

(2024-09-01; PRE-ARXIV) EFFICIENT 1-BIT TENSOR APPROXIMATIONS

ALEX W. N. RIASANOVSKY AND SARAH EL KAZDADI

ABSTRACT. We present a spatially efficient decomposition of matrices and arbitrary-order tensors as linear combinations of tensor products of $\{-1, 1\}$ -valued vectors. For any matrix $A \in \mathbb{R}^{m \times n}$,

$$A - R_w = S_w C_w T_w^\top = \sum_{j=1}^w c_j \cdot \mathbf{s}_j \mathbf{t}_j^\top$$

is a w -width signed cut decomposition of A . Here $C_w = \text{diag}(\mathbf{c}_w)$ for some $\mathbf{c}_w \in \mathbb{R}^w$, and S_w, T_w , and the vectors $\mathbf{s}_j, \mathbf{t}_j$ are $\{-1, 1\}$ -valued.

To store (S_w, T_w, C_w) , we may pack $w \cdot (m + n)$ bits, and require only w floating (or fixed-) point numbers. As a function of w , $\|R_w\|_F$ decays exponentially when applied to $\mathbf{f32}$ matrices with i.i.d. $\mathcal{N}(0, 1)$ entries. Choosing w so that (S_w, T_w, C_w) has the same memory footprint as a $\mathbf{f16}$ or $\mathbf{bf16}$ matrix, the error is comparable.

As a first application, we approximate the weight matrices in the open **Mistral-7B-v0.1** Large Language Model to a 50% spatial compression. Remarkably, all 226 remainder matrices have a relative error $< 6\%$ and the expanded model closely matches **Mistral-7B-v0.1** on the **HuggingFace** leaderboard. Benchmark performance degrades slowly as we reduce the spatial compression from 50% to 25%.

Our algorithm yields efficient signed cut decompositions in 20 lines of pseudocode. It reflects a simple modification from a 1999 paper of Frieze and Kannan. This paper, touted as a foundation of the field of combinatorial limit theory, appears underutilized in randomized numerical linear algebra.

We optimize our open source **Rust** implementation with **SIMD** instructions on **avx2** and **avx512** architectures. We also extend our algorithm from matrices to tensors of arbitrary order and use it to compress a picture of the first author’s cat Angus.

1. INTRODUCTION

In a recent survey [1], Murray et al offer RNLA as a treatment for large-scale problems in high performance computing (HPC) and machine learning (ML).

A dire situation. While communities that rely on NLA now vary widely, they share one essential property: a ravenous appetite for solving larger and larger problems.

— Murray et al [1]

Throughout the survey, the authors emphasize the emergence of structure at scale, which randomized algorithms are uniquely adept at exploiting.

The story in combinatorial limit theory is similar. Here, the *cut norm* $\|\cdot\|_\square$, first coined in 1999 in [2] by Frieze and Kannan, plays a vital role in exposing emergent structure in large networks. Using the cut norm, they approximate matrices with linear combinations of rank-1 matrices with entries in $\{0, 1\}$:

$$\|A\|_{\square} := \max_{S,T} \left| \sum_{(i,j) \in S \times T} a_{ij} \right| \rightsquigarrow A \approx \sum_j d_j \mathbf{1}_{S_j \times T_j}.$$

Following [2], researchers at the Theory Group of Microsoft Research applied these so-called *Frieze-Kannan decompositions* to connect several central, but superficially different, topics in combinatorics.

Although Frieze-Kannan decompositions are spatially efficient, it is challenging to efficiently minimize their errors (c.f. Section 2.1). As a remedy, we base our methods around the signed cut norm $\|\cdot\|_{\blacksquare}$ instead of $\|\cdot\|_{\square}$, and instead approximate with outer products of $\{-1, 1\}$ -valued vectors:

$$\|A\|_{\blacksquare} := \max_{s,t} \langle sA, t \rangle \rightsquigarrow A \approx \sum_j c_j s_j t_j^{\top}.$$

Summarizing Section 4, we argue that these *signed cut decompositions* produce state-of-the-art approximations. As noted in Section 3, the main algorithms we use to find them fit in 20 lines of pseudocode.

1.1. Comparison to bf16 quantization. Our signed cut decompositions trade accuracy for space at a better exchange rate than **bf16** quantization offers for larger matrices, as shown in Figure 1. Moreover, we may make this tradeoff dynamically at runtime to suit the needs of our application.

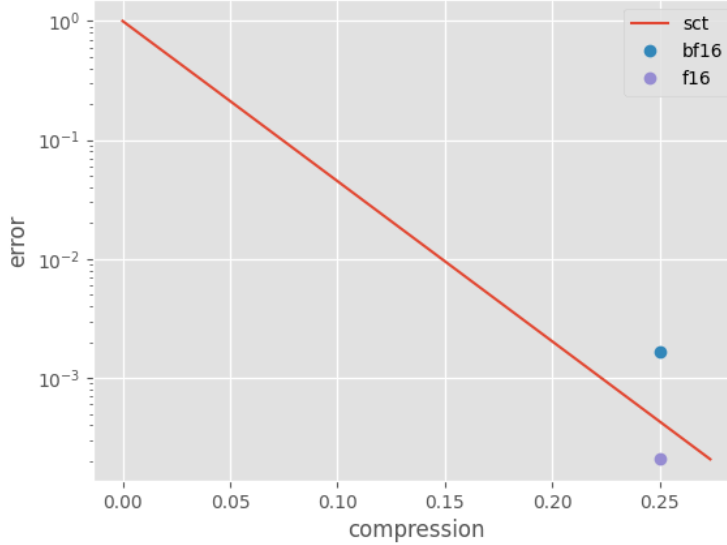


FIGURE 1. Relative error (log-scaled) and relative size of a signed cut decompositions of a 4096×4096 **f64** matrix with i.i.d. standard normal entries. For comparison, we note the tradeoff made with **f16** and **bf16**.

1.2. Approximating Mistral-7B-v0.1. When the original matrix is more structured, signed cut decompositions converge faster. In Figure 2, we show the effect of

signed cut decompositions on the open **Mistral-7B-v0.1** model’s benchmark performance.

