

# Day-14

## Cheatsheet

### 부동소수점 - Cheatsheet

▪ NaN (mantissa $\neq 0$ )	*	11111111	*****
▪ $\pm$ infinity	*	11111111	000000000000000000000000
▪ Lowest/Largest ( $\pm 3.40282 * 10^{+38}$ )	*	11111110	111111111111111111111111
▪ Minimum (normal) ( $\pm 1.17549 * 10^{-38}$ )	*	00000001	000000000000000000000000
▪ Denormal number ( $< 2^{-126}$ )(minimum: $1.4 * 10^{-45}$ )	*	00000000	*****
▪ $\pm 0$	*	00000000	000000000000000000000000

	E4M3	E5M2	half
<b>Exponent</b>	4 [0*-14] (no inf)	5-bit [0*-30]	
<b>Bias</b>	7	15	
<b>Mantissa</b>	4-bit	2-bit	10-bit
<b>Largest (±)</b>	$1.75 \times 2^8$ 448	$1.75 \times 2^{15}$ 57,344	$2^{16}$ 65,536
<b>Smallest (±)</b>	$2^{-6}$ 0.015625	$2^{-14}$ 0.00006	
<b>Smallest (denormal*)</b>	$2^{-9}$ 0.001953125	$2^{-16}$ $1.5258 \times 10^{-5}$	$2^{-24}$ $6.0 \times 10^{-8}$
<b>Epsilon</b>	$2^{-4}$ 0.0625	$2^{-2}$ 0.25	$2^{-10}$ 0.00098

	bf16	float	double
<b>Exponent</b>	8-bit [0*-254]	11-bit [0*-2046]	
<b>Bias</b>	127	1023	
<b>Mantissa</b>	7-bit	23-bit	52-bit
<b>Largest (±)</b>	$2^{128}$ $3.4 \times 10^{38}$	$2^{1024}$ $1.8 \times 10^{308}$	
<b>Smallest (±)</b>	$2^{-126}$ $1.2 \times 10^{-38}$	$2^{-1022}$ $2.2 \times 10^{-308}$	
<b>Smallest (denormal*)</b>	/	$2^{-149}$ $1.4 \times 10^{-45}$	$2^{-1074}$ $4.9 \times 10^{-324}$
<b>Epsilon</b>	$2^{-7}$ 0.0078	$2^{-23}$ $1.2 \times 10^{-7}$	$2^{-52}$ $2.2 \times 10^{-16}$

## 부동소수점 - Limits

```
#include<limits>
// T: float or double

std::numeric_limits<T>::max(); // 최대값
std::numeric_limits<T>::lowest(); // 최솟값 (C++11)
```

```

std::numeric_limits<T>::min(); // 가장 작은 값
std::numeric_limits<T>::denorm_min(); // 가장 작은(비정규) 값
std::numeric_limits<T>::epsilon(); // 엡실론 값
std::numeric_limits<T>::infinity(); // infinity
std::numeric_limits<T>::quiet_NaN(); // NaN(Not a Number)

```

## 부동 소수점 - 유용한 함수들

```

#include<cmath> // C++11
using namespace std;
// T : float or double

bool isnan(T value); // value가 NaN인지 확인
bool isinf(T value); // value가 +-infinity인지 확인
bool isfinite(T value); // value가 NaN or +-infinity가 아닌지 확인

bool isnormal(T value); // value가 정규수인지 확인

T ldexp(T x, p); // 지수 shift  $x \times 2^p$ 
int ilogb(T value); // value의 지수를 출력

```