

Pipelined Implementation of a Real Time Programmable Encoder for Low Density Parity Check Code on a Reconfigurable Instruction Cell Architecture

Zahid Khan, Tughrul Arslan
School of Engineering and Electronics,
The University of Edinburgh, Mayfield Road, Edinburgh, EH9 3JL
z.khan@ed.ac.uk

Abstract

This paper presents pipelined implementation of a real time programmable irregular Low Density Parity Check (LDPC) Encoder as specified in the IEEE P802.16E/D7 standard. The encoder is programmable for frame sizes from 576 to 2304 and for five different code rates. H matrix is efficiently generated and stored for a particular frame size and code rate. The encoder is implemented on Reconfigurable Instruction Cell Architecture which has recently emerged as an ultra low power, high performance, ANSI-C programmable embedded core. Different general and architecture specific optimization techniques are applied to enhance the throughput. With the architecture, a throughput from 10 to 19 Mbps has been achieved. The maximum throughput achieved with pipelining/ multi-core is 78 Mbps.

1. Introduction

Low Density Parity Check (LDPC) codes are attributed to Gallager who proposed them in his 1960 PhD dissertation [1]. The research was lying in dormant due to high complex computation for encoders, introduction of the Reed-Solomon codes and the concatenation of RS and convolutional codes were considered perfectly suitable for error control coding. In 1980, Tanner [2] introduced graphical representation for LDPC codes well known as Tanner graphs. In mid-1990s, Mackay, Luby and others [3], [4], and [5] resurrected them and noticed their importance apparently independently of the work of Gallager.

In this paper, we present implementation of a real time programmable LDPC encoder that can support code lengths from 576 to 2304 and five different code rates which are $\frac{1}{2}$, 2/3A, 2/3B, 3/4A and 3/4B. Each code rate and code length are supported by different parity check matrices that are computed in real time from the model matrices stored in the memory. The proposed architecture can be implemented in ASIC, FPGA or DSP. We

implemented it on Reconfigurable Instruction Cell Architecture [6]. This architecture belongs to the emerging field of Reconfigurable Computing and is an effort to combine the flexibility and programmability of DSP, performance of FPGA and low power consumption of ASIC in one unified core so that the core can meet the requirement of next generation mobile systems. This paper implements the encoder using the established and architecture specific optimization techniques to exploit the parallelism identified inside the algorithm [6].

The rest of the paper is organized such that section 2 presents the real time programmable encoder, section 3 discusses the implementation and optimization on the architecture while section 4 concludes the paper.

2. Real Time Programmable LDPC Encoder

It is often necessary in any type of wireless communication system to adapt its transmission and reception to the varying channel conditions for maintaining a good QoS. For such adaptation, it is necessary that both transmitting and receiving sections can configure themselves in real time. Such configuration or adaptation can be with respect to code length, code rate, modulation scheme and or different encoding/decoding algorithms. Inside a particular Forward Error Correction (FEC) module, the adaptation is with respect to code length and code rate. LDPC is an example of FEC and the IEEE standard defines a specific range of both code lengths and rates which a WiMax system must support.

Such a real time adaptive architecture for IEEE specified LDPC Encoder is presented in Figure 1. This encoder consists of two parts, H matrix calculation, represented in Figure 1(B), and the actual encoding which is shown in Figure 1(A).

In previous work, H matrix calculation is done offline and the H matrix is then stored in the memory to be used for encoding. These designs can support only one code

length as well as code rate and hence are not adaptable to real time systems.

