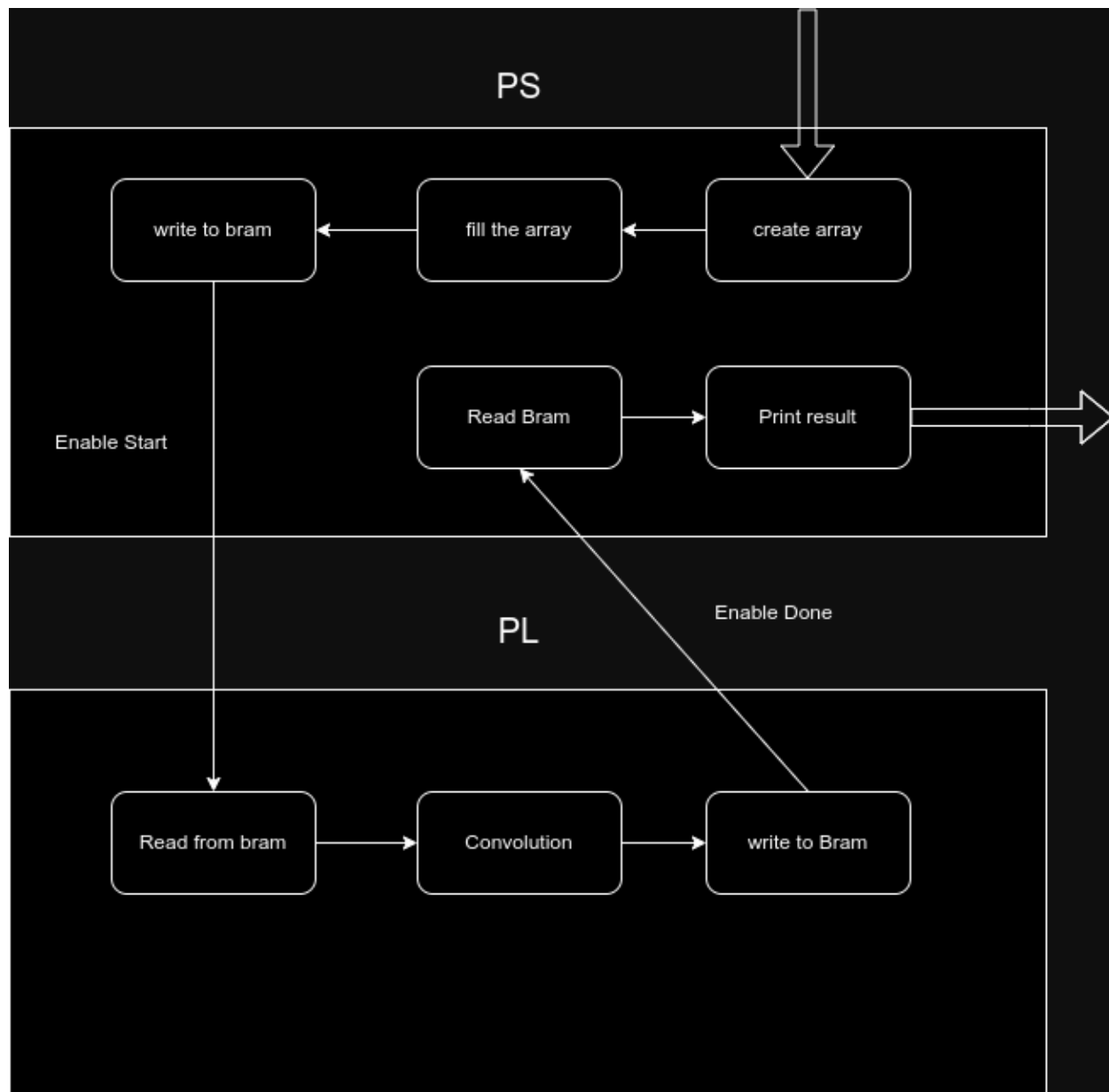


به نام خدا
تهیه کننده: ابراهیم صدیقی شماره دانشجویی: ۹۹۳۱۰۹۸

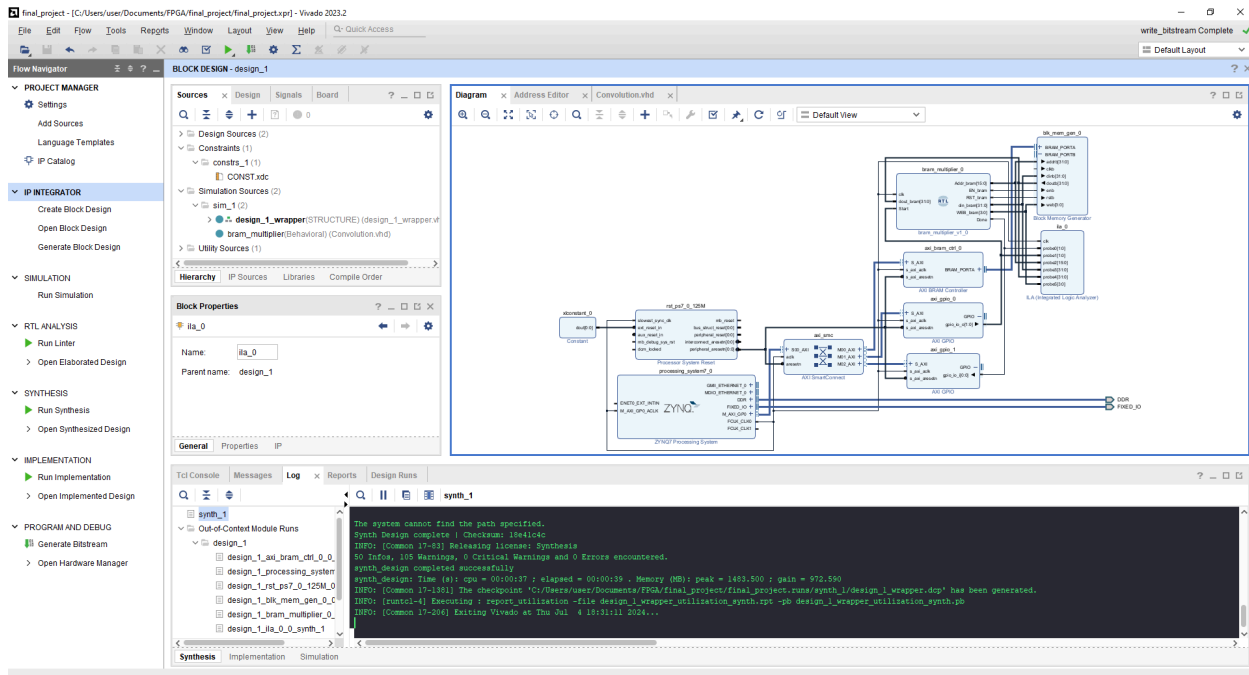
در ابتدا برای پیاده سازی یک شتاب دهنده ابتدا باید با استفاده از co-design قسمت های مختلف PL و PS را طراحی کنیم. نمودار کلی این طراحی به شکل زیر است:



شکل بالا دو قسمت PS و PL را اینگونه پیاده سازی کرده است که ابتدا یک آرایه اصلی و یک آرایه kernel در قسمت PS ساخته می شود، سپس این آرایه ها مقاردهی می شوند و بعد از آن، آرایه ها به داخل BRAM نوشته می شوند. در قسمت PL نیز برنامه در صورت فعال شدن پین start شروع به خواندن BRAM میکند، آرایه ها در قسمت PL خوانده می شوند و عملیات

convolution بر روی آرایه اصلی اجرا می‌شود. خروجی به BRAM نوشته می‌شود و پین done را فعال می‌کند و سپس PS آن را می‌خواند و خروجی را باز می‌گرداند.

برای پیاده‌سازی این سیستم، ابتدا یک پروژه جدید در Vivado ایجاد می‌کنیم. برد EBAZ4205 Development Board را به پروژه خود اضافه می‌کنیم. سپس Block Design را باز می‌کنیم و سه IP با نام‌های AXI BRAM Controller، ZYNQ 7010 و Block Memory Generator را اضافه می‌کنیم. حال یک فایل source جدید باز می‌کنیم و کد VHDL مربوط به پیاده‌سازی Convolution را در آن قرار می‌دهیم. سپس با زدن دکمه Auto Generate مشاهده می‌کنیم که IP ها به هم متصل می‌شوند.



حال به کد vhdل می پردازیم :

```

1  library IEEE;
2  use IEEE.std_logic_1164.all;
3  use IEEE.STD_LOGIC_UNSIGNED.ALL;
4  use IEEE.numeric_std.all;
5
6  entity bram_multiplier is
7  Port (
8      clk : in STD_LOGIC;
