

Flexible High-Level Synthesis Library for Linear Transformations

Wuqiong Zhao, Changan Li, Zhenhao Ji, Zhichen Guo, Xuanbo Chen, You You, Yongming Huang,
Xiaohu You, *Fellow, IEEE*, and Chuan Zhang, *Senior Member, IEEE*

Abstract—Despite decades pursuing efficient hardware design for signal processing based on linear algebra, traditional hardware description languages (HDLs)-based design workflows remain challenging and time-consuming. High-level synthesis (HLS) provides an easier approach but still requires thorough designs of basic modules concerning linear transformations to achieve acceptable hardware efficiency. To simplify the HLS workflow, we propose the FLAMES library, which provides efficient ready-to-use linear transformation modules. Users can implement algorithms with significantly higher code-writing efficiency via the FLAMES library. We demonstrate its effectiveness by implementing the orthogonal matching pursuit list (OMPL) algorithm for compressed sensing in FPGA, achieving 1.56× and 1.12× throughput/slice compared with traditional HLS for the sequential and parallel architecture, respectively.

Index Terms—High-level synthesis (HLS), linear transformations, compressed sensing, hardware implementation, field programmable gate array (FPGA).

I. INTRODUCTION

HIGH-LEVEL synthesis (HLS) for hardware implementation can transcompile high-level programming languages like C/C++ into register-transfer level (RTL) designs [1]–[3]. HLS simplifies the hardware design workflow and lowers the barrier to obtaining efficient architectures [4]. Moreover, it can explore design space to achieve better performance with Pareto-optimal solutions [5]. Despite the promising future of HLS, limitations remain: users must adhere to strict coding guidelines to optimize hardware design [6], conflicting with the HLS concept of a simplified hardware design paradigm.

To alleviate these limitations of HLS, higher-level libraries are constructed for complexity reduction and performance enhancement. The proposed template-based method in [7] shows its superiority in several cases. FBLAS [8] ports the basic linear algebra subprograms (BLAS) library to HLS, promoting HLS design productivity. However, the complicated interfaces of [7], [8] impose inflexibility for users in practical scenarios. We take advantage of `class` and `template` of C++ with Vitis HLS [9] to further reduce the gap between HLS design and software algorithms by providing a user-friendly library with concise and customizable interfaces like Armadillo C++ library [10]. Algorithms can be efficiently implemented with the matrix-based library since an algorithm can be viewed as matrices and their transformations. The optimized hardware can be obtained by simply mapping the matrix transformations into HLS code without the need for extensive design. Our contributions are as follows:

- 1) We present the novel matrix-based HLS library concept FLAMES (Flexible Linear Algebra with Matrix-Empowered Synthesis) for fast algorithm implementation, which is designed as a header-only C++ library for Vitis HLS, simplifying the design workflow of signal processing modules based on linear transformations;
- 2) With our designed HLS library, efficient architectures for orthogonal matching pursuit list (OMPL) compressed sensing (CS) are implemented as a verification. The sequential and parallel designs achieve throughput 94.81 Mb/s and 184.5 Mb/s (1.56× and 1.12× throughput/slice compared with traditional HLS), respectively, with significantly reduced design complexity.

Notations: In this paper, lower-case and upper-case boldface letters \mathbf{a} and \mathbf{A} stand for vector and matrix, respectively. \mathbf{A}^T and \mathbf{A}^{-1} represent the transpose and inverse of matrix \mathbf{A} , respectively. $\|\mathbf{a}\|_1$ and $\|\mathbf{a}\|_2$ take the ℓ_1 - and ℓ_2 -norm of vector \mathbf{a} , respectively. $\langle \mathbf{a}, \mathbf{b} \rangle$ denotes the inner product of vector \mathbf{a} and \mathbf{b} . $\lfloor x \rfloor$ is the floor function returning the largest integer smaller or equal to x . Finally, $\mathcal{N}(\mu, \sigma^2)$ denotes the Gaussian distribution with mean μ and variance σ^2 .

II. HLS LIBRARY DESIGN

A. Framework and Syntax

Despite the power of HLS tools, designing hardware-efficient solutions still requires substantial knowledge and experience in hardware design. Existing HLS libraries like [7], [8] are not as user-friendly as popular linear algebra libraries for general C++ programming like Armadillo [10]. To address the issue, we propose a matrix-based library design where matrix operations and functions are at the center of the design, which enables signal processing algorithms to be expressed in their original form with minimal adjustments without the need for data decomposition or pipeline designs. It supports most basic matrix operations in Armadillo [11]. For synthesis, operations requiring dynamic memory allocation are not supported. Table I provides several examples of library syntax. One notable feature of the FLAMES library is its unified user interface, which does not rely on low-level variable types such as C arrays, C++ standard containers, or streams. It makes the HLS design better organized and less redundant. Furthermore, the library is designed to be header-only, making it easy to install and fully portable. Users only need to include the required header file to access the classes and utilities provided by the FLAMES library, which is open source at <https://github.com/autohdw/flames>.

