

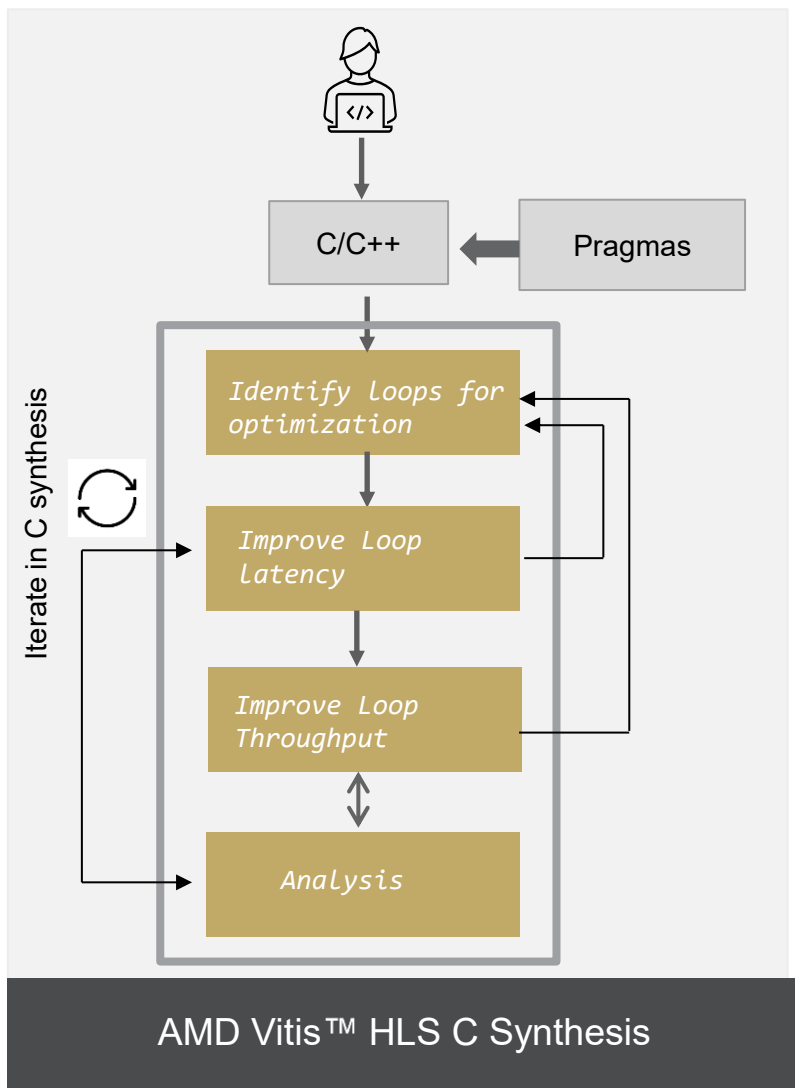
Streamline the AMD Vitis™ HLS Workflow with the Performance Pragma

AMD Vitis HLS

Agenda

- Traditional HLS Optimization
- Challenges of Traditional HLS Optimization
- Simplifying Optimization with Performance Pragma
- Performance Pragma Methodology
 - Key differences with the traditional methodology
- Performance Pragma in Action: Convolution Design
- IDE Enhancements
- Summary & References
- Q&A

Traditional HLS Optimization



- Optimization driven by manual pragma insertion
- Map system performance goals with pragmas inserted manually
- Top-level/system throughput inferred (not explicitly modeled)

