K}^k(B,b)} \|B(x - x^0) - b\|^2 = \operatorname{argmin}_{x \in x^0 + \mathcal{K}^k(B,b)} \|B(x - x^*)\|,$$

where $x^* := B^{-1}b + x^0$ is the optimal solution of the linear system $B(x - x^0) = b$. Next, for $k = 0, \ldots, m$, we set $y^* := B^{\frac{1}{2}}x^*$, $\bar{b} := B^{\frac{1}{2}}b$ and $y^k := B^{\frac{1}{2}}x_R^k$. Then, by definition, we obtain:

$$y^k = \operatorname{argmin}_{y \in y^0 + \mathcal{K}^k(B,\bar{b})} \ (y - y^*)^\top B(y - y^*).$$

According to [18, Theorem 2], this means that $y^k$ coincides with the $k$-th iteration of CG applied to the linear system $B(y - y^*) = 0$ with initial value $y^0 = B^{\frac{1}{2}}x^0$. Moreover, in this case, it follows $(x_R^k - x^*)^\top B(x_R^k - x^*) = \|y^k - y^*\|^2$. Based on this observation and connection between the CR- and CG-iterates, we now want to apply classical techniques for CG, [20], to study the behavior of the distance $\|y^k - y^*\|$ as the iteration $k$ increases. Our goal is to then transfer the obtained results back to CR and AA-R. As usual, we define the following terms: $\bar{y}^k = y^k - L^{-1}B(y^k - y^*)$,

$$\tilde{y}^k = \bar{y}^k - L^{-1}B(\bar{y}^k - y^*), \quad \text{and} \quad \psi(y) = \frac{1}{2}(y - y^*)^\top B(y - y^*).$$

Notice that the introduced linear transformations also preserve the latter gradient descent steps, i.e., it holds that $\bar{y}^k = B^{\frac{1}{2}}\bar{x}_R^k$ and $\tilde{y}^k = B^{\frac{1}{2}}\tilde{x}_R^k$. Here, the point $\bar{y}^k$ is obtained by applying one gradient step (for the objective function $\psi$) with stepsize $L^{-1}$ on the CG-iteration $y^k$ and $\tilde{y}^k$ results from applying two gradient steps with stepsize $L^{-1}$ on $y^k$. Next, we collect several results from [20] for convenience and to fix the notations. The full CG algorithm is shown in Algorithm 3.1.

---

**Algorithm 3.1** CG for the linear system $B(y - y^*) = 0$.

---

1: Choose an initial point $y^0 \in \mathbb{R}^n$ and set $p^0 = r^0 = -B(y^0 - y^*)$.
2: **for** $i = 0, 1, \ldots, n$ **do**
3:     **if** $\|r^i\| = 0$ **then** Break; **end if**
4:     $a_i = \frac{\|r^i\|^2}{\langle p^i, Bp^i \rangle}$.
5:     $y^{i+1} = y^i + a_i p^i$.
6:     $r^{i+1} = r^i - a_i Bp^i$.
7:     $b_i = \frac{\|r^{i+1}\|^2}{\|r^i\|^2}$.
8:     $p^{i+1} = r^{i+1} + b_i p^i$.
9: **end for**

---

PROPOSITION 3.13. *Let the sequence $\{y^k\}_k$ be generated by CG and let $y^{(k)}$ denote the projection of $y^*$ onto the affine space $\mathcal{A}^k := y^0 + \operatorname{span}\{y^1 - y^0, \ldots, y^k - y^0\}$. Then, the following properties are satisfied:*

(i) *([20, Theorem 6.5]) $y^{(k+1)} = y^{k+1} + \frac{2\psi(y^{k+1})}{\|r^k\|^2} p^k$.*

(ii) *([20, Equation (5:3a)]) For all $i \neq j$: $\langle r^i, r^j \rangle = 0$.*

(iii) *([20, Equation (5:3c)]) For all $i < j$, we have $\langle p^i, r^j \rangle = 0$ and for all $i \geq j$, it holds that $\langle p^i, r^j \rangle = \|r^i\|^2$.*

(iv) *([20, Equation (5:6b)]) For all $0 \leq i \leq n$: $\|p^i\|^2 = \|r^i\|^4 \sum_{j=0}^{i} \frac{1}{\|r^j\|^2}$.*

(v) ([20, Equation (5:11)]) For all $0 \leq i \leq n-1$: $\langle r^{i+1}, Br^i \rangle = \langle r^{i+1}, Bp^i \rangle = -a_i^{-1} \|r^{i+1}\|^2$.

(vi) ([20, Equation (5:6a)]) For all $0 \leq i \leq j \leq n$: $\langle p^i, p^j \rangle = \frac{\|r^j\|^2 \|p^i\|^2}{\|r^i\|^2}$.

(vii) ([20, Equation (5:4b)]) For all $i \neq j$: $\langle p^i, Bp^j \rangle = 0$.

(viii) ([20, Equation (5:3d)]) For all $i \neq j, i \neq j+1$: $\langle r^i, Bp^j \rangle = 0$.

(ix) ([20, Equation (5:12)]) We have $a_0 = \frac{\|r^0\|^2}{\langle r^0, Br^0 \rangle}$ and $\frac{\|p^i\|^2}{\langle p^i, Bp^i \rangle} > a_i > \frac{\|r^i\|^2}{\langle r^i, Br^i \rangle}$ for all $i > 0$.

(x) ([20, Equation (5:8b)]) For all $i \geq 1$: $r^{i+1} = (1 + b'_{i-1})r^i - a_i Br^i - b'_{i-1} r^{i-1}$, where $b'_{i-1} = \frac{a_i}{a_{i-1}} b_{i-1} = \frac{a_i}{a_{i-1}} \frac{\|r^i\|^2}{\|r^{i-1}\|^2}$.

The properties stated in Proposition 3.13 will be referred to as Property (i)–(x) in the following. Before studying the convergence behavior of CG in terms of $\|\bar{y}^k - y^*\|$ and $\|\tilde{y}^{k-1} - y^*\|$, let us briefly discuss our underlying motivation.

Our final aim is to prove $f(x^{k+1}) \leq f(g(x^k))$ (including some potential higher-order error terms), where $\{x^k\}_k$ is generated by AA-R. Applying Proposition 3.8, the AA step $x^{k+1}$ is equal to $\bar{x}_G^k$, which is close to $\bar{x}_R^k$. On the other hand, we have $g(x^k) = \bar{g}(x^k) = \bar{g}(\bar{x}_G^{k-1}) = \tilde{x}_G^{k-1}$, which is close to $\tilde{x}_R^{k-1}$. Hence, up to certain error terms, the descent condition "$f(x^{k+1}) \leq f(g(x^k))$" can now be formulated as follows:

$$f(\bar{x}_R^k) \leq f(\tilde{x}_R^{k-1}).$$

Expanding $f$ at $x^0$ (and again ignoring higher-order error terms), this can be further rewritten as:

$$\nabla f(x^0)^\top (\bar{x}_R^k - x^0) + \frac{1}{2}(\bar{x}_R^k - x^0)^\top \nabla^2 f(x^0)(\bar{x}_R^k - x^0)$$
$$\leq \nabla f(x^0)^\top (\tilde{x}_R^{k-1} - x^0) + \frac{1}{2}(\tilde{x}_R^{k-1} - x^0)^\top \nabla f(x^0)(\tilde{x}_R^{k-1} - x^0).$$

Noticing $B = \nabla^2 f(x^0)$, $b = -\nabla f(x^0)$, and $x^* = B^{-1}b + x^0$, this is equivalent to

$$(\bar{x}_R^k - x^*)^\top B(\bar{x}_R^k - x^*) \leq (\tilde{x}_R^{k-1} - x^*)^\top B(\tilde{x}_R^{k-1} - x^*),$$

which, by the previously introduced transformation, can be expressed as $\|\bar{y}^k - y^*\|^2 \leq \|\tilde{y}^{k-1} - y^*\|^2$. This is exactly what we want to show in Theorem 3.14. We note that the proof of Theorem 3.14 would be significantly easier if the stepsize in the gradient mapping $g$ is sufficiently small (potentially much smaller than $L^{-1}$). Here, we provide a general result covering the core case $g(x) = x - \frac{1}{L}\nabla f(x)$.

THEOREM 3.14. *Suppose that $\{y^k\}_k$ is generated by* CG *applied to the linear system $B(y - y^*) = 0$, where $B \in \mathbb{R}^{n \times n}$ is symmetric, positive definite with $\nu := \frac{L}{\|B\|} \geq 1$. Then, we have:*

$$(3.14) \quad \|\bar{y}^k - y^*\|^2 + \left[2\nu + \frac{1}{\nu^2} - 3\right]\frac{\|r^k\|^2}{L^2} + \left[\nu + \frac{1}{\nu} - 2\right]^2 \frac{\|r^{k-1}\|^2}{L^2} \leq \|\tilde{y}^{k-1} - y^*\|^2.$$

*Proof.* First, by [18, Equation (21)] and [26, Theorem 5.3], we have $\mathcal{A}^k = y^0 + \mathcal{K}^k(B, r^0)$ and $y^k \in y^0 + \mathcal{K}^k(B, r^0)$. Hence, both $\bar{y}^k$ and $\tilde{y}^{k-1}$ belong to the affine space $y^0 + \mathcal{K}^{k+1}(B, r^0) = \mathcal{A}^{k+1}$. Furthermore, by the definition of $y^{(k+1)}$, we can derive the following decomposition properties:

$$(3.15) \quad \begin{aligned} \|\bar{y}^k - y^*\|^2 &= \|y^{(k+1)} - \bar{y}^k\|^2 + \|y^{(k+1)} - y^*\|^2, \\ \|\tilde{y}^{k-1} - y^*\|^2 &= \|y^{(k+1)} - \tilde{y}^{k-1}\|^2 + \|y^{(k+1)} - y^*\|^2. \end{aligned}$$

Therefore, it holds that:

$$\|\tilde{y}^{k-1} - y^*\|^2 - \|\bar{y}^k - y^*\|^2 = \|y^{(k+1)} - \tilde{y}^{k-1}\|^2 - \|y^{(k+1)} - \bar{y}^k\|^2$$

(3.16)
$$= \|\tilde{y}^{k-1} - \bar{y}^k\|^2 + 2\langle \tilde{y}^{k-1} - \bar{y}^k, \bar{y}^k - y^{(k+1)}\rangle.$$

