

Fixed-Point Glue Setting

	Section	Page
Introduction	1	1
The problem and a solution	4	2
Glue multiplication	8	4
Glue setting	12	5
Glue-set printing	15	6
The driver program	20	7
Index	27	9

The preparation of this report was supported in part by the National Science Foundation under grants IST-7921977 and MCS-7723728; by Office of Naval Research grant N00014-81-K-0330; and by the IBM Corporation. ‘TeX’ is a trademark of the American Mathematical Society.

1. Introduction. If $\text{T}_{\text{E}}\text{X}$ is being implemented on a microcomputer that does 32-bit addition and subtraction, but with multiplication and division restricted to multipliers and divisors that are either powers of 2 or positive integers less than 2^{15} , it can still do the computations associated with the setting of glue in a suitable way. This program illustrates one solution to the problem.

Another purpose of this program is to provide the first “short” example of the use of **WEB**.

2. The program itself is written in standard Pascal. It begins with a normal program header, most of which will be filled in with other parts of this “web” as we are ready to introduce them.

```
program GLUE(input, output);  
  type  $\langle$ Types in the outer block 6 $\rangle$   
  var  $\langle$ Globals in the outer block 8 $\rangle$   
  procedure initialize; { this procedure gets things started }  
    var  $\langle$ Local variables for initialization 9 $\rangle$   
    begin  $\langle$ Set initial values 10 $\rangle$ ;  
  end;
```

3. Here are two macros for common programming idioms.

```
define incr(#)  $\equiv$  #  $\leftarrow$  # + 1 { increase a variable by unity }  
define decr(#)  $\equiv$  #  $\leftarrow$  # - 1 { decrease a variable by unity }
```

4. The problem and a solution. We are concerned here with the “setting of glue” that occurs when a \TeX box is being packaged. Let x_1, \dots, x_n be integers whose sum $s = x_1 + \dots + x_n$ is positive, and let t be another positive integer. These x_i represent scaled amounts of glue in units of sp (scaled points), where one sp is 2^{-16} of a printer’s point. The other quantity t represents the total by which the glue should stretch or shrink. Following the conventions of \TeX 82, we will assume that the integers we deal with are less than 2^{31} in absolute value.

After the glue has been set, the actual amounts of incremental glue space (in sp) will be the integers $f(x_1), \dots, f(x_n)$, where f is a function that we wish to compute. We want $f(x)$ to be nearly proportional to x , and we also want the sum $f(x_1) + \dots + f(x_n)$ to be nearly equal to t . If we were using floating-point arithmetic, we would simply compute $f(x) \equiv (t/s) \cdot x$ and hope for the best; but the goal here is to compute a suitable f using only the fixed-point arithmetic operations of a typical “16-bit microcomputer.”

The solution adopted here is to determine integers a, b, c such that

$$f(x) = \lfloor 2^{-b} c \lfloor 2^{-a} x \rfloor \rfloor$$

if x is nonnegative. Thus, we take x and shift it right by a bits, then multiply by c (which is 2^{15} or less), and shift the product right by b bits. The quantities a, b , and c are to be chosen so that this calculation doesn’t cause overflow and so that $f(x_1) + \dots + f(x_n)$ is reasonably close to t .

The following method is used to calculate a and b : Suppose

$$y = \max_{1 \leq i \leq n} |x_i|.$$

Let d and e be the smallest integers such that $t < 2^d s$ and $y < 2^e$. Since s and t are less than 2^{31} , we have $-30 \leq d \leq 31$ and $1 \leq e \leq 31$. An error message is given if $d + e \geq 31$; in such a case some x_m has $|x_m| \geq 2^{e-1}$ and we are trying to change $|x_m|$ to $|(t/s)x_m| \geq 2^{d+e-2} \geq 2^{30}$ sp, which \TeX does not permit. (Consider, for example, the “worst case” situation $x_1 = 2^{30} + 1$, $x_2 = -2^{30}$, $t = 2^{31} - 1$; surely we need not bother trying to accommodate such anomalous combinations of values.) On the other hand if $d + e \leq 31$, we set $a = e - 16$ and $b = 31 - d - e$. Notice that this choice of a guarantees that $\lfloor 2^{-a} |x_i| \rfloor < 2^{16}$. We will choose c to be at most 2^{15} , so that the product will be less than 2^{31} .

