

# Selected Fixed Point Problems and Algorithms

*Ch. Booniasirivat, K. Sikorski,<sup>a</sup> Ch. Xiong*

UUCS-07-007

<sup>a</sup>Partially supported by DOE under the C-SAFE center

School of Computing  
University of Utah  
Salt Lake City, UT 84112 USA

March 6, 2007

## *Abstract*

We present a new version of the almost optimal Circumscribed Ellipsoid Algorithm (CEA) for approximating fixed points of nonexpanding Lipschitz functions. We utilize the absolute and residual error criteria with respect to the second norm. The numerical results confirm that the CEA algorithm is much more efficient than the simple iteration algorithm whenever the Lipschitz constant is close to 1. We extend the applicability of the CEA algorithm to larger classes of functions that may be globally expanding, however are nonexpanding/contracting in the direction of fixed points. We also develop an efficient hyper-bisection/secant hybrid method for combustion chemistry fixed point problems.

Keywords: fixed point problems, optimal algorithms, nonlinear equations, ellipsoid algorithm, computational complexity.

## 1 Introduction

We want to approximate a fixed point  $\mathbf{x}^*$  that satisfies the equation  $\mathbf{x}^* = f(\mathbf{x}^*)$  for a given function  $f$ . Many problems can be formulated as fixed-point problems. For example, a root-finding problem for nonlinear equations  $f(\mathbf{x}^*) = 0$  can be rearranged as

$g(\mathbf{x}^*) = f(\mathbf{x}^*) + \mathbf{x}^* = \mathbf{x}^*$ , which is a fixed-point problem. The applications of fixed-point problems include economic equilibria Scarf and Hansen (1973); Border (1985), game theory Scarf (1967b); Garcia et al. (1973); Scarf and Hansen (1973); Border (1985); Yang (1999), boundary value problems Jeppson (1972); Allgower (1977); Yamamoto (1998); Andreas (2003), and chaos and dynamical systems Corless et al. (1990); Szrednicki (2000).

A number of fixed-point theorems have been derived in the last century. Each theorem is focused on a specific class of functions. Banach's fixed-point theorem Banach (1922) states that any contractive function of  $\mathbb{R}^n$  into itself has a unique fixed point. Banach demonstrated that the simple iteration algorithm (SIA),  $\mathbf{x}_{i+1} = f(\mathbf{x}_i)$ , generates a Cauchy sequence that converges to the fixed point of any such function. More general fixed-point theorems focused on continuous functions only. Brouwer demonstrated in Brouwer (1912) that any continuous function from a nonempty, convex and compact subset of  $\mathbb{R}^n$  into itself has at least one fixed point. However, Brouwer's proof using topological arguments was nonconstructive. In 1967 Scarf (1967a) developed a simplicial algorithm for approximating fixed points for Brouwer's maps from a simplex into itself. This was the first constructive proof of Brouwer's fixed-point theorem. Since then, fixed-point computation has become an intensive research area and many new algorithms have been proposed including restart methods Merrill (1972a,b), homotopy methods Eaves (1972); Eaves and Saigal (1972); Allgower and Georg (1990) and ellipsoid algorithms Sikorski et al. (1993); Tsay (1994); Huang et al. (1999). As we mentioned earlier, a fixed-point problem  $\mathbf{x}^* = f(\mathbf{x}^*)$  can be transformed into a root-finding problem  $g(\mathbf{x}^*) = f(\mathbf{x}^*) - \mathbf{x}^* = 0$ . Therefore, effective root-finding algorithms such as Newton type methods can also be used to solve fixed-point problems (see Allgower and Georg (1990)).

Algorithms dealing with nonlinear fixed-point computation are iterative methods. We consider algorithms that are based on function evaluations, and assume that at least one function evaluation is required at each iteration. Most of CPU time is usually consumed by function evaluations. Therefore, we can define the worst case *cost* of an algorithm, for a class of functions  $\mathcal{F}$ , as the maximum number of function evaluations required to achieve a prescribed precision for all functions in  $\mathcal{F}$ . The problem *complexity* is defined as the minimal cost among all possible algorithms for the class  $\mathcal{F}$ . An algorithm is almost *optimal* if and only if its cost is almost minimal. The complexity depends on the class of functions and the selected error criterion, as shown in Sikorski and Woźniakowski (1987); Hirsch et al. (1989); Sikorski (1989, 2001). The contractive functions with contraction factor  $\rho$  close to 1 and nonexpanding functions ( $\rho = 1$ ) are considered first, and then a larger class of functions that are contractive/nonexpanding only in the direction of fixed points (but may be globally expanding), see Vassin and Eremin (2005). The contraction/nonexpansion property is defined with respect to the *Euclidean* norm. The  $\varepsilon$ -approximations of fixed points are obtained by using the *absolute* or *residual* error criteria.

In this paper we present a new version of the Circumscribed Ellipsoid Algorithm (CEA), that does not require the dimensional deflation procedure of Tsay (1994), and as a consequence has the cost lower by a factor of  $n$ . Our CEA algorithm enjoys almost optimal cost  $\approx 2n^2 \ln \frac{1}{\varepsilon}$  for obtaining residual solutions  $\mathbf{x} : \|f(\mathbf{x}) - \mathbf{x}\|_2 \leq \varepsilon$  of nonexpanding functions, and exhibits the cost  $\approx 2n^2(\ln \frac{1}{\varepsilon} + \ln \frac{1}{1-\rho})$  for obtaining absolute solutions  $\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\|_2 \leq \varepsilon$  for contractive functions with factor  $\rho < 1$ . We stress that these bounds hold for both of the considered classes of functions.

We also introduce a univariate hyper-bisection/secant (HBS) method for approximating fixed points of certain combustion chemistry problems. That algorithm enjoys the average case cost  $O(\log \log(1/\varepsilon))$  for computing  $\varepsilon$ -absolute solutions.

We stress that for the infinity-norm case, Shellman and Sikorski (2002) developed a Bisection Envelope Fixed point algorithm (BEFix). Shellman and Sikorski (2003a) also developed a Bisection Envelope Deep-cut Fixed point algorithm (BEDFix) for computing fixed points of two-dimensional nonexpanding functions. Those algorithms exhibit the optimal complexity of  $2\lceil \log_2(1/\varepsilon) \rceil + 1$ . They also developed a recursive fixed point algorithm (PFix) for computing fixed points of  $n$ -dimensional nonexpanding functions with respect to the infinity norm (see Shellman and Sikorski (2003b, 2005)). We note that the complexity of finding  $\varepsilon$ -residual solutions for expanding functions with  $\rho > 1$  is exponential  $O((\rho/\varepsilon)^{(n-1)})$ , as  $\varepsilon \rightarrow 0$ , Hirsch et al. (1989); Deng and Chen (2005).

## 2 Classes of Functions

Given the domain  $B^n = \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\| \leq 1\}$ , the  $n$ -dimensional unit ball, we consider the class of Lipschitz continuous functions

$$\mathcal{B}_{\rho \leq 1}^n \equiv \{f : B^n \rightarrow B^n : \|f(\mathbf{x}) - f(\mathbf{y})\| \leq \rho \|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in B^n\}, \quad (1)$$

where  $n \geq 2$ ,  $\|\cdot\| = \|\cdot\|_2$ , and  $0 < \rho \leq 1$ . In the case when  $0 < \rho < 1$ , the class of functions is denoted by  $\mathcal{B}_{\rho < 1}^n$ . The existence of fixed points of functions in  $\mathcal{B}_{\rho \leq 1}^n$  is assured by the Brouwer's theorem.

In this article, we present the circumscribed ellipsoid algorithm (CEA) that, for every  $f \in \mathcal{B}_{\rho < 1}^n$ , computes an *absolute*  $\varepsilon$ -approximation  $\mathbf{x}$  to  $\mathbf{x}^*$ ,  $\|\mathbf{x} - \mathbf{x}^*\| \leq \varepsilon$ , and for every  $f \in \mathcal{B}_{\rho \leq 1}^n$ , computes a *residual*  $\varepsilon$ -approximation  $\mathbf{x}$ ,  $\|f(\mathbf{x}) - \mathbf{x}\| \leq \varepsilon$ .

We also extend the applicability of the CEA algorithm to larger classes of functions considered by Vassin and Eremin (2005). We indicate that the complexity bounds for the CEA

algorithm do not change in this case. Those larger classes were investigated for problems defined by differential and integral equations originating in geophysics, atmospheric research, material science, and image deblurring Vassin (2005); Vassin and Ageev (1995); Ageev et al. (1997); Vassin and Sereznikova (2004); Perestonina et al. (2003). These problems were effectively solved by Feyer type iterative methods and/or some general optimization techniques, however no formal complexity bounds were derived. These classes are defined by:

$$\mathcal{B}_{\rho \leq 1}^\alpha \equiv \{f : B^n \rightarrow B^n : \text{the set of fixed points of } f, S(f) \neq \emptyset, \quad (2)$$

$$\text{and } \forall \mathbf{x} \in B^n, \text{ and } \forall \alpha \in S(f), \text{ we have } \|f(\mathbf{x}) - f(\alpha)\| \leq \rho \|\mathbf{x} - \alpha\|\},$$

where  $n \geq 2$ ,  $\|\cdot\| = \|\cdot\|_2$ , and  $\rho \leq 1$ . We note that the functions in  $\mathcal{B}_{\rho \leq 1}^\alpha$  may be expanding globally, and therefore the class  $\mathcal{B}_{\rho \leq 1}^n$  is a proper subclass of  $\mathcal{B}_{\rho \leq 1}^\alpha$ .

