



bfloat16 floating-point format

The **bfloat16 (brain floating point)**^{[1][2]} floating-point format is a computer number format occupying 16 bits in computer memory; it represents a wide dynamic range of numeric values by using a floating radix point. This format is a shortened (16-bit) version of the 32-bit IEEE 754 single-precision floating-point format (binary32) with the intent of accelerating machine learning and near-sensor computing.^[3] It preserves the approximate dynamic range of 32-bit floating-point numbers by retaining 8 exponent bits, but supports only an 8-bit precision rather than the 24-bit significand of the binary32 format. More so than single-precision 32-bit floating-point numbers, bfloat16 numbers are unsuitable for integer calculations, but this is not their intended use. Bfloat16 is used to reduce the storage requirements and increase the calculation speed of machine learning algorithms.^[4]

The bfloat16 format was developed by Google Brain, an artificial intelligence research group at Google. It is utilized in many CPUs, GPUs, and AI processors, such as Intel Xeon processors (AVX-512 BF16 extensions), Intel Data Center GPU, Intel Nervana NNP-L1000, Intel FPGAs,^{[5][6][7]} AMD Zen, AMD Instinct, NVIDIA GPUs, Google Cloud TPUs,^{[8][9][10]} AWS Inferentia, AWS Trainium, ARMv8.6-A,^[11] and Apple's M2^[12] and therefore A15 chips and later. Many libraries support bfloat16, such as CUDA,^[13] Intel oneAPI Math Kernel Library, AMD ROCm,^[14] AMD Optimizing CPU Libraries, PyTorch, and TensorFlow.^{[10][15]} On these platforms, bfloat16 may also be used in mixed-precision arithmetic, where bfloat16 numbers may be operated on and expanded to wider data types.

bfloat16 floating-point format

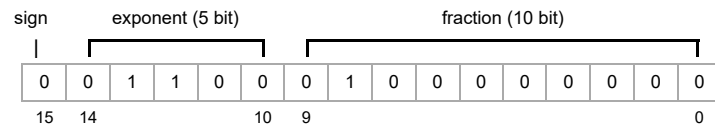
bfloat16 has the following format:

- Sign bit: 1 bit
- Exponent width: 8 bits
- Significand precision: 8 bits (7 explicitly stored, with an implicit leading bit), as opposed to 24 bits in a classical single-precision floating-point format

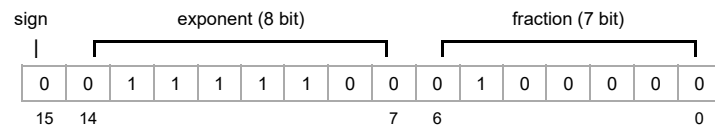
The bfloat16 format, being a shortened IEEE 754 single-precision 32-bit float, allows for fast conversion to and from an IEEE 754 single-precision 32-bit float; in conversion to the bfloat16 format, the exponent bits are preserved while the significand field can be reduced by truncation (thus corresponding to round toward 0) or other rounding mechanisms, ignoring the NaN special case. Preserving the exponent bits maintains the 32-bit float's range of $\approx 10^{-38}$ to $\approx 3 \times 10^{38}$.^[16]

The bits are laid out as follows:

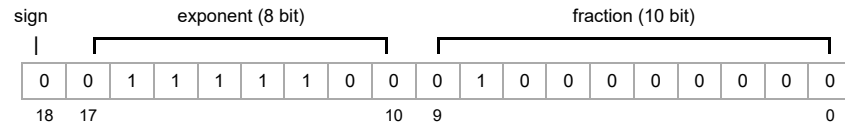
IEEE half-precision 16-bit float



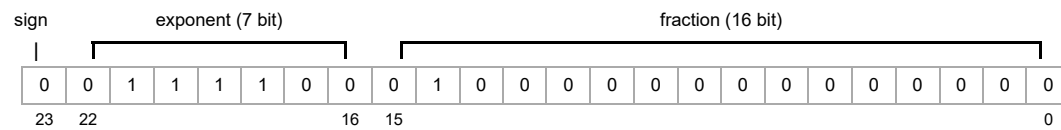
bfloat16



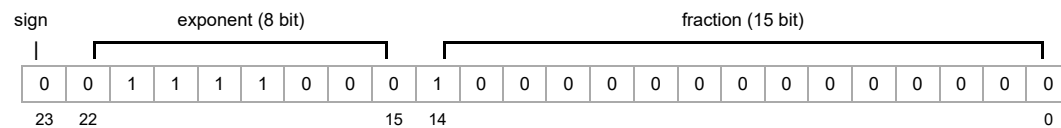
Nvidia's TensorFloat-32 (19 bits)



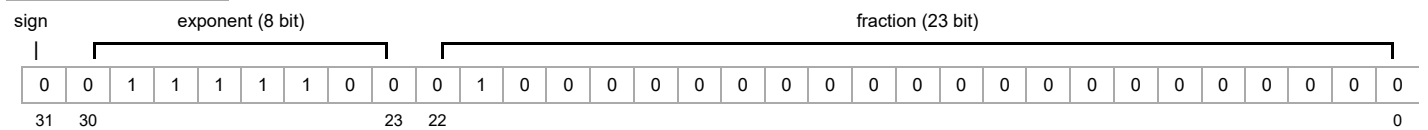
ATI's fp24 format [17]



Pixar's PXR24 format



IEEE 754 single-precision 32-bit float



Exponent encoding

The bfloat16 binary floating-point exponent is encoded using an offset-binary representation, with the zero offset being 127; also known as exponent bias in the IEEE 754 standard.

- $E_{\min} = 01_H - 7F_H = -126$
- $E_{\max} = FE_H - 7F_H = 127$
- Exponent bias = $7F_H = 127$

Thus, in order to get the true exponent as defined by the offset-binary representation, the offset of 127 has to be subtracted from the value of the exponent field.

The minimum and maximum values of the exponent field (00_H and FF_H) are interpreted specially, like in the IEEE 754 standard formats.

Exponent	Significand zero	Significand non-zero	Equation
00_H	<u>zero</u> , -0	<u>subnormal numbers</u>	$(-1)^{\text{signbit}} \times 2^{-126} \times 0.\text{significandbits}$
$01_H, \dots, FE_H$	normalized value		$(-1)^{\text{signbit}} \times 2^{\text{exponentbits} - 127} \times 1.\text{significandbits}$
FF_H	$\pm\text{infinity}$	<u>NaN</u> (quiet, signaling)	

The minimum positive normal value is $2^{-126} \approx 1.18 \times 10^{-38}$ and the minimum positive (subnormal) value is $2^{-126-7} = 2^{-133} \approx 9.2 \times 10^{-41}$.

Rounding and conversion

The most common use case is the conversion between IEEE 754 binary32 and bfloat16. The following section describes the conversion process and its rounding scheme in the conversion. Note that there are other possible scenarios of format conversions to or from bfloat16. For example, int16 and bfloat16.

- From binary32 to bfloat16. When bfloat16 was first introduced as a storage format,^[15] the conversion from IEEE 754 binary32 (32-bit floating point) to bfloat16 is truncation (round toward 0). Later on, when it becomes the input of matrix multiplication units, the conversion can have various rounding mechanisms depending on the hardware platforms. For example, for Google TPU, the rounding scheme in the conversion

is round-to-nearest-even;^[18] ARM uses the non-IEEE Round-to-Odd mode;^[19] for NVIDIA, it supports converting float number to bfloat16 precision in round-to-nearest-even mode.^[20]

- From bfloat16 to binary32. Since binary32 can represent all exact values in bfloat16, the conversion simply pads 16 zeros in the significand bits.^[18]