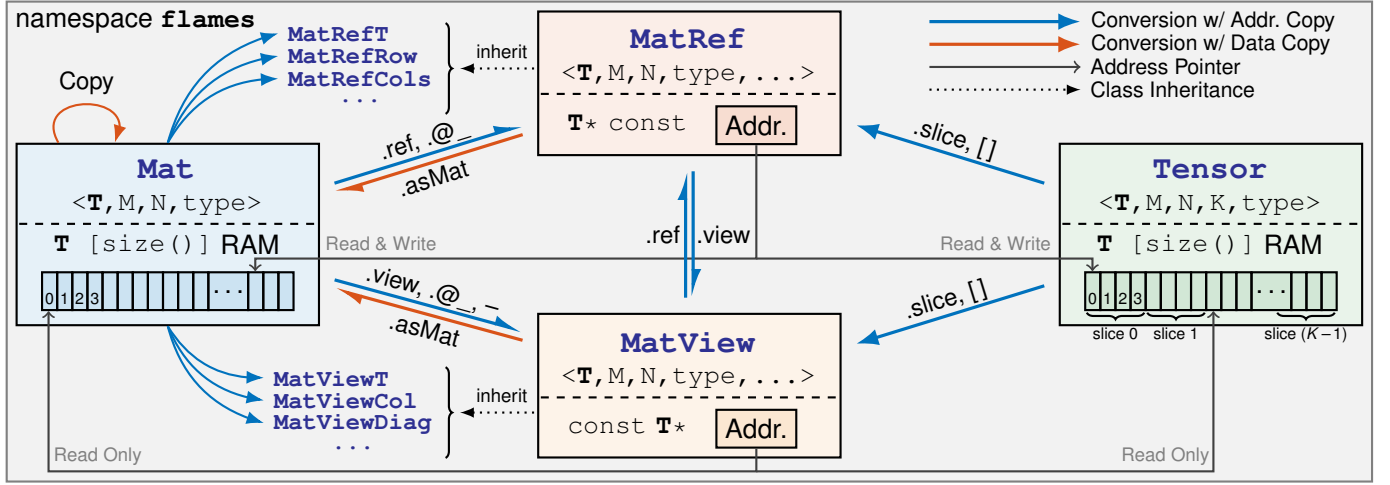


Fig. 1. Classes and their relationships in the FLAMES library, where '@' is a placeholder for function names, such as 't' and 'col'. \mathbf{T} is the data type.

TABLE I
LIBRARY SYNTAX EXAMPLES

HLS Code	Description
$A + B$ and $A - B$	addition/subtraction
$A * B$ and $A \% B$	multiplication/element-wise multiplication
$A += B$ or $A.add(B)$	self adding another matrix
$A.t()$ and $-A$	transpose/negation
$A(2,3)$ or $A.at(2,3)$	accessing element of row 2 column 3
$A.row(1)$	accessing row 1
$A.cols(\{2,3,4\})$	accessing columns 2, 3 and 4
$A.diagMat()$	diagonal as matrix

A and B are both matrices.

B. Basic Structure

The FLAMES library provides several primary classes for matrix and tensor operations, depicted in Fig. 1. These classes include **Mat** for matrices, **MatRef** for matrix references, **MatView** for read-only matrix views, and **Tensor** for tensors. **MatRef** and **MatView** can be used for simple matrix transformations without data copying. To be specific, **MatRef** has a constant pointer to the original data, while **MatView** has a pointer to constant values with read-only access. In addition, **MatView** and **MatRef** can be converted back to **Mat** using the `.asMat` method to allocate the required RAM. The **Tensor** class allocates RAM slice by slice, and slices can be accessed using a **MatRef** or **MatView** with the `.slice` method or the `[]` operator by specifying the slice index. Matrices and tensors have different types that determine data arrangement. To simplify the use of FLAMES, the `auto` keyword can be leveraged to automatically infer **MatView** and **MatRef** types without explicit declaration.

C. Hardware-friendly Design

In addition to the simple user interface, the FLAMES library is carefully designed for efficient hardware implementation.

