

Optimizing Arrays for Performance Demo Script

Introduction

This demonstration script provides high-level instructions on how to remove bottlenecks caused by arrays in a C design.

Preparation:

- Required files: Necessary files are located at *C:\training\optimize_array_performance\demo*
- Required hardware: None
- Supporting materials: None

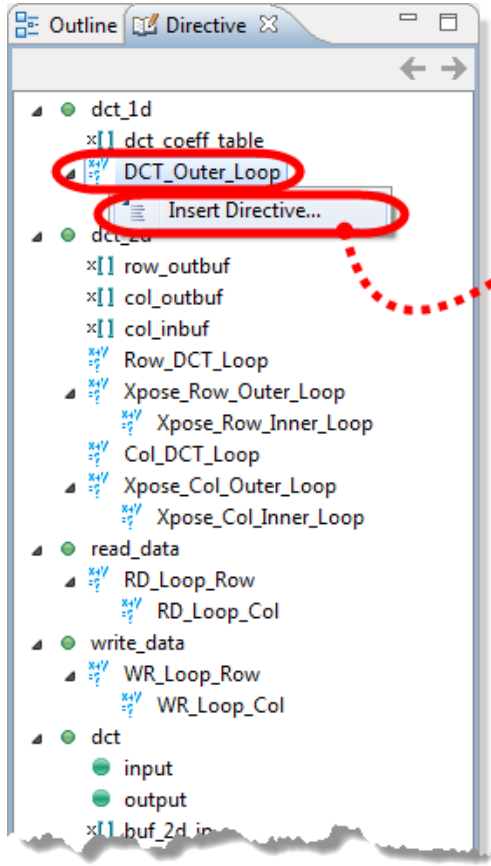
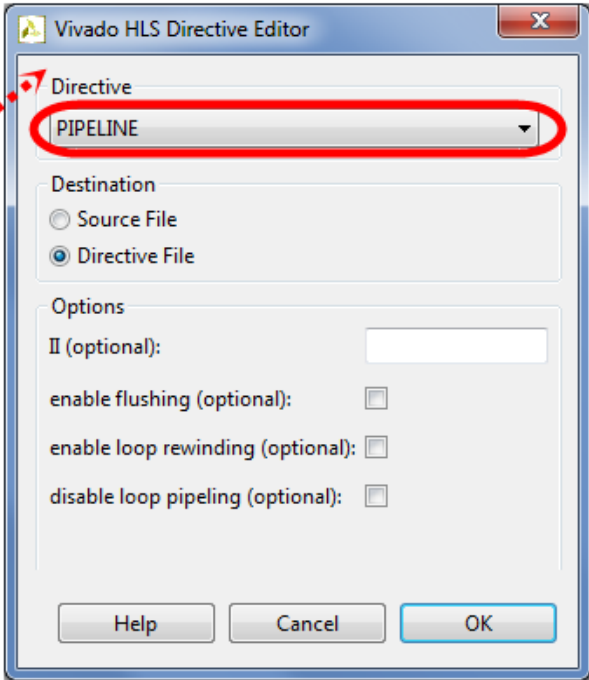
Optimizing Arrays for Performance

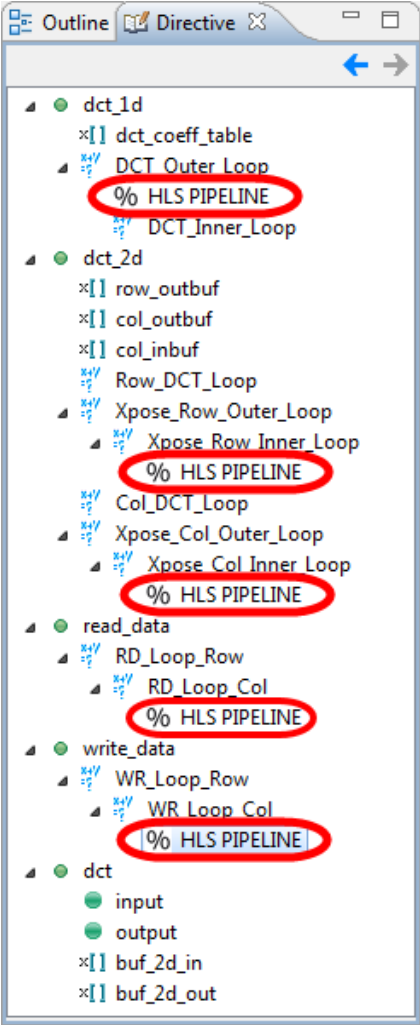
Action with Description	Point of Emphasis and Key Takeaway
<ul style="list-style-type: none">• Launch the Vivado® HLS tool.• Open the provided dct_prj Vivado HLS tool project located at: <i>C:\training\optimize_array_performance\demo\dct_prj</i>	You can open existing Vivado HLS tool projects from the Vivado HLS tool Welcome page.

