# FLAMES Insight

Wuqiong Zhao

September 18, 2023

This document provides additional insights into the FLAMES (**F**lexible **L**inear **A**lgebra with **M**atrix-**E**mpowered **S**ynthesis) library for Vitis HLS [1]. The FLAMES library is open source at https://github.com/autohdw/flames, and the C++ API documentation is available at https://flames.autohdw.com.

# 1   Coverage Comparisons With Armadillo

Though we employ the concept of class-based interfaces from Armadillo [2], the syntax and coverage of FLAMES is not directly comparable with Armadillo.

## 1.1   Similarities and Differences

### 1.1.1   Similarities

1. Both the FLAMES library and the Armadillo library utilize modern C++ features in their code implementation;

2. Both the FLAMES library and the Armadillo library primarily focus on constructing classes and provide numerous member functions (methods) which significantly **simplify code** and **enhance readability**;

3. FLAMES has implemented matrix basics in Armadillo, including classes, operators, member functions, and other related matrix operations.

### 1.1.2   Differences

1. FLAMES is a C++ library for HLS implementation of hardware, while Armadillo is a C++ library for software run on CPU;

2. Though FLAMES is written in C++, it must be synthesizable, i.e., supported by the HLS tool (Vitis HLS [1]);

3. Writing HLS in C++ restricts the usage of class inheritance and virtual functions, and there is a lack of return value optimization (RVO) in Vitis HLS;

4. FLAMES does not include functions that potentially require dynamic memory allocation, such as the sparse matrix class (`SpMat`), as dynamic memory can not be synthesized by Vitis HLS.

## 1.2   HLS C++ Obstacles

Writing HLS in C++ is *substantially different* from C++ software programming with the CPU target. These differences hinder any HLS library implementation to cover most of the Armadillo library for software designs.

The library could have been more easily implemented with class inheritance and virtual functions (especially for `MatView`), however, the synthesis support for virtual functions has been removed in recent Vitis HLS releases. Another headache is the missing return value optimization (RVO) in Vitis HLS, therefore, returning a `Mat` object will inevitably lead to an unneeded copy process.

## 1.3   Coverage Comparisons

FLAMES has implemented almost all basics in Armadillo[1], including classes, operators, member functions, and other matrix operations. Nevertheless, FLAMES is still readily extendible and can work with existing algorithm implementations. The contents marked red in the following tables are later added, thanks to the Reviewer.

Table 1 shows the basic class implementations. The sparse matrix class (`SpMat`) is not implemented here, because it can potentially require dynamic memory (for example, change the number of non-zero elements), which is not synthesizable. It is recommended that users design the implementation to avoid use a sparse matrix representation. This feature may be added in the future should users make a reasonable request.

---

[1]Documentation website: https://arma.sourceforge.net/docs.html.

**Table. 1.** Coverage comparisons of matrix, vector, cube classes.

| Armadillo [2] | FLAMES | Remarks |
|:---:|:---:|:---|
| Mat | Mat | Dense matrix class. |
| Col | Col | Dense column vector class. |
| Row | Row | Dense row vector class. |
| Cube | Row | Dense cube class. |
| SpMat | N/A | Sparse matrix class. |
| const Mat& | MatView | Read-only access to a matrix. |
| Mat& | MatRef | Writable reference to a matrix. |

Table 2 shows the defined operators. It is worth noting that `==` and `!=` are not used for the matrix equality check. Additionally, we also provided the `<<` operator to print the matrix (only in C++ simulation, not used in synthesis).

**Table. 2.** Coverage comparisons of operators.

| Armadillo [2] | FLAMES | Remarks |
|:---:|:---:|:---|
| + (Unary) | + (Unary) | Positive sign. |
| + (Binary) | + (Binary) | Addition of two objects. |
| – (Unary) | – (Unary) | Negative sign. |
| – (Binary) | – (Binary) | Subtraction of two objects. |
| * | * | Matrix multiplication. |
| % | % | Matrix element-wise (Hadamard) multiplication. |
| / | / | Matrix element-wise division. |
| == | == | Element-wise equality evaluation. |
| != | != | Element-wise inequality evaluation. |
| >=, <=, >, < | >=, <=, >, < | Element-wise comparisons. |
| &&, \|\| | &&, \|\| | Element-wise AND and OR logic. |

Methods for main matrix member functions are shown in Table 3. Since FLAMES is for hardware implementation, where dynamic memory allocation is not allowed, all operations that modify the size of a matrix are not supported.