1) *Optimized RAM Usage*: Matrices are often partitioned in algorithm implementations for performance enhancement, imposing difficulties for hardware designers [12]–[14]. In contrast, the FLAMES library offers a neat matrix-based interface while optimizing RAM usage with different **MatTypes**. For instance, diagonal matrices (`MatType::DIAGONAL`) and upper triangular matrices (`MatType::UPPER`) are efficiently

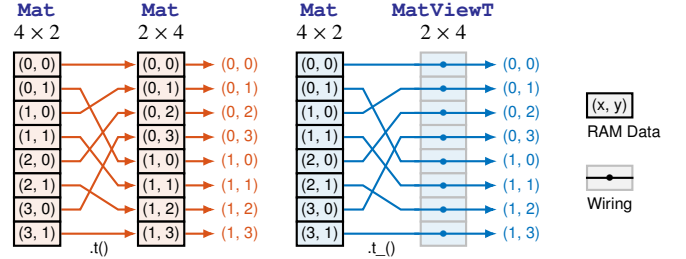


Fig. 2. Permutation network for 4×2 transposed matrix.

stored exploiting their structural properties, while providing a unified user interface. Furthermore, **MatView** avoids unwanted intermediate RAMs due to the lack of return value optimization (RVO) in Vitis HLS. For example, the permutation of the transposed matrix (Fig. 2) is achieved using the **MatViewT** class inherited from **MatView**, which connects wires to the original data instead of copying and allocating RAM for the transposed matrix.

2) *Configurable Parallelism*: The library employs the HLS pragma internally provided by Vitis HLS, controlling the pipelining, unrolling, loop flattening, etc. The parallelism for matrix multiplication, matrix copying, etc. can be configured separately, by defining the corresponding macro such as `FLAMES_MAT_TIMES_UNROLL_FACTOR`. Similarly, the data array partitioning can also be configured.

3) *Optimized Matrix Operations*: The FLAMES library optimizes matrix operations by leveraging the distinctive properties of different matrices. To achieve case-specific hardware optimization, the library employs function overloading. As shown in section III, users do not need to distinguish over 50 overloading functions for matrix multiplication, as optimization is applied automatically. Additionally, FLAMES offers commonly used functions in linear transformations and signal processing. For example, for a diagonally dominant matrix \mathbf{A} , its inverse can be computed efficiently with Neumann series approximation (NSA) [15] as

$$\mathbf{A}^{-1} = \lim_{n \rightarrow \infty} \sum_{i=0}^n (-\mathbf{D}^{-1}\mathbf{E})^i \mathbf{D}^{-1}, \quad (1)$$

where $\mathbf{A} = \mathbf{D} + \mathbf{E}$, $\mathbf{D} \triangleq \mathbf{A} \circ \mathbf{I}$ is the diagonal part while \mathbf{E} is the off-diagonal part (\circ represents the Hadamard product).

TABLE II
NSA ALGORITHM AND FLAMES CODE

Step	Algorithm	FLAMES C++ Implementation
1	$\mathbf{D} = \mathbf{A} \circ \mathbf{I}$	<code>auto D = mat.diagMat_();</code>
2	$\mathbf{E} = \mathbf{A} - \mathbf{D}$	<code>auto E = mat.offDiag_();</code>
3	$\mathbf{D}_I = \mathbf{D}^{-1}$	<code>auto D_I = D.inv();</code>
4	$\mathbf{P} = -\mathbf{D}_I \mathbf{E}$	<code>auto P = -D_I * E;</code>
5	$\mathbf{X} = \mathbf{P}$ (Iter. 1)	<code>auto X = P_ = P;</code>
6	for $i = 2, \dots, n$	for (int i = 2; i <= n; ++i) {
7	$\mathbf{P}^i = \mathbf{P}^{i-1} \mathbf{P}$	<code>P_ *= P;</code>
8	$\mathbf{X} = \mathbf{X} + \mathbf{P}^i$	<code>X += P_;</code>
9	end	}
10	$\mathbf{A}^{-1} = \mathbf{X} \mathbf{D}_I + \mathbf{D}_I$	<code>A_inv = X * D_I + D_I;</code>

Algorithm 1 OMPL with square-root-free QR decomposition.

Input: $\Phi \in \mathbb{R}^{M \times N}$, $\mathbf{y} \in \mathbb{R}^{M \times 1}$.