The IEEE P802.16e/D7 specification defines LDPC code which is based on a set of one or more fundamental LDPC codes. Each of the fundamental codes supports code lengths from 576 to 2304 with code rates 1/2, 2/3A, 2/3B, 3/4A and 3/4B. Each LDPC code in the set of LDPC codes is defined by a matrix H of size m-by-n where n is the length of the code and m is the number of parity check bits in the code. The matrix H is constructed from a set of z-by-z permutation matrices or z-by-z zero matrix. The matrix H is expanded from a binary model matrix H_{bm} of size m_b -by- n_b where $n=z*n_b$ and $m=z*m_b$ where z is an integer ≥ 1 . The model matrix has integer entries ≥ -1 . H matrix is constructed by replacing each '-1' in the model matrix by z-by-z zero matrix. The '0' entry is replaced by z-by-z identity matrix while any other entry greater than 0 is replaced by a z-by-z permutation matrix. The permutations used are circular right shifts and the set of permutation matrices contains the z-by-z identity matrix and its circular right shifted versions.

The H matrix generate section in Figure 1(B) consists of three modules and memory arrays. One memory array stores the model matrices defined in the IEEE standard. Each code rate is represented by a unique model matrix. Five model matrices are stored that correspond to five code rates (1/2, 2/3A, 2/3B, 3/4A, 3/4B). The entire memory consumed by the five model matrices is $12*24 + 8*24 + 8*24 + 6*24 + 6*24 = 960$ bytes. All model matrices have 24 columns whereas the rows depend upon code rates with 12 rows for 1/2, 8 rows for 2/3A, 2/3B, and 6 rows for 3/4A, 3/4B. The purpose of the Configuration Block is to compute the size of the H matrix, and the spreading factor 'z' for a particular code length and code rate. The spreading factor 'z' determines the size of the square matrix which is used to replace each entry in the model matrix to construct the H matrix.

Each model matrix stores values that correspond to H matrix construction for the maximum code length of 2304. For any other code length, the entries in the model matrix is modified according to (1) and (2).

$$p(f, i, j) = \begin{cases} p(i, j), & p(i, j) \leq 0 \\ \left\lfloor \frac{p(i, j) * z}{z_0} \right\rfloor, & p(i, j) > 0 \end{cases} \quad (1)$$

$$p(f, i, j) = \begin{cases} p(i, j), & p(i, j) \leq 0 \\ \text{mod}(p(i, j), z), & p(i, j) > 0 \end{cases} \quad (2)$$

where z and z_0 are the spreading factors for the H matrices of the desired code length and maximum code length respectively. The shift factor $p(f, i, j)$ represents the shift sizes for the corresponding code size and $p(i, j)$ is the shift size for code length of 2304 and this is the i, j th entry in the model matrix stored initially. The $\lfloor x \rfloor$ is flooring function giving the nearest integer towards $-\infty$.

The purpose of the Base Matrix Generate module is to take a model matrix and spreading factor corresponding to a particular code size and then establish a base matrix of the same size as the size of the model matrix but with entries defined by equations 1 and 2. The output of this module is the base matrix which is the basis for the H matrix generation. This matrix is applied to the next module that

