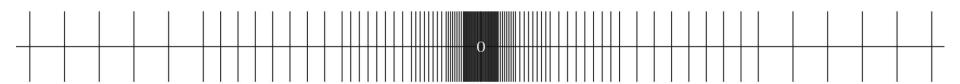
Tools to Detect and Diagnose Floating-Point Errors in Heterogeneous Computing Hardware and Software

A Half-Day Tutorial at Supercomputing 2025

Tut 122: Sunday, 16 November 2025, 8:30am - 12:00pm CST : Location 121









Tools to Detect and Diagnose Floating-Point Errors in Heterogeneous Computing Hardware and Software

A Half-Day Tutorial at Supercomputing 2025

Presenters

Ignacio Laguna





Mark Baranowski and Ganesh Gopalakrishnan





Dolores Miao and Cindy Rubio-González









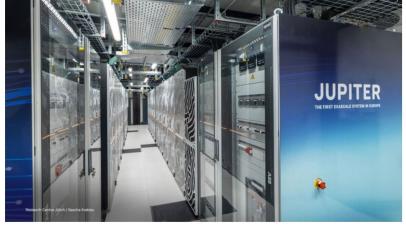
Agenda, and a single QR code to access all our material

- 08:30 : Introduction
- 08:35 : FPChecker
- 09:35 : NixNan (Part-1)
- 09:55 : Plz Fill Survey
- 10:00 : Coffee Break
- 10:30 : Brief Recap
- 10:35 : NixNan (Part-2)
- 11:05 : FloatGuard
- 11:50 : Closing Remarks
- 11:55 : Plz Fill Survey



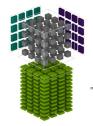
Floating-Point Computations are fundamental to CPUs and GPUs

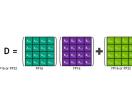






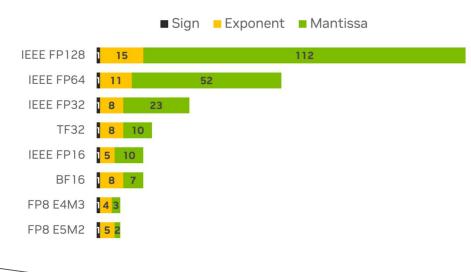






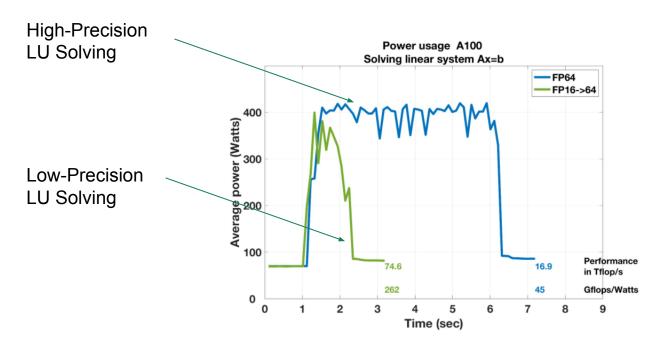
Floating-Point Numbers: A compromise, with many good properties, that can span a large range with a small number of bits Half the representable numbers lie between -1 and +1

Form at	Prec.	Ехр.	Sm. norm.	Lrgst. Norm.	ULP @ emax	eql. 2'c bits	Redn. factor		
ED40	10	+	0.40	0.55	_	44	0.5	IEEE FP128	15
FP16 16 bits	10	5	6.10 E-5	6.55 E4	5	41	2.5	IEEE FP64	11
BF16	7	8	1.18	3.39	120	262	16.3	IEEE FP32	8 23
16 bits			E-38	E38				TF32 i	8 10
TF32 19 bits	10	8	1.18 E-38	3.40 E38	117	265	13.9	IEEE FP16	5 10
19 0115			E-30	E30				BF16	8 7
FP32 32 bits	23	8	1.18 E-38	3.40 E38	104	254	7.9	FP8 E4M3	
FP64 64 bits	52	11	2.23 E-308	1.80 E308	971	2099 bits	32.7	FP8 E5M2	5 2



FP allows us to get away moving 64 bits instead of moving 2099 bits (in 2's comp.)

Lower precision Floating-Point saves energy (and time)



From "Hardware Trends Impacting Floating-Point Computations In Scientific Applications," Dongarra et al.

While increasing precision almost always improves accuracy, there are rare cases where this does not happen.

- The loop in FP32 iterates 10 times (as per real semantics)
- In FP64, it iterates11 times

```
#include <stdio.h>
int main() {
 int cnt=0;
 for (float i = 0.0f; i < 1.0f; i += 0.1) {
   cnt++:
    printf("%.1f ", i);
  printf("\n");
  printf("cnt=%d\n",cnt);
 return 0;
-UU-:--- F1 floatCnt.c
                             Top
                                   L5
#include <stdio.h>
int main() {
 int cnt=0;
 for (double i = 0.0f; i < 1.0f; i += 0.1) {
    cnt++;
    printf("%.1f ", i);
 printf("\n");
  printf("cnt=%d\n",cnt);
 return 0;
-UU-:--- F1 doubCnt.c
                             Top
                                   L1
                                          (C/*l Abbrev)
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$ gcc -o fc floatCnt.c
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$ gcc -o dc doubCnt.c
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$ ./fc
0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9
cnt=10
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$ ./dc
0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
cnt=11
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$
```