Initialization: $\hat{\mathbf{x}} = \mathbf{0}_{N \times 1}$.

```

1: for  $i = 0, 1, \dots, L-1$  do
2:   for  $j = 0, 1, \dots, n-1$  do ▷ BRANCH.
3:      $\tilde{\mathbf{s}}_i^{jn}, \tilde{\mathbf{s}}_i^{jn+1}, \dots, \tilde{\mathbf{s}}_i^{(j+1)n-1} = \arg \max_k^{(n)} |(\mathbf{r}^j, \Phi_{:,k})|;$ 
4:     for  $k = 0, 1, \dots, n-1$  do ▷ SUBBRANCH.
5:        $\mathbf{w} = \Phi_{:, \tilde{\mathbf{s}}_i^{jn+k}};$ 
6:       for  $q = 0, 1, \dots, i-1$  do
7:          $\tilde{\mathbf{R}}_{q,i}^{jn+k} = (\tilde{\mathbf{Q}}_{:,q}^{jn+k})^\top \mathbf{w};$ 
8:          $\mathbf{w} = \mathbf{w} - (\tilde{\mathbf{R}}_{q,i}^{jn+k} / \tilde{\mathbf{R}}_{q,q}^{jn+k}) \tilde{\mathbf{Q}}_{:,q}^{jn+k};$ 
9:          $\tilde{\mathbf{R}}_{i,i}^{jn+k} = \|\mathbf{w}\|_2, \tilde{\mathbf{Q}}_{:,i}^{jn+k} = \mathbf{w};$ 
10:         $\tilde{\mathbf{r}}^{jn+k} = \tilde{\mathbf{r}}^{jn+k} - ((\tilde{\mathbf{Q}}_{:,i}^{jn+k}, \tilde{\mathbf{r}}^{jn+k}) / \tilde{\mathbf{R}}_{i,i}^{jn+k}) \tilde{\mathbf{Q}}_{:,i}^{jn+k};$ 
11:       $\mathbf{a} = \mathbf{u} \arg \min_k^{(n)} \|\tilde{\mathbf{r}}^k\|_1;$  ▷ MERGE.
12:      for  $j = 0, 1, \dots, n-1$  do
13:         $\mathbf{r}^j = \tilde{\mathbf{r}}^{\mathbf{a}^j}, \mathbf{Q}^j = \tilde{\mathbf{Q}}^{\mathbf{a}^j}, \mathbf{R}^j = \tilde{\mathbf{R}}^{\mathbf{a}^j}, \mathbf{s}^j = \tilde{\mathbf{s}}^{\mathbf{a}^j};$ 
14:      for  $j = 0, 1, \dots, n^2-1$  do
15:         $\tilde{\mathbf{r}}^j = \mathbf{r}^{\lfloor j/n \rfloor}, \tilde{\mathbf{Q}}^j = \mathbf{Q}^{\lfloor j/n \rfloor}, \tilde{\mathbf{R}}^j = \mathbf{R}^{\lfloor j/n \rfloor}, \tilde{\mathbf{s}}^j = \mathbf{s}^{\lfloor j/n \rfloor};$ 
16:       $m = \arg \min_k \|\tilde{\mathbf{r}}^k\|_1; \mathbf{Q} = \mathbf{Q}^m, \mathbf{R} = \mathbf{R}^m, \mathbf{s} = \mathbf{s}^m;$ 
17:    for  $l = 0, 1, \dots, L-1$  do ▷ SOLVE.
18:       $\hat{\mathbf{x}}_{s_{L-l-1}} = ((\tilde{\mathbf{Q}}_{:,L-l-1}, \mathbf{y}) - \mathbf{R}_{L-l-1, L-l-1} \hat{\mathbf{x}}_s) / \mathbf{R}_{L-l-1, L-l-1};$ 

```

Output: Estimated L -sparse signal $\hat{\mathbf{x}}$.

As shown in Table II, it is simple to map an algorithm to the FLAMES code.

III. CASE STUDY AND VERIFICATION

Compressed sensing [16] is a powerful technique with numerous applications, including millimeter wave (mmWave) channel estimation in the angular domain for multiple-input multiple-output (MIMO) systems [17], [18]. It is modeled as

$$\mathbf{y} = \Phi \mathbf{x} + \mathbf{n}, \quad (2)$$

where $\mathbf{y} \in \mathbb{R}^{M \times 1}$ is the compressed signal, $\Phi \in \mathbb{R}^{M \times N}$ ($M < N$) is the sensing matrix, $\mathbf{x} \in \mathbb{R}^{N \times 1}$ is the original L -sparse signal, $\mathbf{n} \sim \mathcal{N}(0, \sigma_n^2)$ is the additive white Gaussian noise (AWGN), and σ_n^2 is the noise variance. The OMP algorithm in [19] can be employed to achieve efficient sparse signal recovery. However, the performance of OMP is not satisfactory enough, and an efficient OMPL algorithm in [20] can be applied to enhance the recovery performance. The OMPL algorithm extends the OMP iteration to n lists, which are branched and merged in each iteration. Every set of support

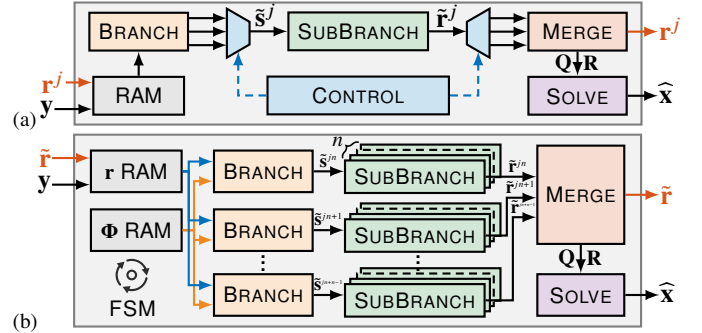


Fig. 3. Hardware architecture of OMPL. (a) Sequential. (b) Parallel.

