

# 1 Universal Numbers Library: Multi-format Variable 2 Precision Arithmetic Library

3 E. Theodore L. Omtzigt <sup>1\*</sup> and James Quinlan <sup>2\*¶</sup>

4 <sup>1</sup> Stillwater Supercomputing, Inc, USA <sup>2</sup> School of Mathematical and Physical Sciences, University of  
5 New England, USA ¶ Corresponding author \* These authors contributed equally.

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

## Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Open Journals](#) 

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

## 6 Summary

7 *Universal Numbers Library*, or *Universal* for short, is a self-contained C++ header-only template  
8 library that contains implementations of many number representations and standard arithmetic  
9 on arbitrary configuration integer and real numbers (Omtzigt et al., 2020). In particular,  
10 the library includes integers, decimals, fixed-points, rationals, linear floats, tapered floats,  
11 logarithmic, SORNs, interval, level-index, and adaptive-precision binary and decimal integers  
12 and floats, each offering a verification suite.

13 The primary pattern using a posit number type as example, is:

```
#include <universal/number/posit/posit.hpp>

template<typename Real>
Real MyKernel(const Real& a, const Real& b) {
    return a * b; // replace this with your kernel computation
}

constexpr double pi = 3.14159265358979323846;

int main() {
    using Real = sw::universal::posit<32,2>;

    Real a = sqrt(2);
    Real b = pi;
    std::cout << "Result: " << MyKernel(a, b) << std::endl;
}
```

14 *Universal* delivers software and hardware co-design capabilities to develop low and mixed-  
15 precision algorithms for reducing energy consumption in signal processing, Industry 4.0, machine  
16 learning, robotics, and high-performance computing applications (Omtzigt & Quinlan, 2022).  
17 The package includes command-line tools for visualizing and interrogating numeric encodings,  
18 an interface for setting and querying bits, and educational examples showcasing performance  
19 gain and numerical accuracy with the different number systems. In addition, a Docker container  
20 is available to experiment without cloning and building from the source code.

```
21 $ docker pull stillwater/universal
22 $ docker run -it --rm stillwater/universal bash
```

23 *Universal* started in 2017 as a bit-level arithmetic reference implementation of the evolving  
24 unum Type III (posit and valid) standard. However, the demands for supporting various number  
25 systems, such as adaptive-precision integers to solve large factorials, adaptive-precision floats  
26 to act as Oracles, or comparing linear and tapered floats provided the opportunity to create a

27 complete platform for numerical analysis and computational mathematics. As a result, several  
28 projects have leveraged *Universal*, including Matrix Template Library (MTL4), Geometry +  
29 Simulation Modules (G+SMO), Bembel, a fast IGA BEM solver, and the Odeint ODE solver.

30 The default build configuration will produce the command line tools, a playground, and  
31 educational and application examples. It is also possible to construct the full regression suite  
32 across all the number systems. For instance, the shortened output for the commands `single`  
33 and `single 1.23456789` are below.

```
$ single
min exponent          -125
max exponent          128
radix                  2
radix digits           24
min                   1.17549e-38
max                   3.40282e+38
lowest                -3.40282e+38
epsilon (1+1ULP-1)   1.19209e-07
round_error            0.5
denorm_min            1.4013e-45
infinity               inf
quiet_NAN              nan
signaling_NAN         nan
...
```

```
$ single 1.23456789
scientific   : 1.2345679
triple form  : (+,0,0b00111100000011001010010)
binary form  : 0b0.0111'1111.001'1110'0000'0110'0101'0010
color coded  : 0b0.0111'1111.001'1110'0000'0110'0101'0010
```