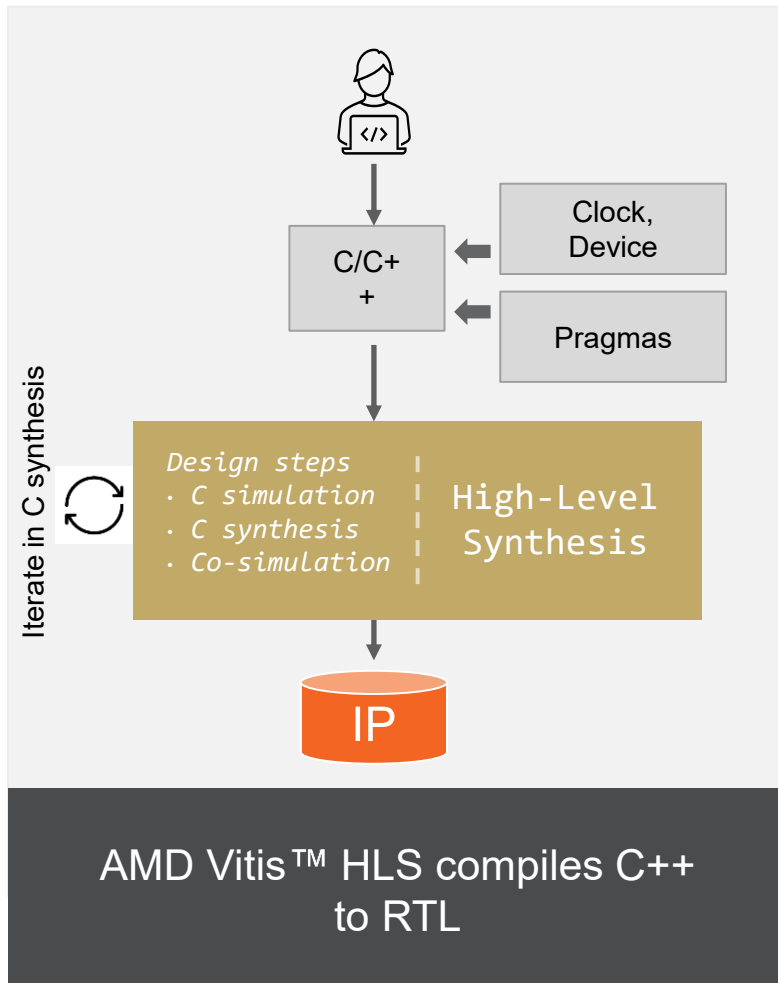
```
void xfv(...)
...
#pragma HLS ARRAY_PARTITION variable=row_ind complete
#pragma HLS ARRAY_PARTITION variable=OutputValues complete
#pragma HLS ARRAY_PARTITION variable=OutputValues1 complete
#pragma HLS ARRAY_PARTITION variable=src_buf complete
#pragma HLS ARRAY_PARTITION variable=src_buf complete
#pragma HLS ARRAY_PARTITION variable=buf complete
for (...)
#pragma HLS pipeline
#pragma HLS LOOP_FLATTEN
for (col = 0; col<img_width>>+)
#pragma HLS UNROLL
for (..)
#pragma HLS UNROLL
```

MODULES & LOOPS	ISSUE	LATENCY(CYCLES)	ITERATION LATENCY	INTERVAL	TRIP COUNT	PIPELINED
Filter2DKernel (5)		2,087,211		2,087,068		dataflow
entry_proc		0		0		no
ReadFromMem (2)	Hi	2,074,007		2,074,007		no
Window2D (1)		2,087,067		2,087,067		no
Filter2D (2)	Hi	2,073,858		2,073,858		no
WriteToMem (1)	Hi	2,073,676		2,073,676		no

The screenshot shows a list of memory accesses on the left, including 'LineBuffer_13_addr', 'LineBuffer_13_load', 'icmp_In158', 'col_ptr_1', 'col_ptr_2', 'col_ptr_write_In123', 'new_pixel', 'LineBuffer_12_addr', 'LineBuffer_12_load', 'reuse_addr_reg_load', 'addr_cmp', 'LineBuffer_13_addr_write_In155', 'reuse_addr_reg_write_In161', 'LineBuffer_11_addr', and 'LineBuffer_11_load'. To the right is a timing diagram with horizontal bars representing dataflow and vertical lines representing control signals.

Optimization Driven by Manual Pragma Insertion!!!

Challenges of Traditional HLS Optimization



- In the regular workflow, achieving performance goals **requires expertise in choosing the right combination of pragmas...**

- Pragma-based designs can be **less flexible and sensitive to changes**
- Vision library color detection uses over 40+ classic pragmas
 - unroll, pipeline, flatten...
- Any changes to throughput goals could affect these 40+ pragmas...

Choosing the Right Combination of Pragmas can be a Challenge!!!

