

https://github.com/NVIDIA/trt-samples-for-hackathon-cn/tree/master/cookbook/50-Resource

浮点数

数据类型	FP64	FP32	TF32	FP16	BF16	FP8e5m2	FP8e4m3
符号位数	1	1	1	1	1	1	1
指数位数 (\$k\$)	11	8	8	5	8	5	4
尾数位数 (\$n\$)	52	23	10	10	7	2	3

最大值（最小值）

数据类型	FP64	FP32	TF32	FP16	BF16
最大值符号位 (\$s_{\bar{2}}\$)	0	0	0	0	0
指数位 (\$e_{\bar{2}}\$)	11111111110	11111110	11111110	11110	11111110
尾数位 (\$m_{\bar{2}}\$)	111...1	111...1	111...1	111...1	1111111
$E=2^{\{e\}-\left(2^{k-1}\right)}$	$\$1023\$$	$\$127\$$	$\$127\$$	$\$15\$$	$\$127\$$
$M=\sum_{i=1}^n\frac{1}{2^{\{i\}}}$	$\$1-\frac{1}{2^{\{52\}}\$$	$\$1-\frac{1}{2^{\{23\}}\$$	$\$1-\frac{1}{2^{\{10\}}\$$	$\$1-\frac{1}{2^{\{10\}}\$$	$\$1-\frac{1}{2^{\{7\}}\$$
$\max=\left(-1\right)^{\{s\}2^{\{E\}}\left(1+M\right)}$	$\$1.798\times10^{\{308\}}\$$	$\$3.403\times10^{\{38\}}\$$	$\$3.401\times10^{\{38\}}\$$	$\$65504.\$$	$\$3.390\times10^{\{38\}}\$$
最小值符号位 (\$s_{\bar{2}}\$)	1	1	1	1	1
$\min=\left(-1\right)^{\{s\}2^{\{E\}}\left(1+M\right)}$	$\$-1.798\times10^{\{308\}}\$$	$\$-3.403\times10^{\{38\}}\$$	$\$-3.401\times10^{\{38\}}\$$	$\$-65504.\$$	$\$-3.390\times10^{\{38\}}\$$

绝对最小值

数据类型	FP64	FP32	TF32	FP16	BF16
符号位 (\$s_{\bar{2}}\$)	0	0	0	0	0
指数位 (\$e_{\bar{2}}\$)	00000000000	00000000	00000000	00000	00000000
尾数位 (\$m_{\bar{2}}\$)	000...01	000...01	000...01	000...01	0000001
$E=1-\left(2^{k-1}\right)$	$\$-1022\$$	$\$-126\$$	$\$-126\$$	$\$-14\$$	$\$-126\$$
$M=\frac{1}{2^{\{n\}}}$	$\$\frac{1}{2^{\{52\}}\$$	$\$\frac{1}{2^{\{23\}}\$$	$\$\frac{1}{2^{\{10\}}\$$	$\$\frac{1}{2^{\{10\}}\$$	$\$\frac{1}{2^{\{7\}}\$$
$\text{value}=\left(-1\right)^{\{s\}2^{\{E\}}M\$$	$\$4.941\times10^{\{-324\}}\$$	$\$1.401\times10^{\{-45\}}\$$	$\$1.148\times10^{\{-41\}}\$$	$\$5.960\times10^{\{-8\}}\$$	$\$9.184\times10^{\{-41\}}\$$

其他值

数据类型	$\$+\infty\$$	$\$-\infty\$$	NaN
符号位 (\$s_{\bar{2}}\$)	0	1	0 或 1
指数位 (\$e_{\bar{2}}\$)	全 1	全 1	全 0
尾数位 (\$m_{\bar{2}}\$)	全 0	全 0	非全 0