9      Addr_bram : out std_logic_vector(15 downto 0);
10     EN_bram : out STD_LOGIC := '0';
11     dout_bram : in std_logic_vector(31 downto 0);
12     RST_bram : out STD_LOGIC := '1';
13     din_bram : out std_logic_vector(31 downto 0);
14     WEB_bram : out std_logic_vector(3 downto 0);
15     Done : out STD_LOGIC := '0';
16     Start : in STD_LOGIC
17 );
18 end bram_multiplier;
19
20 architecture Behavioral of bram_multiplier is
21     constant MATRIX_SIZE : integer := 28;
22     constant KERNEL_SIZE : integer := 3;
23
24     signal adr_s : std_logic_vector(15 downto 0) := (others => '0');
25     signal RData1 : std_logic_vector(31 downto 0) := (others => '0');
26     signal RData2 : std_logic_vector(31 downto 0) := (others => '0');
27     signal WEB_S : std_logic_vector(3 downto 0) := (others => '0');
28     signal product : std_logic_vector(31 downto 0) := (others => '0');
29     signal Is_Run : STD_LOGIC := '0';
30     signal cnt : integer := 0;
31
32     type matrix_28x28 is array (0 to MATRIX_SIZE-1, 0 to MATRIX_SIZE-1) of std_logic_vector(31 downto 0);
33     type matrix_3x3 is array (0 to KERNEL_SIZE-1, 0 to KERNEL_SIZE-1) of std_logic_vector(31 downto 0);
34
35     signal input_matrix : matrix_28x28 := (others => (others => (others => '0')));
36     signal kernel_matrix : matrix_3x3 := (others => (others => (others => '0')));
37     signal output_matrix : matrix_28x28 := (others => (others => (others => '0')));