Simplifying Optimization with Performance Pragma



AMD Performance Pragma simplifies HLS optimization



Allows users to define a high-level throughput goal



Shifts manual pragma selection optimization burden to the compiler



Enables Automatic Pragma Generation: No more manual pragma guessing!



Intelligently infers and applies optimizations (pipelining, unrolling, etc.)



Offers flexible throughput control via target specification



Tool automatically determines optimal pragma configuration



Represents a new, higher-level way to constrain design throughput



Provides a more intuitive and efficient path to desired performance

Performance Pragma

Performance Pragma can be applied to a top-level function and/or individual loops

Top-Level Performance Pragma

- Defines a design-wide throughput goal
- Guides the compiler to optimize the entire design
- Automatically infers and applies loop-level pragmas based on analysis

Loop-Level Performance Pragma

- Targets specific loops for local control
- Can be automatically inferred based on top-level performance pragma or manually applied
- Enables fine-grain optimization, infers classic pragmas (pipeline, unroll, etc...)

Benefits

New top-level performance goal representation	Helps achieve top-level performance pragma targets
Design-wide analysis to validate performance goals	Precise control of loop behavior

Performance Pragma Methodology

Offers a streamlined approach to guide the HLS optimization process to achieve optimal performance

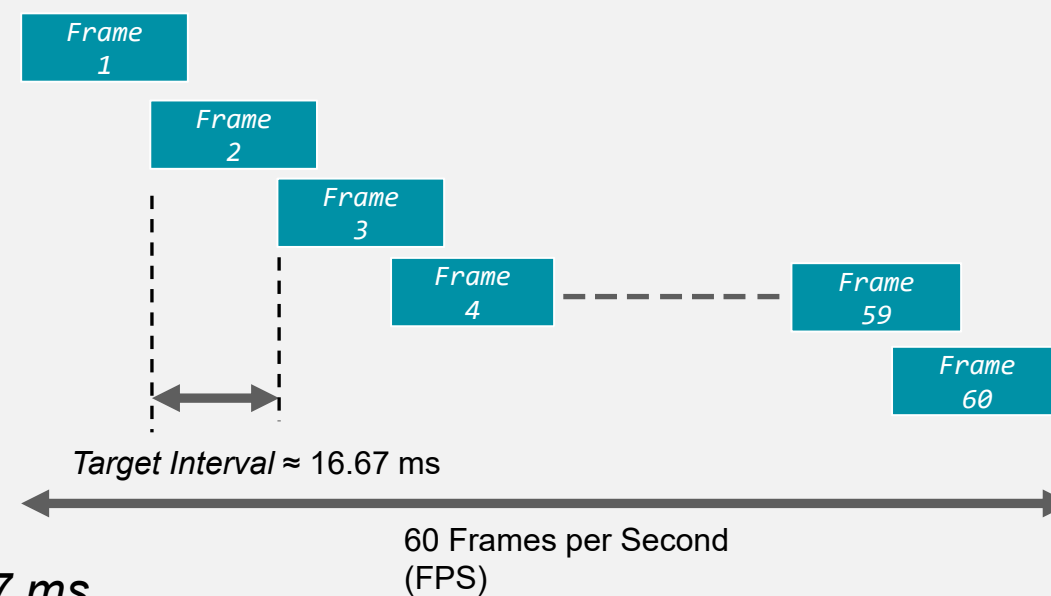


Consider a video application aiming for a frame rate of 60 frames per second (60 fps)

Determine the top-level performance target (`target_ti`)

- To achieve the 60 FPS, the top-level function must be ready to process a new frame within 1/60th of a second
- Implies:

$$\text{target_ti} = 1 / \text{throughput}(\text{FPS}) = 1 / 60 \text{ seconds} \approx 16.67 \text{ ms}$$



Performance Pragma Methodology

Offers a streamlined approach to guide the HLS optimization process to achieve optimal performance



Re-architect the Code for Dataflow

- Design necessitates a re-architecture for an explicit load-compute-store (LCS) modeling compatible with the dataflow pragma
- Allows AMD Vitis™ HLS to effectively:
 - Optimize the code
 - Exploit task-level parallelism

Performance Pragma Methodology

Offers a streamlined approach to guide the HLS optimization process to achieve optimal performance in the user's hardware implementation