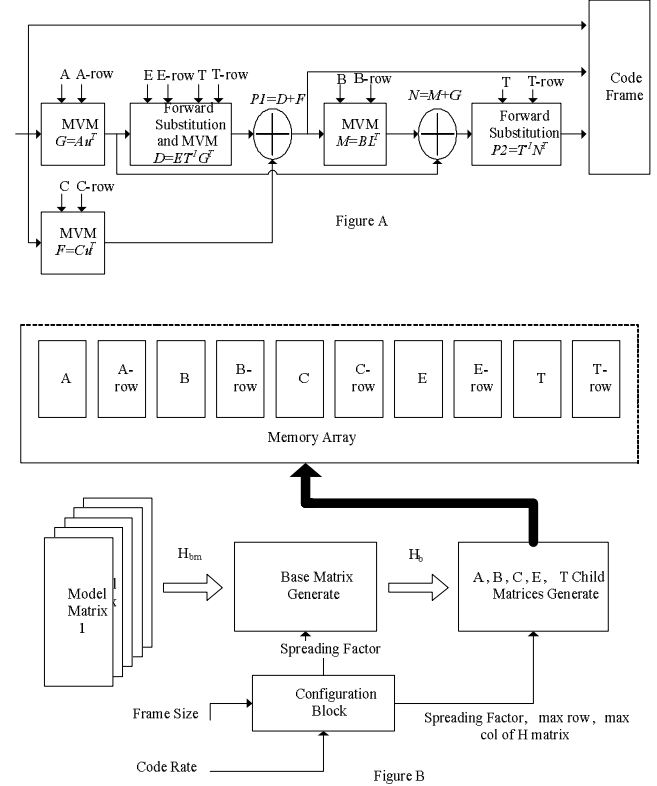


Figure 1: (Real Time Programmable LDPC Encoder)

generates the child matrices A, B, C, E, and T. These children of the H matrix are generated directly instead of the H matrix. The dimension of the child matrices are defined in IEEE 802.16E. If the child matrices are stored conventionally then the total memory consumed by all these matrices for a code size of 2304 and code rate 1/2 is $1152*2304 = 2.53\text{Mbits} = 316\text{Kbytes}$, where 1152 is the row size and 2304 is the column size of the H matrix. This is a very huge memory and dedicating this much memory to the H matrix alone is not recommend to be feasible for either DSP, FPGA or even ASIC implementation. This is for the reason that model matrices and intermediate vectors also require storage and the overall storage then becomes prohibitively high. The total memory required for an example code length of 2304 and code rate of 1/2 can be $= 2.53\text{M(H matrix)} + 1152(\text{information bits}) + (1152+1152+96)(\text{for storing intermediate results as well as parity bits}) = 324.4\text{Kbytes}$.

This huge memory storage requirement can be reduced by adopting the methodology in [8] for storing the child matrices. Since H is sparse regarding number of 1's, a

huge memory reduction is possible if only the indexes of the 1's inside the H matrix or the child matrices are stored. As an example, the total number of 1's in the H matrix for 2304 code length and $\frac{1}{2}$ code rate is 8024 or approximately 8K. If indexes of only the 1's inside the H matrix are stored then only 8Kwords or 16Kbytes of memory is needed in stead of 316Kbytes which is equivalent to a reduction of about 20 times in memory. The authors in [8] have stored the indexes of only 1's inside the child matrices. Each child matrix is accompanied by another array which stores a 1 for the index of H and zero for the end of the row. Therefore, a zero in this array will indicate the end of the row of the associated child matrix. This is an extra overhead which is a 1-bit array of size 8K.

Since the encoder is designed to support all code rates and code lengths, the maximum memory sizes are computed based on exhaustive MatLab simulation of the maximum code length and all code rates. The maximum memory sizes for A, B, C, E and T matrices are determined and then used in the design.

The 'A, B, C, E, T' child matrices Generate' module generates the child matrices and stores the indexes of the 1's inside them in the corresponding memory array. The memory array x-row stores '0' for each entry and '1' for the end of each row. Here x represents a child matrix. This array is used as pointer to indicate the end of row for matrix-vector computation.

After child matrices generation, they are applied to the encoder (figure 1(A)) together with the information bits. The processing in each block is similar to the processing defined in [8]. Here MVM is the matrix-vector-multiplication with a vector output. Forward Substitution is used to obtain $y^* T^T = x$ using $y = T^* x$.

The first MVM module in Figure 1(A) generate $F=C*u^T$ and $G = A*u^T$. The vector G is applied to the next module to generate $D=E*T^T*G^T$. D and F are vector added in the next VA module shown by the encircled plus sign. The output of the VA module is PI which is a subset of the parity check bits that contribute to the code bits. The next MVM module generates the vector $M=B*PI^T$. The VA adder module then adds using modulo-2 addition and generates $N=M+G$. The last module uses forward substitution to generate $P2=T^T*N^T$ and the frame is then given by $c = [u \ P1 \ P2]$.

The serial implementation can be implemented on either ASIC, FPGA or DSP. For speed improvement on either ASIC or FPGA a pipelined version of the encoder is also presented as shown in Figure 2.

