

Fixed Point Arithmetic and the Packed Decimal Format

This set of lectures discusses the idea of **fixed point arithmetic** and its implementation by Packed Decimal Format on the IBM Mainframes.

The topics include

1. A review of IBM S/370 floating point formats, focusing on the precision with which real numbers can be stored.
2. The difference between the precision requirements of business applications and those of scientific applications.
3. An overview of the Packed Decimal format, as implemented on the IBM S/370 and predecessors.

We begin with a review of IBM notation for denoting bit numbers in a byte or word. From our viewpoint, it is a bit non-standard.

IBM S/370: Terminology and Notation

The IBM 370 is a byte-addressable machine; each byte has a unique address.

The standard storage sizes on the IBM 370 are byte, halfword, and fullword.

Byte	8 binary bits	
Halfword	16 binary bits	2 bytes
Fullword	32 binary bits	4 bytes.

In IBM terminology, the leftmost bit is bit zero, so we have the following.

Byte

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

Halfword

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
---	---	---	---	---	---	---	---	---	---	----	----	----	----	----	----

Fullword

0 – 7	8 – 15	16 – 23	24 – 31
-------	--------	---------	---------

Comment: The IBM 370 seems to be a “big endian” machine.

S/370: Available Floating Point Formats

There are three available formats for representing floating point numbers.

Single precision	4 bytes	32 bits: 0 – 31
Double precision	8 bytes	64 bits: 0 – 63
Extended precision	16 bytes	128 bits; 0 – 127.

The standard representation of the fields is as follows.

Format	Sign bit	Exponent bits	Fraction bits
Single	0	1 – 7	8 – 31
Double	0	1 – 7	8 – 63
Extended	0	1 – 7	8 – 127

NOTE: Unlike the IEEE–754 format, greater precision is not accompanied by a greater range of exponents.

The precision of the format depends on the number of bits used for the fraction.

Single precision	24 bit fraction	1 part in 2^{24}	7 digits precision *
Double precision	56 bit fraction	1 part in 2^{56}	16 digits precision **

* $2^{24} = 16,777,216$ ** $2^{56} \approx (10^{0.30103})^{56} \approx 10^{16.86} \approx 7 \cdot 10^{16}$.

Precision Example: Slightly Exaggerated

Consider a banking problem. Banks lend each other money overnight.

At 3% annual interest, the overnight interest on \$1,000,000 is \$40.492.

Suppose my bank lends your bank \$10,000,000 (ten million).

You owe me \$404.92 in interest; \$10,000,404.92 in total.

With seven significant digits, the amount might be calculated as \$10,000,400.

My bank loses \$4.92.

I want my books to balance to the penny. I do not like floating point arithmetic.

TRUE STORY

When DEC (the Digital Equipment Corporation) was marketing their PDP-11 to a large New York bank, it supported integer and floating point arithmetic.

At this time, the PDP-11 did not support decimal arithmetic.

The bank told DEC something like this:

“Add decimal arithmetic and we shall buy a few thousand. Without it – no sale.”

What do you think that DEC did?

Precision Example: Weather Modeling

Suppose a weather model that relies on three data to describe each point.

1. The temperature in Kelvins. Possible values are 200 K to 350 K.
2. The pressure in millibars. Typical values are around 1000.
3. The percent humidity. Possible ranges are 0.00 through 1.00 (100%).

Consider the errors associated with single precision floating point arithmetic. The values are precise to 1 part in 10^7 .

The maximum temperature errors resulting at any step in the calculation would be:

$3.5 \cdot 10^{-5}$ Kelvins

$1.0 \cdot 10^{-4}$ millibars

$1.0 \cdot 10^{-5}$ percent in humidity.

As cumulative errors tend to cancel each other out, it is hard to imagine the cumulated error in the computation of any of these values as becoming significant to the result.

Example: A weather prediction accurate to ± 1 Kelvin is considered excellent.

Encoding Decimal Digits

There are ten decimal digits. As $2^3 < 10 \leq 2^4$, we require four binary bits in order to represent each of the digits.