Run C Simulation and Determine Loop Trip Counts

- That metric is crucial since the performance pragma algorithm requires loop budgeting...
- By default, variable loop bounds are assumed to be "1024," which can be inaccurate...
- For variable loop bounds, provide trip count information via the `loop_tripcount max=N pragma...`

Performance Pragma Methodology

Offers a streamlined approach to guide the HLS optimization process to achieve optimal performance in the user's hardware implementation



Add the Top-Level Performance Target

- Apply the desired performance goal using the top-level performance pragma

```
#pragma HLS performance target_ti = 16.67 ms
```

Add/Update Local Performance Targets

- Identify performance-critical loops after C synthesis
- Specify loop-level performance targets for these critical loops
- Iteratively refine & re-synthesize until the design meets performance

Performance Pragma Methodology: Key Differences

Comparing the traditional approach (manual pragma insertion) to the top-level performance pragma workflow...

Top-Level Throughput Constraint

Unlike starting optimization at individual loop level, you define a system-wide performance target first

Needs Trip Count for Dynamic Loops

Provides the tool with crucial information for accurate performance estimation, especially for loops with variable iterations

May also need Loop-Level Performance Pragma

While the tool automates, you can still fine-tune specific bottlenecks for more granular control

Performance Pragma in Action: Convolution Design

Convolution 2D: Calculate the Performance Target

Step 1



Step 2



Step 3



Step 4

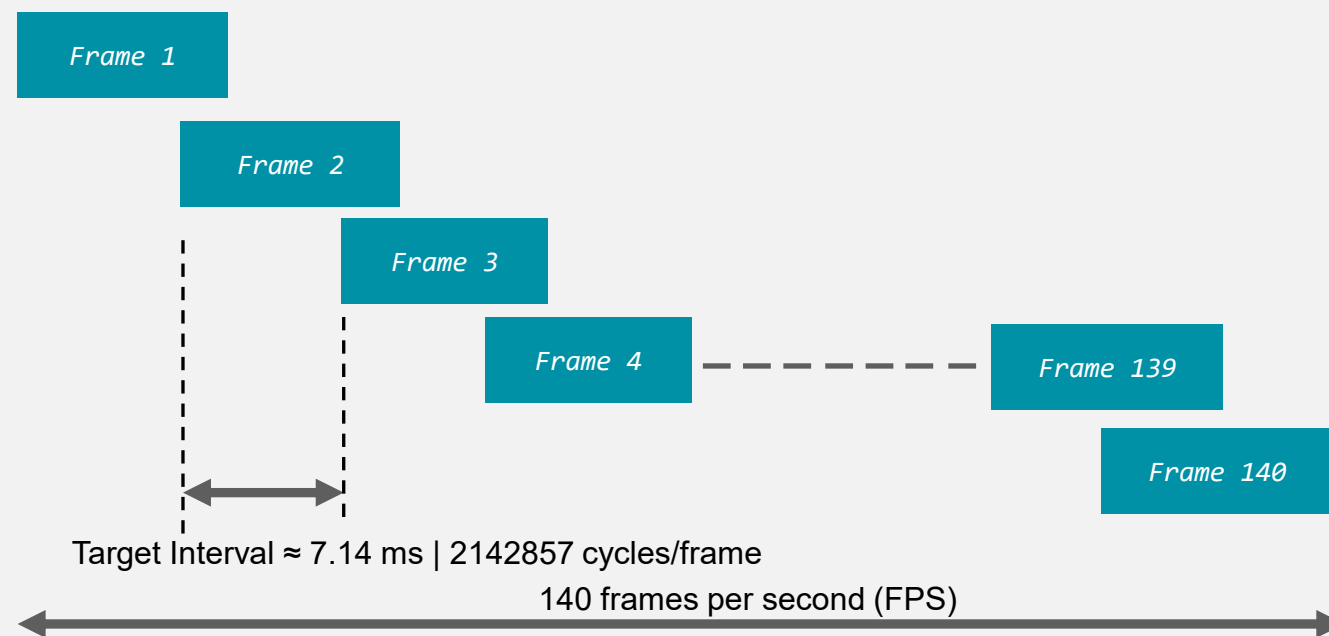
Calculate the Performance Target Based on Throughput Goal

- Goal: Convolution function process an HD 140 frames per second @ 300 MHz clock...