The computation of c is the tricky part. The “ideal” value for c would be $\rho = 2^{a+b}t/s$, since $f(x)$ should be approximately $(t/s) \cdot x$. Furthermore it is better to have c slightly larger than ρ , instead of slightly smaller, since the other operations in $f(x)$ have a downward bias. Therefore we shall compute $c = \lceil \rho \rceil$. Since $2^{a+b}t/s < 2^{a+b+d} = 2^{15}$, we have $c \leq 2^{15}$ as desired.

We want to compute $c = \lceil \rho \rceil$ exactly in all cases. There is no difficulty if $s < 2^{15}$, since c can be computed directly using the formula $c = \lfloor (2^{a+b}t + s - 1)/s \rfloor$; overflow will not occur since $2^{a+b}t < 2^{15}s < 2^{30}$.

Otherwise let $s = s_1 2^l + s_2$, where $2^{14} \leq s_1 < 2^{15}$ and $0 \leq s_2 < 2^l$. We will essentially carry out a long division. Let t be “normalized” so that $2^{30} \leq 2^h t < 2^{31}$ for some h . Then we form the quotient and remainder of $2^h t$ divided by s_1 ,

$$2^h t = q s_1 + r_0, \quad 0 \leq r_0 < s_1.$$

It follows that $2^{h+l}t - q s = 2^l r_0 - q s_2 = r$, say. If $0 \geq r > -s$ we have $q = \lceil 2^{h+l}t/s \rceil$; otherwise we can replace (q, r) by $(q \pm 1, r \mp s)$ repeatedly until r is in the correct range. It is not difficult to prove that q needs to be increased at most once and decreased at most seven times, since $2^l r_0 - q s_2 < 2^l s_1 \leq s$ and since $q s_2/s \leq (2^h t/s_1)(s_2/2^l s_1) < 2^{31}/s_1^2 \leq 8$. Finally, we have $a + b - h - l = -1$ or -2 , since $2^{28+l} \leq 2^{14}s = 2^{a+b+d-1}s \leq 2^{a+b}t < 2^{a+b+d}s = 2^{15}s < 2^{30+l}$ and $2^{30} \leq 2^h t < 2^{31}$. Hence $c = \lceil 2^{a+b-h-l}q \rceil = \lceil \frac{1}{2}q \rceil$ or $\lceil \frac{1}{4}q \rceil$.

An error analysis shows that these values of a, b , and c work satisfactorily, except in unusual cases where we wouldn’t expect them to. When $x \geq 0$ we have

$$\begin{aligned} f(x) &= 2^{-b}(2^{a+b}t/s + \theta_0)(2^{-a}x - \theta_1) - \theta_2 \\ &= (t/s)x + \theta_0 2^{-a-b}x - \theta_1 2^a t/s - 2^{-b}\theta_0\theta_1 - \theta_2 \end{aligned}$$

where $0 \leq \theta_0, \theta_1, \theta_2 < 1$. Now $0 \leq \theta_0 2^{-a-b}x < 2^{e-a-b} = 2^{d+e-15}$ and $0 \leq \theta_1 2^a t/s < 2^{a+d} = 2^{d+e-16}$, and the other two terms are negligible. Therefore $f(x_1) + \dots + f(x_n)$ differs from t by at most about $2^{d+e-15}n$. Since 2^{d+e} is larger than $(t/s)y$, which is the largest stretching or shrinking of glue after expansion, the error is at worst about $n/32000$ times as much as this, so it is quite reasonable. For example, even if fill glue is being used to stretch 20 inches, the error will still be less than $\frac{1}{1600}$ of an inch.

5. To sum up: Given the positive integers s , t , and y as above, we set

$$a \leftarrow \lfloor \lg y \rfloor - 15, \quad b \leftarrow 29 - \lfloor \lg y \rfloor - \lfloor \lg t/s \rfloor, \quad \text{and} \quad c \leftarrow \lceil 2^{a+b} t/s \rceil.$$

The implementation below shows how to do the job in Pascal without using large numbers.

6. T_EX wants to have the glue-setting information in a 32-bit data type called *glue_ratio*. The Pascal implementation of T_EX82 has *glue_ratio* = *real*, but alternative definitions of *glue_ratio* are explicitly allowed.