Action with Description	Point of Emphasis and Key Takeaway
<ul style="list-style-type: none"> Access and review the source files (<i>dct.c</i> and <i>dct.h</i>) from the Explorer pane. 	<p>This C design uses a discrete cosine transformation (DCT). The function implements a 2D DCT algorithm by first processing each row of the input array via a 1D DCT, then processing the columns of the resulting array through the same 1D DCT. It calls the <i>read_data</i>, <i>dct_2d</i>, and <i>write_data</i> functions.</p> <p>The <i>read_data</i> function is defined at line 54 and consists of two loops: RD_Loop_Row and RD_Loop_Col.</p> <p>The <i>write_data</i> function is defined at line 66 and consists of two loops to perform writing the result. The <i>dct_2d</i> function, defined at line 23, calls the <i>dct_1d</i> function and performs transpose.</p> <p>Finally, the <i>dct_1d</i> function, defined at line 4, uses <i>dct_coeff_table</i> and performs the required function by implementing a basic iterative form of the 1D Type II DCT algorithm.</p>

Action with Description	Point of Emphasis and Key Takeaway
<ul style="list-style-type: none"> Run C synthesis. Review the Synthesis report. 	<p>Once synthesis completes, the Synthesis report will open in the main viewing area.</p> <p>Notice that the estimated clock period is within the requested clock period.</p> <p>The Synthesis report also contains latency and throughput information of the design. The results correspond to the default solution (without any directives). You can further reduce the numbers down by specifying your requirements via directives.</p> <p>The Synthesis log is available in the Console pane.</p>
<p>What is the worst-case latency of the design?</p> <p>Answer: 2935</p>	
<ul style="list-style-type: none"> Create a new solution named <i>solution2</i> (select Project > New Solution). Accept the default settings and click Finish. 	<p>Creating a new solution allows for different optimizations to be compared easily.</p> <p>There is no need to copy directives from the previous solution because the previous solution does not have any directives. Even if the directives are copied (the default setting), there would be no impact to the demo since there are no directives in the initial solution.</p>
<p>As part of the optimization process, as seen in the "Pipeline for Performance" demo, you will begin by pipelining the loops in the design.</p>	