Determine the Top-Level Performance Target (`target_ti`)

Target Interval (`target_ti`) can be expressed in time (ms) or number of cycles:

- Time: 140 frames/sec hence: `target_ti`
(milliseconds/frame) = $1000 / 140 = 7.14$ ms
- Cycles: 140 frames/sec at 300 MHz hence:
 - $\text{target_ti} = (\text{kernel freq.}) / (\text{throughput}) =$
 $(300 * 10^6 \text{ cycles/second}) / 140 =$
 $2,142,857.14$ cycles/frame



Convolution 2D: Re-architect the Code for Dataflow

Step 1



Step 2



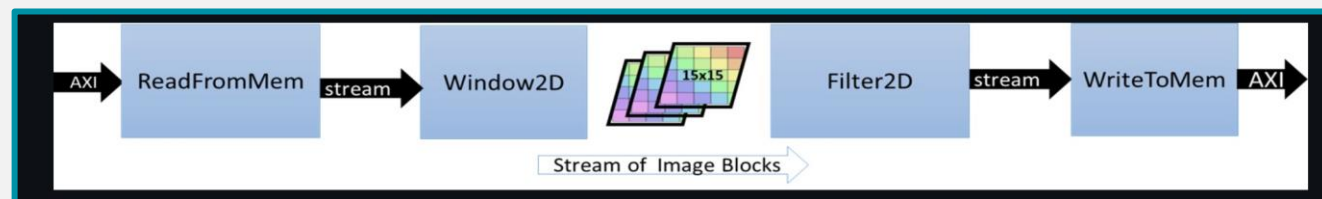
Step 3



Step 4

Re-architect the Code for Dataflow

- Design necessitates a re-architecture to make load-compute-store explicit and apply the dataflow pragma



```

void Filter2DKernel(
    const char      coeffs[256],
    float           factor,
    short           bias,
    unsigned short   width,
    unsigned short   height,
    unsigned short   stride,
    const unsigned char src[MAX_IMAGE_WIDTH*MAX_IMAGE_HEIGHT],
    unsigned char     dst[MAX_IMAGE_WIDTH*MAX_IMAGE_HEIGHT])
{
    #ifdef STEP1
    #pragma HLS performance target_ti = 2142857
    #endif
    #pragma HLS dataflow
  
```

Convolution 2D: Run C Simulation with Code Analyzer

Step 1**Step 2****Step 3****Step 4**

Run C Simulation and Determine Loop Trip Counts

- Run C simulation with Code Analyzer
- Add trip counts for dynamic variable loops

```
    for (int x = 0; x < width; x++)  
    {  
#pragma HLS LOOP_TRIPCOUNT max=1920  
        // Read a 2D window of pixels  
        window w = window_stream.read();
```

Convolution 2D: Add the Top/Loop-Level Performance Target (1/4)

Step 1



Step 2



Step 3



Step 4

Add the Top-Level Performance Target

- Apply the top-level performance pragma using target_ti

```
void Filter2DKernel(  
    const char        coeffs[256],  
    float             factor,  
    short             bias,  
    unsigned short    width,  
    unsigned short    height,  
    unsigned short    stride,  
    const unsigned char src[MAX_IMAGE_WIDTH*MAX_IMAGE_HEIGHT],  
    unsigned char      dst[MAX_IMAGE_WIDTH*MAX_IMAGE_HEIGHT])  
{  
  
    #ifdef STEP1  
    #pragma HLS performance target_ti = 2142857  
    #endif  
    #pragma HLS dataflow
```


Convolution 2D: Add the Top/Loop-Level Performance Target (2/4)



Identify Bottleneck Loops (if any)

- Here Window2D does not meet the performance target of 2,142,857.14 cycles/frame