**Table. 3.** Coverage comparisons of main matrix member functions.

| Armadillo [2] | FLAMES | Remarks |
|---|---|---|
| .zeros | .setZero | Set all elements to zero. |
| .fill | .setValue | Set all elements to a specified value. |
| .memptr | .rawDataPtr | Raw pointer to memory. |
| [] | [] | Access an element (viewed on 1D). |
| (,) | (,) | Access an element (viewed on 2D). |
| .col | .col/.col_ | Column (as a copy/view). |
| .cols | .cols/.cols_ | Columns (as a copy/view). |
| .row | .row/.row_ | Row (as a copy/view). |
| .rows | .rows/.rows_ | Rows (as a copy/view). |
| .t | .t/.t_ | Transpose (as a copy/view). |
| .i | .invNSA | Matrix inverse (NSA is hardware friendly)[2]. |
| .diagvec | .diagVec | Take the diagonal and return a vector. |
| .diagmat | .diagMat | Generate a diagonal matrix. |
| .norm(,"fro") | .power | Power of a matrix (square of $\ell_2$-norm). |
| .print | .print | Print the matrix (not for synthesis). |

Basic matrix operations, including `mul`, `add`, `sub`, `innerProd` are also supported by FLAMES.

FLAMES does not currently cover other functions, including decomposition, factorization, and statics. However, an **interface** to existing algorithm IPs can be easily implemented, by getting the raw data pointer with the `.rawDataPtr` method.

# 2 FLAMES Simplicity

## 2.1 Code Length Comparisons

We use a design methodology for general-purpose DSP as a baseline [3] (the HLS book: https://kastner.ucsd.edu/hlsbook/) and compare it with our code implementation using the FLAMES library. There is a mapping relationship between the design methodology in the reference and our baseline (implementations of operations, e.g., matrix multiplication,

and addition, can be found in [3]). It is worth noting that through our comparison, achieving the same hardware efficiency and performance as the FLAMES library often requires longer code implementation in the baseline approach.

### 2.1.1 Matrix-Vector Multiplication

For multiplication between matrix `A` and vector `b`, with the result as `c`, FLAMES has shorter code than the traditional HLS [3]. This example can also be accessed at https://github.com/autohdw/flames/tree/master/examples/mat-vec-multiplication.

With FLAMES.

```
1   void top(const M& A, const V& b, V& c) { c = A * b; }
```

Without FLAMES. [3]

```
1   void top(dtype A[4][4], dtype b[4], dtype c[4]) {
2   #pragma HLS ARRAY_PARTITION variable = A complete
3   #pragma HLS ARRAY_PARTITION variable = b complete
4       for (size_t i = 0; i != 4; ++i) {
5           for (size_t j = 0; j != 4; ++j) {
6   #pragma HLS UNROLL
7               c[i] = A[i][j] * b[j];
8           }
9       }
10  }
```

### 2.1.2 Neumann Series Approximation (NSA) Inverse

This example can also be accessed at https://github.com/autohdw/flames/tree/master/examples/mat-inv-nsa.

FLAMES provides a built-in Neumann series approximation (NSA) inversion function `.invNSA()`.

With FLAMES.

```
1   #include "flames/flames.hpp"
2
3   using dtype = FxP<8, 8>;
4   using M     = Mat<dtype, 4, 4>;
5
6   M top(const M& A) { return A.invNSA(); }
7
8   int main() {
9       M A{ 10, -2, 1, 0, 1, -8, 2, 0, 0, 0, 11, -1, 0, 1, 2, 4 };
10      A.print("A = ");
11      M A_inv = top(A);
```

```
12      A_inv.print("A_inv = ");
13      return 0;
14  }
```

Even without using the `.invNSA()` function, the implementation of the algorithm can be done with only a few lines of code.