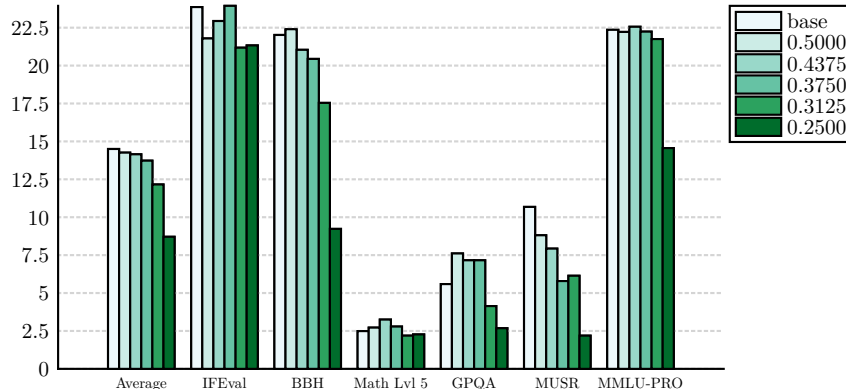


FIGURE 2. Benchmark results formed by approximating the weight matrices of **Mistral-7B-v0.1** with signed cut decompositions and expanding them at different widths. Here, 0.500 reflects a 2-fold spatial compression, and 0.2500 reflects a 4-fold spatial compression.

1.3. Higher order tensors. All of our techniques extend from matrices to tensors of arbitrary order, as show in Section 5. In Figure 3, we compress a picture the first author’s cat by treating the implicit array of **RGB** values as an order-3 tensor.

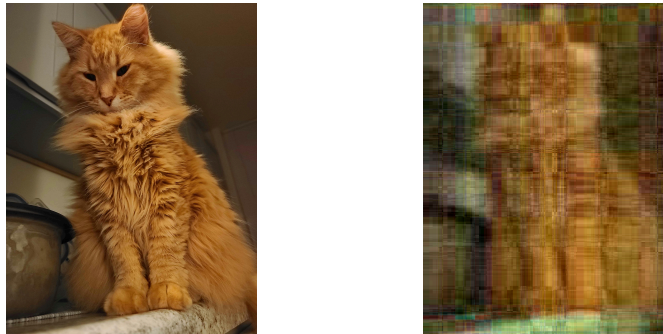


FIGURE 3. A picture of Angus and a very thin signed cut decomposition.

2. CUT NORMS AND DECOMPOSITIONS

2.1. The unsigned cut norm and Frieze-Kannan decompositions. We adapt techniques from combinatorial limit theory, a subfield of combinatorics cultivated by the Theory Group of Microsoft Research in the early 2000s. The germ of this theory is the *cut norm*, coined originally by Frieze and Kannan [2], which connects cut-set problems, matrix approximations, and the celebrated Regularity Lemma of Szemerédi. For any matrix $A = (a_{ij}) \in \mathbb{R}^{m \times n}$, the *cut norm* of A is

$$\|A\|_{\square} := \max_{S,T} \left| \sum_{(i,j) \in S \times T} a_{ij} \right| = \max_{S,T} |\langle \mathbf{1}_{S \times T}, A \rangle|,$$

where S (T) is any set of rows (columns) of A . In other words, $\|\cdot\|_{\square}$ measures the largest absolute sum over all submatrices of A . Although $\|A\|_{\square}$ is **MAX SNP hard**, it can theoretically be approximated efficiently in polynomial time (c.f. [3], [4]). As a motivating question, suppose we wish to optimize Frieze-Kannan decompositions on error, space, and time. As a function of the width w , we define the program

$$\text{FK}_w(A) : \begin{cases} \min & \|A - \sum_{j=1}^w d_j \cdot \mathbf{1}_{S_j \times T_j}\|_F \\ \text{s.t.} & d_1, \dots, d_w \in \mathbb{R} \\ & S_1, \dots, S_w \subseteq \{1, \dots, m\} \\ & T_1, \dots, T_w \subseteq \{1, \dots, n\} \end{cases}.$$

In [2], the authors approach $\text{FK}_w(A)$ with randomized algorithms which estimate the cut norm of a residual matrix and use the result to perform a rank-1 update on the remainder matrix. Their algorithms are complicated by two unhelpful facts:

1. $\text{FK}_1(A)$ is not in general solved by calculating $\|A\|_{\square}$. The optimal approximation $A \approx d \cdot \mathbf{1}_{S \times T}$ has $d = \langle \mathbf{1}_{S \times T}, A \rangle \cdot (|S| \cdot |T|)^{-1}$ for some other S, T .
2. Pairs of matrices of the form $\mathbf{1}_{S \times T}$ are typically far from orthogonal. If we reuse our solution to $\text{FK}_w(A)$ to solve $\text{FK}_{w+1}(A)$, then these correlations force us in general to apply large updates to the scalars d_1, \dots, d_w .

2.2. Signed cut decompositions and the signed cut norm. By instead basing our algorithms around the *signed cut norm* $\|\cdot\|_{\blacksquare}$, we mitigate the issues from Section 2.1 in practice. For any matrix $A \in \mathbb{R}^{m \times n}$, let

$$\|A\|_{\blacksquare} := \max_{s,t} \langle s, At \rangle = \max_{x,y} \langle x, Ay \rangle = \|A\|_{\infty \rightarrow 1}$$

where s, t are $\{-1, 1\}$ -valued and x, y are $[-1, 1]$ -valued. As show in [3], calculating $\|\cdot\|_{\blacksquare}$ is also **MAX SNP hard**. We approximate it with a randomized greedy algorithm. We define a *signed cut decomposition of width w* as any expression

$$A - R_w = S_w C_w T_w^\top = \sum_{j=1}^w c_j \cdot \mathbf{s}_j \mathbf{t}_j^\top$$

where $R_w \in \mathbb{R}^{m \times n}$, $S_w \in \{-1, 1\}^{m \times w}$, $C_w = \text{diag}(\mathbf{c}_w)$ for some $\mathbf{c}_w \in \mathbb{R}^w$, and $T_w \in \{-1, 1\}^{n \times w}$. As noted in Table 1, signed cut decompositions maintain the space efficiency of Frieze-Kannan decompositions while more closely emulating singular value decompositions (see Section 4).

