

VLSI Architectures for Communications and Signal Processing

Kiran Gunnam

IEEE SSCS Distinguished Lecturer
Director of Engineering, Violin Memory

Outline

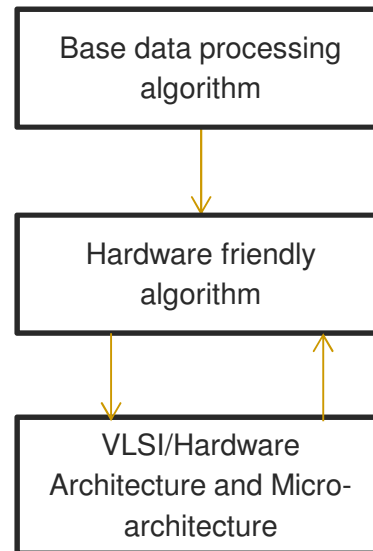
Part I

- Trends in Computation and Communication
- Basics
- Pipelining and Parallel Processing
- Folding, Unfolding, Retiming, Systolic Architecture Design

Part II

- LDPC Decoder
- Turbo Equalization, Local-Global Interleaver and Queuing
- Error Floor Mitigation (Brief)
- T-EMS (Brief)

VLSI Architectures for Communications and Signal Processing



A systematic design technique is needed to transform the communication and signal processing algorithms to practical VLSI architecture.

- Performance of the base algorithm has to be achieved using the new hardware friendly algorithm
- Area, power, speed constraints govern the choice and the design of hardware architecture.
- Time to design is increasingly becoming important factor: Configurable and run-time programmable architectures
- More often, the design of hardware friendly algorithm and corresponding hardware architecture involves an iterative process.

Communication and Signal Processing applications

- Wireless Personal Communication – 3G,B3G,4G,...etc.
– 802.16e,802.11n,UWB,...etc.
- Digital Video/Audio Broadcasting
– DVB-T/H, DVB-S,DVB-C, ISDB-T,DAB,...etc.
- Wired Communications
– DSL, HomePlug, Cable modem, etc.
- Storage
- -Magnetic Read Channel, Flash read channel
- Video Compression
- TV setup box

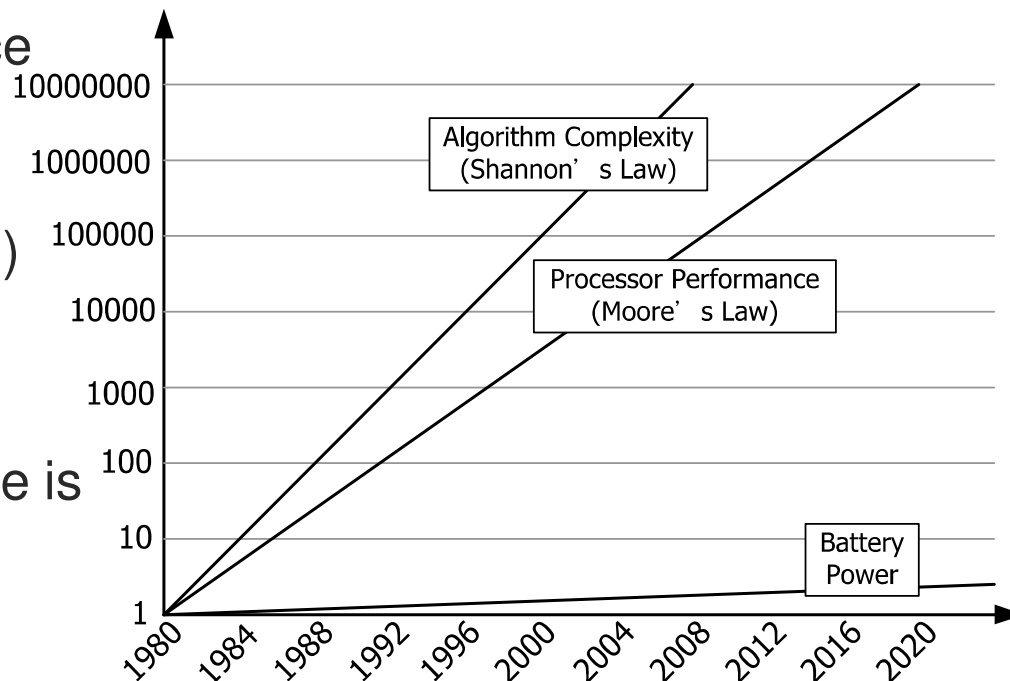
Convergence of Communications and Semiconductor technologies

- High system performance
- – Increase Spectrum efficiency of modem (in bits/sec/Hz/m³)
 - Multi-antenna diversity
 - Beamforming
 - Multi-user detection
 - Multi-input Multi-output (MIMO) Systems • Etc.
- • High silicon integrations
 - Moore's Law
 - High-performance silicon solutions
 - Low power and cost
 - Mobile devices getting more computation, vision and graphics capabilities

Challenges in VLSI for Communication and Signal Processing

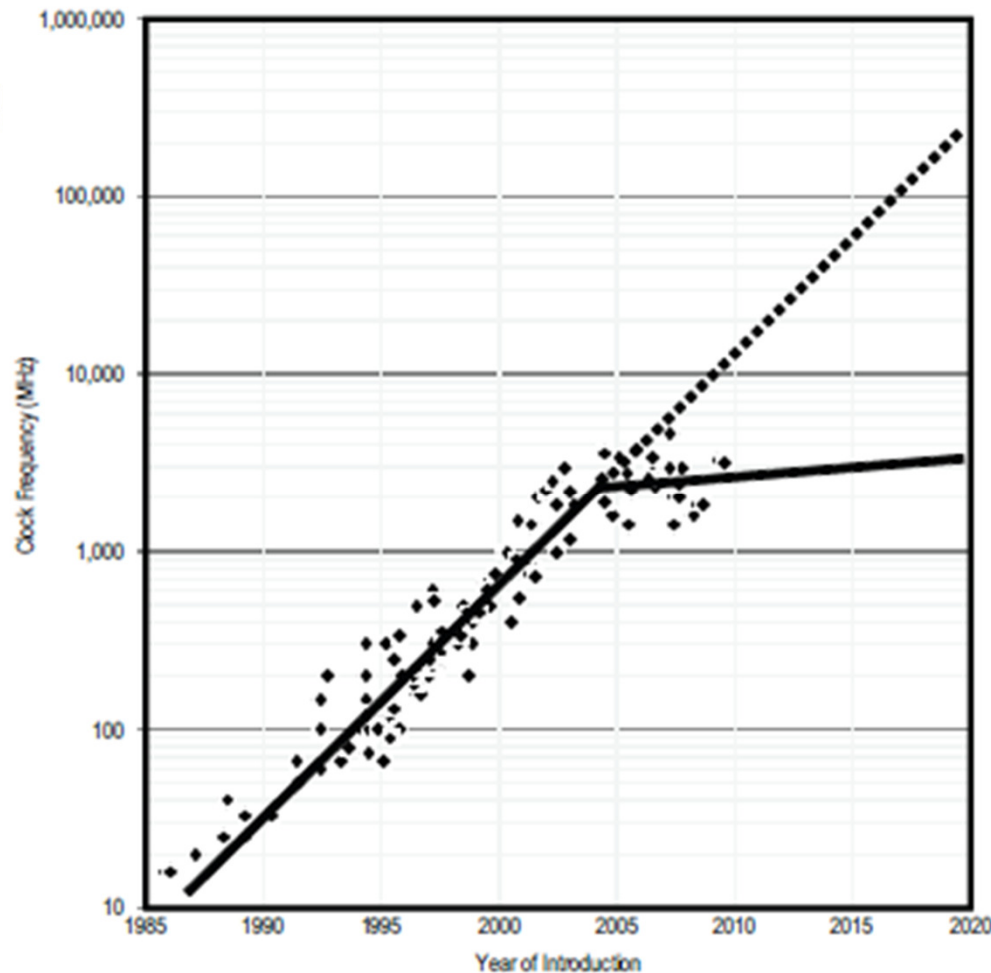
- How to bridge the gap between communication algorithms and IC capabilities.
- Efficient and Flexible DSP VLSI methods considering communication algorithmic requirements
- - High performance
 - Flexibility
 - Low energy
 - Low cost (design)
 - Low cost (area)

