

# LOCAL IMPROVEMENT RESULTS FOR ANDERSON ACCELERATION WITH INACCURATE FUNCTION EVALUATIONS\*

ALEX TOTH<sup>†</sup>, J. AUSTIN ELLIS<sup>‡</sup>, TOM EVANS<sup>§</sup>, STEVEN HAMILTON<sup>§</sup>,  
C. T. KELLEY<sup>‡</sup>, ROGER PAWLOWSKI<sup>†</sup>, AND STUART SLATTERY<sup>§</sup>

**Abstract.** We analyze the convergence of Anderson acceleration when the fixed point map is corrupted with errors. We consider uniformly bounded errors and stochastic errors with infinite tails. We prove local improvement results which describe the performance of the iteration up to the point where the accuracy of the function evaluation causes the iteration to stagnate. We illustrate the results with examples from neutronics.

**Key words.** nonlinear equations, Anderson acceleration, local improvement

**AMS subject classifications.** 65H10, 82D75

**DOI.** 10.1137/16M1080677

**1. Introduction.** Anderson acceleration [3] is an iterative method for fixed point problems

$$(1) \quad u = G(u),$$

where  $u \in R^N$  and  $G : R^N \rightarrow R^N$ . The method was designed to accelerate Picard or fixed point iteration

$$(2) \quad u_{k+1} = G(u_k).$$

Anderson acceleration was originally designed for electronic structure computations and is now very common in that field. Other methods used in electronic structure computations, such as Pulay mixing (or Direct Inversion on the Iterative Subspace) [18, 20, 26, 27, 28, 30] or nonlinear GMRES [5, 22, 24, 39] are essentially the same as Anderson acceleration. Other applications include nuclear reactor design [14, 36] and hydrology [21].

The Anderson iteration maintains a history of residuals

$$F(u) = G(u) - u$$

of size at most  $m + 1$ , where the *depth*  $m$  is an algorithmic parameter. When  $m$  is important, we will call the iteration Anderson( $m$ ). Anderson(0) is fixed point iteration by definition.

\*Received by the editors July 11, 2016; accepted for publication (in revised form) October 26, 2016; published electronically October 26, 2017.

<http://www.siam.org/journals/sisc/39-5/M108067.html>

**Funding:** The research reported in this paper has been partially supported by the Consortium for Advanced Simulation of Light Water Reactors (<http://www.casl.gov>), an Energy Innovation Hub (<http://www.energy.gov/hubs>) for Modeling and Simulation of Nuclear Reactors under U.S. Department of Energy contract DE-AC05-00OR22725, and National Science Foundation grants DMS-1406349 and SI2-SSE-1339844.

<sup>†</sup>Sandia National Laboratories, MS 0316, P.O. Box 5800, Albuquerque, NM 87185 (artoth@sandia.gov, rppawlo@sandia.gov).

<sup>‡</sup>Department of Mathematics, North Carolina State University, Box 8205, Raleigh, NC 27695-8205 (jaellis2@ncsu.edu, Tim.Kelley@ncsu.edu).

<sup>§</sup>Oak Ridge National Laboratory, Oak Ridge, TN 37831 (evanstm@ornl.gov, hamiltonsp@ornl.gov, slatterysr@ornl.gov).

The formal description in Algorithm 1 is most convenient for analysis and is the one we will use in this paper. Algorithm 1 is not an optimal way to implement Anderson acceleration. We refer the reader to [7, 34, 35, 37, 38] for examples of efficient implementations.

---

**Algorithm 1.** Anderson acceleration.

---

**anderson**( $u_0, G, m$ )

$u_1 = G(u_0); F_0 = G(u_0) - u_0$

**for**  $k = 1, \dots$  **do**

$m_k = \min(m, k)$

$F_k = G(u_k) - u_k$

Minimize  $\|\sum_{j=0}^{m_k} \alpha_j^k F_{k-m_k+j}\|$  subject to

$\sum_{j=0}^{m_k} \alpha_j^k = 1.$

$u_{k+1} = (1 - \beta_k) \sum_{j=0}^{m_k} \alpha_j^k u_{k-m_k+j} + \beta_k \sum_{j=0}^{m_k} \alpha_j^k G(u_{k-m_k+j})$

**end for**

---

In a recent paper [35] we showed that Anderson acceleration with  $\beta_k \equiv 1$  converged when the fixed point map  $G$  was a contraction and that the rate of convergence was no worse than Picard iteration. There were several variations on that theme in [35] on the depth  $m$ , the choice of norm, and whether or not  $G$  was affine.

We will denote the fixed point by  $u^* = G(u^*)$ , the Anderson iteration for  $G$  by  $\{u_k\}_{k=0}^\infty$ , and we use the conventional notation

$$e = u - u^*.$$

We will need contractivity of  $G$  in this paper as well and state that as a formal assumption.

ASSUMPTION 1.1. *There are  $c \in (0, 1)$ ,  $\gamma > 1$ , and a neighborhood*

$$\mathcal{B}(\rho_G) = \{u \mid \|e\| \leq \rho_G\}.$$

*$G$  is Lipschitz continuously differentiable in  $\mathcal{B}(\rho_G)$  with Lipschitz constant  $\gamma - 1$ , and*

$$\|G(u) - G(v)\| \leq c\|u - v\|$$

*for all  $u, v \in \mathcal{B}(\rho_G)$ .*

Note that if  $\beta_k \equiv \beta$  is fixed throughout the iteration, then Anderson( $m$ ) with  $\beta \neq 0$  is the same iteration as Algorithm 1 as applied to the damped fixed point map

$$G_\beta(u) = (1 - \beta)u + \beta G(u).$$

Hence, if damped Picard converges with damping parameter  $\beta$ , so will Anderson( $m$ ) with mixing parameter  $\beta$ . Given this equivalence, analysis with  $\beta = 1$  extends to the case with arbitrary nonzero mixing parameter by considering the map  $G_\beta$ . We will then set  $\beta = 1$  in the remainder of this paper.

The purpose of this paper is to extend the main result in [35] to the case where there are errors in the evaluation of  $G$ . The theorem in [35], which we state below, requires an additional assumption on the coefficients  $\{\alpha_i^{(k)}\}$ .

ASSUMPTION 1.2. *There exists some constant  $M_\alpha$  such that for all  $0 \leq k \leq m$*

$$(3) \quad \sum_{i=0}^{m_k} |\alpha_i^{(k)}| \leq M_\alpha.$$

The result from [35] of interest here is the following theorem.

**THEOREM 1.** *Assume that Assumptions 1.1 and 1.2 hold. Let  $\{u_k\}$  be the Anderson iterations. Let  $\hat{c} \in (c, 1)$  be given. Then there is  $\rho_A$  such that for all  $u_0 \in \mathcal{B}(\rho_A)$*

$$(4) \quad \|e_k\| \leq \hat{c}^k \frac{1+c}{1-c} \|e_0\| \leq \rho_A$$

and

$$(5) \quad \|F_k\| \leq \hat{c}^k \|F_0\|.$$

One may interpret (5) as the statement that Anderson performs “no worse” than Picard if the initial iterate is sufficiently good. The interpretation of (4) is in line with the standard results that relate the nonlinear residual to the error [15]. If one estimates the condition number of the Jacobian of  $F$  by  $u \in \mathcal{B}(\rho_G)$  by

$$\kappa(F'(u)) \leq \frac{1+c}{1-c},$$

then (4) and (5) lead to

$$\frac{\|e_k\|}{\|e_0\|} \leq \kappa(F'(u^*)) \hat{c}^k \approx \kappa(F'(u^*)) \frac{\|F(u_k)\|}{\|F(u_0)\|}.$$

Many of the special cases considered in [35] did not require Assumption 1.2. However, that assumption plays a central role in the analysis in this paper. Hence, the results in this paper do not extend to the all of the special cases we considered in [35].

Our results are in the spirit of the local improvement ideas from [9], the mesh-independence results in [1, 2, 11, 16], and, in the case of stochastic errors, the tracking theorems from [41].

**1.1. Classes of errors.** We will consider two cases, one (deterministic) where there is a uniform bound on the errors and a second (stochastic) where the error can only be estimated probabilistically. This latter case is motivated by applications where there is a Monte Carlo simulation embedded in the evaluation of the fixed point map. In both cases, one computes  $G(u)$  inaccurately and obtains

$$\hat{G}(u) = G(u) + \epsilon(u).$$

The error  $\epsilon$  represents effects such as rounding error, an inaccurate matrix-vector product for linear problems, the effect of internal tolerances from single-physics solves in multiphysics coupling problems, and a Monte Carlo simulation embedded in the function evaluation. Our results describe the effects of the errors in the function evaluation on the performance of  $\text{anderson}(u_0, \hat{G}, m)$  in terms of properties of the error  $\epsilon(u)$ .

