DWConv & Standard Conv



PWConv & Matrix-Matrix Multiplication



PE

**Contributions:**

1- Reconfigurable PE (**DSP size**, Systolic array, **Column based**)

**6 times more 8x8** multiplier comparing to a DSP Block (two 8x8),

**RS Data flow**, High frequency, flexible data movement. Great for **SConv, DWConv, PWConv, Matrix-Matrix Multiplication**
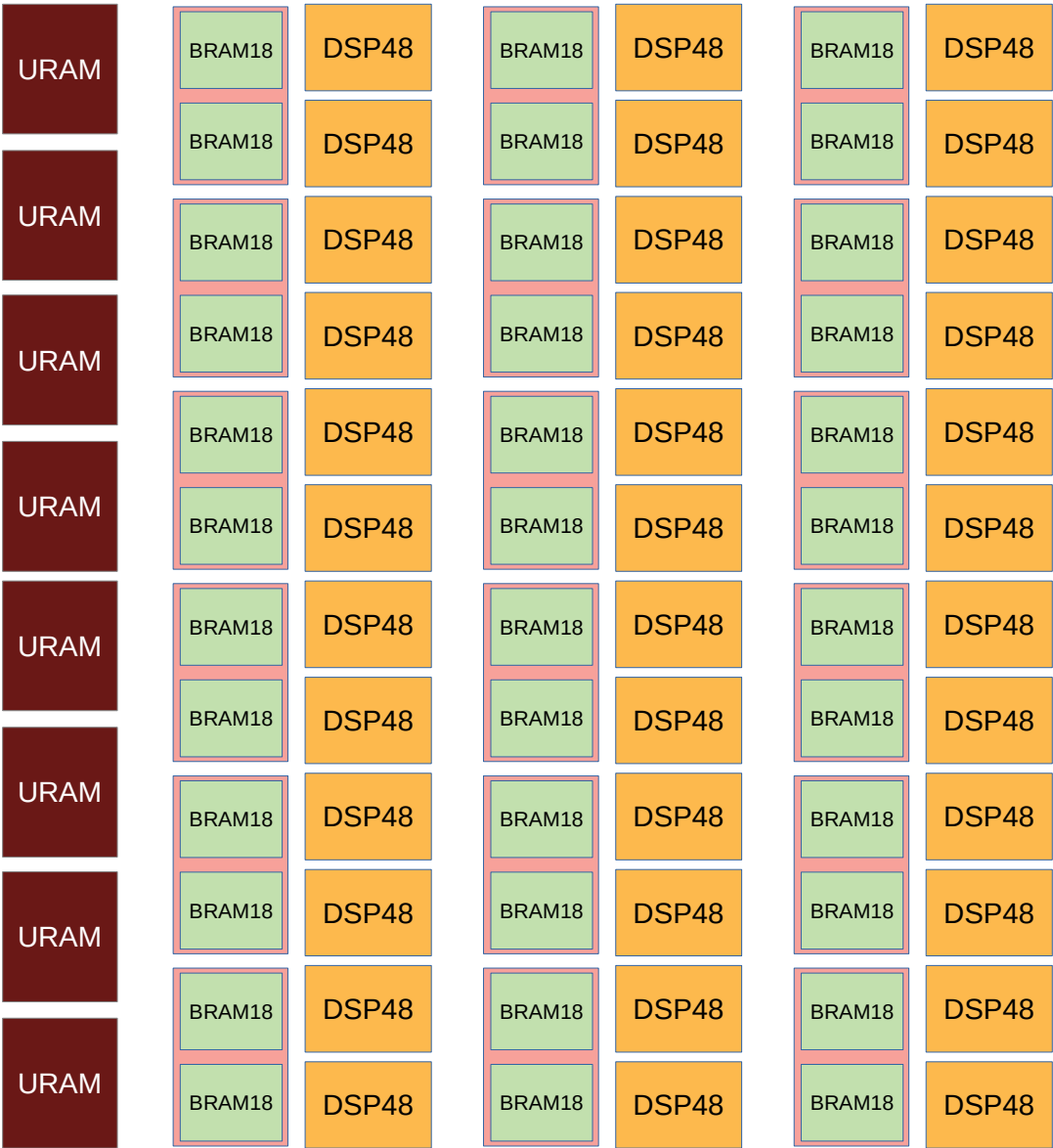
DPS – BRAM ratio 1/1 (same as Ultrascale+ arch), Low number of intermediary outputs in practice

