

Remaining Data Types

CT Team: A. Gieriet, J. Gruber, R. Gübeli, M. Meli, M. Rosenthal, A. Rüst, J. Scheier, M. Thaler

■ Real numbers

- Single vs. double precision

■ Arrays

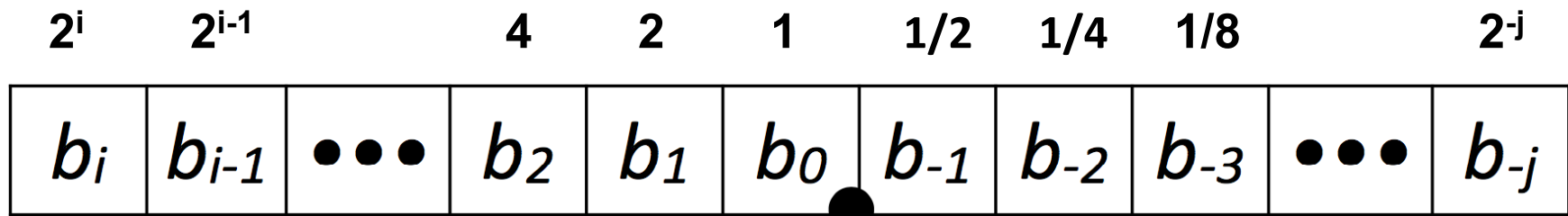
- One dimensional
- Pointer and Reference
- Multi dimensional

■ Strings

■ Structs

- At the end of this lesson you will be able
 - to explain the structure and notation of real numbers in decimal and binary
 - to translate real numbers between decimal and binary
 - to explain how arrays are stored and how to access to single elements
 - to discuss how structs and strings are stored

■ Binary representation



■ Value

$$\sum_{k=-j}^i b_k \times 2^k$$

■ Examples

- $5 \frac{3}{4}$ 101.11_2
- $2 \frac{7}{8}$ 10.111_2

■ Real numbers: general notation

$$\text{Value} = \text{Sign} \cdot \text{Fraction} \cdot \text{Base}^{\text{Exponent}}$$

■ Normalized scientific notation

- Single non-zero digit to the left of the decimal (binary) point

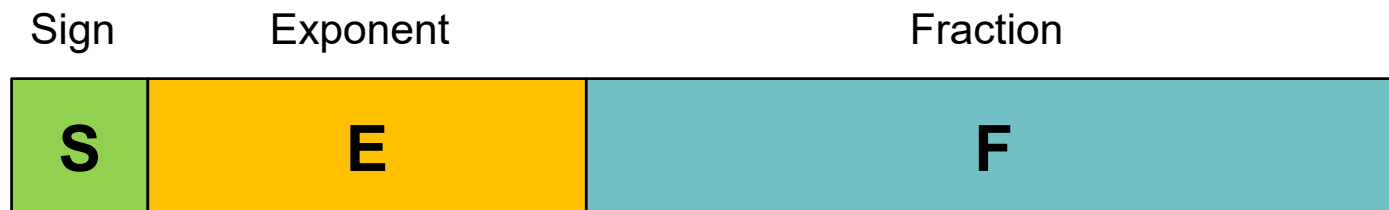
■ Decimal vs. binary

- Decimal: $0.00002796 = 2.796 \cdot 10^{-5}$
- Binary: $0.000010110111 = 1.0110111 \cdot 2^{-5}$

■ Defines how floating point numbers are represented

- Easy exchange of data between machines
- Simplifies hardware algorithms
- Exists since 1980

■ Binary notation



■ Fraction



- Binary notation
- Representing normalized numbers \rightarrow number of form **1**.xxxx..
- In IEEE 754 standard, the **1** is implicit

$$\text{Fraction value} = (1 + F)$$

- Example

$$1.265625 \text{ d} = 1.010001 \text{ b} \rightarrow F = 010001$$

Remark: $1.010001 \text{ b} = (1 + 0 \times 2^{-1} + 1 \times 2^{-2} + \dots + 1 \times 2^{-6}) \text{ d}$