Using Property (i) and the definition of the CG-step, we have $y^{(k+1)} = y^{k+1} + \frac{2\psi(y^{k+1})}{\|r^k\|^2}p^k$ and $y^{k+1} = y^k + a_k p^k$. Consequently, setting $\gamma_k = 2\psi(y^{k+1})/\|r^k\|^2 + a_k$ and applying $r^k = -B(y^k - y^*)$, we obtain

(3.17) $$y^{(k+1)} - \bar{y}^k = \left[\frac{2\psi(y^{k+1})}{\|r^k\|^2} + a_k\right]p^k + [y^k - \bar{y}^k] = \gamma_k p^k - \frac{1}{L}r^k.$$

Moreover, we have $y^k - \bar{y}^{k-1} = a_{k-1}p^{k-1} - \frac{1}{L}r^{k-1}$ and

(3.18) $$\bar{y}^k - \tilde{y}^{k-1} = (I - L^{-1}B)(y^k - \bar{y}^{k-1}) = (I - L^{-1}B)(a_{k-1}p^{k-1} - L^{-1}r^{k-1}).$$

We now consider the first term in (3.16):

$$\|\tilde{y}^{k-1} - \bar{y}^k\|^2 = \|(I - L^{-1}B)(a_{k-1}p^{k-1} - L^{-1}r^{k-1})\|^2 = T_1 - 2L^{-1}T_2 + L^{-2}T_3,$$

where $T_1 = a_{k-1}^2\|(I - L^{-1}B)p^{k-1}\|^2$, $T_2 = \langle a_{k-1}(I - L^{-1}B)p^{k-1}, (I - L^{-1}B)r^{k-1}\rangle$, and $T_3 = \|(I - L^{-1}B)r^{k-1}\|^2$. The update rule for $r^k$ yields

(3.19) $$a_{k-1}Bp^{k-1} = r^{k-1} - r^k.$$

We first expand the term $T_1$:

$$T_1 = \|a_{k-1}p^{k-1} - L^{-1}(r^{k-1} - r^k)\|^2$$
$$= a_{k-1}^2\|p^{k-1}\|^2 - 2a_{k-1}L^{-1}\langle p^{k-1}, r^{k-1} - r^k\rangle + L^{-2}\|r^{k-1} - r^k\|^2.$$

Applying Property (ii) and (iii), it holds that:

$$\|r^{k-1} - r^k\|^2 = \|r^{k-1}\|^2 + \|r^k\|^2, \quad \langle p^{k-1}, r^{k-1} - r^k\rangle = \|r^{k-1}\|^2,$$

and thus, it follows $T_1 = a_{k-1}^2\|p^{k-1}\|^2 - \frac{2a_{k-1}}{L}\|r^{k-1}\|^2 + \frac{1}{L^2}(\|r^{k-1}\|^2 + \|r^k\|^2)$. Next, we estimate the term $T_2$:

$$T_2 = \langle a_{k-1}p^{k-1}, r^{k-1}\rangle - 2L^{-1}\langle a_{k-1}Bp^{k-1}, r^{k-1}\rangle + L^{-2}\langle a_{k-1}Bp^{k-1}, Br^{k-1}\rangle.$$

By Property (iii), we have $\langle a_{k-1}p^{k-1}, r^{k-1}\rangle = a_{k-1}\|r^{k-1}\|^2$. Furthermore, applying (3.19) and Property (ii), we obtain $\langle a_{k-1}Bp^{k-1}, r^{k-1}\rangle = \langle r^{k-1} - r^k, r^{k-1}\rangle = \|r^{k-1}\|^2$ and $\langle a_{k-1}Bp^{k-1}, Br^{k-1}\rangle = \langle r^{k-1} - r^k, Br^{k-1}\rangle$. Utilizing Property (v), we can infer:

$$\langle a_{k-1}Bp^{k-1}, Br^{k-1}\rangle = \langle r^{k-1} - r^k, Br^{k-1}\rangle = \|r^{k-1}\|_B^2 + a_{k-1}^{-1}\|r^k\|^2.$$

Substituting these expressions yields $T_2 = a_{k-1}\|r^{k-1}\|^2 - \frac{2}{L}\|r^{k-1}\|^2 + \frac{1}{L^2}\|r^{k-1}\|_B^2 + \frac{1}{L^2 a_{k-1}}\|r^k\|^2$. Finally, let us consider the term $T_3$; we have:

$$T_3 = \|r^{k-1}\|^2 - 2L^{-1}\|r^{k-1}\|_B^2 + L^{-2}\|Br^{k-1}\|^2.$$

Together, this establishes the following representation of $\|\tilde{y}^{k-1} - \bar{y}^k\|^2$:

$$(3.20) \quad \|\tilde{y}^{k-1} - \bar{y}^k\|^2 = a_{k-1}^2\|p^{k-1}\|^2 + \left[\frac{6}{L^2} - \frac{4a_{k-1}}{L}\right]\|r^{k-1}\|^2$$

$$- \frac{4}{L^3}\|r^{k-1}\|_B^2 + \frac{1}{L^4}\|Br^{k-1}\|^2 + \left[\frac{1}{L^2} - \frac{2}{L^3 a_{k-1}}\right]\|r^k\|^2.$$

We continue with the inner product term $\langle \tilde{y}^{k-1} - \bar{y}^k, \bar{y}^k - y^{(k+1)}\rangle$. By (3.17) and (3.18), we have:

$$\langle \tilde{y}^{k-1} - \bar{y}^k, \bar{y}^k - y^{(k+1)}\rangle = \langle (I - L^{-1}B)(a_{k-1}p^{k-1} - L^{-1}r^{k-1}), \gamma^k p^k - L^{-1}r^k\rangle$$

$$= Q_1 - L^{-1}Q_2,$$

where $Q_1 = \langle a_{k-1}p^{k-1} - \frac{1}{L}r^{k-1}, \gamma_k p^k - \frac{1}{L}r^k\rangle$ and $Q_2 = \langle a_{k-1}Bp^{k-1} - \frac{1}{L}Br^{k-1}, \gamma_k p^k - \frac{1}{L}r^k\rangle$. Applying Property (vi), (ii) and (iii), it holds that:

$$\langle p^{k-1}, p^k\rangle = \frac{\|r^k\|^2\|p^{k-1}\|^2}{\|r^{k-1}\|^2}, \quad \langle r^{k-1}, p^k\rangle = \|r^k\|^2, \quad \langle p^{k-1}, r^k\rangle = \langle r^k, r^{k-1}\rangle = 0,$$

which implies $Q_1 = a_{k-1}\gamma_k\langle p^{k-1}, p^k\rangle - \frac{\gamma_k}{L}\langle r^{k-1}, p^k\rangle - \frac{a_{k-1}}{L}\langle p^{k-1}, r^k\rangle + \frac{1}{L^2}\langle r^k, r^{k-1}\rangle = a_{k-1}\gamma_k\frac{\|r^k\|^2\|p^{k-1}\|^2}{\|r^{k-1}\|^2} - \frac{\gamma_k}{L}\|r^k\|^2$. Similarly, we can expand $Q_2$ as follows:

$$Q_2 = a_{k-1}\gamma_k\langle Bp^{k-1}, p^k\rangle - \frac{\gamma_k}{L}\langle Br^{k-1}, p^k\rangle - \frac{a_{k-1}}{L}\langle Bp^{k-1}, r^k\rangle + \frac{1}{L^2}\langle Br^k, r^{k-1}\rangle.$$

Applying Property (vii), (viii), (ii), (iii), (v), and (3.19), we can infer $\langle Bp^{k-1}, p^k\rangle = 0$, $\langle Br^{k-1}, p^k\rangle = 0$, $\langle a_{k-1}Bp^{k-1}, r^k\rangle = \langle r^{k-1} - r^k, r^k\rangle = -\|r^k\|^2$, and $\langle Br^k, r^{k-1}\rangle = -\frac{\|r^k\|^2}{a_{k-1}}$, which yields $Q_2 = \frac{1}{L}(1 - \frac{1}{La_{k-1}})\|r^k\|^2$. Therefore, the inner product term $\langle \tilde{y}^{k-1} - \bar{y}^k, \bar{y}^k - y^{(k+1)}\rangle$ is given by:

$$\langle \tilde{y}^{k-1} - \bar{y}^k, \bar{y}^k - y^{(k+1)}\rangle = a_{k-1}\gamma_k\frac{\|r^k\|^2\|p^{k-1}\|^2}{\|r^{k-1}\|^2} - \left[\frac{\gamma_k}{L} + \frac{1}{L^2} - \frac{1}{L^3 a_{k-1}}\right]\|r^k\|^2.$$

Summing (3.20) and the previous expression, we obtain:

$$\|\tilde{y}^{k-1} - \bar{y}^k\|^2 + 2\langle \tilde{y}^{k-1} - \bar{y}^k, \bar{y}^k - y^{(k+1)}\rangle$$

$$= a_{k-1}^2\|p^{k-1}\|^2 + \left[\frac{6}{L^2} - \frac{4a_{k-1}}{L}\right]\|r^{k-1}\|^2 - \frac{4}{L^3}\|r^{k-1}\|_B^2 + \frac{1}{L^4}\|Br^{k-1}\|^2$$

$$(3.21) \qquad + \left[2a_{k-1}\gamma_k\frac{\|r^k\|^2\|p^{k-1}\|^2}{\|r^{k-1}\|^2} - \frac{2\gamma_k}{L}\|r^k\|^2 - \frac{1}{L^2}\|r^k\|^2\right].$$

We continue with two sub-cases.