While chip performance is increasing, algorithm complexity for new systems is outpacing it.



Courtesy: Ravi Subramanian (Morphics)

Single Processor Performance Trends



Courtesy: NAE Report, "The Future of Computing Performance: Game Over or Next Level?"

8/18/2013

FIGURE S.1 Historical growth in single-processor performance and a forecast of processor performance to 2020, based on the ITRS roadmap.

The dashed line represents expectations if single-processor performance had continued its historical trend. The vertical scale is logarithmic. A break in the growth rate at around 2004 can be seen.

Before 2004, processor performance was growing by a factor of about 100 per decade; since 2004, processor performance has been growing and is forecasted to grow by a factor of only about 2 per decade.

In 2010, this expectation gap for single-processor performance is about a factor of 10; by 2020, it will have grown to a factor of 100.

Note that this graph plots processor clock rate as the measure of processor performance. Other processor design choices impact processor performance, but clock rate is a dominant processor performance determinant.

Scaling Trends

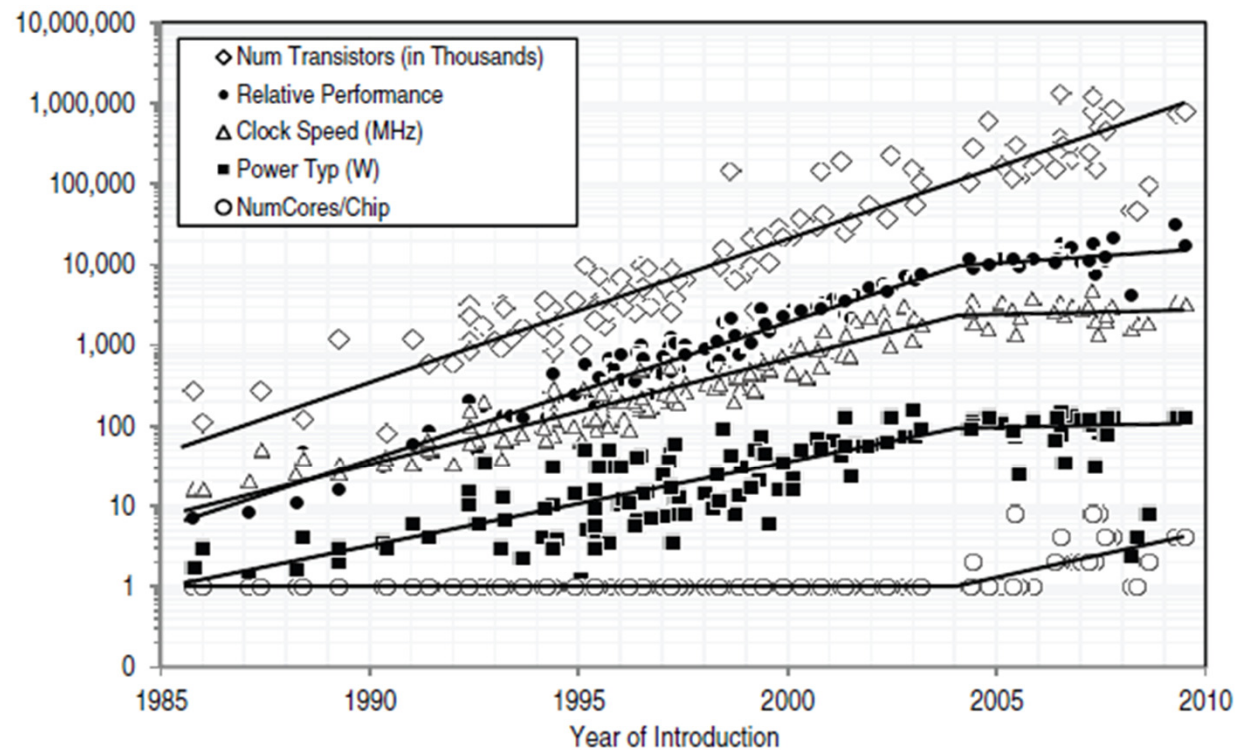


FIGURE 2.1 Transistors, frequency, power, performance, and cores over time (1985-2010). The vertical scale is logarithmic. Data curated by Mark Horowitz with input from Kunle Olukotun, Lance Hammond, Herb Sutter, Burton Smith, Chris Batten, and Krste Asanović.

Courtesy: NAE Report, "The Future of Computing Performance: Game Over or Next Level?"

Why Dedicated Architectures?

