

SUPPLEMENTARY MATERIALS: ALGEBRAIC ERROR ANALYSIS FOR MIXED-PRECISION MULTIGRID SOLVERS*

STEPHEN F. MCCORMICK[†], JOSEPH BENZAKEN[‡], AND RASMUS TAMSTORF[‡]

SM1. Double relaxation sweeps. The theory in the main paper is restricted to $V(1, 0)$ -cycles, meaning that each \mathcal{V} uses one *pre-smoothing* sweep on the way down through the coarse grids and no *post-smoothing* sweeps on the way back up to the finest. This specific cycle simplifies the analysis because post-smoothing sweeps substantially complicate the estimates by the accumulation of errors in the transfer of residuals to the coarse levels and, to a lesser extent, because a single pre-smoothing sweep is easier to analyze than multiple sweeps. On the other hand, while the analysis of a general V -cycle would be too complicated to present here, we can more easily analyze multiple pre-smoothing sweeps as we illustrate now for a $V(2, 0)$ -cycle.

A relatively simple way to handle multiple sweeps is to combine them into one. To this end, for each $j \in \{1, 2\}$, consider a monotonically energy-convergent stationary linear iteration $x \leftarrow x - M^{(j)}(Ax - b)$, where $M^{(j)} \in \mathbb{R}^{n \times n}$, and let α_j be a constant such that computing $M^{(j)}z$ for a vector $z \in \mathbb{R}^n$ in $\hat{\epsilon}$ -precision yields the result $M^{(j)}z + \delta$, $\|\delta\| \leq \alpha_j \hat{\epsilon} \|z\|$. Then the key point here is that the error propagation matrix for relaxation with preconditioner $M^{(1)}$ followed by relaxation with preconditioner $M^{(2)}$ can be written as $(I - M^{(2)}A)(I - M^{(1)}A) = I - M^{(2)}A$, where $M^{(2)} = M^{(1)} + M^{(2)} - M^{(2)}AM^{(1)}$. We can therefore think of, and implement, two relaxation sweeps as just the *one* sweep $y \leftarrow y - M^{(2)}(Ay - r)$, which means that we can analyze a $V(2, 0)$ -cycle as just a $V(1, 0)$ -cycle with this $M^{(2)}$. The implication is that we just need to provide estimates for $\|M^{(2)}\|$ and a constant $\alpha_{M^{(2)}}$ such that computing $M^{(2)}z$ for any vector z in $\hat{\epsilon}$ -precision yields the result $M^{(2)}z + \delta_{M^{(2)}}$, $\|\delta_{M^{(2)}}\| \leq \alpha_{M^{(2)}} \hat{\epsilon} \|z\|$. This is done in our next theorem.

THEOREM SM1.1. Double Sweeps. *For a double sweep in the ordering specified by $M^{(2)}z = ((M^{(2)}z) + (M^{(1)}z) - (M^{(2)}(A(M^{(1)}z))))$, the constant $\alpha_{M^{(2)}}$ can be chosen as follows:*

$$\begin{aligned} \alpha_{M^{(2)}} &= \|M^{(2)}\| + (1 + \hat{\epsilon})[(\|M^{(2)}\| + \alpha_2 \hat{\epsilon})(\|A\| \alpha_1 + \psi \dot{m}_A \alpha_1 \hat{\epsilon} + \psi \dot{m}_A \|M^{(1)}\|)] \\ \text{(SM1.1)} \quad &+ \|A\| \cdot \|M^{(1)}\| \alpha_2 + \|M^{(1)}\| + \|M^{(2)}\| + 2\alpha_1 + 2\alpha_2. \end{aligned}$$

Moreover, in general, $\|M^{(2)}\| \leq \|M^{(1)}\| + \|M^{(2)}\| + \|A\| \cdot \|M^{(1)}\| \cdot \|M^{(2)}\|$. The special case $M^{(1)} = M^{(2)} = \frac{\omega}{\|A\|}I$, $0 < \omega < 2$ yields the sharp estimate $\|M^{(2)}\| \leq \frac{2\omega}{\|A\|}$.

Proof. The following is meant to clarify the stages for computing $M^{(2)}z$, with the

*Received by the editors July 1, 2020; accepted for publication (in revised form) February 18, 2021; published electronically June 24, 2021.

[†]University of Colorado at Boulder, Boulder, CO (stephen.mccormick@colorado.edu).