Decomposition	Rank-1 terms	Vector entries	Size (bits)
Singular value	$\sigma_j \cdot \mathbf{u}_j \mathbf{v}_j^\top$	\mathbb{R}	$(m + n + 1) \cdot f$
Frieze-Kannan	$d_j \cdot \mathbf{1}_{S_j \times T_j}$	$\{0, 1\}$	$m + n + f$
Signed cut	$c_j \cdot \mathbf{s}_j \mathbf{t}_j^\top$	$\{-1, 1\}$	$m + n + f$

TABLE 1. A comparison of three matrix decompositions equivalent to finite sums of rank-1 matrices. Here, a floating point number occupies f bits in memory.

3. METHODS

For convenience, we let $\sigma(n) := \{-1, 1\}^n$. Our analog to $\text{FK}_w(A)$ is the program

$$\text{SC}_w(A) : \begin{cases} \min & \|A - \sum_{j=1}^w c_j \cdot \mathbf{s}_j \mathbf{t}_j^\top\|_F \\ \text{s.t.} & c_1, \dots, c_w \in \mathbb{R} \\ & \mathbf{s}_1, \dots, \mathbf{s}_w \in \sigma(m) \\ & \mathbf{t}_1, \dots, \mathbf{t}_w \in \sigma(n) \end{cases}.$$

Our two main algorithms are straightforward and greedy. We include a few optimizations within each subsection and defer others (e.g., **SIMD**, bitset storage, and alignment) to Section 3.4. Algorithm 4 estimates $\|\cdot\|_{\blacksquare}$ iteratively. Algorithm 5 bootstraps these estimates to extend approximate solutions from $\text{SC}_k(A)$ to $\text{SC}_{k+1}(A)$. Finally, Algorithm 6 adjusts the coefficients to improve the regression. In Section 5, we generalize to tensors.

3.1. Naive signed cut-sets. Note that $\|A\|_{\blacksquare}$ is the solution to the program

$$\text{SgnCut}(A) : \begin{cases} \max & \langle \mathbf{s}, A\mathbf{t} \rangle \\ \text{s.t.} & \mathbf{s} \in \sigma(m) \\ & \mathbf{t} \in \sigma(n) \end{cases}.$$

In [2], the authors explore the analogous program for the cut norm. By replacing \mathbf{s} with $\text{sgn}(A\mathbf{t})$, or \mathbf{t} with $\text{sgn}(A^\top \mathbf{s})$, we see that $\text{Cut}(A)$ is equivalent to the programs

$$\text{LeftSgnCut}(A) : \begin{cases} \max & \|A^\top \mathbf{s}\|_1 \\ \text{s.t.} & \mathbf{s} \in \sigma(m) \end{cases} \quad \text{RightSgnCut}(A) : \begin{cases} \max & \|A\mathbf{t}\|_1 \\ \text{s.t.} & \mathbf{t} \in \sigma(n) \end{cases}.$$

Algorithm 4 exploits this to approach a locally maximal cut by flipping signs.

GREEDY SIGNED CUT($A \in \mathbb{R}^{m \times n}$):

```

1  sample  $\mathbf{s}_0 \in \sigma(m)$  uniformly
2  sample  $\mathbf{t}_0 \in \sigma(n)$  uniformly
3  let  $c_0 := -\infty$ 
4  for  $j \in \{1, 2, \dots\}$ :
5      let  $\mathbf{s}_j := \text{sgn}(A \cdot \mathbf{t}_{j-1})$ 
6      let  $\mathbf{t}_j := \text{sgn}(A^\top \cdot \mathbf{s}_j)$ 
7      let  $c_j := \langle \mathbf{s}_j, A \cdot \mathbf{t}_j \rangle$ 
8      if  $c_j \leq c_{j-1}$ :
9          yeet  $(c_j, \mathbf{s}_j, \mathbf{t}_j)$ 
```

ALGORITHM 4. A randomized greedy algorithm for $\text{Cut}(A)$.

By caching the intermediate vectors $A \cdot \mathbf{t}_j, A^\top \mathbf{s}_j$, we may sparsify the matrix-vector products. To see this, note that at iteration $j \in \{2, 3, \dots\}$,

$$\begin{aligned} \mathbf{s}_j &= \text{sgn}(A \cdot \mathbf{t}_{j-1}) = \text{sgn}(A \cdot \mathbf{t}_{j-2} + A \cdot (\mathbf{t}_{j-1} - \mathbf{t}_{j-2})), \text{ and} \\ \mathbf{t}_j &= \text{sgn}(A^\top \cdot \mathbf{s}_{j-1}) = \text{sgn}(A^\top \cdot \mathbf{s}_{j-2} + A^\top \cdot (\mathbf{s}_{j-1} - \mathbf{s}_{j-2})). \end{aligned}$$

In practice, as j increases, the $\{-2, 0, 2\}$ -valued vectors $\mathbf{s}_j - \mathbf{s}_{j-1}$ and $\mathbf{t}_j - \mathbf{t}_{j-1}$ become increasing sparse. We store the variables as follows:

- **sign bitsets:** the $\{-1, 1\}$ -valued vectors \mathbf{s}_j and \mathbf{t}_j are stored with unsigned integers. Here, each 1 bit indicating a -1 . At iteration j , we also store \mathbf{s}_{j-1} and \mathbf{t}_{j-1} .
- **contiguous, aligned vectors:** the real-valued vectors $A \cdot \mathbf{t}_{j-1}$ and $A^\top \cdot \mathbf{s}_j$ are stored with aligned slices of `f32` values. The matrices A and A^\top are aligned, padded, and in column-major.