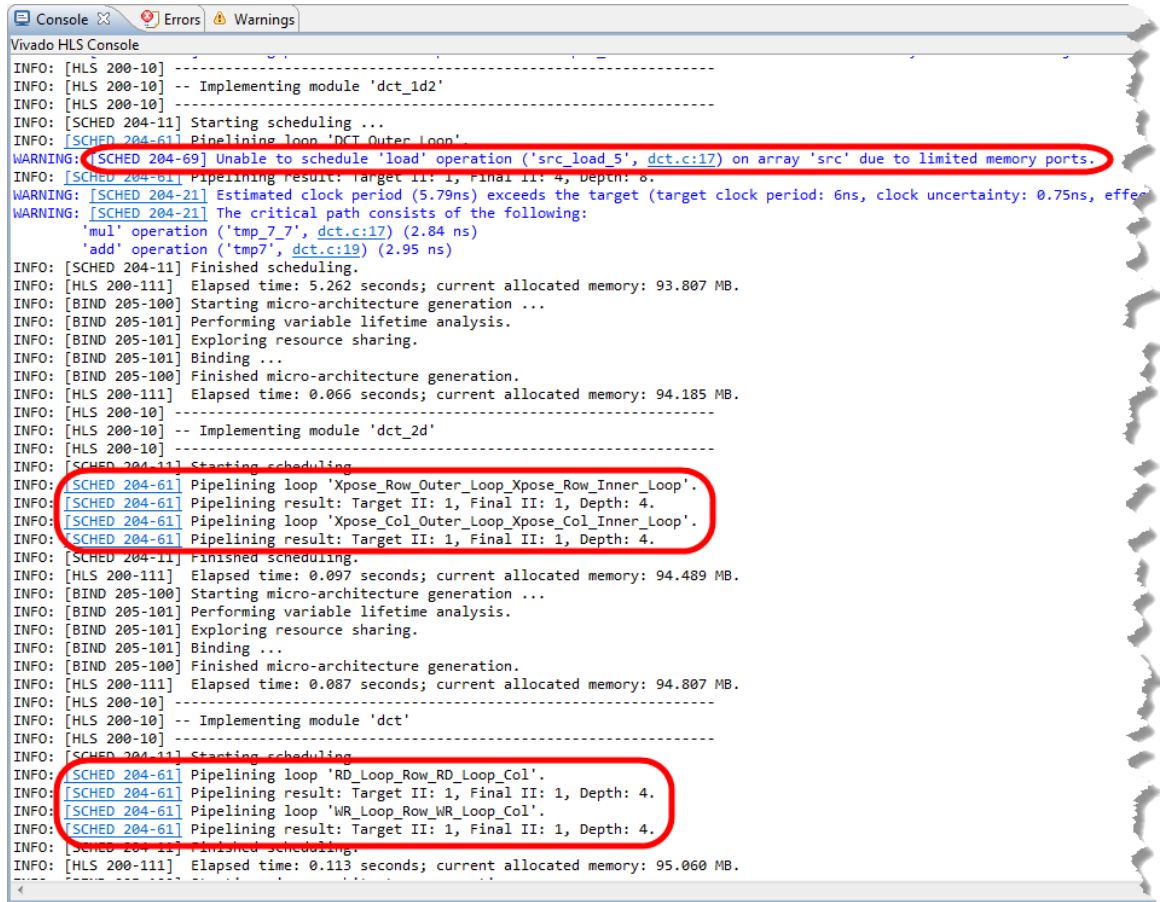
Action with Description	Point of Emphasis and Key Takeaway
<ul style="list-style-type: none"> Apply the PIPELINE directive on <i>DCT_Outer_Loop</i> of the <i>dct_1d</i> function (shown below). 	<p>Apply the PIPELINE directive on the outer loop.</p> <p>Moving the PIPELINE directive from the inner loop to the outer loop of <i>dct_1d</i> will lead to more parallelism of the multiply and add operations.</p> <p>That is, eight (8) multiply and add operations are performed concurrently, thus minimizing the number of cycles required to compute each value in the output array.</p> <p>Leave the II field blank since the design tries to target II as 1; i.e., it will try to optimize the loop to accept a new input for every cycle.</p> <p>You will find the directive written to the <i>directive.tcl</i> file under the <i>dct_prj</i> > <i>solution2</i> > <i>constraints</i> folder.</p> <pre>set_directive_pipeline "dct_1d/DCT_Outer_Loop"</pre>

Action with Description	Point of Emphasis and Key Takeaway
 <p>The screenshot shows the Vivado Outline window with a project hierarchy. The 'DCT_Outter_Loop' is highlighted with a red circle, and the 'Insert Directive...' option is also highlighted with a red circle. A red dotted arrow points from the 'Insert Directive...' option to the 'Vivado HLS Directive Editor' window.</p>	 <p>The screenshot shows the Vivado HLS Directive Editor window. The 'Directive' dropdown menu is set to 'PIPELINE' and is highlighted with a red circle. The 'Destination' section shows 'Directive File' selected. The 'Options' section has three checkboxes: 'enable flushing (optional)', 'enable loop rewinding (optional)', and 'disable loop pipeling (optional)', all of which are unchecked.</p>