Encoding of special values

Positive and negative infinity

Just as in IEEE 754, positive and negative infinity are represented with their corresponding sign bits, all 8 exponent bits set (FF_{hex}) and all significand bits zero. Explicitly,

```
val    s_exponent_signcnd
+inf = 0_11111111_00000000
-inf = 1_11111111_00000000
```

Not a Number

Just as in IEEE 754, NaN values are represented with either sign bit, all 8 exponent bits set (FF_{hex}) and not all significand bits zero. Explicitly,

```
val    s_exponent_signcnd
+NaN = 0_11111111_klmnopq
-NaN = 1_11111111_klmnopq
```

where at least one of k, l, m, n, o, p , or q is 1. As with IEEE 754, NaN values can be quiet or signaling, although there are no known uses of signaling bfloat16 NaNs as of September 2018.

Range and precision

Bfloat16 is designed to maintain the number range from the 32-bit IEEE 754 single-precision floating-point format (binary32), while reducing the precision from 24 bits to 8 bits. This means that the precision is between two and three decimal digits, and bfloat16 can represent finite values up to about 3.4×10^{38} .

Examples

These examples are given in bit *representation*, in hexadecimal and binary, of the floating-point value. This includes the sign, (biased) exponent, and significand.

```
3f80 = 0 01111111 00000000 = 1
c000 = 1 10000000 00000000 = -2
```

```
7f7f = 0 11111110 11111111 =  $(2^8 - 1) \times 2^{-7} \times 2^{127} = 3.38953139 \times 10^{38}$  (max finite positive value in bfloat16 precision)
0080 = 0 00000001 00000000 =  $2^{-126} \approx 1.175494351 \times 10^{-38}$  (min normalized positive value in bfloat16 precision and single-precision floating point)
```

The maximum positive finite value of a normal bfloat16 number is $3.38953139 \times 10^{38}$, slightly below $(2^{24} - 1) \times 2^{-23} \times 2^{127} = 3.402823466 \times 10^{38}$, the max finite positive value representable in single precision.

Zeros and infinities

```
0000 = 0 00000000 00000000 = 0
8000 = 1 00000000 00000000 = -0
```

```
7f80 = 0 11111111 00000000 = infinity
ff80 = 1 11111111 00000000 = -infinity
```

Special values

```
4049 = 0 10000000 1001001 =  $3.140625 \approx \pi$  ( pi )
3eab = 0 01111101 0101011 =  $0.333984375 \approx 1/3$ 
```

NaNs

```
ffc1 = x 11111111 1000001 => qNaN
ff81 = x 11111111 0000001 => sNaN
```

See also

- Half-precision floating-point format: 16-bit float w/ 1-bit sign, 5-bit exponent, and 11-bit significand, as defined by [IEEE 754](#)
- [ISO/IEC 10967](#), Language Independent Arithmetic
- [Primitive data type](#)
- [Minifloat](#)
- [Google Brain](#)
- [Lawsuit against Google for its use of bfloat16 in TPU](#)

