

# A Decimal Floating-Point Specification

<http://www2.hursley.ibm.com/decimal/>

Eric M. Schwarz

Michael F. Cowlishaw (IBM UK),

Ronald M. Smith, Charles F. Webb

eServer Microprocessor Development

IBM Corporation

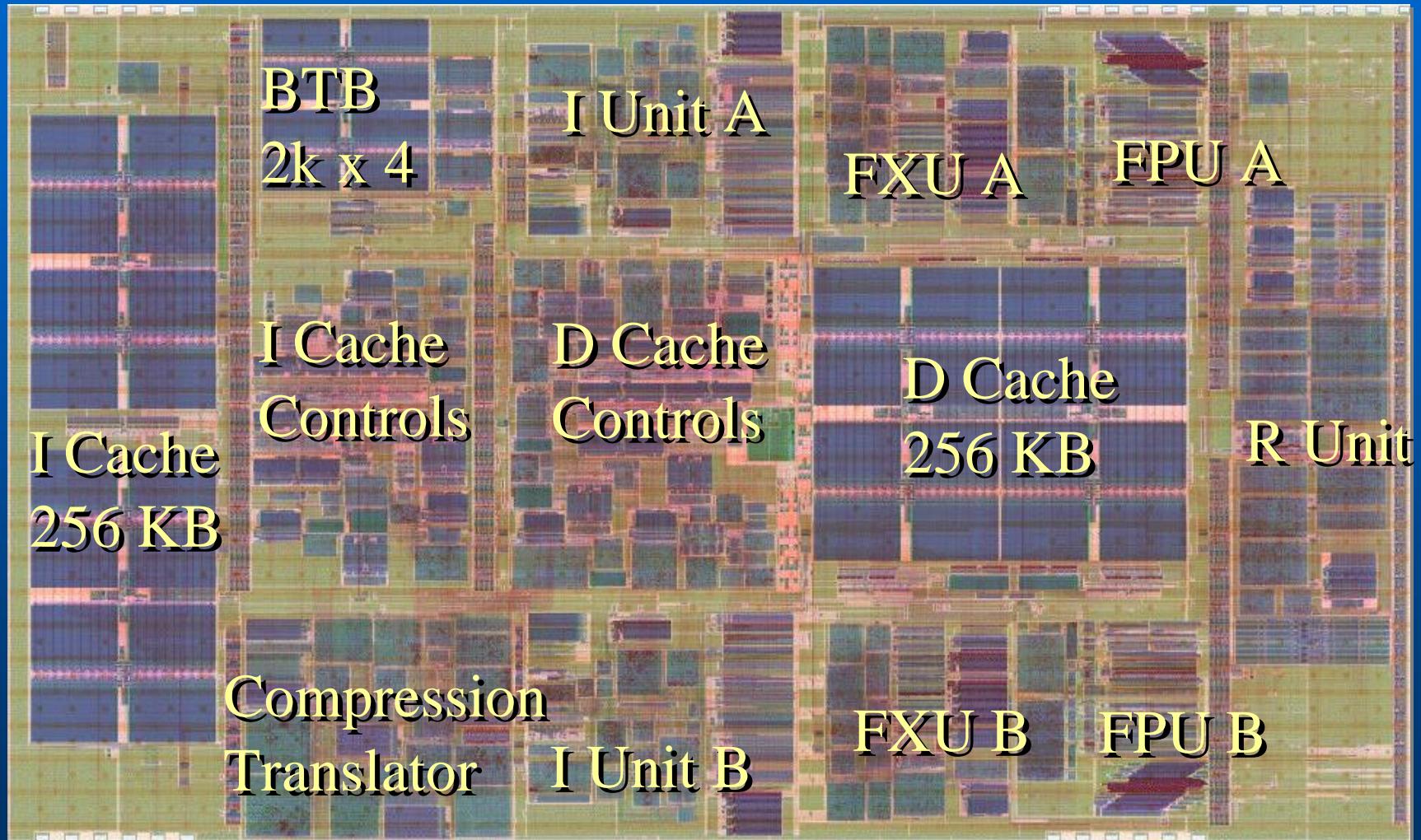


# Motivation

- People think in terms of decimals
- BFP can't represent decimal fractions such as 0.1
- Scaling of BFP requires rounding
- Decimal data is common, of all numeric data from commercial databases – 55% decimal, 26.4% small int, 17.3% integer, and 1.4% binary float
- How do we get decimal arithmetic into hardware?