3.2. Simple signed cut decompositions. Algorithm 5 builds signed cut decompositions by repeatedly invoking Algorithm 4.

GREEDY DECOMPOSITION($A \in \mathbb{R}^{m \times n}$, $w \in \mathbb{N}$):

```

1  let  $R_0 := A$ 
2  let  $S_0 := 0 \in \{-1, 1\}^{m \times 0}$ 
3  let  $\mathbf{c}_0 := 0 \in \mathbb{R}^0$ 
4  let  $T_0 := 0 \in \{-1, 1\}^{n \times 0}$ 
5  for  $k \in \{0, \dots, w-1\}$ :
6      let  $(\mathbf{c}_{k+1}, \mathbf{s}_{k+1}, \mathbf{t}_{k+1}) = \text{greedy\_signed\_cut}(R_k)$ 
7      form  $\mathbf{c}_{k+1}$  by extending  $\mathbf{c}_k$  with  $\mathbf{c}_{k+1}$ 
8      form  $S_{k+1}$  by extending  $S_k$  with  $\mathbf{s}_{k+1}$ 
9      form  $T_{k+1}$  by extending  $T_k$  with  $\mathbf{t}_{k+1}$ 
10     let  $R_{k+1} := R_k - \frac{\mathbf{c}_{k+1}}{mn} \cdot \mathbf{s}_{k+1} \mathbf{t}_{k+1}^\top$ 
11 yeet  $(S_w, \text{diag}(\mathbf{c}_w), T_w)$ 
```

ALGORITHM 5. A “greedy” least-squares approach to signed cut decompositions. The remainder R_{k+1} is the orthogonal complement of R_k with respect to the 1-dimensional subspace of $\mathbb{R}^{m \times n}$ spanned by $\mathbf{s}_{k+1} \mathbf{t}_{k+1}^\top$.

By inspection,

$$\|R_{k+1}\|_F^2 = \|R_k\|_F^2 - \frac{c_{k+1}^2}{mn}.$$

for all $0 \leq k < w$. If c_k is proportional to $\|R_{k-1}\|_F$, then $\|R_k\|_F$ decreases exponentially. Empirically, results in Section 4 affirm this hypothesis.

To amortize the cost of the rank-1 updates to R_k , we delay the `matmul` accumulation in Algorithm 5 by storing each R_k implicitly as (R', S, C, T) where $R_k = R' + SCT^\top$. When (S, C, T) reaches a carefully chosen width (32), we flush SCT^\top into R' to form R_k . We store all real-valued vectors and matrices (except the diagonal matrices) as in Section 3.1. For the $\{-1, 1\}$ -valued matrices, we also use bitsets with the 1 bit indicating a -1 sign.

3.3. Least square corrections.

There is a simple improvement that can be made to Algorithm 5 based on the method of least squares. In the k -th iteration of Algorithm 6, we treat S_k and T_k as fixed and allow the diagonal matrix C_k to vary so as to minimize $\|A - S_k C_k T_k^\top\|_F$.

LEAST SQUARES DECOMPOSITION($A \in \mathbb{R}^{m \times n}$, $w \in \mathbb{N}$):

```

1  let  $S_0 := 0 \in \{-1, 1\}^{m \times 0}$ 
2  let  $c_0 := 0 \in \mathbb{R}^0$ 
3  let  $T_0 := 0 \in \{-1, 1\}^{n \times 0}$ 
4  for  $k \in \{0, \dots, w-1\}$ :
5      let  $R_k := S_k \text{diag}(c_k) T_k^\top$ 
6      let  $(\_, s_{k+1}, t_{k+1}) = \text{greedy\_signed\_cut}(R_k)$ 
7      form  $S_{k+1}$  by extending  $S_k$  with  $s_{k+1}$ 
8      form  $T_{k+1}$  by extending  $T_k$  with  $t_{k+1}$ 
9      let  $c_{k+1}$  minimize  $\|A - S_{k+1} \text{diag}(c_{k+1}) T_{k+1}^\top\|_F$ 
10 yeet  $(S_w, \text{diag}(c_w), T_w)$ 

```

ALGORITHM 6. An improvement to Algorithm 5. Here, R_{k+1} is the orthogonal complement of A with respect to the span of $s_1 t_1^\top, \dots, s_{k+1} t_{k+1}^\top$.

To calculate c_{k+1} in Algorithm 6, we formulate the associated least-squares problem and solve its *normal equation*

$$X_k^\top X_{k+1} \cdot c_{k+1} = X_{k+1}^\top \cdot A.$$

Here, X_{k+1} is the linear operator from \mathbb{R}^k to $\mathbb{R}^{m \times n}$ defined by sending j -th standard basis vector to $s_j t_j^\top$.

In our experiments, the improvement from swapping Algorithm 5 for Algorithm 6 is not significant for many matrices. We offer two possible explanations:

1. If only c_{k+1} varies, then Algorithm 5 and Algorithm 6 are equivalent.
2. The more orthogonal the matrices $s_k t_k^\top$ are, the less Algorithm 6 helps.

A lengthier analysis goes beyond the scope of this paper.

3.4. SIMD Optimizations. We improve the implementation in [5] by leveraging SIMD intrinsics on both `avx2` and `avx512` architectures. Interestingly, neither Algorithm 4 nor Algorithm 5 makes a single call to `fused multiply add` instructions. To demonstrate this, suppose we wish to compute the inner product $\langle s, x \rangle$ where

$$s = \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{bmatrix} \in \sigma(n) \quad \text{and} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n.$$

Assuming $n = 8\ell$ for some ℓ , we apply the isometry $\mathbb{R}^n \cong \bigoplus_{j=1}^{\ell} \mathbb{R}^8$ and write

$$\begin{aligned}
\mathbf{s} &= \mathbf{s}_1 \oplus \cdots \oplus \mathbf{s}_\ell && \text{where each } \mathbf{s}_j \in \sigma(8), \text{ and} \\
\mathbf{x} &= \mathbf{x}_1 \oplus \cdots \oplus \mathbf{x}_\ell && \text{where each } \mathbf{x}_j \in \mathbb{R}^8, \text{ so that} \\
\langle \mathbf{s}, \mathbf{x} \rangle &= \langle \mathbf{s}_1, \mathbf{x}_1 \rangle + \cdots + \langle \mathbf{s}_\ell, \mathbf{x}_\ell \rangle.
\end{aligned}$$

Since our signed sets are stored as bitsets, we need only use the bits of each \mathbf{s}_j to flip the sign of each \mathbf{x}_j entry and then sum. On `avx512` instructions (and the upcoming `avx10` instructions), we may offload this task to `mask` instructions with 512-bit registers. On `avx2` and `aarch64`, we may instead broadcast 32 bits to smaller `f32` registers, then use shifts instructions to move bits of interest to flip sign bits.

