# ICM – Computer Science Major M1 Cyber Physical and Social Systems Course unit on Data Interoperability and Semantics

#### **Tutorial 1**

#### 1) Bluetooth Mesh Model

The definitions below are extracts of a Bluetooth standard. Complete the tables.

Bluetooth Mesh Model specification v1.0.1 Section 3.1.6.1 Generic Battery Level (modified)

The Generic Battery Level state is a value ranging from 0 percent through 100 percent. The values must be divided by two (ex., 0x03 stands for 1.5 %). Values above 0xC8 (100 %) are forbidden.

	encoded	decoded
8-bit unsigned integer (UINT8)	0xA7 (hexa)	
UINT8	0b01010110 (binary)	

#### Bluetooth Mesh Model specification v1.0.1 Section 3.1.7.6 Local Altitude

The Local Altitude field determines the altitude of the device relative to the Generic Location Global Altitude. This is a 16-bit signed integer in decimeters. The valid range is from -32768 decimeters (0x8000) through 0 decimeters (0x0000) up to 32765 decimeters (0x7FFD).

The following formula can be used to decode data:  $B2T(X) = -x_{w-1} \cdot 2^{w-1} + \sum_{i=0}^{w-2} x_i \cdot 2^i$ 

	encoded	decoded
INT16 - Big Endian (AB)	0x8001	
INT16 - Little Endian (BA)	0x8001	

#### 2) IEEE 754

Decode the following hexadecimal strings using the Single precision IEEE 754 floating-point standard:

- a) 0x3F000000
- b) 0x44800000

Encode the following numbers using the Single precision IEEE 754 floating-point standard:

- c) -3.0
- d) 24.0

## 3) UTF-8

Convert the following strings to UTF-8:

- a) cs
- b) ¥
- c) €
- d) 😩

#### 4) CBOR

Below is a document encoded using the Concise Binary Object Representation.

Write a JSON equivalent of this document.

Spaces and new lines in the representation are just there to help you.

82 A3 61 76 FA 3F800000 62 6363 63 455552 62 6373 63 E282AC

A3 61 76 FA 40F80000 62 6363 63 434E59 62 6373 62 C2A5

#### 5) Base32 decoding

Below is the RFC 4648 Base32 Alphabet

Value	Encoding	Value	Encoding	Value	Encoding	Value	Encoding
0	A	9	J	18	S	27	3
1	В	10	K	19	T	28	4
2	C	11	L	20	U	29	5
3	D	12	M	21	V	30	6
4	E	13	N	22	W	31	7
5	F	14	O	23	X		
6	G	15	P	24	Y	(pad)	=
7	H	16	Q	25	Z		
8	I	17	R	26	2		

Decode the following string: JVSWK5BAJJXWKICAGRYG2II=

Sequence of values:

Sequence of 5 bits:

Decoded string:

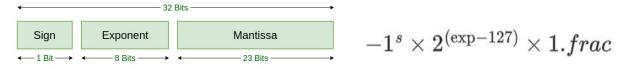
#### 6) CRLF

Different OSs represent ends of line differently.

Microsoft Windows uses CRLF, Unix uses only LF, Classic MacOS used CR, then LF since Mac OS X.

- a) What is the full name of CR and LF?
- b) What are their ASCII codes?
- c) With what escape characters can these characters be included in a JSON string?
- d) Why is it problematic that different OSs use different end-of-line characters?

#### Appendix A: Single precision IEEE 754 floating-point standard



As an example, the value for number 0xbfc00000 is -1.5

#### Appendix B: UTF-8 Encoding

UTF-8 encodes code points in one to four bytes, depending on the value of the code point. In the following table, the x characters are replaced by the bits of the code point:

Code point $\leftrightarrow$ UTF-8 conversion					
First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4
U+0000	U+007F	0xxxxxxx			
U+0080	U+07FF	110xxxxx	10xxxxxx		
U+0800	U+FFFF	1110xxxx	10xxxxxx	10xxxxxx	
U+10000	U+10FFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

The first 128 code points (ASCII) need one byte. The next 1,920 code points need two bytes to encode, which covers the remainder of almost all Latin-script alphabets, and also IPA extensions, Greek, Cyrillic, Coptic, Armenian, Hebrew, Arabic, Syriac, Thaana and N'Ko alphabets, as well as Combining Diacritical Marks. Three bytes are needed for the remaining 61,440 code points of the Basic Multilingual Plane (BMP), including most Chinese, Japanese and Korean characters. Four bytes are needed for the 1,048,576 code points in the other planes of Unicode, which include emoji (pictographic symbols), less common CJK characters, various historic scripts, and mathematical symbols.

A "character" can take more than 4 bytes because it is made of more than one code point. For instance a national flag character takes 8 bytes since it is "constructed from a pair of Unicode scalar values" both from outside the BMP.

### **Encoding process**

In these examples, character \( '\) indicate how bits from the code point are separated and distributed among the UTF-8 bytes. Additional bits added by the UTF-8 encoding process are underlined.