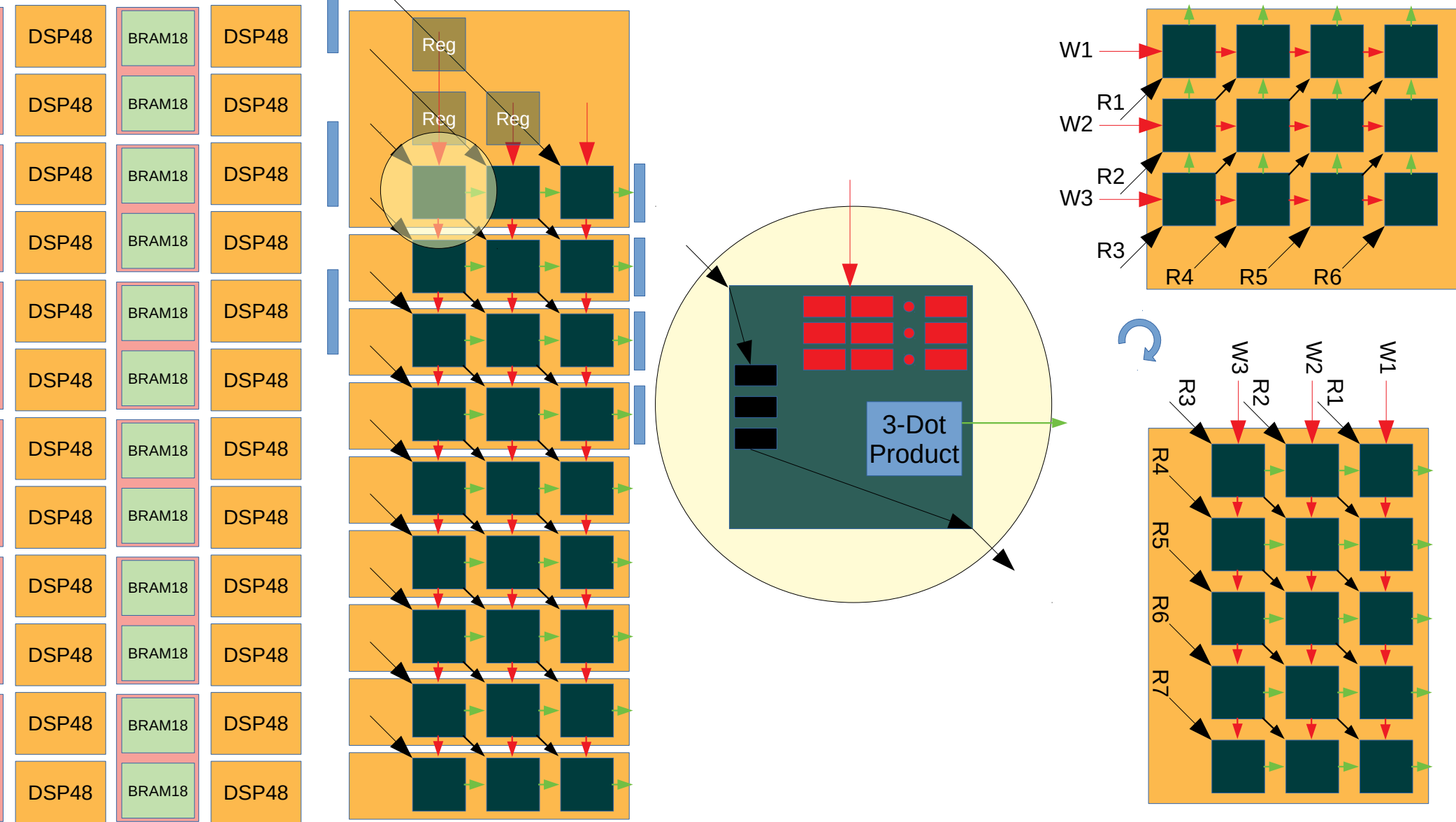Parameterized (for any budget limitation) – can integrate multi precision idea

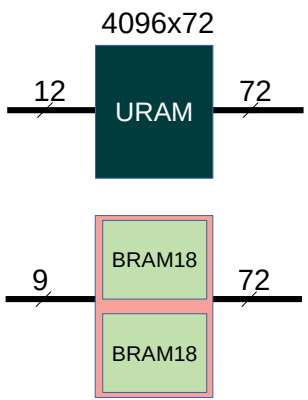2- Compare with cascade paper (Prof. Nachiket)