MODULES & LOOPS	ISSUE	TARGET TI(CYCLES)	ESTIMATED TI(CYCLES)
<div> <div>Filter2DKernel (5)</div> <div> <div>entry_proc</div> <div>ReadFromMem (2)</div> <div> <div>Window2D (1)</div> <div> <div>update_window (2)</div> <div>Outline_VITIS_LOOP_141_1</div> <div>Window2D_Pipeline_VITIS_LOOP_149_3 (1)</div> <div>Filter2D (2)</div> <div>WriteToMem (1)</div> </div> </div> </div> </div>	<div> <div>ti</div> <div>ti</div> <div></div> <div></div> <div></div> <div>ti</div> <div>ti</div> </div>	<div> <div>2,142,857</div> <div>2,073,873</div> <div>532,196,988</div> <div>532,196,985</div> <div>1</div> <div></div> <div>575,394</div> <div>2,073,612</div> </div>	<div> <div>45,915,038</div> <div>2,074,007</div> <div>45,915,037</div> <div>45,915,034</div> <div>1</div> <div></div> <div>2,073,858</div> <div>2,073,676</div> </div>

```

void Window2D(
    unsigned short width,
    unsigned short height,
    hls::stream<U8> &pixel_stream,
    hls::stream<window> &window_stream)
{
    // Line buffers - used to store [FILTER_V_SIZE-1] entire
    U8 LineBuffer[FILTER_V_SIZE-1][MAX_IMAGE_WIDTH];

    // Sliding window of [FILTER_V_SIZE][FILTER_H_SIZE] pixels
    window Window;

    unsigned col_ptr = 0;
    unsigned ramp_up = width*((FILTER_V_SIZE-1)/2)+(FILTER_V_SIZE-1);
    unsigned num_pixels = width*height;
    unsigned num_iterations = num_pixels + ramp_up;
    const unsigned max_iterations = MAX_IMAGE_WIDTH*MAX_IMAGE_HEIGHT;

    // Iterate until all pixels have been processed
    update_window: for (int n=0; n<num_iterations; n++)
    {
        #pragma HLS LOOP_TRIPCOUNT max=2087047
    }
}
    
```

Convolution 2D: Add the Top/Loop-Level Performance Target (3/4)

Step 1



Step 2



Step 3



Step 4

Add/Update Local Performance Targets

- Apply loop-level performance pragmas
 - Specify loop-level performance pragmas for the Window2D function to achieve target_ti

```
update_window: for (int n=0; n<num_iterations; n++)
{
    #pragma HLS LOOP_TRIPCOUNT max=2087047

#ifdef STEP2
    MET=no TARGET=2142857, ESTIMATE=45915034 (CYCLES)
    #pragma HLS performance target_ti = 2142857
#endif
```

Convolution 2D: Add the Top/Loop-Level Performance Target (4/4)

Target Still Not Reached

MODULES & LOOPS	ISSUE	TARGET TI(CYCLES)	ESTIMATED TI(CYCLES)	PERFORMANCE CONSTRAINT
Filter2DKernel (5)		2,142,857	45,915,038	pragma
entry_proc				
> ReadFromMem (2)		2,073,873	2,074,007	
Window2D (1)		532,196,988	45,915,037	
> update_window (2)		2,142,857	45,915,034	pragma
Outline_VITIS_LOOP_141_1		1	1	
Window2D_Pipeline_VITIS_LOOP_149_3 (1)				
VITIS_LOOP_149_3		26	13	inferred
> Filter2D (2)		575,394	2,073,858	
> WriteToMem (1)		2,073,612	2,073,676	

