
EFloat: Entropy-coded Floating Point Format for Deep Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In a large class of deep learning models, specifically vector embedding models
2 in NLP, we observe that floating point exponent values tend to cluster around
3 few unique values, presenting entropy encoding opportunities. The proposed
4 EFloat floating point number format encodes frequent exponent values and signs
5 with Huffman codes to minimize the average exponent field width while keeping
6 the original exponent range unchanged. Saved bits then become available to the
7 significand increasing the EFloat numeric precision on average by 4.3 bits compared
8 to other low-precision floating point formats of equal bit budget. The EFloat format
9 makes 8-bit and smaller floats practical by preserving the full exponent range
10 of a 32-bit floating point representation. We currently use the EFloat format for
11 compressing and saving memory used in large NLP deep learning models while
12 I/O and memory bandwidth savings in GPUs and AI accelerators are also possible.
13 Using RMS-error as a precision metric, we demonstrate that EFloat provides more
14 accurate floating point representation than other formats with the same bit budget.
15 EF12 with 12-bit budget has less end-to-end application error than the 16-bit
16 BFloat16. EF16 RMS-error is 17 to 35 times less than BFloat16 RMS-error for a range
17 of datasets. Using the NDCG metric for evaluating ranked results of similarity and
18 dissimilarity queries in NLP, we demonstrate that EFloat matches the result quality
19 of other floating point representations with larger bit budgets.

20 1 Introduction

21 As natural language processing (NLP) models expand their capabilities, complexity, and training costs,
22 the model sizes have been increasing dramatically. For example, state-of-the-art transformer-based
23 NLP models such as BERT (Vaswani et al. (2017)), Megatron-LM (Shoeybi et al. (2020)), Open AI
24 GPT-3 (Brown et al. (2020)), or Google Switch-C Transformers (Fedus et al. (2021)), contain from
25 hundreds of millions, to even trillion parameters (Hoefler (2020); Fedus et al. (2021)). Although
26 NLP model compression is a very active area of research (Section 6), its current focus is on model
27 inference scenarios, in which reduced precision and integer quantization are commonly used, given
28 that the original model need not be restored.

29 The primary goal for this work is to explore compression strategies for large vector embedding models
30 such that one can recover or minimize the loss in the original model, or use the same compressed
31 model in the inference phase. The *database embedding (db2Vec)*, a vector embedding technique
32 designed to develop semantic models from multi-modal relational database tables (Bordawekar and
33 Shmueli (2016, 2017)), forms the impetus behind this exploration. Db2Vec differs from its NLP
34 counterparts, such as Word2Vec (Mikolov et al. (2013)) and GloVe (Pennington et al. (2014b)), in
35 that its source data follows the relational data model (Date (1982)) (the source data is not a natural
36 language document but a relational database table). Considering the relational database tables can

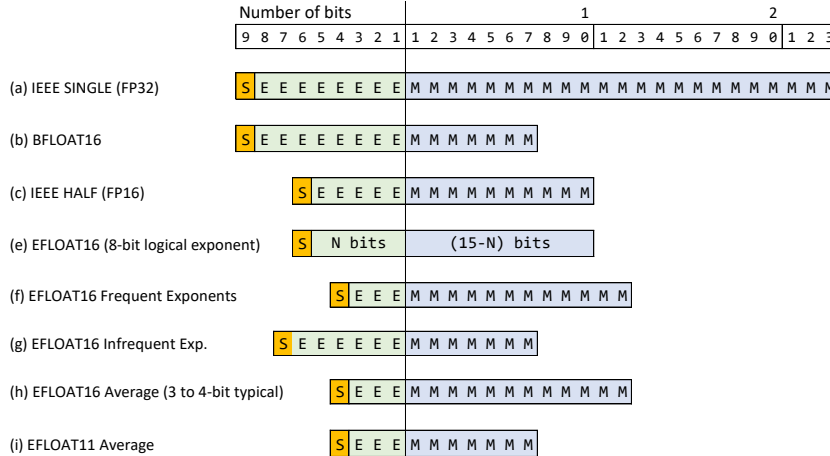


Figure 1: Floating point formats are compared. EFloat has a fixed total width, but the boundary between the exponent and the significand is variable (e). The exponent is entropy coded, providing an average of 4.3 extra bits of precision to the significand (e.g., (h)), while keeping the logical exponent range at 8 bits, same as that of FP32. EFloat has greater precision and range than the existing FP formats having the same bit budget.

37 be very large (e.g., billions of rows in a table) with a large number of unique tokens, it leads to a
 38 much larger vocabulary than a traditional natural language document. As a result, trained db2Vec
 39 models can be very large. Any trained vector embedding model is a snapshot of it’s weight matrices
 40 and consists of weight values represented typically with IEEE 32-bit single-precision floating point
 41 (**FP32**) format. Therefore, we focus on compression approaches that exploit different low-precision
 42 floating point formats.

43 Existing low-precision floating-point(FP) formats make a tradeoff between the number of exponent
 44 and significand bits. An FP number is of the form

$$-1^{signbit} \times 2^{exponent-bias} \times significand$$

45 The exponent largely determines the range of minimum and maximum values representable by the
 46 format and the significand width determines the precision (A constant bias is added to exponents to
 47 make them all positive integers which simplifies magnitude comparisons.) For example, the BFloat16
 48 (**BF16**) format with an 8-bit exponent and 7-bit significand has a wide range but low precision when
 49 compared to FP32 (Wang and Kumar (2019)) and IEEE 754-2019 Half-precision (**FP16**) as illustrated
 50 in Figures 1(a,b,c). On the other hand, FP16 with a 5-bit exponent and 10-bit significand has a greater
 51 precision but a tighter range than BF16 (IEEE (2019)).

52 In this work, we introduce a new low-precision FP format, *EFloat* (*EFn*), that uses entropy-coded
 53 variable-width exponent and a variable-width significand with a total FP bit budget n , as illustrated in
 54 Fig.1(e). Our design is motivated by a key pattern that we observed across a wide range of vector
 55 embedding models: post-training, these models use only few of the $2^8 = 256$ unique exponents
 56 available in FP32 and with a bell-shaped distribution caused by a certain class of non-linear activation
 57 functions used in model training. The EFloat design exploits this behavior and assigns the least
 58 number of exponent bits to most common exponent values, without losing the exponent range of the
 59 original floating point value.

60 The proposed EFloat format has the following benefits:

- 61 • Reduced-bit representation of *any* floating point format (e.g., FP32, FP16), by using fewer
 62 exponent bits to map the **same** exponent range as the original value.
- 63 • For a given bit budget (e.g., 16), EFloat provides more accurate representation of the FP32
 64 values than BF16 and FP16 by using fewer exponent bits to capture the same range as before,
 65 and then using the remaining bits to increase significand precision.
- 66 • The format is suitable for both memory and bandwidth **compression** and reduced-bit
 67 **computations** over *pre-trained* vector embedding models. Software implementation trades

68 compute cycles with capacity and I/O bandwidth savings. A factor of 3 reduction in memory
 69 footprint is achieved converting FP32 values to EF11. Hardware conversion is possible for
 70 FPn to EFn and vice versa with simple Static RAM based lookup tables.

