

```

else if  $\alpha_i \cdot x_{i+1} + \beta_i \cdot B^{m_{i+1}}$  overflows, then
     $x_i = \alpha_i \cdot x_{i+1} / B + \beta_i \cdot B^{m_{i+1}-1}$  (carefully!)
     $m_i = m_{i+1} - 1$ 
else if  $\alpha_i \cdot x_{i+1} + \beta_i \cdot B^{m_{i+1}}$  underflows, then
     $x_i = \alpha_i \cdot x_{i+1} \cdot B + \beta_i \cdot B^{m_{i+1}+1}$  (carefully!)
     $m_i = m_{i+1} + 1$ 
endif
endfor

```

The true value of x_i is forgotten from the computer, but it is the same way as gotten from \tilde{m}_i and m_i as described above: If \tilde{m}_i is less than 1 in magnitude, multiply it by B and add 1 to it, then subtract the floor of \tilde{m}_i from it so the largest is not zero, and then change \tilde{m}_i to $B^{-\tilde{m}_i}$.

I have not debugged this pseudocode, but I think it is basically correct. Alan Edelman points out that these can all be blocked in straight code.

6. Compute $y = \prod_{i=1}^n d_i$ for all i in one parallel operation.

To protect against over/underflow, one can modify this scheme in one of the following ways. The easiest way is to compute \log of the moment \log of $\sum_{j=1}^i \log g_j$ using an add-scan, and define $d_i = \exp(\log g_i - \bar{d}_{i-1})$. Let $\bar{d} = \max \log g_i$ and $e = \max \log g_i$. Replace d_i by $\exp(\log g_i - \bar{d})$ and replace \bar{d} by $\log g_i - e$. Exponentiate the new d_i and \bar{d} to get scaled values, do the largest of each of which is 1. Now perform steps (4) through (6) of the above algorithm.

This has the added cost of a logarithm and exponent, and loses a little of them too. But it will not overflow and almost certainly not underflow. (There are other algorithms that might be marginally safer against underflow signs of them be accumulated and applied with a multiply scan operation only ± 1 's.)

A better way to protect against overflow is to modify the multiply-scan as follows. Instead of computing $d_i = d_{i-1} \cdot f_i$ and integer m_i such that d_i cannot overflow or underflow, compute $\tilde{d}_i = d_{i-1} \cdot f_i / B^{m_i}$ where B is a big power of the radix near overflow threshold. Here is a code, which can obviously be "scanned" for m_i :

```

 $\tilde{d}_0 = 1; m_0 = 0$ 
for  $j = 1, n$ 
  if  $\tilde{d}_{i-1} \cdot f_i$  neither overflows nor underflows then
     $\tilde{d}_i = \tilde{d}_{i-1} \cdot f_i$ 
     $m_i = m_{i-1}$ 
  else  $\tilde{d}_{i-1} \cdot f_i$  would overflow then
     $\tilde{d}_i = \tilde{d}_{i-1} \cdot f_i / B$  (computed carefully!)
     $m_i = m_{i-1} + 1$ 
  else  $\tilde{d}_{i-1} \cdot f_i$  would underflow then
     $\tilde{d}_i = \tilde{d}_{i-1} \cdot f_i \cdot B$  (computed carefully!)
     $m_i = m_{i-1} - 1$ 
  endif
endifor

```

If B is close to overflow, at this point \tilde{d}_i loses its magnitude, multiply them by B and subtract 1 from the result. When the magnitude is small, the largest one can be subtracted from all of them (so the largest is $d_i = \tilde{d}_i B^{m_i}$ computed without fear of overflow).

The third approach is to run $x_i = \alpha x_{i+1} + \beta_i$ sequentially, scaling if necessary as one goes. The code is the same as the last code above:

```

 $\tilde{x}_n = \beta_n; m_n = 0$ 
for  $i = n - 1$  downto 1
  if  $\alpha x_{i+1} + \beta_i \cdot B^{m_{i+1}}$  neither over/underflows,
     $x_i = \alpha \cdot x_{i+1} + \beta_i \cdot B^{m_{i+1}}$ 
     $m_i = m_{i+1}$ 

```

3 Exploiting extra range

If we have extra exponent range available, we can greatly diminish spent testing and scaling. If the data is in IEEE single precision, (normalized!) numbers can be computed without over/underflow in IEEE and products of 128 in IEEE extended. If the data is IEEE double then p be computed in IEEE extended. Thus, any scaling tests would only have 8, 16 or 128 products.