- 1. The Unicode code point for the euro sign € is U+20AC.
- 2. As this code point lies between U+0800 and U+FFFF, this will take three bytes to encode.
- 3. Hexadecimal 20AC is binary 0010 | 0000 10 | 100. The two leading zeros are added because a three-byte encoding needs exactly sixteen bits from the code point.
- 4. Because the encoding will be three bytes long, its leading byte starts with three 1s, then a 0 (1110...)
- 5. The four most significant bits of the code point are stored in the remaining low order four bits of this byte (11100010), leaving 12 bits of the code point yet to be encoded (...0000 10|10 1100).
- 6. All continuation bytes contain exactly six bits from the code point. So the next six bits of the code point are stored in the low order six bits of the next byte, and 10 is stored in the high order two bits to mark it as a continuation byte (so 1000010).
- 7. Finally the last six bits of the code point are stored in the low order six bits of the final byte, and again 10 is stored in the high order two bits (10101100).

The three bytes  $\underline{1110}0010\ \underline{10}000010\ \underline{10}101100$  can be more concisely written in hexadecimal, as E2 82 AC.

# Appendix C: Some Unicode code points

	000	001	002	003	004	005	006	007
0	NUL 8800	<b>DLE</b>	SP 0020	0030	@ 0040	P 0050	0060	<b>p</b>
1	SOH 0001	DC1	0021	1	A 0041	Q 0051	<b>a</b>	<b>q</b>
2	STX 0002	DC2	0022	2	B 0042	<b>R</b>	<b>b</b>	<b>r</b>
3	ETX 0003	DC3	# 0023	3	C 0043	S 0053	C 0063	S 0073
4	EOT 0004	DC4	\$	4	D 0044	T	d	t
5	ENQ 0005	NAK 0015	% 0025	5	E 0045	U 0055	e 0065	u 0075
6	ACK 0006	SYN 0016	& 0026	6	F 0046	V 0056	<b>f</b>	<b>V</b>
7	BEL 0007	ETB 0017	0027	7	<b>G</b>	W 0057	g 0067	W 0077
8	BS 0008	CAN 0018	0028	8	H 0048	X 0058	h	X 0078
9	HT 0009	E M 0019	0029	9	I 0049	Y 0059	i 0069	<b>y</b>
Α	LF 000A	SUB	* 002A	003A	<b>J</b>	Z 005A	<b>j</b>	<b>Z</b>
В	VT 0008	ESC 001B	+ 002B	• • • • •	K 0048		k	{ 007B
С	FF 000C	FS 001C	• 002C	< 003C	L	005C	1	007C
D	CR 0000	GS 001D		003D	M 004D	0050	m 0060	} 007D
Е	80 000E	RS 001E	• 002E	> 003E	N 004E	<b>∧</b>	n 006E	~ 007E
F	SI	US 001F	/ 002F	?	O 004F	005F	O 006F	DEL 007F

Basic Latin Range: U+0000 – U+007F

	008	009	00A	00B	00C	00D	00E	00F
0	0080	DCS 0090	NB SP	0080	À	$\mathbf{\tilde{D}}_{0000}$	à	ð
1	0081	PU1 0091	00A1	<u>+</u>	Á	$ ilde{\mathbf{N}}_{_{0001}}$	á	ñ
2	BPH 0082	PU2 0092	<b>¢</b>	2	Â	Ò	â	Ò 00F2
3	NBH 0083	STS 0093	£	3	Ã	Ó	ã	<b>Ó</b>
4	IND 0084	CCH 0094	<b>Q</b>	0084	Ä	Ô	ä	ô
5	NEL 0085	MW 0095	¥	μ 0085	Å	Õ	å	Õ 00F5
6	SSA 0086	SPA 0096	00046	¶ 0086	Æ	Ö	æ	Ö 00F6
7	ESA 0087	<b>EPA</b>	<b>§</b>	0087	Ç	0007	Ç 00E7	00F7
8	HTS	SOS 0098	00A8	5 0088	È	Ø 0008	è	Ø 00F8
9	HTJ 0089	XXX 0099	6V8	1 0089	É	Ù	é	ù
Α	VTS 008A	SCI 009A	<u>a</u>	<b>Q</b> 008A	Ê	Ú	ê	ú ODFA
В	PLD 008B	CS1 0098	<b>≪</b>	>>> 0088	Ë	Û	ë	û
С	PLU	ST 009C		1/4	Ì	Ü	ì	ü
D	RI 0080	OSC 009D	SHY	1/2	Í	Ý	í	ý
Ε	SS2 008E	PM 009E	® ODAE	3/4 008E	Î	Þ	î	þ
F	SS3	APC 009F	00AF	¿ OOBF	Ï 00CF	ß	i OOEF	ÿ