- 71 • For a given dataset, many different FP to EF conversion tables are possible. Tables may be
 72 optimized for maximum significant width (highest exponent compression) at the expense of
 73 worse precision for few floats with infrequent exponents (less significant bits for outliers)
 74 and vice versa.
- 75 • Since vector embedding models are used in a wide array of NLP transformer architectures,
 76 the EFloat format can be used for a much wider (and more space consuming) class of NLP
 77 models.

78 In Section 2, we first present the analysis of various vector embedding models. The EFloat format
 79 is presented in Section 3. Section 4 describes key steps in conversion between EFloat and other FP
 80 formats. Section 5 presents an error analysis of various EFloat widths (EFn) against BF16 and FP16.
 81 In Section 6, a review of related work on model compression and floating-point formats for deep
 82 learning is presented. Finally, Section 7 presents conclusions and outlines future directions.

83 2 Analyzing vector embedding models

84 Vector embedding models are extensively used in natural language processing (NLP) to capture
 85 and exploit semantic relationships of word entities (e.g., words, sentences, phrases, paragraphs, or
 86 documents). A trained vector embedding model consists of a set of vectors, each vector encoding
 87 a *distributed* representation of inferred semantics of a word entity, i.e., a single vector captures
 88 different attributes of the inferred semantics (Hinton et al. (1986)), created, in part, by contributions
 89 by other word entities. Every vector embedding model implements some variant of the *log-bilinear*
 90 language (LBL) model that predicts the probability of the next word w_i given the previous words
 91 (*context*) (Hinton (2013); Almeida and Xexeo (2019); Bender and Koller (2020)). The LBL model
 92 first predicts a real-valued vector representation of a word by *linearly* combining the real-valued vector
 93 representations of its context words. Then the distributed representation of the word is computed
 94 based on the similarity between the predicted representation and the representations of all words in
 95 the vocabulary. This step is accomplished using the *normalized exponential* or *Softmax* function over
 96 the associated feature vectors. The output of the Softmax function is the probability distribution over
 97 V different possible outcomes, where V is the vocabulary size.

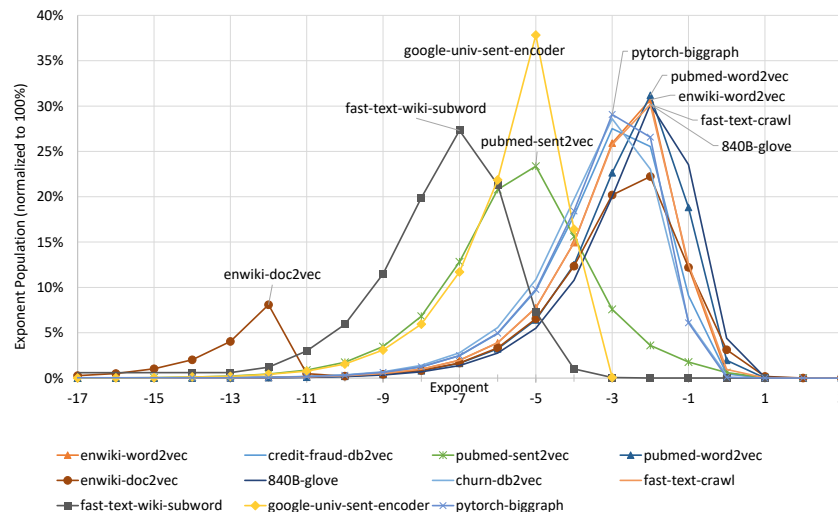


Figure 2: Histogram of the exponent fields of 32-bit floating-point (FP32) values found in vector-embedding and related NLP models. Only the db2Vec, word2vec, doc2vec, and sentence-encoder models were generated. Others were downloaded as publically available pretrained models.

98 Figure 2 presents histograms of exponent values in multiple pre-trained vector embedding models,
 99 where the X-axis represents exponent values (from the 8-bit exponent portion of a 32-bit IEEE 754

100 floating point value). For each X-axis value, the Y-axis represents normalized number of occurrences
 101 of that exponent value, i.e., a histogram. Vector embedding and related NLP models presented in Fig. 2
 102 include word embedding (word2Vec) (Mikolov et al. (2013); Zhang et al. (2019)), sentence (sent2Vec)
 103 and document embedding (doc2Vec) (Le and Mikolov (2014); Chen et al. (2019)), GloVe (Pennington
 104 et al. (2014a)), subword embedding (FastText) (Bojanowski et al. (2017a,b)), database embedding
 105 (db2Vec), graph embedding (PyTorch BigGraph Lerer et al. (2019)), and Google’s transformer-based
 106 universal sentence encoder (Cer et al. (2018); Google (2021)) using the Brown corpus (Brown-
 107 corpus (2021)). All these models implement different variations of the LBL model. The word2Vec
 108 based models, e.g., word2Vec, sent2Vec, doc2Vec, db2Vec, and FastText, use a neural network with
 109 different versions of Softmax as the activation function. GloVe, on the other hand, is a count-based
 110 optimization approach that uses a word co-occurrence matrix and weighted least-square as the
 111 optimization function. The FastText subword model (Joulin et al. (2016); Bojanowski et al. (2017a))
 112 assigns a vector for every character n -gram, using an extended skip-gram model (Mikolov et al.
 113 (2013)) and then, words are represented as the sum of these representations. The universal sentence
 114 encoder generates embedding vectors for sentences using a standard Transformer architecture that
 115 takes word embedding vectors as input and uses a Softmax function to compute attention (Vaswani
 116 et al. (2017)). Irrespective of the model type, we observe that exponent values cluster around a
 117 certain range of values, and display a distinct *peak*. The only exception is the doc2Vec model that
 118 exhibits two peaks as the doc2Vec first builds fine-grained embeddings for words and then uses them
 119 to build embeddings for coarser-grained entities such as paragraphs via concatenating and averaging
 120 individual word vectors which results in a smaller second peak as observed in Fig. 2.

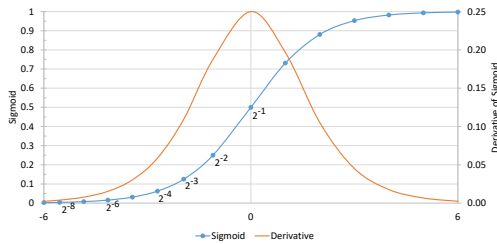


Figure 3: The Sigmoid $\sigma(x)$ curve and its gradient. The floating-point (FP32) exponent of few neural weights are overlaid on $\sigma(x)$.