C-synth report with top and loop-level Performance pragma

Solution

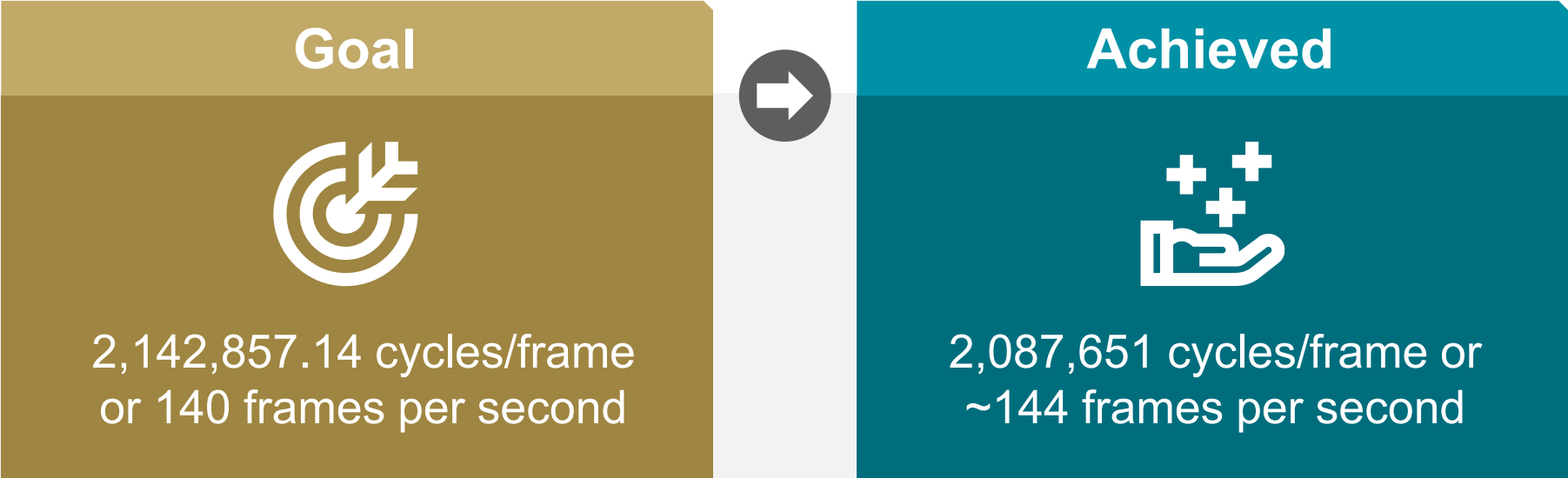
Set a target loop TI = 1 for this shift loop to nudge the compiler to execute the loop in a single cycle

```

for(int i = 0; i < FILTER_V_SIZE-2; i++) {
#ifdef STEP3
    #pragma HLS performance target_ti= 1
#endif
    LineBuffer[i][col_ptr] = LineBuffer[i+1][col_ptr];
}
    
```

Voila!! Meets performance!

Convolution 2D: Results



Category	Using Classic Pragmas	Using Performance Pragma
Target Interval	2,087,651	2,087,651
Optimizations Pragmas in the Design	8	3 (2.6X fewer pragmas)

Limitations of Performance Pragma in v2025.1

Features	Behavior/Limitation
<code>ap_cint</code>	The tool exits with an explicit warning message
<code>ap_(u)int</code> / <code>ap_(u)fixed</code>	<ul style="list-style-type: none"> No performance models are inaccurate
Big constant arrays of <code>ap_int</code> / <code>ap_fixed</code>	<ul style="list-style-type: none"> Using them with macro <code>NON_C99STRING</code> will result in a compilation error
<code>std::complex<ap_fixed></code>	Lead to a compiler error on Windows
HLS IP blocks (FFT, FIR, ...), <code>hls_math.h</code> , <code>ap_wait</code> , <code>hls::vector</code> , <code>ap_axis/ap_axiu</code>	<ul style="list-style-type: none"> No performance models are inaccurate
<code>hls::stream_of_blocks</code> , <code>hls::task</code> , <code>hls::split</code> / <code>hls::merge</code> , <code>hls::print</code> , <code>hls::half</code> , <code>ap_utils.h</code> , <code>hls_fpo.h</code> , <code>ap_float</code> , RTL Blackboxes, OpenCL, <code>hls::burst_maxi</code> , <code>hls::fence</code> , <code>hls::directio</code>	The tool exits with an explicit warning message
Deprecated HLS pragmas	Assertion failure
Function pipeline with sub loop(s)	Assertion failure

IDE Enhancements

Redesigned Table for Synthesis Report

The screenshot shows a table titled 'Performance & Resource Estimates' with a toolbar above it. The toolbar includes icons for expand/collapse, list view, hierarchy view, info, warning, error, refresh (annotated with '2'), percentage, and search (annotated with '1'). A filter icon and a 'Modules' checkbox (annotated with '3') are on the right. The table has columns: MODULES & LOOPS, ISSUE, TARGET TI(CYCLES), ESTIMATED TI(CYCLES), PERFORMANCE CONSTRAINT, LATENCY(CYCLES), ITERATION LATENCY, INTERVAL, TRIP COUNT, PIPELINED, and STRUCTURE. The data rows are: Filter2DKernel (5), entry_proc, ReadFromMem (2) (annotated with '4'), Window2D (1), Filter2D (2), and WriteToMem (1). The 'PERFORMANCE CONSTRAINT' column shows a 'pragma' for the top row.