■ Exponent



- Binary excess notation

$$\text{Exponent value} = (E - \text{Bias})$$

- Bias: 127 (float) / 1023 (double)
- Example (float):

$$0111'1010b = 122 \rightarrow 122 - 127 = -5 \rightarrow \text{Value} = 2^{-5}$$

$$\text{Value} = 2^{-5} \rightarrow -5 + 127 = 122 = 0111'1010b$$

■ Exponent Single precision (float)



E		Meaning
0000 0000		Reserved
0000 0001		-126_{10}
0000 0010		-125_{10}
0111 1111		0_{10}
1111 1110		127_{10}
1111 1111		Reserved

→ Unnormalized
→ See later...

→ Not a Number
→ See later...

■ Sign



- 1 bit
- Sign value = $(-1)^S$

■ Sign and Magnitude Representation



$$\text{Value} = (-1)^S \cdot (1 + F) \cdot 2^{(E - \text{Bias})}$$

- More exponent bits → wider range of numbers
- More fraction bits → higher precision

Exercise: Single Precision (float)

- **Decimal 7.5 d = ?**

Exercise: Single Precision (float)

■ Binary 10111111'01010000'00000000'00000000 = ?

Scalar Types: Single Precision

■ Special Values

E	F	S	Result	
255	$\neq 0$		Not a Number (NaN)	
255	0	1	$-\infty$	
255	0	0	∞	
$0 < E < 255$			$(-1)^S \cdot (1 + F) \cdot 2^{(E-Bias)}$	Normalized
0	$\neq 0$		$(-1)^S \cdot (0 + F) \cdot 2^{(1-Bias)}$	Unnormalized
0	0	1	-0	
0	0	0	$+0$	

Ranges (Single Precision)

■ **Normalized (positive):** $(1 + F) \cdot 2^{(E-Bias)}$

- Smallest

0 00000001 000000000000000000000000

 $+2^{-126}(1+0) = 2^{-126}$
- Largest

0 11111110 111111111111111111111111

 $+2^{127}(1+(1-2^{-23}))$

■ **Unnormalized (positive):** $(0 + F) \cdot 2^{(1-Bias)}$

- Smallest

0 00000000 000000000000000000000001

 $+2^{-126}(2^{-23}) = 2^{-149}$
- Largest

0 00000000 111111111111111111111111

 $+2^{-126}(1-2^{-23})$



■ Single precision (float)

- 32 bits
- Exponent bias: 127
- Dynamics: $-3.4 \cdot 10^{38} \dots -1.17 \cdot 10^{-38}$, 0, $1.17 \cdot 10^{-38} \dots 3.4 \cdot 10^{38}$
- Resolution: $2^{-24} = 0.6 \cdot 10^{-7} \rightarrow 7$ digits



■ Double precision (double)

- 64 bits
- Exponent bias: 1023
- Dynamics: $-1.8 \cdot 10^{308} \dots -2.2 \cdot 10^{-308}$, 0, $2.2 \cdot 10^{-308} \dots 1.8 \cdot 10^{308}$
- Resolution: $2^{-53} = 0.11 \cdot 10^{-15} \rightarrow 15$ digits



■ Precision:

```
int main(void) {  
    float a = 1.8;  
    float b = 18;  
    if(a*10 == b) printf("\nNumbers are the same\n\n");  
    else          printf("\nNumbers differ!!\n\n");  
}
```

- Output: Numbers are the same

```
int main(void) {  
    float a = 1.81;  
    float b = 18.1;  
    if(a*10 == b) printf("\nNumbers are the same\n\n");  
    else          printf("\nNumbers differ!!\n\n");  
}
```

- Output Numbers differ!!