For our purposes we shall let *glue_ratio* be a record that is packed with three fields: The *a_part* will hold the positive integer $a + 16$, the *b_part* will hold the nonnegative integer b , and the *c_part* will hold the nonnegative integer c . When the formulas above tell us to take $b > 30$, we might as well set $c \leftarrow 0$ instead, because $f(x)$ will be zero in all cases when $b > 30$. Note that we have only about 25 bits of information in all, so it should fit in 32 bits with ease.

```

⟨Types in the outer block 6⟩ ≡
  glue_ratio = packed record a_part: 1 .. 31; { the quantity  $e = a + 16$  in our derivation }
               b_part: 0 .. 30; { the quantity  $b$  in our derivation }
               c_part: 0 .. '100000; { the quantity  $c$  in our derivation }
  end;
  scaled = integer; { this data type is used for quantities in sp units }

```

This code is used in section 2.

7. The real problem is to define the procedures that T_EX needs to deal with such *glue_ratio* values: (a) Given scaled numbers s , t , and y as above, to compute the corresponding *glue_ratio*. (b) Given a nonnegative scaled number x and a *glue_ratio* g , to compute the scaled number $f(x)$. (c) Given a *glue_ratio* g , to print out a decimal equivalent of g for diagnostic purposes.

The procedures below can be incorporated into T_EX82 via a change file without great difficulty. A few modifications will be needed, because T_EX's *glue_ratio* values can be negative in unusual cases—when the amount of stretchability or shrinkability is less than zero. Negative values in the *c_part* will handle such problems, if proper care is taken. The error message below should either become a warning message or a call to T_EX's *print_err* routine; in the latter case, an appropriate help message should be given, stating that glue cannot stretch to more than 18 feet long, but that it's OK to proceed with fingers crossed.

8. Glue multiplication. The easiest procedure of the three just mentioned is the one that is needed most often, namely, the computation of $f(x)$.

Pascal doesn't have built-in binary shift commands or built-in exponentiation, although many computers do have this capability. Therefore our arithmetic routines use an array called '*two_to_the*', containing powers of two. Divisions by powers of two are never done in the programs below when the dividend is negative, so the operations can safely be replaced by right shifts on machines for which this is most appropriate. (Contrary to popular opinion, the operation ' $x \text{ div } 2$ ' is not the same as shifting x right one binary place, on a machine with two's complement arithmetic, when x is a negative odd integer. But division *is* equivalent to shifting when x is nonnegative.)

⟨Globals in the outer block 8⟩ \equiv
two_to_the: **array** [0 .. 30] **of** *integer*; { *two_to_the*[k] = 2^k }

See also sections 15 and 20.

This code is used in section 2.

9. ⟨Local variables for initialization 9⟩ \equiv
 k : 1 .. 30; { an index for initializing *two_to_the* }

This code is used in section 2.

10. ⟨Set initial values 10⟩ \equiv
two_to_the[0] \leftarrow 1;
for $k \leftarrow 1$ **to** 30 **do** *two_to_the*[k] \leftarrow *two_to_the*[$k - 1$] + *two_to_the*[$k - 1$];

This code is used in section 2.

11. We will use the abbreviations *ga*, *gb*, and *gc* as convenient alternatives to Pascal's **with** statement. The glue-multiplication function f , which replaces several occurrences of the '*float*' macro in T_EX82, is now easy to state:

```
define ga  $\equiv$  g.a_part
define gb  $\equiv$  g.b_part
define gc  $\equiv$  g.c_part
function glue_mult(x : scaled; g : glue_ratio): integer; { returns  $f(x)$  as above, assuming that  $x \geq 0$  }
begin if ga > 16 then  $x \leftarrow x \text{ div } \textit{two\_to\_the}[\textit{ga} - 16]$  { right shift by  $a$  places }
else  $x \leftarrow x * \textit{two\_to\_the}[16 - \textit{ga}]$ ; { left shift by  $-a$  places }
glue_mult  $\leftarrow (x * \textit{gc}) \text{ div } \textit{two\_to\_the}[\textit{gb}]$ ; { right shift by  $b$  places }
end; { note that  $b$  may be as large as 30 }
```

12. Glue setting. The *glue_fix* procedure computes a , b , and c by the method explained above. \TeX does not normally compute the quantity y , but it could be made to do so without great difficulty.