We finally introduce a univariate combustion chemistry fixed point problem defined in the class:

$$G = \{g : [a, b] \rightarrow [a, b] : \quad (3)$$

$$g(t) = C_4 + \frac{C_5}{\left(\sqrt{\frac{C_1 t^2}{t-C_2} \cdot e^{-\frac{E_c}{Rt}} + C_3} + \sqrt{\frac{C_1 t^2}{t-C_2} \cdot e^{-\frac{E_c}{Rt}}}\right)^2}\},$$

where the interval  $[a, b]$ , and the constants  $C_1, \dots, C_5, R$  and  $E_c$  are defined later in section 6.

We solve this problem in the equivalent zero finding formulation  $f(x) = 0$ , where  $f(x) = x - g(x)$ . We derive the HBS method that is a modification of the bisection/regula-falsi/secant (BRS) method which was proven almost optimal in the average case setting Novak et al. (1995), with the complexity  $O(\log \log(1/\varepsilon))$ . To get  $\varepsilon$ -absolute approximations with  $\varepsilon = 10^{-4}$  we need at most 6 function evaluations in the class  $G$ . Since the number of function evaluations in the HBS method is essentially bounded by the number of function evaluations in the BRS method, we conclude that the average case complexity of HBS is at most  $O(\log \log(1/\varepsilon))$ .

### 3 Constructive Lemmas and Cost Bounds

The following lemma is the basis of the *circumscribed ellipsoid* algorithm (CEA) Sikorski (1989); Sikorski et al. (1993); Tsay (1994); Sikorski (2001). The proof of this lemma can be found in Sikorski et al. (1993); Sikorski (2001).

**Lemma 3.1** *Let  $f \in \mathcal{B}_{\rho < 1}^n$ . Suppose that  $A \subseteq B^n$  contains the fixed point  $\mathbf{x}^*$ . Then, for every  $\mathbf{x} \in A$ ,  $\mathbf{x}^* \in A \cap B^n(\mathbf{c}, \gamma)$ , where*

$$\mathbf{c} = \mathbf{x} + \frac{1}{1 - \rho^2}(f(\mathbf{x}) - \mathbf{x})$$

and

$$\gamma = \frac{\rho}{1 - \rho^2} \|f(\mathbf{x}) - \mathbf{x}\|.$$

The above lemma yields

**Corollary 3.2** *Let  $f \in \mathcal{B}_{\rho < 1}^n$ . If*

$$\|f(\mathbf{x}) - \mathbf{x}\| \leq \frac{1 - \rho^2}{\rho} \varepsilon, \tag{4}$$

*then  $\mathbf{x} + (f(\mathbf{x}) - \mathbf{x})/(1 - \rho^2)$  is an absolute  $\varepsilon$ -approximation to the fixed point  $\mathbf{x}^*$ .*

The following lemma and corollary exhibit upper bounds on the number of iterations of the CEA algorithm.

**Lemma 3.3** *For any  $\delta \in (0, 1)$ , and  $f \in \mathcal{B}_{\rho \leq 1}^n$ , the CEA algorithm requires at most  $i = \lceil 2n(n+1) \cdot \ln\left(\frac{2+\delta}{\delta}\right) \rceil$  iterations to compute  $\mathbf{x}_i \in \mathbb{R}^n$  such that  $\|f(\mathbf{x}_i) - \mathbf{x}_i\| \leq \delta$ , as  $\delta \rightarrow 0$ .*

*Proof:*

We give a sketch of the proof, since it is similar to the proof of Lemma 2.2 of Huang et al. (1999).

The upper bound on the number of function evaluations of the CEA algorithm is obtained by replacing the volume reduction constant  $0.843\gamma^{-2}$  from the Interior Ellipsoid Algorithm of Huang et al. (1999), by the volume-reduction constant of the CEA algorithm given by  $\exp(-1/(2(n+1))) < 1$  Khachiyan (1979); Sikorski (2001). Following the formula (2.3) of the proof in Huang et al. (1999) we get

$$\frac{\delta^n}{(2+\delta)^n} \leq e^{-\frac{i}{2(n+1)}}. \quad (5)$$

The number of iterations that guarantees to obtain an  $\delta$ -residual approximation  $\mathbf{x}_i$  is the smallest  $i$  for which this inequality is violated. Therefore, we get:

$$n \ln \left( \frac{2+\delta}{\delta} \right) \leq i \frac{1}{2(n+1)}, \quad (6)$$

that yields

$$i = \lceil 2n(n+1) \ln \left( \frac{2+\delta}{\delta} \right) \rceil,$$

and completes the proof.

We conclude that the worst case cost of the CEA algorithm is

$$i = O(2n^2 \ln(2/\delta)), \text{ as } \delta \rightarrow 0. \quad (7)$$

We also remark that the bound obtained in the above Lemma is better by a factor of  $n$  than the bound obtained in Tsay (1994).

As a direct corollary from this lemma we get:

**Corollary 3.4** *If  $\rho < 1$ , the CEA algorithm finds an  $\varepsilon$ -approximation  $\mathbf{x}_i$  of the fixed point  $\mathbf{x}^*$  in the absolute sense,  $\|\mathbf{x}_i - \mathbf{x}^*\|_2 \leq \varepsilon$ , within  $i \cong 2n^2(\ln(1/\varepsilon) + \ln(1/(1-\rho)))$  iterations.*

Finally, we derive a constructive lemma for the larger class  $B_{\rho \leq 1}^\alpha$ . This lemma is similar in nature to Lemma 3.1, since after each function evaluation it enables us to locate the set of fixed points in the intersection of the previous set with a certain half space (in Lemma 3.1 it was the intersection of the previous set with a ball). This lemma implies that the CEA algorithm can be also applied to functions in  $\mathcal{B}_{\rho \leq 1}^\alpha$  yielding the same complexity bounds as in the smaller classes  $\mathcal{B}_{\rho \leq 1}^n$ , since all of the arguments in the proof of Lemma 3.3 and Lemma 2.2 of Huang et al. (1999) hold in the classes  $\mathcal{B}_{\rho \leq 1}^\alpha$ .

**Lemma 3.5** *Let  $f \in \mathcal{B}_{\rho \leq 1}^\alpha$  and  $A \subset B^n$  be such that the set of fixed points  $S(f) \subset A$ . Then:*

$$\forall \mathbf{x} \in A, \quad S(f) \subset A \cap H_{\mathbf{x}}, \quad (8)$$

where the halfspace  $H_{\mathbf{x}} = \{\mathbf{y} \in \mathbb{R}^n : (\mathbf{y} - \mathbf{c})^T (f(\mathbf{x}) - \mathbf{c}) \geq 0\}$ , for  $\mathbf{c} = (f(\mathbf{x}) + \mathbf{x})/2$ .

*Proof:*

Suppose on the contrary that there exists  $\alpha \in S(f)$  such that  $(\alpha - \mathbf{c})^T (f(\mathbf{x}) - \mathbf{c}) < 0$ . Then, since  $f(\mathbf{x}) - \mathbf{c} = \mathbf{c} - \mathbf{x}$ , we get:

$$\begin{aligned} \|f(\mathbf{x}) - \alpha\|^2 &= (f(\mathbf{x}) - \mathbf{c} + \mathbf{c} - \alpha)^T (f(\mathbf{x}) - \mathbf{c} + \mathbf{c} - \alpha) \\ &= \|f(\mathbf{x}) - \mathbf{c}\|^2 + \|\mathbf{c} - \alpha\|^2 + 2(\mathbf{c} - \alpha)^T (f(\mathbf{x}) - \mathbf{c}) \\ &> \|f(\mathbf{x}) - \mathbf{c}\|^2 + \|\mathbf{c} - \alpha\|^2 \\ &> \|\mathbf{x} - \mathbf{c}\|^2 + \|\mathbf{c} - \alpha\|^2 - 2(\mathbf{c} - \alpha)^T (\mathbf{c} - \mathbf{x}) \\ &= \|\mathbf{x} - \mathbf{c} + \mathbf{c} - \alpha\|^2 = \|\mathbf{x} - \alpha\|^2, \end{aligned} \quad (9)$$

which contradicts that  $f \in \mathcal{B}_{\rho \leq 1}^\alpha$  and completes the proof.

## 4 Circumscribed Ellipsoid Algorithm

The circumscribed ellipsoid algorithm (CEA) originally designed by Khachiyan (1979) for Linear Programming Problems was developed by Sikorski et al. (1993) for fixed point problems. The details of the CEA algorithm can be found in Tsay (1994) and in a forthcoming (numerically stable) implementation paper Sikorski et al. (2006). The idea behind this algorithm is to construct a sequence of decreasing volume ellipsoids,  $E(\mathbf{x}_k)$ , centered at the evaluation points  $\mathbf{x}_k$ ,  $k = 0, \dots$ , with  $E(\mathbf{x}_0) = B^n$ ,  $\mathbf{x}_0 = \mathbf{0}$ , that enclose the set of fixed points. The volume of an ellipsoid is decreased at each step by a factor

of  $\exp(-1/(2(n+1))) < 1$ , Khachiyan (1979); Sikorski (2001). Both, Lemma 3.3 and Lemma 3.5 allow us to locate the set of fixed points  $S(f)$  in the intersection  $I(k)$  of the current ellipsoid  $E(\mathbf{x}_k)$  with the halfspace  $H(k) = \{\mathbf{y} \in \mathfrak{R}^n : (\mathbf{y} - \mathbf{x}_k)^T (f(\mathbf{x}_k) - \mathbf{x}_k) \geq 0\}$ . Indeed, this is the case, since  $H(k)$  contains as a subset the ball  $B(\mathbf{c}, \gamma)$  of Lemma 3.3 as well as the halfspace  $H_{\mathbf{x}_k}$  of Lemma 3.5. Then, the smallest volume ellipsoid  $E(\mathbf{x}_{k+1})$  circumscribing the set  $I(k)$  is constructed. This process is continued until one of the stopping criteria is satisfied and/or the center of the last ellipsoid is an  $\varepsilon$ -approximation of the fixed point.