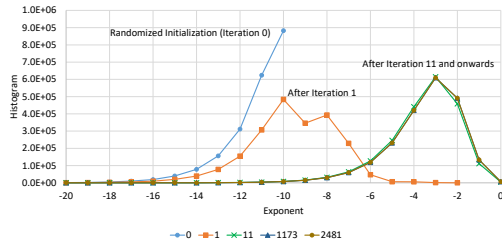


Figure 4: Exponent population (FP32) of model weights during the training of churn-db2vec model, for iterations 0 to through 2481

121 The Softmax family of activations functions used in vector embedding models is responsible for the
 122 clustering behavior of exponents (Figure 2). To understand the reasons, let us delve deeper into the
 123 training of an embedding model. For illustration purposes, we use database embedding (db2Vec)
 124 of the Telcom Churn data (IBM (2020)) as an example. db2Vec is an adaptation of the word2Vec
 125 approach, and has been designed to build an embedding model from structured database tables that
 126 adhere to the relational data model. Like word2Vec, db2Vec also uses Skipgram with Negative
 127 Sampling (SGNS) as the training approach. The SGNS approach uses a binary classifier based on
 128 the logistic (Sigmoid) function instead of using the Softmax-based predictor. The overall training
 129 process involves multiple back-propagation iterations to update model weights using the gradient of
 130 the Sigmoid function. Weights get updated iteratively during the back-propagation process by the
 131 error computed for that iteration. Practically, the error is computed using the gradient of the activation
 132 function. During model training, we observe that the weights rapidly converge (Fig.4) to their final
 133 values. Their exponents are substantially clustered at the slope of the Sigmoid curve, the 2^{-8} to 2^0
 134 output range of Sigmoid, as evidenced by Figures 3 and 4. Training eliminates smaller exponents
 135 from the model because the activation function output is practically zero for any input value when
 136 weights are small, and large exponents are non-existent of normalization of weights.

137 3 The EFloat format: EF_n

138 The key idea behind the EFloat format is the variable-width encoding of exponents using the well-
 139 known Huffman algorithm. Frequency of unique exponent values in the dataset determine the
 140 coded-exponent widths which may vary between as small as 1-bit and some software configurable
 141 maximum, e.g., 8-bit (Figure 1(e,f,g)). Thanks to the entropy coding of the Huffman algorithm,

142 frequent exponent values are coded with fewer bits and infrequent exponents are coded with more
 143 bits as observed in Fig.5.

144 Bits saved from the exponent become part of the significand, therefore increasing the floating point
 145 precision compared to other float formats with the same bit budget. An N -bit coded-exponent in
 146 an EF16 float results in a $(15-N)$ -bit significand as shown in Fig.1(e). Since Efloats with frequent
 147 exponents have wide significands, the entire dataset has a greater precision on average. Efloats with
 148 infrequent exponents have narrow significands. But, their contribution is relatively small in the
 149 common calculations used in model training and inferencing, such as dot-products, vector-sums, and
 150 cosine-similarity (EFloat precision is quantified and compared to prior formats in Section 5.).

151 EFloat on average have greater precision and range than any other fixed-field FP format with the
 152 same bit budget. For example, EF16 with a 3-bit coded-exponent has 12-bits of significand compared
 153 to the 7-bit significand found in a BF16 (Figures 1(h,b)). EFloat exponent's logical width is *always*
 154 8-bit, the same as for FP32 and BF16, irrespective of EFloat width. Even for extremely narrow floats
 155 such as EF8, the logical exponent width can be 8-bit since encoding compresses the exponent field.

156 The EFloat format compresses special values of IEEE 754, such as signed zeros and infinities
 157 losslessly. NaN are semantically compressed losslessly: converting a NaN to and from FP32 to EFn
 158 and vice-versa still results in a NaN. Denormal floats may round to zero since least significant bits of
 159 significands are truncated during encoding.

160 4 EFloat encoding and decoding

161 **The Huffman algorithm:** is a popular lossless compression algorithm used in many compression
 162 tools and compressed data formats (Salomon (2004)). Data symbols are encoded with variable-length
 163 binary codes whose length are determined by the symbol probabilities in the data stream. The
 164 algorithm builds a binary tree with each leaf assigned a symbol. Higher probability symbols are
 165 closer to the tree root than others. The path from the tree root to the leaf is the binary coding of the
 166 symbol. To demonstrate with a trivial example, the letters A, B, C occurring with probabilities of
 167 0.5, 0.25, and 0.25 may be encoded with the bit patterns 0, 10, and 11, respectively. The algorithm
 168 yields 1.5-bit/symbol compression efficiency, better than 8-bits/symbol using an ASCII representation
 169 or 2-bits/symbol using a simplistic mapping of the 3 letters to 2-bit integers. Huffman coding is
 170 optimal when symbol probabilities are negative powers of 2. However, it is an effective compression
 171 method even for non-optimal data distributions. Fig.5 shows the Huffman coded exponent widths as
 172 a function of exponent frequencies of a word2vec trained model.

173 Huffman codes have the *prefix* property which states that no code is a prefix of a longer code (due
 174 its tree structure.) As a result, the Huffman code not only encodes the original symbol but the
 175 code-length as well. We use this property to locate the bit position of the movable boundary between
 176 the exponent and the significand fields when decoding EFloats(Fig.1(e)).

177 **Length-Limiting:** The basic Huffman algorithm, depending on probabilities, may produce extremely
 178 wide codes consuming the entire width of EFloats and more. We use the *Length-Limiting* variant of
 179 the Huffman algorithm to set a maximum coded-exponent width (Abali et al. (2020)). In Fig.5, the
 180 maximum code-width is set to 8-bits resulting in infrequent exponents encoded with that maximum.

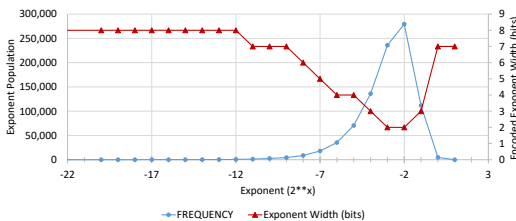


Figure 5: EFloat variable exponent widths are a function of the exponent population (word2vec)

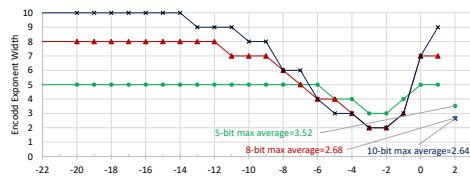


Figure 6: EFloat uses Length-Limited Huffman algorithm to set the maximum width of coded-exponents to 5, 8, or 10-bits

181 The software-defined limit presents an opportunity to tune the EFloat precision to a particular NLP
 182 application requirements. Figure 6 shows the effect of limiting maximum width of coded-exponents to

183 5, 8, and 10-bits. As the limit is increased, the least frequent exponents are coded with the maximum-
184 width codes, therefore their respective significands lose precision. On the other hand, with increasing
185 limits the most frequent exponents are coded with fewer bits reducing both the minimum and the
186 average coded-exponent widths, therefore their respective significands gain precision. Therefore,
187 EFloat not only compresses the regular floats but for a given EFn budget of n -bits the application can
188 optimize the end-to-end precision by adjusting the maximum coded-exponent width.

