



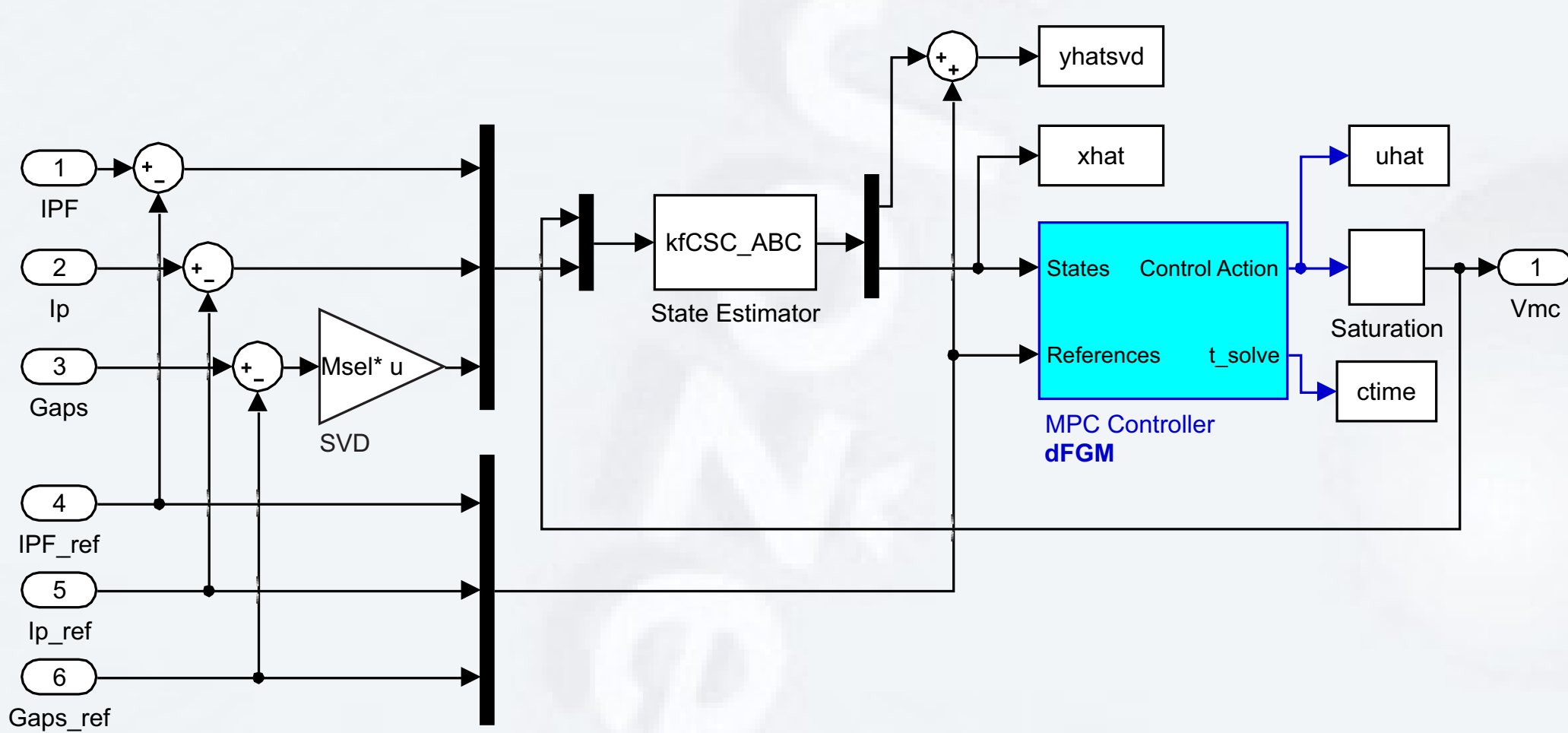
FPGA acceleration of Model Predictive Control for ITER Plasma current and shape control

Samo Gerksič^a, Boštjan Pregelj^a, Matija Perne^a

^a Jozef Stefan Institute, Jamova cesta 39, 1000 Ljubljana, Slovenia

Overview

- Plasma Current and Shape Controller (PCSC) for ITER
- Singular Value Decomposition (SVD) -based Model Predictive Control (MPC) using a dual Fast Gradient Method (dFGM) quadratic programming (QP) solver
- FPGA acceleration with a High-Level Synthesis (HLS) approach



SVD-based fast MPC PCSC block scheme

Model Predictive Control (MPC)

MPC is an advanced model-based process control technique, based on on-line optimization of predicted future courses of the process signals. Typically, a simplified control model in the discrete-time linear state space form is used, and the MPC controller is used in conjunction with a Kalman state estimator. It is closely related to LQ (linear quadratic) optimal control, and is convenient for control of multivariable processes such as PCSC. It is efficient in handling of constraints - in the case of ITER PCSC:

- input constraints of the power supply voltages V_{PF}
- output constraints of the superconductive coil currents I_{PF} .