The CEA algorithm is summarized in Fig. 1. In each iteration we construct the minimum volume ellipsoid  $E(\mathbf{x}_k)$  that encloses the set of fixed points. The worst case cost of the CEA algorithm with the  $\delta$ -residual error termination is  $i \cong 2n^2 \ln(2/\delta)$  as  $\delta \rightarrow 0$ . This and Corollary 3.4 indicate that the CEA algorithm is more efficient than the SIA algorithm whenever  $n$  is not too large and  $\rho$  is close to 1, since the number of iterations of the SIA algorithm to get an  $\varepsilon$ -absolute solution is  $\lceil \ln(1/\varepsilon) / \ln(1/\rho) \rceil$ .

We now explain all three stopping criteria used in the CEA algorithm. The first one denoted as (1) in Fig. 1 is an absolute error criterion  $\|\mathbf{x}_k - \mathbf{x}^*\| \leq \varepsilon$  where  $\mathbf{x}_k$  is the computed  $\varepsilon$ -approximation. This stopping criterion occurs when the radius of the minimal-volume ellipsoid  $E(\mathbf{x}_k)$  is at most  $\varepsilon$ . The second one denoted as (2) is also an absolute error criterion implied by Corollary 3.2 when the condition  $\|f(\mathbf{x}_k) - \mathbf{x}_k\| \leq (1 - \rho^2)\varepsilon/\rho$  is satisfied. By using this criterion we may reduce the number of function evaluations in some test cases. The last denoted by (3) is a residual error criterion  $\|\mathbf{x}_k - f(\mathbf{x}_k)\| \leq \varepsilon$ . We use it only when  $\rho = 1$ , since for  $\rho < 1$  it may result in premature termination especially if we want to get  $\varepsilon$ -absolute solutions.

We note that each iteration of this algorithm requires construction of the minimum volume ellipsoid that is equivalent (in numerically stable implementation) to computing a positive definite matrix  $A_k$  defining that ellipsoid as well as its eigenvalues  $\lambda_i(A_k)$ ,  $0 < \lambda_1(A_k) \leq \dots \leq \lambda_n(A_k)$  and eigenvectors. The direct construction of the ellipsoids as outlined in the above figure is numerically unstable.

An efficient, stable implementation of this algorithm will appear in a forthcoming paper Sikorski et al. (2006).



```

input  $\{\rho \leq 1; \varepsilon > 0; k_{\max}$  (maximum number of iterations);
function  $f \in B_{\rho \leq 1}^n$  or  $f \in B_{\rho \leq 1}^\alpha$ ;
 $k := 0; \mathbf{x}_0 := \mathbf{0}; A_0 := I_{n \times n}$  (initial ellipsoid  $E(\mathbf{x}_0) = B^n$ );
while  $k \leq k_{\max}$  do begin
    if  $\sqrt{\lambda_n(A_k)} \leq \varepsilon$  then (1)
        return  $\mathbf{x}_k$  as an absolute  $\varepsilon$ -approximation;
     $\mathbf{a} := \mathbf{x}_k - f(\mathbf{x}_k)$ ;
    if  $\|\mathbf{a}\| \leq (1 - \rho^2)\varepsilon/\rho$  then (2)
        return  $\mathbf{x}_k - \mathbf{a}/(1 - \rho^2)$  as an absolute  $\varepsilon$ -approximation;
    if  $\|\mathbf{a}\| \leq \varepsilon$  then (3)
        return  $\mathbf{x}_k$  as a residual  $\varepsilon$ -approximation;
     $\xi := \mathbf{a}^T \mathbf{a} / ((1 + \rho)\sqrt{\mathbf{a}^T A_k \mathbf{a}})$ ;
     $\alpha := n(1 - \xi)/(n + 1)$ ;
     $\beta := \sqrt{\frac{n^2}{n^2 - 1}(1 - \xi^2)}$ ;
     $\gamma := (n\xi + 1)/(n + 1)$ ;
     $\mathbf{z} := A_k \mathbf{a} / \sqrt{\mathbf{a}^T A_k \mathbf{a}}$ ;
     $\mathbf{x}_{k+1} := \mathbf{x}_k - \gamma \mathbf{z}$ ;
     $A_{k+1} := \beta^2(A_k - (1 - \alpha^2/\beta^2)\mathbf{z}\mathbf{z}^T)$ ;
     $k := k + 1$ ;
end
if  $k = k_{\max}$  then
    return failed to compute  $\varepsilon$ -approximation in  $k_{\max}$  iterations

```

Figure 1: The circumscribed ellipsoid algorithm (CEA).

## 5 Numerical Results

In this section we compare numerical results for several test functions with utilizing numerically stable implementations of the CEA and SIA algorithms. We exhibit the total CPU time (in seconds), number of iterations and the number of the stopping criterion that resulted in termination. We also exhibit (in the parentheses) the upper bounds (see Lemma 3.3, where  $\delta = \varepsilon(1 - \rho)$  for the absolute termination case) on the number of iterations of the CEA and SIA algorithms. The speedup factors that represent the ratios of CPU time when using the CEA algorithm instead of the SIA algorithm are also included. All tests are carried out with the Linux operating system on Intel Pentium IV 2.4 GHz based machine.

The user specified termination parameter  $\varepsilon$  is modified in our code according to the follow-

ing rules. If the function evaluations are carried out in single precision, then

$$\varepsilon = \max(\varepsilon, \text{Macheps}(1)), \quad (10)$$

and if in double precision, then

$$\varepsilon = \max(\varepsilon, \text{Macheps}(2)), \quad (11)$$

where  $\text{Macheps}(j)$ ,  $j = 1, 2$ , are respectively machine precision in single and double floating point representations.

If the absolute termination is to be used (case  $\rho < 1$ ) the user may request to modify the  $\varepsilon$  by

$$\varepsilon = \max\left(\varepsilon, \frac{\text{Macheps}(j)}{1 - \rho}\right). \quad (12)$$

This is justified by the absolute sensitivity of the problem that is characterized by the condition number equal to  $\frac{1}{1-\rho}$ .

In the tests of functions 2 and 3, some initial balls are not the unit balls. In these cases, the problems are defined on a general ball  $B^n(\mathbf{c}, \gamma)$ . They can be transformed to the unit ball  $B^n(0, 1)$  as follows.

Let  $f : B^n(\mathbf{c}, \gamma) \rightarrow B^n(\mathbf{c}, \gamma)$  denote the original function defined on a ball  $B^n(\mathbf{c}, \gamma)$ . The modified function  $\bar{f} : B^n(0, 1) \rightarrow B^n(0, 1)$  defined on a unit ball  $B^n(0, 1)$  is

$$\bar{f}(\mathbf{x}) = \frac{1}{\gamma} (f(\mathbf{x} \cdot \gamma + \mathbf{c}) - \mathbf{c}).$$

It turns out that if  $f$  is  $\rho$ -Lipschitz on the ball  $B^n(\mathbf{c}, \gamma)$ , then  $\bar{f}$  is  $\rho$ -Lipschitz on the ball  $B^n(0, 1)$ . Indeed, if  $\mathbf{x}_1, \mathbf{x}_2 \in B^n(0, 1)$ , then

$$\begin{aligned} \|\bar{f}(\mathbf{x}_1) - \bar{f}(\mathbf{x}_2)\| &= \frac{1}{\gamma} \|f(\mathbf{x}_1 \cdot \gamma + \mathbf{c}) - f(\mathbf{x}_2 \cdot \gamma + \mathbf{c})\| \\ &\leq \frac{1}{\gamma} \cdot \rho \|\gamma \cdot (\mathbf{x}_1 - \mathbf{x}_2)\| \\ &= \rho \|\mathbf{x}_1 - \mathbf{x}_2\|, \end{aligned}$$

i.e., they satisfy Lipschitz condition with the same  $\rho$  and

$$\|\bar{f}(\mathbf{x})\| = \frac{1}{\gamma} \|f(\mathbf{x} \cdot \gamma + \mathbf{c}) - \mathbf{c}\| \leq \frac{1}{\gamma} \cdot \gamma = 1,$$

i.e.,  $\bar{f} : B^n(0, 1) \rightarrow B^n(0, 1)$ .

In all the tables we exhibit the user specified  $\varepsilon$ . In all cases these  $\varepsilon$ 's were then modified according to equations 10-12, however in only a few cases their values were changed.

**Test 1.** This test function is a simple affine mapping  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  given by

$$T_1 : \quad f(\mathbf{x}) = \rho \mathbf{x} + (1 - \rho) \mathbf{s},$$

where the constant vector  $\mathbf{s}$  is chosen from  $B^n(0, 1)$ . Obviously,  $\mathbf{s} = \mathbf{x}^*$  is the unique fixed point of  $f$  for  $\rho < 1$ . The results for  $n = 5$  are exhibited in Table 1. Table 2 shows the results with varying  $n$  from 2 to 5.

Table 1: T1:  $n = 5, \varepsilon = 10^{-6}$  and  $\mathbf{x}^* = [0.1, 0.3, 0.4, 0.1, 0.2]^T$ .