4. RESULTS

4.1. Approximating a random matrix. In this section, we evaluate the quality of signed cut decompositions by approximating a random matrix $A \in \mathbb{R}^{4096 \times 4096}$ with independent $\mathcal{N}(0, 1)$ entries. As an array of `f64` numbers, this occupies $8 \cdot 4096 \cdot 4096$ bytes, or $64 \cdot 4096 \cdot 4096$ bits, in memory. In Figure 7, we consider the intermediate signed cut decompositions $A = R_k + \sum_{j=1}^k c_j \mathbf{s}_j \mathbf{t}_j^\top$. We store each triple $(c_j, \mathbf{s}_j, \mathbf{t}_j)$ so the scalar has `f64` precision and the signed vectors $\mathbf{s}_j, \mathbf{t}_j$ as a bitset, as in Section 3. This way, each triple can be stored with $64 + m + n$ bits in memory. For each k , we plot (p_k, r_k) where

$$p_k := \frac{k \cdot (64 + 4096 + 4096)}{64 \cdot 4096 \cdot 4096} \quad \text{and} \quad r_k := \frac{\|R_k\|_F}{\|A\|_F}.$$

This way, p_k is the *compression rate* of the k -width approximation of A . For comparison, we plot $(0.25, r)$ where r is the relative error from downcasting A to `bf16` or `f16` precision.

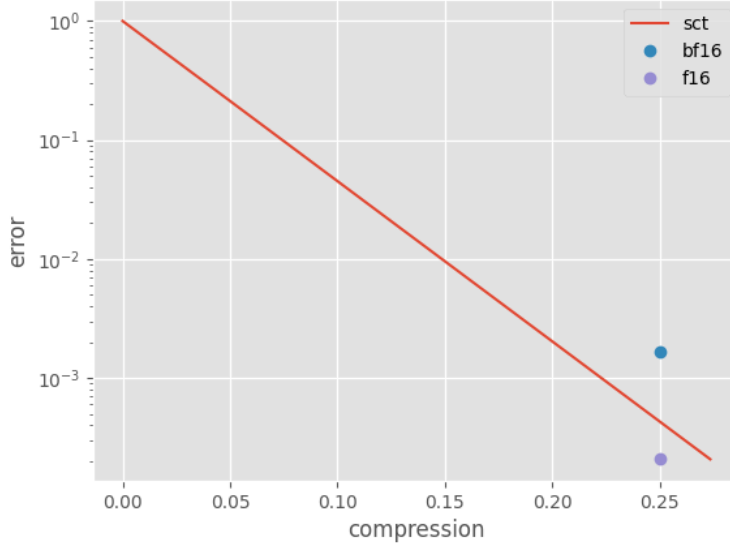


FIGURE 7. Relative error (log-scale) plotted against spatial compression rate for a standard normal 4096×4096 matrix with **f64** precision. Here, **bf16** has error ≈ 0.00166 and **f16** has error ≈ 0.000208 . Signed cut decompositions breaks even on error at 0.2064 and 0.2734, respectively.

4.2. **Approximating Mistral-7B-v0.1.** In this section, we apply signed cut decompositions to approximate the open **Mistral-7B-v0.1** model.

4.2.1. *Frobenius norm error.* We repeat our experiment in Section 4.1 by forming signed cut decompositions of the 226 weight matrices from the open **Mistral-7B-v0.1** model. Since these matrices are stored with **bf16** precision, each such $A \in \mathbb{R}^{m \times n}$ occupies $16mn$ bits in the original **safetensors** file. In the approximation $A \approx \sum_{j=1}^k c_j \mathbf{s}_j \mathbf{t}_j^\top$, we store each scalar c_j with **f32** precision and store each vector sign vector $\mathbf{s}_j, \mathbf{t}_j$ as a bitset, as described in Section 3. So the summands together occupy $32 + m + n$ bits in memory.

Targeting a 2-fold spatial compression, our approximation $A \approx \sum_{j=1}^w c_j \mathbf{s}_j \mathbf{t}_j^\top$ has width

$$w \approx \frac{1}{2} \cdot \frac{16mn}{32 + m + n}.$$

In Figure 8, we plot the relative error as a function of the compression rate.

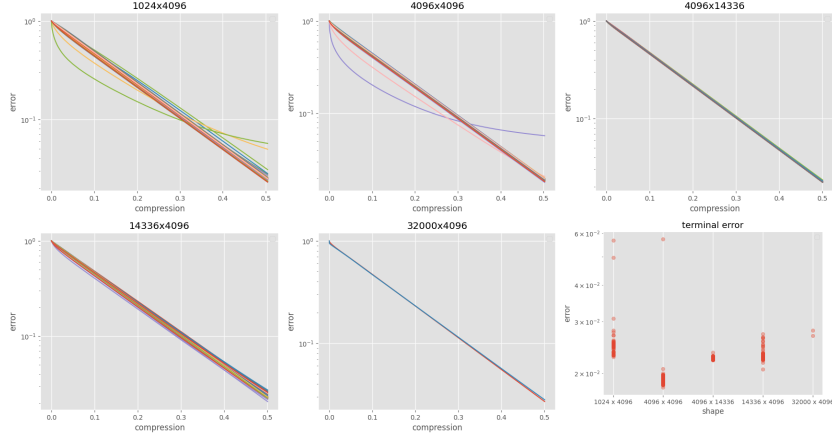


FIGURE 8. Relative error $\frac{\|R_k\|_F}{\|A\|_F}$ (log-scale) of signed cut decompositions of every matrix in **Mistral-7B-v0.1** after running Algorithm 5. Here, we group by matrix shape. The final width corresponds to a 50% spatial compression from the original **bf16** matrix.