Solving QP problems for MPC using dFGM

In MPC, a QP optimisation problem must be solved in each time step of the algorithm, which is difficult with large-scale multivariable systems with fast dynamics. Using complexity reduction techniques for MPC and C code optimisation, with the dFGM algorithm the required accuracy was achieved in 500 iterations with peak computation times 3 ms using a standard Intel x86 CPU using a single thread. This is considered sufficiently fast for ITER, but not for experimentation on smaller tokamaks with faster dynamics. The execution of the dFGM algorithm may be accelerated by parallelisation of matrix-vector numerical operations within iterations, but this is not well-suited for the standard CPU architecture with the thread scheduler timescale 10 μ s, thread synchronisation takes more time than computation.

FPGA implementation

FPGA implementation supports small-scale parallelism, but they have other issues regarding the implementation of complex algorithms: relatively low clock frequencies, limited resources for massive parallelism, specific programming. Aiming to avoid manual recoding as much as possible, the High-Level Synthesis (HLS) approach of Xilinx Vivado HLS was adopted:

- C code (simplifications required: static allocation of variables, direct implementation of functions)
- emulated single float precision (multiple cycles/operation)
- conversion directives (pragma commands) to the VHDL compiler for more optimized implementation: pipeline, loop unroll, loop merge, array partition
- hardware: Xilinx Zedboard, Zc706

The dFGM main iteration loop must be optimized for speed. Its computational load is dominated by two matrix-vector multiplications (33x99 * 99 and 99x33 * 33), while other vector operations represent less than 10% of the workload.

MPC QP, hard constraints

$$\min_z \frac{1}{2} z^T H z + c^T z$$

subject to $Cz \leq b$

MPC QP, soft state constraints

$$\min_{z,s} \frac{1}{2} z^T H z + c^T z + \frac{1}{2} s^T W s + w^T s$$

subject to $C_x z \leq b_x + s$
 $C_u z \leq b_u$
 $s \leq 0$

$$c = \begin{bmatrix} C_x \\ C_u \end{bmatrix}, \quad b = \begin{bmatrix} b_x \\ b_u \end{bmatrix}$$

Dual FGM algorithm

$$v^k = v^k + \beta^k (v^k - v^{k-1})$$

$$y^k = -H^{-1} (C^T v^k + c)$$

$$v^{k+1} = v^k + C y^k - \overline{\rho} \overline{\alpha}_{h,w,w} (v^k + C y^k)$$

$$\overline{\rho} \overline{\alpha}_{h,w,w} (t_i) = \begin{cases} t_i & \text{if } t_i \leq b_i \\ b_i & \text{if } t_i > b_i \text{ and } i \text{ hard} \\ \frac{t_i + W_{ii} b_i - w_i}{W_{ii} + 1} & \text{if } b_i + w_i \geq t_i > b_i \text{ and } i \text{ soft} \\ \frac{t_i + W_{ii} b_i - w_i}{W_{ii} + 1} & \text{if } b_i + w_i < t_i \text{ and } i \text{ soft} \end{cases}$$