189 That the minimum and average code widths are inversely proportional to the maximum code-width
190 may appear counterintuitive (Figure 6). Let the maximum coded-exponent width be K -bits ($K = 5$,
191 $K = 8$, etc.). K -bits are sufficient to represent N unique exponents in the dataset provided
192 $\lceil \log_2(N) \rceil \leq K$ holds. When K is chosen much larger than the minimum $\lceil \log_2(N) \rceil$, some codes
193 can have fewer than K -bits. Consider a short L -bit code such that $L \leq K$ (e.g., assume $K = 8$,
194 $L = 2$ and the short-code = 00.) Due to the prefix property the L -bit code consumes 2^{K-L} bit
195 patterns out of the 2^K maximum possible (e.g., all 8-bit patterns whose 2-bit prefix is 00, 64 patterns
196 in total, are consumed by the code 00.) The remaining $N - 1$ exponents can still be encoded if
197 $N - 1 \leq 2^K - 2^{K-L}$ holds. We observed in practice that increasing the maximum code-width by 2
198 to 3 bits over the minimum $\lceil \log_2(N) \rceil$ gives a good compression ratio.

199 **EFloat Encoding and the Code-Table:** During the conversion from FP32 to an EFn (e.g., EF16),
200 exponents in the original dataset are histogrammed first, e.g., Fig.2. The histogram representing
201 probabilities of the exponents is the input to the Length-Limiting Huffman algorithm. The output is a
202 256-entry (2^8) code-table indexed by the original 8-bit exponent. Each table entry contains a pair, the
203 variable-width coded-exponent and its width. Note that the code-table is quite small, tens of bytes in
204 practice, since few unique exponents are present in most NLP datasets as Fig.2 shows.

205 When the sign bits have a skewed distribution, e.g., if they are substantially positive, then the sign
206 bit and the 8-bit exponent may be treated as a single 9-bit integer when histogramming. Thanks to
207 Huffman coding, a skewed sign bit distribution may provide up to one additional bit of precision to
208 the significand.

209 Using the code-table, the entire dataset is converted from FP32 to the chosen EFn width (e.g., EF16)
210 replacing original exponents with coded-exponents. Least significant bits of the FP32 significand are
211 truncated to match the EFn width. For example, in Fig.5, the algorithm encodes the most frequent
212 exponent with 2-bits. Accounting for the sign bit, this yields a 13-bit significand in EF16 by truncating
213 the bottom 10-bit of the 23-bit significand of FP32. We use the *round-to-nearest* method to provide
214 on average 0.5 bits of additional precision: when the leading bit of the truncated part is 1 the upper
215 part of the significand is incremented +1 provided it doesn't overflow in to the exponent field.

216 For large datasets, a statistically representative subset may also be used to reduce histogram collection
217 time. When the histogram is known in advance, a pre-built code-table may be used. Pre-built
218 code-tables eliminate the overhead of executing the Huffman algorithm. During training exponents
219 rapidly converge to their final values as observed in Fig.4. The exponent distribution is practically
220 identical for all iterations 11 to 2481, Therefore, a single pre-built code-table optimized for final
221 iterations may serve for all iterations start to finish. The same pre-built table, although suboptimal
222 for early iterations, may be used because significand precision is not as important at that point in
223 time; model weight updates are dominated by exponent updates. Once exponents settled to their final
224 values the significand precision becomes important since model weights updates progressively get
225 smaller.

226 **EFloat Decoding:** For EFloat to FP32 conversion (i.e., decoding) we use a inverse mapping of the
227 code-table described earlier. A decoder-table indexed by the coded-exponent may be used to decode
228 the original exponent value and the significand's leading bit position in constant time. Each table entry
229 contains the original exponent and width of the coded-exponent. To index the decoder-table with
230 variable-width codes many entries are filled with duplicates. For example, the 2-bit coded-exponent
231 00 is duplicated 64 times in the table at locations 00000000 through 00111111 with each location
232 containing the pair (original exponent and code-width= 2). Duplicating entries is equivalent to having
233 *logical don't care* bits in the index which is a useful in hardware based decoder implementations.

234 The second element of each table entry contains the EFloat significand width. Since the significand
235 was truncated earlier during the FP32 to EFloat conversion, the missing least significant bits must be
236 padded with zeros to match the original FP32 width.

237 **5 Evaluating the EFloat representation**

238 In this section, we evaluate the efficacy of the EFloat format using two sets of experiments. The
 239 first set measures the loss of precision in representing FP32 data in BF16, FP16, and EFloat formats
 240 with bit budgets from 16 to 8 bits. The second set of experiments compares the quality of ranked
 241 results for *similarity* and *dissimilarity* queries using the *Normalized Discounted Cumulative Gain*
 242 (*NDCG*) score for BF16, FP16, and various EFloat formulations. Table 1 presents the list of models
 243 used in these experiments, along with their characteristics: model types, model size (stored using
 244 FP32), number of unique exponents, range of exponent bits generated by the Huffman algorithm, the
 245 average count of exponent bits, and minimum and maximum average count of significant bits. For
 246 EF16, the average significand length is 4.3 bits more than BF16 (with 7-bit significand) and 1.2 bits
 247 more than FP16 (with 10-bit significand).

Table 1: EFloat characteristics from EF16 to EF8 for different datasets

Model	Type	Size	Unique exponents	EFn exponent bits			EFn significand bits (Avg.)	
				Min	Max	Avg.	Max (EF16)	Min (EF8)
churn	db2vec	20 MB	23	3	5	3.6	11.4	3.4
crawl	fast-text	4.3 GB	30	3	5	3.4	11.6	3.6
enwiki	word2vec	9.6 GB	27	4	5	4.2	10.8	4.8
MDM	db2vec	14 GB	24	3	6	3.6	11.4	3.4
840B	GloVe	5.3 GB	35	3	6	3.5	11.5	3.5
wiki-sw	fast-text	2.2 GB	22	3	5	3.6	10.5	3.4
virginia	db2Vec	222 MB	24	3	5	3.7	11.3	3.4

248 The first set of experiments compares the loss of precision due to the least significant significand
 249 bits being truncated during conversion from FP32 to various lower-precision formats. Given a
 250 low-precision format (e.g., EF16 or BF16), the values are converted back to FP32, and the arithmetic
 251 difference, $f^o - f^c$, of the original FP32 value, f^o , and the regenerated FP32 value, f^c , is computed.
 252 This difference represents the precision loss due to conversion. Root Mean Square Error (RMSE)
 253 metric is then used to summarize the loss of precision across a dataset of N floats as:

$$RMSE = \sqrt{\frac{1}{N} \sum_k (f_k^o - f_k^c)^2}$$

254 We then compare the errors of BF16/FP16 and EFn by dividing $RMSE_{BF16/FP16}$ by $RMSE_{EFn}$
 255 in Table 2. Ratios greater than 1.0 indicate that the EFloat error is less than BF16 or FP16 errors.
 256 For EF16, across all models, we observe an average RMSE error ratio of 24.1 for BF16, and 1.5 for
 257 FP16. Note that for these experimental results, the datasets were encoded with a minimum of 3-bit
 258 and a maximum of 6-bit coded-exponents resulting in an average width in the range of 3.4 to 4.2-bits
 259 (Table 1). Accordingly, for EF16, the *minimum* significand width is 10-bit which is 3-bit wider than
 260 BF16, and of the same length as FP16. Therefore, EF16 has significantly higher precision against
 261 BF16 than FP16. Also, Table 2 shows that EF12 has the same to slightly better RMSE than BF16
 262 since the RMSE ratios are in the 1.0 to 2.2 range. Thus, EF12 uses 25% less bandwidth and memory
 263 capacity than BF16 for similar floating-point precision.