This procedure replaces several occurrences of the ‘*unfloat*’ macro in \TeX 82. It would be written as a function that returns a *glue_ratio*, if Pascal would allow functions to produce records as values.

```

procedure glue_fix( $s, t, y$  : scaled; var  $g$  : glue_ratio);
var  $a, b, c$  : integer; { components of the desired ratio }
     $k, h$  : integer; {  $30 - \lfloor \lg s \rfloor, 30 - \lfloor \lg t \rfloor$  }
     $s0$  : integer; { original (unnormalized) value of  $s$  }
     $q, r, s1$  : integer; { quotient, remainder, divisor }
     $w$  : integer; {  $2^l$ , where  $l = 16 - k$  }
begin  $\langle$  Normalize  $s$ ,  $t$ , and  $y$ , computing  $a$ ,  $k$ , and  $h$  13  $\rangle$ ;
if  $t < s$  then  $b \leftarrow 15 - a - k + h$  else  $b \leftarrow 14 - a - k + h$ ;
if  $(b < 0) \vee (b > 30)$  then
    begin if  $b < 0$  then write_ln('!_Excessive_glue. '); { error message }
     $b \leftarrow 0$ ;  $c \leftarrow 0$ ; { make  $f(x)$  identically zero }
    end
else begin if  $k \geq 16$  then { easy case,  $s_0 < 2^{15}$  }
     $c \leftarrow (t \text{ div } \text{two\_to\_the}[h - a - b] + s0 - 1) \text{ div } s0$  { here  $1 \leq h - a - b \leq k - 14 \leq 16$  }
    else  $\langle$  Compute  $c$  by long division 14  $\rangle$ ;
    end;
 $ga \leftarrow a + 16$ ;  $gb \leftarrow b$ ;  $gc \leftarrow c$ ;
end;

```

```

13.  $\langle$  Normalize  $s$ ,  $t$ , and  $y$ , computing  $a$ ,  $k$ , and  $h$  13  $\rangle \equiv$ 
begin  $a \leftarrow 15$ ;  $k \leftarrow 0$ ;  $h \leftarrow 0$ ;  $s0 \leftarrow s$ ;
while  $y < '10000000000$  do {  $y$  is known to be positive }
    begin decr( $a$ );  $y \leftarrow y + y$ ;
    end;
while  $s < '10000000000$  do {  $s$  is known to be positive }
    begin incr( $k$ );  $s \leftarrow s + s$ ;
    end;
while  $t < '10000000000$  do {  $t$  is known to be positive }
    begin incr( $h$ );  $t \leftarrow t + t$ ;
    end;
end { now  $2^{30} \leq t = 2^h t_0 < 2^{31}$  and  $2^{30} \leq s = 2^k s_0 < 2^{31}$ , hence  $d = k - h$  if  $t/s < 1$  }

```

This code is used in section 12.

```

14.  $\langle$  Compute  $c$  by long division 14  $\rangle \equiv$ 
begin  $w \leftarrow \text{two\_to\_the}[16 - k]$ ;  $s1 \leftarrow s0$  div  $w$ ;  $q \leftarrow t \text{ div } s1$ ;  $r \leftarrow ((t \text{ mod } s1) * w) - ((s0 \text{ mod } w) * q)$ ;
if  $r > 0$  then
    begin incr( $q$ );  $r \leftarrow r - s0$ ;
    end
else while  $r \leq -s0$  do
    begin decr( $q$ );  $r \leftarrow r + s0$ ;
    end;
if  $a + b + k - h = 15$  then  $c \leftarrow (q + 1) \text{ div } 2$  else  $c \leftarrow (q + 3) \text{ div } 4$ ;
end

```

This code is used in section 12.