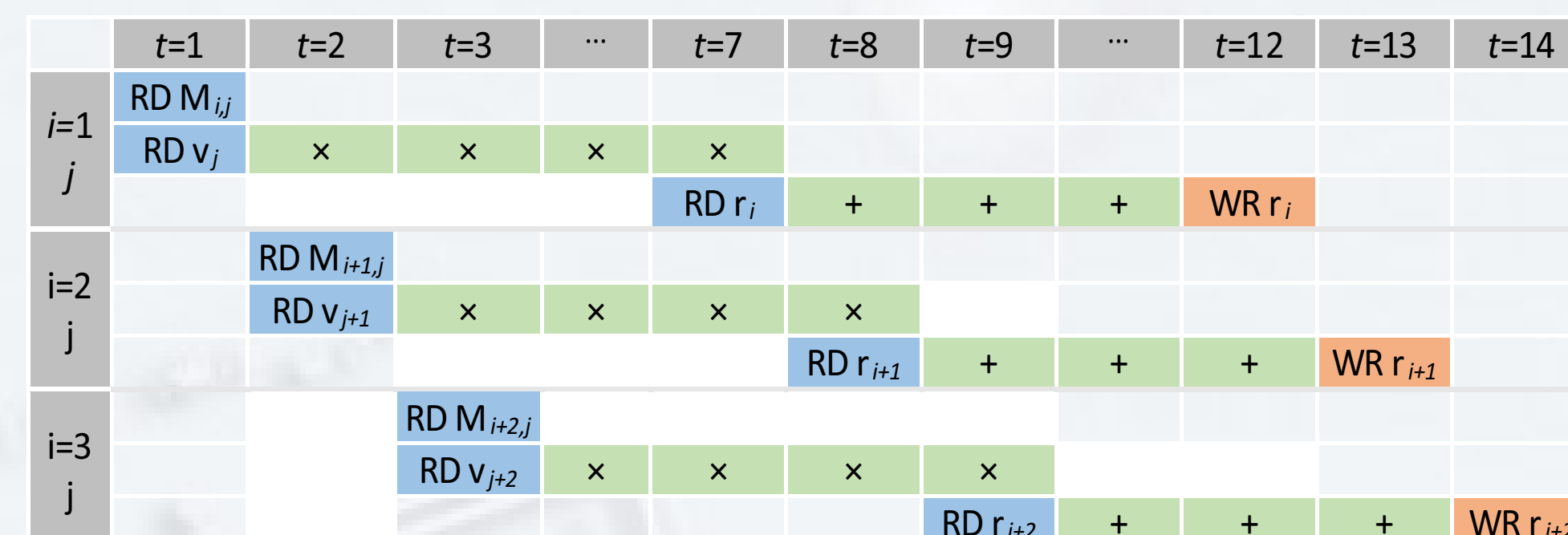
Matrix-vector multiplication implementation

$$r_{[m \times 1]} = M_{m \times n} \times v_{[n \times 1]}$$

Approach 1 (initial)

$$r_i = r_i + M_{ij} \times v_j$$

```
// mat_vec_mult_full (e>M, ws->v, ws->tmp_var_n);
mv_mult_M_x_v: for (j = 0; j < n_constraint; j++)
{
    vv_mult_M_x_v: for (i = 0; i < n_opt_var; i++)
    {
        tmp_var_n[i] += M[i][j] * v[j];
    }
}
```