Energy Efficiency

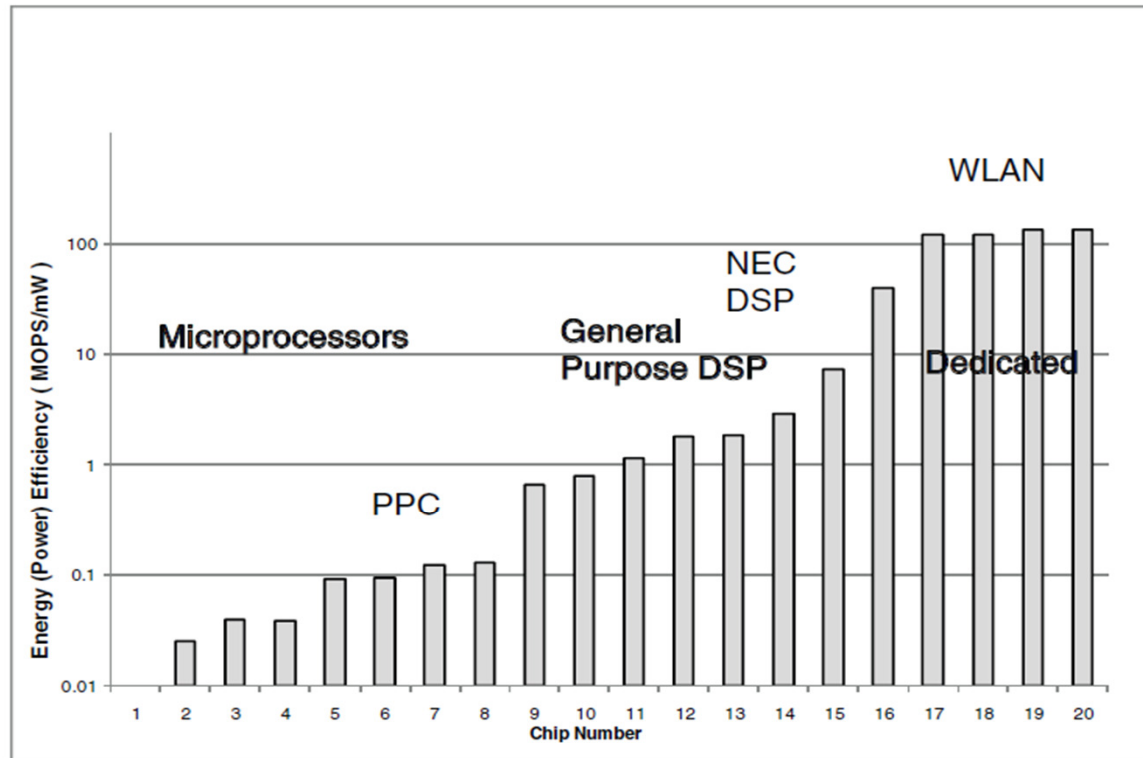


FIGURE 3.5 Energy efficiency comparison of CPUs, DSPs, and ASICs. SOURCE: Robert Brodersen of the University of California, Berkeley, and Teresa Meng of Stanford University. Data published at International Solid-State Circuits Conference (0.18- to 0.25- μm).

Courtesy: NAE Report, "The Future of Computing Performance: Game Over or Next Level?"

Why Dedicated Architectures?

Area Efficiency

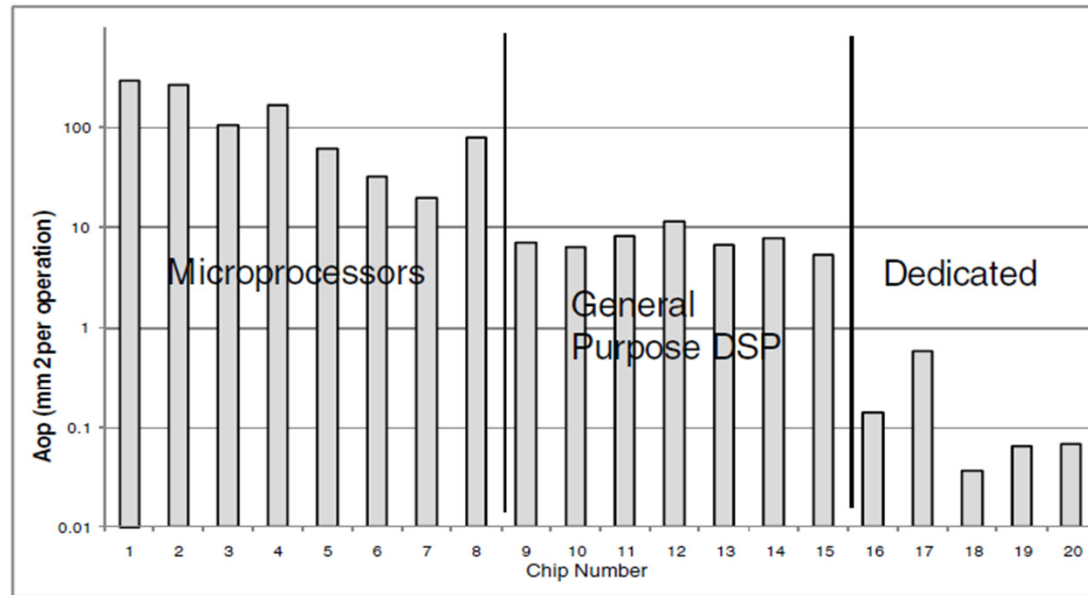


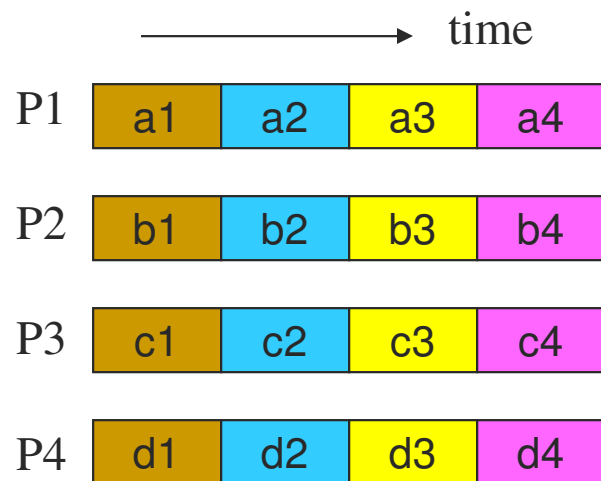
FIGURE 3.6 Area efficiency comparison of CPUs, DSPs, and ASICs. SOURCE: Robert Brodersen of the University of California at Berkeley and Teresa Meng of Stanford University.

NAE Report Recommendation: “Invest in research and development of parallel architectures driven by applications, including enhancements of chip multiprocessor systems and conventional data-parallel architectures, cost effective designs for application-specific architectures, and support for radically different approaches.”

Courtesy: NAE Report, “The Future of Computing Performance: Game Over or Next Level?”