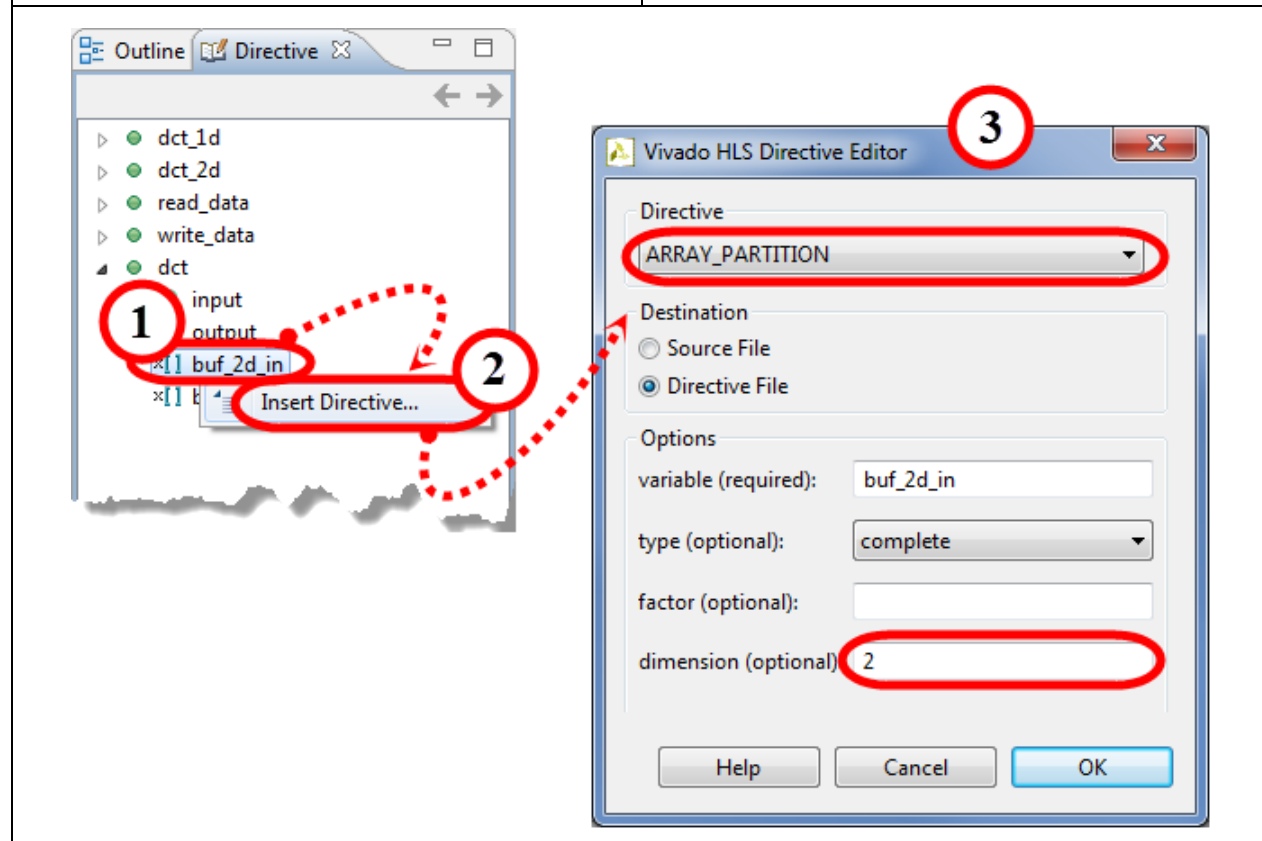
Action with Description	Point of Emphasis and Key Takeaway
<ul style="list-style-type: none"> Similarly, apply the PIPELINE directive to the following loops: <ul style="list-style-type: none"> <i>Xpose_Row_Inner_Loop</i> of the <i>dct_2d</i> function <i>Xpose_Col_Inner_Loop</i> of the <i>dct_2d</i> function <i>RD_Loop_Col</i> of the <i>read_data</i> function <i>WR_Loop_Col</i> of the <i>write_data</i> function 	<p>The Directive tab should look like the figure below after you finish applying the PIPELINE directive.</p>  <p>The screenshot shows the 'Directive' tab in the IDE. The code hierarchy is as follows:</p> <ul style="list-style-type: none"> dct_1d <ul style="list-style-type: none"> dct_coeff_table DCT_Outer_Loop <ul style="list-style-type: none"> % HLS PIPELINE DCT_Inner_Loop dct_2d <ul style="list-style-type: none"> row_outbuf col_outbuf col_inbuf Row_DCT_Loop Xpose_Row_Outer_Loop <ul style="list-style-type: none"> Xpose_Row_Inner_Loop <ul style="list-style-type: none"> % HLS PIPELINE Col_DCT_Loop Xpose_Col_Outer_Loop <ul style="list-style-type: none"> Xpose_Col_Inner_Loop <ul style="list-style-type: none"> % HLS PIPELINE read_data <ul style="list-style-type: none"> RD_Loop_Row <ul style="list-style-type: none"> RD_Loop_Col <ul style="list-style-type: none"> % HLS PIPELINE write_data <ul style="list-style-type: none"> WR_Loop_Row <ul style="list-style-type: none"> WR_Loop_Col <ul style="list-style-type: none"> % HLS PIPELINE dct <ul style="list-style-type: none"> input output buf_2d_in buf_2d_out
<ul style="list-style-type: none"> Run C synthesis. 	<p>Once synthesis completes, the Synthesis report will open in the main viewing area.</p>
<ul style="list-style-type: none"> Compare the results of two solutions (<i>solution1</i> and <i>solution2</i>). 	<p>This allows you to compare the different optimizations of the project.</p> <p>You should see the comparison report as shown below.</p>