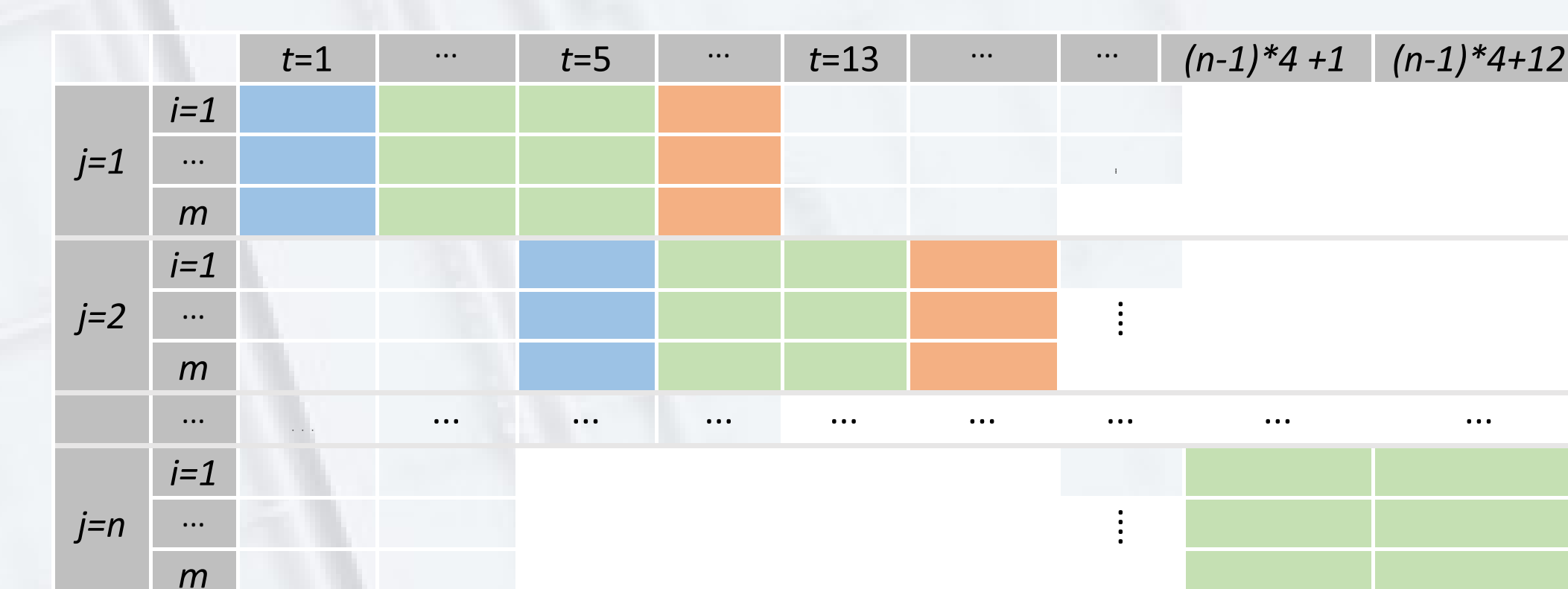
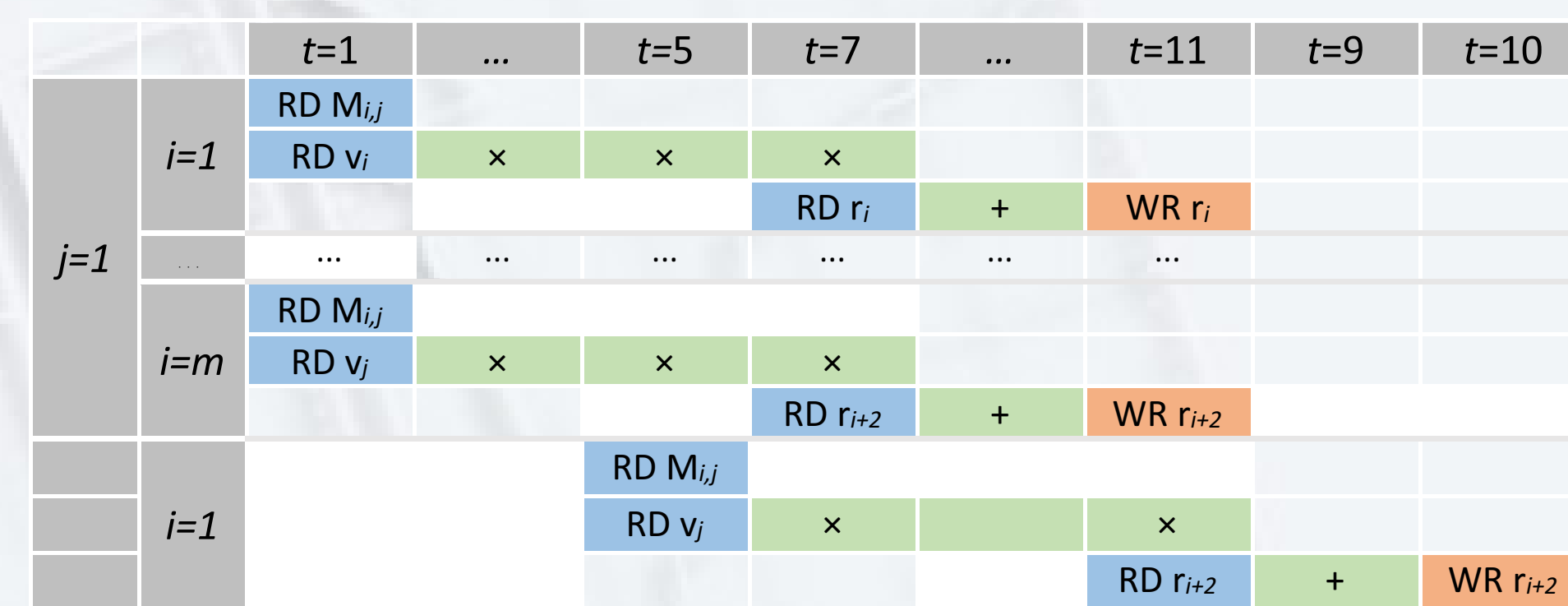


Time schedule for parts $M_{ij} \times v_j$ of $M \times v$ product.