In the deterministic case we assume that there is  $\epsilon_0$  such that

$$(6) \quad \|\epsilon(u)\| \leq \epsilon_0$$

for all  $u$ . For the deterministic case we prove a local improvement result of the type proposed in [9]. Here the convergence proceeds like the error-free case until the error reaches a lower bound that prohibits a proof of asymptotic convergence.

We will express the stochastic case as in [41] which considered the effects of an embedded Monte Carlo simulation on Newton and JFNK iterations. We assume that  $\hat{G}$  contains a Monte Carlo simulation with  $N_{MC}$  trials, where  $N_{MC}$  can be adjusted to improve accuracy. Following [41], we make the following assumption.

**ASSUMPTION 1.3.** *There is  $\rho_{MC}$  and a function  $c_G$  such that if  $\|u - u^*\| \leq \rho_{MC}$  and  $\delta > 0$ , then*

$$(7) \quad \text{Prob} \left( \|\epsilon(u, N_{MC})\| > \frac{c_G(\delta)}{\sqrt{N_{MC}}} \right) < \delta.$$

When the number of trials is important we will make the dependence on  $N_{MC}$  explicit,

$$\hat{G}(u, N_{MC}) = G(u) + \epsilon(u, N_{MC}),$$

and suppress it when only  $\epsilon_0$  is relevant.

**2. Theory.** In this section we prove two theorems on local improvement. The deterministic case (see subsection 2.1) is based on the analysis from [35]. This work is related to our previous work on transition from deterministic errors to stochastic errors in the context of Jacobian-Free Newton Krlov (JFNK) methods from [41, 42]. The JFNK analysis was made particularly challenging by the difficulty Krylov methods [31, 32, 33] have with error accumulation due to inaccurate matrix-vector products. Anderson acceleration does not have that problem, and the transition from the deterministic case to the stochastic one is much simpler.

### 2.1. Uniformly bounded errors.

**THEOREM 2.** *Assume that Assumptions 1.1 and 1.2 hold. Let  $\hat{c} \in (c, 1)$ . Then for  $\rho_D$  and  $\epsilon_0$  sufficiently small, (6) and  $u_0 \in \mathcal{B}(\rho_D)$  imply that for all  $k \geq 0$ ,*

$$(8) \quad \|e_k\| \leq \hat{c}^k \frac{1+c}{1-c} \|e_0\| + \left(1 + \frac{\gamma \rho_D}{2(1-c)}\right) (M_\alpha + c) \frac{1-\hat{c}^k}{(1-c)(1-\hat{c})} \epsilon_0 \leq \rho_D,$$

and

$$(9) \quad \|F_k\| \leq \hat{c}^k \|F_0\| + \left(1 + \frac{\gamma \rho}{2(1-c)}\right) (M_\alpha + c) \frac{1-\hat{c}^k}{1-\hat{c}} \epsilon_0.$$

*Proof.* First note that Assumption 1.1 implies that  $F'$  is Lipschitz continuous in  $\mathcal{B}(\rho_G)$  with Lipschitz constant  $\gamma$ .

Let  $\rho_0 \leq \rho_G$  be small enough so that

$$(10) \quad \left(1 + \frac{\gamma \rho}{2(1-c)}\right) \left(c + \frac{\gamma \rho M_\alpha}{2(1-c)} \hat{c}^{-m}\right) \leq \hat{c}.$$

Now set

$$\rho_D = \min \left\{ \rho_0, \left(1 + \frac{\gamma \rho_0}{2(1-c)}\right) \right\}.$$

Assume that

$$(11) \quad \|e_0\| \leq \rho_D$$

and

$$(12) \quad \epsilon_0 \leq \frac{(1-c)(1-\hat{c})}{\left(1 + \frac{\gamma \rho_D}{2(1-c)}\right)(M_\alpha + c)} \rho_D.$$

We proceed by induction. The result is clear for  $k = 0$ . To see this note that (11) implies that  $u_0 \in \mathcal{B}(\rho_D)$ . Inequalities (8) and (9) hold trivially.

Next, let  $k \geq 0$  and suppose that  $\|e_j\| \leq \rho$  and the bounds (8) and (9) hold for each  $0 \leq j \leq k$ . We will first show that  $\|e_{k+1}\| \leq \rho_D$ . Lipschitz continuity of  $F'$  implies that

$$F_{k-m_k+i} = F'(u^*)e_{k-m_k+i} + \Delta_{k-m_k+i}, \quad 0 \leq i \leq m_k,$$

where

$$(13) \quad \|\Delta_{k-m_k+i}\| \leq \frac{\gamma}{2} \|e_{k-m_k+i}\|^2.$$

Our assumptions imply that  $\|F'(u^*)^{-1}\| \leq 1 - c$ . Hence, (13) implies that

$$(14) \quad \begin{aligned} e_{k-m_k+i} &= F'(u^*)^{-1}(F_{k-m_k+i} - \Delta_{k-m_k+i}) \\ &= F'(u^*)^{-1}(\hat{F}_{k-m_k+i} - \Delta_{k-m_k+i} - \epsilon(u_{k-m_k+i})). \end{aligned}$$

The formula for  $u_{k+1}$  and the constraint  $\sum_{i=0}^{m_k} \alpha_i^{(k)} = 1$  imply that

$$u_{k+1} = \sum_{i=0}^{m_k} \alpha_i^{(k)} \hat{G}(u_{k-m_k+i}) = \sum_{i=0}^{m_k} \alpha_i^{(k)} [u_{k-m_k+i} + \hat{F}(u_{k-m_k+i})]$$

and

$$(15) \quad e_{k+1} = \sum_{i=0}^{m_k} \alpha_i^{(k)} [e_{k-m_k+i} + \hat{F}(u_{k-m_k+i})].$$

We substitute (14) into (15) to obtain

$$(16) \quad \begin{aligned} e_{k+1} &= \sum_{i=0}^{m_k} \alpha_i^{(k)} [F'(u^*)^{-1}(\hat{F}_{k-m_k+i} - \Delta_{k-m_k+i} - \epsilon(u_{k-m_k+i})) + \hat{F}(u_{k-m_k+i})] \\ &= F'(u^*)^{-1} \sum_{i=0}^{m_k} \alpha_i^{(k)} [(I + F'(u^*))\hat{F}_{k-m_k+i} - \Delta_{k-m_k+i} - \epsilon(u_{k-m_k+i})] \\ &= F'(u^*)^{-1} \left( G'(u^*) \sum_{i=0}^{m_k} \alpha_i^{(k)} \hat{F}_{k-m_k+i} - \bar{\Delta}_k - E_k \right), \end{aligned}$$

where

$$\bar{\Delta}_k = \sum_{i=0}^{m_k} \alpha_i^{(k)} \Delta_{k-m_k+i} \quad \text{and} \quad E_k = \sum_{i=0}^{m_k} \alpha_i^{(k)} \epsilon(u_{k-m_k+i}).$$

The minimization condition on the coefficients implies that

$$\left\| \sum_{i=0}^{m_k} \alpha_i^{(k)} \hat{F}_{k-m_k+i} \right\| \leq \|\hat{F}_k\|.$$

Hence

$$(17) \quad \|e_{k+1}\| \leq \frac{1}{1-c} \left( c\|\hat{F}_k\| + \|\bar{\Delta}_k\| + \|E_k\| \right).$$

Assumption 1.2 implies that

$$(18) \quad \|E_k\| \leq \sum_{i=0}^{m_k} |\alpha_i^{(k)}| \epsilon_0 \leq M_\alpha \epsilon_0.$$

Since  $\|e_j\| \leq \rho_D$  for each  $0 \leq j \leq k$ , (13) implies that

$$(19) \quad \|\bar{\Delta}_k\| \leq \frac{\gamma}{2} \sum_{i=0}^{m_k} |\alpha_i^{(k)}| \|e_{k-m_k+i}\|^2 \leq \frac{\gamma \rho_D}{2} \sum_{i=0}^{m_k} |\alpha_i^{(k)}| \|e_{k-m_k+i}\|.$$

The induction hypothesis implies that

$$(20) \quad \begin{aligned} \|e_{k-m_k+i}\| &\leq \frac{\|F_{k-m_k+i}\|}{1-c} \\ &\leq \frac{1}{1-c} \left( \hat{c}^{k-m_k+i} \|F_0\| + \left(1 + \frac{\gamma \rho_D}{2(1-c)}\right) \frac{(M_\alpha + c)(1 - \hat{c}^{k-m_k+i})}{1 - \hat{c}} \epsilon_0 \right) \\ &\leq \frac{1}{1-c} \left( \hat{c}^{k-m} \|F_0\| + \left(1 + \frac{\gamma \rho_D}{2(1-c)}\right) \frac{(M_\alpha + c)(1 - \hat{c}^k)}{1 - \hat{c}} \epsilon_0 \right). \end{aligned}$$