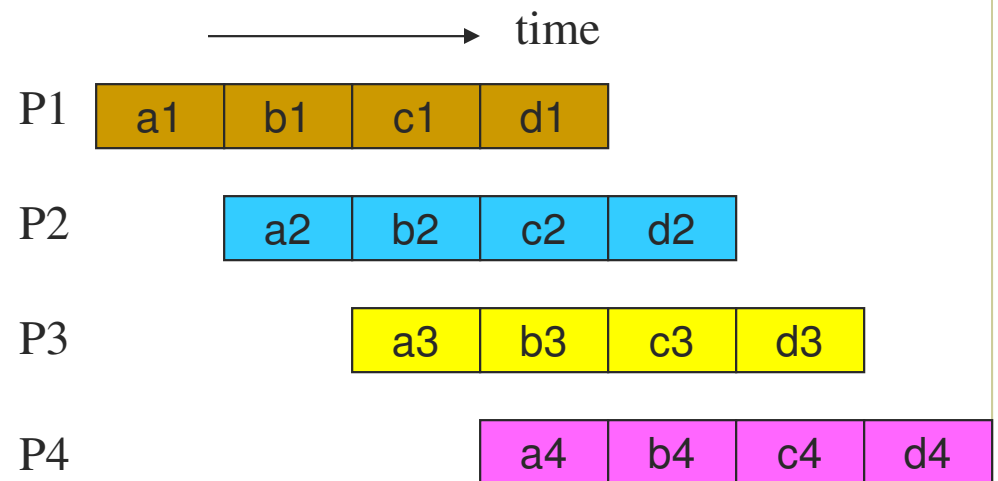
Basic Ideas

■ Parallel processing



Less inter-processor communication
Complicated processor hardware

■ Pipelined processing



More inter-processor communication
Simpler processor hardware

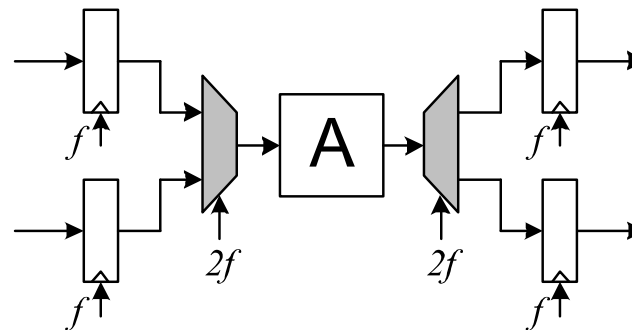
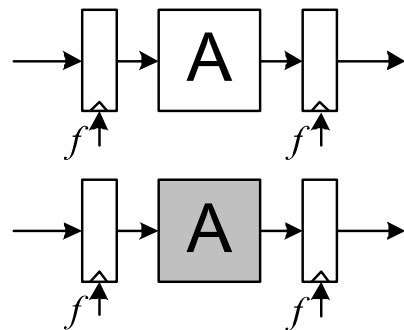
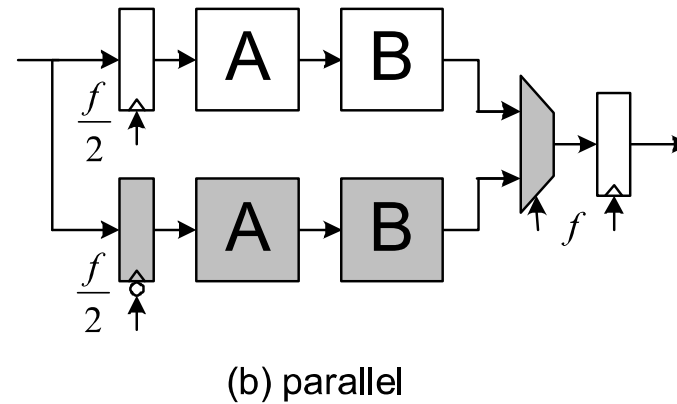
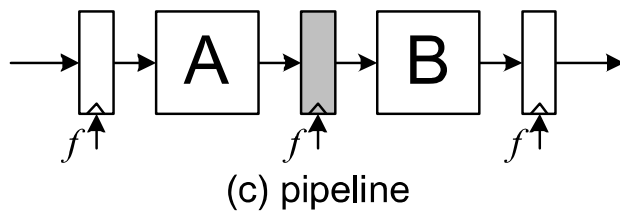
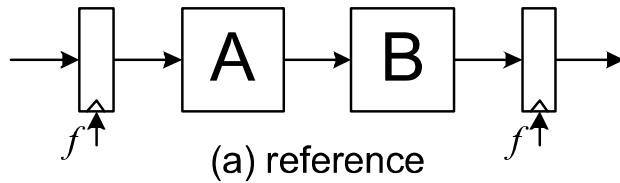
Colors: different types of operations performed

a, b, c, d: different data streams processed

Can combine parallel processing and pipelining-will have 16 processors instead of 4.

Courtesy: Yu Hen Hu

Basic Ideas

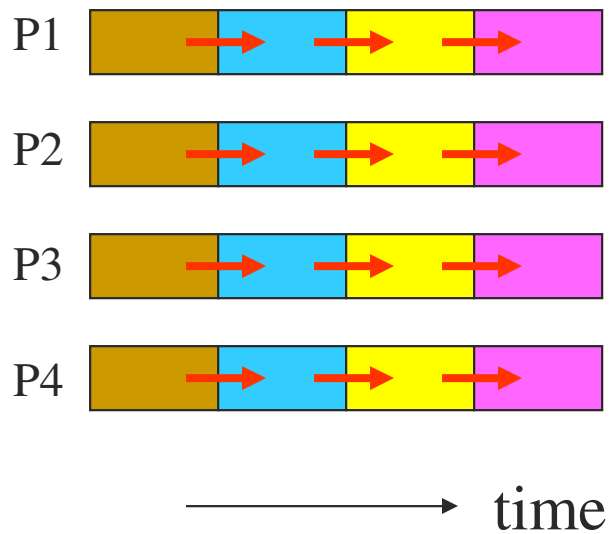


Basic micro-architectural techniques: reference architecture (a), and its parallel (b) and pipelined (c) equivalents. Reference architecture (d) for time-multiplexing (e). Area overhead is indicated by shaded blocks.

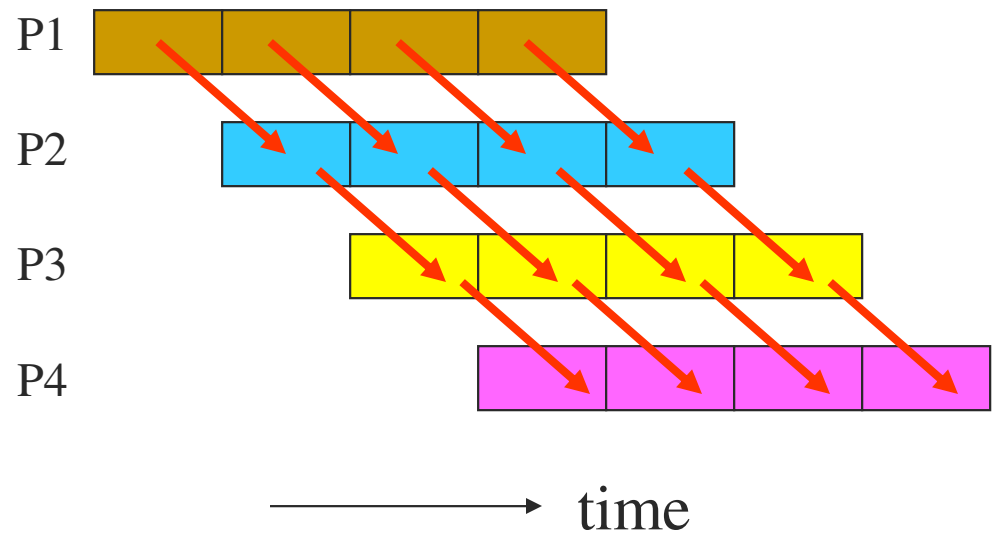
Bora et. al, "Power and Area Efficient VLSI Architectures for Communication Signal Processing", ICC 2006