$\rho$	CEA	SIA	Speedup
$1 - 10^{-1}$	0.00276, 287 (1008), 2	0.000017, 119 (132), 2	0.00616
$1 - 10^{-2}$	0.00282, 302 (1146), 2	0.000173, 1247 (1375), 2	0.06135
$1 - 10^{-3}$	0.00283, 311 (1284), 2	0.001746, 12531 (13809), 2	0.61696
$1 - 10^{-4}$	0.00300, 321 (1423), 2	0.017640, 125361 (138149), 2	5.88000
$1 - 10^{-5}$	0.00252, 296 (1561), 2	0.175800, 1253671 (1381545), 2	69.7619
$1 - 10^{-6}$	0.00230, 304 (1699), 2	1.758000, 12536780 (13815504), 2	764.348

Table 2: T1:  $\rho = 1 - 10^{-6}, \varepsilon = 10^{-6}, \mathbf{x}_1^* = [0.1, 0.3]^T, \mathbf{x}_2^* = [0.1, 0.3, 0.4]^T, \mathbf{x}_3^* = [0.1, 0.3, 0.4, 0.1]^T$ , and  $\mathbf{x}_4^* = [0.1, 0.3, 0.4, 0.1, 0.2]^T$ .

$n$	$\mathbf{x}^*$	CEA	SIA	Speedup
2	$\mathbf{x}_1^*$	0.000132, 86 (339), 2	0.987, 11971078 (13815504), 2	7477.27
3	$\mathbf{x}_2^*$	0.000823, 185 (679), 2	1.367, 12448820 (13815504), 2	1661.00
4	$\mathbf{x}_3^*$	0.000826, 187 (1132), 2	1.462, 12467678 (13815504), 2	1769.98
5	$\mathbf{x}_4^*$	0.002300, 304 (1699), 2	1.758, 12536780 (13815504), 2	764.35

**Test 2.** This test function is a complex function from Howland and Vaillancourt (1985), given by

$$T_2 : \quad f(\mathbf{z}) = g(g(z)),$$

where

$$g(z) = \frac{z^2 + c \cos^2 z}{z + \sin z \cos z},$$

$z$  is a complex variable and  $c$  is a complex constant. We consider this test function as a two-dimensional real function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , i.e.,  $n = 2$ , but we evaluate the function as a one-dimensional complex function. The problem is tested with two values of the constant  $c$  :  $c = 1.025$  and  $c = \pi/4 + 1.2 + i(\pi - 1.17)$ . The fixed points of this problem are  $(0, 0.69029)^T$  and  $(2.14062, -2.50683)^T$ , respectively. The results are exhibited in Table 3.

Table 3: T2:  $B_1 = B^2([0.0, 0.1]^T, 1)$ ,  $B_2 = B^2([2.2, -2.2]^T, 1)$ ,  $c_1 = 1.025$ ,  $\rho_1 = 0.9989885$ ,  $c_2 = \pi/4 + 1.0 + i(\pi - 1.17)$ , and  $\rho_2 = 0.9984$ .

$\varepsilon$	$\rho$	$c$	Ball	CEA	SIA	Speedup
$10^{-2}$	$\rho_1$	$c_1$	$B_1$	0.000017, 3 (146), 2	0.00206, 583 (4551), 2	121.18
$10^{-3}$	$\rho_1$	$c_1$	$B_1$	0.000034, 7 (173), 2	0.00823, 2290 (6826), 2	242.06
$10^{-4}$	$\rho_1$	$c_1$	$B_1$	0.000063, 14 (201), 2	0.01452, 4102 (9102), 2	230.48
$10^{-5}$ *	$\rho_1$	$c_1$	$B_1$	0.000060, 14 (207), 2	0.01440, 4250 (9624), 2	240.00
$10^{-6}$ *	$\rho_1$	$c_1$	$B_1$	0.000060, 14 (207), 2	0.01450, 4250 (9624), 2	241.67
$10^{-2}$	$\rho_2$	$c_2$	$B_2$	0.000065, 14 (140), 1	0.00946, 2876 (2876), 1	145.54
$10^{-3}$	$\rho_2$	$c_2$	$B_2$	0.000092, 20 (168), 1	0.01339, 4314 (4314), 1	145.54
$10^{-4}$	$\rho_2$	$c_2$	$B_2$	0.000113, 25 (196), 1	0.01804, 5752 (5752), 1	159.65
$10^{-5}$ **	$\rho_2$	$c_2$	$B_2$	0.000119, 28 (207), 1	0.01916, 6369 (6369), 1	161.01
$10^{-6}$ **	$\rho_2$	$c_2$	$B_2$	0.000118, 28 (207), 1	0.01916, 6369 (6369), 1	162.37

\* these input  $\varepsilon$ 's were changed to  $5.89 \times 10^{-5}$  via equation 12.

\*\* these input  $\varepsilon$ 's were changed to  $3.72 \times 10^{-5}$  via equation 12.

**Test 3.** This test function is a parabolic, periodical function from Sikorski et al. (1993), given by

$$T_3 : \quad f(x_1, x_2) = [f_1(x_1, x_2), f_2(x_1, x_2)],$$

where

$$f_i(x_1, x_2) = \frac{\rho}{2}(x_i - 2m)^2 + 1 - \frac{\rho}{2}$$

for  $i = 1, 2$  and  $2m - 1 \leq x_i \leq 2m + 1$  and  $m$  is an arbitrary integer. This two-dimensional contractive function has a unique fixed point at  $(1, 1)^T$ . The results are exhibited in Table 4. It is worth noting that the CEA algorithm can terminate with absolute error criterion even when  $\rho = 1 - 10^{-15}$ ; i.e., when the function is almost nonexpanding.

**Test 4.** This test function is a saw-like, periodical function from Sikorski et al. (1993), given by

$$T_4 : \quad f(x_1, x_2) = \begin{pmatrix} \sqrt{3}/2 & -1/2 \\ 1/2 & \sqrt{3}/2 \end{pmatrix} \begin{pmatrix} f_1(x_1, x_2) \\ f_2(x_1, x_2) \end{pmatrix}$$

Table 4: T3:  $B_1 = B^2([0, 0]^T, 2)$ ,  $B_2 = B^2([0.1, 0.2]^T, 2)$ ,  $\rho_1 = 1 - 10^{-3}$ ,  $\rho_2 = 1 - 10^{-5}$ , and  $\rho_3 = 1 - 10^{-15}$ .

$\varepsilon$	$\rho$	Ball	CEA	SIA	Speedup
$10^{-3}$	$\rho_1$	$B_1$	0.000062, 34 (182), 2	0.000101, 1120 (7598), 2	8.42
$10^{-3}$	$\rho_2$	$B_1$	0.000065, 45 (237), 2	0.00106, 11874 (760087), 2	70.67
$10^{-4}$	$\rho_1$	$B_2$	0.000068, 47 (210), 2	0.00029, 2768 (9899), 2	4.40
$10^{-4}$	$\rho_2$	$B_2$	0.000077, 54 (265), 2	0.00336, 37376 (990344), 2	40.98
$10^{-6}$	$\rho_2$	$B_2$	0.000114, 79 (320), 2	0.0249, 277744 (1450859), 2	228.44
$10^{-6*}$	$\rho_3$	$B_2$	0.000154, 87 (596), 1	75.35, $1.8 \times 10^8$ ( $1.6 \times 10^{16}$ ), 2	489286

\* In these tests we removed the modification of  $\varepsilon$  according to equation 12, to test the limits of numerical precision of our algorithm.

where

$$f_i(x_1, x_2) = \min_{j=1,99} (\rho|x_i - m - 10^{-2}j| + i/3)$$

for  $i = 1, 2$  and  $m \leq x_i \leq m + 1$  and  $m$  is an arbitrary integer. This function has a unique fixed point at  $(-0.01946, 0.75933)^T$ . The results are exhibited in Table 5. It is worth noting that the speedup factor is increased by 4 orders of magnitude when  $\rho$  is closer to 1 by 4 orders of magnitude.

Table 5: T4:  $\varepsilon = 10^{-6}$ ,  $x^* = [-0.04, 0.74]^T$ ,  $B_1 = B^2([0, 0]^T, 1)$ ,  $B_2 = B^2([0, 0]^T, 2)$ , and  $B_3 = B^2([0.1, 0.2]^T, 2)$ .

$\rho$	Ball	CEA	SIA	Speedup
$1 - 10^{-2}$	$B_1$	0.000152, 36 (229), 1	0.00302, 1109 (1375), 2	19.87
$1 - 10^{-2}$	$B_2$	0.000169, 40 (237), 1	0.00306, 1109 (1444), 2	18.11
$1 - 10^{-2}$	$B_3$	0.000172, 41 (237), 1	0.00307, 1141 (1444), 2	17.85
$1 - 10^{-6}$	$B_1$	0.000152, 36 (339), 1	37.49, 13815504 (13815504), 1	246644
$1 - 10^{-6}$	$B_2$	0.000166, 41 (348), 1	40.40, 14508651 (14508651), 1	243373
$1 - 10^{-6}$	$B_3$	0.000168, 41 (348), 1	40.38, 14508651 (14508651), 1	240357

**Test 5.** This test function is a nine-dimensional function from Xu and Shi (1992), given by

$$T_5 : \quad f(\mathbf{x}) = x - tF(\mathbf{x})$$

where

$$F(\mathbf{x}) = A\mathbf{x} + G(\mathbf{x}) - \mathbf{b},$$

$$A = \begin{bmatrix} B & -I & 0 \\ -I & B & -I \\ 0 & -I & B \end{bmatrix},$$

$$B = \begin{bmatrix} 4 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 4 \end{bmatrix},$$

$$G(x) = \frac{1}{8} \left[ g\left(\frac{1}{4}, \frac{1}{4}, x_1\right), g\left(\frac{2}{4}, \frac{2}{4}, x_2\right), \dots, g\left(\frac{3}{4}, \frac{3}{4}, x_9\right) \right]^T,$$

and

$$\mathbf{b} = \left[ \phi\left(0, \frac{1}{4}\right) + \phi\left(\frac{1}{4}, 0\right), \phi\left(\frac{2}{4}, 0\right), \phi\left(\frac{3}{4}, 0\right) + \phi\left(1, \frac{1}{4}\right), \phi\left(0, \frac{2}{4}\right), 0, \phi\left(1, \frac{2}{4}\right), \right. \\ \left. \phi\left(0, \frac{3}{4}\right) + \phi\left(\frac{1}{4}, 1\right), \phi\left(\frac{2}{4}, 1\right), \phi\left(\frac{3}{4}, 1\right) + \phi\left(1, \frac{3}{4}\right) \right]^T,$$

with

$$g(s, t, u) = \frac{t}{s} + |u - 1|^3, \quad \phi(s, t) = \sin(0.5\pi st).$$

We choose the value of  $t = \sqrt{L^2(1 - \rho^2)}/L^2$  with  $L = 0.8234$  and select the Lipschitz constant  $\rho$  close to 1. This choice implies that  $f$  is contractive Xu and Shi (1992); Tsay (1994). The results are exhibited in Table 6. It turns out that the SIA algorithm is faster than the CEA algorithm when  $\rho$  is not too close to 1. This is because the cost of each iteration of the CEA algorithm increases with the dimension as  $O(n^3)$ . However, when  $\rho$  is close to 1, the CEA algorithm is much faster than the SIA algorithm.