Assumption 1.2 and (20) imply that

$$(21) \quad \|\bar{\Delta}_k\| \leq \frac{\gamma \rho_D M_\alpha}{2(1-c)} \left( \hat{c}^{k-m} \|F_0\| + \left(1 + \frac{\gamma \rho_D}{2(1-c)}\right) \frac{(M_\alpha + c)(1 - \hat{c}^k)}{1 - \hat{c}} \epsilon_0 \right).$$

The induction hypothesis for  $\|F_k\|$ , (18), and (21) imply that

$$\begin{aligned} c\|\hat{F}_k\| + \|\bar{\Delta}_k\| + \|E_k\| &\leq c \left( \hat{c}^k \|F_0\| + \left(1 + \frac{\gamma \rho_D}{2(1-c)}\right) (M_\alpha + c) \frac{1 - \hat{c}^k}{1 - \hat{c}} \epsilon_0 \right) \\ &\quad + \frac{\gamma \rho_D M_\alpha}{2(1-c)} \left( \hat{c}^{k-m} \|F_0\| + \left(1 + \frac{\gamma \rho_D}{2(1-c)}\right) \frac{(M_\alpha + c)(1 - \hat{c}^k)}{1 - \hat{c}} \epsilon_0 \right) + (M_\alpha + c) \epsilon_0. \end{aligned}$$

Hence,

$$(22) \quad \begin{aligned} c\|\hat{F}_k\| + \|\bar{\Delta}_k\| + \|E_k\| &\leq \left( c + \frac{\gamma \rho_D M_\alpha}{2(1-c)} \hat{c}^{-m} \right) \hat{c}^k \|F_0\| \\ &\quad + \left[ \left(1 + \frac{\gamma \rho_D}{2(1-c)}\right) \left( c + \frac{\gamma \rho_D M_\alpha}{2(1-c)} \right) (1 - \hat{c}^k) + 1 - \hat{c} \right] \frac{M_\alpha + c}{1 - \hat{c}} \epsilon_0. \end{aligned}$$

We substitute (10) into (22) and obtain

$$(23) \quad \begin{aligned} c\|\hat{F}_k\| + \|\bar{\Delta}_k\| + \|E_k\| &\leq \frac{\hat{c}}{1 + \frac{\gamma \rho_D}{2(1-c)}} \hat{c}^k \|F_0\| + (\hat{c}(1 - \hat{c}^k) + 1 - \hat{c}) \frac{M_\alpha + c}{1 - \hat{c}} \epsilon_0 \\ &= \frac{\hat{c}^{k+1}}{1 + \frac{\gamma \rho_D}{2(1-c)}} \|F_0\| + (1 - \hat{c}^{k+1}) \frac{M_\alpha + c}{1 - \hat{c}} \epsilon_0. \end{aligned}$$

We have, combining (23), (17), (11), and (12),

$$\begin{aligned} \|e_{k+1}\| &\leq \frac{1}{1-c} \left( \frac{\hat{c}^{k+1}}{1 + \frac{\gamma \rho_D}{2(1-c)}} \|F_0\| + (1 - \hat{c}^{k+1}) \frac{M_\alpha + c}{1 - \hat{c}} \epsilon_0 \right) \\ &\leq \frac{(1+c)\hat{c}^{k+1}}{(1-c)(1 + \frac{\gamma \rho_D}{2(1-c)})} \|e_0\| + (1 - \hat{c}^{k+1}) \frac{M_\alpha + c}{(1-c)(1 - \hat{c})} \epsilon_0 \\ &\leq \hat{c}^{k+1} \rho_D + (1 - \hat{c}^{k+1}) \rho_D = \rho_D. \end{aligned}$$

Hence  $u_{k+1} \in \mathcal{B}(\rho_D)$ . We will complete the proof by showing that (8) and (9) hold for  $k+1$ . Since

$$\begin{aligned} F_{k+1} &= \int_0^1 F'(u^* + te_{k+1})e_{k+1} dt \\ &= \int_0^1 F'(u^* + te_{k+1})F'(u^*)^{-1} dt \left( G'(u^*) \sum_{i=0}^{m_k} \alpha_i^{(k)} \hat{F}_{k-m_k+i} - \bar{\Delta}_k - E_k \right), \end{aligned}$$

we have

$$(24) \quad \|F_{k+1}\| \leq \int_0^1 \|F'(u^* + te_{k+1})F'(u^*)^{-1}\| dt \left( c\|\hat{F}_k\| + \|\bar{\Delta}_k\| + \|E_k\| \right).$$

Lipschitz continuity of  $F'$  in  $\mathcal{B}(\rho_D)$  implies that

$$\begin{aligned} \|F'(u^* + te_{k+1})F'(u^*)^{-1}\| &= \|I - [F'(u^*) - F'(u^* + te_{k+1})]F'(u^*)^{-1}\| \\ &\leq 1 + \|F'(u^*) - F'(u^* + te_{k+1})\| \|F'(u^*)^{-1}\| \\ &\leq 1 + \frac{\gamma t \|e_{k+1}\|}{1-c}. \end{aligned}$$

Therefore,

$$(25) \quad \int_0^1 \|F'(u^* + te_{k+1})F'(u^*)^{-1}\| dt \leq 1 + \frac{\gamma \|e_{k+1}\|}{2(1-c)} \leq 1 + \frac{\gamma \rho_D}{2(1-c)}.$$

We combine (25), (24), and (23) to obtain

$$\begin{aligned} \|F_{k+1}\| &\leq \left( 1 + \frac{\gamma \rho_D}{2(1-c)} \right) (c\|F_k\| + \|\bar{\Delta}_k\| + \|E_k\|) \\ &\leq \left( 1 + \frac{\gamma \rho_D}{2(1-c)} \right) \left( \frac{\hat{c}^{k+1}}{1 + \frac{\gamma \rho_D}{2(1-c)}} \|F_0\| + (1 - \hat{c}^{k+1}) \frac{M_\alpha + c}{1 - \hat{c}} \epsilon_0 \right) \\ (26) \quad &= \hat{c}^{k+1} \|F_0\| + \left( 1 + \frac{\gamma \rho_D}{2(1-c)} \right) (M_\alpha + c) \frac{1 - \hat{c}^{k+1}}{1 - \hat{c}} \epsilon_0. \end{aligned}$$

Thus (9) holds for  $k+1$ , and (8) follows directly from (9) and

$$\begin{aligned} \|e_{k+1}\| &\leq \frac{\|F_{k+1}\|}{1-c} \\ &\leq \hat{c}^{k+1} \frac{\|F_0\|}{1-c} + \left( 1 + \frac{\gamma \rho_D}{2(1-c)} \right) (M_\alpha + c) \frac{1 - \hat{c}^{k+1}}{(1-c)(1-\hat{c})} \epsilon_0 \\ (27) \quad &\leq \hat{c}^{k+1} \frac{1+c}{1-c} \|e_0\| + \left( 1 + \frac{\gamma \rho_D}{2(1-c)} \right) (M_\alpha + c) \frac{1 - \hat{c}^{k+1}}{(1-c)(1-\hat{c})} \epsilon_0. \end{aligned}$$

This completes the induction.  $\square$

Compare the factor  $(1-c)^{-2}$  in the second term on the left side of (8) to the single power of  $(1-c)^{-1}$  in (9). This indicates that ill conditioning in the presence of a large error bound  $\epsilon_0$  in  $G$  can lead to a poor result even if the residual is small. Contrast this with Theorem 1, where  $(1-c)^{-1}$  is in (4) and not in (5) at all.

**2.2. Random errors.** The result in this section is a direct consequence of Theorem 2 and an application of the logic in [42].