4 Extensions

Of course, it would be nice if the user could supply his or her own data tive operator, and have the system perform parallel prefix. Given the handling as described above, it would be nice to have these basic floa be done automatically, rather than expecting the user to handle them

Appendix Comments on Inverse Iteration on the CM-2

1

The task at hand is to solve $Bx = y$ where B is an n by n upper bidiagonal x and y are n -vectors. Let d be the diagonal entries, of B , and b be the superdiagonal entries. The usual way to solve $Bx = y$ is to solve for $i = n$ down to $i = 1$, $x_{i+1} = 0$ and $x_i = (y_i - b x_{i+1}) / d_i$. I can think of three ways to solve this, with various kinds of immunity against overflow. I have not analyzed all these (they are not all equivalent to evaluating the recurrence nearly perfect backward error), but I don't see any obvious dangers in the scan operation to the extent possible, assuming only scans for floating-point multiply (although the best solution would involve no scan).

The first two, and most parallel, methods involve the following factorization $B = D_1 D_2$, where D_1 and D_2 are diagonal, and E is bidiagonal with 1 on the diagonal. $D_2 = \text{diag}(d_1, \dots, d_{n-1})$ is given by d and $d \prod_{j=1}^i (q_j / b_j)$. $D_1 = \text{diag}(e_1, \dots, e_{n-1})$ is given by $e = y / d$. One can also view E as upper triangular with all ones on and above the diagonal. Thus, computing $D_2 E^{-1} D_1 x = B^{-1} y$ involves the following:

1. Compute $e_i \neq 1$, $i = 1, \dots, n-1$ for all i in one parallel step.
2. Compute $q_i = d \prod_{j=1}^i (q_j / b_j)$ for all i using a multiply-scan operation.
3. Compute $e_i \neq 1$, $i = 1, \dots, n-1$ for all i in two parallel steps.
4. Compute $q_i \neq D_1 x$ for all i in one parallel operation.
5. Compute $e_i \neq E^{-1} y_1$ for all i using an add-scan operation (which is all by E^{-1} is).

¹These are some early notes written for J.-P. Brunet.

Addition/Subtraction: compute $z^{k_r} = (x \cdot r) \pm (y \cdot r)$. Statements in braces are unnecessary on machines that returned wrapped results on over/underflow. It assumes there are sticky overflow and underflow flags as in IEEE arithmetic. Multiplications and divisions by r can be done by modifying the exponent directly. It assumes round to nearest mode and flush to zero underflow (i.e. not gradual underflow), although the changes to account for other assumptions are simple. Besides r the machine constant $t = \text{underflow threshold/machine epsilon}$ will be used. I have arranged the "if" tests in decreasing order of likelihood of their being executed.

```

i f ( m = n ) t h e n
    z = x ± y
    k = m
    i f ( o v e r f l o w ) t h e n
        { z = ( x / r ) ± ( y / r ) }
        k = k + 1
    e l s e i f ( u n d e r f l o w ) t h e n
        { z = ( x · r ) ± ( y · r ) }
        k = k - 1
    e n d i f
e l s e i f ( m = n - 1 ) t h e n
    z = x ± ( y / r ) /* no overflow possible if round to nearest */
    k = n
    i f 0 < |z| < t t h e n
        z = ( x · r ) ± y
        k = m
    e n d i f
e l s e i f ( m = n + 1 ) t h e n
    z = y ± ( x / r ) /* no overflow possible if round to nearest */
    k = m
    i f 0 < |z| < t t h e n
        z = ( y · r ) ± x
        k = n
    e n d i f
e l s e i f ( m < n - 1 ) t h e n
    z = x
    k = n
e l s e i f ( n < m - 1 ) t h e n
    z = y
    k = m
e n d i f

```

It would be interesting to benchmark this parallel prefix operation the protection against over/underflow I propose here, to see how much us.