Data Dependence

- Parallel processing requires NO data dependence between processors

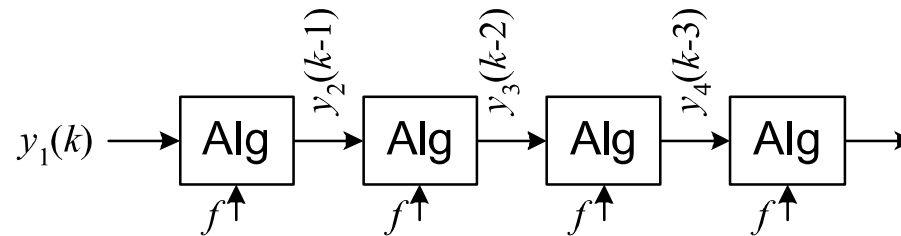


- Pipelined processing will involve inter-processor communication

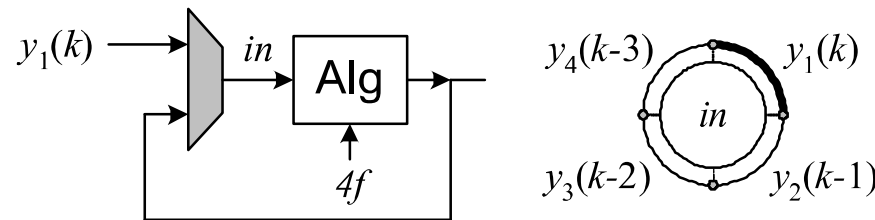


Courtesy: Yu Hen Hu

Folding



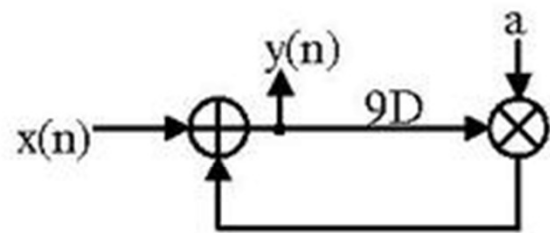
(a) reference



(b) folding

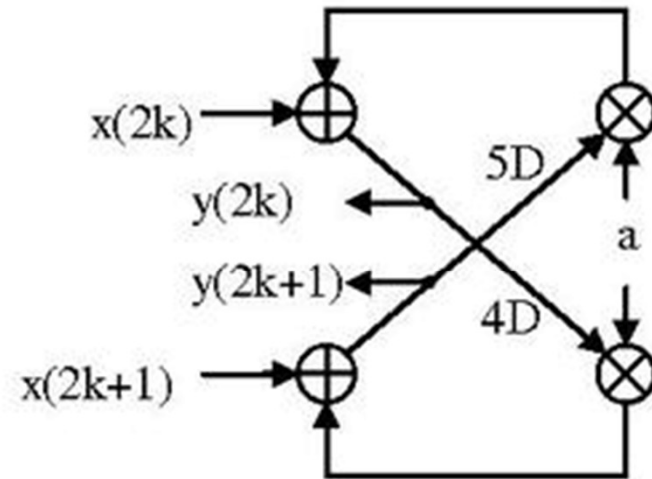
Concept of folding: (a) time-serial computation, (b) operation folding. Block *Alg* performs some algorithmic operation.

Unfolding



$$y(n) = ay(n-9) + x(n)$$

unfolding



$$y(2k) = ay(2k-9) + x(2k)$$

$$y(2k+1) = ay(2k-8) + x(2k+1)$$

transform the dfg of 1 input and 1 output into dfg that receives 2 inputs and produce 2 outputs at each time.

Courtesy: Yu Hen Hu

Block Processing

- One form of vectorized parallel processing of DSP algorithms. (Not the parallel processing in most general sense)
- Block vector: $[x(3k) \ x(3k+1) \ x(3k+2)]$
- Clock cycle: can be 3 times longer
- Original (FIR filter):

$$y(n) = a \cdot x(n) + b \cdot x(n-1) + c \cdot x(n-2)$$

Courtesy: Yu Hen Hu

- Rewrite 3 equations at a time:

$$\begin{bmatrix} y(3k) \\ y(3k+1) \\ y(3k+2) \end{bmatrix} = a \begin{bmatrix} x(3k) \\ x(3k+1) \\ x(3k+2) \end{bmatrix} + b \begin{bmatrix} x(3k-1) \\ x(3k) \\ x(3k+1) \end{bmatrix} + c \begin{bmatrix} x(3k-2) \\ x(3k-1) \\ x(3k) \end{bmatrix}$$

- Define block vector $\mathbf{x}(k) = \begin{bmatrix} x(3k) \\ x(3k+1) \\ x(3k+2) \end{bmatrix}$
- Block formulation:

$$\mathbf{y}(k) = \begin{bmatrix} a & 0 & 0 \\ b & a & 0 \\ c & b & a \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} 0 & c & b \\ 0 & 0 & c \\ 0 & 0 & 0 \end{bmatrix} \mathbf{x}(k-1)$$

Systolic Architectures

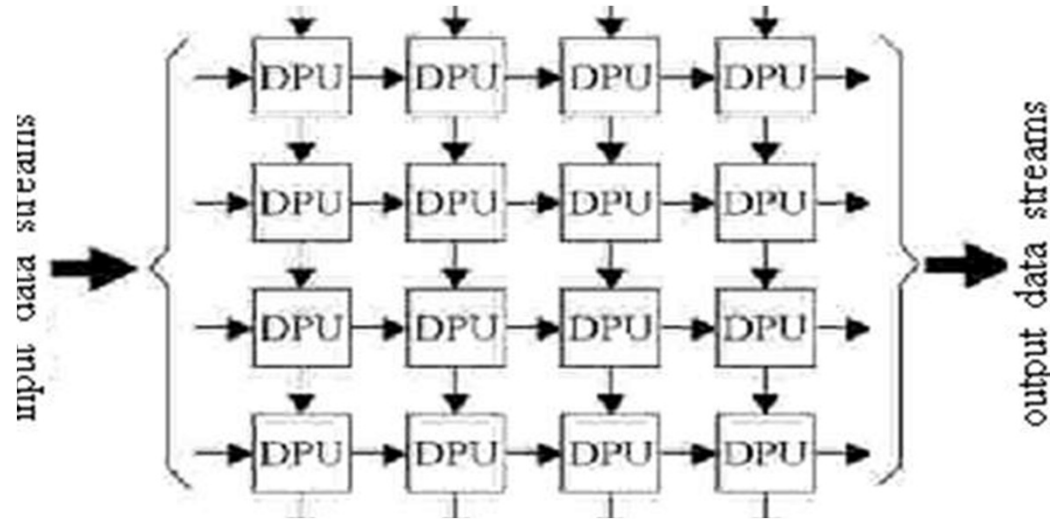


Figure by Rainier

Matrix-like rows of data processing units called cells.

Transport Triggered.

Matrix multiplication $C=A*B$.

A is fed in a row at a time from the top of the array and is passed down the array,

B is fed in a column at a time from the left hand side of the array and passes from left to right.

Dummy values are then passed in until each processor has seen one whole row and one whole column.

The result of the multiplication is stored in the array and can now be output a row or a column at a time, flowing down or across the array.

LDPC DECODER

8/18/2013

Requirements for Wireless systems and Storage Systems

Magnetic Recording systems

- Data rates are 3 to 5 Gbps.
- Real time BER requirement is $1e-10$ to $1e-12$
- Quasi real-time BER requirement is $1e-15$ to $1e-18$
- Main Channel impairments: ISI+ data dependent noise (jitter)
+ erasures
- Channel impairments are getting worse with the increasing recording densities.

Wireless Systems:

- Data rates are 0.14 Mbps (CDMA 2000) to 326.4 Mbps (LTE UMTS/4GSM) .
- Real time BER requirement is $1e-6$
- Main Channel impairments: ISI (frequency selective channel)
 - + time varying fading channel
 - + space selective channel
 - + deep fades
- Increasing data rates require MIMO systems and more complex channel estimation and receiver algorithms

In general the algorithms used in wireless systems and magnetic recording systems are similar. The increased complexity in magnetic recording system stems from increased data rates while the SNR requirements are getting tighter.

For ISI channels, the near optimal solution is turbo equalization using a detector and advanced ECC such as LDPC.