**THEOREM 3.** *Assume that Assumptions 1.1, 1.2, and 1.3 hold. Let an integer  $K > 0$  and  $\omega \in (0, 1)$  be given; then if  $N_{MC}$  is sufficiently large,  $\rho_S$  sufficiently small, and  $u_0 \in \mathcal{B}(\rho_S)$ , the conclusions of Theorem 2 hold with probability no less than  $1 - \omega$  for all  $0 \leq k \leq K$ .*

*Proof.* Let  $\rho_D$  and  $\epsilon_0$  be small enough so that the conclusions of Theorem 2 hold. Let

$$\rho_S = \min(\rho_D, \rho_{MC}),$$

where  $\rho_{MC}$  is from Assumption 1.3. We will simultaneously tune  $\delta$  and  $N_{MC}$  by

$$\delta = 1 - (1 - \omega)^{1/K} \quad \text{and} \quad N_{MC} = \frac{\epsilon_0}{c_G(\delta)}.$$

Then (7) implies that, for any given  $k$ ,

$$\text{Prob}(\|\epsilon(u_k, N_{MC})\| \leq \epsilon_0) \geq 1 - \delta = (1 - \omega)^{1/K}.$$

Hence,

$$\text{Prob}(\|\epsilon(u_k, N_{MC})\| \leq \epsilon_0 \text{ for all } 0 \leq k \leq K - 1) \geq (1 - \delta)^K = (1 - \omega).$$

Hence we may now invoke Theorem 2 to complete the proof.  $\square$

In the next section we will present three examples that illustrate the principal consequence of the theory, that the iteration history will track that of the error-free iteration until the errors cause the iteration to stagnate. In the example in subsection 3.1 we artificially control the error to illustrate this point. The subsequent examples are more realistic.

**3. Examples.** In this section we consider three examples. The first is a simple example where we can directly control the stochastic errors. The second two are more realistic. One has deterministic error, and the other has an embedded Monte Carlo simulation in the fixed point map.

**3.1. Chandrasekhar H-equation.** This example is based on one from [35]. We will solve a discretization of the Chandrasekhar H-equation [4, 6]

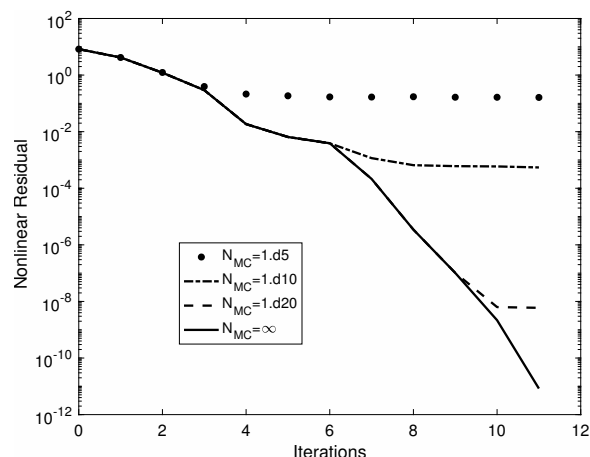
$$(28) \quad H(\mu) = G(H) \equiv \left( 1 - \frac{\omega}{2} \int_0^1 \frac{\mu}{\mu + \nu} H(\nu) d\nu \right)^{-1}.$$

In (28)  $\omega \in [0, 1]$  is a parameter, and one seeks a solution  $H^* \in C[0, 1]$ .

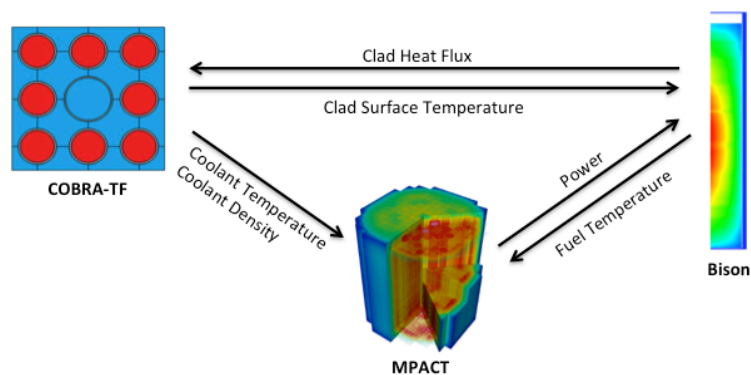
In this example we set  $\omega = .99$  and discretize with the composite midpoint rule on an equally spaced mesh with 500 intervals. We introduce errors by applying a relative perturbation to the output of the fixed point map that is normally distributed with mean zero and standard deviation  $1/\sqrt{N_{MC}}$ . We solve the problem with Anderson(1). In Figure 1 we plot the convergence history for four values of  $N_{MC}$ , where  $N_{MC} = \infty$  is the error-free case.

In the figure one can clearly see that the iterations for the problem with errors track the iteration history for the exact problem until the errors in the fixed point map cause stagnation. This is what the theory predicts. In this example, the rate of convergence of Picard iteration is  $\omega = .99$ , and, as reported in [35], Anderson acceleration does much better than that.



FIG. 1. *Anderson with errors.*

**3.2. Deterministic coupling with Tiamat.** Here we consider the effect of varying internal solve tolerances in the context of coupled multiphysics problems. These errors are deterministic. The specific coupling we consider is called Tiamat [25]. This is a nuclear reactor physics code designed for fuel rod analysis. The codes participating in this coupling, and the dependencies between these codes, are illustrated in Figure 2. Tiamat can solve this coupled problem by either Picard iteration or Anderson acceleration, and we will examine the effect that varying the quality of the solution produced by these single-physics codes has on the behavior of these coupled solution methods.

FIG. 2. *Participating codes and intercode dependencies for the Tiamat coupling.*

**3.2.1. Single-physics codes.** We will first overview the single-physics codes in this coupling. First, Bison is a single-rod fuel performance code which models the mechanical, thermal, and chemical behavior within a fuel rod [13]. This code is built upon the MOOSE framework [12], and it utilizes finite element discretization and solves the resulting system by JFNK. The system of equations that this code solves is necessarily time dependent, so we denote the set of equations which Bison solves

within a given time step by

$$(29) \quad f_B(x_B, T_c, q) = 0.$$

This equation represents coupled equations for heat conduction and equilibrium mechanics. This is solved for the state variables  $x_B$ , which represent the displacement field and temperature distribution. The other quantities in this equation are coupling parameters which depend on the solutions to the other sets of physics and affect the solution to this equation.  $T_c$  is the cladding surface temperature and  $q$  is the fission heat generation rate, which act as a Dirichlet boundary condition and source in the heat equation, respectively. As Bison models a single fuel rod, a separate instance of Bison is used to model each rod in the reactor core.

Next, the neutronics code MPACT models neutron transport in the reactor core and the resulting power generation due to fission [8]. This behavior is governed by the neutron transport equation. MPACT approximates the three-dimensional (3D) transport equation by decomposing it into a coupled set of two-dimensional (2D) equations in the radial direction, which are solved by the method of characteristics, and one-dimensional (1D) equations in the axial direction. We denote the discretized form of the transport equation solved by MPACT as

$$(30) \quad f_M(x_M, T_f, T_w, \rho_w) = 0.$$

Unlike Bison, this equation is steady-state. In this, the state variables  $x_M$  represent the dominant eigenvalue and the scalar flux, which is a measure of neutron intensity. The quantities  $T_f$ ,  $T_w$ , and  $\rho_w$  are the fuel temperature, coolant temperature, and coolant density, respectively. This equation's dependence on these material properties is due to their effect on coefficients in the transport equation called cross sections.

Finally, the thermal-hydraulics code COBRA-TF models coolant flow through the reactor core [29]. This code uses a two-fluid, three-field representation of two-phase flow. The fields modeled are liquid, vapor, and entrained droplets. The governing equations are a coupled system representing conservation of mass, momentum, and energy in each of the fields. Like MPACT, we seek a steady-state solution to this system. While this code always solves the time dependent form of the equations, it includes a "pseudo-steady-state" mode which employs a time marching scheme until the solution has been judged to have sufficiently approached steady-state. We denote the steady-state system of equations that COBRA-TF solves by

$$(31) \quad f_C(x_C, q'') = 0.$$

In this,  $x_C$  represents the density, enthalpy, and velocity for each field and the pressure.  $q''$  is the heat flux from the fuel to the coolant, so this represents a source in the energy equations.

It then remains to describe how the coupling parameter vectors are computed. Again, these quantities indicate dependencies between the codes, represented as shown in Figure 2. Notationally, we represent computation of coupling parameter vectors by

$$(32) \quad T_f = r_{M,B}(x_B), \quad \begin{pmatrix} T_w \\ \rho_w \end{pmatrix} = \begin{pmatrix} r_{M,C,T}(x_C) \\ r_{M,C,\rho}(x_C) \end{pmatrix}, \quad T_c = r_{B,C}(x_C),$$

$$q = r_{B,M}(x_M), \quad q'' = r_{C,B}(x_B).$$

These transfer functions simply map responses computed from the solution to one set of physics to inputs for another.