With FLAMES but not directly use `.invNSA()`.

```
1   #include "flames/flames.hpp"
2
3   using dtype = FxP<8, 8>;
4   using M     = Mat<dtype, 4, 4>;
5
6   M top(const M& A) {
7       M A_inv;
8       const auto D = A.diagMat_(); // diagonal part
9       const auto E = A.offDiag_(); // off-diagonal part
10      Mat<dtype, 4, 4, MatType::DIAGONAL> D_inv;
11      D_inv.invDiag(D); // inverse of diagonal part
12      Mat<dtype, 4, 4, MatType::NORMAL> product = (-D_inv) * E;
13      Mat<dtype, 4, 4, MatType::NORMAL> sum_tmp = A_inv = product; // the first
        ↪ iteration
14      Mat<dtype, 4, 4, MatType::NORMAL> tmp;
15      const size_t iter = 4;
16  MAT_INV_NSA:
17      for (size_t i = 1; i < iter; ++i) {
18          tmp.mul(A_inv, product);
19          A_inv = tmp;
20          sum_tmp += tmp;
21      }
22      A_inv.mul(sum_tmp, D_inv);
23      return A_inv += D_inv;
24  }
25
26  int main() {
27      M A{ 10, -2, 1, 0, 1, -8, 2, 0, 0, 0, 11, -1, 0, 1, 2, 4 };
28      A.print("A = ");
29      M A_inv = top(A);
30      A_inv.print("A_inv = ");
31      return 0;
32  }
```

Directly using built-in functions provided by Vitis HLS [1] and the design method in [3] without FLAMES, the required code is much longer, and the readability is much poorer.

Without FLAMES.

```
1   #include "ap_fixed.h"
2   #include <iostream>
3   #include <string>
```

6

```cpp
4
5  using dtype = ap_fixed<17, 8>;
6
7  void print(dtype A[4][4], std::string s = "") {
8  #ifndef __SYNTHESIS__
9      std::cout << s << "[";
10     for (size_t i = 0; i + 1 < 4; ++i) {
11         std::cout << "[";
12         for (size_t j = 0; j + 1 < 4; ++j) std::cout << A[i][j] << ", ";
13         std::cout << A[i][3] << "]," << std::endl;
14     }
15     std::cout << "[";
16     for (size_t j = 0; j + 1 < 4; ++j) std::cout << A[3][j] << ", ";
17     std::cout << A[3][3] << "]]" << std::endl;
18 #endif
19 }
20
21 void mat_copy(dtype from[4][4], dtype to[4][4]) {
22 #pragma HLS INLINE
23 #pragma HLS ARRAY_PARTITION variable = from complete
24 #pragma HLS ARRAY_PARTITION variable = to complete
25     for (size_t i = 0; i != 4; ++i) {
26 #pragma HLS UNROLL
27         for (size_t j = 0; j != 4; ++j) {
28 #pragma HLS LOOP_FLATTEN
29             to[i][j] = from[i][j];
30         }
31     }
32 }
33
34 void top(dtype A[4][4], dtype A_inv[4][4]) {
35 #pragma HLS ARRAY_PARTITION variable = A complete
36 #pragma HLS ARRAY_PARTITION variable = A_inv complete
37     dtype product[4][4];
38     dtype D_inv[4];
39     for (size_t i = 0; i != 4; ++i) {
40 #pragma HLS UNROLL
41         D_inv[i] = static_cast<dtype>(1) / A[i][i];
42     }
43 MAT_DIAG_TIMES_MAT_NORMAL:
44     for (size_t j = 0; j != 4; ++j) {
45 #pragma HLS UNROLL
46         for (size_t i = 0; i != 4; ++i) {
47 #pragma HLS LOOP_FLATTEN
48             product[i][j] = i == j ? static_cast<dtype>(0) :
49                 static_cast<dtype>(-D_inv[i] * A[i][j]);
            }
50     }
51     dtype sum_tmp[4][4], tmp[4][4];
52     mat_copy(product, sum_tmp);
```

```
53    mat_copy(product, A_inv);
54    const size_t iter = 4;
55    for (size_t i = 1; i < iter; ++i) {
56    // tmp = A_inv * product;
57    GEMM:
58        for (size_t _i = 0; _i != 4; ++_i) {
59        GEMM_r:
60            for (size_t r = 0; r != 4; ++r) {
61 #pragma HLS UNROLL
62            GEMM_c:
63                for (size_t c = 0; c != 4; ++c) {
64 #pragma HLS LOOP_FLATTEN
65                    if (_i == 0) tmp[r][c] = static_cast<dtype>(0);
66                    tmp[r][c] += A_inv[r][_i] * product[_i][c];
67                }
68            }
69        }
70        mat_copy(tmp, A_inv);
71        // sum_tmp += tmp;
72        for (size_t _i = 0; _i != 4; ++_i) {
73 #pragma HLS UNROLL
74            for (size_t j = 0; j != 4; ++j) {
75 #pragma HLS LOOP_FLATTEN
76                sum_tmp[_i][j] += tmp[_i][j];
77            }
78        }
79    }
80 // A_inv = sum_tmp * D_inv;
81 MAT_NORMAL_TIMES_MAT_DIAG:
82    for (size_t i = 0; i != 4; ++i) {
83 #pragma HLS UNROLL
84        for (size_t j = 0; j != 4; ++j) {
85 #pragma HLS LOOP_FLATTEN
86            A_inv[i][j] = sum_tmp[i][j] * D_inv[j];
87        }
88    }
89    // A_inv += D_inv;
90    for (size_t i = 0; i != 4; ++i) {
91 #pragma HLS UNROLL
92        A_inv[i][i] += D_inv[i];
93    }
94 }
95
96 int main() {
97    dtype A[4][4] = { { 10, -2, 1, 0 }, { 1, -8, 2, 0 }, { 0, 0, 11, -1 }, {
       ↪ 0, 1, 2, 4 } };
98    print(A, "A = ");
99    dtype A_inv[4][4];
100    top(A, A_inv);
101    print(A_inv, "A_inv = ");
```

```
102        return 0;
103    }
```