```

```

32 type matrix_28x28 is array (0 to MATRIX_SIZE-1, 0 to MATRIX_SIZE-1) of std_logic_vector(31 downto 0);
33 type matrix_3x3 is array (0 to KERNEL_SIZE-1, 0 to KERNEL_SIZE-1) of std_logic_vector(31 downto 0);
34
35 signal input_matrix : matrix_28x28 := (others => (others => (others => '0')));
36 signal kernel_matrix : matrix_3x3 := (others => (others => (others => '0')));
37 signal result_matrix : matrix_28x28 := (others => (others => (others => '0')));
38
39 signal current_row : integer := 0;
40 signal current_col : integer := 0;
41 signal kernel_row : integer := 0;
42 signal kernel_col : integer := 0;
43 begin
44
45     Addr_bram <= adr_s;
46     RST_bram <= '0';
47     EN_bram <= '1';
48     WEB_bram <= WEB_S;
49
50     process(clk)
51     begin
52         if rising_edge(clk) then
53             if (Start='1') and (Is_Run = '0') then
54                 Is_Run <= '1';
55                 Done <= '0';
56                 cnt <= 0;
57                 WEB_S <= "0000";
58                 adr_s <= (others => '0');
59                 current_row <= 0;
60                 current_col <= 0;
61                 kernel_row <= 0;
62                 kernel_col <= 0;
63             elsif (Is_Run = '1') then
64                 if cnt < (MATRIX_SIZE * MATRIX_SIZE) then
65                     if cnt < (MATRIX_SIZE * MATRIX_SIZE) then
66                         input_matrix(current_row, current_col) <= dout_bram;
67                         adr_s <= std_logic_vector(to_unsigned(cnt + 1, 16));
68                         cnt <= cnt + 1;
69                         current_col <= (current_col + 1) mod MATRIX_SIZE;
70                         if current_col = 0 then
71                             current_row <= current_row + 1;
72                         end if;
73                     elsif cnt < (MATRIX_SIZE * MATRIX_SIZE + KERNEL_SIZE * KERNEL_SIZE) then
74                         kernel_matrix(kernel_row, kernel_col) <= dout_bram;
75                         adr_s <= std_logic_vector(to_unsigned(cnt + 1, 16));
76                         cnt <= cnt + 1;
77                         kernel_col <= (kernel_col + 1) mod KERNEL_SIZE;
78                         if kernel_col = 0 then
79                             kernel_row <= kernel_row + 1;
80                         end if;
81                     else
82                         -- Perform convolution
83                         for i in 0 to MATRIX_SIZE-KERNEL_SIZE loop
84                             for j in 0 to MATRIX_SIZE-KERNEL_SIZE loop
85                                 product <= (others => '0');
86                                 for ki in 0 to KERNEL_SIZE-1 loop
87                                     for kj in 0 to KERNEL_SIZE-1 loop
88                                         product <= std_logic_vector(
89                                             resize(
90                                                 unsigned(product) + unsigned(input_matrix(i+ki, j+kj)) * unsigned(kernel_matrix(ki, kj)), 32
91                                             );
92                                     end loop;
93                                 end loop;
94                                 result_matrix(i, j) <= product;
95                             end loop;
96                         end loop;
97                     ..
98                 end loop;
99
100                 -- Write result back to BRAM
101                 if cnt < (MATRIX_SIZE * MATRIX_SIZE + KERNEL_SIZE * KERNEL_SIZE + MATRIX_SIZE * MATRIX_SIZE) then
102                     adr_s <= std_logic_vector(to_unsigned(cnt + 1, 16));
103                     din_bram <= result_matrix(current_row, current_col);
104                     WEB_S <= "1111";
105                     cnt <= cnt + 1;
106                     current_col <= (current_col + 1) mod MATRIX_SIZE;
107                     if current_col = 0 then
108                         current_row <= current_row + 1;
109                     end if;
110                 else
111                     Done <= '1';
112                     if Start='0' then
113                         Is_Run <= '0';
114                     end if;
115                 end if;
116             end if;
117         end process;
118     end Behavioral;
119
120

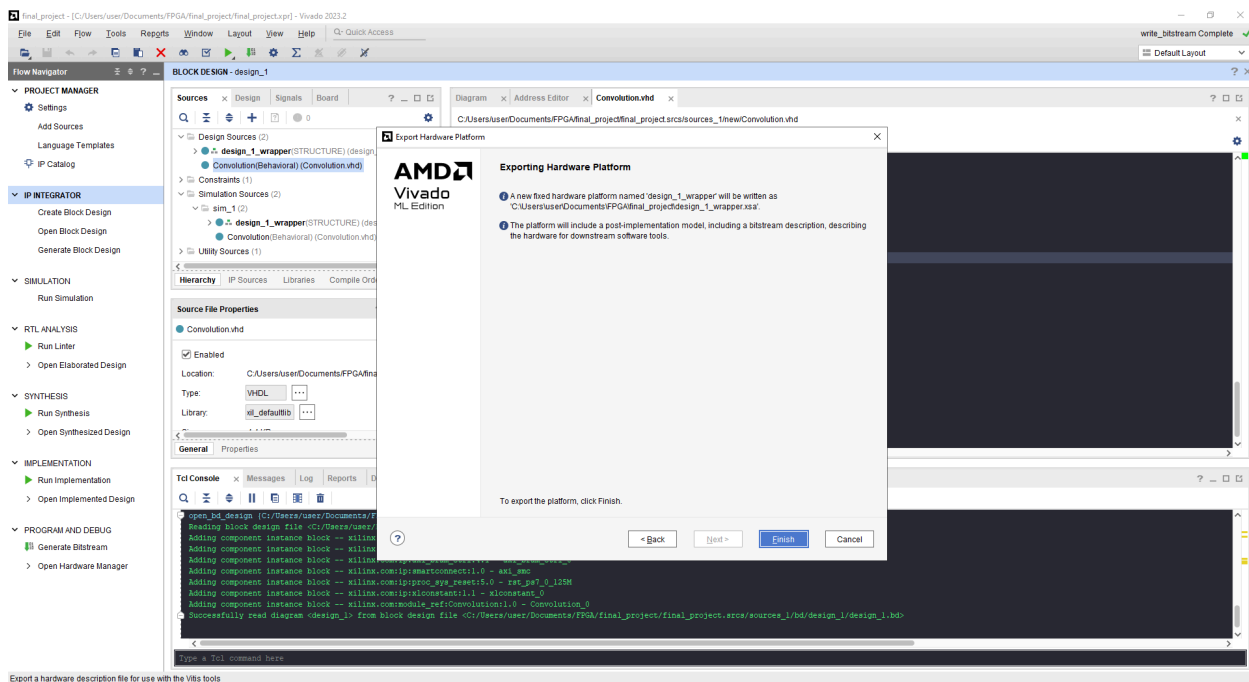
```