3- new suggestion to use each 18KBRAM as 36bit streamer using external controler circuit (delivering 662MHz) (in cascade paper: 18bit)
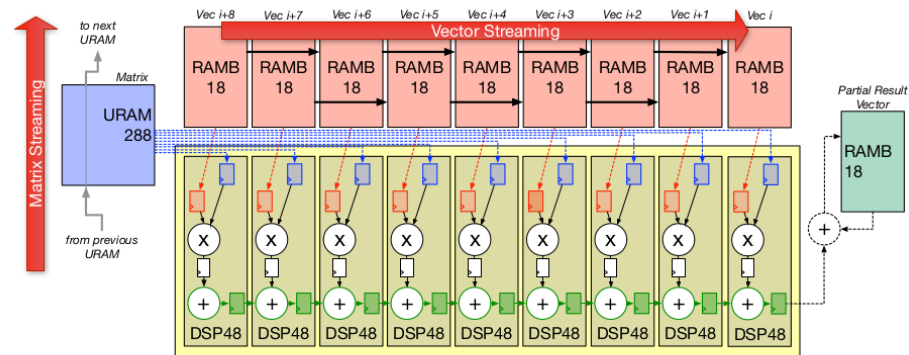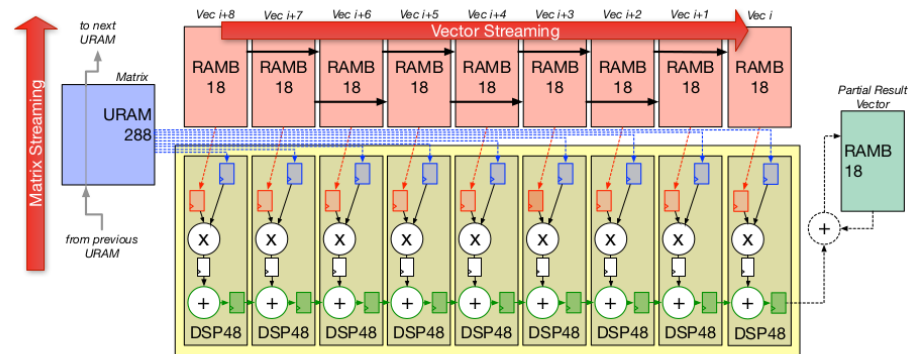
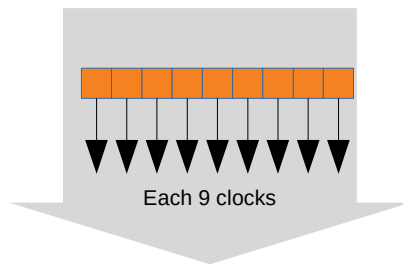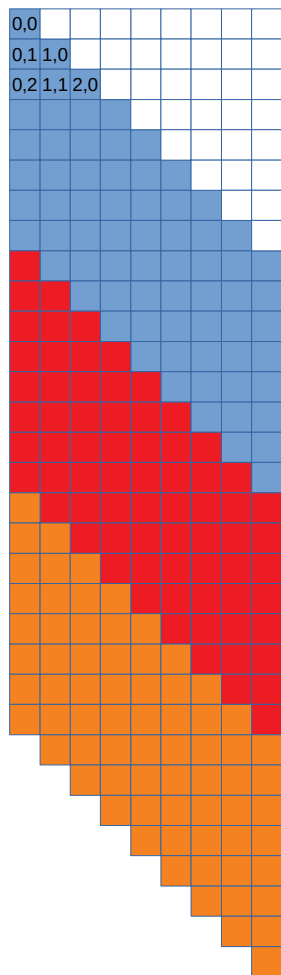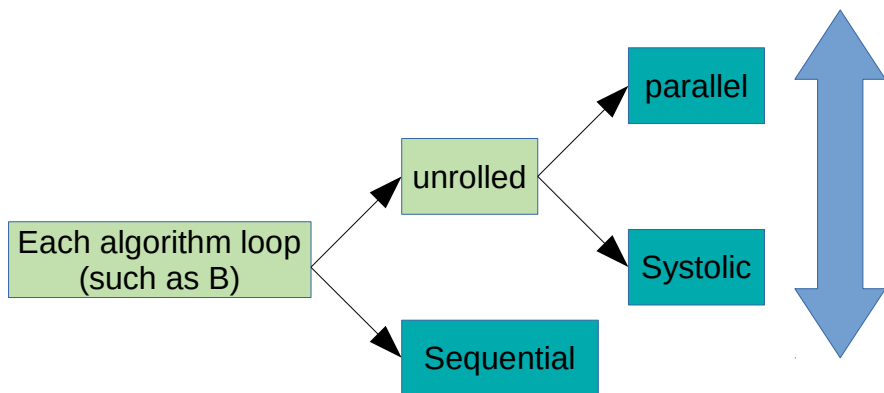UltraScale+ architecture distribution:

# My amassing MLBlocks world

DSP48 BRAM18 DSP48

Reg Reg Reg

3-Dot Product

W1 R1 W2 R2 W3 R3 R4 R5 R6

R3 W3 R2 W2 R1 W1 R4 R5 R6 R7

4096x72

12 URAM 72

9 BRAM18 72

BRAM18

Each 9 clocks

Each algorithm loop
(such as B)

unrolled

parallel

Systolic

Sequential

Dis Parallel:

1- more fan in and outs (since we are talking about small Pes it is fine)