This is a four stage pipelined encoder which promises to improve the speed by approximately 2.5 times at the expense of increased memory of size $(1152+1152)*3$ bits. The 2.5 times is because of the second stage as the second stage is the most time consuming block and cannot be divided into smaller blocks. The architecture can implement both due to its flexibility.

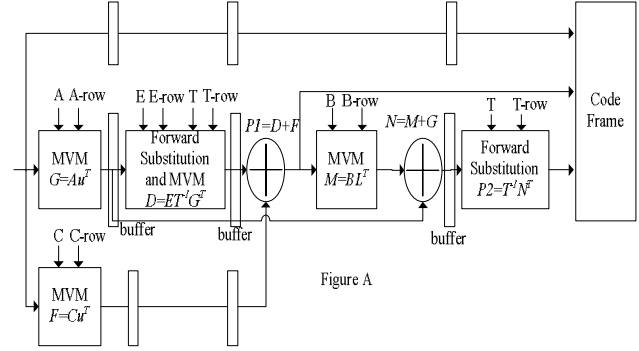


Figure 2: (Pipelined Real Time Encoder)

3. Implementation on Reconfigurable Instruction Cell based Architecture

The architecture [6] is a dynamic reconfigurable fabric which has coarse grained heterogeneous functional units (cells) connected to each other through a reconfigurable interconnect structure. The functional units support primitive operators that can perform addition/subtraction, multiplication, logic, multiplex, shift and register operations. Additional functional units (cells) are provided to handle control/branch operations. Each functional unit operates on 32-bit operands. The current implementation of the architecture is structured so that it can have 8 memory read and 8 memory write interfaces and 16 memory banks. The memory width of each bank is 8-bit and 16 banks provide four 32-bit memory accesses in parallel.

There are three major optimization techniques for increasing the throughput [9]. These are classified as Loop Unrolling, Split Computation and Multi Sampling. Loop unrolling is applied to independent iterations to use the resources as much as possible.

With each memory interface, 32-bit word can be accessed thus allowing four 32-bit memory access operations in parallel. The architecture has other optimizations like converting small conditional branches into simple multiplexing operations, storing intermediate results in distributed registers to avoid more costly memory access and providing more functional units for parallel execution of operations. The functional units and interconnects in the architecture allow the possible exploitation of instruction level parallelism inside the code to as much as the number of functional units. In this paper, the objective is to exploit generic and architecture specific optimizations for increased throughput.

The code for the encoder consists of Vector_Add, Matrix_Vector_Mult, Forward_Substitution, Base_Matrix_Generate, and H_Matrix_Generate. Of all these modules Base_Matrix_Generate and H_Matrix_Generate are called once at start up or whenever the parameters on which H matrix generation depends get changed. These parameters are the code rate

and the number of information bits used to produce code word of a particular length. The code is first simulated without applying any optimization techniques. The simulated result shows only 3.5 Mbps throughput using code length of 2304 with code rate $\frac{1}{2}$. The throughput has been increased to 10 Mbps for the same rate and code length by applying generic and architecture specific optimization to different modules. These are discussed below:

3.1 Optimization specific to algorithm

This is carried out in the matrix A that has different distribution for 1's for different code rates. Separate functions have been written for all code rates as well as code lengths. The LDPC encoder supports 5 different code rates and code lengths from 576 to 2304. This is carried out in order to help the reconfigurable fabric optimize the code efficiently. By doing so, many if statements have been eliminated that would cause unnecessary jumps and more time in execution. The resources which mainly consist of logic gates for these functions are increased to unroll the loop by four times. This has caused significant reduction in execution time.