264 Note that the RMSE method amplifies larger errors due to the squaring of differences. EFloat coded
 265 floating point values with short significands (i.e., those with infrequent exponents) are disproportion-
 266 ately represented in the RMSE summation. However, the true measure of error for vector embedding
 267 models will be the evaluation of ranked results for similarity queries for different floating point
 268 formats. Unlike the binning in traditional classification inference tasks, ranked results from similarity
 269 queries are far more sensitive to numerical precision. We use the *Normalized Discounted Cumulative*
 270 *Gain (NDCG)* metric (Järvelin and Kekäläinen (2002); Wang et al. (2013)), to evaluate the quality
 271 of ranked results for different floating point formats. NDCG is widely used in information retrieval
 272 and web search to evaluate the relevance of retrieved documents. NDCG is a normalization of the
 273 Discounted Cumulative Gain (DCG) measure. DCG is calculated as a weighted sum of the degree of
 274 relevancy of the ranked items, where the weight is a decreasing function of the position of an item.
 275 NDCG is computed by normalizing DCG by IDCg, which is the DCG measure for a perceived ideal
 276 ranking result. Thus, the NDCG measure always lies within [0.0,1.0].

Table 2: BFloat16 (BF16), IEEE Half (FP16), and EF16–8 precision comparisons using RMSE-with-FP32 ratio. Higher is better.

Model	EF16		EF14		EF12		EF10		EF8	
	BF16	FP16	BF16	FP16	BF16	FP16	BF16	FP16	BF16	FP16
churn	22.5	1.4	5.6	0.4	1.4	0.09	0.3	0.02	0.08	0.005
crawl	34.6	2.2	8.6	0.5	2.2	0.1	0.5	0.03	0.1	0.008
enwiki	16.9	1.0	4.2	0.3	1.0	0.07	0.3	0.02	0.06	0.004
MDM	27.9	1.8	6.9	0.4	1.8	0.1	0.4	0.03	0.1	0.007
840B	25.0	1.6	6.3	0.4	1.6	0.09	0.4	0.02	0.09	0.006
wiki-sw	22.0	1.4	5.5	0.3	1.2	0.08	0.3	0.02	0.08	0.005
virginia	19.6	1.2	4.9	0.3	1.2	0.08	0.3	0.02	0.07	0.004

277 For a given vector embedding model, we choose $q = 20$ randomly selected distinct query points.
 278 For each query point, we compute similar and dissimilar points by computing cosine similarities
 279 over the corresponding vectors. For similarity queries, the result contains a list of points sorted in
 280 decreasing order of their similarity scores (most similar pair of items will have score closer to 1.0),
 281 and for dissimilarity queries, the result list is sorted in increasing order of their similarity scores (most
 282 dissimilar pair of items will have score closer to -1.0). For each query point, we run similarity and
 283 dissimilarity queries for different floating point formats, and use the top $k = 10$ results for each test
 284 to compute the NDCG score, (**NDCG@10**). In our evaluation, we use the ranked results for FP32
 285 as the baseline for calculating the IDCG. For each model, we report the average NDCG@10 score
 286 computed over 20 query points using BF16, FP16, and various EFn from EF16 to EF8.

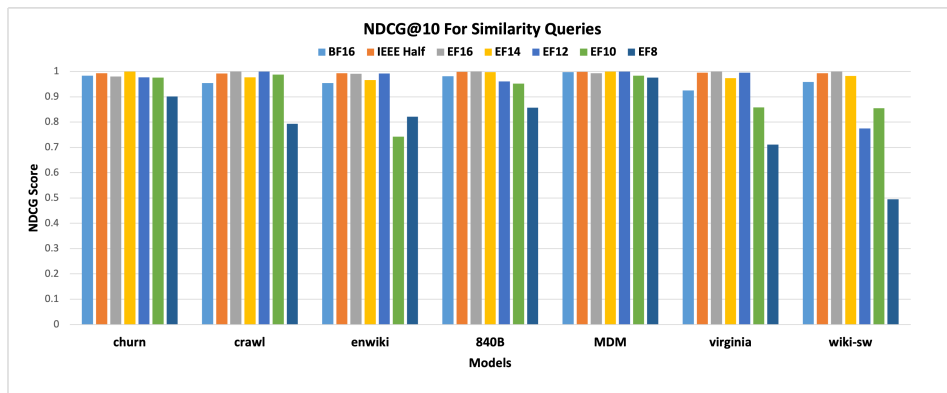


Figure 7: Evaluation of similarity query accuracy using NDCG score across different floating point formats. Higher score (closer to 1.0) is better.

287 Figure 7 presents NDCG10 results for similarity queries, and Figure 8 presents NDCG10 results for
 288 dissimilarity queries. For both similarity and dissimilarity queries, EF16 matches or exceeds the
 289 quality of BF16 or FP16 (in particular, among the three formats, BF16 provides the worst quality
 290 results). Furthermore, EF14 and EF12 provide similar quality results as EF16 in many instances. The
 291 two lower-precision EFn, EF10 and EF8, consistently generate the least quality results.

292 In summary, results from the two sets of experiments (Table 2, and Figures 7 and 8), conclusively
 293 demonstrate that: (1) Given a bit budget, EFloat has higher accuracy than other formats, (2) In many
 294 scenarios, EFn with reduced bit budget (e.g., EF14 or EF12) provides results of quality comparable to
 295 higher precision formats, e.g., BF16, and FP16. These results validate the design of the EFloat format,
 296 and demonstrate that EFloats can be used for compressing and computing using vector embedding
 297 models.

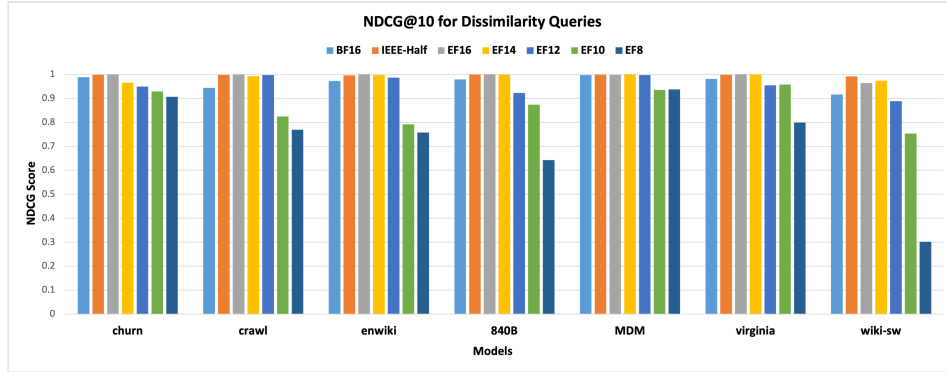


Figure 8: Evaluation of dissimilarity query accuracy using NDCG score across different floating point format. Higher score (closer to 1.0) is better.