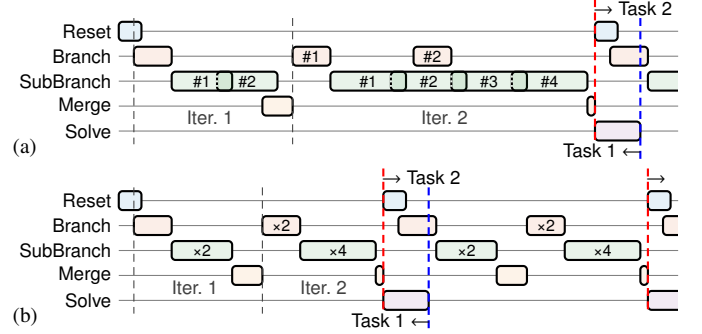


Fig. 4. Timeline of (a) sequential and (b) parallel designs with $L = 2$, $n = 2$.

indices is a candidate in OMPL. With the square-root-free QR decomposition [21] method for OMP proposed in [22], the OMPL algorithm can be represented as Alg. 1.

The hardware architecture of OMPL consists of 4 modules. For a given candidate, the BRANCH module computes n sub-candidates. The SUBBRANCH module updates the QR decomposition result and residual of a sub-candidate. The MERGE module selects the best n candidates among the n^2 sub-candidates according to their residual. The SOLVE module solves the least square (LS) problem. The sequential architecture (Fig. 3(a)) reuses one BRANCH module and one SUBBRANCH module, whereas the parallel architecture (Fig. 3(b)) contains n BRANCH modules and n^2 SUBBRANCH modules. In each iteration, the parallel architecture processes all the sub-candidates in parallel. As is shown in the timeline traces (Fig. 4), task-level pipelining is automatically achieved for both architectures. With matrix-based operations, FLAMES provides a clearer task flow, making it easier for users to fully exploit the task-level pipelining optimization provided by HLS tools. Due to the iterative nature of OMPL, pipelining is limited in this case. To demonstrate the availability of task-level pipelining optimization, an additional example is provided in [11].

The main iterations of OMPL compressed sensing are between step 1 and 15, where BRANCH and MERGE operations are repeated.

BRANCH (step 3). For each candidate, BRANCH selects n new support indices as sub-candidates according to the correlation of Φ and \mathbf{r}^j as illustrated in Fig. 5. A binary sensing matrix with 1-bit quantization is used [23]–[25] to simplify the correlation calculation, resulting in Φ being `Mat<bool, M, N>` in FLAMES. Simply writing $*$, the binary matrix-vector multiplication is specially optimized in

FLAMES with configurable parallelism p : a multiplexer array followed by an adder and accumulator array is employed. This corresponds to the internal implementation of $+=$ following an `if` condition. The accumulator array outputs the $\Phi^T \mathbf{r}^j$ vector every $\frac{M}{p}$ clock cycles, which is then processed with an absolute value calculating array before the sorting network locates the n largest elements of the vector. Each of the sub-candidates leads to one SUBBRANCH. Notably, only one BRANCH operation is required in the first iteration.

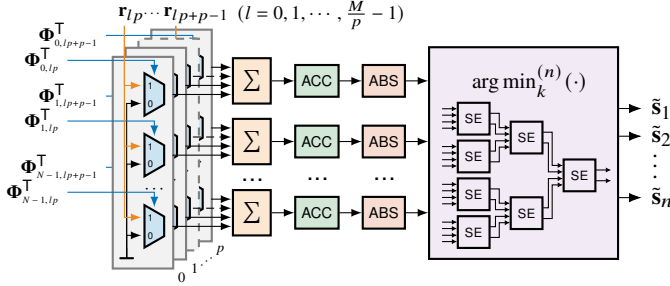
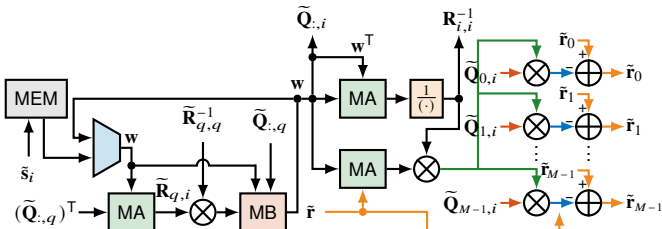
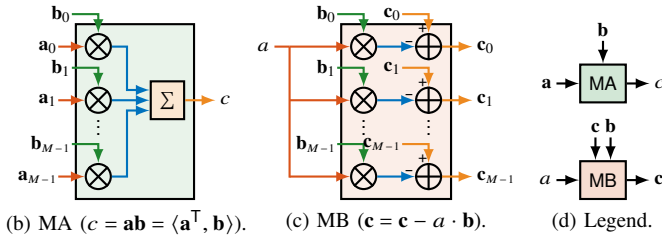


Fig. 5. Hardware architecture of the BRANCH module.