**3.2.2. Coupled problem formulation and solution.** Given the notation introduced above, we will now define the coupled problem that is solved in Tiamat and show how it is solved. At each time step in the simulation, we are interested in finding solutions to the single-physics systems such that each residual equation is simultaneously solved. That is, we seek to solve the monolithic system

$$(33) \quad F \begin{pmatrix} x_B \\ x_C \\ x_M \end{pmatrix} = \begin{pmatrix} f_B(x_B, T_c, q) \\ f_C(x_C, q'') \\ f_M(x_M, T_f, T_w, \rho_w) \end{pmatrix} = 0,$$

such that the constraints imposed by the transfer functions are satisfied.

---

**Algorithm 2.** Block Gauss–Seidel nonlinear solve for Tiamat.

---

- 1: Given  $x_B^0, x_C^0$ , and  $x_M^0$ .
  - 2: **for**  $k = 0, 1, \dots$  until converged **do**
  - 3:   Transfer Bison to MPACT,  $T_f^k = r_{M,B}(x_B^k)$ .
  - 4:   Transfer COBRA-TF to MPACT,  $T_w^k = r_{M,C,T}(x_C^k)$  and  $\rho_w^k = r_{M,C,\rho}(x_C^k)$ .
  - 5:   Solve  $f_M(x_M, T_f^k, T_w^k, \rho_w^k) = 0$  for  $x_M^{k+1}$ .
  - 6:   Transfer MPACT to Bison,  $q^{k+1} = r_{B,M}(x_M^{k+1})$ .
  - 7:   Transfer COBRA-TF to Bison,  $T_c^k = r_{B,C}(x_C^k)$ .
  - 8:   Solve  $f_B(x_B, T_c^k, q^{k+1}) = 0$  for  $x_B^{k+1}$ .
  - 9:   Transfer Bison to COBRA-TF,  $q''_{k+1} = r_{C,B}(x_B^{k+1})$ .
  - 10:   Solve  $f_C(x_C, q''_{k+1}) = 0$  for  $x_C^{k+1}$ .
  - 11: **end for**
- 

Couplings of this type are frequently solved by Picard iteration [14, 23, 43]. In Algorithm 2 we present a block Gauss–Seidel Picard iteration scheme for solving this coupled system in which we repeatedly solve the single-physics systems in a given order and transfer updated coupling parameters as they are obtained. To apply Anderson acceleration to this iteration scheme, the fixed point map  $G$  must be explicitly stated. We apply Anderson acceleration to a fixed point problem defined in terms of the transferred coupling parameter data. Denoting the solution to (29) given  $T_c$  and  $q$  as  $x_B(T_c, q)$ , the solution to (30) given  $T_f, T_w$ , and  $\rho_w$  as  $x_M(T_f, T_w, \rho_w)$ , and the solution to (31) given  $q''$  as  $x_C(q'')$ , we define

$$(34) \quad G \begin{pmatrix} T_f \\ T_w \\ \rho_w \\ T_c \end{pmatrix} = \begin{pmatrix} r_{M,B}(x_B(T_c, r_{B,M}(x_M(T_f, T_w, \rho_w)))) \\ r_{M,C,T}(x_C(r_{C,B}(x_B(T_c, r_{B,M}(x_M(T_f, T_w, \rho_w)))))) \\ r_{M,C,\rho}(x_C(r_{C,B}(x_B(T_c, r_{B,M}(x_M(T_f, T_w, \rho_w)))))) \\ r_{B,C}(x_C(r_{C,B}(x_B(T_c, r_{B,M}(x_M(T_f, T_w, \rho_w)))))) \end{pmatrix}.$$

By finding a fixed point of this map, the solutions  $x_B$ ,  $x_C$ , and  $x_M$  that are computed internally are a solution to the fully coupled system (33). Evaluating this map involves solving the single-physics systems in the order MPACT, Bison, and then COBRA-TF, so given initial iterates  $T_f^0 = r_{M,B}(x_B^0)$ ,  $T_w^0 = r_{M,C,T}(x_C^0)$ ,  $\rho_w^0 = r_{M,C,\rho}(x_C^0)$ , and  $T_c^0 = r_{B,C}(x_C^0)$ , the iteration  $u_{k+1} = G(u_k)$  produces the same sequence as Algorithm 2.

As the components of the fixed point map  $G$  represent various physical quantities, they may exist on vastly different scales. In particular, with the chosen units, temperature unknowns are several orders of magnitude larger than the density unknowns. As a result of this, small changes in temperature unknowns may contribute

much more significantly to the least-squares problem in Anderson acceleration than relatively larger changes in density. To address this issue, rather than the original problem  $u = G(u)$ , we attempt to solve the scaled fixed point problem

$$(35) \quad v = MG(M^{-1}v) \equiv H(v),$$

where  $v = Mu$  are scaled variables. We define the diagonal scaling matrix  $M$  by

$$(36) \quad M = \begin{pmatrix} \text{diag}(T_f^0)^{-1} & & & \\ & \text{diag}(T_w^0)^{-1} & & \\ & & \text{diag}(\rho_w^0)^{-1} & \\ & & & \text{diag}(T_c^0)^{-1} \end{pmatrix}.$$

As a result of an initialization process, described in [25], we have fairly good initial iterates, so this should provide reasonably good scaling.

**3.2.3. Single-physics tolerances.** Convergence in Tiamat is judged on a local and global level. By global convergence, we refer to convergence of the coupled system. This convergence is judged based on small changes in various responses between coupled iterations. The specific criteria are described in more detail in [25]. Local convergence refers to the convergence of a single solve of an individual physics code. Because the individual sets of physics are only solved approximately, the fixed point map  $G$  is never evaluated exactly. The size of the error in the evaluation of  $G$  is directly related to the size of the error in the computed solutions for each of the single-physics systems, and each code features several parameters involved in judging convergence. Bison solves its system by JFNK, so it judges convergence based on the nonlinear residual  $f_B$ . Determinations of convergence are based upon a combination of relative and absolute tolerances. MPACT features an inner fixed point iteration to solve for the scalar flux and dominant eigenvalue. The convergence of this code is based upon small relative changes in the scalar flux distribution and absolute changes in the eigenvalue between iterations. Finally, with COBRA-TF operating in pseudo-steady-state mode, it determines steady-state convergence based on sufficient smallness of several quantities. The only criterion we will consider in the following section is called the global energy balance tolerance. This criterion measures the difference between energy gains and losses in the system as a percentage of energy gains, which should be zero at steady-state. A complete description of the criteria for achieving steady-state can be found in [29].

**3.2.4. Problem setup.** We now consider numerical tests consisting of simulation of a single fuel rod for one time step. We compare the behavior of Anderson acceleration and Picard with various levels of accuracy in the three single-physics codes. In the following cases we have a single  $\text{UO}_2$  fuel rod at full power. The rod has height 3.66m and operates at a power level of 67kW. The inlet coolant has temperature 559F, 600ppm dissolved boron, and pressure 2250psi. These physical parameters are chosen to be realistic for a typical fuel rod in a pressurized water reactor. The coupling parameter vectors are volume-averaged over 49 axial regions, so Anderson is solving a fixed point problem for 196 unknowns. The individual codes may use a finer mesh internally. Damping is required for Picard to converge, so we choose the mixing parameter  $\beta = 0.5$  for both Picard and Anderson. For Anderson acceleration we used a storage depth of  $m = 2$ . We will analyze convergence behavior for both methods in terms of the scaled fixed point residual  $F(v) = H(v) - v$ , with  $H$  defined as in (35). While this is a fairly small problem, it has been observed that for this sort of problem

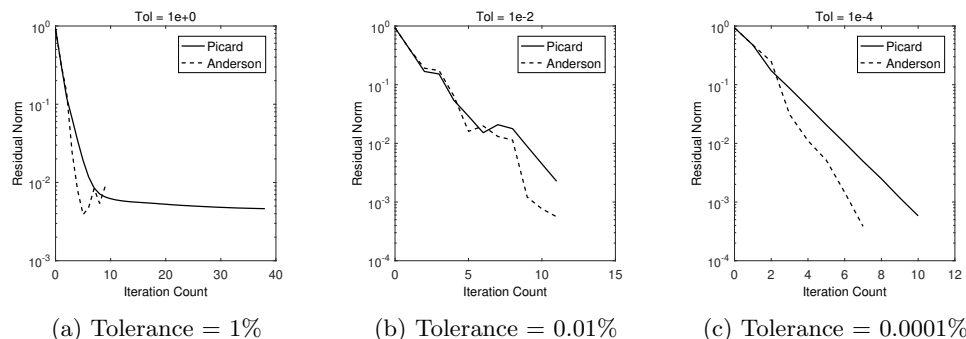


FIG. 3. Varying CTF global energy balance tolerance in Tiamat single-rod tests.

single-rod simulations are often sufficient for analyzing coupled convergence behavior, as convergence difficulties occur primarily due to behavior in the axial direction [36].