298 6 Related Work

299 Model quantization is widely used to compress *pre-trained* models for the inference phase (Gupta
300 and Agrawal (2020)). Quantization covers two broad approaches: the first represents a full-precision
301 (e.g., 32-bit) floating point weight value using reduced (e.g., BF16 or FP16) or mixed precision
302 floats, and the second converts full-precision floating point values into integer values with fewer bits
303 (e.g., INT8, INT4, and INT1 (Migacz (2017); Wu et al. (2020); TensorFlow Documentation (2020);
304 Jacob et al. (2017))). In conjunction with the model compression work, there has been significant
305 work in devising reduced-precision floating point formats tuned for broader machine learning and
306 HPC applications (Sapunov (2020); Abdelfattah et al. (2020)). Unlike the inference-focused model
307 compression work, reduced-precision floating points are designed to work for both model training and
308 inference phases. The most common reduced-precision floating point formats use 16 bits. Current
309 16-bit implementations include IEEE 754 half-precision (FP16); Brain Floating Point, BFloat16
310 (Wang and Kumar (2019); Kalamkar et al. (2019)); and Deep Learning Float (DLFloat) (Agrawal
311 et al. (2019)), with 1 sign bit, 6 exponent bits, and 7 fraction bits. TensorFloat-32 (TF32) from
312 Nvidia is a 19-bit format that combines 8 exponent bits from BFLOAT16 and 10 exponent bits from
313 IEEE FP16 (Kharya (2020)). Hybrid Block Floating Point (HBFP) (Drumond et al. (2018)), Intel
314 Nervana’s Flexpoint (Koster et al. (2017)), and Microsoft MSFP (Rouhani et al. (2020)) formats
315 combine the advantages of fixed point and floating point formats by splitting up the significand and
316 the exponent part which is shared across multiple numeric values. Recent research proposals have
317 described training of key deep learning models using even reduced precision floating point values
318 (8- and 4-bit floats) (Sun et al. (2020); Wang et al. (2018); Cambier et al. (2020); Mellempudi et al.
319 (2019)). Recently proposed AdaptiveFloat (Tambe et al. (2020)) is an inference-targeted floating-point
320 format which maximizes its dynamic range at a network layer granularity by dynamically shifting
321 its exponent range via modifications to the exponent bias and by optimally clipping (quantizing) its
322 representable datapoints. Our proposed EFloat design practically achieves the same result without
323 altering the exponent range and quantizing full-precision values.

324 7 Conclusion

325 We introduced EFloat, a novel entropy-coded variable length floating point format for deep learning
326 applications. This format can be used for compressing a trained deep learning model, as well as for
327 enabling more accurate model representations using reduced-precision floating point formats. While
328 our intended use cases were initially for the database embedding (db2Vec) workloads, we demonstrate
329 that the proposed format works effectively for other vector embedding models, and can be used for a
330 much broader class of NLP models including transformer-based models. Broadly, EFloat may be
331 used in deep learning applications where tradeoffs need to be made between range, precision, memory
332 capacity and bandwidth savings. As a future work, we plan to explore the Benford distribution
333 pattern (Benford (1938); Newcomb (1881)) exhibited by significands of vector embedding models
334 (Appendix A in the supplementary document) and investigate its application in rounding EFloat
335 values. A follow-up study on 8-bit floats and integers is being considered as well.

References

- 336
- 337 Bulent Abali, Bart Blaner, John Reilly, Matthias Klein, Ashutosh Mishra, Craig B Agricola, Bedri
338 Sendir, Alper Buyuktosunoglu, Christian Jacobi, William J Starke, et al. 2020. Data Compression
339 Accelerator on IBM POWER9 and z15 Processors. In *2020 ACM/IEEE 47th Annual International*
340 *Symposium on Computer Architecture (ISCA)*. IEEE Computer Society, 1–14.
- 341 Ahmad Abdelfattah, Hartwig Anzt, Erik G. Boman, Erin Carson, Terry Cojean, Jack Dongarra, Mark
342 Gates, Thomas Grützmacher, Nicholas J. Higham, Sherry Li, Neil Lindquist, Yang Liu, Jennifer
343 Loe, Piotr Luszczek, Pratik Nayak, Sri Pranesh, Siva Rajamanickam, Tobias Ribizel, Barry Smith,
344 Kasia Swirydowicz, Stephen Thomas, Stanimire Tomov, Yaohung M. Tsai, Ichitaro Yamazaki, and
345 Urike Meier Yang. 2020. A Survey of Numerical Methods Utilizing Mixed Precision Arithmetic.
346 arXiv:2007.06674 [cs.MS]
- 347 A. Agrawal, S. M. Mueller, B. M. Fleischer, X. Sun, N. Wang, J. Choi, and K. Gopalakrishnan. 2019.
348 DLFloat: A 16-b Floating Point Format Designed for Deep Learning Training and Inference. In
349 *2019 IEEE 26th Symposium on Computer Arithmetic (ARITH)*. 92–95.
- 350 Felipe Almeida and Geraldo Xexeo. 2019. Word Embeddings: A Survey. arXiv:1901.09069 [cs.CL]
- 351 Emily M. Bender and Alexander Koller. 2020. Climbing towards NLU: On Meaning, Form, and
352 Understanding in the Age of Data. In *Proceedings of the 58th Annual Meeting of the Association*
353 *for Computational Linguistics*. Association for Computational Linguistics, Online, 5185–5198.
354 <https://doi.org/10.18653/v1/2020.acl-main.463>
- 355 Frank Benford. 1938. The law of anomalous numbers. In *Proc. American Philosophical Society*,
356 Vol. 78. 551–572.
- 357 Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017a. Enriching Word
358 Vectors with Subword Information. *Transactions of the Association for Computational Linguistics*
359 5 (2017), 135–146.
- 360 Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017b. En-
361 riching Word Vectors with Subword Information. arXiv:1607.04606 [cs.CL]
362 <https://fasttext.cc/docs/en/english-vectors.html> MIT License.
363 <https://github.com/facebookresearch/fastText/blob/master/LICENSE>.
- 364 Rajesh Bordawekar and Oded Shmueli. 2016. Enabling Cognitive Intelligence Queries in Relational
365 Databases using Low-dimensional Word Embeddings. *CoRR* abs/1603.07185 (2016). <http://arxiv.org/abs/1603.07185>
- 366
- 367 Rajesh Bordawekar and Oded Shmueli. 2017. Using Word Embedding to Enable Semantic Queries
368 in Relational Databases. In *Proceedings of the 1st Workshop on Data Management for End-to-End*
369 *Machine Learning (Chicago, IL, USA) (DEEM'17)*. ACM, New York, NY, USA, Article 5, 4 pages.
370 <https://doi.org/10.1145/3076246.3076251>
- 371 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-
372 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,
373 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
374 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin,
375 Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford,
376 Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. *CoRR*
377 abs/2005.14165 (2020). arXiv:2005.14165 <https://arxiv.org/abs/2005.14165>
- 378 Brown-corpus 2021. Brown Corpus. [http://www.sls.hawaii.edu/bley-vroman/brown_](http://www.sls.hawaii.edu/bley-vroman/brown_nolines.txt)
379 [nolines.txt](http://www.sls.hawaii.edu/bley-vroman/brown_nolines.txt)
- 380 Leopold Cambier, Anahita Bhiwandiwala, Ting Gong, Mehran Nekuii, Oguz H Elibol, and Hanlin
381 Tang. 2020. Shifted and Squeezed 8-bit Floating Point format for Low-Precision Training of Deep
382 Neural Networks. arXiv:2001.05674 [cs.LG]

383 Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah
384 Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Brian Strope, and Ray Kurzweil.
385 2018. Universal Sentence Encoder for English. In *Proceedings of the 2018 Conference on*
386 *Empirical Methods in Natural Language Processing: System Demonstrations*. Association for
387 Computational Linguistics, Brussels, Belgium, 169–174. [https://doi.org/10.18653/v1/](https://doi.org/10.18653/v1/D18-2029)
388 D18-2029 Apache-2.0 License.