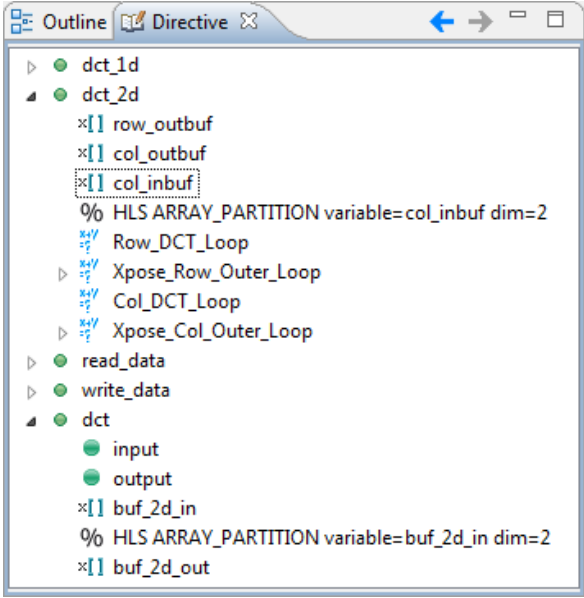
Action with Description	Point of Emphasis and Key Takeaway																																																			
<div><div>Performance Estimates</div><div><div>Timing (ns)</div><table><tr><td>Clock</td><td></td><td>solution2</td><td>solution1</td></tr><tr><td>ap_clk</td><td>Target</td><td>10.00</td><td>10.00</td></tr><tr><td></td><td>Estimated</td><td>7.769</td><td>3.770</td></tr></table></div><div><div>Latency (clock cycles)</div><table><tr><td></td><td></td><td>solution2</td><td>solution1</td></tr><tr><td>Latency</td><td>min</td><td>843</td><td>2935</td></tr><tr><td></td><td>max</td><td>843</td><td>2935</td></tr><tr><td>Interval</td><td>min</td><td>843</td><td>2935</td></tr><tr><td></td><td>max</td><td>843</td><td>2935</td></tr></table></div><div>Utilization Estimates</div><table><tr><td></td><td>solution2</td><td>solution1</td></tr><tr><td>BRAM_18K</td><td>5</td><td>5</td></tr><tr><td>DSP48E</td><td>8</td><td>1</td></tr><tr><td>FF</td><td>491</td><td>246</td></tr><tr><td>LUT</td><td>1364</td><td>964</td></tr><tr><td>URAM</td><td>0</td><td>0</td></tr></table></div>		Clock		solution2	solution1	ap_clk	Target	10.00	10.00		Estimated	7.769	3.770			solution2	solution1	Latency	min	843	2935		max	843	2935	Interval	min	843	2935		max	843	2935		solution2	solution1	BRAM_18K	5	5	DSP48E	8	1	FF	491	246	LUT	1364	964	URAM	0	0	
Clock		solution2	solution1																																																	
ap_clk	Target	10.00	10.00																																																	
	Estimated	7.769	3.770																																																	
		solution2	solution1																																																	
Latency	min	843	2935																																																	
	max	843	2935																																																	
Interval	min	843	2935																																																	
	max	843	2935																																																	
	solution2	solution1																																																		
BRAM_18K	5	5																																																		
DSP48E	8	1																																																		
FF	491	246																																																		
LUT	1364	964																																																		
URAM	0	0																																																		
<ul style="list-style-type: none">What is the worst-case latency of the design?	Answer: 843																																																			
<ul style="list-style-type: none">Go to the Utilization Estimates section and note the number of DSP48E and block RAMs used to implement <i>solution2</i>.	Answer: Number of BRAM_18K: 5 Number of DSP48E: 8																																																			

Action with Description	Point of Emphasis and Key Takeaway
<ul style="list-style-type: none">Select the Console tab and review the synthesis information.	<p>From the Synthesis log, note that the design was not able to achieve the requested II on <i>DCT_Outer_Loop</i> because of the limited memory ports on the <i>src</i> element.</p> <p>The <i>src</i> is input to the <i>dct_1d</i> function and <i>dct_1d</i> is called twice in the <i>dct_2d</i> function (line 33 and line 44).</p> <p>At line 33, <i>in_block</i> is accessed via the <i>src</i> element in <i>dct_1d</i>. At line 44, <i>col_inbuf</i> is accessed via the <i>src</i> element in <i>dct_1d</i>.</p> <p>Therefore, you will need to partition both the <i>col_inbuf</i> and <i>in_block</i> arrays to achieve a throughput of 1.</p>