We consider variation in the global energy balance tolerance for COBRA-TF, the scalar flux tolerance for MPACT, and the relative residual tolerance for Bison. When varying the tolerance for one code, the other tolerances for that code are set to levels such that only the specified criterion is determining convergence, and the tolerances for the other codes are set to their tightest values.

**3.2.5. Numerical results.** Results from varying the COBRA-TF global energy balance tolerance are given in Figure 3. Figure 3a shows results from utilizing a very loose 1% tolerance. During the coupled solve, this resulted in COBRA-TF declaring convergence after only one time step in each solve. With such a large error, both Picard and Anderson approach a stagnation point, and the residuals for both methods seem to stagnate about approximately the same level. While the theory present above suggests that Anderson may stagnate at a larger residual norm than Picard, this is not observed here. The residuals for Anderson simply appear more jagged about the same stagnation point as Picard. We lastly note that Anderson declares convergence, but it does so at a point far from the actual solution. This is because global convergence is based upon small changes in responses between coupled iterations, so false convergence may be declared if the iteration stagnates. A more accurate judgment of global convergence could be made based upon the true residuals (29)–(31), but obtaining this data would have required significant, invasive code modification. This illustrates a potential danger of this black-box approach to code coupling. In the remaining figures, we reduce the tolerance to 0.01% and 0.0001%. With the moderately loose 0.01% tolerance, Anderson and Picard achieve coupled convergence in the same number of iterations, and with the tightest 0.0001% tolerance Anderson significantly outperforms Picard. For both methods there is an improvement as the tolerance is reduced from 0.01% to 0.0001%, but the improvement is far more significant for Anderson.

Figure 4 shows results from solving the coupled problem with both Picard and Anderson with various MPACT scalar flux convergence tolerances. The results are fairly similar to what was observed when varying the COBRA-TF convergence criteria. In Figure 4a we utilize a loose tolerance of  $1e-1$ , and we observe convergence up until a stagnation point is reached for both methods. As before, this stagnation point occurs at roughly the same residual norm size for both Picard and Anderson, so it does not seem that Anderson amplifies the error in the function evaluation significantly more than Picard. In the remaining figures we consider an intermediate tolerance

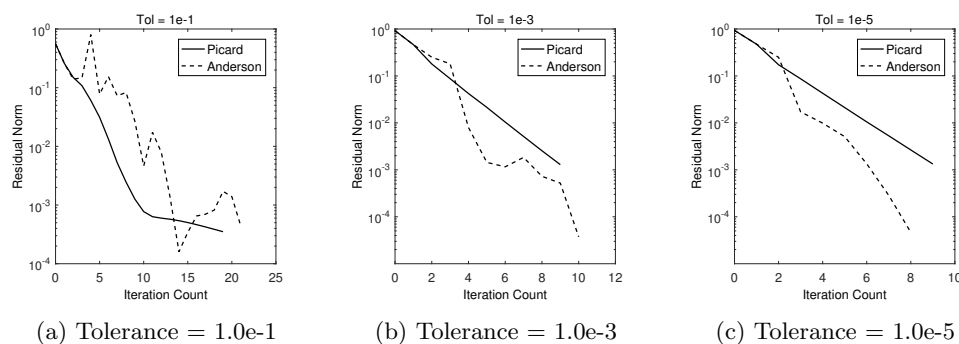


FIG. 4. Varying MPACT scalar flux tolerance in Tiamat single-rod tests.

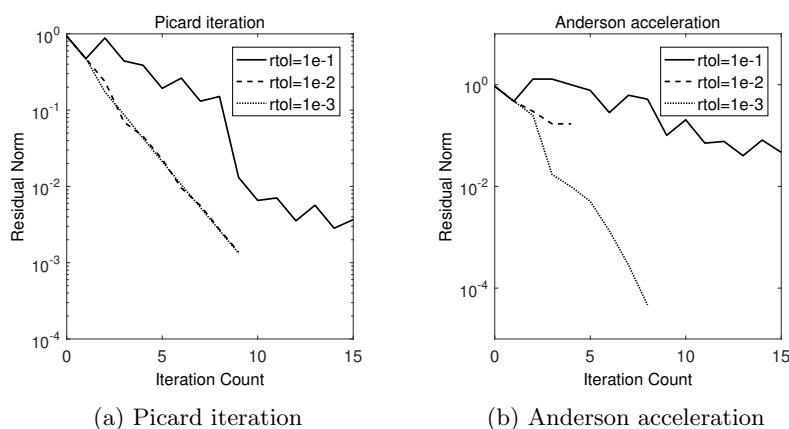


FIG. 5. Varying Bison relative residual tolerance in Tiamat single-rod tests.

of  $10^{-3}$  and a tight tolerance of  $10^{-5}$ . With tolerance  $10^{-3}$ , we see that Anderson and Picard perform comparably, and with tolerance  $10^{-5}$  Anderson again performs noticeably better than Picard. The Picard curves in these two figures are nearly identical, so it seems that with tolerance  $10^{-3}$  the fixed point map error introduced by MPACT is small enough to be negligible. Conversely, Anderson achieves a significant improvement as the tolerance is reduced from  $10^{-3}$  to  $10^{-5}$ . This again suggests that Anderson may require a smaller fixed point map error than Picard to perform well.

Finally, Figure 5 shows results from solving the coupled system while varying the Bison relative residual tolerance. We observe that with the loosest tolerance of  $10^{-1}$  both Picard and Anderson perform quite poorly. Anderson seems to make no progress, while Picard achieves some reduction in the residual and then begins to stagnate. Reducing the relative tolerance to  $10^{-2}$ , Picard performs significantly better. In fact, there is little visible difference between this curve and the tighter tolerance of  $10^{-3}$ , so it seems that even with a fairly loose relative residual tolerance of  $10^{-2}$  the Bison solution is accurate enough so as to not affect Picard significantly. Conversely, for Anderson there is a significant difference between the curves for each of the tolerances. It initially seems as though Anderson is performing better with relative tolerance  $10^{-2}$  than with  $10^{-1}$ , but it quickly reaches a stagnation point and declares false convergence. Only with the tolerance reduced to  $10^{-3}$  does Anderson

perform well, and little difference is observed if the tolerance is reduced beyond this point. Hence, we again see that Anderson may be more noticeably affected than Picard by the presence of large errors in the evaluation of the fixed point map.

**3.3. Neutron transport.** In this section we compare Anderson acceleration with Picard iteration applied to the linear Boltzmann equation for neutron transport. In particular we consider the case where the transport equation is solved using a Monte Carlo approach and is therefore subject to stochastic errors as described in subsection 2.2.

**3.3.1. Theory.** The steady-state, energy-independent, first-order linear Boltzmann transport equation for isotropic scattering is

$$(37) \quad \mu \frac{\partial \psi}{\partial x} + \Sigma_t \psi(x, \mu) = \frac{1}{2} [\Sigma_s \phi(x) + q(x)],$$

where  $\psi$  is the neutron angular flux,  $\phi$  is the scalar flux given by

$$(38) \quad \phi(x) = \int_{-1}^1 \psi(x, \mu') d\mu',$$

$\mu$  is the angular variable,  $\Sigma_t$  is the total cross section,  $\Sigma_s$  is the scattering cross section, and  $q$  is a source term. Boundary conditions are given by  $\psi(0, \mu) = \psi(L, -\mu) = 0$  for  $\mu > 0$ .

We consider two different problem formulations. The first is obtained by inverting the operator on the left-hand side of (37) and integrating over angle to produce an equation purely in terms of  $\phi$ :

$$(39) \quad \phi = \int_{-1}^1 d\mu \left[ \mu \frac{\partial}{\partial x} + \Sigma_t \right]^{-1} \frac{1}{2} (\Sigma_s \phi + q),$$

which suggests the fixed point map

$$(40) \quad \phi^{k+1} = \int_{-1}^1 d\mu \left[ \mu \frac{\partial}{\partial x} + \Sigma_t \right]^{-1} \frac{1}{2} (\Sigma_s \phi^k + q).$$

This iterative process is often referred to as source iteration [19]. This problem can also be written as a nonlinear equation for  $\phi$ :

$$(41) \quad F(\phi) = \phi - \int_{-1}^1 d\mu \left[ \mu \frac{\partial}{\partial x} + \Sigma_t \right]^{-1} \frac{1}{2} (\Sigma_s \phi + q).$$

If a Monte Carlo method is used for the operator inversion and integration over angle, no angular discretization is necessary.