Safe Limits for Exponentiation			
IEEE Single		$r = 2$	$r = 2^{92}$
$f = 2^{28}$	$b = 16$	255	49150
	$b = 32$	1.67 · 10 ³	2.2 · 10 ⁹
$f = 2^{126}$	$b = 16$	260	49930
	$b = 32$	1.70 · 10 ³	2.7 · 10 ⁹
$f = 2^{149}$	$b = 16$	219	42223
	$b = 32$	1.44 · 10 ²	7.6 · 10 ⁹
IEEE Double		$r = 2$	$r = 2^{536}$
$f = 2^{1024}$	$b = 16$	31	49150
	$b = 32$	2.09 · 10 ³	2.2 · 10 ⁹
$f = 2^{1022}$	$b = 16$	32	49246
	$b = 32$	2.10 · 10 ³	2.2 · 10 ⁹
$f = 2^{1074}$	$b = 16$	30	46775
	$b = 32$	1.99 · 10 ³	0.6 · 10 ⁹

From this table, we see the limiting case is taking power of the safe number. When $b = 16$ we must clearly take the larger of the two r values to get a large safe exponent, and even then it is less than 50000. Choosing b is reasonable. Should we choose $r = 1$ or larger r ? I believe the larger r is better because of the larger parallel prefix operations it allows, and because it is easier to implement, since the representation of a number is almost unique: n ways to store a nonzero number in the form $f \cdot r^n$.

In order to implement our two operations it suffices to explain how to multiply numbers in this scaled format.

Multiplication: compute $z^k = (x \cdot r^m) \times (y \cdot r^n)$. Statements in braces are unnecessary on machines that returned wrapped results on over/underflow. It assumes there are sticky overflow and underflow flags as in IEEE arithmetic. Multiplications and divisions can be done by modifying the exponent directly.

```

z = x · y
k = m + n
if (overflow) then
    { z = ( x / r ) · ( y / r ) }
    k = k + 1
else if (underflow) then
    { z = ( x · r ) · ( y · r ) }
    k = k - 1
endif

```

recurrences are enough to do 2 term recurrences). At this time, however, there is no evidence that we frequently need to solve 3 or more term recurrences.

This leads us to propose the following building blocks for parallel prefix:

1. Scalar multiply parallel prefix with scaling to avoid over/underflow
2. 2 by 2 matrix multiply parallel prefix with scaling to avoid over/underflow

2 Specifications for Parallel Prefix

Basically, each floating point number is replaced by a pair of a floating point number and an integer, with the pair representing a value of the form $m \cdot 2^e$. The first problem is to choose r and the number of bits to store e for easy implementation and the ability to do very large parallel prefix operations without fear of over/underflow. So given r , the number of bits s in the exponent, an integer, and the largest and smallest positive possible values of a floating point number, we will ask how high a power of 2 the largest and smallest floating point numbers can be raised to without over/underflow in the form $m \cdot 2^e$.

We will only consider IEEE single and double precision formats. The numbers are given in the following table:

	IEEE Single	IEEE Double
Approximate overflow threshold	2^{128}	2^{1024}
Underflow threshold	2^{-126}	2^{-1022}
Smallest subnormal number	2^{-149}	2^{-1074}

Reasonable values for r from IEEE single precision are 1536 or 1537, and from IEEE double precision are 1536 or 1537. The source for these last two values is the wrapped exponential arithmetic: If overflow is trapped, the floating point unit is supposed to return the true answer in single precision and the true answer in double precision. Similarly, if underflow is trapped the value returned is supposed to be the true value times the true value.

Reasonable values for b , the number of bits in the mantissa, are 23 or 24.

The following table enumerates the approximate highest powers to which a number f can be safely raised using the scaled format as a function of n .

So we need to compute $N = (PM, P)$ where each is a 2 by 2 matrix, and \cdot is matrix multiply. Alternatively, we could use the following, equivalent

$$\begin{aligned} f &= \text{Par Pr } (\frac{\tau}{f}, x) \\ g &= \tau / f \quad (\text{componentwise vector division}) \\ h &= \text{Par Pr } (fg, x) \\ x &= h \cdot f \quad (\text{componentwise vector multiplication}) \end{aligned}$$

Unfortunately, this is very nonrobust because f frequently overflows. Even in IEEE double precision, with 308 bits, it does not take many consecutive floating point multiplies to get a number out of range. This can be partially mitigated by doing a Par (Progfix) operation, but this is not a satisfactory solution because it is slower and less accurate (see the appendix).