**Case 1:** $k = 1$. Using the fact $r^0 = p^0$, Property (ii), and (3.19), it follows $\|Br^0\|^2 = \|Bp^0\|^2 = a_0^{-2}\|r^0 - r^1\|^2 = a_0^{-2}(\|r^0\|^2 + \|r^1\|^2)$ and $\langle r^0, Br^0\rangle = \langle r^0, Bp^0\rangle = a_0^{-1}\langle r^0, r^0 - r^1\rangle = a_0^{-1}\|r^0\|^2$. Using these two equalities, we can simplify (3.21) to:

$$\|\tilde{y}^0 - \bar{y}^1\|^2 + 2\langle \tilde{y}^0 - \bar{y}^1, \bar{y}^1 - y^{(2)}\rangle$$

$$= a_0^2\|r^0\|^2 + \left[\frac{6}{L^2} - \frac{4a_0}{L}\right]\|r^0\|^2 - \frac{4}{L^3 a_0}\|r^0\|^2 + \frac{1}{L^4 a_0^2}(\|r^0\|^2 + \|r^1\|^2)$$

$$+ \left[2a_0\gamma_1\|r^1\|^2 - 2\gamma_1 L^{-1}\|r^1\|^2 - L^{-2}\|r^1\|^2\right]$$

$$= L^{-2}(Q_3\|r^0\|^2 + Q_4\|r^1\|^2),$$

where $Q_3$ and $Q_4$ are defined as $Q_3 = a_0^2 L^2 + 6 - 4a_0 L - 4(a_0 L)^{-1} + (a_0 L)^{-2}$ and $Q_4 = 2a_0\gamma_1 L^2 - 2\gamma_1 L - 1 + (a_0 L)^{-2}$. By Property (ix), we have $a_0 = \|r^0\|^2/\langle r^0, Br^0 \rangle \geq \frac{1}{\|B\|} \geq \frac{1}{L}$ and $a_1 > \|r^1\|^2/\langle r^1, Br^1 \rangle \geq \frac{1}{\|B\|} \geq \frac{1}{L}$. Hence, by the definition of $\gamma_1$, we can infer $\gamma_1 L \geq a_1 L > L\|B\|^{-1} = \nu \geq 1$ and $a_0 L \geq \nu \geq 1$. Therefore, it holds that:

$$Q_4 = 2\gamma_1 L(a_0 L - 1) - 1 + (a_0 L)^{-2} \geq 2a_0 L + (a_0 L)^{-2} - 3 \geq 2\nu + \nu^{-2} - 3 \geq 0,$$

where – in the last equality – we used the fact that the function $x \mapsto 2x + x^{-2}$ is monotonically increasing for $x \geq 1$. Concerning $Q_3$, we notice:

$$Q_3 = (a_0 L + (a_0 L)^{-1} - 2)^2 \geq 0.$$

Since $x \mapsto x + \frac{1}{x} - 2$ is monotonically increasing and nonnegative for $x \in [1, \infty)$, we can further infer $Q_3 = (a_0 L + (a_0 L)^{-1} - 2)^2 \geq (\nu + \nu^{-1} - 2)^2$, Together, we obtain $\|\bar{y}^1 - y^*\|^2 + (\nu + \nu^{-1} - 2)^2 \frac{\|r^0\|^2}{L^2} + (2\nu + \nu^{-2} - 3)\frac{\|r^1\|^2}{L^2} \leq \|\tilde{y}^0 - y^*\|^2$.

**Case 2:** $k \geq 2$. We first utilize Property (x): $a_{k-1}Br^{k-1} = (1 + b'_{k-2})r^{k-1} - r^k - b'_{k-2}r^{k-2}$. Along with Property (ii), this allows to calculate $\|r^{k-1}\|_B^2$ and $\|Br^k\|^2$:

$$\langle r^{k-1}, Br^{k-1}\rangle = \frac{1}{a_{k-1}}\langle r^{k-1}, (1 + b'_{k-2})r^{k-1} - r^k - b'_{k-2}r^{k-2}\rangle = \frac{1 + b'_{k-2}}{a_{k-1}}\|r^{k-1}\|^2,$$

$$\|Br^{k-1}\|^2 = \langle Br^{k-1}, Br^{k-1}\rangle = \frac{(1 + b'_{k-2})^2}{a_{k-1}^2}\|r^{k-1}\|^2 + \frac{1}{a_{k-1}^2}\|r^k\|^2 + \frac{(b'_{k-2})^2}{a_{k-1}^2}\|r^{k-2}\|^2.$$

Therefore, the term (3.21) can be decomposed as follows: $\|\tilde{y}^{k-1} - \bar{y}^k\|^2 + 2\langle \tilde{y}^{k-1} - \bar{y}^k, \bar{y}^k - y^{(k+1)}\rangle = Q_5 + Q_6$, where

$$Q_5 = a_{k-1}^2\|p^{k-1}\|^2 + \left[\frac{6}{L^2} - \frac{4a_{k-1}}{L} - \frac{4(1+b'_{k-2})}{L^3 a_{k-1}} + \frac{(1+b'_{k-2})^2}{L^4 a_{k-1}^2}\right]\|r^{k-1}\|^2 + \frac{(b'_{k-2})^2}{a_{k-1}^2 L^4}\|r^{k-2}\|^2$$

$$Q_6 = 2a_{k-1}\gamma_k\frac{\|r^k\|^2\|p^{k-1}\|^2}{\|r^{k-1}\|^2} - \frac{2\gamma_k}{L}\|r^k\|^2 - \frac{1}{L^2}\|r^k\|^2 + \frac{1}{L^4 a_{k-1}^2}\|r^k\|^2.$$

We start with bounding $Q_6$. First, by Property (iv), it holds that:

$$2a_{k-1}\gamma_k\frac{\|r^k\|^2\|p^{k-1}\|^2}{\|r^{k-1}\|^2} \geq 2a_{k-1}\gamma_k\frac{\|r^k\|^2\|r^{k-1}\|^2}{\|r^{k-1}\|^2} = 2a_{k-1}\gamma_k\|r^k\|^2.$$

Thus, we have $Q_6 \geq \frac{1}{L^2}[2a_{k-1}\gamma_k L^2 - 2\gamma_k L - 1 + (a_{k-1}L)^{-2}]\|r^k\|^2$. The coefficient in the parentheses can be shown to be larger or equal than $2\nu + \nu^{-2} - 3$ by using the same strategy as in **Case 1** for $Q_4$. This yields $Q_6 \geq (2\nu + \nu^{-2} - 3)L^{-2}\|r^k\|^2$. Next, recalling the definition of $b'_{k-2}$ (see Property (x)), we obtain:

$$b'_{k-2} = \frac{a_{k-1}}{a_{k-2}}\frac{\|r^{k-1}\|^2}{\|r^{k-2}\|^2} \quad\Longrightarrow\quad \frac{(b'_{k-2})^2}{a_{k-1}^2 L^4}\|r^{k-2}\|^2 = \frac{\|r^{k-1}\|^2}{a_{k-2}^2 L^4\|r^{k-2}\|^2}\|r^{k-1}\|^2.$$

In addition, by Property (iv), it follows:

$$a_{k-1}^2\|p^{k-1}\|^2 = a_{k-1}^2\|r^{k-1}\|^4\sum_{j=0}^{k-1}\frac{1}{\|r^j\|^2} \geq a_{k-1}^2\|r^{k-1}\|^2 + a_{k-1}^2\frac{\|r^{k-1}\|^2}{\|r^{k-2}\|^2}\|r^{k-1}\|^2.$$

Setting $Q_7 = a_{k-1}^2 L^2 + 6 - 4a_{k-1}L - 4(a_{k-1}L)^{-1} + (a_{k-1}L)^{-2} = (a_{k-1}L + (a_{k-1}L)^{-1} - 2)^2$ and using the previous inequalities, we can lower bound $Q_5$ by:

$$Q_5 \geq \left[Q_7 - \frac{4b'_{k-2}}{a_{k-1}L} + \frac{(1+b'_{k-2})^2 - 1}{a_{k-1}^2 L^2} + \left[a_{k-1}^2 L^2 + \frac{1}{a_{k-2}^2 L^2}\right]\frac{\|r^{k-1}\|^2}{\|r^{k-2}\|^2}\right]\frac{\|r^{k-1}\|^2}{L^2}.$$

Let us denote the term in parentheses by $Q_8$. It suffices to show that $Q_8$ is nonnegative. In particular, it holds that:

$$Q_8 \geq Q_7 + \left[a_{k-1}L - \frac{1}{a_{k-2}L}\right]^2 \frac{\|r^{k-1}\|^2}{\|r^{k-2}\|^2} + \frac{2a_{k-1}}{a_{k-2}} \frac{\|r^{k-1}\|^2}{\|r^{k-2}\|^2} - \frac{4b'_{k-2}}{a_{k-1}L} + \frac{2b'_{k-2}}{a_{k-1}^2 L^2}$$

$$\geq Q_7 + 2b'_{k-2} - \frac{4b'_{k-2}}{a_{k-1}L} + \frac{2b'_{k-2}}{a_{k-1}^2 L^2} \geq \left[\nu + \nu^{-1} - 2\right]^2 + 2b'_{k-2}\left[1 - (a_{k-1}L)^{-1}\right]^2,$$

where we again used $a_{k-1}L \geq \nu$. This finally establishes (3.14), which concludes the proof of Theorem 3.14. $\qquad\square$

The previous result shows that performing two gradient steps on $y^{k-1}$ achieves less progress in terms of the distance to the optimal solution compared to performing one gradient step on $y^k$. In fact, we are able to prove a similar result for CG, which is of independent interest. More precisely, CG can decrease the distance to the optimal solution no slower than the gradient method with stepsize $\frac{1}{L}$. Hence, $y^k$ can provide more progress than $\bar{y}^{k-1}$. The proof is much easier and is deferred to Appendix B.

THEOREM 3.15. *Let $\{y^k\}_k$ be generated by* CG *applied to the system $B(y - y^*) = 0$, where $B \in \mathbb{R}^{n \times n}$ is symmetric, positive definite with $\|B\| \leq L$. Then, we have:*

$$\|y^{k+1} - y^*\|^2 \leq \|\bar{y}^k - y^*\|^2.$$

As discussed at the beginning of this section, the sequences $\{x_R^k\}_k$ and $\{y^k\}_k$ are equivalent up to a linear transformation, i.e., it holds that $y^k = B^{\frac{1}{2}} x_R^k$. This allows to transfer our obtained results back to the CR method. We summarize our observations for CR in the following corollary.

COROLLARY 3.16. *Let $B \in \mathbb{R}^{n \times n}$ be a symmetric, positive definite matrix with $\|B\| \leq L$ and let $x^0 \in \mathbb{R}^n$ be given. Suppose that $\{x_R^k\}_k$ is generated by the* CR *method to solve the linear system $B(x - x^0) = b$. Then, we have:*

$$\varphi(\bar{x}_R^k) \leq \varphi(\tilde{x}_R^{k-1}) \quad and \quad \varphi(x_R^k) \leq \varphi(\bar{x}_R^{k-1}),$$

*where $\bar{x}_R^k$ and $\tilde{x}_R^k$ are defined in (3.9) and $\varphi(x) := \frac{1}{2}(x - x^0)^\top B(x - x^0) - b^\top(x - x^0)$.*

**3.6. Proof of Theorem 3.6.** In this subsection, we combine our obtained results and show that AA-R locally decreases the objective function no slower than a gradient descent step with stepsize $\frac{1}{L}$ (up to a certain higher-order error term).