We note that in all our tests the CEA algorithm terminated with the absolute error criterion (1) or (2) much faster than indicated by the upper bounds. More numerical tests and comparisons with Newton type methods will be carried out in a forthcoming paper Sikorski et al. (2006).

## 6 Combustion Chemistry Fixed Point Problem

In this section we focus on the design of nearly optimal fixed point solvers that are applied to univariate fixed point problems originating in modeling combustion of energetic materials. In particular, our solvers are able to efficiently approximate fixed points of nonlinear equations modeling burn rates of HMX type materials. These fixed point calculations are utilized in large scale transient combustion simulations. They are repeated at every grid cell of a very large model and at every time step. This is why they have to be extremely fast and sufficiently accurate.

Table 6: T5:  $\varepsilon = 10^{-6}$  and  $\rho$  is varied from  $1 - 10^{-2}$  to  $1 - 10^{-15}$ .

$\rho$	CEA	SIA	Speedup
$1 - 10^{-2}$	0.0601, 1630 (3440), 2	0.000088, 93 (1375), 2	0.00146
$1 - 10^{-3}$	0.0765, 2138 (3854), 1	0.000310, 332 (13809), 2	0.00405
$1 - 10^{-4}$	0.0818, 2288 (4269), 1	0.001044, 1122 (138149), 2	0.01276
$1 - 10^{-5}$	0.0834, 2333 (4683), 1	0.00351, 3775 (1381545), 2	0.04209
$1 - 10^{-6}$	0.0847, 2364 (5098), 1	0.01166, 12527 (13815504), 2	0.13766
$1 - 10^{-7}$	0.0847, 2349 (5512), 1	0.03830, 41440 (138155099), 2	0.45218
$1 - 10^{-8}$	0.0846, 2344 (5927), 1	0.1277, 136790 ( $\approx 1.4 \times 10^9$ ), 2	1.50946
$1 - 10^{-9}$ *	0.0851, 2358 (6341), 1	0.418, 450622 ( $\approx 1.4 \times 10^{10}$ ), 2	4.91187
$1 - 10^{-10}$ *	0.0847, 2337 (6612), 1	1.374, 1478283 ( $\approx 1.4 \times 10^{11}$ ), 2	16.2220
$1 - 10^{-11}$ *	0.0843, 2359 (7170), 1	4.49, 4824754 ( $\approx 1.4 \times 10^{12}$ ), 2	53.2622
$1 - 10^{-12}$ *	0.0843, 2330 (7585), 1	13.93, 14999188 ( $\approx 1.4 \times 10^{13}$ ), 2	165.243
$1 - 10^{-13}$ *	0.0846, 2348 (7999), 1	42.26, 45615144 ( $\approx 1.4 \times 10^{14}$ ), 2	499.527
$1 - 10^{-14}$ *	0.0855, 2330 (8414), 1	169.1, 138613488 ( $\approx 1.4 \times 10^{15}$ ), 2	1977.19
$1 - 10^{-15}$ *	0.0853, 2348 (8828), 1	523.9, 420287875 ( $\approx 1.4 \times 10^{16}$ ), 2	6141.85

\* In these tests we removed the modification of  $\varepsilon$  according to equation 12, to test the limits of numerical precision of our algorithm.

We derive an almost optimal (on the average) hyper-bisection/secant (HBS) modification of a hybrid bisection-regula falsi-secant (BRS) method Novak et al. (1995); Sikorski (2001) to solve a nonlinear fixed point problem that is derived from the Ward, Son, Brewster (WSB) combustion model Ward et al. (1998a,b); Ward (1997). The WSB model assumes a simple overall kinetic scheme to represent the complex chemical reactions in condensed phase and gas phase: a zero-order mildly exothermic decomposition reaction with high activation energy in condensed phase is followed by a second-order highly exothermic reaction with zero activation energy in gas phase. Given the values of initial solid temperature  $T_0$  and gas pressure  $P$ , the WSB model defines the functions for the condensed phase mass flux  $m$  and the burning surface temperature  $T_s$ :  $m(T_s)$  and  $T_s(m)$ . The problems in condensed phase and gas phase are solved separately. Solving for the condensed phase problem, we obtain a solution for the condensed phase mass flux Ward et al. (1998a,b); Ward (1997):

$$m(T_s) = \sqrt{\frac{A_c R T_s^2 k_c \rho_c}{E_c (c_p (T_s - T_0) - \frac{Q_c}{2})}} \cdot e^{-\frac{E_c}{R T_s}} \quad (13)$$

where  $m$  is the condensed phase mass flux ( $\frac{kg}{m^2 s}$ ),  $T_s$  is the burning surface temperature,  $T_0$  is the initial solid temperature,  $c_p$  is the specific heat,  $k_c$  is the thermal conductivity,  $A_c$  is the pre-exponential factor,  $Q_c$  is the chemical heat release per unit mass,  $\rho_c$  is the solid

density,  $E_c$  is the solid phase apparent activation energy and  $R$  is the gas constant (Table 7). Solving for the gas phase problem, we obtain a solution for the burning surface temperature Ward et al. (1998a,b); Ward (1997):

$$T_s(m) = T_0 + \frac{Q_c}{c_p} + \frac{Q_g}{\frac{2mc_p}{\sqrt{m^2 + \frac{\Delta k_g B_g P^2 W^2}{c_p R^2}} - m} + c_p} \quad (14)$$

Table 7: List of constant values

$c_p = 1.4 \times 10^3 J \cdot kg^{-1} \cdot K^{-1}$	$R = 8.314 J \cdot K^{-1} \cdot mol^{-1}$
$k_c = 0.2 W \cdot m^{-1} \cdot K^{-1}$	$k_g = 0.07 W \cdot m^{-1} \cdot K^{-1}$
$A_c = 1.637 \times 10^{15} s^{-1}$	$B_g = 1.6 \times 10^{-3} m^3 \cdot kg^{-1} \cdot s^{-1}$
$Q_c = 4.0 \times 10^5 J \cdot kg^{-1}$	$Q_g = 3.018 \times 10^6 J \cdot kg^{-1}$
$\rho_c = 1.8 \times 10^3 kg \cdot m^{-3}$	$W = 3.42 \times 10^{-2} kg \cdot mol^{-1}$
$E_c = 1.76 \times 10^5 J \cdot mol^{-1}$	

where, subscript  $g$  denotes variables for the gas phase,  $W$  is the overall molecular weight of gas product,  $P$  is the pressure and  $B_g$  is the gas phase pre-exponential factor (Table 7).

With the known equations for  $T_s(m)$  and  $m(T_s)$  outlined above, we want to compute the burn rate  $m$  and the burning surface temperature  $T_s$ . The  $T_s$  and  $m$  are dependent on the local pressure  $P$  and the initial solid temperature  $T_0$ . The nonlinear dependence between  $m(T_s)$  and  $T_s(m)$  implies the need for an iterative solution.

Since  $A_c$ ,  $k_c$ ,  $\rho_c$ ,  $E_c$ ,  $Q_c$ ,  $c_p$ ,  $k_g$ ,  $B_g$ ,  $Q_g$ ,  $W$  and  $R$  are all constant physical properties of materials, then for a given values of initial solid temperature  $T_0$  and gas phase pressure  $P$ , we can simplify the equations for  $m(T_s)$  and  $T_s(m)$  to:

$$m(T_s) = \sqrt{\frac{C_1 T_s^2}{T_s - C_2}} \cdot e^{-\frac{E_c}{RT_s}} \quad (15)$$

and



$$T_s(m) = C_4 + \frac{C_5}{(\sqrt{m^2 + C_3} + m)^2} \quad (16)$$

where,  $C_1 = \frac{A_c R K_c \rho_c}{E_c c_p}$ ,  $C_2 = T_0 + \frac{Q_c}{2c_p}$ ,  $C_3 = \frac{4k_g B_g P^2 W^2}{c_p R^2}$ ,  $C_4 = T_0 + \frac{Q_c}{c_p}$  and  $C_5 = \frac{C_3 Q_g}{c_p}$ . When plugging  $m(T_s)$  into  $T_s(m)$ , we get the fixed point equation  $T_s = G(T_s)$ , where

$$G(T_s) = C_4 + \frac{C_5}{\left( \sqrt{\frac{C_1 T_s^2}{T_s - C_2} \cdot e^{-\frac{E_c}{RT_s}} + C_3} + \sqrt{\frac{C_1 T_s^2}{T_s - C_2} \cdot e^{-\frac{E_c}{RT_s}}} \right)^2}$$

Now, the problem can be formulated as a zero finding problem: given gas phase pressure  $P$ , initial solid temperature  $T_0$  and a small positive number  $\varepsilon$ , we want to find an  $\varepsilon$ -approximation  $T_\varepsilon$  of the exact zero  $T_{s\_true}$  with respect to the root criterion  $\frac{|T_\varepsilon - T_{s\_true}|}{|T_{max} - T_{min}|} \leq \varepsilon$  by solving the nonlinear equation  $f(T_s) = 0$ , where  $[T_{min}, T_{max}]$  is the interval containing the solution  $T_{s\_true}$  and  $f(T_s) = T_s - G(T_s)$ . Figure 2 depicts a number of functions with various  $T_0$  and  $P$ .