[‡]Walt Disney Animation Studios, Burbank, CA (Joseph.Benzaken@disneyanimation.com, Rasmus.Tamstorf@disneyanimation.com).

subscripted w 's standing for the indicated quantities computed in $\dot{\epsilon}$ -precision:

$$M^{(2)}z = \underbrace{\underbrace{(M^{(2)}z)}_{w_4} + \underbrace{(M^{(1)}z)}_{w_1}}_{w_5} - \underbrace{\underbrace{(M^{(2)}(A(M^{(1)}z)))}_{w_1}}_{w_2}.$$

Thus, by definition,

$$(SM1.2) \quad w_1 = M^{(1)}z + \delta_1, \quad \|\delta_1\| \leq \alpha_1 \dot{\epsilon} \|z\|.$$

We then use (2.3) in the main paper to obtain the slight overestimate

$$w_2 = AM^{(1)}z + A\delta_1 + \delta_2, \quad \|\delta_2\| \leq \psi \dot{m}_A \dot{\epsilon} \|w_1\| = \psi \dot{m}_A \dot{\epsilon} \|M^{(1)}z + \delta_1\|.$$

Letting $\delta_3 = A\delta_1 + \delta_2$, where $\|\delta_3\| \leq \|A\| \cdot \|\delta_1\| + \psi \dot{m}_A \dot{\epsilon} (\|\delta_1\| + \|M^{(1)}z\|)$, yields

$$w_2 = AM^{(1)}z + \delta_3, \quad \|\delta_3\| \leq \left(\|A\| + \psi \dot{m}_A \dot{\epsilon} \right) \alpha_1 + \psi \dot{m}_A \|M^{(1)}\| \dot{\epsilon} \|z\|.$$

Similarly,

$$w_3 = M^{(2)}AM^{(1)}z + M^{(2)}\delta_3 + \delta_4, \quad \|\delta_4\| \leq \alpha_2 \dot{\epsilon} \|w_2\| = \alpha_2 \dot{\epsilon} \|AM^{(1)}z + \delta_3\|.$$

Letting $\delta_5 = M^{(2)}\delta_3 + \delta_4$, where $\|\delta_5\| \leq \|M^{(2)}\| \cdot \|\delta_3\| + \alpha_2 \dot{\epsilon} (\|AM^{(1)}z\| + \|\delta_3\|)$, yields $w_3 = M^{(2)}AM^{(1)}z + \delta_5$, $\|\delta_5\| \leq \Upsilon \dot{\epsilon} \|z\|$, where

$$(SM1.3) \quad \Upsilon = (\|M^{(2)}\| + \alpha_2 \dot{\epsilon}) (\|A\| \alpha_1 + \psi \dot{m}_A \alpha_1 \dot{\epsilon} + \psi \dot{m}_A \|M^{(1)}\|) + \|A\| \cdot \|M^{(1)}\| \alpha_2 \dot{\epsilon} \|z\|.$$

We can now use the estimate $w_4 = M^{(2)}z + \delta_6$, $\|\delta_6\| \leq \alpha_2 \dot{\epsilon} \|z\|$ together with (SM1.2) and (SM1.3) to obtain

$$\begin{aligned} w_5 &= (w_1 + w_4 - w_3) + \delta_7 + \delta_8 = M^{(2)}z + \delta_1 - \delta_5 + \delta_6 + \delta_7 + \delta_8, \\ \|\delta_7\| &\leq \dot{\epsilon} \|w_1 + w_4\| \leq (\|M^{(1)}\| + \|M^{(2)}\| + \alpha_1 + \alpha_2) \dot{\epsilon} \|z\|, \\ \|\delta_8\| &\leq \dot{\epsilon} \|w_1 + w_4 - w_3 + \delta_7\| \leq \dot{\epsilon} (\|M^{(2)}\| \|z\| + \|\delta_1\| + \|\delta_5\| + \|\delta_6\| + \|\delta_7\|). \end{aligned}$$

Letting $\delta = \delta_1 - \delta_5 + \delta_6 + \delta_7 + \delta_8$, then $w_5 = M^{(2)}z + \delta$, where

$$\begin{aligned} \|\delta\| &\leq \|\delta_1\| + \|\delta_5\| + \|\delta_6\| + \|\delta_7\| + \|\delta_8\| \\ &\leq \left(\dot{\epsilon} \|M^{(2)}\| + (1 + \dot{\epsilon}) (\|\delta_1\| + \|\delta_5\| + \|\delta_6\| + \|\delta_7\|) \right) \|z\| \\ &\leq \left(\|M^{(2)}\| + (1 + \dot{\epsilon}) [\alpha_1 + \alpha_2 + \Upsilon + \alpha_1 + \alpha_2] \right) \dot{\epsilon} \|z\|, \end{aligned}$$

thus establishing (SM1.1). The estimates for $\|M^{(2)}\|$ are straightforward. \square

SM2. Second-order Chebyshev iteration. The equation in (SM1.1) can be used in a recursive way to analyze any Krylov method, where the error propagation matrix is a polynomial in A . For example, it is fairly straightforward to show that $\alpha_{M^{(2)}} = \mathcal{O}\left(\frac{\dot{m}_A}{\|A\|}\right)$ for the \mathcal{K}^{th} -order Chebyshev relaxation (cf., [SM1]), although the constant in this order bound depends exponentially on \mathcal{K} . On the other hand, a more

direct approach can achieve a somewhat tighter bound, as illustrated for the case $\mathcal{K} = 2$ in our next theorem.