Throughout this section, we will work with the following choices $B = \nabla^2 f(x^0)$, $b = -\nabla f(x^0)$, and $A = B + E_m$, where $E_m$ is defined in (3.7).

*Proof of Theorem 3.6.* Clearly, (3.6) holds for $k = 0$. Furthermore, we only need to verify (3.6) for one cycle of AA-R as all assumptions and results will also hold for subsequent cycles, since all the subsequent iterations would also belong to $U$ by Proposition 3.5. Let $U = S_\epsilon$ be the neighborhood defined in Proposition 3.5. Then, for all $k \in \mathbb{N}$, we have:

$$\|h(x^{k+1})\| \leq (1 - (2\kappa_r)^{-1})\|h(x^k)\|.$$

Proposition 3.2 establishes $\kappa(H_k^\top H_k) \leq M_H^2$ for some $M_H > 0$ and by Lemma 3.4, we can infer $\|\hat{g}^k - x^0\| = \mathcal{O}(\|b\|)$. Due to $x^{k+1} = \hat{g}^k$, this just means $\|x^{k+1} - x^0\| = \mathcal{O}(\|b\|)$. Notice that this estimate holds for every $k = 0, 1, \ldots, m$ and therefore, it follows $M_x^m = \mathcal{O}(\|b\|)$. Furthermore, using (3.2), we obtain $\|E_m\| = \mathcal{O}(\|b\|)$. Reducing $\epsilon$ if

necessary, we may assume that $\|A - B\| = \|E_m\| < \mu$, which ensures the invertibility of $A$ as shown in the proof of Lemma 3.11. Now, let $\{x_G^k\}_k$ and $\{x_R^k\}_k$ be generated by GMRES and CR applied to the linear systems $A(x - x^0) = b$ and $B(x - x^0) = b$ with $x_G^0 = x_R^0 = x^0$, respectively. By Proposition 3.8, we have $\bar{x}_G^k = x^{k+1}$ for all $k = 0, \ldots, m$ and $\kappa((\bar{X}_G^k)^\top \bar{X}_G^k) \leq M^2$ for all $k \in [m]$. Moreover, since the perturbed gradient mapping $\bar{g}$ is exact at each $x^k$, $k = 0, \ldots, m$, it holds that

$$g(x^k) = \bar{g}(x^k) = \bar{g}(\bar{x}_G^{k-1}) = \tilde{x}_G^{k-1} \quad \forall \, k = 1, \ldots, m.$$

In addition, we have $\|\bar{x}_G^k - x^0\| = \|x^{k+1} - x^0\| = \mathcal{O}(\|b\|)$. Reducing $\epsilon$ — if necessary — we may assume $\epsilon \leq \epsilon_\sharp$, where $\epsilon_\sharp$ was introduced in the proof of Theorem 3.12. Thus, all conditions in Theorem 3.12 are satisfied and it follows

$$(3.22) \qquad \|\bar{x}_G^k - \bar{x}_R^k\| = \mathcal{O}(\|b\|^2) \quad \text{and} \quad \|\tilde{x}_G^k - \tilde{x}_R^k\| = \mathcal{O}(\|b\|^2) \quad \forall \, k = 0, \ldots, m.$$

Moreover, since $g$ is a contraction on $\mathbb{B}_r(x^\star)$ and due to $\|x^k - x^0\| = \mathcal{O}(\|b\|)$, we have $\|g(x^k) - x^0\| \leq \|g(x^k) - g(x^0)\| + \|g(x^0) - x^0\| = \mathcal{O}(\|b\|)$. Reusing the notation from Corollary 3.16, the Lipschitz continuity of the Hessian $\nabla^2 f$ then implies
(3.23)
$$f(x^{k+1}) = f(x^0) + \varphi(x^{k+1}) + \mathcal{O}(\|x^{k+1} - x^0\|^3) = f(x^0) + \varphi(\bar{x}_G^k) + \mathcal{O}(\|b\|^3),$$
$$f(g(x^k)) = f(x^0) + \varphi(g(x^k)) + \mathcal{O}(\|g(x^k) - x^0\|^3) = f(x^0) + \varphi(\tilde{x}_G^{k-1}) + \mathcal{O}(\|b\|^3),$$

see, e.g., [25, Lemma 4.1.1]. Since the mapping $\varphi$ is quadratic, we can further write:

$$\varphi(\bar{x}_G^k) = \varphi(\bar{x}_R^k) + \nabla\varphi(\bar{x}_R^k)^\top(\bar{x}_G^k - \bar{x}_R^k) + \tfrac{1}{2}(\bar{x}_G^k - \bar{x}_R^k)^\top B(\bar{x}_G^k - \bar{x}_R^k),$$
$$\varphi(\tilde{x}_G^{k-1}) = \varphi(\tilde{x}_R^{k-1}) + \nabla\varphi(\tilde{x}_R^{k-1})^\top(\tilde{x}_G^{k-1} - \tilde{x}_R^{k-1}) + \tfrac{1}{2}(\tilde{x}_G^{k-1} - \tilde{x}_R^{k-1})^\top B(\tilde{x}_G^{k-1} - \tilde{x}_R^{k-1}).$$

Next, applying Lemma 3.11 for the case $A = B$, it holds that:

$$\|\nabla\varphi(\bar{x}_R^k)\| = \|B(\bar{x}_R^k - x^0) - b\| \leq \|B(x_R^k - x^0) - b\| \leq \|B(x^0 - x^0) - b\| = \|b\|,$$

where we used Proposition 3.10 in the last step. Similarly, we can show $\|\nabla\varphi(\tilde{x}_R^{k-1})\| \leq \|b\|$. Thus, combining (3.22) and the representations of $\varphi(\bar{x}_G^k)$ and $\varphi(\tilde{x}_G^{k-1})$, we obtain

$$|\varphi(\bar{x}_G^k) - \varphi(\bar{x}_R^k)| \leq \|\nabla\varphi(\bar{x}_R^k)\|\|\bar{x}_G^k - \bar{x}_R^k\| + \tfrac{L}{2}\|\bar{x}_G^k - \bar{x}_R^k\|^2 = \mathcal{O}(\|b\|^3),$$
$$|\varphi(\tilde{x}_G^{k-1}) - \varphi(\tilde{x}_R^{k-1})| \leq \|\nabla\varphi(\tilde{x}_R^{k-1})\|\|\tilde{x}_G^{k-1} - \tilde{x}_R^{k-1}\| + \tfrac{L}{2}\|\tilde{x}_G^{k-1} - \tilde{x}_R^{k-1}\|^2 = \mathcal{O}(\|b\|^3),$$

Using these estimates in (3.23), we can infer

$$f(x^{k+1}) = f(x^0) + \varphi(\bar{x}_R^k) + \mathcal{O}(\|b\|^3), \quad f(g(x^k)) = f(x^0) + \varphi(\tilde{x}_R^{k-1}) + \mathcal{O}(\|b\|^3).$$

The conclusion then follows immediately from Corollary 3.16. $\qquad\qquad\square$

**4. A Function Value-Based Globalization for AA-R.** Based on the local descent properties established in the last section, we now propose a globalization mechanism for AA-R. We prove global convergence and provide simple global-to-local transition results for the globalized AA-R algorithm. To the best of our knowledge, this is the first function value-based globalization of AA-R that achieves both global and local convergence. Previously, only heuristic strategies seem to be available, see [31, 27, 42].

The full procedure is presented in Algorithm 4.1. Our core idea is to check whether the AA step $x_{\mathsf{AA}}^k$ satisfies a sufficient decrease condition

$$(4.1) \quad f(x_{\mathsf{AA}}^k) \leq f(x^k) - \gamma\|\nabla f(x^k)\|^2 + \min\{c_1\|\nabla f(x^{k-\hat{m}})\|^\nu, c_2\|\nabla f(x^{k-\hat{m}})\|^2, c_3\},$$

**Algorithm 4.1** A Globalized AA Scheme with Restarting

---

1: Choose an initial point $x^0 \in \mathbb{R}^n$, the memory parameter $m$, and constants $\gamma, \nu, c_1, c_2, c_3 > 0$. Set $\hat{m} = 0$.
2: **for** $k = 0, 1, \dots$ **do**
3:   Set $\hat{m} = \mathrm{mod}(k, m+1)$.
4:   **if** $\hat{m} = 0$ **then**
5:     Set $x^{k+1} = g(x^k)$.
6:   **else**
7:     Calculate the coefficient $\alpha^k$ based on the sequence $\{h(x^k), \dots, h(x^{k-\hat{m}})\}$ by solving (1.3) and set $x_{\mathsf{AA}}^k = g^{k-\hat{m}} + G_k \alpha^k$.
8:     **if** $f(x_{\mathsf{AA}}^k) > f(x^k) - \gamma\|\nabla f(x^k)\|^2 + \min\{c_1\|\nabla f(x^{k-\hat{m}})\|^\nu, c_2\|\nabla f(x^{k-\hat{m}})\|^2, c_3\}$ **then**
9:       Set $x^{k+1} = g(x^k)$.
10:     **else**
11:       Set $x^{k+1} = x_{\mathsf{AA}}^k$.
12:     **end if**
13:   **end if**
14: **end for**

---

where $\gamma$, $\nu$, and $c_i$, $i = 1, 2, 3$, are given parameters. We accept the AA step as new iterate $x^{k+1} = x_{\mathsf{AA}}^k$ if condition (4.1) holds. Otherwise, a gradient step $x^{k+1} = g(x^k)$ is performed (which ensures decrease of the objective function values). We summarize several basic convergence properties of Algorithm 4.1 in the following theorem.