Approach 2 (inner loop unroll)

- Unroll the inner i loop (dimension m)
- Partition M (dimension m)
- Partition r (dimension m)
- Pipeline the outer j loop

```
#pragma HLS ARRAY_PARTITION variable=M complete dim=1
#pragma HLS ARRAY_PARTITION variable=n complete
mv_mult_M_x_v: for (j = 0; j < n_constraint; j++)
{
    #pragma HLS PIPELINE
    vv_mult_M_x_v: for (i = 0; i < n_opt_var; i++)
    {
        #pragma HLS UNROLL
        tmp_var_n[i] += M[i][j] * v[j];
    }
}
```



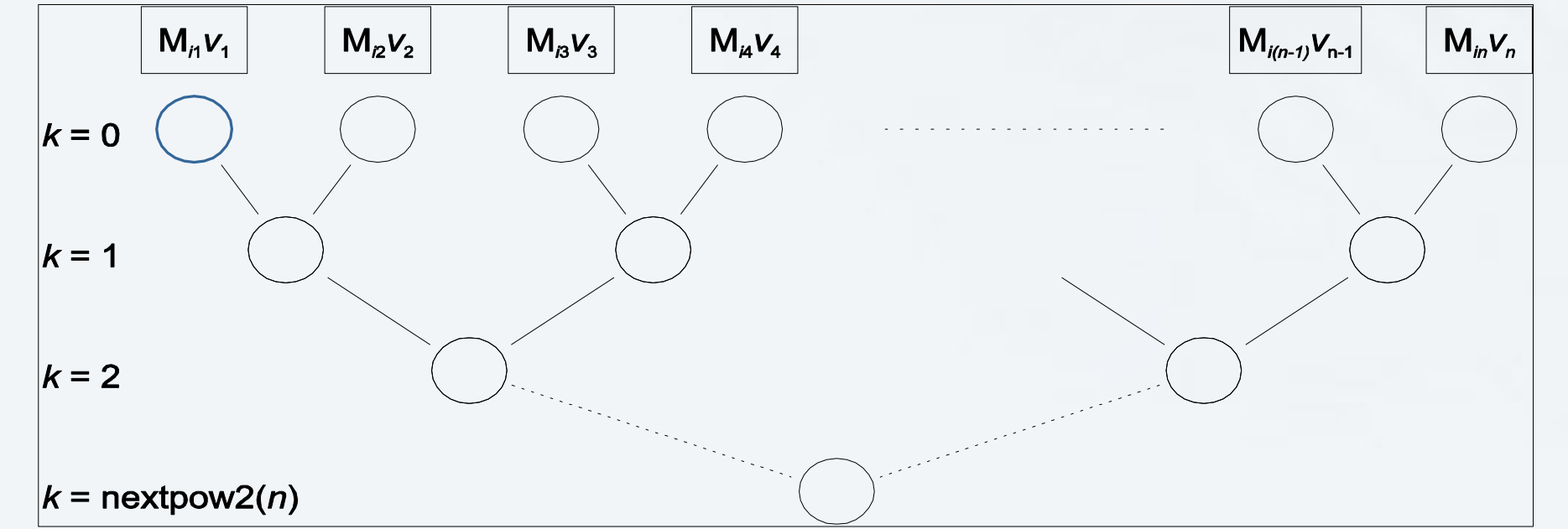
Timing schedule of $M \times v$ product, with parallelized inner loop.