Second-order Chebyshev relaxation can be formed from two sweeps of Richardson iteration with error propagation factors of the form $I - s_j A$, $j = 1, 2$. Assume that the coefficients in the Chebyshev factors are chosen to so that $0 < s_j = \mathcal{O}(\frac{1}{\|A\|})$, $j = 1, 2$. This assumption would generally hold in the multigrid context when the smoothing interval is chosen to be a fixed percentage of the upper spectrum of A . We can thus write Chebyshev iteration in the form $y \leftarrow y - M_C(Ay - r)$, where $M_C = \omega_1 I - \omega_2 A$, $\omega_1 = \mathcal{O}(\frac{1}{\|A\|}) > 0$, and $\omega_2 = \mathcal{O}(\frac{1}{\|A\|^2}) > 0$. Define $\dot{m}_A = \frac{m_A}{1 - m_A \dot{\epsilon}}$, where m_A is the maximum number of nonzeros in the rows of A . Finally, suppose that computing $M_C z$ in $\dot{\epsilon}$ -precision for any vector z yields $M_C z + \delta_{M_C}$, $\|\delta_{M_C}\| \leq \alpha_{M_C} \dot{\epsilon} \|z\|$, for some constant α_{M_C} .

THEOREM SM2.1. Chebyshev. *For one second-order Chebyshev iteration in the ordering specified by $M_C z = (\omega_1 z) - (\omega_2 (Az))$, we can choose*

$$\alpha_{M_C} = \|M_C\| + (\omega_1 + (1 + \dot{\epsilon})\omega_2 \psi \dot{m}_A + \omega_2 \|A\|)(1 + \dot{\epsilon}).$$

Note that $\|M_C\| = \mathcal{O}(\frac{1}{\|A\|})$ and, if $\psi \approx \mathcal{O}(\|A\|)$, then $\alpha_{M_C} = \mathcal{O}(\frac{\dot{m}_A}{\|A\|})$.

Proof. Computing $M_C z$ according to $M_C z = \underbrace{(\omega_1 z)}_{w_1} - \underbrace{(\omega_2 (Az))}_{w_2}$, we have

$$\underbrace{\underbrace{w_1}_{w_3}}_{w_4}$$

$$\begin{aligned} w_1 &= \omega_1 z + \delta_1, & \|\delta_1\| &\leq \omega_1 \dot{\epsilon} \|z\|, \\ w_2 &= Az + \delta_2, & \|\delta_2\| &\leq \psi \dot{m}_A \dot{\epsilon} \|z\|, \\ w_3 &= -\omega_2 w_2 + \delta_3, & \|\delta_3\| &\leq \omega_2 \dot{\epsilon} \|Az + \delta_2\| \leq \omega_2 \dot{\epsilon} (\|A\| + \psi \dot{m}_A \dot{\epsilon}) \|z\|, \\ w_4 &= w_1 + w_3 + \delta_4, & \|\delta_4\| &\leq \dot{\epsilon} \|w_1 + w_3\| = \dot{\epsilon} \|M_C z + \delta_1 + \omega_2 \delta_2 + \delta_3\|. \end{aligned}$$

This implies that $w_4 = M_C z + \delta_C$, $\delta_C = \delta_1 - \omega_2 \delta_2 + \delta_3 + \delta_4$. The theorem now follows from noting that

$$\begin{aligned} \|\delta_C\| &\leq \|\delta_1 - \omega_2 \delta_2 + \delta_3\| + \dot{\epsilon} \|M_C z + \delta_1 + \omega_2 \delta_2 + \delta_3\| \\ &\leq \|M_C\| \|z\| \dot{\epsilon} + (\|\delta_1\| + \|\omega_2 \delta_2\| + \|\delta_3\|)(1 + \dot{\epsilon}). \end{aligned}$$

Remark SM2.2. Chebyshev iteration based on A preconditioned by its diagonal D can be formed from two sweeps of damped Jacobi with error propagation operators $I - s_j D^{-1} A = I - M^{(j)} A$, where $M^{(j)} = M_{ii}^{(j)} = \mathcal{O}(\frac{\kappa(D)}{\|A\|})$, $j = 1, 2$. To estimate α_{M_C} for this case, we can mimic the proof of Theorem SM2.1, but with the understanding now that the ω_j are matrices: $\omega_1 = M^{(2)} + M^{(1)}$ and $\omega_2 = M^{(2)} A M^{(1)}$. If the diagonal matrices $M^{(j)}$ have been formed accurately beforehand, perhaps in the setup phase at higher precision, then the line of reasoning is much the same as the above proof with two extra steps to account for the increased complexity of M_C . The resulting estimate is of the same order as that in Theorem SM2.1 with the exception that a power of $\kappa(D)$ appears in the implied constant. When the diagonal entries of A are widely varying, it may be more effective to construct $D^{-1} A$ in higher precision before it is used in V-cycles in order to avoid an explicit dependence on the condition number of D appearing in the estimate for α_{M_C} .

REFERENCES

- [SM1] M. ADAMS, M. BREZINA, J. HU, AND R. TUMINARO, *Parallel Multigrid Smoothing: Polynomial Versus Gauss–Seidel*, *Journal of Computational Physics*, 188 (2003), pp. 593–610, [https://doi.org/10.1016/S0021-9991\(03\)00194-3](https://doi.org/10.1016/S0021-9991(03)00194-3).