THEOREM 4.1. *Suppose that* (A.1) *holds and let $f$ be bounded from below. Let the sequence $\{x^k\}_k$ be generated by Algorithm 4.1 with $\gamma, c_1, c_3 > 0$, $c_2 < \frac{1}{2mL}$, and $\nu \in (2, 3)$. Then, we have*

$$\sum\nolimits_{k=0}^{\infty} \|\nabla f(x^k)\|^2 < \infty \quad and \quad \lim_{k \to \infty} \|\nabla f(x^k)\| = 0.$$

*In addition, if $\gamma < \frac{1}{2L}$ and if the conditions* (A.2)–(A.4) *are satisfied with $x^\star$ being an accumulation point of $\{x^{k(m+1)}\}_k$, then we have $x^k \to x^\star$ and all AA steps will be eventually accepted, i.e., Algorithm 4.1 locally turns into Algorithm 2.1.*

*Proof.* Notice that the $k$-th cycle starts at iteration $(k-1)(m+1)$ and ends at iteration $k(m+1)$ (with $x^{k(m+1)}$ serving as initial point for the next cycle). In order to keep the notation simple, we concentrate on the first cycle. Since the first iteration within each cycle is a gradient descent step, i.e., $x^1 = g(x^0)$, we can deduce $f(x^1) \le f(x^0) - \frac{1}{2L}\|\nabla f(x^0)\|^2$. For all $k = 1, \dots, m$, the iterate $x^{k+1}$ either results from a gradient descent step or an AA step satisfying the acceptance criterion:

$$f(x^{k+1}) \le f(x^k) - \gamma\|\nabla f(x^k)\|^2 + c_2\|\nabla f(x^0)\|^2.$$

Hence, each update $k = 1, \dots, m$ satisfies $f(x^{k+1}) \le f(x^k) - \min\{\frac{1}{2L}, \gamma\}\|\nabla f(x^k)\|^2 + c_2\|\nabla f(x^0)\|^2$. Summing these estimates from 1 to $m$, we obtain

$$f(x^{m+1}) \le f(x^0) - \min\left\{\tfrac{1}{2L}, \gamma\right\} \sum\nolimits_{k=1}^{m} \|\nabla f(x^k)\|^2 - \left[\tfrac{1}{2L} - mc_2\right]\|\nabla f(x^0)\|^2.$$

Defining $\sigma := \min\{\frac{1}{2L}, \gamma, \frac{1}{2L} - mc_2\} > 0$, this result holds for every cycle of Algorithm 4.1, i.e., we have

$$f(x^{(k+1)(m+1)}) \le f(x^{k(m+1)}) - \sigma \sum\nolimits_{i=k(m+1)}^{k(m+1)+m} \|\nabla f(x^i)\|^2 \quad \forall\, k \in \mathbb{N}.$$

Summing this inequality for all $k \in \mathbb{N}$ and noticing that $f$ is bounded from below, it follows $\sum_{i=0}^{\infty} \|\nabla f(x^i)\|^2 < \infty$ which readily implies $\|\nabla f(x^i)\| \to 0$. Next, let $x^\star$ be an accumulation point of $\{x^{k(m+1)}\}_k$ satisfying (A.2)–(A.4). By Theorem 3.6, there is a neighborhood $U$ of $x^\star$ such that if $y^0 \in U$, then the sequence $\{y^k\}_k$ generated by Algorithm 2.1 satisfies

$$f(y^{k+1}) \leq f(g(y^k)) + \mathcal{O}(\|\nabla f(y^{k-\hat{m}})\|^3) \leq f(y^k) - \tfrac{1}{2L}\|\nabla f(y^k)\|^2 + \mathcal{O}(\|\nabla f(y^{k-\hat{m}})\|^3).$$

Thus, by shrinking $U$ if necessary and using $\gamma < \frac{1}{2L}$, we can assume

$$(4.2) \qquad f(y^{k+1}) \leq f(y^k) - \gamma\|\nabla f(y^k)\|^2 + c_1\|\nabla f(y^{k-\hat{m}})\|^\nu \quad \forall\, k.$$

Since $x^\star$ is an accumulation point of $\{x^{k(m+1)}\}_k$, there exists $s$ with $x^{s(m+1)} \in U$. We now set $y^k := x^{k+s(m+1)}$. Due to $y^0 \in U$, $\|\nabla f(x^i)\| \to 0$, and since the conditions (A.1)–(A.4) are satisfied, we can inductively infer that every AA step fulfills (4.2) and is accepted as new iterate, i.e., we have $y^{k+1} = x^{s(m+1)+k+1} = x_{\mathsf{AA}}^{s(m+1)+k}$, $k \geq 1$. Convergence of $\{x^k\}_k$ then follows from Proposition 3.5 and (A.2). □

**5. Numerical Experiments.** In this section, we conduct preliminary numerical experiments to illustrate the performance and convergence behavior of AA-R and to empirically verify the descent properties of Algorithm 4.1[1].

**5.1. Nonconvex Classification.** We consider two nonconvex classification problems, namely a nonlinear least-squares problem and a student's-$t$ problem. A detailed introduction of the tested problems is deferred to the subsequent paragraphs. We will compare Algorithm 4.1 with four different methods:

(1) The gradient descent method (GD) with fixed step size $\frac{1}{L}$. This is the original Picard iteration (1.2) and serves as a baseline.
(2) Pure AA with and without restarting [47]. Pure AA does not use any globalization strategy, i.e., in each step, we perform an AA iteration.
(3) AA with residual-based globalization. Our implementation is based on A2DR [14, Algorithm 3] and we consider two variants with and without restarting. A2DR uses a residual-based acceptance mechanism and we apply the method with the default parameters as suggested in [14].
(4) L-BFGS. We implement L-BFGS with Wolfe conditions as in [26, Algorithm 7.5]. The line-search parameter is set to $10^{-4}$ and the parameter in Wolfe's condition is set to 0.9 as suggested in [26]. The maximum number of line-search iterations is set to $1,000$.

We note that the line search procedure in L-BFGS can contain many function and gradient evaluations per iteration. Therefore, it is improper to give comparisons solely based on the number of iterations. In our plots, the $x$-axis is typically set as the number of oracle calls, which appears to be a more appropriate and fair measure when comparing AA algorithms and L-BFGS. Specifically, the computation of each function value is counted as one oracle call and every gradient evaluation contributes as an additional oracle call. The $y$-axis is set as $(f(x^k) - f^\star)/\max\{f^\star, 1\}$ or $\|\nabla f(x^k)\|$, respectively, where $f^\star$ denotes the best objective function value among all algorithms over the maximum oracle calls. In the figures, we will add a special mark "$\star$" once the current AA step is rejected and a gradient step is performed in Algorithm 4.1.

We continue with the description of the utilized training datasets and several universal implementational details. We use the CIFAR10 dataset [21], which contains

---

[1]Code available under https://github.com/yangliu-op/AndersonAcceleration
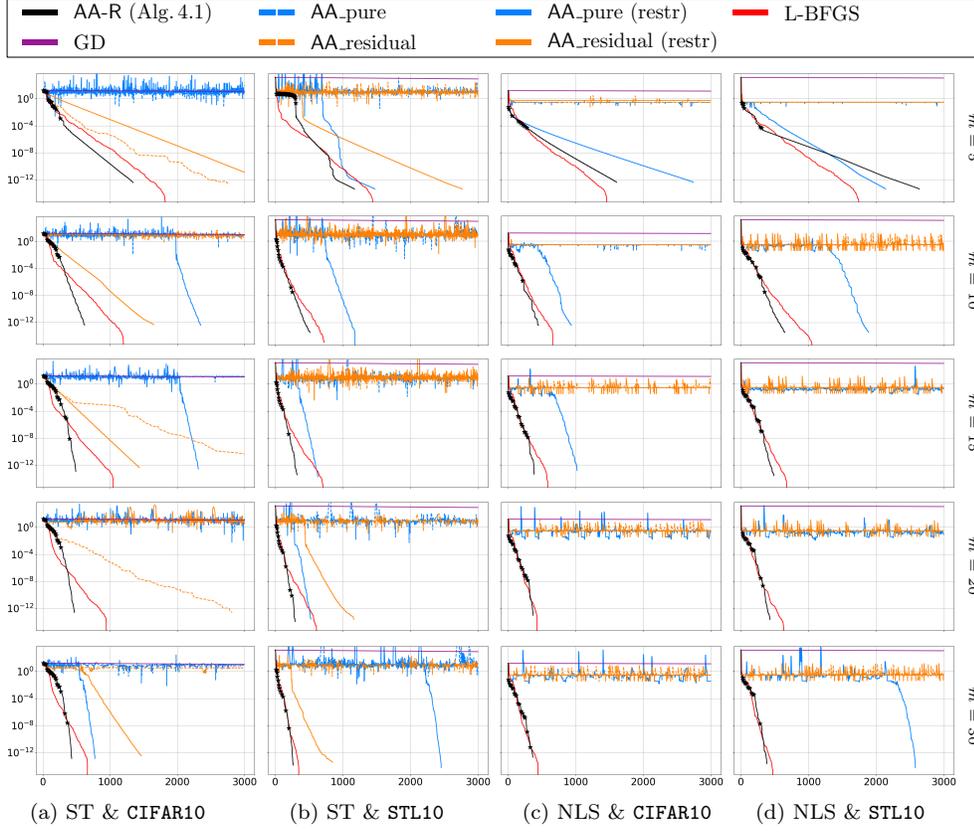
FIG. 1. *Relative error* $(f(x^k) - f^*)/\max\{f^*, 1\}$ *vs. Oracle calls for the student's t (ST) and nonlinear least-squares (NLS) problem and the datasets* **CIFAR10** *and* **STL10**. *The plots in each column are generated using the identical initial point* $x^0 \sim \mathcal{N}^d(0, 1)$. *In each row, the different* AA *methods and* **L-BFGS** *are executed using the same memory parameter* $m \in \{5, 10, 15, 20, 30\}$. *The x-axes of each plot have the same scaling* $0 - 3,000$ *(as shown in the bottom row).*