SUBBRANCH (step 4 to 10). With the square-root-free QR decomposition [22], \mathbf{w} , $\mathbf{R}_{:,i}$ and $\mathbf{Q}_{:,i}$ are calculated iteratively by reusing linear transformation units including vector-vector multiplication and vector-scalar multiplication. For hardware efficiency, $\mathbf{R}_{:,i}^{-1}$ is stored instead of $\mathbf{R}_{:,i}$. In Fig. 6(e), `dtype` is the data type of \mathbf{w} which writes into $\mathbf{Q}_{:,i}$ via a `MatRefCol` created by `.col_`, and the `innerProd` function calculates the inner product of two vectors.



(a) Overall architecture. The superscript $jn+k$ is omitted for simplicity.



```

auto w = Phi.col<dtype>(cand_s[j*n+k]);
SB_LOOP: for (int p = 0; p + 1 < i; ++p)
  w -= (Ri(p, i) = innerProd(Qi.col_(p), w))
        / Ri(p, p) * Qi.col_(p);
auto wp = Ri(i, i) = w.power(); /* l2-norm square */
Qi.col_(i) = w;
ri -= innerProd(w, ri) / wp * w;

```

(e) Corresponding C++ code.

Fig. 6. Hardware architecture and C++ code of the SUBBRANCH module.

MERGE (step 11 to 16). To seek the n least residual support sets in n^2 candidates with unique indices forming

\mathbf{a} , denoted as $\text{u arg min}_k^{(n)}$, ℓ_1 -norm is used instead of ℓ_2 -norm for hardware simplicity. The corresponding matrices are updated until step 15. MERGE is implemented with a multi-layer sorting tree utilizing `flames/sort.hpp`. Each layer contains sorting nodes that select the n smallest values of its input. Finally, the selected set is obtained by finding the set with minimum residual of the n lists.

SOLVE (step 17 to 18). It performs matrix-vector multiplication, backward substitution, and location mapping to solve the LS problem. With simple $\mathbf{Q} * \mathbf{y}$, the matrix-vector multiplication is auto implemented like that in the BRANCH module, but with a multiplier array replacing the multiplexer array. Backward substitution is then performed with a PE array of size L (a for loop in FLAMES), where PE is a 3-input, 1-output processing unit shown in Fig. 7 (* for multiplier and += for ACC in FLAMES). Finally, $\tilde{\mathbf{x}}_l$ is mapped to its location s_l via the mapping unit, implemented by assigning to the corresponding elements using the operator `[]`.

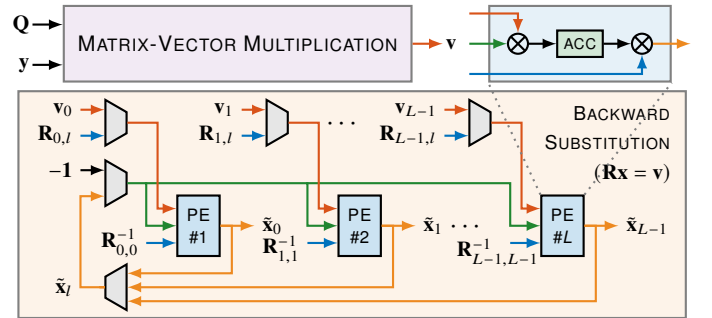


Fig. 7. Hardware architecture of the SOLVE module (mapping unit omitted).

IV. EXPERIMENTAL RESULTS

A. Performance Verification

Fig. 8 compares the normalized mean square error (NMSE) performance defined as $\mathbb{E}[\|\tilde{\mathbf{x}} - \mathbf{x}\|_2 / \|\mathbf{x}\|_2]$ versus the signal-to-noise ratio (SNR), where the OMPL list size is $n = 2$. The LS method is compared as a benchmark with no dimension compression. OMPL outperforms OMP [19] using both random binary ('-R' suffix) and designed binary sensing matrices ('-D' suffix) in Fig. 8(a), where $M = 64$, $N = 128$, $L = 8$. In Fig. 8(b), N varies among 24, 32 and 40, and $M = 16$, $L = 2$.