Non-intuitive behavior caused by Rounding!

```
#include <stdio.h>
int main() {
 int cnt=0;
 for (float i = 0.0f; i < 1.0f; i += 0.1) {
   cnt++;
   printf("%.1f ", i):
 printf("\n");
 printf("cnt=%d\n",cnt);
 return 0:
                                 L5
UU-:--- F1 floatCnt.c
                                         (C/*l Abbrev) ------
#include <stdio.h>
int main() {
 int cnt=0:
 for (double i = 0.0f; i < 1.0f; i += 0.1) {
   cnt++;
   printf("%.1f ", i);
 printf("\n");
 printf("cnt=%d\n",cnt);
 return 0:
-UU-:--- F1 doubCnt.c
                            Top L1
                                         (C/*l Abbrev) ------
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$ gcc -o fc floatCnt.c
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$ gcc -o dc doubCnt.c
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$ ./fc
0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9
cnt=10
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$ ./dc
0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
cnt=11
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$
```

```
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$ ./
(sign: 0, exp: 0, mantissa: 0x000000)
Value: 0.1000000014901161193848
(sign: 0, exp: 123, mantissa: 0x4CCCCD)
Value: 0.2000000029802322387695
(sign: 0, exp: 124, mantissa: 0x4CCCCD)
Value: 0.3000000119209289550781
(sign: 0, exp: 125, mantissa: 0x19999A)
Value: 0.4000000059604644775391
(sign: 0, exp: 125, mantissa: 0x4CCCD)
Value: 0.5000000000000000000000000
(sign: 0, exp: 126, mantissa: 0x000000)
Value: 0.6000000238418579101562
(sign: 0, exp: 126, mantissa: 0x19999A)
Value: 0.7000000476837158203125
(sign: 0, exp: 126, mantissa: 0x333334)
Value: 0.8000000715255737304688
(sign: 0, exp: 126, mantissa: 0x4CCCCE)
Value: 0.9000000953674316406250
(sign: 0, exp: 126, mantissa: 0x666668)
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$ ./
(sign: 0, exp: 0, mantissa: 0x0000000000000)
Value: 0.1000000000000000055511
(sign: 0, exp: 1019, mantissa: 0x99999999999A)
Value: 0.20000000000000000111022
(sign: 0, exp: 1020, mantissa: 0x999999999999A)
Value: 0.30000000000000000444089
(sign: 0, exp: 1021, mantissa: 0x33333333333334)
Value: 0.40000000000000000222045
(sign: 0, exp: 1021, mantissa: 0x99999999999A)
Value: 0.500000000000000000000000
(sign: 0, exp: 1022, mantissa: 0x0000000000000)
Value: 0.599999999999999777955
(sign: 0, exp: 1022, mantissa: 0x33333333333333)
Value: 0.699999999999999555911
(sign: 0, exp: 1022, mantissa: 0x6666666666666)
Value: 0.799999999999999333866
(sign: 0, exp: 1022, mantissa: 0x999999999999)
Value: 0.8999999999999999111822
(sign: 0, exp: 1022, mantissa: 0xCCCCCCCCCCCC)
Value: 0.999999999999998889777
(sign: 0, exp: 1022, mantissa: 0xFFFFFFFFFFFF)
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$
(base) ganesh@ganesh-XPS-17-9710:~/parfloat-cs6961/c-prec-expts$ ☐
```

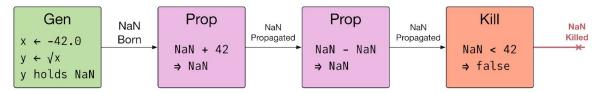
base, gallesil@gallesil-Ars-17-9/10.-/pai/toat-cs0901/c-piec-expts\$ gc

Focus of this tutorial

- Excessive Rounding Error (not addressed in this tutorial)
- Exhausting the range (addressed)
 - Detection of Exceptions (main focus)

FP Exceptions (Exceptional Values) are Undesirable (often are bugs)

- FP32: $(1.9 \times 10^{19})^2 \rightarrow INF$... (in E4M3, $16^2 = INF$) \leftarrow a popular ML low-prec format \circ Easily generated during dot-product
- Of the FP Exceptions, we focus on INF and NaN
 - NaNs can skew control-flow
- This macro is buggy
 - o #define MAX(x,y) ((x >= y) ? x : y)
 hint : consider x == NaN or y == NaN
- Exceptions often vanish harmlessly (e.g. in ML codes)
 - But it is difficult to know which control-flows they have affected meanwhile



Tools Presented



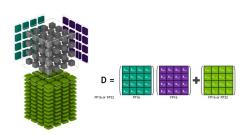
GPU-FPX + NixNaN

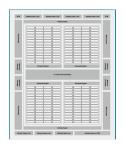
FloatGuard

NVIDIA GPUs









AMD GPUs

















Agenda

- 08:30 : Introduction
- 08:35 : FPChecker
- 09:35 : NixNan (Part-1)
- 09:55 : Fill Survey
- 10:00 : Coffee Break
- 10:30 : Brief Recap
- 10:35 : NixNan (Part-2)
- 11:05 : FloatGuard
- 11:50 : Closing Remarks
- 11:55 : Fill Survey