3.2 Optimization of Vector_Add

The aim of this module is to provide modulo-2 addition of two input vectors. Since the addition is bit wise, enough parallelism is present inside the module. This parallelism is exploited for increased throughput through simple loop unrolling. The original code is as:

```
for (i=0;i<zf;i++) { // zf is spreading factor and is 24 to 96
*(Ptrtovout+i) = *(Ptrtovin1+i) ^ *(Ptrtovin2+i); }
```

The above code has $2*zf$ read and zf write memory operations. It will take $3*zf$ cycles in any DSP processor and if 4 nsec is taken as the memory access time then its execution will take $12*zf$ nsec on a DSP Processor. If the code is expanded as below then all the four memory reads from the same memory array can be done in just one 4 nsec time and all the four writes to the same memory array can also be completed in just one 4 nsec time. The memory access time is then $(zf/4 \text{ (read)} + zf/4 \text{ (read)} + zf/4 \text{ (write)}) * 4 = 3*zf$ which is equivalent to a 4 times reduction in memory access time.

```
for (i=0;i<zf;i+=4) {
*(Ptrtovout+i+0)= *(Ptrtovin1+i+0) ^ *(Ptrtovin2+i+0);
*(Ptrtovout+i+1)= *(Ptrtovin1+i+1) ^ *(Ptrtovin2+i+1);
*(Ptrtovout+i+2)= *(Ptrtovin1+i+2) ^ *(Ptrtovin2+i+2);
*(Ptrtovout+i+3)= *(Ptrtovin1+i+3) ^ *(Ptrtovin2+i+3); }
```

In the case of packed computation in which four independent data bytes are accessed as one data pack and processed as separate and independent 4 data bytes, only two read interfaces and one write interfaces are required

thus reducing the burden on the resources. The architecture also supports 8 adders and 8 XOR operations. This implies that address calculation and XORing can be done in parallel thus saving more execution time.

3.3 Replacing Jumps with Multiplexing

The architecture executes the code in steps. A step is defined as combination of instructions that can be executed in the fabric provided by the architecture. The group of instructions need not be necessarily independent of each other. A step is determined by the number of available resources, conditional branch and the length of the critical path.

The architecture is structured to support only one jump per step. If a function or a piece of code generates more than one jumps then the architecture places them in separate steps. The main idea behind optimization is to place as much code inside one step as possible. If the entire function is placed in one step then maximum utilization of resources can be achieved and the time of execution will most probably be less than the case in which the same function takes more than one step. For the jump to be replaced with multiplexing, the variables and the paths for the execution must be independent of each other and must not include a memory read in its conditioning. This can be explained from the following code.

```
if(*(PtrtoMrow+j+4))
{
*(Ptrtovout+i) = (*(PtrtoM+j+4)+Ptrtovin) ^
*(Ptrtovout+i);
j++; } }
```

The above code will definitely produce a jump as the memory access in the code inside the parenthesis is made dependent upon the 'if' condition. If the 'if' condition is true then there will be memory access otherwise not. The jump produced by this code will cause underutilization of resources. This code can be structured like the following to avoid the instantiation of the jump cell.

```
temp1 = *(PtrtoMrow+j+4)
temp2 = (*(PtrtoM+j+4)+Ptrtovin) ^ *(Ptrtovout+i);
temp3 = *(Ptrtovout+i);

*(Ptrtovout+i) = temp1 ? temp2 : temp3;
```

By storing the intermediate values in the registers and making the variables independent of each other, the instantiation of the jump cell has been avoided. The same strategy has been adopted in all other places inside the code to either reduce or eliminate the jumping and using multiplexers to perform the same functionality as is performed with jump cell. Loop unrolling has been used in addition to the above mentioned techniques and the execution time and throughput has been calculated as follows:

One time execution of LDPC Encoder = 554.926

Number of steps taken:	106129
Two times execution of LDPC Encoder	= 666.232 μ sec
Number of steps taken:	123674
Execution time of actual encoding to be used in real time	= 666.232 – 554.926 = 111.3 μ sec
Number of steps of actual encoding	= 123674 – 106129 = 17545
Execution time per bit	= 111.3/1152 = 96.61 nsec/bit
Throughput	= 1/96.61 = 10.4 Mbps ($\frac{1}{2}$ rate).