■ Precision:

```
int main(void) {  
    int x    = 33554431;  
    float y = 33554431;  
    printf("int: %d\n", x );  
    printf("float %f\n", y );  
}
```


- Output:

```
int: 33554431  
float: 33554432.000000
```

- Single Precision float of 33554431:

Binary32: 4C000000

Status	Sign [1]	Exponent [8]	Significand [23]
Normal	0 (+)	10011000 (+25)	1.00000000000000000000000 (1.0)


$$2^{25} = 33554432$$

<https://babbage.cs.qc.cuny.edu/IEEE-754>

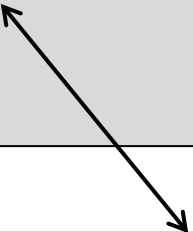
■ Patriot Missile Error

- During the Gulf War in 1991, a U.S. Patriot missile failed to intercept an Iraqi Scud missile, and 28 Americans were killed.
- The problem was caused by the inaccuracy of the binary representation of 0.10.
 - The Patriot incremented a counter once every 0.10 seconds.
 - It multiplied the counter value by 0.10 to compute the actual time.
- However, the (24-bit) binary representation of 0.10 actually corresponds to 0.099999904632568359375, which is off by 0.000000095367431640625.
- The system was running 100 hours, the time ended up being off by 0.34 seconds - enough time for a Scud to travel 500 meters!



■ Pointer to first element of array

```
int main(void) {  
    int j, sum = 0;  
    int array[5];  
    int *ptr;                                // pointer to integer  
    for (j = 0; j < 5; j++) {  
        array[j] = j;  
    }  
    ptr = array;                             // copy pointer  
    for (j = 0; j < 5; j++) {  
        sum = sum + *ptr;                    // add value ptr points  
        ptr++;                               // point to next element  
    }  
}
```



```
for (j = 0; j < 5; j++) {  
    sum = sum + array[j];  
}
```

■ Multidimensional arrays are "Arrays of Arrays"

```
void main(void) {
    int array [3][3];
    int *ptr;
    int i, j, k;
    j = 0;
    for (i = 0; i < 3; i++) {
        for (k = 0; k < 3; k++) {
            array [i][k] = j++;
        }
    }
    ptr = array [0];
    for (i = 0; i < 9; i++) {
        printf("index: %d, value: %d\n", i, *ptr);
        ptr++;
    }
}
```

- Array a[3][3]:
$$\begin{bmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{11} & a_{12} \\ a_{20} & a_{21} & a_{22} \end{bmatrix}$$
- Is stored as: $[a_{00} \ a_{01} \ a_{02} \ a_{10} \ a_{11} \ a_{12} \ a_{20} \ a_{21} \ a_{22}]$

Arrays: Who Cares?

■ Access to arrays: g++, Linux

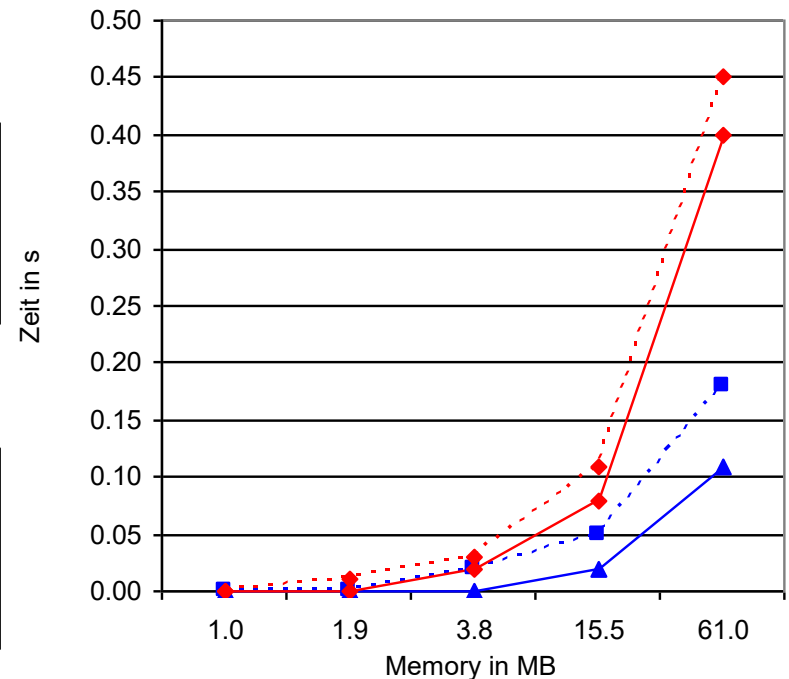
- long int array [N][N]
- N = 500, 700, 1000, 2000, 4000, 8000, 12000, 16000, 20000
M = 1, 1.9, 3.8, 15.3, 61.0, 244, 550, 980, 1500 MB