$60,000$ images with $32 \times 32$ colored pixels and the **STL10** dataset [8], which consists of $5,000$ colored images of size $96 \times 96$. Given that both datasets contain 10 classes, we split the data into even and odd classes to allow binary classification. We use $\{u_i, v_i\}_{i=1}^n$ to denote the training samples, where $u_i \in \mathbb{R}^d$ represents the training image and $v_i \in \{0, 1\}$ is the associated label. We set $U = \{u_1, u_2, \ldots, u_n\}^\top \in \mathbb{R}^{n \times d}$. We terminate the algorithms once $\|\nabla f(x^k)\| \leq 10^{-7}$ or the number of oracle calls exceeds $3,000$. The memory parameter $m$ is chosen from $m \in \{5, 10, 15, 20, 30\}$ for all AA-based methods and **L-BFGS**. The regularization parameter $\lambda$ in (5.1) and (5.2) is set to $10^{-2}$ for **CIFAR10** and to $10^{-1}$ for **STL10**. The initial points $x^0 \sim \mathcal{N}^d(0, 1)$ are generated following a normal distribution. Finally, in Algorithm 4.1, we utilize the default parameters: $\gamma = \frac{0.01}{2L}$, $c_1 = c_3 = 1$, $c_2 = \frac{0.99}{2mL}$, and $\nu = 2.1$. Let us briefly motivate this default choice. In order to promote acceptance of AA steps (and to ensure potential acceleration), the descent condition (4.1) should not be too strict. This can be achieved when $\gamma$ is small and when the min-term in (4.1) is large. Hence, we set $\gamma$ fairly small, $c_2$ close to the theoretical threshold, and $\nu$ close to 2. Furthermore, we have found that the simple choice $c_1 = c_3 = 1$ works well enforcing sufficient progress during the first iterations. An additional ablation study for $c_1$, $c_2$, $c_3$, and $\gamma$ is
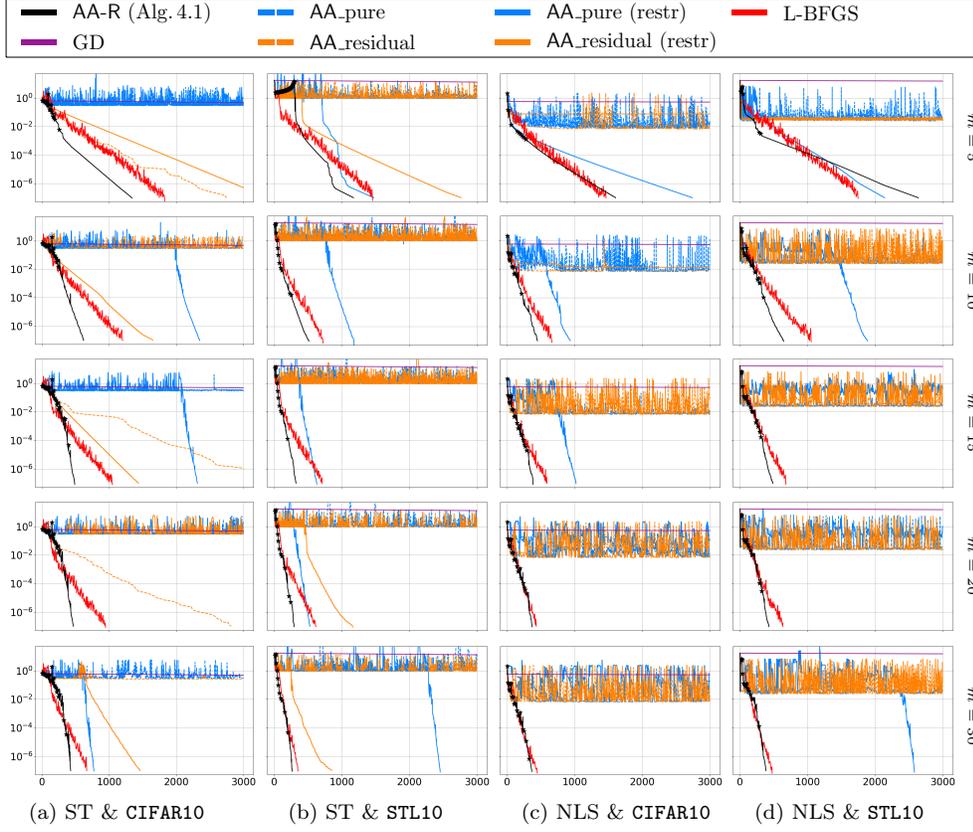
Fig. 2. $\|\nabla f(x^k)\|$ vs. Oracle calls for the student's t (ST) and nonlinear least-squares (NLS) problem and the datasets CIFAR10 and STL10. The plots in each column are generated using the identical initial point $x^0 \sim \mathcal{N}^d(0, 1)$. In each row, the different AA methods and L-BFGS are executed using the same memory parameter $m \in \{5, 10, 15, 20, 30\}$. The x-axes of each plot have the same scaling $0 - 3,000$ (as shown in the bottom row).

discussed in subsection 5.3. We use the LSQR method [29] to solve the AA subproblem (1.3) and to compute $\alpha^k = -H_k^\dagger h^{k-\hat{m}} = -(H_k^\top H_k)^{-1} H_k^\top h^{k-\hat{m}}$. (Here, $H_k^\dagger$ represents the Moore-Penrose pseudo-inverse of $H_k$). The termination condition of LSQR is set to $\|H_k^\top (H_k \alpha + h^{k-\hat{m}})\| < 10^{-16}$.

Next, we present the classification models used in our numerical comparison:

• **Student's-t Loss with $\ell_2$-Regularization (ST).** We consider the following classification problem with student's-t loss, [3, 2],

$$(5.1) \qquad f(x) = \frac{1}{n}\sum\nolimits_{i=1}^{n} \log\left(1 + (u_i^\top x - v_i)^2/\mu\right) + \frac{\lambda}{2}\|x\|^2.$$

The Lipschitz constant of $\nabla f$ is given by $L = \frac{2}{\mu n}\|U\|^2 + \lambda$ and we set $\mu = 20$.

• **Nonlinear Least-Squares Problem with $\ell_2$-Regularization (NLS).** As a second example, we consider a nonlinear least-squares problem, [50],

$$(5.2) \qquad f(x) = \frac{1}{n}\sum\nolimits_{i=1}^{n} (\psi(u_i^\top x) - v_i)^2 + \frac{\lambda}{2}\|x\|^2,$$

where $\psi(z) = 1/(1 + e^{-z})$ is the sigmoid function. The Lipschitz constant of

$\nabla f$ is given by $L = \frac{1}{6n}\|U\|^2 + \lambda$.

The initial points for all algorithms and $m \in \{5, 10, 15, 20, 30\}$ are identical for each tested dataset and classification model. Figures 1 and 2 illustrate that AA-R with function value-based globalization (4.1) is a competitive solver. Specifically, Algorithm 4.1 requires the least amount of oracle calls to satisfy the stopping criterion when $m \in \{10, 15, 20, 30\}$. However, in the low memory case $m = 5$, Algorithm 4.1 and the restarting strategy seem less effective (especially for the nonlinear least-squares problem). The plots in Figures 1–2 generally underline the potential of function value- and descent-based globalization mechanisms for AA schemes. Rejections predominantly occur in the early stage of the iterative process to ensure global convergence and progress of Algorithm 4.1. In addition, transition to a pure AA-R phase with accelerated convergence is maintained — as indicated by our theoretical results.

As the applications tested in this section are nonconvex, we have recorded the smallest eigenvalues of the respective Hessians in the last iterations of AA-R for each of the problems and datasets. We have observed that that these eigenvalues are all approximately equal to $\lambda > 0$ and hence, assumption (A.2) is locally satisfied.

**5.2. Descent Properties.** In Figure 3, we plot the measure

$$\rho_k := \max\{f(x_{\mathsf{AA}}^k) - f(g(x^k)), 0\}/\|\nabla f(x^{k-\hat{m}})\|^3$$

versus the number of iterations $k$ to further visualize and verify the descent properties derived in Theorem 3.6. If the AA step achieves descent, $f(x_{\mathsf{AA}}^k) \leq f(g(x^k))$, then we have $\rho_k = 0$ and we locally expect $\rho_k \approx 0$ for all $k$ sufficiently large. The special marks "$\star$" in Figure 3 again indicate that an AA step did not pass the descent condition (4.1). Figure 3 illustrates that $\rho_k$ indeed stays zero eventually and that no AA steps are rejected locally. This observation is slightly less pronounced on CIFAR10 when $m = 30$.

**5.3. Ablation Study.** Finally, we provide an additional ablation study for the parameters $c_1$, $c_2$, $c_3$, and $\gamma$ used in Algorithm 4.1 and in the definition of the descent condition (4.1). Based on Theorem 4.1, $c_1$, $c_2$, $c_3$, and $\gamma$ need to satisfy the conditions $0 < \gamma < \frac{1}{2L}$, $c_1, c_3 > 0$, and $0 < c_2 < \frac{1}{2mL}$. We compare our default choice with the following extreme sets of parameters:

$$(5.3) \qquad \gamma = \frac{1}{2L},\ c_1 = c_2 = c_3 = 0, \quad \text{and} \quad \gamma = 0,\ c_1 = c_3 = 10^{10},\ c_2 = \frac{1}{2mL}.$$

These two choices correspond to highly strict and loose acceptance criteria for the AA step $x_{\mathsf{AA}}^k$. Since $\nu \in (2, 3)$ has only limited influence, we omit an explicit ablation study for $\nu$ and use the default value $\nu = 2.1$. Figure 4 demonstrates that Algorithm 4.1 is robust with respect to the choice of $c_1$, $c_2$, $c_3$, and $\gamma$. In particular, performance is only affected marginally when using the more extreme parameters (5.3).

**6. Conclusion.** In this work, we study descent properties of an Anderson accelerated gradient method with restarting. We first show that the iterates generated by AA-R are equivalent to the iterates generated by GMRES after an additional gradient step within each restarting cycle. Based on the symmetry of the underlying system matrix, we then analyze the error between the iterates generated by GMRES and CR and verify that this error is controllable and related to some higher-order perturbation terms. After connecting CR and CG, the desired descent property for AA-R can be expressed in terms of distances to the respective optimal solution for the iterates generated by CG. We establish such a convergence result for CG utilizing classical techniques. Combining these different observations, we prove that AA-R can decrease
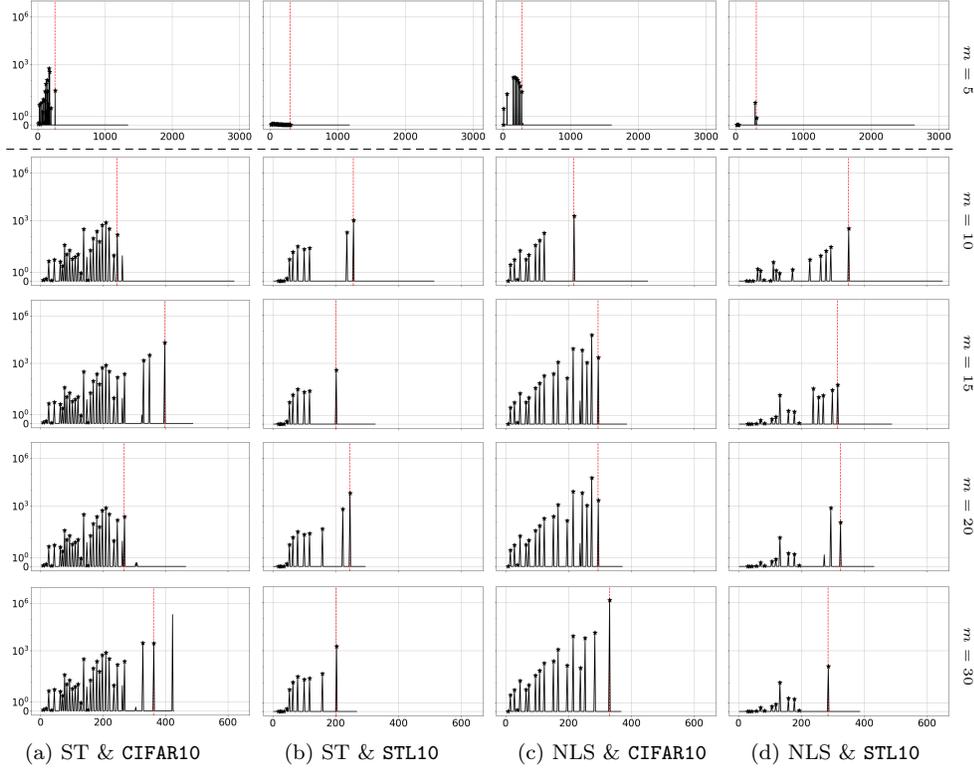
FIG. 3. *Plot of $\rho_k = \max\{f(x_{\mathsf{AA}}^k) - f(g(x^k)), 0\}/\|\nabla f(x^{k-\hat{m}})\|^3$ vs. Oracle calls for Algorithm 4.1. The marks "$\star$" indicate rejected AA steps. After the last rejected AA step (red dashed vertical line), $\rho_k$ mostly stays $0$, which verifies and illustrates Theorem 3.6. The x-axes of each plot in the rows $m \in \{10, 15, 20, 30\}$ have the same scaling $0\,$–$\,700$. For $m = 5$, the scaling is $0\,$–$\,3,000$.*