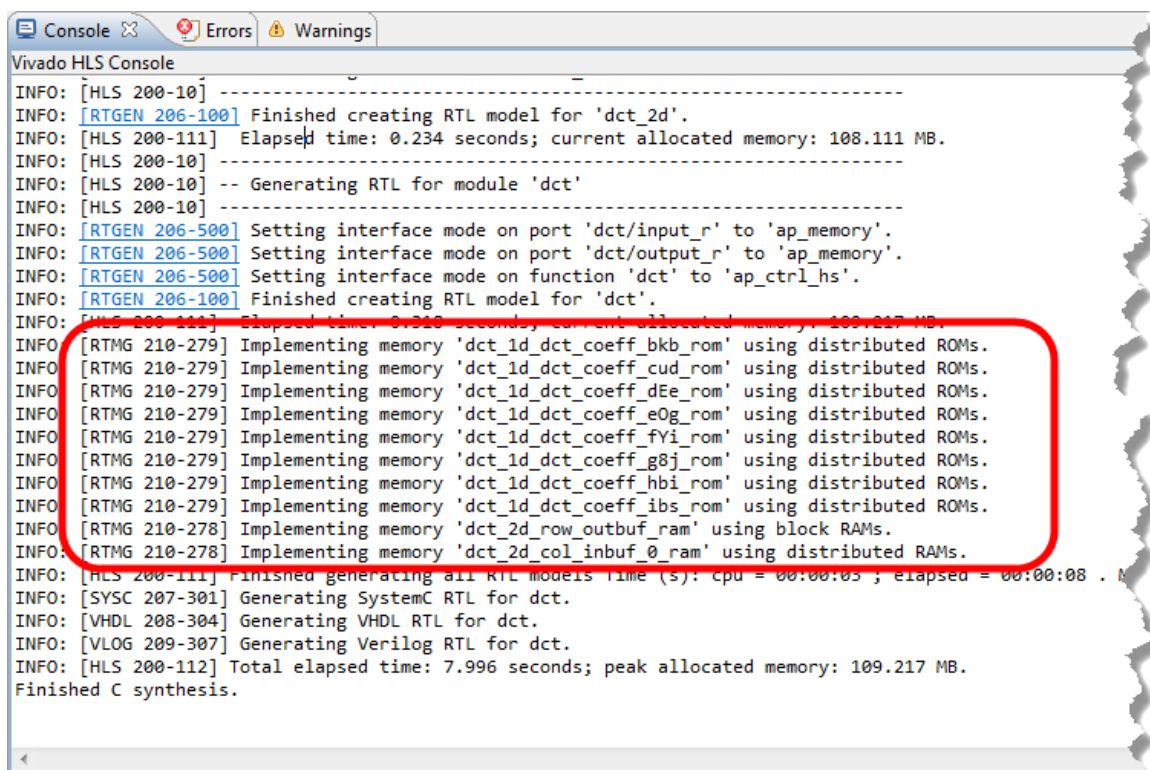
Action with Description	Point of Emphasis and Key Takeaway
 <pre> Vivado HLS Console INFO: [HLS 200-10] ----- INFO: [HLS 200-10] -- Implementing module 'dct_1d2' INFO: [HLS 200-10] ----- INFO: [SCHED 204-11] Starting scheduling ... INFO: [SCHED 204-61] Pipelining loop 'DCT Outer Loop'. WARNING: [SCHED 204-69] Unable to schedule 'load' operation ('src_load_5', dct.c:17) on array 'src' due to limited memory ports. INFO: [SCHED 204-61] Pipelining result: Target II: 1, Final II: 4, Depth: 8. WARNING: [SCHED 204-21] Estimated clock period (5.79ns) exceeds the target (target clock period: 6ns, clock uncertainty: 0.75ns, effective clock period: 5.04ns). WARNING: [SCHED 204-21] The critical path consists of the following: 'mul' operation ('tmp_7_7', dct.c:17) (2.84 ns) 'add' operation ('tmp7', dct.c:19) (2.95 ns) INFO: [SCHED 204-11] Finished scheduling. INFO: [HLS 200-111] Elapsed time: 5.262 seconds; current allocated memory: 93.807 MB. INFO: [BIND 205-100] Starting micro-architecture generation ... INFO: [BIND 205-101] Performing variable lifetime analysis. INFO: [BIND 205-101] Exploring resource sharing. INFO: [BIND 205-101] Binding ... INFO: [BIND 205-100] Finished micro-architecture generation. INFO: [HLS 200-111] Elapsed time: 0.066 seconds; current allocated memory: 94.185 MB. INFO: [HLS 200-10] ----- INFO: [HLS 200-10] -- Implementing module 'dct_2d' INFO: [HLS 200-10] ----- INFO: [SCHED 204-11] Starting scheduling ... INFO: [SCHED 204-61] Pipelining loop 'Xpose_Row_Outer_Loop_Xpose_Row_Inner_Loop'. INFO: [SCHED 204-61] Pipelining result: Target II: 1, Final II: 1, Depth: 4. INFO: [SCHED 204-61] Pipelining loop 'Xpose_Col_Outer_Loop_Xpose_Col_Inner_Loop'. INFO: [SCHED 204-61] Pipelining result: Target II: 1, Final II: 1, Depth: 4. INFO: [SCHED 204-11] Finished scheduling. INFO: [HLS 200-111] Elapsed time: 0.097 seconds; current allocated memory: 94.489 MB. INFO: [BIND 205-100] Starting micro-architecture generation ... INFO: [BIND 205-101] Performing variable lifetime analysis. INFO: [BIND 205-101] Exploring resource sharing. INFO: [BIND 205-101] Binding ... INFO: [BIND 205-100] Finished micro-architecture generation. INFO: [HLS 200-111] Elapsed time: 0.087 seconds; current allocated memory: 94.807 MB. INFO: [HLS 200-10] ----- INFO: [HLS 200-10] -- Implementing module 'dct' INFO: [HLS 200-10] ----- INFO: [SCHED 204-11] Starting scheduling ... INFO: [SCHED 204-61] Pipelining loop 'RD_Loop_Row_RD_Loop_Col'. INFO: [SCHED 204-61] Pipelining result: Target II: 1, Final II: 1, Depth: 4. INFO: [SCHED 204-61] Pipelining loop 'WR_Loop_Row_WR_Loop_Col'. INFO: [SCHED 204-61] Pipelining result: Target II: 1, Final II: 1, Depth: 4. INFO: [SCHED 204-11] Finished scheduling. INFO: [HLS 200-111] Elapsed time: 0.113 seconds; current allocated memory: 95.060 MB. </pre>	
<p>You will now solve the <i>dct_1d</i> pipeline II problem by increasing the memory bandwidth available to it.</p> <p>This will be done by partitioning the arrays from which the <i>dct_1d</i> inner loops read data (<i>in_block</i> in <i>Row_DCT_Loop</i> and <i>col_inbuf</i> in <i>Col_DCT_Loop</i>).</p>	
<ul style="list-style-type: none"> • Create a new solution named <i>solution3</i>. • Accept the default settings and click Finish. 	<p>In this solution, you will apply the ARRAY_PARTITION directive to <i>buf_2d_in</i> of the <i>dct</i> function and <i>col_inbuf</i> of the <i>dct2d</i> function.</p>