Here are the standard encodings. Note the resemblance to hexadecimal, except that the non-decimal hexadecimal digits are not shown.

0	0000	5	0101
1	0001	6	0110
2	0010	7	0111
3	0011	8	1000
4	0100	9	1001

Note that the binary codes 1010, 1011, 1100, 1101, 1110, 1111 for hexadecimal digits A B C D E F are not used to represent decimal digits.

Packed decimal representation will have other uses for these binary codes.

Packed Decimal Data

Packed decimal representation makes use of the fact that only four binary bits are required to represent any decimal digit.

Numbers can be **packed**, two digits to a byte.

How do we represent signed numbers in this representation?

The answer is that we must include a sign “half byte”.

The IBM format for packed decimal allows for an arbitrary number of digits in each number stored. The range is from 1 to 31, inclusive.

After adding the sign “half byte”, the number of hexadecimal digits used to represent any number ranges from 2 through 32.

Numbers with an odd digit count are converted to proper format by the addition of a leading zero; “1234” becomes “01234”, and “1” becomes “01”.

The system must allocate an integral number of bytes to store each datum, so each number stored in Packed Decimal format must have an odd number of digits.

Packed Decimal: The Sign “Half Byte”

In the S/370 Packed Decimal format, the sign is stored “to the right” of the string of digits as the least significant hexadecimal digit.

The standard calls for use of all six non–decimal hexadecimal digits as sign digits. The standard is as follows:

Binary	Hex	Sign	Comments
1010	A	+	
1011	B	–	
1100	C	+	The standard plus sign
1101	D	–	The standard minus sign
1110	E	+	
1111	F	+	A common plus sign, resulting from a shortcut in translation from Zoned format.

Packed Decimal Data

The packed decimal format is the one preferred by business for financial use. Zoned decimal format is not used for arithmetic, but just for conversions.

As is suggested by the name, the packed format is more compact.

Zoned format one digit per byte

Packed format two digits per byte (mostly)

In the packed format, the rightmost byte stores the sign in its rightmost part, so the rightmost byte of packed format data contains only one digit.

All other bytes in the packed format contain two digits, each with value in 0 – 9. This implies that each packed constants always has an odd number of digits.

A leading 0 may be inserted, as needed.

The standard sign fields are: negative X'D'
 non-negative X'C'

The length may be from 1 to 16 bytes, or 1 to 31 decimal digits.

Examples +7 | 7C |
 - 13 | 01 | 3D |

Example: Addition of Two Packed Decimal Values

Consider two integer values, stored in packed decimal format.

Note that 32-bit two's complement representation would limit the integer to a bit more than nine digits: $-2,147,483,648$ through $2,147,483,647$.

Integers stored as packed decimals can have 31 digits and vary between $-9,999,999,999,999,999,999,999,999,999,999,999$ and $+9,999,999,999,999,999,999,999,999,999,999,999$ ($9.99 \cdot 10^{30}$)

Consider the addition of the numbers -97 and $12,541$.

-97 would be expanded to -097 and stored as **097D**.

$12,541$ would be stored as **12541C**.

The CPU does what we would do. It first pads the smaller number with 0's.

00097D

12541C

It then performs the operation, denoted as " $12541 - 97$ ", and stores the result as 12444 , or **12444C** in packed decimal format.

Fixed Point: Where is the Decimal Point?

Consider the following addition problem: $1.23 + 10.11405$.

The answer is obviously 11.34405.

But the Packed Decimal format does not store the decimal point, so

1.23 is stored as **123C** and 10.11405 is stored as **1011405C**.

Using our previous logic, we line up the numbers

	0000123C
	1011405C

The result is denoted **10111528C**, which might be 10.111528.

It cannot represent the correct answer.

This is a “feature” of **fixed point arithmetic**, in which all numbers are stored with the decimal point in the same place. In this system, the code must guarantee that 1.23 is stored as 0123000C.

The result is now

0123000C
1011405C , stored as 1134405C.

The code must interpret and output this as the value 11.34405.