این کد VHDL یک ماژول به نام `bram_multiplier` را پیاده‌سازی می‌کند که وظیفه انجام ضرب ماتریسی (کانولوشن) را با استفاده از حافظه بلوکی (BRAM) بر عهده دارد. این ماژول شامل پورت‌هایی برای سیگنال ساعت (`clk`)، آدرس BRAM، سیگنال فعال‌سازی (`EN_bram`)، داده‌های خروجی از `BRAM (dout_bram)`، سیگنال بازنشانی (`RST_bram`)، داده‌های ورودی به `BRAM (din_bram)`، سیگنال‌های نوشتن به `BRAM (WEB_bram)`، سیگنال اتمام عملیات (`Done`) و سیگنال شروع عملیات (`Start`) است.

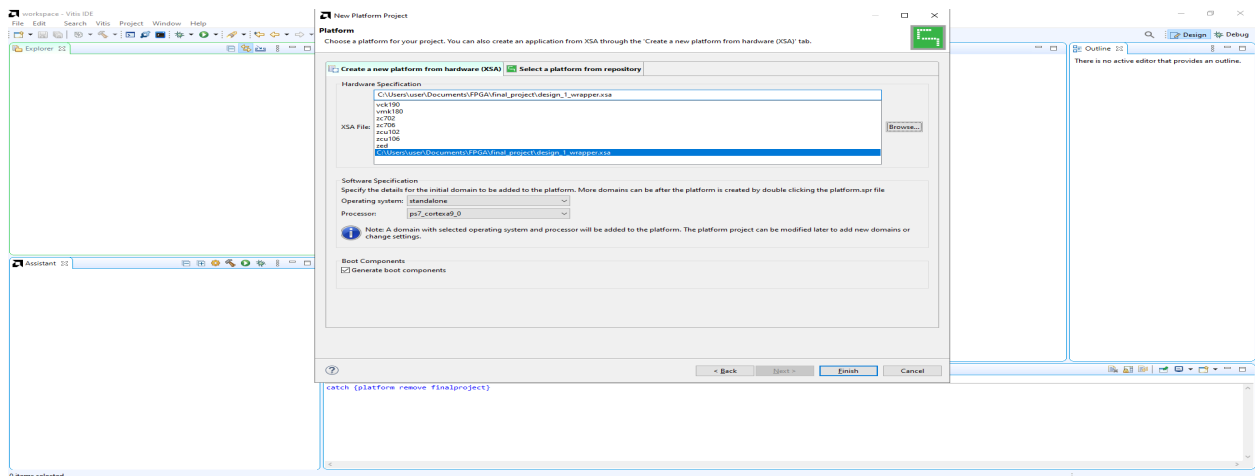
در معماری Behavioral این ماژول، ثابت‌هایی برای اندازه ماتریس ورودی (28×28) و اندازه کرنل (3×3) تعریف شده‌اند. سیگنال‌های داخلی برای آدرس‌دهی، داده‌های خوانده شده، داده‌های نوشته شده، سیگنال‌های نوشتن، محصول نهایی، وضعیت اجرای عملیات و شمارنده تعریف شده‌اند. همچنین، انواع داده‌های ماتریسی برای ذخیره ماتریس ورودی، کرنل و نتیجه عملیات نیز تعریف شده‌اند.

فرآیند اصلی با فعال شدن سیگنال `Start` آغاز می‌شود. در این مرحله، داده‌های ماتریس ورودی و کرنل از BRAM خوانده می‌شوند. سپس، عملیات کانولوشن بر روی ماتریس ورودی با استفاده از کرنل انجام می‌شود و نتیجه در ماتریس ذخیره می‌گردد. این عملیات شامل جمع حاصلضرب عناصر ماتریس ورودی با کرنل است. پس از اتمام کانولوشن، نتایج به BRAM نوشته می‌شوند و سیگنال `Done` فعال می‌شود تا نشان دهد که عملیات به پایان رسیده است.