4.2.2. *Large language model benchmarks.* Using $\frac{j}{16}$ for $j \in \{4, 5, 6, 7, 8\}$ as our target compression rate, we select appropriate widths based on the matrix dimensions m, n as in Table 2.

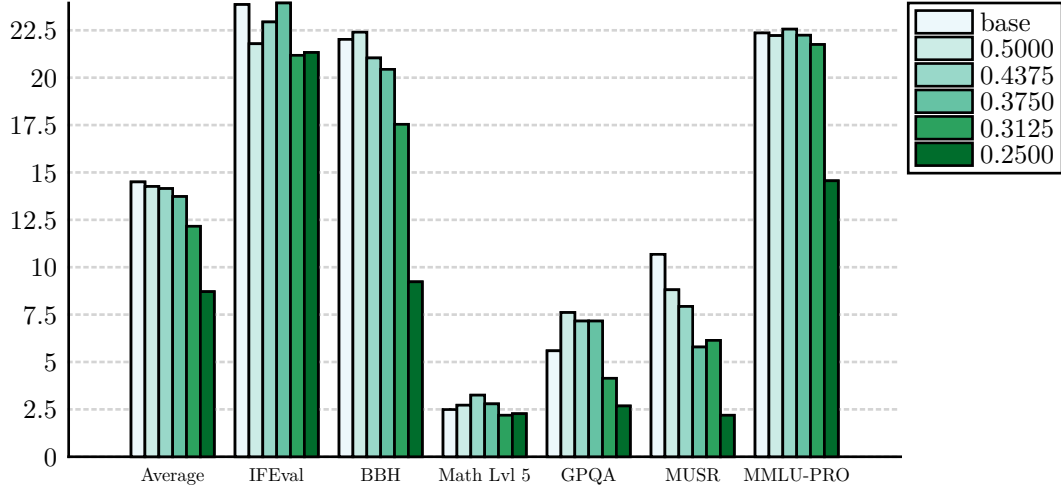
Shape	Count	Width by compression				
		0.5000	0.4375	0.3750	0.3125	0.2500
1024×4096	64	6512	5698	4884	4070	3256

Shape	Count	Width by compression				
		0.5000	0.4375	0.3750	0.3125	0.2500
4096×4096	64	16320	14280	12240	10200	8160
4096×14336	32	25442	22261	19081	15901	12721
14336×4096	64	25442	22261	19081	15901	12721
32000×4096	2	29023	25395	21767	18139	14511

TABLE 2. Widths matching each shape to each desired compression rate.

We expand the corresponding truncated signed cut approximations and write new **safetensors** files, reusing all other values. In order to measure the quality of signed cut decompositions as a vehicle for model quantization, we truncate the weight matrix approximations $A \approx \sum_{j=1}^w c_j \mathbf{s}_j \mathbf{t}_j^\top$ at different weights corresponding to our target compression rates. With m, n fixed, a w -width signed cut approximation of at **bf16** $m \times n$ matrix has compression rate

$$p_{m,n}(w) = \frac{w \cdot (32 + m + n)}{16mn}.$$

FIGURE 9. Benchmark performance of signed cut decompositions of the open **Mistral-7B-v0.1** model at compression rates between 0.5 and 0.25.

5. GENERALIZATIONS

As mentioned in Section 1, we motivate the tensor signed cut decomposition by approximating a picture of the first author’s cat Angus, as in Figure 10.

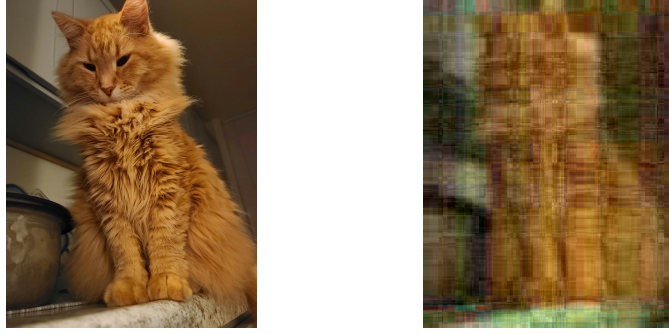


FIGURE 10. A picture of Angus and a very thin signed cut decomposition.

To begin, we convert this `jpg` image to a 4000×3000 array of `RGB` triples, with each value stored as an unsigned 8-bit integer. As a tensor, we may write $\mathbf{a} \in \mathbb{R}^{4000 \times 3000 \times 3}$, with axes corresponding to rows, columns, and colors, respectively. Our goal is to optimize an approximation of the form

$$\mathbf{a} \approx \mathbf{a}_w = \sum_{j=1}^w c_j \cdot \mathbf{x}_j \otimes \mathbf{y}_j \otimes \mathbf{k}_j.$$

By using `f32` precision for the scalars, and bitsets for the sign vectors, each quadruple $(c_j, \mathbf{x}_j, \mathbf{y}_j, \mathbf{k}_j)$, occupies $4 + 500 + 375 + 0.375 = 879.375$ bytes of memory.