# z900 Microprocessor



17.9mm x 9.9mm – 47 million transistors

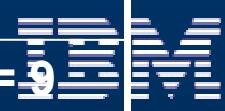
IBM

# Goals of Decimal Floating Point Specification

- It should allow efficient hardware or software implementation
- It should support numbers to be used in ANSI X3.274-1996 (REXX)
- It should support numbers to be used in ANSI/IEEE 854-1987
- It should allow efficient use of existing databases
- It should allow for future expansion



# Software Requirements

<b>C &amp; C++</b>	<b>S/390</b> <b>AS/400</b> <b>Others</b>	<b>Fixed</b>	<b>31</b> <b>31</b> <b>38</b>	<b>C only</b> <b>C only</b> <b>In libraries</b>
<b>COBOL</b>	<b>All</b>	<b>Fixed</b>	<b>31</b>	<b>32 dig Float</b>
<b>C#</b>	<b>All</b>	<b>Floating</b>	<b>28</b>	
<b>Java</b>	<b>All</b>	<b>Floating</b>	<b>Infinite</b>	<b>BigDecimal</b>
<b>OS/400 CL</b>	<b>AS/400</b>	<b>Fixed</b>	<b>15</b>	<b>Scale &lt;= 9</b>
<b>PL/I</b>	<b>S/390</b> <b>AS/400</b> <b>Others</b>	<b>Fixed</b>	<b>15</b> <b>15</b> <b>31</b>	<b>Decimal</b> <b>Float is actually binary</b>
<b>PSM</b>	<b>All</b>	<b>Fixed</b>	<b>31</b>	<b>Impl in C</b>
<b>Rexx</b>	<b>All</b>	<b>Floating</b>	<b>Infinite</b>	
<b>RPG</b>	<b>AS/400</b>	<b>Fixed</b>	<b>30</b>	<b>Scale &lt;= 9</b> 

# Decimal Arithmetic Specification

- Operations – add, subtract, plus, minus, multiply, divide, power, divide-integer, remainder, compare
- Format – scientific, engineering, numeric
- Rounding modes – half up, half down, half-even
- Precision control
- Flags - inexact
- Machine (concrete) representation



# Machine Representation

- Sign
- Integer
- Exponent
- $(-1)^{\text{sign}} * (\text{integer}) * 10^{\text{exponent}}$
- Common data lengths are 4B, 8B, 10B, 16B - choose 8B and 16B.



# Integer Representation

- **Binary format**
  - Compact
  - Fast multiplication and division
  - Scaling, normalization, threshold detection slow
  - Requires conversion from stored BCD data
- **Binary Coded Decimal**
  - No conversion needed
  - Rounding, normalization easy
  - Moderately slower arithmetic than binary
  - Wasteful, uses only 62.5% space
- **Compressed BCD**
  - Minor conversion stage needed
  - Advantages of BCD for scaling and of Binary for compactness



# Compressed BCD Format

- 10 binary bits for 3 decimal digits
- 7 bits for 2 digits

Decimal Digits      ->    binary bits

1        2        3

Abcd efgh ijkl -> p qrs tuv wxy



# Compression example

Abcd efgh ijk1 -> p qrs tuv wxy

<b>aei</b>	<b>P qrs tuv wxy</b>
<b>000</b>	<b>0 bcd fgh jkl</b>
<b>100</b>	<b>1 00d fgh jkl</b>
<b>010</b>	<b>1 01d bch jkl</b>
<b>001</b>	<b>1 10d fgh bcl</b>
<b>011</b>	<b>1 11d 00h bcl</b>
<b>101</b>	<b>1 11d 01h fgl</b>
<b>110</b>	<b>1 11d 10h jkl</b>
<b>111</b>	<b>1 11d 11h 00l</b>