389 Qingyu Chen, Yifan Peng, and Zhiyong Lu. 2019. BioSentVec: creating sentence embed-
390 dings for biomedical texts. In *The 7th IEEE International Conference on Healthcare Infor-*
391 *matics*. <https://github.com/ncbi-nlp/BioSentVec> License: [https://github.com/ncbi-](https://github.com/ncbi-nlp/BioSentVec/blob/master/LICENSE.txt)
392 [nlp/BioSentVec/blob/master/LICENSE.txt](https://github.com/ncbi-nlp/BioSentVec/blob/master/LICENSE.txt).

393 C. J. Date. 1982. A Formal Definition of the relational Model. *SIGMOD Rec.* 13, 1 (1982), 18–29.
394 <https://doi.org/10.1145/984514.984515>

395 Mario Drumond, Tao Lin, Martin Jaggi, and Babak Falsafi. 2018. Training DNNs with Hybrid
396 Block Floating Point. In *Proceedings of the 32nd International Conference on Neural Information*
397 *Processing Systems (NIPS'18)*. Curran Associates Inc., Red Hook, NY, USA, 451–461.

398 William Fedus, Barret Zoph, and Noam Shazeer. 2021. Switch Transformers: Scaling to Trillion
399 Parameter Models with Simple and Efficient Sparsity. arXiv:2101.03961 [cs.LG]

400 Google. 2021. TensorFlow Universal Sentence Encoder. [https://tfhub.dev/google/](https://tfhub.dev/google/universal-sentence-encoder/4)
401 [universal-sentence-encoder/4](https://tfhub.dev/google/universal-sentence-encoder/4) Apache-2.0 License.

402 Manish Gupta and Puneet Agrawal. 2020. Compression of Deep Learning Models for Text: A Survey.
403 arXiv:2008.05221 [cs.CL]

404 Geoffrey Hinton. 2013. Learning Distributed Representations For Statistical Language Modelling.
405 <https://www.cs.toronto.edu/~hinton/csc2535/note/h1b1.pdf>.

406 G. E. Hinton, J. L. McClelland, and D. E. Rumelhart. 1986. Distributed representations. *Parallel*
407 *distributed processing: explorations in the microstructure of cognition 1* (1986), 77–109.

408 Torsten Hoefler. 2020. OpenAI's GPT-3 took more than 314 Zettaflop (mixed fp16/fp32) to train!
409 <https://twitter.com/thoefler/status/1268291652701077504?s=20>.

410 IBM. 2020. Predict Customer Churn using Watson Machine Learning and Jupyter Notebooks
411 on Cloud Pak for Data. <https://github.com/IBM/telco-customer-churn-on-icp4d>
412 Apache-2.0 License.

413 IEEE. 2019. "IEEE Standard for Floating-Point Arithmetic". , 84 pages. [https://doi.org/10.](https://doi.org/10.1109/IEEESTD.2019.8766229)
414 [1109/IEEESTD.2019.8766229](https://doi.org/10.1109/IEEESTD.2019.8766229)

415 Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig
416 Adam, and Dmitry Kalenichenko. 2017. Quantization and Training of Neural Networks for
417 Efficient Integer-Arithmetic-Only Inference. arXiv:1712.05877 [cs.LG]

418 Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated Gain-Based Evaluation of IR Techniques.
419 *ACM Transactions on Information Systems* 20, 4 (2002). [https://doi.org/10.1145/582415.](https://doi.org/10.1145/582415.582418)
420 582418

421 Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016. Bag of Tricks for
422 Efficient Text Classification. arXiv:1607.01759 [cs.CL]

423 Dhiraj Kalamkar, Dheevatsa Mudigere, Naveen Mellempudi, Dipankar Das, Kunal Banerjee,
424 Sasikanth Avancha, Dharma Teja Vooturi, Nataraj Jammalamadaka, Jianyu Huang, Hector Yuen,
425 Jiyan Yang, Jongsoo Park, Alexander Heinecke, Evangelos Georganas, Sudarshan Srinivasan, Ab-
426 hisek Kundu, Misha Smelyanskiy, Bharat Kaul, and Pradeep Dubey. 2019. A Study of BFLOAT16
427 for Deep Learning Training. arXiv:1905.12322 [cs.LG]