The second problem formulation is based on the Nonlinear Diffusion Acceleration (NDA) form of the neutron transport equation [17, 40, 41, 42]. Here, the “high-order” problem given by (37) is accelerated using the “low-order” diffusion equation

$$(42) \quad \frac{d}{dx} \left[ -\frac{1}{3\Sigma_t} \frac{d\phi}{dx} + \hat{D}^{HO} \phi \right] + (\Sigma_t - \Sigma_s) \phi = q(x).$$

The coefficient  $\hat{D}$  is given by

$$(43) \quad \hat{D} = \frac{J^{HO} + \frac{1}{3\Sigma_t} \frac{d\phi^{HO}}{dx}}{\phi^{HO}},$$

where  $\phi^{HO}$  is the scalar flux of the high-order problem and  $J^{HO}$  is the current of the high-order problem given by

$$(44) \quad J^{HO}(x) = \int_{-1}^1 \psi(x, \mu') \mu' d\mu'.$$

This selection for the diffusion coefficient forces consistency between the high-order and low-order formulations. We can express the problem as a nonlinear equation for  $\phi$ :

$$(45) \quad F(\phi) = \frac{d}{dx} \left[ -\frac{1}{3\Sigma_t} \frac{d\phi}{dx} + \hat{D}^{HO}(\phi)\phi \right] + (\Sigma_t - \Sigma_s)\phi - q.$$

We write  $\hat{D}^{HO}(\phi)$  to demonstrate the dependence of  $\hat{D}^{HO}$  on  $\phi$  as is seen in (43), in which  $\phi^{HO}$  and  $J^{HO}$  are recovered from the solution to (37). We discretize the low-order problem with second-order central differences and transfer the output of the Monte Carlo solve of the high-order problem by tallying fluxes and currents within each spatial cell. We refer the reader to [17, 40, 42] for the details of the discretization. As with the original fixed point map, the high-order problem is solved using a Monte Carlo approach. For the numerical results, the Monte Carlo solution is done with the Profugus code [10], an open source, multigroup Monte Carlo transport solver developed at Oak Ridge National Laboratory.

**3.3.2. Numerical results.** We consider two monoenergetic fixed source problems in single material slab geometry with isotropic scattering [19] to illustrate the theory. We consider the moderate and high scattering cases where the scattering ratio  $c = \Sigma_s/\Sigma_t = .80$  and  $.99$ . We apply source iteration and nonlinear diffusion acceleration to each case with Profugus as the high-order solve. We then accelerate both of these methods using Anderson(2). The iterations' residual behaviors are then observed.

In each of the Profugus computations we fix the number of histories. For  $c = .80$ , we have  $N_{MC} = 1e6$  during each step of the iteration. The same is repeated for  $c = .99$  but instead with  $N_{MC} = 4e6$  histories. The high scattering case requires a higher number of particles because the mean spectral radius is near unity. For a deterministic transport sweep the spectral radius of source iteration is bounded above by the scattering ratio. Moderate stochastic error in the function evaluations can lead to poor stability or even divergence. This occurs because we may violate Assumption 1.1 if the error in the stochastic transport sweep is too large, ultimately leading to loss of contractivity.

For the computations reported in Table 1 and Figure 6 we use the parameter values from [41] along with the moderate scattering case.

TABLE 1  
Problem data from [41].

Parameter	Moderate scattering	High scattering
$\Sigma_t$	10	10
$\Sigma_s$	8.0	9.9
$\tau$	1	1
$q$	.5	.5
Spatial cells	50	50

In Figures 6a and 6b, we plot the average relative residuals over 10 independent runs on the  $y$ -axis and the cumulative number of particle histories on the  $x$ -axis, with



each marker indicating an iterate's residual. We ran each of the iterations until the error in the function evaluation stagnated the residual reduction.

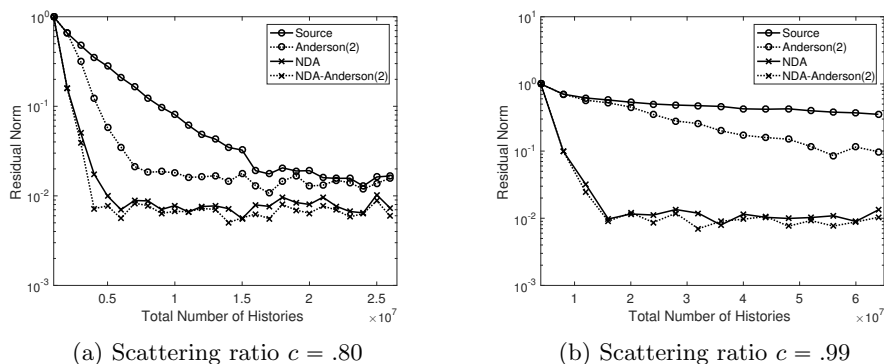


FIG. 6. Residual histories of Profugus dependent fixed source calculations.

As is shown, Anderson(2) was able to accelerate source iteration and further provide speed up of nonlinear diffusion acceleration, though marginal. In Figure 6a, Anderson(2) quickly reduces the residual relative to source iteration until the noise in the function evaluation dominates. Source iteration reaches the same stagnating residual as Anderson(2), but the convergence is slower, as expected. Comparing NDA and NDA-Anderson(2), we observe that the two methods have roughly the same residual history and stagnate simultaneously. As the theory states, Anderson does no worse than Picard iteration if the noise is well bounded with high probability.

In Figure 6b, source iteration now with the higher scattering ratio has a far harder time reducing the residual each iteration. Anderson(2) lowers the spectral radius of source iteration and achieves faster convergence. Though not plotted in the high scattering case, source iteration and its Anderson accelerated version stagnate at roughly 70 and 30 iterations, respectively. As for the previous case, NDA and NDA-Anderson(2) follow the same convergence and both stagnate at the same residual.

Finally, we observed negligible differences in iteration histories storing more than two Anderson vectors and only slight incremental improvement on intermediate iterates using Anderson(2) over Anderson(1). We also have observed that once stagnation is reached, the Anderson accelerated versions are more stable.

**4. Conclusions.** We prove local improvement theorems for Anderson acceleration for problems in which the fixed point map is corrupted by deterministic or stochastic noise. Our results predict stagnation in the iteration once the residuals are at the level of the noise. We present numerical results that support the predictions of the theory.

#### REFERENCES

- [1] E. L. ALLGOWER AND K. BÖHMER, *Application of the mesh independence principle to mesh refinement strategies*, SIAM J. Numer. Anal., 24 (1987), pp. 1335–1351, <https://doi.org/10.1137/0724086>.
- [2] E. L. ALLGOWER, K. BÖHMER, F. A. POTRA, AND W. C. RHENBOLDT, *A mesh-independence principle for operator equations and their discretizations*, SIAM J. Numer. Anal., 23 (1986), pp. 160–169, <https://doi.org/10.1137/0723011>.