734d = 0111 0011 0100  
-> 0 111 011 100



# Compression logic

Compress 12 bit BCD to 10 bit binary encoding:

$$p = a \mid e \mid i$$

$$q = (b \ e') \mid i \mid (a \ e)$$

$$r = (c \ i') \mid e \mid (a \ i)$$

$$s = d$$

$$t = (a \ e) \mid (f \ (a' \mid i')) \mid (b \ e \ i')$$

$$u = (a \ i) \mid (c \ e \ i') \mid (g \ e')$$

$$v = h$$

$$w = j \mid (b \ i) \mid (f \ a \ i)$$

$$x = k \mid (c \ i) \mid (g \ a \ i)$$

$$y = l$$



# Expansion Logic

Expand 10 bit binary encoding to 12 bit BCD:

$$a = (p \ q' \ r') \mid (p \ q \ r \ (t \mid u))$$

$$b = (q \ p') \mid (t \ p \ q' \ r) \mid (w \ q \ (r' \mid (t' \ u')))$$

$$c = (r \ (p' \mid (q' \ u))) \mid (x \ p \ q \ (r' \mid (t' \ u')))$$

$$d = s$$

$$e = (p \ r \ (q' \mid u' \mid t))$$

$$f = (t \ (p' \mid r')) \mid (p \ q \ r \ t' \ u \ w)$$

$$g = (u \ (p' \mid r')) \mid (p \ q \ r \ t' \ u \ x)$$

$$h = v$$

$$i = (p \ q \ (r' \mid t' \mid u))$$

$$j = (w \ (p' \mid q' \mid (r \ t)))$$

$$k = (x \ (p' \mid q' \mid (r \ t)))$$

$$l = y$$



# Exponent Representation

- **Binary two's complement**
  - Allows for fast addition
  - Compact
- **Binary unsigned with bias**
  - Similar advantage to above
  - Similar to IEEE BFP
- **BCD**
  - No conversion
  - Inefficient and slow



# Special Values

- Reserved Exponent Values
  - Similar to IEEE BFP
  - Compact
- Using separate bits
  - Easier to detect
  - wasteful



# Length of Exponent

- Must meet IEEE 854 standard
  - $(E_{max} - E_{min}) > 5 \times \text{Digits}$  prefer  $10 \times D$
  - Double precision  $8 \times E_{max}$  single + 7
  - Recommend  $E_{min} = -E_{max}$
- Short (8B) – 11 bit (-999 to +999-(D-1))
- Long(16B) – 15 bit (-9999 to +9999-(D-1))



# Decimal Digits Possible

	BCD	7:2	10:3	Binary
<b>Single (52 bits)</b>	<b>13 digits (0)</b>	<b>14 digits (3)</b>	<b>15 digits (2)</b>	<b>15 digits (2)</b>
<b>Double (112 bits)</b>	<b>28 digits (0)</b>	<b>32 digits (0)</b>	<b>33 digits (2)</b>	<b>33 digits (2)</b>



# Proposed Decimal Format

Sign	Exponent	Filler	Integer
1 bit	11 bit unsigned binary w/bias	2 bit	15 digit (Cmpr BCD 10:3)
1 bit	15 bit unsigned binary w/bias	2 bit	33 digit (Cmpr BCD 10:3)



# Remarks

- Determined a possible machine representation
- Would like multiple manufactures to support a representation
- Need to keep software standards inline with hardware limitations
- Encourage hardware implementation of decimal floating point
- <http://www2.hursley.ibm.com/decimal/>