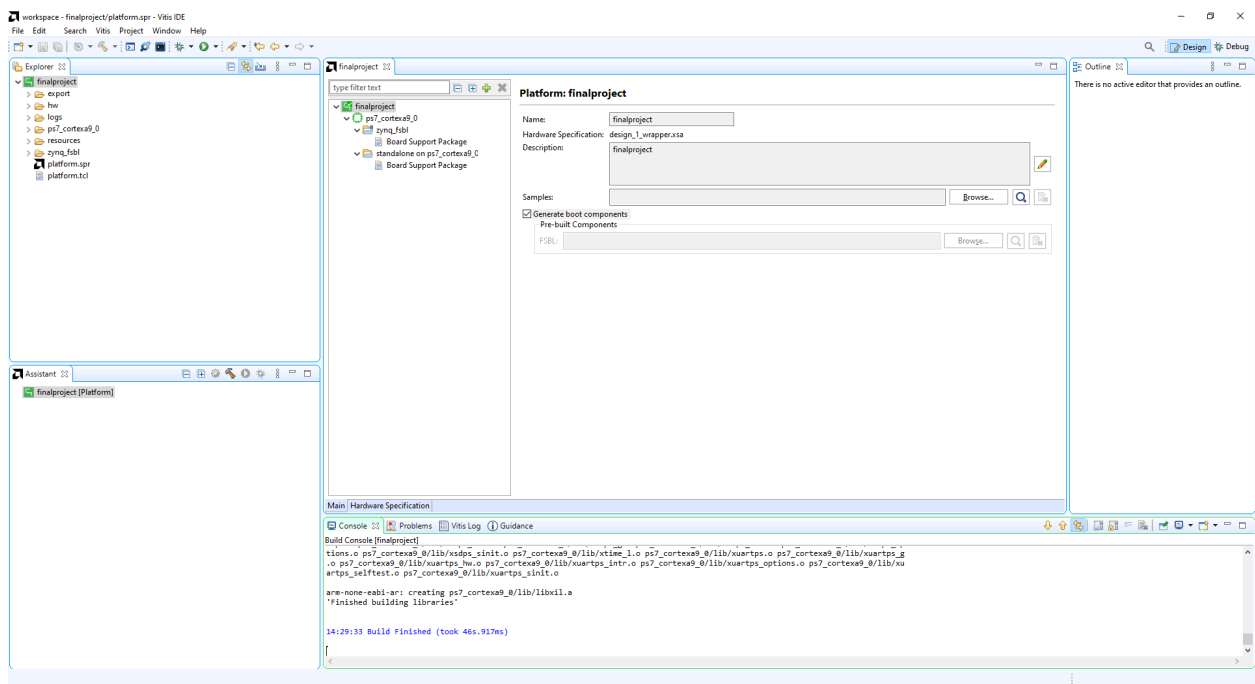
پروژه را سنتز میکنیم و سخت افزار را export میکنیم.