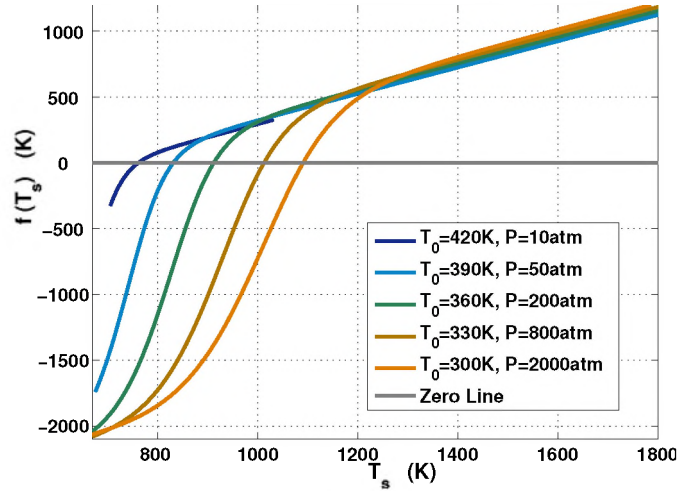


Figure 2: Functions  $f(T_s)$  with various  $T_0$  and  $P$

## 7 Description of the Algorithm

For the univariate zero finding problem, the solution can be efficiently obtained by using a hybrid bisection-secant method. The reasons for adopting a hybrid method are:

1. Bisection Method was proved to be optimal in the worst case in the class of all methods using arbitrary sequential linear information Sikorski (2001, 1982).
2. A hybrid bisection-secant type method is more efficient than plain bisection in the average case and exhibits almost optimal average case complexity Novak et al. (1995); Sikorski (2001).

## 7.1 Solution Interval

Given the function  $f(T_s)$ , we can identify a solution interval  $[T_{min}, T_{max}]$  such that  $f(T_s)$  has different signs at its endpoints. Since  $f(T_s)$  is continuous, it follows that the exact zero  $T_{s\_true}$  such that  $f(T_{s\_true}) = 0$  is in  $[T_{min}, T_{max}]$ . Since the constants  $C_3, C_4, C_5$  and variable  $m$  are all positive, then  $T_s(m) > C_4$ , which implies  $G(T_s) > C_4$  and  $f(C_4) = C_4 - G(C_4) < 0$ . Thus, we can define the left endpoint of the interval as  $T_{min} = C_4$ . To estimate the maximum of the function  $G(T_s)$ , we need to solve for  $\frac{dG(T_s)}{dT_s} = 0$ , which gives the local maximum of  $G(T_s)$  at  $T_{s\_max} = C_2 - \frac{E_c}{2R} + \sqrt{C_2^2 + \frac{E_c^2}{4R^2}} > C_2$  for  $T_s > 0$ , where  $C_2 > 0$  (see Appendix). Since the function has negative derivative for  $T_s > T_{s\_max}$  (see Appendix), we define  $T_{max} = G(T_{s\_max})$  if  $T_{min} < T_{s\_max}$ , or  $T_{max} = G(T_{min})$  otherwise. In other words,  $T_{max} = G(\text{Max}(T_{min}, T_{s\_max}))$ . This choice of  $T_{max}$  implies that  $G(T_{max}) < T_{max}$  and therefore  $f(T_{max}) > 0$ . In the majority of tests,  $T_{min}$  was greater than  $T_{s\_max}$ . An example is shown in Figures 3 and 4.

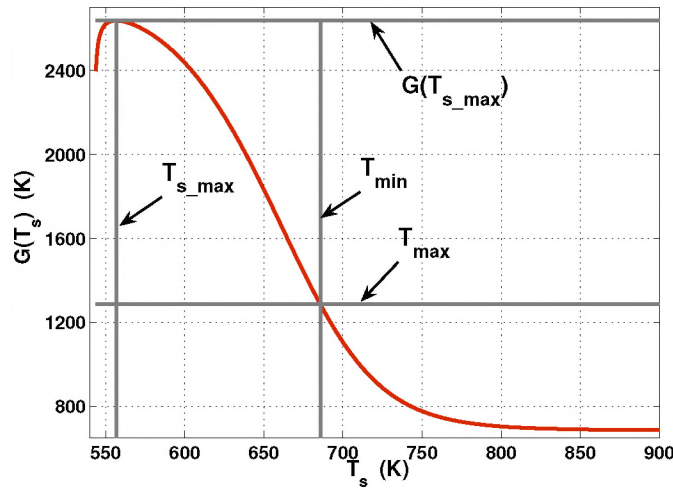


Figure 3: Function  $G(T_s)$  with  $T_0 = 400K$  and  $P = 10atm$  ( $T_{min} > T_{s\_max}$ )

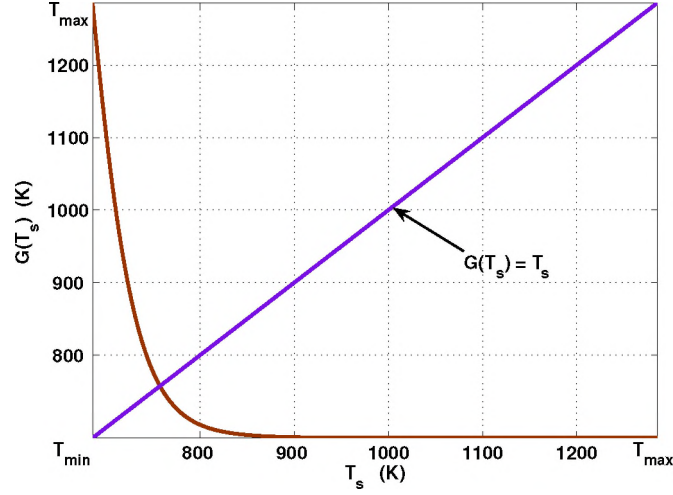


Figure 4: Zoom-in of Figure 3 with  $G(T_s)$ ,  $T_s \in [T_{min}, T_{max}]$

## 7.2 The BRS Method

We now outline the hybrid bisection-regula falsi-secant (BRS) method investigated in Novak et al. (1995). This method computes points  $T_{s,i} \in [T_{min}, T_{max}]$ , at which the function  $f(T_s)$  is evaluated, and a subinterval  $[l_i, r_i]$  containing a zero of  $f(T_s)$ . We always have  $f(l_i) \leq 0 \leq f(r_i)$ ,  $T_{s,i} \in [l_{i-1}, r_{i-1}]$  and  $[l_i, r_i] \subset [l_{i-1}, r_{i-1}]$ . Each function evaluation is preceded by the verification of the stopping rule. At the beginning and later after each bisection step, the method takes two steps of the regula falsi method, starting from the endpoints of the interval  $[l_{i-1}, r_{i-1}]$ . Then, the secant steps are performed until they are well defined, result in points in the current interval and reduce the length of the current interval by at least one half in every three steps. If any of these conditions is violated, a bisection step takes place. We utilize the absolute error termination criterion given by:

$$Stop_i = \begin{cases} 1, & \text{if } (r_i - l_i)/2 \leq \varepsilon \cdot (T_{max} - T_{min}), \text{ then } T_{s,i} = (l_i + r_i)/2 \\ 0, & \text{otherwise} \end{cases}$$

After the termination ( $Stop_i = 1$ ), we have  $\frac{|T_{s,i} - T_{s,true}|}{|T_{max} - T_{min}|} \leq \varepsilon$ , since the exact solution  $T_{s,true} \in [l_i, r_i]$ . We further define the points

$$Secant(u, w) = \begin{cases} u - \frac{f(u) \cdot (u-w)}{f(u) - f(w)}, & \text{if } f(u) \neq f(w) \\ \text{undefined}, & \text{otherwise} \end{cases}$$

for the secant method with  $u, w \in [T_{min}, T_{max}]$ , and

$$Bisection_i = \frac{l_i + r_i}{2}$$

for the bisection method. In each step, a new interval  $I_i$  containing a zero of  $f(T_s)$  is computed,

$$I_i = \begin{cases} [l_{i-1}, x_i], & \text{if } f(x_i) > 0 \\ [x_i, r_{i-1}], & \text{if } f(x_i) < 0 \\ [x_i, x_i], & \text{if } f(x_i) = 0 \end{cases}$$

The complete method is summarized in the flowchart in Figure 5.