Approach 3 (scalar product tree-sum)

Multiplication as series of scalar products:

$$r = M \times v = \begin{bmatrix} M_{1,1}v_1 + M_{1,2}v_2 + \dots + M_{1,n}v_n \\ M_{2,1}v_1 + M_{2,2}v_2 + \dots + M_{2,n}v_n \\ \vdots \\ M_{m,1}v_1 + M_{m,2}v_2 + \dots + M_{m,n}v_n \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^n (M_{1,j}v_j) \\ \sum_{j=1}^n (M_{2,j}v_j) \\ \vdots \\ \sum_{j=1}^n (M_{m,j}v_j) \end{bmatrix} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_m \end{bmatrix}$$

Binary-tree-wise parallel summation:



Parallelized calculation of vector scalar product (mat $M \times v$)

Approach 4 (tall matrix truncation)

Vertically divide M into p wide matrices M_1, \dots, M_p and process them in parallel

$$r = M \times v = \begin{bmatrix} M_1 \\ M_2 \\ \vdots \\ M_p \end{bmatrix} \times v = \begin{bmatrix} M_1 \times v \\ M_2 \times v \\ \vdots \\ M_p \times v \end{bmatrix} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_p \end{bmatrix}$$

```
mv_mult_C_x_x: for (i = 0; i < 33; i++) // (i = 0; i < n_constraint; i++)
{
    #pragma HLS PIPELINE
    vv_mult_C_x_x: for (j = 0; j < n_opt_var; j++)
    {
        #pragma HLS UNROLL //factor=33
        //tmp_var_p4[i] += C[i][j] * x[j];
        tmp11[j] = C[i][j] * x[j];
        tmp21[j] = C[i+33][j] * x[j];
        tmp31[j] = C[i+66][j] * x[j];
    }
    reset_tmp_rest_m2: for (j=n_opt_var;j<128;j++)
    {
        #pragma HLS UNROLL
        tmp11[j] = 0.0;
        tmp21[j] = 0.0;
        tmp31[j] = 0.0;
    }
    sum2m2: for (j=0;j<32;j++)
    {
        #pragma HLS UNROLL
        tmp12[j] = tmp11[j]+tmp11[j+32];
        tmp22[j] = tmp21[j]+tmp21[j+32];
        tmp32[j] = tmp31[j]+tmp31[j+32];
    }
    sum3m2: for (j=0;j<16;j++)
    {
        #pragma HLS UNROLL
        tmp11[j] = tmp12[j]+tmp12[j+16];
        tmp21[j] = tmp22[j]+tmp22[j+16];
        tmp31[j] = tmp32[j]+tmp32[j+16];
    }
    sum4m2: for (j=0;j<8;j++)
    {
        #pragma HLS UNROLL
        tmp12[j] = tmp11[j]+tmp11[j+8];
        tmp22[j] = tmp21[j]+tmp21[j+8];
        tmp32[j] = tmp31[j]+tmp31[j+8];
    }
    sum5m2: for (j=0;j<4;j++)
    {
        #pragma HLS UNROLL
        tmp11[j] = tmp12[j]+tmp12[j+4];
        tmp21[j] = tmp22[j]+tmp22[j+4];
        tmp31[j] = tmp32[j]+tmp32[j+4];
    }
    tmp12[0] = tmp11[0]+tmp11[2];
    tmp12[1] = tmp11[1]+tmp11[3];
    tmp_var_p4[i] = tmp12[0] + tmp12[1];
    tmp22[0] = tmp21[0]+tmp21[2];
    tmp22[1] = tmp21[1]+tmp21[3];
    tmp_var_p4[i+33] = tmp22[0] + tmp22[1];
    tmp32[0] = tmp31[0]+tmp31[2];
    tmp32[1] = tmp31[1]+tmp31[3];
    tmp_var_p4[i+66] = tmp32[0] + tmp32[1];
}
```

Matrix-vector multiplication operation count (in clock cycles)

Status	Iteration* latency (Lat)	Initiation* Interval (II)	Iteration* count (It)	Total latency (It*II + Lat)
Initial	m	$m+3$	n	$n(m+3) + n^2$
Inner unroll	12	4	n	$n^2 + 4 + 12$
Scalar product tree-sum	32	1	m	$m + 32$
Tall matrix truncation	32	1	m/p	$m/p + 32$

* The parameters are stated for the outer loop

** The additional n clocks are required to enter the inner loop n -times. This is handled using loop flatten

Conclusions

The automatic conversion does not yield useful results from the original C code, because the code structure prevents automated application of conversion optimization routines, so that the execution time is 45 ms, which is much longer than the CPU implementation, with low use of resources

After a series of manual modifications of the code with pragma directions, the computation time for Xilinx ZC706 is reduced to 1 ms (3x faster than CPU).

Further acceleration is presumably possible by using a FPGA with hardware flow-point multipliers, or by using fixed-point arithmetics (lower dynamic range)