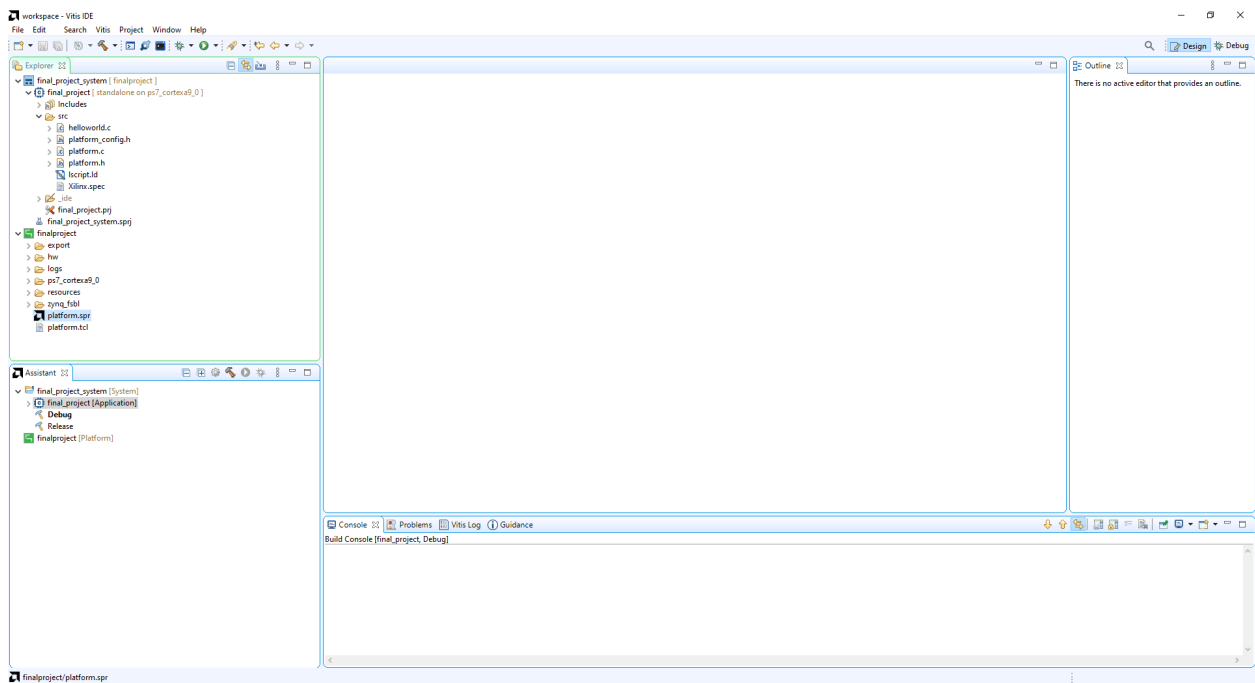
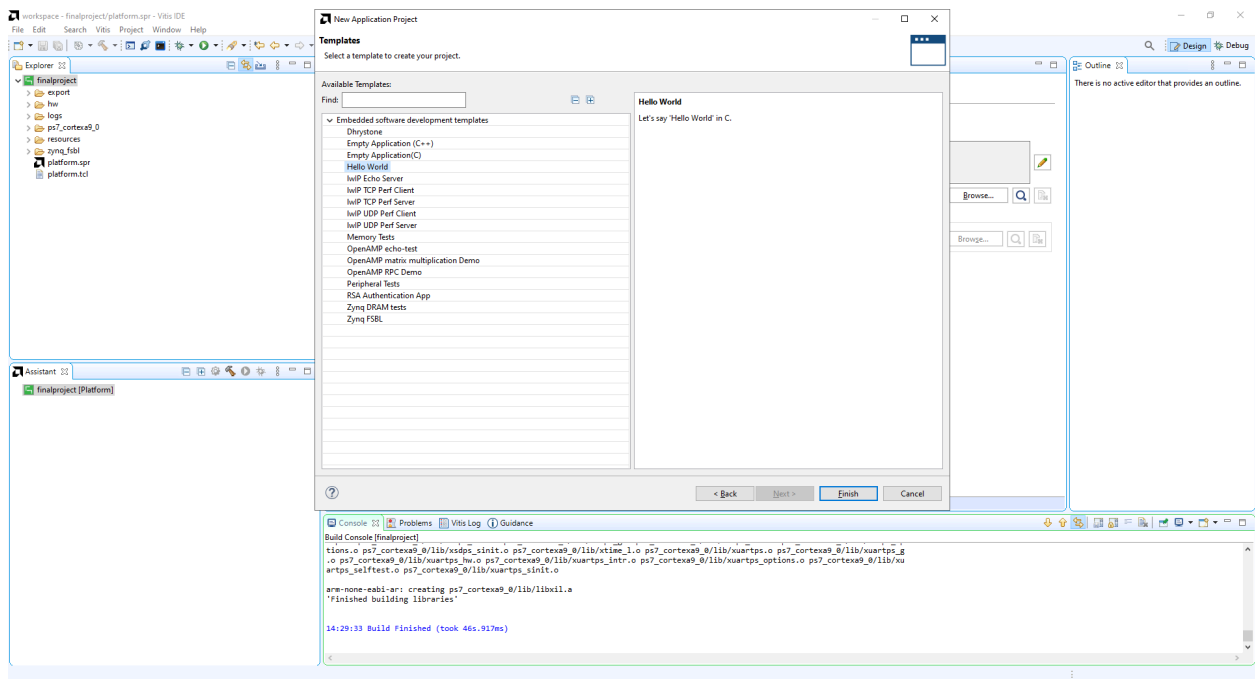
در مرحله بعد `vitis classic` را باز میکنیم. در بخش `new` یک `platform project` را ایجاد میکنیم و فایل سخت افزار `export` شده را به نرم افزار میدهیم.



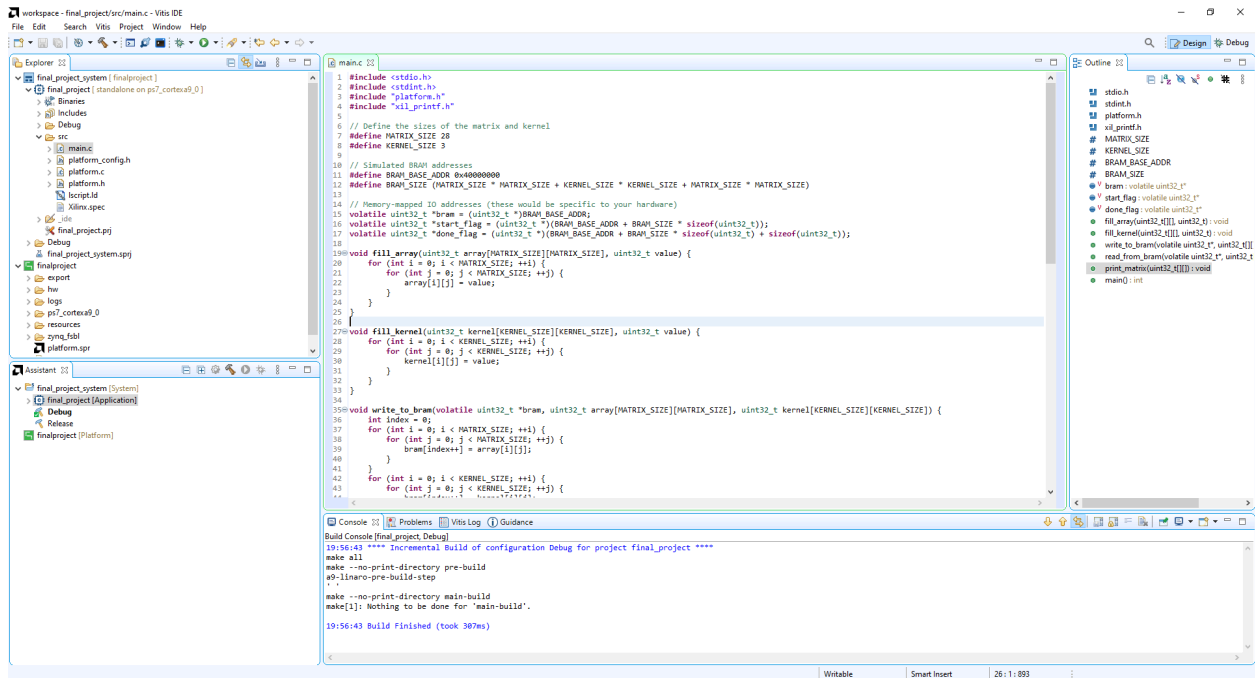
پروژه را build میکنیم.