Introduction to Channel Coding

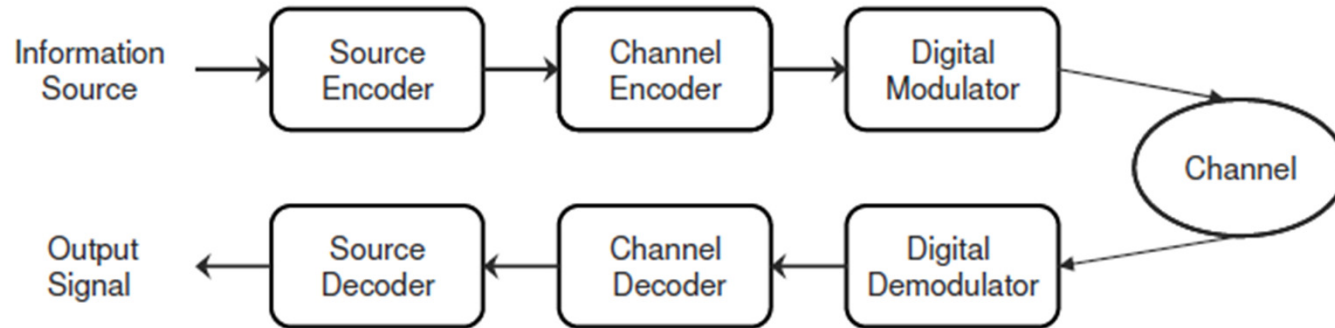


Figure 1: Block Diagram of Communication System

- ❑ Channel Encoder introduces redundancy in transmitted bit stream
- ❑ Channel decoder use the redundancy to correct the errors due to channel impairments and noise.
- ❑ Low-Density Parity-Check (LDPC) Code is the best available error correction code.

Some Notation and Terminology

- k - number of input (information) bits.
 n – number of output (coded) bits
- Code is a set of 2^k vectors of length n
- Encoder is a specific mapping of 2^k inputs to codewords
- Decoder tries to recover the information bits from the received code word that is corrupted with channel impairments
- d_{\min} – minimum distance of the code
Smallest Hamming distance between any two codewords.

Courtesy: Dr. Krishna Narayanan (Texas A&M)

Shannon Capacity and Channel Codes

- The **Shannon limit** or **Shannon capacity** of a communications channel is the theoretical maximum information transfer rate of the channel, for a particular noise level.
- Random and long code lengths achieve channel capacity.
- To construct a random code, pick 2^k codewords of length n at random. Code is guaranteed to be good as k !
- Decoding random codes, require storage of 2^k codewords
- There are only about 10^{82} ($\sim 2^{276}$) atoms in the universe.
- Encoding/Decoding complexities don't increase drastically with k
- Storage does not increase drastically with k
- Randomness Vs Structure
- Random codes are good
- But structure is needed to make it practical

Courtesy: Dr. Krishna Narayanan (Texas A&M)

Coding Theory Advances

- There are two kinds of codes: Block Codes and Convolutional codes
- Block Codes: In an (n,k) block code, k bits are encoded in to n bits. Block code is specified by $k \times n$ generator matrix G or an $(n-k) \times n$ parity check matrix H
- Examples: Hamming, BCH, Reed Solomon Codes. Hard decision decoding is used. Soft decoding possible- but complex.
- Convolutional codes: Can encode infinite sequence of bits using shift registers. Soft decision decoding such as viterbi can achieve optimal maximum likelihood decoding performance.
- Turbo Codes (1993): Parallel concatenated convolutional codes.
- Rediscovery: **LDPC** Block code(**1962, 1981, 1998**). Near shannon limit code, Efficient soft decoding (message passing) and with iterations.

Progress in Error Correction Systems

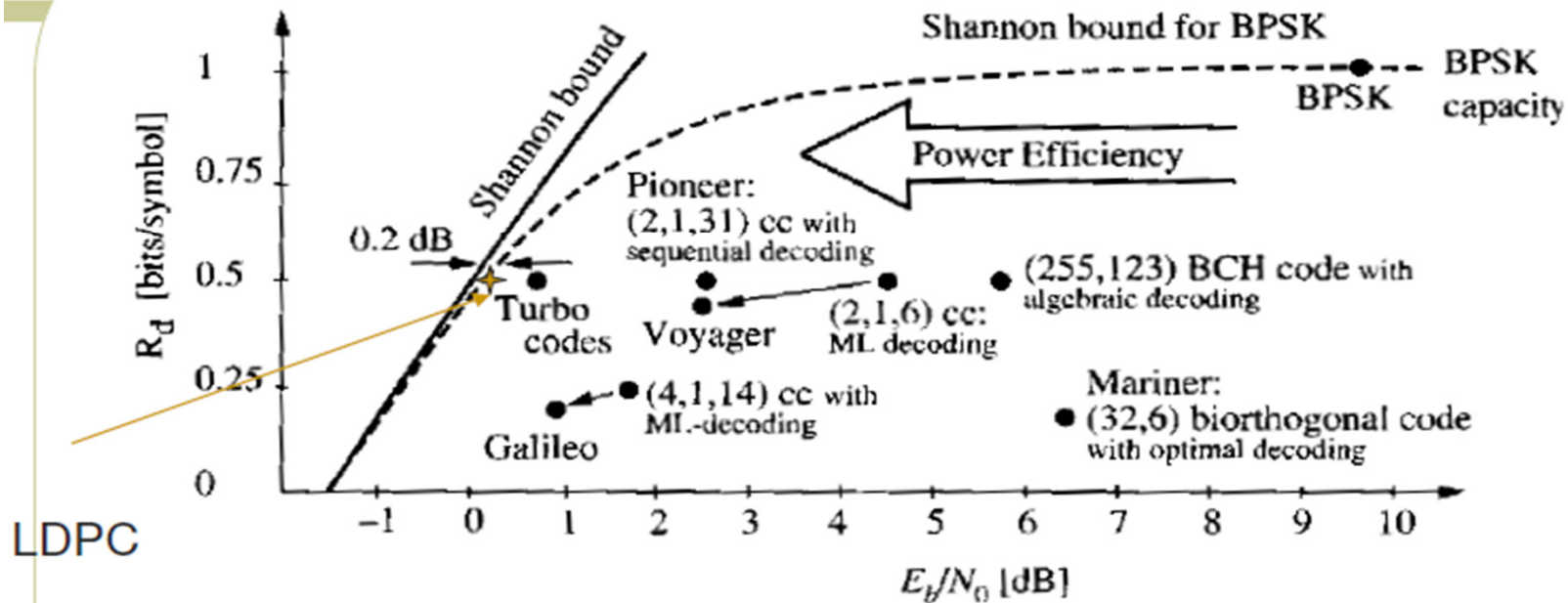
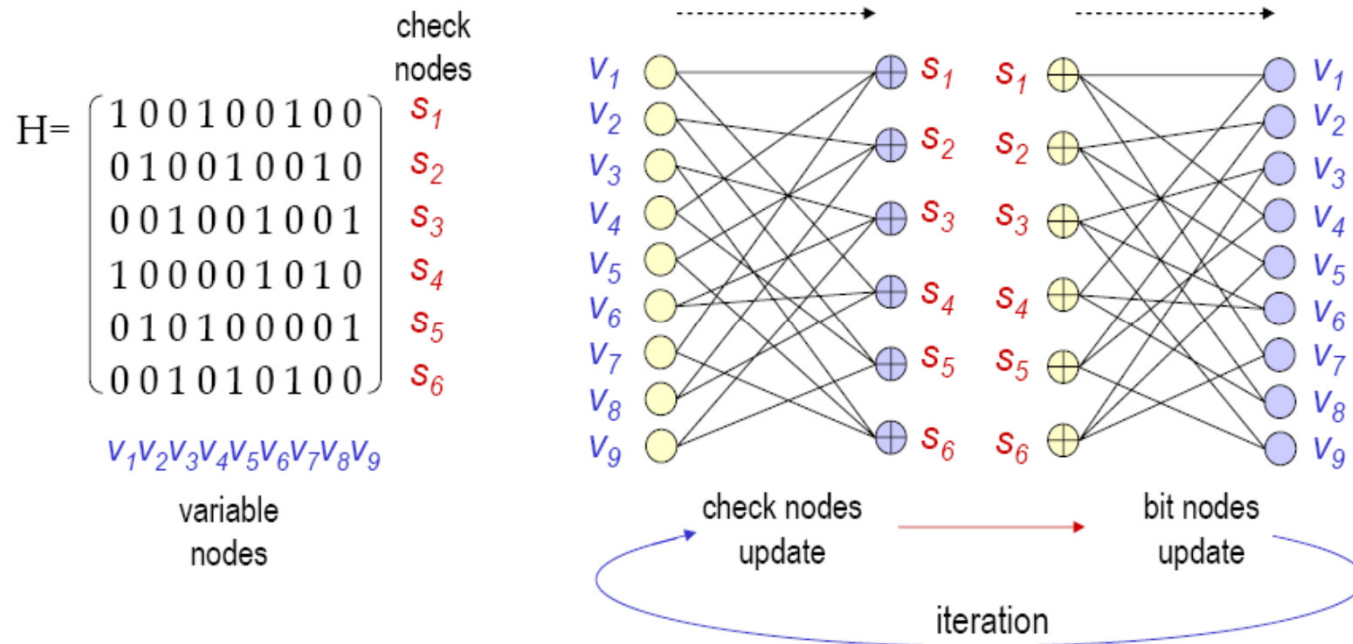