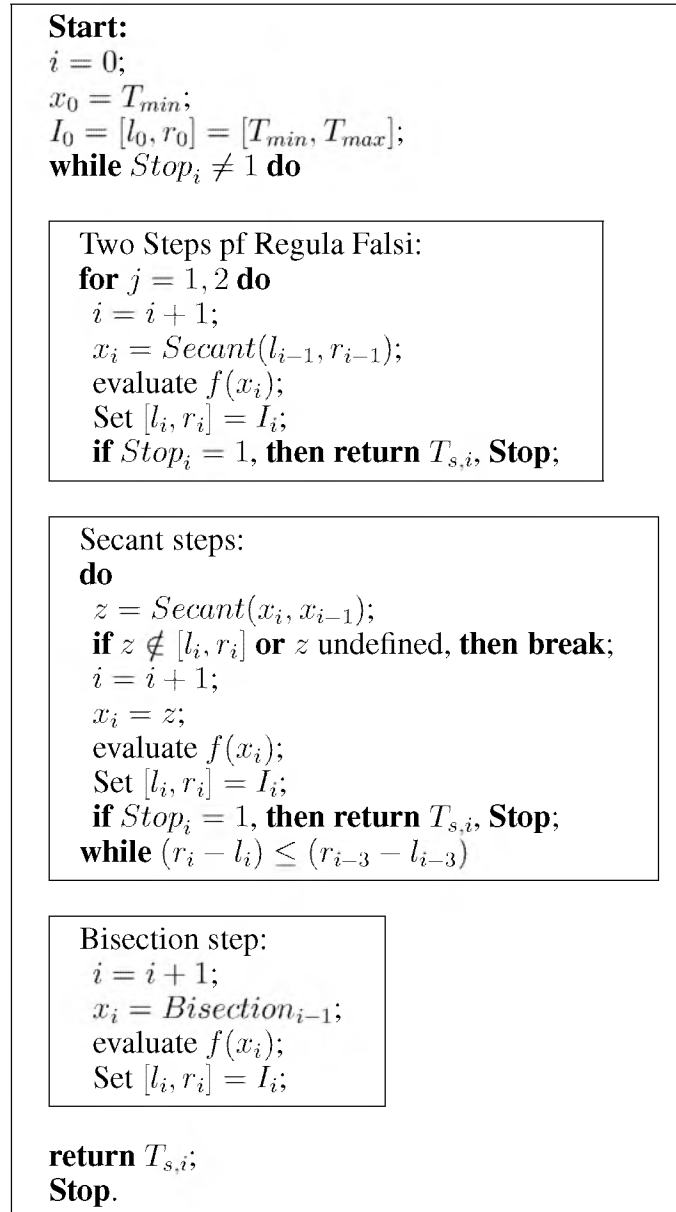


Figure 5: Flowchart of the BRS method

The BRS method is almost optimal on the average Novak et al. (1995); Sikorski (2001). The average number  $m_{aver}$  of function evaluations for finding an  $\varepsilon$ -approximation to the

solution is bounded by:

$$m_{aver} \leq \frac{1}{\log\left(\frac{1+\sqrt{5}}{2}\right)} \cdot \log \log\left(\frac{1}{\varepsilon}\right) + A$$

where  $A$  is a constant Novak et al. (1995).

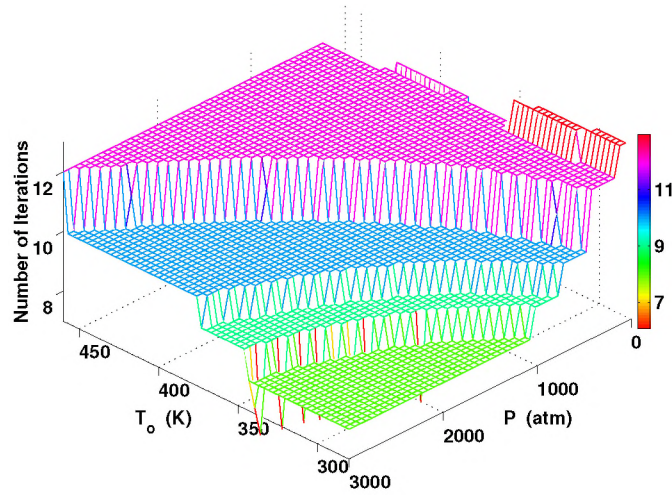


Figure 6: Iteration Graph of the BRS method

For practical combustion simulations, the initial solid bulk temperature  $T_0$  usually varies within the interval  $TE = [280, 460] K$  and the gas phase pressure  $P$  within the interval  $PR = [0, 3000] atm$ . The number of iterations of the BRS method can be visualized in 3-D as a function of  $T_0$  and  $P$  within those intervals. To carry out the tests, we selected  $60 \times 50$  evenly spaced grid nodes for the set of parameters  $TE \times PR$ . By choosing  $\varepsilon = 10^{-4}$ , the average number of iterations is 10.5, where the average is defined as the total number of iterations divided by the number of tested functions. We observed (Figure 6) that for low  $P$  and high  $T_0$  it took about 12-13 iterations to solve the problem. The differences in the numbers of iterations can be explained by the special graph shapes of the functions in different sub-domains of the parameter set  $TE \times PR$ . Figure 2 indicates that the functions have characteristic break-points, where the largest function curvatures occur. The solution of  $f(T_s) = 0$  turns out to be very close to the high curvature break-point for high  $T_0$  and low  $P$ . In this case, more iterations are required for the secant method to converge. We derive a modification to the BRS method in order to lower the average number of iterations.

### 7.3 HBS - a modified BRS method

To derive the HBS method, we first divide the parameter plane  $TE \times PR$  into three subdomains  $D_i, i = 1, 2, 3$  (Figure 7) by two lines  $P = 4 \cdot (T_o - 250)$  and  $P = 15 \cdot (T_o - 250)$ ,

$$D_i = \begin{cases} D_1, & \text{if } P \leq 4 \cdot (T_o - 250) \\ D_2, & \text{if } 4 \cdot (T_o - 250) < P \leq 15 \cdot (T_o - 250) \\ D_3, & \text{if } P > 15 \cdot (T_o - 250) \end{cases}$$

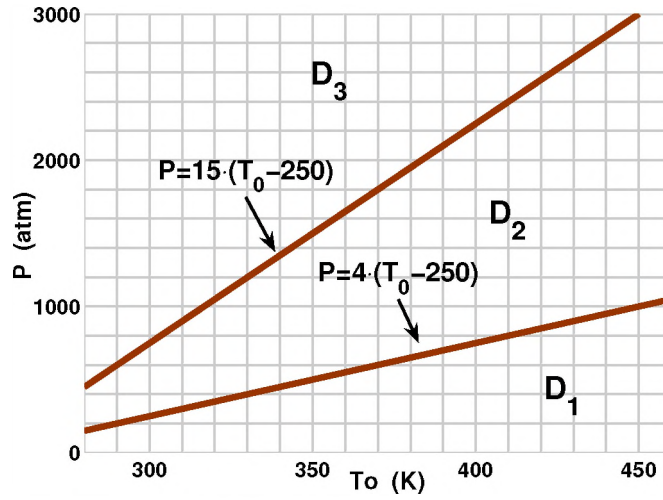


Figure 7:  $TE \times PR$  domain subdivision

For each subdomain, we run two steps of hyper-bisection method defined as:

$$Hyperbis_i = l_i + \xi_i \cdot (r_i - l_i)$$

where  $\xi_i = \frac{T_{s,i} - l_i}{r_i - l_i} \in [0, 1]$ . Numerical experiments indicate that the solutions are usually distributed around the point  $T_{min} + \lambda \cdot (T_{max} - T_{min})$  in the sub-domain  $D_1$ , where  $\lambda = 0.12$ . We therefore utilize  $\xi_1 = \lambda$  for the first step of hyper-bisection and,

$$\xi_2 = \begin{cases} \delta^2, & \text{if function evaluation } f(Hyperbis_1) < 0, \\ 1 - \delta, & \text{otherwise,} \end{cases}$$

where  $\delta = 0.2$ . Those choices guarantee that in most cases the solution will be in the interval  $[Hyperbis_1, Hyperbis_2]$ . The same strategy applies to subdomains  $D_2$  and  $D_3$ , with  $\lambda = 0.18$  for  $D_2$  and  $\lambda = 0.25$  for  $D_3$ . Parameter  $\delta$  equals 0.2 for all the subdomains. Ideally, the solution interval will be reduced to 2% – 5% of its original length after two

steps of the hyper-bisection. Thereafter, the BRS method is used to find the solution. When choosing the same set of test functions and  $\varepsilon = 10^{-4}$ , the average number of iterations of the HBS method is 5.7 as compared to 10.5 of the BRS method. The maximum number of iterations is 6. The iteration graph of the HBS method is shown in Figure 8. Since we ultimately want to compute the burn rate  $m$ , the equation (15) is used to compute  $m = m(T_s)$ . The maximum of the relative error  $e_{rel}$  in the computation of  $m$  is  $1.2 \times 10^{-3}$  in the  $TE \times PR$  domain, where  $e_{rel} = \frac{|m_{calc} - m_{real}|}{|m_{real}|}$  and the real solution  $m_{real}$  is computed by our HBS method with  $\varepsilon = 10^{-15}$ .

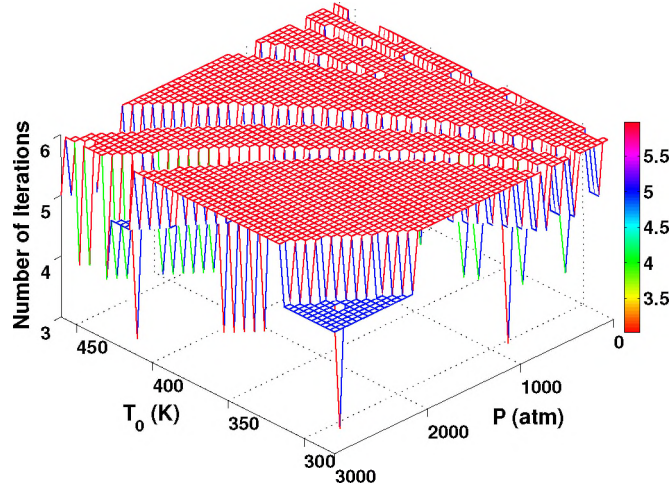


Figure 8: Iteration Graph of the HBS method

## 8 Appendix

We derive formulas for the first derivative of the function  $G(\cdot)$  and its roots that are needed in verification of the definition of initial interval locating the fixed point, as well as in showing existence of exactly one fixed point of  $G$ .