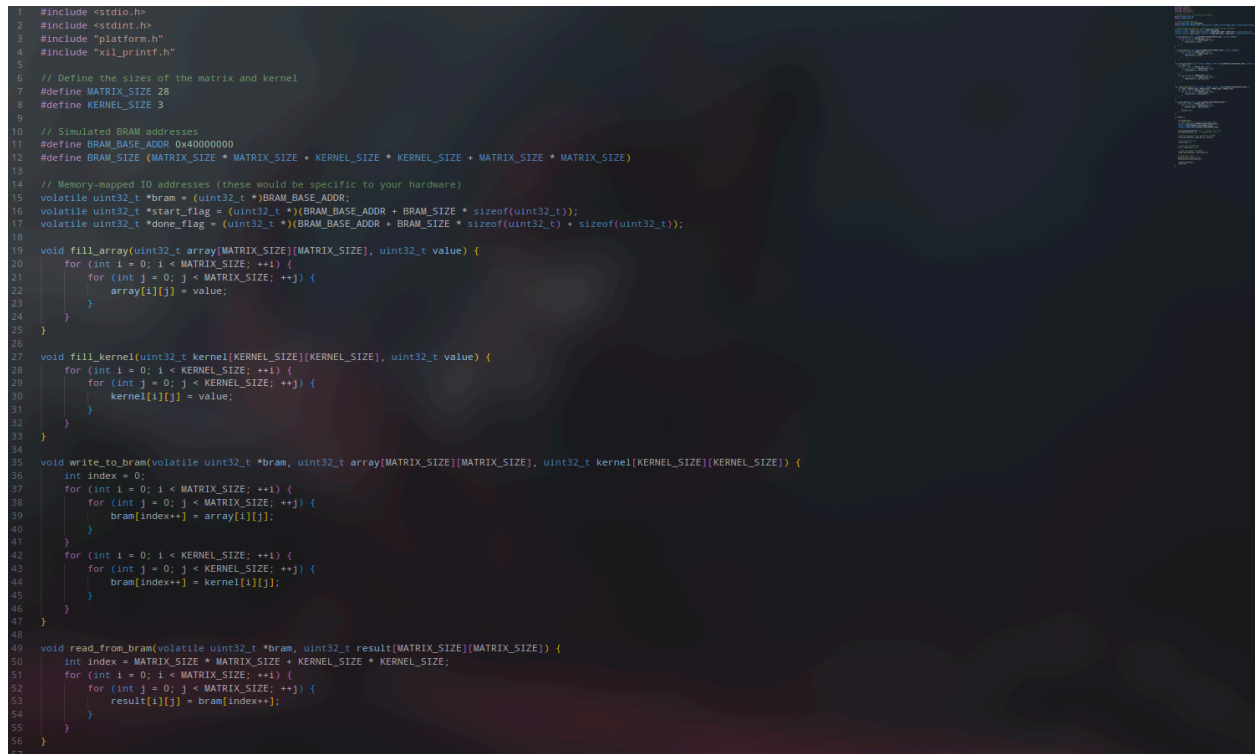
حال یک application project میسازیم. سخت افزار را انتخاب میکنیم و به عنوان template از hello world استفاده میکنیم.



درون پوشه src فایل helloworld.c را به main.c تغییر می دهیم و درون فایل main.c کد مربوط به ps خود را اضافه میکنیم. حال پروژه نرم افزاری را build میکنیم.