References

1. Teich, Paul (2018-05-10). "Tearing Apart Google's TPU 3.0 AI Coprocessor" (<https://www.nextplatform.com/2018/05/10/tearing-apart-google-tpu-3-0-ai-coprocessor/>). *The Next Platform*. Retrieved 2020-08-11. "Google invented its own internal floating point format called "bfloat" for "brain floating point" (after Google Brain)."
2. Wang, Shibo; Kanwar, Pankaj (2019-08-23). "BFLOAT16: The secret to high performance on Cloud TPUs" (<https://cloud.google.com/blog/products/ai-machine-learning/bfloat16-the-secret-to-high-performance-on-cloud-tpus>). *Google Cloud*. Retrieved 2020-08-11. "This custom floating point format is called "Brain Floating Point Format," or "bfloat16" for short. The name flows from "Google Brain", which is an artificial intelligence research group at Google where the idea for this format was conceived."
3. Tagliavini, Giuseppe; Mach, Stefan; Rossi, Davide; Marongiu, Andrea; Benin, Luca (2018). "A transprecision floating-point platform for ultra-low power computing". *2018 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. pp. 1051–1056. arXiv:1711.10374 (<https://arxiv.org/abs/1711.10374>). doi:10.23919/DAT.2018.8342167 (<https://doi.org/10.23919%2FDAT.2018.8342167>). ISBN 978-3-9819263-0-9. S2CID 5067903 (<https://api.semanticscholar.org/CorpusID:5067903>).
4. Dr. Ian Cutress (2020-03-17). "Intel: Cooper lake Plans: Why is BF16 Important?" (<https://web.archive.org/web/20200318044239/https://www.anandtech.com/show/15631/intels-cooper-lake-plans-the-chip-that-wasnt-meant-to-exist-dies-for-you>). Archived from the original (<https://www.anandtech.com/show/15631/intels-cooper-lake-plans-the-chip-that-wasnt-meant-to-exist-dies-for-you>) on March 18, 2020. Retrieved 2020-05-12. "The bfloat16 standard is a targeted way of representing numbers that give the range of a full 32-bit number, but in the data size of a 16-bit number, keeping the accuracy close to zero but being a bit more loose with the accuracy near the limits of the standard. The bfloat16 standard has a lot of uses inside machine learning algorithms, by offering better accuracy of values inside the algorithm while affording double the data in any given dataset (or doubling the speed in those calculation sections)."
5. Khari Johnson (2018-05-23). "Intel unveils Nervana Neural Net L-1000 for accelerated AI training" (<https://venturebeat.com/2018/05/23/intel-unveils-nervana-neural-net-l-1000-for-accelerated-ai-training/>). *VentureBeat*. Retrieved 2018-05-23. "...Intel will be extending bfloat16 support across our AI product lines, including Intel Xeon processors and Intel FPGAs."
6. Michael Feldman (2018-05-23). "Intel Lays Out New Roadmap for AI Portfolio" (<https://www.top500.org/news/intel-lays-out-new-roadmap-for-ai-portfolio/>). *TOP500 Supercomputer Sites*. Retrieved 2018-05-23. "Intel plans to support this format across all their AI products, including the Xeon and FPGA lines"
7. Lucian Armasu (2018-05-23). "Intel To Launch Spring Crest, Its First Neural Network Processor, In 2019" (<https://www.tomshardware.com/news/intel-neural-network-processor-lake-crest,37105.html>). *Tom's Hardware*. Retrieved 2018-05-23. "Intel said that the NNP-L1000 would also support bfloat16, a numerical format that's being adopted by all the ML industry players for neural networks. The company will also support bfloat16 in its FPGAs, Xeons, and other ML products. The Nervana NNP-L1000 is scheduled for release in 2019."
8. "Available TensorFlow Ops | Cloud TPU | Google Cloud" (<https://cloud.google.com/tpu/docs/tensorflow-ops>). *Google Cloud*. Retrieved 2018-05-23. "This page lists the TensorFlow Python APIs and graph operators available on Cloud TPU."