- [3] D. G. ANDERSON, *Iterative procedures for nonlinear integral equations*, J. ACM, 12 (1965), pp. 547–560.
- [4] I. W. BUSBRIDGE, *The Mathematics of Radiative Transfer*, Cambridge Tracts 50, Cambridge University Press, Cambridge, UK, 1960.
- [5] N. N. CARLSON AND K. MILLER, *Design and application of a gradient-weighted moving finite element code I: In one dimension*, SIAM J. Sci. Comput., 19 (1998), pp. 728–765, <https://doi.org/10.1137/S106482759426955X>.
- [6] S. CHANDRASEKHAR, *Radiative Transfer*, Dover, New York, 1960.
- [7] A. M. COLLIER, A. C. HINDMARSH, R. SERBAN, AND C. S. WOODWARD, *User Documentation for KINSOL v2.8.0*, Tech. Report UCRL-SM-208116, Lawrence Livermore National Laboratory, Livermore, CA, 2015.
- [8] B. COLLINS, T. DOWNAR, J. GEHIN, A. GODFREY, A. GRAHAM, D. JABAAY, B. KELLEY, K. CLARNO, K. KIM, B. KOCHUNAS, E. LARSEN, Y. LIU, Z. LIU, W. MARTIN, S. PALMTAG, M. ROSE, T. SALLER, S. STIMPSON, T. TRAHAN, J. WANG, W. WIESELQUIST, M. YOUNG, AND A. ZHU, *MPACT Theory Manual*, Tech. Report CASL-U-2015-0078-000, Consortium for Advanced Simulation of LWRs, 2015.
- [9] J. E. DENNIS AND H. F. WALKER, *Inaccuracy in quasi-Newton methods: Local improvement theorems*, in Mathematical Programming at Oberwolfach II, Mathematical Programming Study 22, North-Holland, Amsterdam, 1984, pp. 70–85.
- [10] T. EVANS, S. HAMILTON, AND S. SLATTERY, *ORNL-CEES/Profugus*, Oak Ridge National Laboratory, 2016, <https://github.com/ORNL-CEES/Profugus>.
- [11] W. R. FENG AND C. T. KELLEY, *Mesh independence of matrix-free methods for path following*, SIAM J. Sci. Comput., 21 (2000), pp. 1835–1850, <https://doi.org/10.1137/S1064827598339360>.
- [12] D. GASTON, C. NEWMAN, AND G. HANSEN, *MOOSE: A parallel computational framework for coupled systems of nonlinear equations*, Nuclear Engineering and Design, 239 (2009), pp. 1768–1778, <https://doi.org/10.1016/j.nucengdes.2009.05.021>.
- [13] J. HALES, S. NOVASCONE, G. PASTORE, D. PEREZ, B. SPENCER, AND R. WILLIAMSON, *BISON Theory Manual*, Tech. Report October, Fuels Modeling and Simulation Department, Idaho National Laboratory, 2013, [https://neup.inl.gov/SiteAssets/FY2014%20Documents/BISON-theory\\_manual.pdf](https://neup.inl.gov/SiteAssets/FY2014%20Documents/BISON-theory_manual.pdf).
- [14] S. HAMILTON, M. BERRILL, K. CLARNO, R. PAWLOWSKI, A. TOTH, C. T. KELLEY, T. EVANS, AND B. PHILIP, *An assessment of coupling algorithms for nuclear reactor core physics simulations*, J. Comput. Phys., 311 (2016), pp. 241–257.
- [15] C. T. KELLEY, *Iterative Methods for Linear and Nonlinear Equations*, Frontiers Appl. Math. 16, SIAM, Philadelphia, 1995.
- [16] C. T. KELLEY AND E. W. SACHS, *Mesh independence of Newton-like methods for infinite dimensional problems*, J. Int. Eq. Appl., 3 (1991), pp. 549–573.
- [17] D. A. KNOLL, H. PARK, AND K. SMITH, *Application of the Jacobian-free Newton-Krylov method to nonlinear acceleration of transport source iteration in slab geometry*, Nuclear Sci. Eng., 167 (2011), pp. 122–132.
- [18] K. N. KUDIN, G. E. SCUSERIA, AND E. CANCEÈ, *A black-box self-consistent field convergence algorithm: One step closer*, J. Chem. Phys., 116 (2002), pp. 8255–8261, <https://doi.org/10.1063/1.1470195>.
- [19] E. LEWIS AND W. MILLER, *Computational Methods of Neutron Transport*, American Nuclear Society, LaGrange Park, IL, 1993.
- [20] L. LIN AND C. YANG, *Elliptic preconditioner for accelerating the self-consistent field iteration in Kohn–Sham density functional theory*, SIAM J. Sci. Comput., 35 (2013), pp. S277–S298, <https://doi.org/10.1137/120880604>.
- [21] P. A. LOTT, H. F. WALKER, C. S. WOODWARD, AND U. M. YANG, *An accelerated Picard method for nonlinear systems related to variably saturated flow*, Adv. Water Res., 38 (2012), pp. 92–101.
- [22] K. MILLER, *Nonlinear Krylov and moving nodes in the method of lines*, J. Comput. Appl. Math., 183 (2005), pp. 275–287.
- [23] L. MONTI AND T. SCHULENBERG, *Coupled ERANOS/TRACE system for HPLWR 3 pass core analyses*, in International Conference on Mathematics, Computational Methods and Reactor Physics, Saratoga Springs, NY, 2009, pp. 1–14.
- [24] C. W. OOSTERLEE AND T. WASHIO, *Krylov subspace acceleration for nonlinear multigrid with application to recirculating flows*, SIAM J. Sci. Comput., 21 (2000), pp. 1670–1690, <https://doi.org/10.1137/S1064827598338093>.

- [25] R. PAWLOWSKI, K. CLARNO, R. MONTGOMERY, R. SALKO, T. EVANS, J. TURNER, AND D. GASTON, *Design of a high fidelity core simulator for analysis of pellet clad interaction*, in Proceedings of the ANS MC2015 – Joint International Conference on Mathematics and Computation (M&C), Supercomputing in Nuclear Applications (SNA) and the Monte Carlo (MC) Method, vol. ANS MC2015 CD, Nashville, TN, 2015.
- [26] P. PULAY, *Convergence acceleration of iterative sequences. The case of SCF iteration*, Chem. Phys. Lett., 73 (1980), pp. 393–398.
- [27] P. PULAY, *Improved SCF convergence acceleration*, J. Comput. Chem., 3 (1982), pp. 556–560.
- [28] 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.
- [29] R. SALKO AND M. AVRAMOVA, *CTF Theory Manual*, Tech. Report CASL-U-2015-0054-000, Consortium for Advanced Simulation of LWRs, 2015.
- [30] R. SCHNEIDER, T. ROHWEDDER, A. NEELOV, AND J. BLAUERT, *Direct minimization for calculating invariant subspaces in density functional computations of the electronic structure*, J. Comput. Math., 27 (2008), pp. 360–387.
- [31] V. SIMONCINI AND D. B. SZYLD, *Flexible inner-outer Krylov subspace methods*, SIAM J. Numer. Anal., 40 (2003), pp. 2219–2239, <https://doi.org/10.1137/S0036142902401074>.
- [32] V. SIMONCINI AND D. B. SZYLD, *Theory of inexact Krylov subspace methods and applications to scientific computing*, SIAM J. Sci. Comput., 25 (2003), pp. 454–477, <https://doi.org/10.1137/S1064827502406415>.
- [33] V. SIMONCINI AND D. B. SZYLD, *Recent computational developments in Krylov subspace methods for linear systems*, Numer. Linear Algebra Appl., 14 (2007), pp. 1–59.
- [34] *SUNDIALS (SUite of Nonlinear and Differential/ALgebraic Solvers)*, <http://www.llnl.gov/casc/sundials>.
- [35] A. TOTH AND C. T. KELLEY, *Convergence analysis for Anderson acceleration*, SIAM J. Numer. Anal., 53 (2015), pp. 805–819, <https://doi.org/10.1137/130919398>.
- [36] A. TOTH, C. T. KELLEY, S. SLATTERY, S. HAMILTON, K. CLARNO, AND R. PAWLOWSKI, *Analysis of Anderson acceleration on a simplified neutronics/thermal hydraulics system*, in Proceedings of the ANS MC2015 – Joint International Conference on Mathematics and Computation (M&C), Supercomputing in Nuclear Applications (SNA) and the Monte Carlo (MC) Method, vol. ANS MC2015 CD, Nashville, TN, 2015.
- [37] A. TOTH AND R. PAWLOWSKI, *NOX::Solver::AndersonAcceleration Class Reference*, [https://trilinos.org/docs/dev/packages/nox/doc/html/classNOX\\_1\\_1Solver\\_1\\_1AndersonAcceleration.html](https://trilinos.org/docs/dev/packages/nox/doc/html/classNOX_1_1Solver_1_1AndersonAcceleration.html) (2015).
- [38] H. W. WALKER AND P. NI, *Anderson acceleration for fixed-point iterations*, SIAM J. Numer. Anal., 49 (2011), pp. 1715–1735, <https://doi.org/10.1137/10078356X>.
- [39] T. WASHIO AND C. OOSTERLEE, *Krylov subspace acceleration for nonlinear multigrid schemes*, Electron. Trans. Numer. Anal., 6 (1997), pp. 271–290.
- [40] J. WILLERT, *Hybrid Deterministic/Monte Carlo Methods for Solving the Neutron Transport Equation and  $k$ -Eigenvalue Problem*, Ph.D. thesis, North Carolina State University, Raleigh, NC, 2013.
- [41] J. WILLERT, X. CHEN, AND C. T. KELLEY, *Newton’s method for Monte Carlo-based residuals*, SIAM J. Numer. Anal., 53 (2015), pp. 1738–1757, <https://doi.org/10.1137/130905691>.
- [42] J. WILLERT, C. T. KELLEY, D. A. KNOLL, AND H. K. PARK, *Hybrid deterministic/Monte Carlo neutronics*, SIAM J. Sci. Comput., 35 (2013), pp. S62–S83, <https://doi.org/10.1137/120880021>.
- [43] J. YAN, B. KOCHUNAS, M. HURSIN, T. DOWNAR, Z. KAROUTAS, AND E. BAGLIETTO, *Coupled computational fluid dynamics and MOC neutronic simulations of Westinghouse PWR fuel assemblies with grid spacers*, in 14th International Topical Meeting on Nuclear Reactor Thermalhydraulics, Toronto, Ontario, 2011.