the objective function $f$ locally. These novel findings can be used in the design of effective, function value-based globalization mechanisms for AA-R approaches. We propose one such possible AA-R globalization and conduct numerical experiments on two large-scale learning problems that illustrate our theoretical results.

**Acknowledgments.** We would like to thank the Associate Editor and three anonymous reviewers for their detailed and constructive comments, which have helped greatly to improve the quality and presentation of the manuscript.

### Appendix A. Proof of Proposition 3.3.

*Proof.* We show that Proposition 3.3 is a direct application of [33, Theorem 5.1]. Due to (2.1), the constant $\kappa_g$ in [33] reduces to $1 - \frac{1}{\kappa_r}$ and we have $\theta_k \leq 1$ and $\beta_k = 1$. Moreover, all the points of interest lie in $\mathbb{B}_r(x^\star)$, so all the expansions of the residuals in [33, Section 3] are legitimate. The core estimate (5.18) in [33, Theorem 5.1] then reduces to (3.4). The proof is complete if all assumptions in [33, Theorem 5.1] hold.

Using $x^i \in \mathbb{B}_r(x^\star)$, $i = k - \hat{m}, \ldots, k$ and (2.1), it follows

$$\|h^i - h^{i-1}\| \geq \|x^i - x^{i-1}\| - \|g(x^i) - g(x^{i-1})\| \geq \kappa_r^{-1}\|x^i - x^{i-1}\| \quad \forall\, i = k - \hat{m} + 1, \ldots, k.$$

This is exactly Assumption 2.3 in [33] (see also [33, Remark 2.1]). Next, we verify the sufficient linear independence condition introduced in [33, Lemma 5.2]. Let us
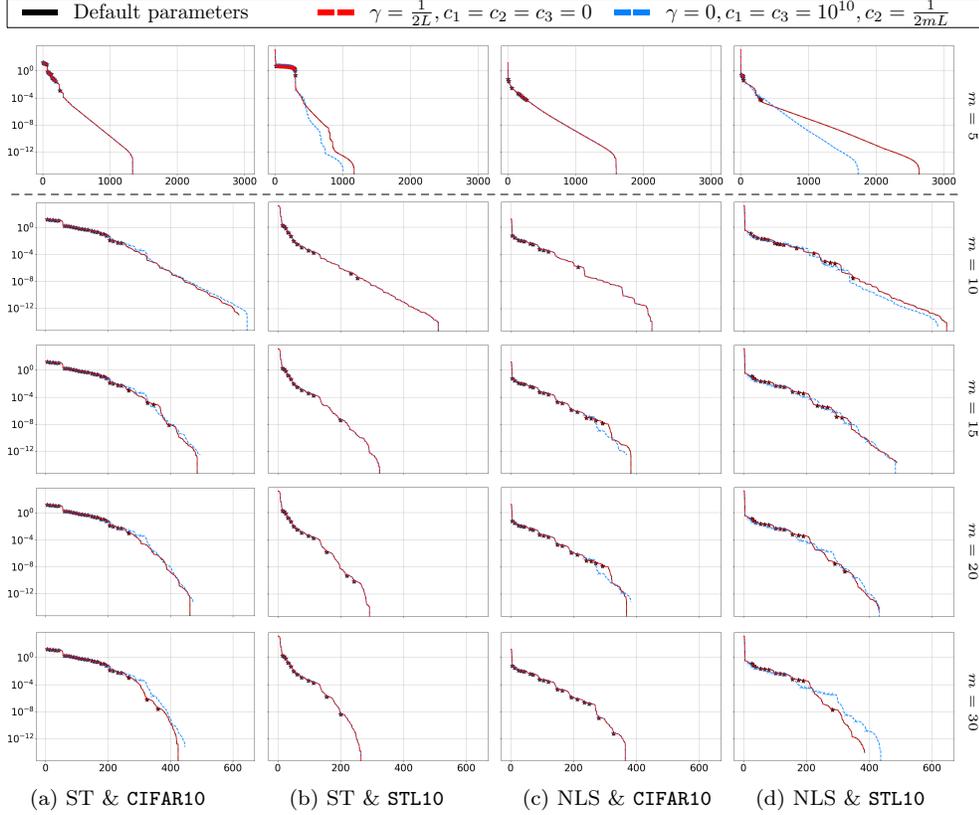
FIG. 4. *Ablation study of Algorithm 4.1 using different $c_1$, $c_2$, $c_3$, and $\gamma$. Each plot depicts $(f(x^k) - f^*)/\max\{f^*, 1\}$ vs. Oracle calls for three different runs of Algorithm 4.1. We compare the default parameters with the extreme choices $\gamma = \frac{1}{2L}$, $c_1 = c_2 = c_3 = 0$ and $\gamma = 0$, $c_1 = c_3 = 10^{10}$, $c_2 = \frac{1}{2mL}$. The x-axes of each plot in the rows $m \in \{10, 15, 20, 30\}$ have the same scaling $0$–$700$. For $m = 5$, the scaling is $0$–$3,000$.*

define $\tilde{H}_k = [h^k - h^{k-1}, \ldots, h^{k-\hat{m}+1} - h^{k-\hat{m}}] =: [v^1, \ldots, v^{\hat{m}}]$. We note that there is a fixed nonsingular matrix $P \in \mathbb{R}^{\hat{m} \times \hat{m}}$ such that $\tilde{H}_k = H_k P$ and $\kappa(\tilde{H}_k^\top \tilde{H}_k) \leq \kappa(H_k^\top H_k) \kappa(P^\top P)$. Therefore, by Proposition 3.2 and using $x^i \in U_1$, $i = k - \hat{m}, \ldots, k$, (A.4) implies that the condition number of $\tilde{H}_k^\top \tilde{H}_k$ is bounded by some $\tilde{M}^2$. Let $\mathcal{V}_i = \text{span}\{v^1, \ldots, v^i\}$ denote the linear subspace spanned by the first $i$ columns of $\tilde{H}_k$ and let $\tilde{H}_k = Q_k R_k$ be the QR decomposition of $\tilde{H}_k$. We then have $\kappa(R_k^\top R_k) = \kappa(\tilde{H}_k^\top \tilde{H}_k) \leq \tilde{M}^2$. Furthermore, let $\{r_{ii}\}_{1 \leq i \leq \hat{m}}$ denote the diagonal entries of $R_k$. By [33, Proposition 5.2], it follows $r_{11}^2 = \|v_1\|^2$ and $r_{ii}^2 = \|v_i\|^2 \sin^2(v_i, \mathcal{V}_{i-1})$ for all $2 \leq i \leq \hat{m}$. Since $R_k$ is upper triangular, the diagonal entries $r_{ii}$, $i = 1, \ldots, \hat{m}$, are exactly the eigenvalues of $R_k$. Consequently, we obtain

$$(\|v_i\|^2/\|v_1\|^2) \cdot \sin^2(v_i, \mathcal{V}_{i-1}) = r_{ii}^2/r_{11}^2 \geq \sigma_{\min}(R_k)^2/\sigma_{\max}(R_k)^2 \geq 1/\tilde{M}^2.$$

In addition, we have $\|v_i\|^2/\|v_1\|^2 \leq \sigma_{\max}(\tilde{H}_k)^2/\sigma_{\min}(\tilde{H}_k)^2 \leq \tilde{M}^2$. Combining these inequalities, this yields $|\sin(v_i, \mathcal{V}_{i-1})| \geq \tilde{M}^{-2}$ which verifies the last remaining assumption in [33, Lemma 5.2 and Theorem 5.1]. This concludes the proof. $\qquad\square$

**Appendix B. Proof of Theorem 3.15.**

*Proof.* Similar to (3.15) and utilizing the projection $y^{(k+1)}$, we have:

$$\|\bar{y}^k - y^*\|^2 - \|y^{k+1} - y^*\| = \|\bar{y}^k - y^{(k+1)}\|^2 - \|y^{k+1} - y^{(k+1)}\|^2$$

$$= \|\bar{y}^k - y^{k+1}\|^2 + 2\langle \bar{y}^k - y^{k+1}, y^{k+1} - y^{(k+1)}\rangle$$

$$= \|a_k p^k - L^{-1}r^k\|^2 + 2\langle a_k p^k - L^{-1}r^k, \gamma_k p^k\rangle$$

$$= a_k^2\|p^k\|^2 - 2a_k L^{-1}\langle p^k, r^k\rangle + L^{-2}\|r^k\|^2 + 2a_k\gamma_k\|p^k\|^2 - 2\gamma_k L^{-1}\langle p^k, r^k\rangle$$

$$= (a_k^2 + 2a_k\gamma_k)\|p^k\|^2 + \tfrac{1}{L}[L^{-1} - 2\gamma_k - 2a_k]\|r^k\|^2$$

$$\geq \left[a_k^2 + 2a_k\gamma_k - 2\gamma_k L^{-1} - 2a_k L^{-1} + L^{-2}\right]\|r^k\|^2$$

$$= \left[2\gamma_k(a_k - L^{-1}) + (a_k - L^{-1})^2\right]\|r^k\|^2 \geq 0,$$

where we have used Property (iii) to show that $\langle p^k, r^k\rangle = \|r^k\|^2$, Property (iv) to show that $\|p^k\|^2 \geq \|r^k\|^2$, and Property (ix) to show that $\gamma_k \geq a_k \geq 1/L$. □

## REFERENCES

[1] D. G. Anderson, *Iterative procedures for nonlinear integral equations*, J. ACM, 12 (1965), pp. 547–560.