Dis Systolic:

1- tougher scheduling, rythmic scheduling
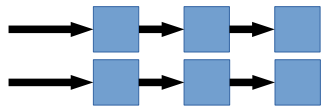
2- prevent circuit fusions (less optimization)

$B = B_{Seq} \times B_{par} \times B_{Sys}$

# of Physical MAC:    $\times B_{par} \times B_{Sys}$
# of Input:                $\times B_{par}$          (without internal serial to parallel)
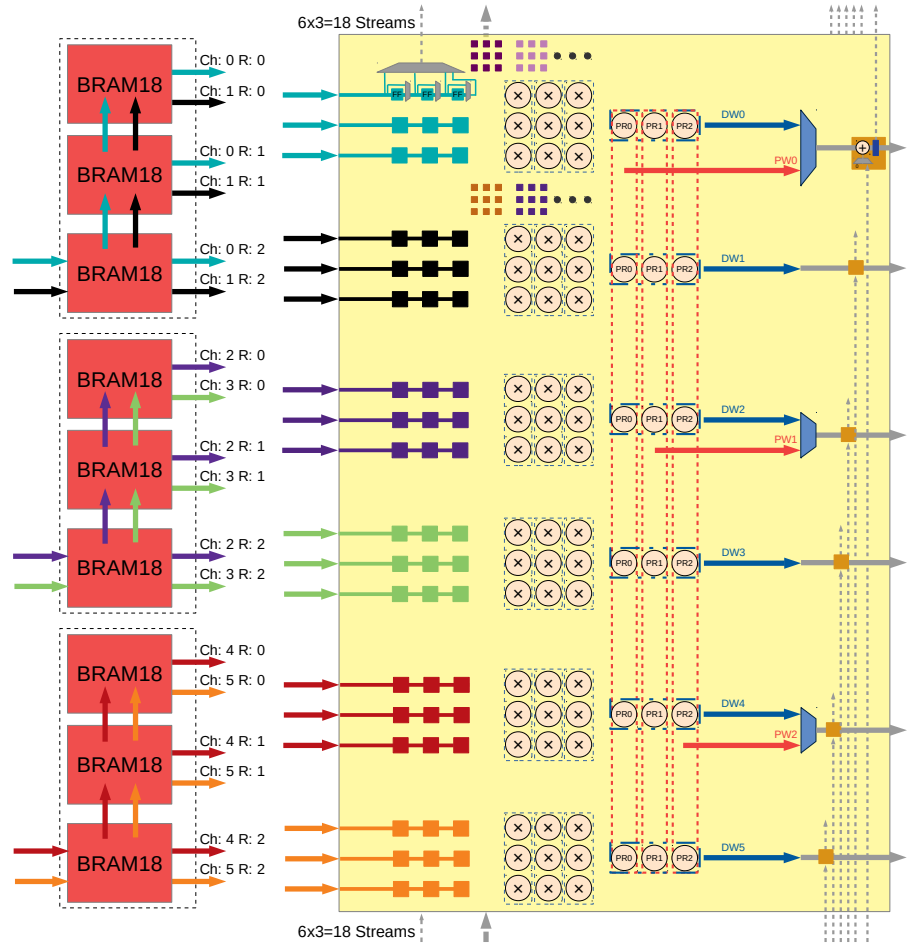# of Output:              $\times B_{par}$

Params = {right side indexes}

for $\quad param^{i}_{sch\_0}$: $\quad 0 \rightarrow sch^{i}\_0$

$\quad$ for $\quad param^{i}_{sch\_1}$: $\quad 0 \rightarrow sch^{i}\_1$

$\quad\quad$ for $\quad param^{i}_{sch\_2}$: $\quad 0 \rightarrow sch^{i}\_2$

$\quad\quad\quad$ for $\quad param^{i}_{seq}$: $\quad 0 \rightarrow comp\_seq^{i}$

$\quad\quad\quad\quad$ for $\quad param^{i}_{un}$: $\quad 0 \rightarrow comp\_un^{i}$

DSP-size (BRAM-size)

BRAM18
Ch: 0 R: 0
Ch: 1 R: 0

BRAM18
Ch: 0 R: 0
Ch: 1 R: 0
Ch: 2 R: 0
Ch: 3 R: 0

8bit
8bit
8bit
8bit
48bit_6x8

RW0
RW1
RW2
RW3

8bit
8bit
8bit
8bit
48bit_6x8

RW0
RW1
RW2
RW3

DSP-size (BRAM-size)

BRAM18

Ch: 0 R: 0
Ch: 1 R: 0

BRAM18

Ch: 0 R: 0
Ch: 1 R: 0
Ch: 2 R: 0
Ch: 3 R: 0

8bit
8bit
8bit
8bit

48bit_6x8

RW0
RW1
RW2
RW3