MODULES & LOOPS	ISSUE	TARGET TI(CYCLES)	ESTIMATED TI(CYCLES)	PERFORMANCE CONSTRAINT	LATENCY(CYCLES)	ITERATION LATENCY	INTERVAL	TRIP COUNT	PIPELINED	STRUCTURE
Filter2DKernel (5)		2,142,857	2,087,068	pragma	2,087,211		2,087,068		dataflow	function
entry_proc					0		0		no	function
ReadFromMem (2)		2,073,873	2,074,007		2,074,007		2,074,007		no	function
Window2D (1)		532,196,988	2,087,067		2,087,067		2,087,067		no	function
Filter2D (2)		575,394	2,073,858		2,073,858		2,073,858		no	function
WriteToMem (1)		2,073,612	2,073,676		2,073,676		2,073,676		no	function

1. Updated table in response to new performance pragma
2. Richer report visualization: Instantly toggle cycles/ns for clearer insights
3. Greater filtering capability
4. Icons and colorization for csynth_design results condensation

Quickly perform detailed performance and resource analysis with new table

Pragma Overlays in Source Editor

- (1) - HLS performance pragma (top-level function)
- (2) – (4) – Inferred pragmas from (1)

```
matrix_vector_base.c > ...
1  #define SIZE 8
2  typedef int BaseType;
3
4  void matrix_vector(BaseType M[SIZE][SIZE], BaseType V_In[SIZE], BaseType V_Out[SIZE]) {
5      BaseType i, j;
6      #pragma HLS performance target_ti=SIZE*SIZE 1
7      data_loop:
8          #pragma HLS performance target_ti=40 | #pragma HLS performance target_ti=40 | #pragma HLS pipeline ii=5
9          for (i = 0; i < SIZE; i++) {
10             BaseType sum = 0;
11             dot_product_loop:
12                 #pragma HLS performance target_ti=1 | #pragma HLS unroll complete
13                 for (j = 0; j < SIZE; j++) {
14                     sum += V_In[j] * M[i][j];
15                 }
16             V_Out[i] = sum;
17         }
18     }
```

Review optimizations made by the performance pragma...

Summary



Performance Pragma simplifies complex HLS optimization by enabling users to define a high-level throughput goal, shifting the optimization burden to the compiler for automatic pragma generation and application of key transformations, offering flexible throughput control and a more efficient path to desired hardware performance.

Q&A

References

1. [Vitis™ HLS Introductory Examples](#)
2. [Vitis HLS Doc](#)
3. Vitis HLS Product Page - [Vitis HLS](#)
4. [Vitis HLS tutorial](#)

Disclaimers

The information contained herein is for informational purposes only and is subject to change without notice. While every precaution has been taken in the preparation of this document, it may contain technical inaccuracies, omissions and typographical errors, and AMD is under no obligation to update or otherwise correct this information. Advanced Micro Devices, Inc. makes no representations or warranties with respect to the accuracy or completeness of the contents of this document, and assumes no liability of any kind, including the implied warranties of noninfringement, merchantability or fitness for particular purposes, with respect to the operation or use of AMD hardware, software or other products described herein. No license, including implied or arising by estoppel, to any intellectual property rights is granted by this document. Terms and limitations applicable to the purchase or use of AMD products are as set forth in a signed agreement between the parties or in AMD's Standard Terms and Conditions of Sale. GD-18u.

© 2025 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, Vitis, and combinations thereof are trademarks of Advanced Micro Devices, Inc. Other product names used in this publication are for identification purposes only and may be trademarks of their respective owners. Certain AMD technologies may require third-party enablement or activation. Supported features may vary by operating system. Please confirm with the system manufacturer for specific features. No technology or product can be completely secure.