428 Paresh Kharya. 2020. TensorFloat-32 in the A100 GPU Accelerates AI Training, HPC up to 20x.
429 Nvidia Blog.

- 430 Urs Koster, Tristan J. Webb, Xin Wang, Marcel Nassar, Arjun K. Bansal, William H. Constable,
431 Oguz H. Elibol, Scott Gray, Stewart Hall, Luke Hornof, Amir Khosrowshahi, Carey Kloss, Ruby J.
432 Pai, and Naveen Rao. 2017. Flexpoint: An Adaptive Numerical Format for Efficient Training of
433 Deep Neural Networks. arXiv:1711.02213 [cs.LG]
- 434 Quoc V. Le and Tomas Mikolov. 2014. Distributed Representations of Sentences and Documents.
435 *CoRR* abs/1405.4053 (2014). <http://arxiv.org/abs/1405.4053>
- 436 Adam Lerer, Ledell Wu, Jiajun Shen, Timothee Lacroix, Luca Wehrstedt, Abhijit Bose, and
437 Alex Peysakhovich. 2019. PyTorch-BigGraph: A Large-scale Graph Embedding System.
438 arXiv:1903.12287 [cs.LG] <https://github.com/facebookresearch/PyTorch-BigGraph>
439 BSD License. <https://github.com/facebookresearch/PyTorch-BigGraph/blob/master/LICENSE.txt>.
- 440 Naveen Mellempudi, Sudarshan Srinivasan, Dipankar Das, and Bharat Kaul. 2019. Mixed Precision
441 Training With 8-bit Floating Point. arXiv:1905.12334 [cs.LG]
- 442 Szymon Migacz. 2017. 8-bit Inference with TensorRT. [https://on-demand.gputechconf.com/
443 gtc/2017/presentation/s7310-8-bit-inference-with-tensorrt.pdf](https://on-demand.gputechconf.com/gtc/2017/presentation/s7310-8-bit-inference-with-tensorrt.pdf).
- 444 Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient Estimation of Word
445 Representations in Vector Space. *CoRR* abs/1301.3781 (2013). [http://arxiv.org/abs/1301.
446 3781](http://arxiv.org/abs/1301.3781)
- 447 Simon Newcomb. 1881. Note on the frequency of use of the different digits in natural numbers.
448 *American Journal of Mathematics* 4, 1/4 (1881), 39–40.
- 449 Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014a. GloVe: Common Crawl
450 840B Dataset. <https://nlp.stanford.edu/projects/glove> License: Open Data Commons
451 Public Domain Dedication and License.
- 452 Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014b. GloVe: Global Vectors for
453 Word Representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural
454 Language Processing*. 1532–1543. <http://aclweb.org/anthology/D/D14/D14-1162.pdf>
- 455 Bitia Rouhani, Daniel Lo, Ritchie Zhao, Ming Liu, Jeremy Fowers, Kalin Ovtcharov, Anna Vino-
456 gradsky, Sarah Massengill, Lita Yang, Ray Bittner, Alessandro Forin, Haishan Zhu, Taesik Na,
457 Prerak Patel, Shuai Che, Lok Chand Koppaka, Xia Song, Subhojit Som, Kaustav Das, Saurabh
458 Tiwary, Steve Reinhardt, Sitaram Lanka, Eric Chung, and Doug Burger. 2020. Pushing the Limits
459 of Narrow Precision Inferencing at Cloud Scale with Microsoft Floating Point. In *NeurIPS 2020*.
460 ACM.
- 461 David Salomon. 2004. *Data compression: the complete reference*. Springer Science & Business
462 Media.
- 463 Grigory Sapunov. 2020. FP64, FP32, FP16, BFLOAT16, TF32, and other members of the ZOO.
464 [https://moocaholic.medium.com/fp64-fp32-fp16-bfloat16-tf32-and-other-members-
465 a1ca7897d407](https://moocaholic.medium.com/fp64-fp32-fp16-bfloat16-tf32-and-other-members-of-the-zoo-a1ca7897d407).
- 466 Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan
467 Catanzaro. 2020. Megatron-LM: Training Multi-Billion Parameter Language Models Using Model
468 Parallelism. arXiv:1909.08053 [cs.CL]
- 469 Xiao Sun, Naigang Wang, Chia-yu Chen, Jia-min Ni, Ankur Agrawal, Xiaodong Cui, Swagath
470 Venkararamini, Kaoutar El Maghraoui, and Vijayalakshmi Srinivasan. 2020. Ultra-Low Precision
471 4-bit Training of Deep Neural Networks. In *Advances in Neural Information Processing Systems*,
472 Vol. 33.
- 473 T. Tambe, E. Y. Yang, Z. Wan, Y. Deng, V. Janapa Reddi, A. Rush, D. Brooks, and G. Y. Wei.
474 2020. Algorithm-Hardware Co-Design of Adaptive Floating-Point Encodings for Resilient Deep
475 Learning Inference. In *2020 57th ACM/IEEE Design Automation Conference (DAC)*. 1–6. [https:
476 //doi.org/10.1109/DAC18072.2020.9218516](https://doi.org/10.1109/DAC18072.2020.9218516)
- 477 TensorFlow Documentation. 2020. TensorFlow Lite 8-bit quantization specification. [https:
478 //www.tensorflow.org/lite/performance/quantization_spec](https://www.tensorflow.org/lite/performance/quantization_spec)

- 479 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz
480 Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. arXiv:1706.03762 [cs.CL]
- 481 Naigang Wang, Jungwook Choi, Daniel Brand, Chia-Yu Chen, and Kailash Gopalakrishnan. 2018.
482 Training Deep Neural Networks with 8-bit Floating Point Numbers. In *Advances in Neural*
483 *Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-
484 Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc., 7675–7684.
- 485 Shibo Wang and Pankaj Kumar. 2019. BFloat16: The secret to high performance on Cloud TPUs.
- 486 Yining Wang, Liwei Wang, Yuanzhi Li, Di He, Tie-Yan Liu, and Wei Chen. 2013. A Theoretical
487 Analysis of NDCG Type Ranking Measures. *CoRR* abs/1304.6480 (2013). [http://arxiv.org/](http://arxiv.org/abs/1304.6480)
488 [abs/1304.6480](http://arxiv.org/abs/1304.6480)
- 489 Hao Wu, Patrick Judd, Xiaojie Zhang, Mikhail Isaev, and Paulius Micikevicius. 2020. In-
490 teger Quantization for Deep Learning Inference: Principles and Empirical Evaluation.
491 arXiv:2004.09602 [cs.LG]
- 492 Yijia Zhang, Qingyu Chen, Zhihao Yang, Hongfei Lin, and Zhiyong Lu. 2019. BioWordVec,
493 improving biomedical word embeddings with subword information and MeSH. *Scientific Data* 6,
494 52 (2019). <https://github.com/ncbi-nlp/BioSentVec> License: [https://github.com/ncbi-](https://github.com/ncbi-nlp/BioSentVec/blob/master/LICENSE.txt)
495 [nlp/BioSentVec/blob/master/LICENSE.txt](https://github.com/ncbi-nlp/BioSentVec/blob/master/LICENSE.txt).

496 **Checklist**

- 497 1. For all authors...
- 498 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
499 contributions and scope? [Yes]
- 500 (b) Did you describe the limitations of your work? [Yes]
- 501 (c) Did you discuss any potential negative societal impacts of your work? [N/A] This work
502 does not deal with societal impact of AI.
- 503 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
504 them? [Yes]
- 505 2. If you are including theoretical results...
- 506 (a) Did you state the full set of assumptions of all theoretical results? [N/A] Experimental
507 work.
- 508 (b) Did you include complete proofs of all theoretical results? [N/A] Experimental work.
- 509 3. If you ran experiments...
- 510 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
511 mental results (either in the supplemental material or as a URL)? [No] The codes used
512 in experiments are proprietary.
- 513 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
514 were chosen)? [No] We used pre-trained models
- 515 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
516 ments multiple times)? [No] Experiments used pre-trained models.
- 517 (d) Did you include the total amount of compute and the type of resources used (e.g.,
518 type of GPUs, internal cluster, or cloud provider)? [No] Experiments used pre-trained
519 models.
- 520 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 521 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 522 (b) Did you mention the license of the assets? [Yes] License information is provided in the
523 citation.
- 524 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
525 Information provided in the citation.
- 526 (d) Did you discuss whether and how consent was obtained from people whose data you're
527 using/curating? [N/A]
- 528 (e) Did you discuss whether the data you are using/curating contains personally identifiable
529 information or offensive content? [N/A]
- 530 5. If you used crowdsourcing or conducted research with human subjects...
- 531 (a) Did you include the full text of instructions given to participants and screenshots, if
532 applicable? [N/A]
- 533 (b) Did you describe any potential participant risks, with links to Institutional Review
534 Board (IRB) approvals, if applicable? [N/A]
- 535 (c) Did you include the estimated hourly wage paid to participants and the total amount
536 spent on participant compensation? [N/A]