9. Elmar Haußmann (2018-04-26). "Comparing Google's TPuv2 against Nvidia's V100 on ResNet-50" (<https://web.archive.org/web/20180426200043/https://blog.riseml.com/comparing-google-tpuv2-against-nvidia-v100-on-resnet-50-c2bbb6a51e5e>). *RiseML Blog*. Archived from the original (<https://blog.riseml.com/comparing-google-tpuv2-against-nvidia-v100-on-resnet-50-c2bbb6a51e5e>) on 2018-04-26. Retrieved 2018-05-23. "For the Cloud TPU, Google recommended we use the bfloat16 implementation from the official TPU repository with TensorFlow 1.7.0. Both the TPU and GPU implementations make use of mixed-precision computation on the respective architecture and store most tensors with half-precision."
10. Tensorflow Authors (2018-07-23). "ResNet-50 using BFloat16 on TPU" (https://github.com/tensorflow/tpu/tree/0ece10f6f4e523eab79aba0247b513fe57d38ae6/models/experimental/resnet_bfloat16). *Google*. Retrieved 2018-11-06.
11. "BFLOAT16 extensions for Armv8-A" (https://community.arm.com/developer/ip-products/processors/b/ml-ip-blog/posts/bfloat16-processing-for-neural-networks-on-armv8_2d00_a). *community.arm.com*. 29 August 2019. Retrieved 2019-08-30.
12. "AArch64: add support for newer Apple CPUs · llvm/llvm-project@677da09" (<https://github.com/llvm/llvm-project/commit/677da09d0259d7530d32e85cb561bee15f0066e2>). *GitHub*. Retrieved 2023-05-08.
13. "CUDA Library bfloat16 Intrinsics" (https://docs.nvidia.com/cuda/cuda-math-api/group_CUDA_MATH_INTRINSIC_BFLOAT16.html#grou_p_CUDA_MATH_INTRINSIC_BFLOAT16).
14. "ROCm version history" (https://github.com/RadeonOpenCompute/ROCm/blob/8bd9a527405cb466d45b3b343a33434e79b1d387/version_history.md#miopen-20). *github.com*. Retrieved 2019-10-23.
15. Joshua V. Dillon, Ian Langmore, Dustin Tran, Eugene Brevdo, Srinivas Vasudevan, Dave Moore, Brian Patton, Alex Alemi, Matt Hoffman, Rif A. Saurous (2017-11-28). TensorFlow Distributions (Report). arXiv:1711.10604 (<https://arxiv.org/abs/1711.10604>). Bibcode:2017arXiv171110604D (<https://ui.adsabs.harvard.edu/abs/2017arXiv171110604D>). Accessed 2018-05-23. "All operations in TensorFlow Distributions are numerically stable across half, single, and double floating-point precisions (as TensorFlow dtypes: tf.bfloat16 (truncated floating point), tf.float16, tf.float32, tf.float64). Class constructors have a validate_args flag for numerical asserts"
16. "Livestream Day 1: Stage 8 (Google I/O '18) - YouTube" (<https://www.youtube.com/watch?v=vm67WcLzfv&t=2555>). *Google*. 2018-05-08. Retrieved 2018-05-23. "In many models this is a drop-in replacement for float-32"
17. Buck, Ian (2005-03-13), "Chapter 32. Taking the Plunge into GPU Computing" (<https://developer.nvidia.com/gpugems/gpugems2/part-iv-general-purpose-computation-gpus-primer/chapter-32-taking-plunge-gpu>), in Pharr, Matt (ed.), *GPU Gems*, Addison-Wesley, ISBN 0-321-33559-7, retrieved 2018-04-05.
18. "The bfloat16 numerical format" (<https://cloud.google.com/tpu/docs/bfloat16>). *Google Cloud*. Retrieved 2023-07-11. "On TPU, the rounding scheme in the conversion is round to nearest even and overflow to inf."
19. "Arm A64 Instruction Set Architecture" (<https://developer.arm.com/documentation/ddi0596/2021-12/SVE-Instructions/BFDDOT--vectors---BFLOAT16-floating-point-dot-product->). *developer.arm.com*. Retrieved 2023-07-26. "Uses the non-IEEE Round-to-Odd rounding mode."
20. "1.3.5. BFloat16 Precision Conversion and Data Movement" (https://docs.nvidia.com/cuda/pdf/CUDA_Math_API.pdf) (PDF). *docs.nvidia.com*. p. 199. Retrieved 2023-07-26. "Converts float number to nv_bfloat16 precision in round-to-nearest-even mode and returns nv_bfloat16 with converted value."