For code rate $\frac{3}{4}$, the throughput is measured to be approximately 19 Mbps. This is the highest code rate that IEEE 802.16 defines for the irregular LDPC codes.

3.4 Hardware Pipelining the code

The hardware pipelining is about pipelining the memory read, computation and memory write operations inside a code. If a loop consists of all these processes and brought inside one step then the step will loop into itself. Such a step can be pipelined and using a two stage pipeline, the execution time can be reduced by about 2.5 times.

An example of a hardware pipelining is shown in Figure 3 and 4. This is the graph of a data path that reads from the memory, adds together the words read from the memory and then writes the results to the same location inside the memory. The complete datapath has been brought into one the architecture step and it loops to itself. The execution time of this step has been measured to be 28 nsec and if the maximum iteration is 120, the execution time of the entire loop is $28 \times 120 = 3.36 \mu$ sec.

Using pipelining, the data path can be divided into three smaller data paths with critical path length not exceeding 10 nsec. The time of execution of the step is now approximately 2.5 times less than the total time and is about $3.36/2.5 = 1.344 \mu$ sec. By using hardware pipelining, the throughput of the encoder is achieved in the range from 26 Mbps for $\frac{1}{2}$ code rate to 47 Mbps for $\frac{3}{4}$ code rate. This is better than an FPGA [7] implementation which achieved 22 Mbps from a real time programmable LDPC encoder for regular codes.

Another solution to increase the throughput is to extend pipeline to the individual modules as shown in Figure 2. There are two possibilities to this type of approach. The first is structuring existing fabric to support inter-module pipeline in addition to the existing pipeline for the intra-module. The second is to use more than one fabrics and build the architecture shown in Figure 2 using these fabrics in a pipelined fashion. We used the second option and simulated the modules in a way as would be executed in a pipeline. We simulated first the MVM that performs multiplication of A with the input information frame. After this, Forward Substitution Module has been simulated. The third simulation was carried out for Vector Add, MVM module for matrix B and the Vector Add after MVM B. The last is carried out for the forward substitution in which matrix T is involved.

This type of build up resulted in 3 times improvement in throughput. The throughput is then 78 Mbps. The fabric

used in this case consists only 64 cells and has much less resources compared to the fabric used for the implementation of a single LDPC encoder.