- Blue: CPU-time

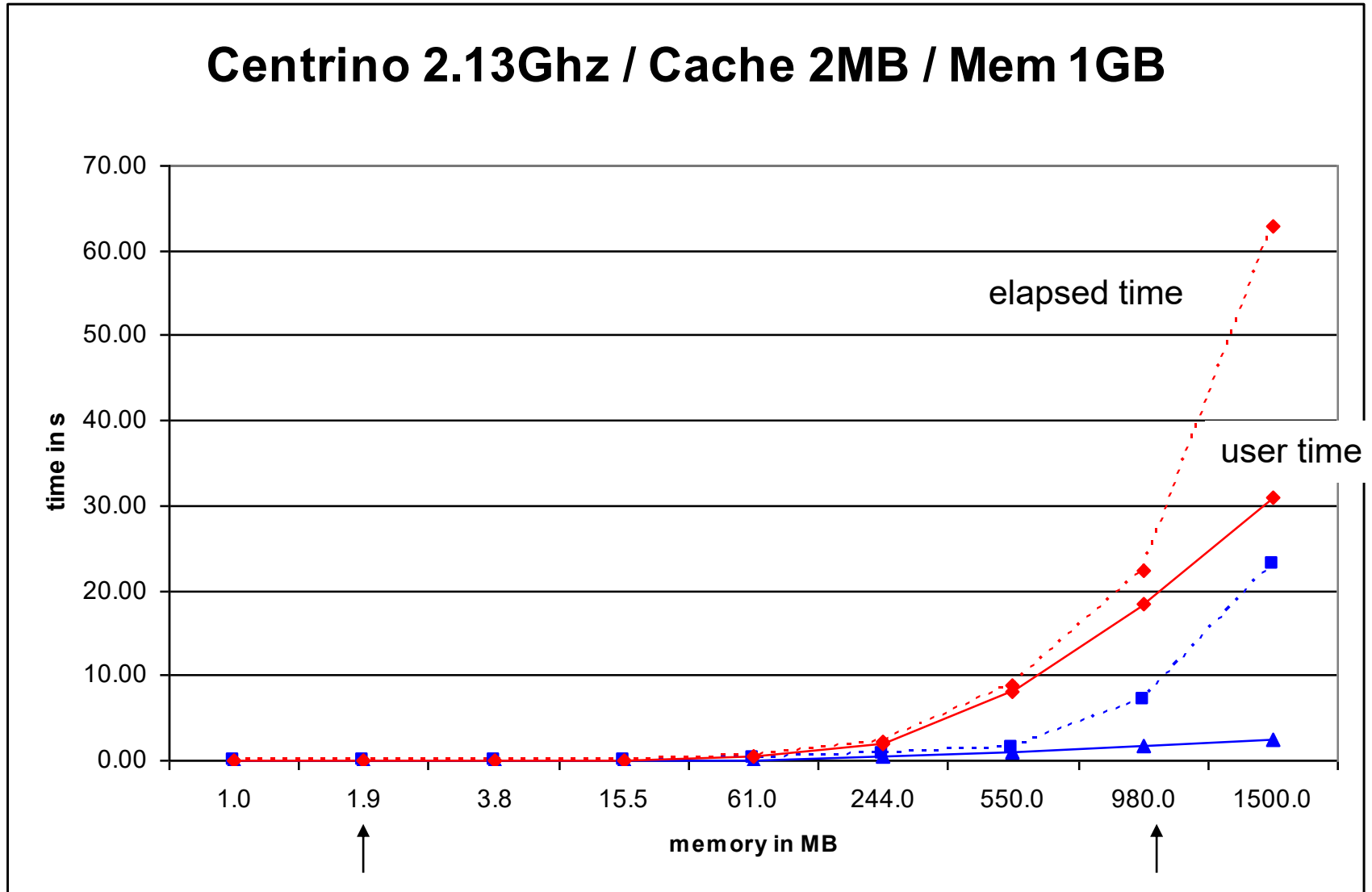
```
for (i = 0; i < s; i++) {  
    for (j = 0; j < s; j++) {  
        array[i][j]=i+j;  
    }  
}
```