حال به کد C می پردازیم:



```

57
58 void print_matrix(uint32_t matrix[MATRIX_SIZE][MATRIX_SIZE]) {
59     for (int i = 0; i < MATRIX_SIZE; ++i) {
60         for (int j = 0; j < MATRIX_SIZE; ++j) {
61             printf("%10u ", matrix[i][j]);
62         }
63         printf("\n");
64     }
65 }
66
67 int main() {
68     init_platform();
69     // Create and fill the input matrix and kernel
70     uint32_t input_matrix[MATRIX_SIZE][MATRIX_SIZE];
71     uint32_t kernel[KERNEL_SIZE][KERNEL_SIZE];
72     uint32_t result_matrix[MATRIX_SIZE][MATRIX_SIZE];
73
74     fill_array(input_matrix, 1); // Example fill value
75     fill_kernel(kernel, 1); // Example fill value
76
77     // Write the input matrix and kernel to BRAM
78     write_to_bram(bram, input_matrix, kernel);
79
80     // Set the start flag
81     *start_flag = 1;
82
83     // Wait for the done flag
84     while (*done_flag == 0);
85
86     // Read the result from BRAM
87     read_from_bram(bram, result_matrix);
88
89     // Print the result
90     printf("Result Matrix:\n");
91     print_matrix(result_matrix);
92
93     cleanup_platform();
94     return 0;
95 }
96
97

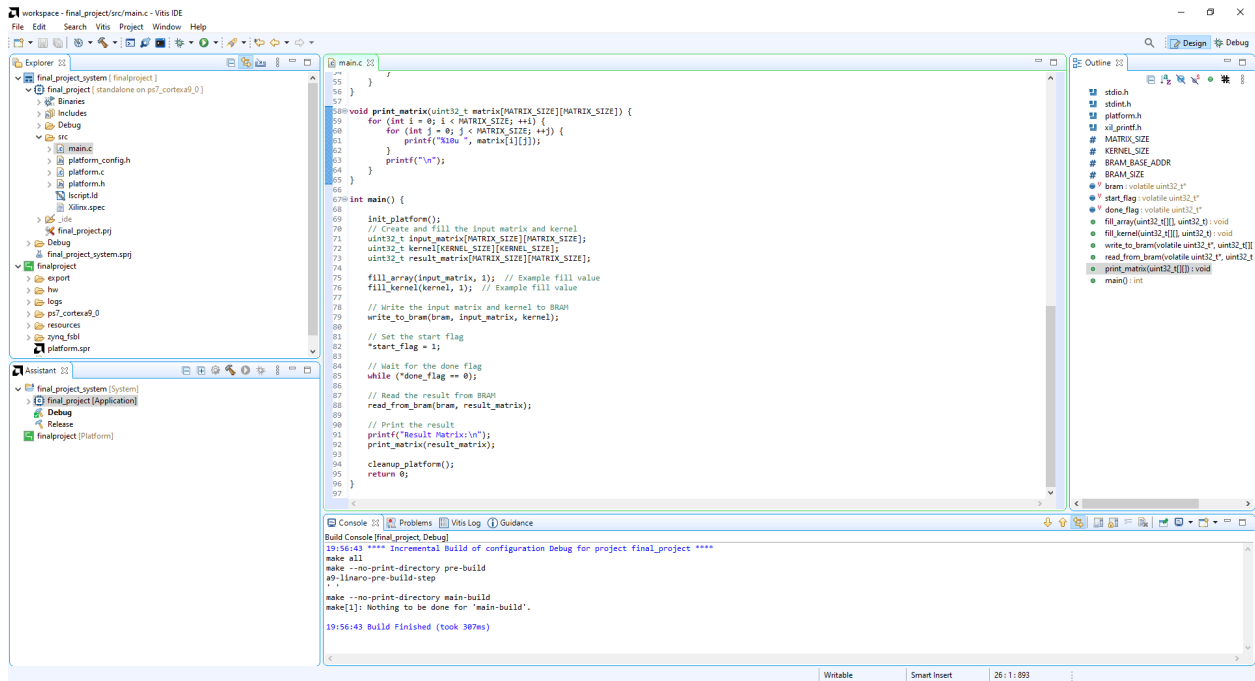
```

این کد یک مثال ساده از استفاده از حافظه‌های مموری مپ شده در برنامه‌نویسی سخت‌افزاری است. در این برنامه، از حافظه Block RAM (BRAM) برای ارتباط بین نرم‌افزار و سخت‌افزار استفاده می‌شود. ابتدا، اندازه و آدرس پایه حافظه BRAM تعریف می‌شوند تا بتوان اطلاعات را به آن نوشته و از آن خواند. سپس، توابع کمکی برای پر کردن ماتریس و کرنل با مقادیر مشخص و نوشتن این اطلاعات به حافظه BRAM تعریف می‌شوند.

در `main`، ابتدا ماتریس و کرنل با مقادیر مورد نظر پر می‌شوند. سپس، این داده‌ها به حافظه BRAM با استفاده از تابع `write_to_bram` نوشته می‌شوند. پس از آن، `flag` شروع (`start_flag`) به 1 تنظیم می‌شود تا سخت‌افزار بفهمد که پردازش باید شروع شود. برنامه در حالی که `flag` انجام (`done_flag`) به 0 است، منتظر اتمام پردازش در سخت‌افزار می‌ماند. وقتی که این `flag` به 1 تغییر کند، نتیجه از حافظه BRAM خوانده شده و در ماتریس `result_matrix` قرار داده می‌شود. در نهایت، ماتریس نتیجه چاپ می‌شود.

این نمونه نشان می‌دهد که چگونه از حافظه‌های مموری مپ شده برای انتقال داده‌ها بین نرم‌افزار و سخت‌افزار استفاده کرد تا فرآیندهای پردازشی مانند کانولوشن به صورت کارآمد اجرا شوند.

به عنوان جمع بندی ما یک شتابدهنده ضرب convolution برای یک ماتریس و یک kernel را به کمک co-design طراحی و پیاده سازی کردیم.



با تشکر