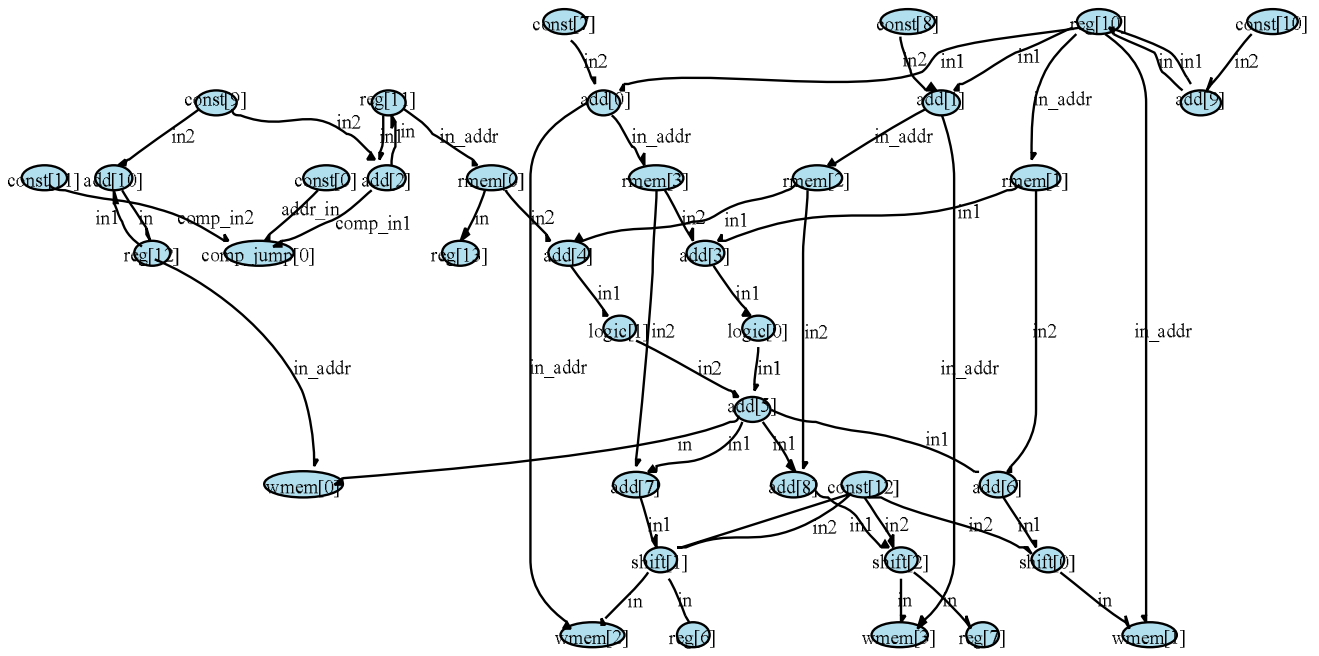
4. Conclusion

In this paper, the authors have presented a novel architecture for the real time LDPC encoder as specified in the IEEE P802.16e/D7 standard targeting WiMax applications and discussed its optimization on a reconfigurable instruction cell architecture (the architecture). Several algorithmic and the architecture specific optimization techniques have been applied. The throughput achieved without pipelining is in the range from 10 to 19 Mbps while with pipelining, it is in the range from 26 to 47Mbps. This is a considerable throughput while providing the flexibility and programmability from a high level 'C' programming language.

The multi-core solution is also provided as an example for throughput enhancement to meet emerging wireless standards. The multicore can provide a maximum throughput of 78 Mbps.

5. References

- [1] R. Gallager, "Low-Density parity-check codes", IRE Trans. Information Theory, pp. 21-28, Jan. 1962
- [2] R.M. Tanner, "A recursive approach to low complexity codes", IEEE Trans. Information Theory, pp. 533-547, Sept. 1981
- [3] D. MacKay and R. Neal, "Good codes based on very sparse matrices", in Cryptography and coding, 5th IMA Conf., C.Boyd, Ed., Lecture Notes in Computer Science, pp. 100-111, Berlin, Germany: Springer, 1995
- [4] D. Mackay, "Good error correcting codes based on very sparse matrices", IEEE Trans. Information Theory, pp. 399-431, March 1999
- [5] N. Alon and M. Luby, "A linear time erasure-resilient code with nearly optimal recovery", IEEE Trans. Information Theory, pp. 1732-1736, Nov. 1996.
- [6] Ying Yi, Nousias, I., Milward, M., Khawam, S. Arslan, T., Lindsay, I., "System-level Scheduling on Instruction Cell Based Reconfigurable Systems" Design, Automation and Test in Europe, 2006. DATE '06. Proceedings, Volume 1, 6-10 March 2006 Page(s):1 - 6 IEEE P802.16E/D7 Specification published in 2006
- [7] L. Dong-U, L. Wayne, W. Connie, J. Christopher, "A flexible hardware encoder for low density parity check codes", Proc. IEEE Symp. Field-Programmable Custom Computing, 2004
- [8] http://uat.freescale.com/files/technology_publications/doc/Papers/Eintell5441.pdf
- [9] SC140 Application Note: How to implement a Viterbi Decoder on the StarCore SC140. Application number ANSC140VIT: Available on www.freescale.com



step[LNewRICA_latest_8bitmemory_sL539] (3)

Figure 3:(Non-pipelined Step of the architecture) (rmem is read memory cell, wmem is write memory cell, const is a cell to store constant values, add is adder cell, reg is register cell and shift is shifter cell, comp_jump is logical cell)