Let T be an n by n symmetric tridiagonal matrix with diagonal entries a_i and off-diagonal entries b_i . In order to find its eigenvalues, we need to solve the term recurrence

$$w_i = (a_i - \sigma) w_{i-1} - b_{i-1}^2 w_{i-2} \quad (2)$$

(The number of sign changes in the sequence is the number of eigenvalues less than σ .) This may be written in terms of parallel prefix as follows

$$\begin{bmatrix} w_i \\ w_{i-1} \end{bmatrix} = \begin{bmatrix} a_i - \sigma & -b_{i-1}^2 \\ 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} w_{i-1} \\ w_{i-2} \end{bmatrix} \equiv P_i \cdot \begin{bmatrix} w_{i-1} \\ w_{i-2} \end{bmatrix} = P_i \cdot P_{i-1} \cdots P_1 \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} \equiv Q_i \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad (3)$$

So we need to compute $Q = (PR, R)$, again a 2 by 2 matrix multiply parallel prefix.

This recurrence suffers from the same sensitivity to overflow/underflow. w_i is the determinant of the leading i by i submatrix of $T - \sigma I$, and so over/underflow even for matrices of modest norm and modest dimension. There is no known way to express its solution using scalar multiply parallel prefix as building blocks. In fact, we strongly suspect that no such exists, although we lack as yet a formal proof.

Furthermore, there is a theorem by H. T. Kung which says that if a recurrence relation where each term is a rational function, then it can be evaluated in faster than $\Omega(n)$ time if and only if it can be evaluated with matrix multiply parallel prefix. It turns out that the term recurrence

$$f_i(x_{i-1}) = \frac{\alpha_i x_{i-1} + \beta_i}{\gamma_i x_{i-1} + \delta_i}$$

which can be parallelized by computing x

$$\begin{bmatrix} u_i \\ v_i \end{bmatrix} = \begin{bmatrix} \alpha_i & \beta_i \\ \gamma_i & \delta_i \end{bmatrix} \cdot \begin{bmatrix} u_{i-1} \\ v_{i-1} \end{bmatrix} \equiv S_i \cdot \begin{bmatrix} u_{i-1} \\ v_{i-1} \end{bmatrix} = S_i \cdots S_1 \cdot \begin{bmatrix} x_0 \\ 1 \end{bmatrix}$$

Thus, 2 by 2 matrix multiply parallel prefix is sufficient (and we believe necessary) to parallelize all parallelizable scalar recurrence relations. On the other hand, it is adequate for 3 or more term recurrences (for the same reason I do

DRAFT

A Specification for Floating Point Parallel Prefix

James Demmel
Mathematics Department and Computer Science Division
University of California
Berkeley, CA 94720

July 3, 1992

Abstract

Parallel prefix is a useful operation for various linear algebra operations, including solving bidiagonal systems of equations and finding the eigenvalues of a symmetric tridiagonal matrix. However, the simplest implementations of parallel prefix for the operations of scalar floating point add and scalar floating point multiply are inadequate to solve these important problems. This is because they are too susceptible to over/underflow, and because they apparently cannot solve the general two term recurrence needed to find eigenvalues. In this note we propose a specification for parallel prefix operations overcoming these drawbacks.

1 Motivation

Our notation for parallel prefix will be as follows. Let r_1, \dots, r_n denote n vectors of data objects, which could be scalars or more complex objects. Let \otimes be an associative operator defined on these objects.

$$s_i = r_1 \otimes \dots \otimes r_i .$$

The most basic numerical parallel prefix operations on r_i and s_i should support floating point addition and floating point multiplication. However, floating point operations are not truly associative, and the impact on numerical analysis we will not pursue here.

Let B be an n by n bidiagonal matrix with diagonal entries τ_i and sub-diagonal entries η_i . To solve the linear system $Bx = y$, we need to solve the recurrence

$$x_i = \frac{\tau_{i-1}}{s_i} x_{i-1} + \frac{y_i}{s_i} \equiv \eta_i x_{i-1} + \tau_i \quad (1)$$

This may be done in two mathematically equivalent ways using parallel prefix. We have

$$\begin{bmatrix} x_i \\ 1 \end{bmatrix} = \begin{bmatrix} \eta_i & \tau_i \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x_{i-1} \\ 1 \end{bmatrix} \equiv M_i \cdot \begin{bmatrix} x_{i-1} \\ 1 \end{bmatrix} = M_i \cdot M_{i-1} \cdots M_1 \cdot \begin{bmatrix} 0 \\ 1 \end{bmatrix} \equiv N_i \cdot \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$