## 34 Statement of need

35 High-performance computing (HPC), machine learning, and deep learning tasks (e.g.,  
36 [Carmichael et al., 2019](#); [Cococcioni et al., 2022](#); [Desrentes et al., 2022](#)) have increased  
37 environmental impacts and financial costs due to massive energy consumption ([Haidar,  
38 Abdelfattah, et al., 2018](#)). These both result from growth requirements in processing and  
39 storage. In addition to redesigning algorithms to minimize data movement and processing,  
40 modern systems increasingly support multi-precision arithmetic in hardware ([Haidar, Tomov,  
41 et al., 2018](#)). Recently, NVIDIA added support for low-precision formats to its top-level GPUs  
42 to perform tensor operations ([Choquette et al., 2021](#)), including a 19-bit format having an  
43 exponent of 8 bits and a mantissa of 10 bits (see also [Intel Corporation (2018); kharya:2020].  
44 In addition, the “Brain Floating Point Format,” commonly referred to as “bfloat16,” is a  
45 format developed by Google that enables the training and operation of deep neural networks  
46 using specialized processors called Tensor Processing Units, or TPUs, at higher performance  
47 and cheaper cost ([Wang & Kanwar, 2019](#)). As a result, we see a trend to redesign many  
48 standard algorithms. In particular, designing fast and energy-efficient linear solvers is an active  
49 area of research where low-precision numerics plays a fundamental role ([Carson & Higham,  
50 2018](#); [Haidar et al., 2017](#); [Haidar, Tomov, et al., 2018](#); [Haidar, Abdelfattah, et al., 2018](#);  
51 [Higham et al., 2019](#)).

52 While the primary motivation for low-precision arithmetic is its high performance and energy  
53 efficiency, mixed-precision algorithm designs aim to identify and exploit opportunities to right-  
54 scale the number of systems used for critical computational paths representing the execution  
55 bottleneck. Furthermore, when these algorithms are incorporated into embedded devices and

56 custom hardware engines, we approach optimal performance and power efficiency. Therefore,  
57 investigations into computational mathematics and measuring mixed-precision algorithms'  
58 accuracy, efficiency, robustness, and stability are needed.

59 Custom number systems that optimize the entire system's performance per watt (W) are  
60 crucial components with the rise of embedded devices demanding intelligent behavior. Likewise,  
61 energy efficiency is an essential differentiator for embedded intelligence applications. Using  
62 the distinct arithmetic requirements of the control and data flow can result in considerable  
63 performance and power efficiency gains when creating unique compute solutions. Even within  
64 the data flow, we observe many requirements for precision and the required dynamic range of  
65 the arithmetic operations.

## 66 Verification Suite

67 Each number system contained within *Universal* is supported by a comprehensive verification  
68 environment testing library class API consistency, logic and arithmetic operators, the standard  
69 math library, arithmetic exceptions, and language features such as compile-time constexpr.  
70 The verification suite is run as part of the `make test` command in the build directory.

71 Due to the size of the library, the build system for *Universal* allows for fine-grain control to  
72 subset the test environment for productive development and verification. There are twelve  
73 core build category flags defined:

- 74     ▪ BUILD\_APPLICATIONS
- 75     ▪ BUILD\_BENCHMARKS
- 76     ▪ BUILD\_CI
- 77     ▪ BUILD\_CMD\_LINE\_TOOLS
- 78     ▪ BUILD\_C\_API
- 79     ▪ BUILD\_DEMONSTRATION
- 80     ▪ BUILD\_EDUCATION
- 81     ▪ BUILD\_LINEAR\_ALGEBRA
- 82     ▪ BUILD\_MIXEDPRECISION\_SDK
- 83     ▪ BUILD\_NUMBERS
- 84     ▪ BUILD\_NUMERICS
- 85     ▪ BUILD\_PLAYGROUND

86 The flags, when set during cmake configuration, i.e. `cmake -DBUILD_CI=ON ..`, enable build  
87 targets specialized to the category. For example, the `BUILD_CI` flag turns on the continuous  
88 integration regression test suites for all number systems, and the `BUILD_APPLICATIONS` flag will  
89 build all the example applications that provide demonstrations of mixed-precision, high-accuracy,  
90 reproducible and/or interval arithmetic.

91 Each build category contains individual targets that further refine the build targets. For  
92 example, `cmake -DBUILD_NUMBER_POSIT=ON -DBUILD_DEMONSTRATION=OFF ..` will build just  
93 the fixed-size, arbitrary configuration posit number system regression environment.