When evaluating this approximation, it is natural to first round the entries of \mathbf{a}_w to the nearest value in $\{0, \dots, 255\}$. To account for this, we define

$$R(i) := \left\{ x \in \mathbb{R} : |x - i| = \min_{j \in \{0, \dots, 255\}} |x - j| \right\} \text{ for all } i \in \{0, \dots, 255\}.$$

In particular, $R(0) = (-\infty, 0.5]$, $R(255) = [254.5, \infty)$, and for all $j \in \{1, \dots, 254\}$, $R(j) = [j - \frac{1}{2}, j + \frac{1}{2}]$. In general, for all `RGB` image arrays $\mathbf{a} \in \{0, \dots, 255\}^{m \times n \times 3}$, the set of tensors which round to \mathbf{a} is

$$\mathcal{R}(\mathbf{a}) := \left\{ \mathbf{b} \in \mathbb{R}^{m \times n \times 3} : \text{for all } (i, j, k) \in [m] \times [n] \times [3], (\mathbf{b})_{ijk} \in R((\mathbf{a})_{ijk}) \right\}.$$

Then for all w , we let

$$\text{RGBSgnCut}_w(\mathbf{a}) : \left\{ \begin{array}{l} \min \left\| \mathbf{b} - \sum_{j=1}^w c_j \cdot \mathbf{x}_j \otimes \mathbf{y}_j \otimes \mathbf{k}_j \right\|_2 \\ \text{s.t. } \mathbf{b} \in \mathcal{R}(\mathbf{a}) \\ c_1, \dots, c_w \in \mathbb{R} \\ \mathbf{x}_1, \dots, \mathbf{x}_w \in \sigma(m) \\ \mathbf{y}_1, \dots, \mathbf{y}_w \in \sigma(n) \\ \mathbf{k}_1, \dots, \mathbf{k}_w \in \sigma(3) \end{array} \right. .$$

Pixel-wise, this means that estimating $i \in \{0, \dots, 255\}$ with $z \in \mathbb{R}$ contributes error

$$\rho(i, x) = \begin{cases} \max(0, x - \frac{1}{2}) & \text{if } i = 0 \\ \max(0, 255 - x - \frac{1}{2}) & \text{if } i = 255 \\ \max(0, |x - i| - \frac{1}{2}) & \text{else} \end{cases}$$

5.1. Signed cuts. For our purposes, an order- k tensor is an array $\mathbf{a} \in \mathbb{R}^{\otimes \mathbf{n}}$ where $\mathbf{n} = (n_1, \dots, n_k) \in \mathbb{N}^k$. We define the signed cut norm over $\mathbb{R}^{\otimes \mathbf{n}}$ as follows. First, let

$$\text{sgn}(\mathbf{n}) := \{-1, 1\}^{\otimes \mathbf{n}} = \bigotimes_{i=1}^k \text{sgn}(n_i).$$

Here, $\text{sgn}(\mathbf{n})$ is a $2^{n_1 + \dots + n_k}$ -element frame spanning $\mathbb{R}^{\otimes \mathbf{n}}$. Then for all $\mathbf{a} \in \mathbb{R}^{\otimes \mathbf{n}}$,

$$\|\mathbf{a}\|_{\blacksquare} := \max_{\mathbf{s} \in \text{sgn}(\mathbf{n})} \langle \mathbf{s}, \mathbf{a} \rangle.$$

Generalizing from Section 3.1, we note that for all $\mathbf{n} \in \mathbb{N}^k$ and all $1 \leq j \leq k$, the programs

$$\begin{aligned} \text{TensorSgnCut}(A) : & \begin{cases} \max & \langle \mathbf{s}, \mathbf{a} \rangle \\ \text{s.t.} & \mathbf{s} \in \sigma(\mathbf{n}) \end{cases} \quad \text{and} \\ \text{AxialCut}_i(A) : & \begin{cases} \max & \|\mathbf{a} \times_i \mathbf{s}\|_1 \\ \text{s.t.} & \mathbf{s} \in \sigma\left(\binom{n_j}{j \neq i}\right) \end{cases} \end{aligned}$$

both solve for $\|\mathbf{a}\|_{\blacksquare}$. Here, $\mathbf{a} \times_i \mathbf{s} \in \mathbb{R}^{n_j}$ is the product of \mathbf{a} and \mathbf{s} along the i -th axis. We extend Algorithm 4 from Section 3.1 to Algorithm 11.

AXIAL GREEDY CUT($\mathbf{a} \in \mathbb{R}^{\otimes \mathbf{n}}$):

- 1 **sample** $\mathbf{s}_0 = \bigotimes_{i \in \mathcal{A}} \mathbf{s}_{0i} \in \sigma(\mathbf{n})$ uniformly
- 2 **let** $c_0 := -\infty$
- 3 **for** $j \in \{1, 2, \dots\}$:
- 4 **for** $i \in \{1, \dots, k\}$:
- 5 **let** $\mathbf{s}_{j \downarrow i} := \bigotimes_{i' < i} \mathbf{s}_{ji}$
- 6 **let** $\mathbf{s}_{j \uparrow i} := \bigotimes_{i' > i} \mathbf{s}_{j-1, i}$
- 7 **let** $\mathbf{s}_{j \uparrow i} := \mathbf{s}_{j \downarrow i} \otimes \mathbf{s}_{j \uparrow i}$
- 8 **let** $\mathbf{s}_{ji} := \text{sgn}(\mathbf{a} \times_i \mathbf{s}_{j \uparrow i})$
- 9 **let** $\mathbf{s}_j := \bigotimes_{i \in \mathcal{A}} \mathbf{s}_{ji}$
- 10 **let** $c_j := \langle \mathbf{s}_j, \mathbf{a} \rangle$
- 11 **if** $c_j \leq c_{j-1}$:
- 12 **yeet** (c_j, \mathbf{s}_j)

ALGORITHM 11. Algorithm 4 for tensors.

5.2. Signed cut decompositions. Here, for any $\mathbf{a} \in \mathbb{R}^{\otimes \mathbf{n}}$, a *signed cut decomposition of width w* is any solution to

$$\mathbf{a} - \mathbf{r}_w = S_w \cdot \mathbf{c}_w = \sum_{j=1}^w c_j \cdot \mathbf{s}_j$$

where each $\mathbf{s}_j \in \sigma(\mathbf{n})$, S_w is a linear operator from \mathbb{R}^w to $\mathbb{R}^{\otimes \mathbf{n}}$ which sends the j -th standard basis vector to \mathbf{s}_j , $\mathbf{c}_w = (c_j) \in \mathbb{R}^w$, and $\mathbf{r}_w \in \mathbb{R}^{\otimes \mathbf{n}}$. Continuing, we extend Algorithm 5 from Section 3.2 to Algorithm 12.