- Red: CPU-time

```
for (i = 0; i < s; i++) {  
    for (j = 0; j < s; j++) {  
        array[j][i]=i+j;  
    }  
}
```



... arrays



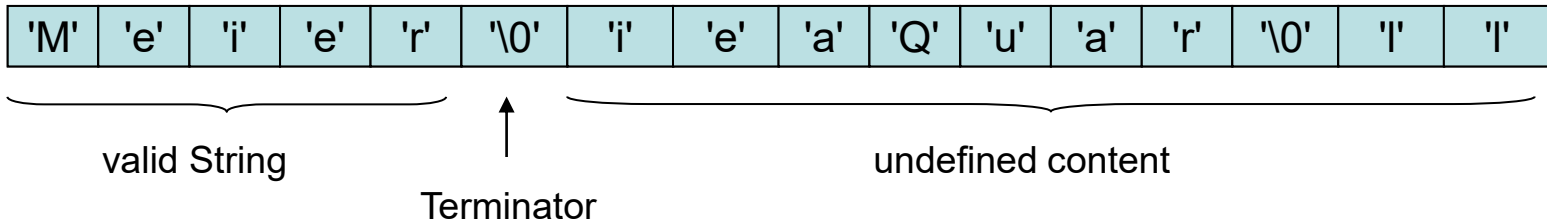
Strings in C

■ Strings in C are not objects

- Declaration with char arrays
- Example

```
char Name[16] = "Meier";
```

- Question: How is the end of a string recognised?
- All Strings in C are 0 – terminated
→ after the last character follows a '0'



- Length of strings not stored → count characters to '\0'

■ Structs contain different data that belong together

- Declaration

```
struct person {  
    char    name[8];  
    char    prename[8];  
    int     pers_nr;  
    unsigned char day_of_birth;  
    unsigned char month_of_birth;  
};
```

- Definition

```
struct person aperson;
```

- Usage

```
strcpy(aperson.Name, "Meier");  
strcpy(aperson.prename, "Peter");
```

Struct Memory Layout

■ Access to struct

```
struct person {  
    char    name[8];  
    char    prename[8];  
    int     pers_nr;  
    unsigned char day_of_birth;  
    unsigned char month_of_birth;  
};
```

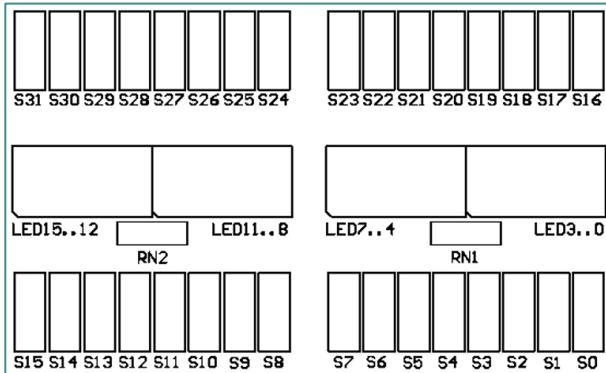
- Using base address

- name: base address + 0
- prename: base address + 8
- pers_nr: base address + 16
- day_of_birth: base address + 18
- month_of_birth: base address + 19