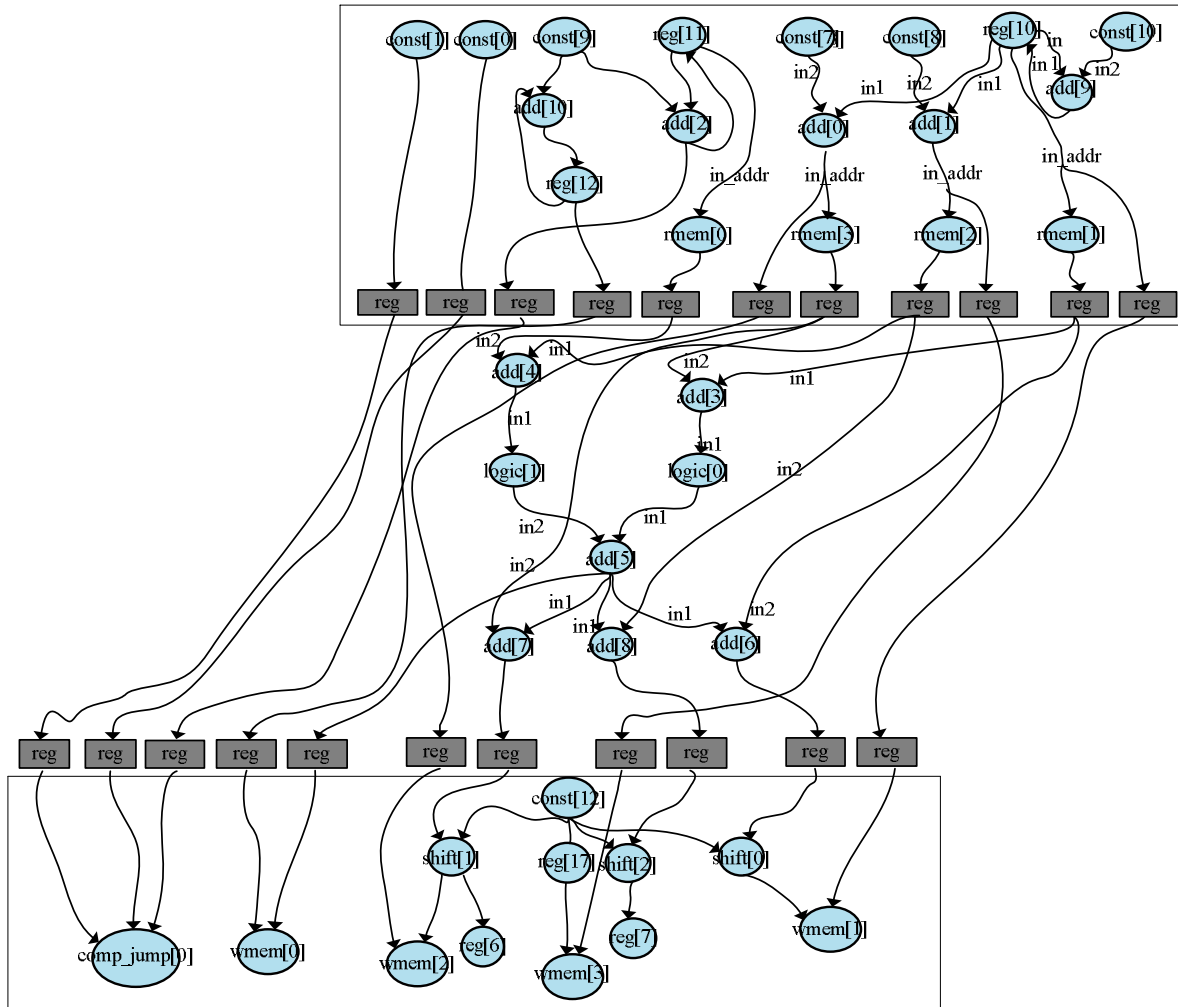


Figure 4: (Pipelined Step of the architecture)