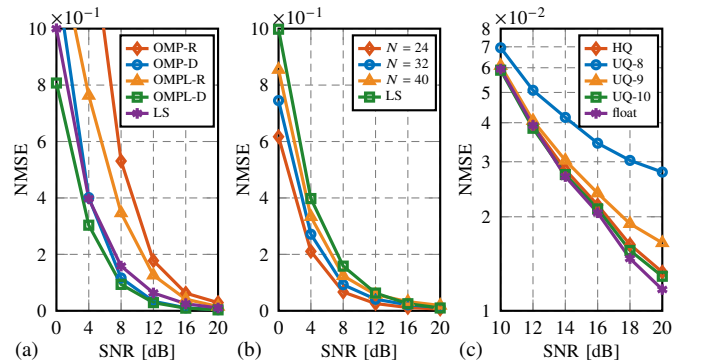


Fig. 8. NMSE performance of: (a) OMP and OMPL, (b) OMPL under different measurement size, (c) OMPL under different quantization schemes.

In Fig. 8(c), we use hybrid quantization (HQ) with an average of around 7 fixed-point bits ($M = 16$, $N = 32$, $L = 2$), which is close to the performance of 10 bits unified quantization (UQ) and float-point precision. In our HQ scheme, \mathbf{r} , $\hat{\mathbf{x}}$, \mathbf{w} , \mathbf{Q} , \mathbf{R} have 8, 10, 3, 7, and 7 bits, respectively.

B. Hardware Implementation

The OMPL algorithm for compressed sensing is synthesized and implemented for Xilinx Zynq-7000 ZC-706 with Vitis HLS 2022.2 [9]. As a baseline, we implemented OMPL following a general HLS design guideline [6], denoted as traditional HLS (T-HLS). Compared with T-HLS, FLAMES-assisted designs achieve 3.96 \times and 1.67 \times latency reduction and 1.56 \times and 1.12 \times throughput/slice for the sequential and parallel designs respectively with fewer than half lines of code, with details in [11]. The FPGA resource consumption and performance comparisons are listed in Table III. The parallel design achieves higher throughput and lower latency than the sequential design at the expense of higher overall resource consumption.

TABLE III
FPGA RESOURCE AND PERFORMANCE COMPARISONS FOR 16 \times 32
COMPRESSED SENSING WITH OMPL

Metric	Sequential		Parallel	
	T-HLS [6]	FLAMES	T-HLS [6]	FLAMES
Slice	2,566	6,912	6,022	8,181
LUT	4,360	15,704	17,129	9,669
FF	5,958	22,114	10,899	39,170
DSP	81	36	75	5
BRAM	14	1	5	1
Latency [μ s]	13.91	3.510	3.113	1.864
Frequency [MHz]	134	104	128	107
Throughput [Mb/s]	23.63	94.81	120.9	184.5
Throughput/slices [Mbps/K slices]	8.819	13.72	20.08	22.55

Remarks: 1) DSPs and BRAMs are disabled in the RTL synthesis phase. However, there is no direct way for the Vitis HLS synthesis phase. 2) The overall resource consumption is reflected by the number of slices. Certain types of resources may differ due to complex HLS optimization strategies. 3) The sequential and parallel designs only differ in the HLS pragma/directive, so the OMPL algorithm can be readily implemented as different architectures. Notably, the FLAMES library excels at achieving high throughput in contrast to HLS's conservative optimization strategy. 4) The FLAMES library significantly simplifies the hardware design with the matrix-based syntax shown in Table I and [11]. Complicated operations like branching and merging in OMPL are implemented easily with matrix-based coding, significantly reducing the complexity and difficulty of hardware design.

V. CONCLUSION

In this paper, we propose FLAMES for matrix-based linear transformations, an open-source high-level synthesis library implemented on Vitis HLS. FLAMES achieves hardware and coding efficiency for fast algorithm implementations. Moreover, we demonstrate the effectiveness of FLAMES through the OMPL algorithm, which verifies the library's capability to synthesize efficient hardware designs.