Action with Description	Point of Emphasis and Key Takeaway
<ul style="list-style-type: none"> Apply the ARRAY_PARTITION directive to <i>buf_2d_in</i> of the <i>dct</i> function as shown in the figure below. 	Partitioning large arrays into multiple smaller arrays or into individual registers can help improve access to data and remove block RAM bottlenecks.



Action with Description	Point of Emphasis and Key Takeaway
<ul style="list-style-type: none"> Similarly, apply the ARRAY_PARTITION directive <i>col_inbuf</i> of the <i>dct_2d</i> function. 	<p>The Directive tab should look like the figure below after you finish applying the ARRAY_PARTITION directive.</p> 
<ul style="list-style-type: none"> Run C synthesis. 	<p>Once synthesis completes, the synthesis report will open in the main viewing area.</p>
<p>Examine the Synthesis log. Has the PIPELINE II directive been met?</p> <p>Yes, the pipeline directive met II=1. The memory bandwidth increased via the array partitioning and thus the design met the requested II value.</p>	
<ul style="list-style-type: none"> Compare the results of the two solutions (<i>solution2</i> and <i>solution3</i>). 	<p>You should see the comparison report as shown below.</p>

Action with Description	Point of Emphasis and Key Takeaway																																																		
<div><div>Performance Estimates</div><div><div>Timing (ns)</div><table><tr><td>Clock</td><td></td><td>solution3</td><td>solution2</td></tr><tr><td>ap_clk</td><td>Target</td><td>10.00</td><td>10.00</td></tr><tr><td></td><td>Estimated</td><td>7.517</td><td>7.769</td></tr></table></div><div><div>Latency (clock cycles)</div><table><tr><td></td><td></td><td>solution3</td><td>solution2</td></tr><tr><td>Latency</td><td>min</td><td>477</td><td>843</td></tr><tr><td></td><td>max</td><td>477</td><td>843</td></tr><tr><td>Interval</td><td>min</td><td>477</td><td>843</td></tr><tr><td></td><td>max</td><td>477</td><td>843</td></tr></table></div><div><div>Utilization Estimates</div><table><tr><td></td><td>solution3</td><td>solution2</td></tr><tr><td>BRAM_18K</td><td>3</td><td>5</td></tr><tr><td>DSP48E</td><td>8</td><td>8</td></tr><tr><td>FF</td><td>844</td><td>491</td></tr><tr><td>LUT</td><td>1879</td><td>1364</td></tr><tr><td>URAM</td><td>0</td><td>0</td></tr></table></div></div>		Clock		solution3	solution2	ap_clk	Target	10.00	10.00		Estimated	7.517	7.769			solution3	solution2	Latency	min	477	843		max	477	843	Interval	min	477	843		max	477	843		solution3	solution2	BRAM_18K	3	5	DSP48E	8	8	FF	844	491	LUT	1879	1364	URAM	0	0
Clock		solution3	solution2																																																
ap_clk	Target	10.00	10.00																																																
	Estimated	7.517	7.769																																																
		solution3	solution2																																																
Latency	min	477	843																																																
	max	477	843																																																
Interval	min	477	843																																																
	max	477	843																																																
	solution3	solution2																																																	
BRAM_18K	3	5																																																	
DSP48E	8	8																																																	
FF	844	491																																																	
LUT	1879	1364																																																	
URAM	0	0																																																	
<p>Latency was reduced to 477. Block RAM usage decreased to 3 from 5.</p> <p>Memory utilization will usually be more after array partitioning. But in this case, some of the memory elements were implemented in the distributed RAMs/ROMS (you can observe this in the Synthesis log in the Console tab) and hence the number of block RAMs was actually reduced compared to the previous solution.</p>																																																			

Action with Description	Point of Emphasis and Key Takeaway
 <pre> Vivado HLS Console INFO: [HLS 200-10] ----- INFO: [RTGEN 206-100] Finished creating RTL model for 'dct_2d'. INFO: [HLS 200-111] Elapsed time: 0.234 seconds; current allocated memory: 108.111 MB. INFO: [HLS 200-10] ----- INFO: [HLS 200-10] -- Generating RTL for module 'dct' INFO: [HLS 200-10] ----- INFO: [RTGEN 206-500] Setting interface mode on port 'dct/input_r' to 'ap_memory'. INFO: [RTGEN 206-500] Setting interface mode on port 'dct/output_r' to 'ap_memory'. INFO: [RTGEN 206-500] Setting interface mode on function 'dct' to 'ap_ctrl_hs'. INFO: [RTGEN 206-100] Finished creating RTL model for 'dct'. INFO: [HLS 200-112] Elapsed time: 0.234 seconds; current allocated memory: 108.111 MB. INFO: [RTMG 210-279] Implementing memory 'dct_1d_dct_coeff_bkb_rom' using distributed ROMs. INFO: [RTMG 210-279] Implementing memory 'dct_1d_dct_coeff_cud_rom' using distributed ROMs. INFO: [RTMG 210-279] Implementing memory 'dct_1d_dct_coeff_dEe_rom' using distributed ROMs. INFO: [RTMG 210-279] Implementing memory 'dct_1d_dct_coeff_eOg_rom' using distributed ROMs. INFO: [RTMG 210-279] Implementing memory 'dct_1d_dct_coeff_fYi_rom' using distributed ROMs. INFO: [RTMG 210-279] Implementing memory 'dct_1d_dct_coeff_g8j_rom' using distributed ROMs. INFO: [RTMG 210-279] Implementing memory 'dct_1d_dct_coeff_hbi_rom' using distributed ROMs. INFO: [RTMG 210-279] Implementing memory 'dct_1d_dct_coeff_ibs_rom' using distributed ROMs. INFO: [RTMG 210-278] Implementing memory 'dct_2d_row_outbuf_ram' using block RAMs. INFO: [RTMG 210-278] Implementing memory 'dct_2d_col_inbuf_0_ram' using distributed RAMs. INFO: [HLS 200-111] Finished generating all RTL models (s): cpu = 00:00:03 ; elapsed = 00:00:08 . INFO: [SYSC 207-301] Generating SystemC RTL for dct. INFO: [VHDL 208-304] Generating VHDL RTL for dct. INFO: [VLOG 209-307] Generating Verilog RTL for dct. INFO: [HLS 200-112] Total elapsed time: 7.996 seconds; peak allocated memory: 109.217 MB. Finished C synthesis. </pre>	

Summary

Memory elements may become bottlenecks when it comes to meeting throughput and latency requirements. In this demo, you learned how to apply a partition directive on arrays in the design and observed its impact on resources and throughput.

References:

- Supporting materials
 - Vivado Design Suite Tutorial: High-Level Synthesis* (UG871)
 - Vivado Design Suite User Guide: High-Level Synthesis* (UG902)