Figure 1.7 Milestones in the drive toward channel capacity achieved by the space systems which evolved over the past 40 years as an answer to the Shannon capacity challenge.

Source :Trellis coding by C. Schlegel, IEEE Press

LDPC Decoding, Quick Recap, 1/5



Variable nodes correspond to the soft information of received bits.

Check nodes describe the parity equations of the transmitted bits.

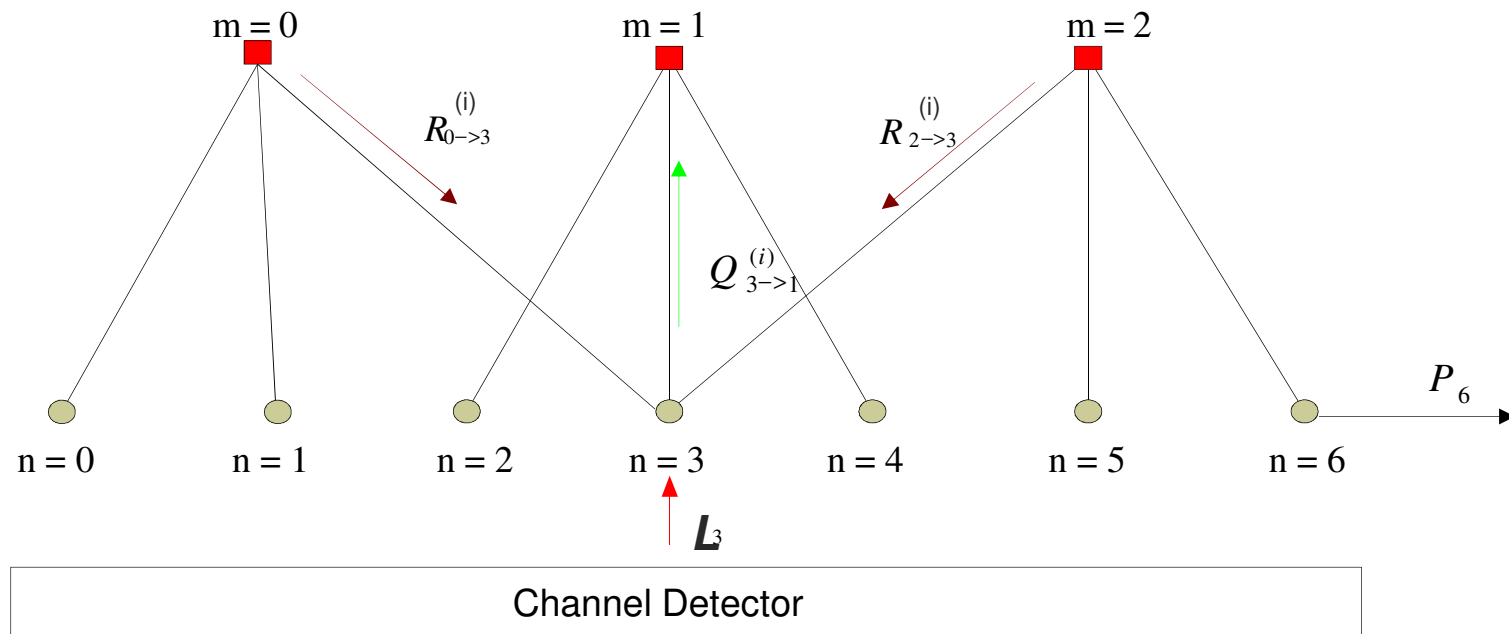
eg. $v_1 + v_4 + v_7 = 0$; $v_2 + v_5 + v_8 = 0$ and so on.

The decoding is successful when all the parity checks are satisfied (i.e. zero).

LDPC Decoding, Quick Recap, 2/5

- There are four types of LLR messages

- Message from the channel to the n -th bit node, L_n
- Message from n -th bit node to the m -th check node $Q_{n \rightarrow m}^{(i)}$ or simply $Q_{nm}^{(i)}$
- Message from the m -th check node to the n -th bit node $R_{m \rightarrow n}^{(i)}$ or simply $R_{mn}^{(i)}$
- Overall reliability information for n -th bit-node P_n



LDPC Decoding, Quick Recap, 3/5

Notation used in the equations

x_n is the transmitted bit n ,

L_n is the initial LLR message for a bit node (also called as variable node) n ,
received from channel/detector

P_n is the overall LLR message for a bit node n ,

\hat{x}_n is the decoded bit n (hard decision based on P_n) ,

[Frequency of P and hard decision update depends on decoding schedule]

$M(n)$ is the set of the neighboring check nodes for variable node n ,

$N(m)$ is the set of the neighboring bit nodes for check node m .

For the i^{th} iteration,

$Q_{nm}^{(i)}$ is the LLR message from bit node n to check node m ,

$R_{mn}^{(i)}$ is the LLR message from check node m to bit node n .