## 2.2   Simplicity Beyond Code Length

Writing HLS C++ code with FLAMES enhances code readability by providing modular and reusable components. FLAMES encapsulate complex functionalities into pre-built functions, allowing designers to focus on higher-level logic and abstraction. This makes the code **better organized**, and **less redundant**. It also facilitates *code maintainability* by making the code easier to comprehend and modify.

# 3   Task-Level Pipelining

Designs **can** leverage task-level pipelining with FLAMES. Auto pipelining is applied in the case study by HLS.

Since the task-level pipelining optimization is provided by the HLS tool (Vitis HLS, for example), there is no significant difference between using FLAMES or not in this term. If anything, FLAMES provides a clearer task flow, where common matrix operations are optimized and managed, which is conducive to task-level pipelining. For instance, the (false) data dependency issue can be avoided by using FLAMES.

To demonstrate the availability of task-level pipelining optimization, an additional example is provided here, which can also be accessed at https://github.com/autohdw/flames/tree/master/examples/task-level-pipelining. In this example, main_task can be pipelined by HLS. It is worth noting that the function main_task needs to be marked as an INLINE function so as to achieve a good pipeline result. The schedule viewer for the top function provided by Vitis HLS 2022.2 is shown in Fig. 1.

```cpp
1  #include "flames/flames.hpp"
2
3  using dtype = FxP<6, 2>;
4  using M = Mat<dtype, 4, 4>;
5  using V = Vec<dtype, 4>;
6
7  void main_task(const M& A, const V& b, V& c) {
8      #pragma HLS INLINE
9      M tmp1;
10     V tmp2;
11     tmp1 = A * A;
12     tmp2 = tmp1 * b;
13     c = tmp2 % tmp2;
14 }
15
16 void top(const M& A1, const M& A2, const M& A3, const V& b, V& c) {
17     V c1, c2;
```

```
18      M tmp;
19      #pragma HLS PIPELINE
20      main_task(A1, b, c1);
21      main_task(A2, c1, c2);
22      main_task(A3, c2, c);
23  }
24
25  int main() {
26      M A1, A2, A3;
27      V b, c;
28      top(A1, A2, A3, b, c);
29      return 0;
30  }
```
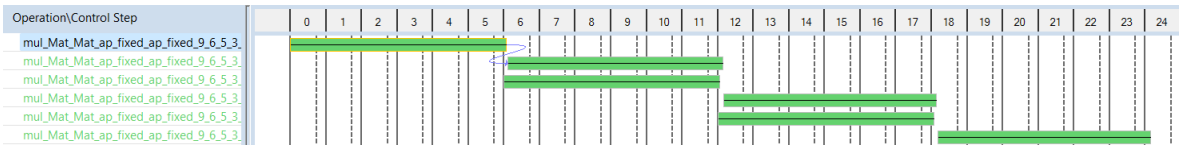


**Fig. 1.** Task-level pipelining result shown in the schedule viewer.

# References

[1]  Xilinx, "Vitis high-level synthesis user guide (UG1399)," Accessed: Feb. 27, 2023, 2023. [Online]. Available: https://docs.xilinx.com/r/en-US/ug1399-vitis-hls.

[2]  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.

[3]  R. Kastner, J. Matai, and S. Neuendorffer, "Parallel programming for FPGAs," *arXiv:1805.03648*, 2018. [Online]. Available: https://arxiv.org/abs/1805.03648.