TENSOR GREEDY DECOMPOSITION($\mathbf{a} \in \mathbb{R}^{\otimes \mathbf{n}}$, $w \in \mathbb{N}$):

```

1  let  $\mathbf{r}_0 := \mathbf{a}$ 
2  let  $S_0 := 0 \in \sigma(\mathbf{n}) \otimes \mathbb{R}^0$ 
3  let  $\mathbf{c}_0 := 0 \in \mathbb{R}^0$ 
4  for  $k \in \{0, \dots, w-1\}$ :
5      let  $(c_{k+1}, \mathbf{s}_{k+1}) = \text{axial\_greedy\_cut}(\mathbf{r}_k)$ 
6      form  $\mathbf{c}_{k+1}$  by extending  $\mathbf{c}_k$  with  $c_{k+1}$ 
7      form  $S_{k+1}$  by extending  $S_k$  with  $\mathbf{s}_{k+1}$ 
8      let  $\mathbf{r}_{k+1} := \mathbf{r}_k - \frac{c_{k+1}}{\prod_i n_i} \cdot \mathbf{s}_{k+1}$ 
9  yeet  $(\mathbf{c}_w, S_w)$ 

```

ALGORITHM 12. Algorithm 5 for tensors.

5.3. Least square corrections. Finally, we extend Algorithm 6 from Section 3.3 to Algorithm 13.

TENSOR LEAST SQUARES($\mathbf{a} \in \mathbb{R}^{\otimes \mathbf{n}}$, $w \in \mathbb{N}$):

```

1  let  $S_0 := 0 \in \sigma(\mathbf{n}) \otimes \mathbb{R}^0$ 
2  let  $\mathbf{c}_0 := 0 \in \mathbb{R}^0$ 
3  for  $k \in \{0, \dots, w-1\}$ :
4      let  $\mathbf{r}_k := \mathbf{a}_k - S_k \cdot \mathbf{c}_k$ 
5      let  $(\_, \mathbf{s}_{k+1}) = \text{axial\_greedy\_cut}(\mathbf{r}_k)$ 
6      form  $S_{k+1}$  by extending  $S_k$  with  $\mathbf{s}_{k+1}$ 
7      let  $\mathbf{c}_{k+1}$  minimize  $\|\mathbf{a} - S_{k+1} \cdot \mathbf{c}_{k+1}\|_2$ 
8  yeet  $(\mathbf{c}_w, S_w)$ 

```

ALGORITHM 13. Algorithm 6 for tensors.

6. DISCUSSION AND FUTURE WORK

We discuss some potential applications of signed cut decompositions for future work.

6.1. High Performance Computing. In Section 4.1, we showed that signed cut decompositions can approximate large random matrices roughly as well as **bf16**

and `f16` quantization while matching the memory footprint. In Section 4.2, we also showed that they approximate highly structured low-precision matrices even after driving down the memory footprint by a factor of 2 or higher. How do these gains in space complexity yield better runtimes? More precisely, we ask:

Question A: How does the runtime of `matvec` (calculating $A \cdot x$) for matrix A and vector x compare when A is replaced by a w -width signed cut decomposition?

At this time, we lack the infrastructure and resources to fully answer this question. As noted in Section 3.4, our `matmul` implementation manipulate bits en lieu of making a call to `FMA` instructions.

6.2. Machine Learning. Recall from Section 4.2 that our experiments replace all 226 weight matrices from the open `Mistral-7B-v0.1` model with the result of expanding signed cut decompositions $\sum_{j=1}^w c_j s_j t_j^\top$ where the width w is chosen so as to emulate a predetermined compression rate. Based on Table 2, we see that our 2-fold compression consists of

$$64 \cdot 6,512 + 64 \cdot 16,320 + 64 \cdot 25,442 + 32 \cdot 25,442 + 2 \cdot 29023 = 3,961,726$$

scalars c_j (and the same number of row and column vectors s_j and t_j). In other words:

Question B: If the s_j and t_j are fixed, can the `Mistral-7B-v0.1` model be efficiently retrained treating only the 3.96 million scalars as variables?

More broadly, we ask:

Question C: What is the best way to build machine learning architecture around signed vectors?

In a feed-forward neural network, some input vector $x \in \mathbb{R}^n$ is passed through hidden layers which pair affine linear transformations with activation functions:

$$\begin{aligned} x &\mapsto Ax + b \\ y &\mapsto \sigma(y). \end{aligned}$$

We may instead calculate $Y CZ \cdot x$ where Y, Z have $\{-1, 1\}$ -valued entries and C is a diagonal matrix of scalars. To justify universal approximation (provided that Z is sufficiently tall), it is sufficient to note that \mathbb{R}^k is spanned by $\sigma(k)$. Can we get away with only the “ $CZ \cdot x$ ” steps, while maintaining universal approximation?

7. ACKNOWLEDGEMENTS

We thank the reviewers Alice, Bob, and Charlie for their helpful comments. Additionally, we thank Bernard Lidický for lending the TODO computing cluster to run our first prototype.

REFERENCES

1. Murray, R., Demmel, J., Mahoney, M., Erichson, N., Melnichenko, M., Malik, O., Grigori, L., Luszczek, P., Dereziński, M., Lopes, M., others: Randomized Numerical Linear Algebra: A

Perspective on the Field With an Eye to Software (2023). DOI: <https://doi.org/10.48550/arXiv.2302>,

2. Frieze, A., Kannan, R.: Quick approximation to matrices and applications. *Combinatorica*. 19, 175–220 (1999). <https://doi.org/10.1007/s004930050052>
3. Alon, N., Naor, A.: Approximating the cut-norm via Grothendieck's inequality. In: *Proceedings of the thirty-sixth annual ACM symposium on Theory of computing*. pp. 72–80 (2004)
4. Lovász, L., Szegedy, B.: Szemerédi's lemma for the analyst. *GAFA Geometric And Functional Analysis*. 17, 252–270 (2007)
5. Neal Riasanovsky, A.: cuts, <https://github.com/ariasanovsky/cuts>

PHILADELPHIA, PA, USA

Email address: a.riasanovsky@gmail.com

URL: <https://riasanovsky.me/>