15. Glue-set printing. The last of the three procedures we need is *print_gr*, which displays a *glue_ratio* in symbolic decimal form. Before constructing such a procedure, we shall consider some simpler routines, copying them from an early draft of the program T_EX82.

```
define unity  $\equiv$  '200000 {  $2^{16}$ , represents 1.0000 }
⟨Globals in the outer block 8⟩ +≡
dig: array [0 .. 15] of 0 .. 9; { for storing digits }
```

16. An array of digits is printed out by *print_digs*.

```
procedure print_digs(k : integer); { prints dig[k - 1] ... dig[0] }
begin while k > 0 do
  begin decr(k); write(chr(ord('0') + dig[k]));
  end;
end;
```

17. A nonnegative integer is printed out by *print_int*.

```
procedure print_int(n : integer); { prints an integer in decimal form }
var k: 0 .. 12; { index to current digit; we assume that  $0 \leq n < 10^{12}$  }
begin k  $\leftarrow$  0;
repeat dig[k]  $\leftarrow$  n mod 10; n  $\leftarrow$  n div 10; incr(k);
until n = 0;
print_digs(k);
end;
```

18. And here is a procedure to print a nonnegative *scaled* number.

```
procedure print_scaled(s : scaled); { prints a scaled real, truncated to four digits }
var k: 0 .. 3; { index to current digit of the fraction part }
begin print_int(s div unity); { print the integer part }
s  $\leftarrow$  ((s mod unity) * 10000) div unity;
for k  $\leftarrow$  0 to 3 do
  begin dig[k]  $\leftarrow$  s mod 10; s  $\leftarrow$  s div 10;
  end;
write(' '); print_digs(4);
end;
```

19. Now we're ready to print a *glue_ratio*. Since the effective multiplier is $2^{-a-b}c$, we will display the scaled integer $2^{16-a-b}c$, taking care to print something special if this quantity is terribly large.

```
procedure print_gr(g : glue_ratio); { prints a glue multiplier }
var j: -29 .. 31; { the amount to shift c }
begin j  $\leftarrow$  32 - ga - gb;
while j > 15 do
  begin write('2x'); decr(j); { indicate multiples of 2 for BIG cases }
  end;
if j < 0 then print_scaled(gc div two_to_the[-j]) { shift right }
else print_scaled(gc * two_to_the[j]); { shift left }
end;
```

20. The driver program. In order to test these routines, we will assume that the *input* file contains a sequence of test cases, where each test case consists of the integer numbers $t, x_1, \dots, x_n, 0$. The final test case should be followed by an additional zero.

⟨Globals in the outer block 8⟩ \equiv
 x : **array** [1 .. 1000] **of** *scaled*; { the x_i }
 t : *scaled*; { the desired total }
 m : *integer*; { the test case number }

21. Each case will be processed by the following routine, which assumes that t has already been read.

procedure *test*; { processes the next data set, given t and m }
 var n : 0 .. 1000; { the number of items }
 k : 0 .. 1000; { runs through the items }
 y : *scaled*; { $\max_{1 \leq i \leq n} |x_i|$ }
 g : *glue_ratio*; { the computed glue multiplier }
 s : *scaled*; { the sum $x_1 + \dots + x_n$ }
 ts : *scaled*; { the sum $f(x_1) + \dots + f(x_n)$ }
 begin *write_ln*(*Test_data_set_number*, m : 1, ^: ^); ⟨Read x_1, \dots, x_n 22⟩;
 ⟨Compute s and y 23⟩;
 if $s \leq 0$ **then** *write_ln*(*Invalid_data(nonpositive sum); this set rejected.* ^)
 else begin ⟨Compute g and print it 24⟩;
 ⟨Print the values of $x_i, f(x_i)$, and the totals 25⟩;
 end;
end;

22. ⟨Read x_1, \dots, x_n 22⟩ \equiv
 begin $n \leftarrow 0$;
 repeat *incr*(n); *read*($x[n]$);
 until $x[n] = 0$;
 decr(n);
 end

This code is used in section 21.

23. ⟨Compute s and y 23⟩ \equiv
 begin $s \leftarrow 0$; $y \leftarrow 0$;
 for $k \leftarrow 1$ **to** n **do**
 begin $s \leftarrow s + x[k]$;
 if $y < \text{abs}(x[k])$ **then** $y \leftarrow \text{abs}(x[k])$;
 end;
 end

This code is used in section 21.

24. ⟨Compute g and print it 24⟩ \equiv
 begin *glue_fix*(s, t, y, g); { set g , perhaps print an error message }
 write(*Glue_ratio_is* ^); *print_gr*(g); *write_ln*(*^ (^, ga - 16 : 1, ^, ^, gb : 1, ^, ^, gc : 1, ^) ^*);
 end

This code is used in section 21.