Latin-1 Supplement: U+0080 – U+00FF

	20A	20B	20C
0	<b>E</b>	<b>S</b>	<u>C</u>
1	<b>₡</b>	<b>₽</b>	
2	<b>C</b>	<b>G</b>	
3	<b>F</b> 20A3	<b>A</b> 2083	
4	€	8	
5	<u>20A4</u>	<u>2084</u>	
6	N N	#t	
7	Pts	S 1086	
8	Rs	20B7	
9	₩ 	₹	
Α	20.49	<b>も</b>	
В	<u>₫</u>	20BA	
С	€	20BB	
D	K	₽	
E	T T	<u>208D</u>	
F	20AE 20AF	20BE <b>B</b> 20BF	

Currency Symbols U+20A0 – U+20C0

	1F60	1F61	1F62	1F63	1F64
0	<b>(#)</b>	•••	<b>(*)</b>	<b>③</b>	
1	1F600	1F610	1F620	1F630	1F640
2	1F601	1F611	1F621	1F631	1F641
_	1F602	1F612	1F622	1F632	1F642
3	1F603	1F613	1F623	1F633	1F643
4	1F604	1F614	1F624	1F634	1F644
5		<b>©</b>	<b>(2)</b>	66	
	1F605	1F615	1F625	1F635	1F645
6	1F606	1F616	1F626	1F636	1F646
7	1F607	1F617	1F627	1F637	1F647
8	1F608	1F618	1F628	) 1F638	1F648
9	<b>(3)</b>	(C)	8		<b>®</b>
Α	1F609	1F619	1F629	1F639	1F649
В	1F60A	1F61A	1F62A	1F63A	1F64A
	1F60B	1F61B	1F62B	1F63B	1F64B
С	<u></u>	<b>®</b>	<b></b>	<b>(1)</b>	
D	1F60C	1F61C	1F62C	1F63C	1F64C
E	1F60D	1F61D	1F62D	1F63D	<u>Ö</u> ,
F	1F60E	1F61E	1F62E	1F63E	1F64E
	1F60F	1F61F	1F62F	1F63F	1F64F

Emoji symbols U+1F600 – U+1F64F

#### Appendix D: RFC 8949 - Concise Binary Object Representation (CBOR)

Concise Binary Object Representation (CBOR) is a binary data serialization format loosely based on JSON authored by C. Bormann. Like JSON it allows the transmission of data objects that contain name—value pairs, but in a more concise manner. This increases processing and transfer speeds at the cost of human readability. It is defined in IETF RFC 8949

#### Specification of the CBOR encoding

CBOR encoded data is seen as a stream of data items. Each data item consists of a header byte containing a 3-bit type and 5-bit short count. This is followed by an optional extended count (if the short count is in the range 24–27), and an optional payload.

For types 0, 1, and 7, there is no payload; the count is the value. For types 2 (byte string) and 3 (text string), the count is the length of the payload. For types 4 (array) and 5 (map), the count is the number of items (pairs) in the payload.

#### **CBOR** cheatsheet

Major type	Description	Binary Shorth	nand
0	an unsigned integer	000_xxxxx	unsigned(#)
1	a negative integer	001_xxxxx	negative(#-1)
2	a byte string	010_xxxxx	bytes(n)
3	a text string (UTF-8)	011_xxxxx	text(n)
4	an array of data items	100_xxxxx	array(n)
5	a map of pairs of data items	101_xxxxx	map(n)
7	Floating-point numbers and values with no content	111_xxxxx	simple(v)

<sup>&</sup>lt;sup>1</sup> For negative(#-1), # represents the negative value minus 1. For example, negative(4) represents a value of -5.

19 0100 .. 19 FFFF

Unsigned intege	er	256 65535	B9 0100 B9 FFFF
Number	Hex	Indefinite	BF
0 15	00 0F	Length represents r	number of key/value pairs.
16 23	10 17		
24 255	18 18 18 FF	Floating point nun	nbers and values with no

#### **Negative integer**

256 .. 65535

Number	Hex
-116	20 2F
-1724	30 37
-25256	38 18 38 FF
-25765536	39 0100 39 FFFF

#### **Array**

Length	Hex
0 15	80 8F
16 23	90 97
24 255	98 18 98 FF
256 65535	99 0100 99 FFFF
Indefinite	9F

Length represents number of key/value pairs.

#### Мар

Length	Hex
0 15	A0 AF
16 23	B0 B7
24 255	B8 18 B8 FF

Floating point numbers and values with no content 5-Bit Value Semantics				
20	false			
21	true			
22	null			
23	undefined			
25	IEEE 754 Half-Precision Float (16 bits follow)			
26	IEEE 754 Single-Precision Float (32 bits			
follow)				
27	IEEE 754 Double-Precision Float (64 bits			
follow) 31	"break" stop code for indefinite-length items			
	, , , , , , , , , , , , , , , , , , , ,			

#### **Examples:**

As JSON:	As CBOR:	
{ "t":2, "h":null }	61 74 02 61	# map(2) # text(1) 1 # "t" # unsigned(2) # text(1) 3 # "h" # primitive(22)

<sup>&</sup>lt;sup>5</sup> For map(n), n represents number of key/value pairs.