REFERENCES

- [1] D. D. Gajski, N. D. Dutt, A. C. Wu, and S. Y.-L. Lin, *High-Level Synthesis: Introduction to Chip and System Design*, D. D. Gajski, Ed. New York, NY: Springer, 1992.
- [2] P. Coussy, D. D. Gajski, M. Meredith, and A. Takach, "An introduction to high-level synthesis," *IEEE Des. Test. Comput.*, vol. 26, no. 4, pp. 8–17, 2009.
- [3] R. Nane, V.-M. Sima, C. Pilato, J. Choi, B. Fort, A. Canis *et al.*, "A survey and evaluation of FPGA high-level synthesis tools," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 35, no. 10, pp. 1591–1604, Oct. 2016.
- [4] Z. Zhong, W. J. Gross, Z. Zhang, X. You, and C. Zhang, "Polar compiler: Auto-generator of hardware architectures for polar encoders," *IEEE Trans. Circuits Syst. I*, vol. 67, no. 6, pp. 2091–2102, Jun. 2020.
- [5] B. C. Schafer and Z. Wang, "High-level synthesis design space exploration: Past, present, and future," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 39, no. 10, pp. 2628–2639, Oct. 2020.
- [6] R. Kastner, J. Matai, and S. Neundorffer, "Parallel programming for FPGAs," 2018, *arXiv:1805.03648*. [Online]. Available: <https://arxiv.org/abs/1805.03648>
- [7] J. Matai, D. Lee, A. Althoff, and R. Kastner, "Composable, parameterizable templates for high-level synthesis," in *Proc. IEEE Design Autom. Test Europe Conf. Exhibit. (DATE)*, Mar. 2016, pp. 744–749.
- [8] T. De Matteis, J. de Fine Licht, and T. Hoefler, "FBLAS: Streaming linear algebra on FPGA," in *Proc. IEEE Int. Conf. HPC Netw. Stor. Anal.*, Nov. 2020, pp. 1–13.
- [9] Xilinx, "Vitis high-level synthesis user guide (UG1399)," 2023, accessed: Feb. 27, 2023. [Online]. Available: <https://docs.xilinx.com/t/en-US/ug1399-vitis-hls>
- [10] C. Sanderson and R. Curtin, "Armadillo: a template-based C++ library for linear algebra," *J. Open Source Softw.*, vol. 1, no. 2, p. 26, 2016.
- [11] W. Zhao, "FLAMES insight," 2023, accessed: Sep. 18, 2023. [Online]. Available: https://flames.autohdw.com/FLAMES_Insight.pdf
- [12] M. Aynala, M. Brown, and K. K. Parhi, "Pipelined parallel FFT architectures via folding transformation," *IEEE Trans. VLSI Syst.*, vol. 20, no. 6, pp. 1068–1081, Jun. 2012.
- [13] B. Yuan and K. K. Parhi, "Architecture optimizations for BP polar decoders," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, May 2013, pp. 2654–2658.
- [14] H. Wang, Y. Ji, Y. Shen, W. Song, M. Li, X. You *et al.*, "An efficient detector for massive MIMO based on improved matrix partition," *IEEE Trans. Signal Process.*, vol. 69, pp. 2971–2986, 2021.
- [15] X. Tan, W. Xu, Y. Zhang, Z. Zhang, X. You, and C. Zhang, "Efficient expectation propagation massive MIMO detector with Neumann-series approximation," *IEEE Trans. Circuits Syst. II*, vol. 67, no. 10, pp. 1924–1928, Oct. 2020.
- [16] D. L. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, Apr. 2006.
- [17] J. Lee, G.-T. Gil, and Y. H. Lee, "Channel estimation via orthogonal matching pursuit for hybrid MIMO systems in millimeter wave communications," *IEEE Trans. Commun.*, vol. 64, no. 6, pp. 2370–2386, Jun. 2016.
- [18] Y. You and L. Zhang, "Bayesian matching pursuit-based channel estimation for millimeter wave communication," *IEEE Commun. Lett.*, vol. 24, no. 2, pp. 344–348, Feb. 2020.
- [19] J. A. Tropp and A. C. Gilbert, "Signal recovery from random measurements via orthogonal matching pursuit," *IEEE Trans. Inf. Theory*, vol. 53, no. 12, pp. 4655–4666, Dec. 2007.
- [20] W. Zhao, Y. You, L. Zhang, X. You, and C. Zhang, "OMPL-SBL algorithm for intelligent reflecting surface-aided mmWave channel estimation," *IEEE Trans. Veh. Technol.*, 2023, to be published.
- [21] S.-F. Hsieh, K. R. Liu, and K. Yao, "A unified square-root-free approach for QRD-based recursive-least-squares estimation," *IEEE Trans. Signal Process.*, vol. 41, no. 3, pp. 1405–1409, Mar. 1993.
- [22] X. Ge, F. Yang, H. Zhu, X. Zeng, and D. Zhou, "An efficient FPGA implementation of orthogonal matching pursuit with square-root-free QR decomposition," *IEEE Trans. VLSI Syst.*, vol. 27, no. 3, pp. 611–623, Mar. 2019.
- [23] S.-T. Xia, X.-J. Liu, Y. Jiang, and H.-T. Zheng, "Deterministic constructions of binary measurement matrices from finite geometry," *IEEE Trans. Signal Process.*, vol. 63, no. 4, pp. 1017–1029, Feb. 2015.
- [24] W. Lu, T. Dai, and S.-T. Xia, "Binary matrices for compressed sensing," *IEEE Trans. Signal Process.*, vol. 66, no. 1, pp. 77–85, Jan. 2018.
- [25] M. Fardad, S. M. Sayedi, and E. Yazdian, "A low-complexity hardware for deterministic compressive sensing reconstruction," *IEEE Trans. Circuits Syst. I*, vol. 65, no. 10, pp. 3349–3361, Oct. 2018.