25. $\langle \text{Print the values of } x_i, f(x_i), \text{ and the totals } 25 \rangle \equiv$

```

begin  $ts \leftarrow 0$ ;
for  $k \leftarrow 1$  to  $n$  do
  begin  $write(x[k] : 20)$ ;
  if  $x[k] \geq 0$  then  $y \leftarrow glue\_mult(x[k], g)$ 
  else  $y \leftarrow -glue\_mult(-x[k], g)$ ;
   $write\_ln(y : 15)$ ;  $ts \leftarrow ts + y$ ;
  end;
 $write\_ln(\text{'Totals', } s : 13, ts : 15, \text{'(versus', } t : 1, \text{' )'})$ ;
end

```

This code is used in section 21.

26. Here is the main program.

```

begin  $initialize$ ;  $m \leftarrow 1$ ;  $read(t)$ ;
while  $t > 0$  do
  begin  $test$ ;  $incr(m)$ ;  $read(t)$ ;
  end;
end.

```

27. Index. Here are the section numbers where various identifiers are used in the program, and where various topics are discussed.

a: [12](#).
a_part: [6](#), [11](#).
abs: [23](#).
b: [12](#).
b_part: [6](#), [11](#).
c: [12](#).
c_part: [6](#), [7](#), [11](#).
chr: [16](#).
decr: [3](#), [13](#), [14](#), [16](#), [19](#), [22](#).
dig: [15](#), [16](#), [17](#), [18](#).
error analysis: [4](#).
error message: [7](#), [12](#).
float: [11](#).
g: [11](#), [12](#), [19](#), [21](#).
ga: [11](#), [12](#), [19](#), [24](#).
gb: [11](#), [12](#), [19](#), [24](#).
gc: [11](#), [12](#), [19](#), [24](#).
GLUE: [2](#).
glue_fix: [12](#), [24](#).
glue_mult: [11](#), [25](#).
glue_ratio: [6](#), [7](#), [11](#), [12](#), [15](#), [19](#), [21](#).
h: [12](#).
hairly mathematics: [4](#).
incr: [3](#), [13](#), [14](#), [17](#), [22](#), [26](#).
initialize: [2](#), [26](#).
input: [2](#), [20](#).
integer: [6](#), [8](#), [11](#), [12](#), [16](#), [17](#), [20](#).
j: [19](#).
k: [9](#), [12](#), [16](#), [17](#), [18](#).
m: [20](#).
main program: [26](#).
n: [17](#), [21](#).
ord: [16](#).
output: [2](#).
print_digs: [16](#), [17](#), [18](#).
print_err: [7](#).
print_gr: [15](#), [19](#), [24](#).
print_int: [17](#), [18](#).
print_scaled: [18](#), [19](#).
program header: [2](#).
q: [12](#).
r: [12](#).
read: [22](#), [26](#).
real: [6](#).
s: [12](#), [21](#).
scaled: [6](#), [11](#), [12](#), [18](#), [20](#), [21](#).
s0: [12](#), [13](#), [14](#).
s1: [12](#), [14](#).
t: [12](#), [20](#).
test: [21](#), [26](#).
ts: [21](#), [25](#).
two_to_the: [8](#), [9](#), [10](#), [11](#), [12](#), [14](#), [19](#).
unfloat: [12](#).
unity: [15](#), [18](#).
w: [12](#).
write: [16](#), [18](#), [19](#), [24](#), [25](#).
write_ln: [12](#), [21](#), [24](#), [25](#).
x: [11](#), [20](#).
y: [12](#).

- ⟨ Compute c by long division 14 ⟩ Used in section 12.
- ⟨ Compute g and print it 24 ⟩ Used in section 21.
- ⟨ Compute s and y 23 ⟩ Used in section 21.
- ⟨ Globals in the outer block 8, 15, 20 ⟩ Used in section 2.
- ⟨ Local variables for initialization 9 ⟩ Used in section 2.
- ⟨ Normalize s , t , and y , computing a , k , and h 13 ⟩ Used in section 12.
- ⟨ Print the values of x_i , $f(x_i)$, and the totals 25 ⟩ Used in section 21.
- ⟨ Read x_1, \dots, x_n 22 ⟩ Used in section 21.
- ⟨ Set initial values 10 ⟩ Used in section 2.
- ⟨ Types in the outer block 6 ⟩ Used in section 2.