1. First derivative of the function  $G(T_s)$  is given by:

$$\frac{dG(T_s)}{dT_s} = \frac{[2RT_s^2 - (2RC_2 + R - E_c)T_s - E_c C_2] \cdot \sqrt{C_1} C_5 \cdot e^{-\frac{E_c}{2RT_s}}}{T_s R \sqrt{C_1 T_s^2 \cdot e^{-\frac{E_c}{2RT_s}} + C_3(T_s - C_2)} \cdot \left[ \sqrt{C_1 T_s^2 \cdot e^{-\frac{E_c}{2RT_s}} + C_3(T_s - C_2)} + T_s \sqrt{C_1} \cdot e^{-\frac{E_c}{2RT_s}} \right]^2}$$



2. Roots of the first derivative  $\frac{dG(T_s)}{dT_s} = 0$  are given by:

$$T_s = C_2 - \frac{E_c}{2R} + \sqrt{C_2^2 + \frac{E_c^2}{4R^2}} > 0 \text{ and } T_s = C_2 - \frac{E_c}{2R} - \sqrt{C_2^2 + \frac{E_c^2}{4R^2}} < 0$$

## References

- AGEEV, A., VASSIN, V., BESSONOVA, E., AND MARKUSEVICH, V. 1997. Radiosounding ionosphere on two frequencies. algorithmic analysis of fredholm-stiltjes integral equation, in: *Theoretical problems in geophysics (in russian)*. 100–118.
- ALLGOWER, E. L. 1977. Application of a fixed point search algorithm to nonlinear boundary value problems having several solutions. In *Fixed Points Algorithms and Applications*, S. Karamardian, Ed. Academic Press, New York, 87–111.
- ALLGOWER, E. L. AND GEORG, K. 1990. *Numerical Continuation Methods*. Springer-Verlag, New York.
- ANDREAS, J. 2003. *Topological Fixed Point Principles for Boundary Value Problems*. Kluwer Academic Publishers.
- BANACH, S. 1922. Sur les opérations dans les ensembles abstraits et leur application aux équations integrales. *Fund. Math.* 3, 133–181.
- BORDER, K. C. 1985. *Fixed point theorems with applications to economics and game theory*. Cambridge University Press, New York.
- BROUWER, L. E. 1912. Über abbildungen von mannigfaltigkeiten. *Mathematische Annalen* 71, 97–115.
- CORLESS, R. M., FRANK, G. W., AND MONROE, J. G. 1990. Chaos and continued fractions. *Physica D* 46, 241–253.
- DENG, X. AND CHEN, X. 2005. On algorithms for discrete and approximate brouwer fixed points. In *Proceedings of the 37th Annual ACM Symposium on Theory of Computing, Baltimore, MD, USA, May 22-24, 2005*, H. N. Gabov and R. Fagin, Eds. ACM.
- EAVES, B. C. 1972. Homotopies for computation of fixed points. *Math. Programming* 3, 1, 1–12.
- EAVES, B. C. AND SAIGAL, R. 1972. Homotopies for computation of fixed points on unbounded regions. *Math. Programming* 3, 2, 225–237.

- GARCIA, C. B., LEMKE, C. E., AND LUETHI, H. 1973. Simplicial approximation of an equilibrium point for non-cooperative n-person games. In *Mathematical Programming*, T. C. Hu and S. M. Robinson, Eds. Academic Press, New York, 227–260.
- HIRSCH, M. D., PAPADIMITRIOU, C., AND VAVASIS, S. 1989. Exponential lower bounds for finding Brouwer fixed points. *Journal of Complexity* 5, 379–416.
- HOWLAND, J. L. AND VAILLANCOURT, R. 1985. Attractive cycle in the iteration of meromorphic functions. *Numer. Math.* 46, 323–337.
- HUANG, Z., KHACHIYAN, L., AND SIKORSKI, K. 1999. Approximating fixed points of weakly contracting mappings. *J. Complexity* 15, 200–213.
- JEPPSON, M. M. 1972. A search for the fixed points of a continuous mapping. In *Mathematical Topics in Economic Theory and Computation*, R. H. Day and S. M. Robinsons, Eds. SIAM, Philadelphia, 122–129.
- KHACHIYAN, L. 1979. Polynomial algorithm in linear programming. *Soviet. Math. Dokl.* 20, 191–194.
- MERRIL, O. H. 1972a. Application and extensions of an algorithm that computes fixed points of a certain upper semi-continuous point of set mapping. Ph.D. thesis, University of Michigan, Ann Arbor, MI.
- MERRIL, O. H. 1972b. A summary of techniques for computing fixed points of continuous mapping. In *Mathematical Topics in Economic Theory and Computation*, R. H. Day and S. M. Robinson, Eds. SIAM, Philadelphia, PA, 130–149.
- NOVAK, E., RITTER, K., AND WOŹNIAKOWSKI, H. 1995. Average case optimality of a hybrid secant-bisection method. *Math. Computation* 64, 1517–1539.
- PERESTONINA, G., PRUTKIN, I., TIMERKHANOVA, L., AND VASSIN, V. 2003. Solving three-dimensional inverse problems of gravimetry and magnetometry for three layer medium (in russian). *Math. Modeling* 15, 2, 69–76.
- SCARF, H. 1967a. The approximation of fixed point of a continuous mapping. *SIAM J. Appl. Math.* 35, 1328–1343.
- SCARF, H. 1967b. The core of an n person game. *Econometrica* 35, 50–69.
- SCARF, H. E. AND HANSEN, T. 1973. *Computation of Economic Equilibria*. Yale University Press.
- SHELLMAN, S. AND SIKORSKI, K. 2002. A two-dimensional bisection envelope algorithm for fixed points. *J. Complexity* 18, 2 (June), 641–659.

- SHELLMAN, S. AND SIKORSKI, K. 2003a. Algorithm 825: A deep-cut bisection envelope algorithm for fixed points. *ACM Trans. Math. Soft.* 29, 3 (Sept.), 309–325.
- SHELLMAN, S. AND SIKORSKI, K. 2003b. A recursive algorithm for the infinity-norm fixed point problem. *J. Complexity* 19, 6 (Dec.), 799–834.
- SHELLMAN, S. AND SIKORSKI, K. 2005. Algorithm 848: A recursive fixed point algorithm for the infinity-norm case. *ACM Trans. Math. Soft.* 31, 4, 580–587.
- SIKORSKI, K. 1982. Bisection is optimal. *Numer. Math* 40, 111–117.
- SIKORSKI, K. 1989. Fast algorithms for the computation of fixed points. In *Robustness in Identification and Control*, R. T. M. Milanese and A. Vicino, Eds. Plenum Press, New York, 49–59.
- SIKORSKI, K. 2001. *Optimal Solution of Nonlinear Equations*. Oxford Press, New York.
- SIKORSKI, K., BOONIASIRIVAT, C., AND TSAY, C.-W. 2006. Algorithm: The circumscribed ellipsoid algorithm for fixed points. *to be submitted to ACM ToMS*.
- SIKORSKI, K., TSAY, C. W., AND WOŹNIAKOWSKI, H. 1993. An ellipsoid algorithm for the computation of fixed points. *J. Complexity* 9, 181–200.
- SIKORSKI, K. AND WOŹNIAKOWSKI, H. 1987. Complexity of fixed points. *J. Complexity* 3, 388–405.
- SZREDNICKI, R. 2000. A generalization of the lefschetz fixed point theorem and detection of chaos. *Proc. Amer. Math. Soc.* 128, 1231–1239.
- TSAY, C. W. 1994. Fixed point computation and parallel algorithms for solving wave equations. Ph.D. thesis, University of Utah, Salt Lake City, UT.
- VASSIN, V. 2005. Ill-posed problems with a priori information: methods and applications. Institute of Mathematics and Mechanics, Russian Academy of Sciences, Ural Subdivision.
- VASSIN, V. AND AGEEV, A. 1995. *Ill-posed problems with a priori information*. VSP, Utrecht, The Netherlands.
- VASSIN, V. AND EREMIN, E. 2005. *Feyer Type Operators and Iterative Processes (in Russian)*. Russian Academy of Sciences, Ural Subdivision, Ekaterinburg.
- VASSIN, V. AND SEREZNIKOVA, T. 2004. Two stage method for approximation of nonsmooth solutions and reconstruction of noisy images (in russian). *Automatica and Telemechanica* 2.

- WARD, M. 1997. A new modeling paradigm for the steady deflagration of homogeneous energetic materials. M.S. thesis, University of Illinois, Urbana-Champaign.
- WARD, M., SON, S., AND BREWSTER, M. 1998a. Role of gas- and condensed-phase kinetics in burning rate control of energetic solids. *Combust. Theory Modeling* 2, 293–312.
- WARD, M., SON, S., AND BREWSTER, M. 1998b. Steady deflagration of hmx with simple kinetics: A gas phase chain reaction model. *Combustion And Flame* 114, 556–568.
- XU, Z. B. AND SHI, X. Z. 1992. A comparison of point and ball iterations in the contractive mapping case. *Computing* 49, 75–85.
- YAMAMOTO, N. 1998. A numerical verification method for solutions of boundary value problems with local uniqueness by banach's fixed-point theorem. *SIAM J. Num. Anal.* 35, 2004–2013.
- YANG, Z. 1999. *Computing equilibria and fixed points: The solution of nonlinear inequalities*. Kluwer Academic Publishers, Dordrecht.