[2] A. Aravkin, M. P. Friedlander, F. J. Herrmann, and T. Van Leeuwen, *Robust inversion, dimensionality reduction, and randomized sampling*, Math. Program., 134 (2012), pp. 101–125.

[3] A. Aravkin, T. Van Leeuwen, and F. Herrmann, *Robust full-waveform inversion using the student's t-distribution*, in SEG Tech. Program Expanded Abstracts, 2011, pp. 2669–2673.

[4] E. Artacho, E. Anglada, O. Diéguez, J. D. Gale, A. García, J. Junquera, R. M. Martin, P. Ordejón, J. M. Pruneda, D. Sánchez-Portal, and J. M. Soler, *The SIESTA method; developments and applicability*, J. Phys.-Condes. Matter, 20 (2008).

[5] Z. Bai, D. Hu, and L. Reichel, *A Newton basis GMRES implementation*, IMA J. Numer. Anal., 14 (1994), pp. 563–581.

[6] W. Bian, X. Chen, and C. Kelley, *Anderson acceleration for a class of nonsmooth fixed-point problems*, SIAM J. Sci. Comput., (2021), pp. S1–S20.

[7] M. Chupin, M.-S. Dupuy, G. Legendre, and E. Séré, *Convergence analysis of adaptive DIIS algorithms with application to electronic ground state calculations*, ESAIM Math. Model. Numer. Anal., 55 (2021), pp. 2785–2825.

[8] A. Coates, A. Ng, and H. Lee, *An analysis of single-layer networks in unsupervised feature learning*, in Proc. Int. Conf. Artif. Intell. Stat. (AISTATS), 2011, pp. 215–223.

[9] M. Ermis and I. Yang, *A3DQN: Adaptive Anderson acceleration for deep Q-networks*, in 2020 IEEE Symposium Series on Computational Intelligence (SSCI), IEEE, 2020, pp. 250–257.

[10] C. Evans, S. Pollock, L. G. Rebholz, and M. Xiao, *A proof that Anderson acceleration improves the convergence rate in linearly converging fixed-point methods (but not in those converging quadratically)*, SIAM J. Numer. Anal., 58 (2020), pp. 788–810.

[11] V. Eyert, *A comparative study on methods for convergence acceleration of iterative vector sequences*, Journal of Computational Physics, 124 (1996), pp. 271–285.

[12] H.-r. Fang and Y. Saad, *Two classes of multisecant methods for nonlinear acceleration*, Numer. Linear Algebra Appl., 16 (2009), pp. 197–221.

[13] D. C.-L. Fong and M. Saunders, *CG versus MINRES: An empirical comparison*, Sultan Qaboos University Journal for Science [SQUJS], 17 (2012), pp. 44–62.

[14] A. Fu, J. Zhang, and S. Boyd, *Anderson accelerated Douglas–Rachford splitting*, SIAM J. Sci. Comput., 42 (2020), pp. A3560–A3583.

[15] M. Geist and B. Scherrer, *Anderson acceleration for reinforcement learning*, arXiv preprint arXiv:1809.09501, (2018).

[16] G. H. Golub and C. F. Van Loan, *Matrix computations*, JHU Press, Baltimore, MD, 2013.

[17] X. Guo, A. Hu, R. Xu, and J. Zhang, *Consistency and computation of regularized mles for multivariate hawkes processes*, arXiv preprint arXiv:1810.02955, (2018).

[18] M. H. Gutknecht, *A brief introduction to Krylov space methods for solving linear systems*, in Front. Comput. Sci., Springer, 2007, pp. 53–62.

[19] N. C. Henderson and R. Varadhan, *Damped Anderson acceleration with restarts and monotonicity control for accelerating EM and EM-like algorithms*, J. Comput. Graph. Stat., 28 (2019), pp. 834–846.

[20] M. R. HESTENES AND E. STIEFEL, *Methods of conjugate gradients for solving linear systems*, J. Res. Natl. Bur. Stand., 49 (1952), pp. 409–435.

[21] A. KRIZHEVSKY, *Learning multiple layers of features from tiny images*, Master's thesis, University of Tront, (2009).

[22] J. LOFFELD AND C. S. WOODWARD, *Considerations on the implementation and use of Anderson acceleration on distributed memory and GPU-based parallel computers*, in Adv. Math. Sci., Springer, 2016, pp. 417–436.

[23] V. MAI AND M. JOHANSSON, *Anderson acceleration of proximal gradient methods*, in Int. Conf. Mach. Learn., PMLR, 2020, pp. 6620–6629.

[24] R. MEYER, *On the convergence of algorithms with restart*, SIAM J. Numer. Anal., 13 (1976), pp. 696–704.

[25] Y. NESTEROV, *Lectures on convex optimization*, vol. 137, Springer, 2018.

[26] J. NOCEDAL AND S. J. WRIGHT, *Numerical optimization*, Springer Series in Operations Research and Financial Engineering, Springer, New York, second ed., 2006.

[27] W. OUYANG, Y. PENG, Y. YAO, J. ZHANG, AND B. DENG, *Anderson acceleration for nonconvex ADMM based on Douglas-Rachford splitting*, Comput. Graph. Forum, 39 (2020), pp. 221–239.

[28] W. OUYANG, J. TAO, A. MILZAREK, AND B. DENG, *Nonmonotone globalization for Anderson acceleration using adaptive regularization*, arXiv preprint arXiv:2006.02559, (2020).

[29] C. C. PAIGE AND M. A. SAUNDERS, *LSQR: An algorithm for sparse linear equations and sparse least squares*, ACM transactions on mathematical software, 8 (1982), pp. 43–71.

[30] A. L. PAVLOV, G. W. OVCHINNIKOV, D. Y. DERBYSHEV, D. TSETSERUKOU, AND I. V. OSELEDETS, *AA-ICP: Iterative closest point with Anderson acceleration*, in IEEE Int. Conf. Robot. Autom. (ICRA), IEEE, 2018, pp. 1–6.

[31] Y. PENG, B. DENG, J. ZHANG, F. GENG, W. QIN, AND L. LIU, *Anderson acceleration for geometry optimization and physics simulation*, ACM Trans. Graph., 37 (2018), p. 42.

[32] X.-H. PHAM, M. ALAMIR, F. BONNE, AND P. BONNAY, *On the use of Anderson acceleration in hierarchical control*, arXiv preprint arXiv:2112.04299, (2021).

[33] S. POLLOCK AND L. G. REBHOLZ, *Anderson acceleration for contractive and noncontractive operators*, IMA J. Numer. Anal., 41 (2021), pp. 2841–2872.

[34] S. POLLOCK, L. G. REBHOLZ, AND M. XIAO, *Anderson-accelerated convergence of Picard iterations for incompressible Navier–Stokes equations*, SIAM J. Numer. Anal., 57 (2019), pp. 615–637.

[35] F. A. POTRA AND H. ENGLER, *A characterization of the behavior of the Anderson acceleration on linear problems*, Linear Alg. Appl., 438 (2013), pp. 1002–1011.

[36] M. J. D. POWELL, *Restart procedures for the conjugate gradient method*, Math. Program., 12 (1977), pp. 241–254.

[37] P. P. PRATAPA AND P. SURYANARAYANA, *Restarted Pulay mixing for efficient and robust acceleration of fixed-point iterations*, Chem. Phys. Lett., 635 (2015), pp. 69–74.

[38] T. ROHWEDDER AND R. SCHNEIDER, *An analysis for the DIIS acceleration method used in quantum chemistry calculations*, J. Math. Chem., 49 (2011), pp. 1889–1914.

[39] Y. SAAD, *Iterative methods for sparse linear systems*, SIAM, 2003.

[40] Y. SAAD AND M. H. SCHULTZ, *GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems*, SIAM J. Sci. Comput., 7 (1986), pp. 856–869.

[41] D. SCIEUR, A. D'ASPREMONT, AND F. BACH, *Regularized nonlinear acceleration*, in Adv. Neural Inf. Process. Syst., 2016, pp. 712–720.

[42] D. SCIEUR, A. D'ASPREMONT, AND F. BACH, *Nonlinear acceleration of stochastic algorithms*, arXiv preprint arXiv:1706.07270, (2017).

[43] W. SHI, S. SONG, H. WU, Y.-C. HSU, C. WU, AND G. HUANG, *Regularized Anderson acceleration for off-policy deep reinforcement learning*, preprint arXiv:1909.03245, (2019).

[44] E. STIEFEL, *Relaxationsmethoden bester Strategie zur Lösung linearer Gleichungssysteme*, Commentarii Mathematici Helvetici, 29 (1955), pp. 157–179.

[45] W. TANG AND P. DAOUTIDIS, *Fast and stable nonconvex constrained distributed optimization: the ellada algorithm*, Optimization and Engineering, 23 (2022), pp. 259–301.

[46] A. TOTH AND C. KELLEY, *Convergence analysis for Anderson acceleration*, SIAM J. Numer. Anal., 53 (2015), pp. 805–819.

[47] H. F. WALKER AND P. NI, *Anderson acceleration for fixed-point iterations*, SIAM J. Numer. Anal., 49 (2011), pp. 1715–1735.

[48] D. WANG, Y. HE, AND H. DE STERCK, *On the asymptotic linear convergence speed of Anderson acceleration applied to ADMM*, J. Sci. Comput., 88 (2021), pp. 1–35.

[49] F. WEI, C. BAO, AND Y. LIU, *Stochastic Anderson mixing for nonconvex stochastic optimization*, in Adv. Neural Inf. Process. Syst., vol. 34, 2021, pp. 22995–23008.

[50] P. Xu, F. Roosta, and M. W. Mahoney, *Second-order optimization for non-convex machine learning: An empirical study*, Proc. SIAM Int. Conf. Data Min., (2020), pp. 199–207.

[51] J. Zhang, B. O'Donoghue, and S. Boyd, *Globally convergent type-I Anderson acceleration for nonsmooth fixed-point iterations*, SIAM J. Optim., 30 (2020), pp. 3170–3197.

[52] J. Zhang, Y. Peng, W. Ouyang, and B. Deng, *Accelerating ADMM for efficient simulation and optimization*, ACM Trans. Graph., 38 (2019), pp. 1–21.