94 It is also possible to run specific test suite components, for example, to validate algorithmic  
95 changes to more complex arithmetic functions, such as square root, exponent, logarithm, and  
96 trigonometric functions. Here is an example, assuming that the logarithmic number system  
97 has been configured during the cmake build generation:

```
98 $ make lns_trigonometry
```

99 The repository's README file has all the details about the build and regression environment  
100 and how to streamline its operation.

## 101 Availability and Documentation

102 *Universal Number Library* is available under the [MIT License](#). The package may be cloned or  
103 forked from the [GitHub repository](#). Documentation is provided via Docs, including a tutorial  
104 introducing primary functionality and detailed reference and communication networks. The  
105 library employs extensive unit testing.

## 106 Acknowledgements

107 We want to acknowledge all code contributions, bug reports, and feedback from numerous  
108 other developers and users.

## 109 References

110 Carmichael, Z., Langroudi, H. F., Khazanov, C., Lillie, J., Gustafson, J. L., & Kudithipudi,  
111 D. (2019). Deep positron: A deep neural network using the posit number system. *2019*  
112 *Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 1421–1426.

113 Carson, E., & Higham, N. J. (2018). Accelerating the solution of linear systems by iterative  
114 refinement in three precisions. *SIAM Journal on Scientific Computing*, *40*(2), A817–A847.

115 Choquette, J., Gandhi, W., Giroux, O., Stam, N., & Krashinsky, R. (2021). NVIDIA A100  
116 tensor core GPU: Performance and innovation. *IEEE Micro*, *41*(2), 29–35.

117 Cococcioni, M., Rossi, F., Emanuele, R., & Saponara, S. (2022). Small reals representations  
118 for deep learning at the edge: A comparison. *Proc. Of the 2022 Conference on Next*  
119 *Generation Arithmetic (CoNGA'22)*.

120 Desrentes, O., Resmerita, D., & Dinechin, B. D. de. (2022). A Posit8 decompression operator  
121 for deep neural network inference. *Next Generation Arithmetic: Third International*  
122 *Conference, CoNGA 2022, Singapore, March 1–3, 2022, Revised Selected Papers*, 13253,  
123 14.

124 Haidar, A., Abdelfattah, A., Zounon, M., Wu, P., Pranesh, S., Tomov, S., & Dongarra,  
125 J. (2018). The design of fast and energy-efficient linear solvers: On the potential of  
126 half-precision arithmetic and iterative refinement techniques. *International Conference on*  
127 *Computational Science*, 586–600.

128 Haidar, A., Tomov, S., Dongarra, J., & Higham, N. J. (2018). Harnessing GPU tensor  
129 cores for fast FP16 arithmetic to speed up mixed-precision iterative refinement solvers.  
130 *SC18: International Conference for High Performance Computing, Networking, Storage*  
131 *and Analysis*, 603–613.

132 Haidar, A., Wu, P., Tomov, S., & Dongarra, J. (2017). Investigating half precision arithmetic  
133 to accelerate dense linear system solvers. *Proceedings of the 8th Workshop on Latest*  
134 *Advances in Scalable Algorithms for Large-Scale Systems*, 1–8.

135 Higham, N. J., Pranesh, S., & Zounon, M. (2019). Squeezing a matrix into half precision, with  
136 an application to solving linear systems. *SIAM Journal on Scientific Computing*, *41*(4),  
137 A2536–A2551.

138 Intel Corporation. (2018). *BFLOAT16 - Hardware Numerics Definition*. <https://www.intel.com/content/dam/develop/external/us/en/documents/bf16-hardware-numerics-definition-white-paper.pdf>  
139  
140

- 141 Omtzigt, E. T. L., Gottschling, P., Seligman, M., & Zorn, W. (2020). Universal Numbers  
142 Library: Design and implementation of a high-performance reproducible number systems  
143 library. *arXiv:2012.11011*.
- 144 Omtzigt, E. T. L., & Quinlan, J. (2022). Universal: Reliable, reproducible, and energy-efficient  
145 numerics. *Conference on Next Generation Arithmetic*, 100–116.
- 146 Wang, S., & Kanwar, P. (2019). BFloat16: The secret to high performance on cloud TPUs.  
147 *Google Cloud Blog*, 4.

DRAFT