sizeof(struct person) = 20

Identifier	Memory Content	Address	
name	M	0x0000'2000	← Base Address
	e	0x0000'2001	
	i	0x0000'2002	
	e	0x0000'2003	
	r	0x0000'2004	
	\0	0x0000'2005	
prename		0x0000'2006	
		0x0000'2007	
	P	0x0000'2008	
	e	0x0000'2009	
	t	0x0000'200A	
	e	0x0000'200B	
pers_ne	r	0x0000'200C	
	\0	0x0000'200D	
		0x0000'200E	
		0x0000'200F	
	156	0x0000'2010	
	034	0x0000'2011	
day_of_birth	15	0x0000'2012	
month_of_birth	8	0x0000'2013	

Union Memory Layout



Memory Map

S 7	S 6	S 5	S 4	S 3	S 2	S 1	S 0	0x6000'0200
S 15	S 14	S 13	S 12	S 11	S 10	S 9	S 8	0x6000'0201
S 23	S 22	S 21	S 20	S 19	S 18	S 17	S 16	0x6000'0202
S 31	S 30	S 29	S 28	S 27	S 26	S 25	S 24	0x6000'0203

```
typedef union {
    struct {
        volatile uint8_t S7_0;
        volatile uint8_t S15_8;
        volatile uint8_t S23_16;
        volatile uint8_t S31_24;
    } BYTE;
    struct {
        volatile uint16_t S15_0;
        volatile uint16_t S31_16;
    } HWORD;
    volatile uint32_t WORD;
} reg_ct_dipsw_t;

#define CT_DIPSW ((reg_ct_dipsw_t *) 0x60000200)
...
sizeof (reg_ct_dipsw_t) = 4
```

WORD	HWORD	BYTE	Address
	s15_0	s7_0	0x6000'0200
		s15_8	0x6000'0201
	s31_16	s23_16	0x6000'0202
		s31_24	0x6000'0203

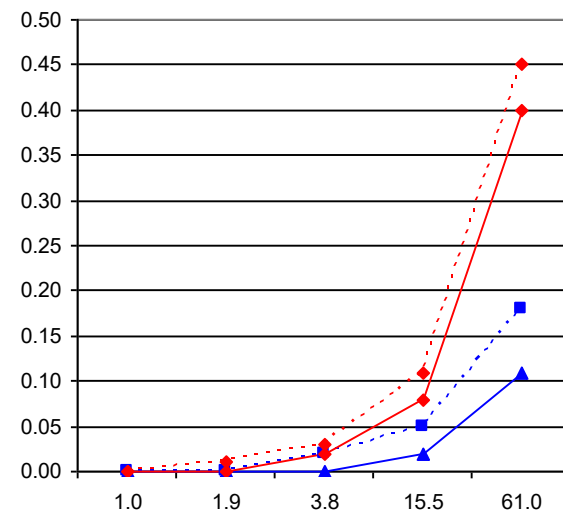
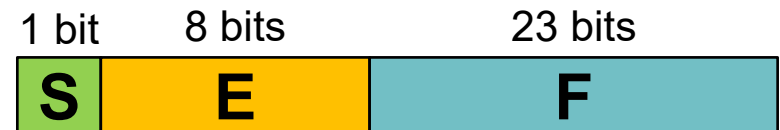
Conclusion

■ Real numbers

- IEEE Standard 754 / 854
- Single vs. double precision

■ Memory Layout of

- Arrays
- Strings
- Structs



'M'	'e'	'i'	'e'	'r'	'\0'	'i'	'e'	'a'	'Q'	'u'	'a'	'r'	'\0'	'l'	'l'
-----	-----	-----	-----	-----	------	-----	-----	-----	-----	-----	-----	-----	------	-----	-----