LDPC Decoding, Quick Recap, 4/5

(A) check node processing: for each m and $n \in \mathbf{N}(m)$,

$$R_{mn}^{(i)} = \delta_{mn}^{(i)} \kappa_{mn}^{(i)} \quad (1)$$

$$\kappa_{mn}^{(i)} = |R_{mn}^{(i)}| = \min_{n' \in \mathbf{N}(m) \setminus n} |Q_{n'm}^{(i-1)}| \quad (2)$$

The sign of check node message $R_{mn}^{(i)}$ is defined as

$$\delta_{mn}^{(i)} = \left(\prod_{n' \in \mathbf{N}(m) \setminus n} \text{sgn}(Q_{n'm}^{(i-1)}) \right) \quad (3)$$

where $\delta_{mn}^{(i)}$ takes value of $+1$ or -1

LDPC Decoding, Quick Recap, 5/5

(B) *Variable-node processing*: for each n and $m \in M(n)$:

$$Q_{nm}^{(i)} = L_n + \sum_{m' \in M(n) \setminus m} R_{m'n}^{(i)} \quad (4)$$

(C) P Update and Hard Decision

$$P_n = L_n + \sum_{m \in M(n)} R_{mn}^{(i)} \quad (5)$$

A hard decision is taken where $\hat{x}_n = 0$ if $P_n \geq 0$, and $\hat{x}_n = 1$ if $P_n < 0$.

If $\hat{x}_n H^T = 0$, the decoding process is finished with \hat{x}_n as the decoder output; otherwise, repeat steps (A) to (C).

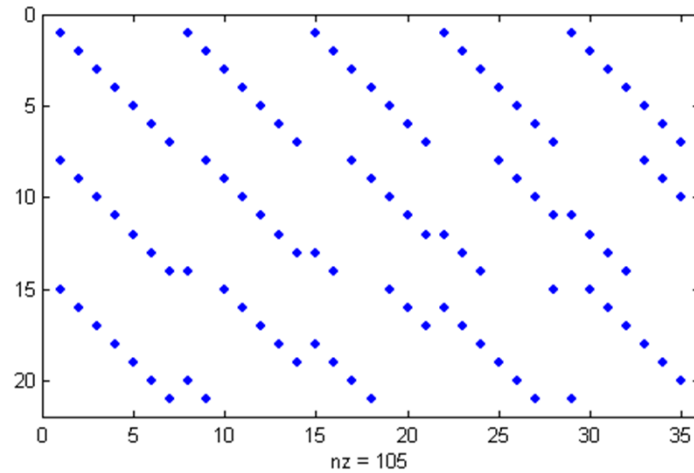
If the decoding process doesn't end within some maximum iteration, stop and output error message.

Scaling or offset can be applied on R messages and/or Q messages for better performance.

Example QC-LDPC Matrix

$$H = \begin{bmatrix} I & I & I & \dots & I \\ I & \sigma & \sigma^2 & \dots & \sigma^{r-1} \\ I & \sigma^2 & \sigma^4 & \dots & \sigma^{2(r-1)} \\ \vdots & & & & \\ I & \sigma^{c-1} & \sigma^{(c-1)2} & \dots & \sigma^{(c-1)(r-1)} \end{bmatrix}$$

$$\sigma = \begin{bmatrix} 0 & 0 & \dots & 0 & 1 \\ 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \vdots & & & & \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix}$$



Example H Matrix, Array
LDPC

r row/ check node degree=5

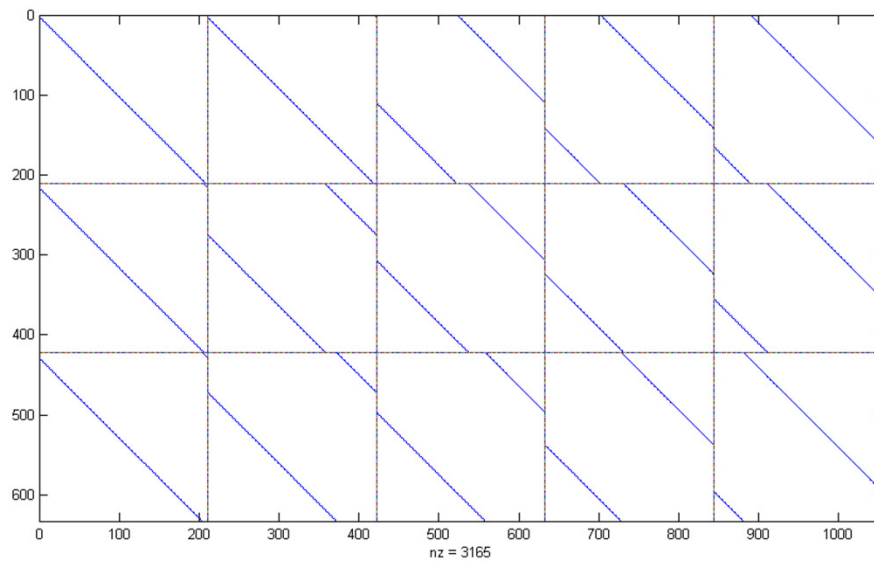
c columns/variable node degree=3

Sc circulant size=7

$N=Sc*r=35$

Example QC-LDPC Matrix

$$S_{3 \times 5} = \begin{bmatrix} 2 & 3 & 110 & 142 & 165 \\ 5 & 64 & 96 & 113 & 144 \\ 7 & 50 & 75 & 116 & 174 \end{bmatrix}$$



Example H Matrix

r row degree=5

c column degree =3

Sc circulant size=211

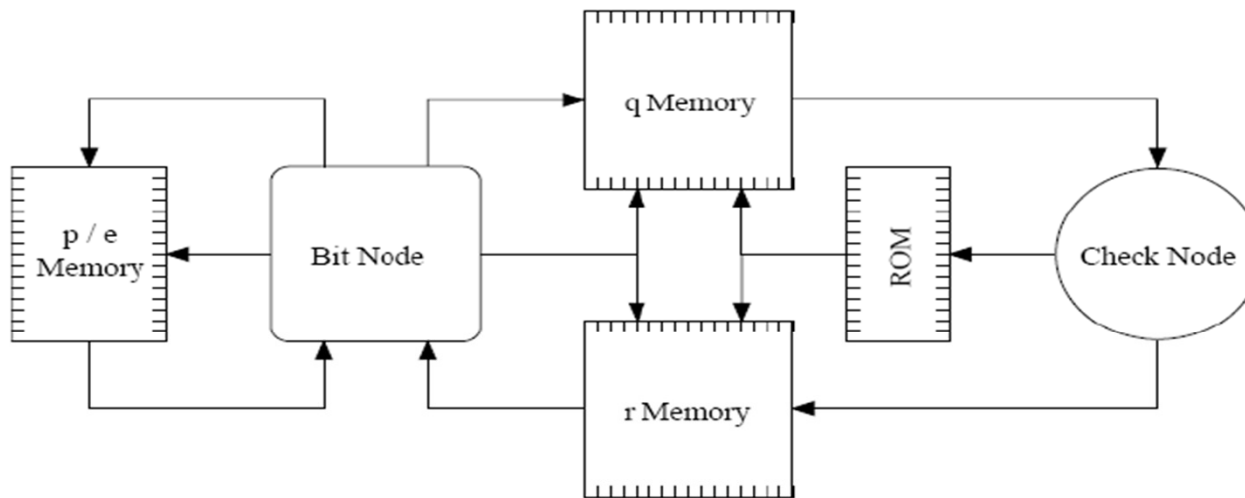
$N=Sc*r=1055$

Decoder Architectures

- Parallelization is good-but comes at a steep cost for LDPC.
- Fully Parallel Architecture:
- All the check updates in one clock cycle and all the bit updates in one more clock cycle.
- Huge Hardware resources and routing congestion.

Decoder Architectures, Serial

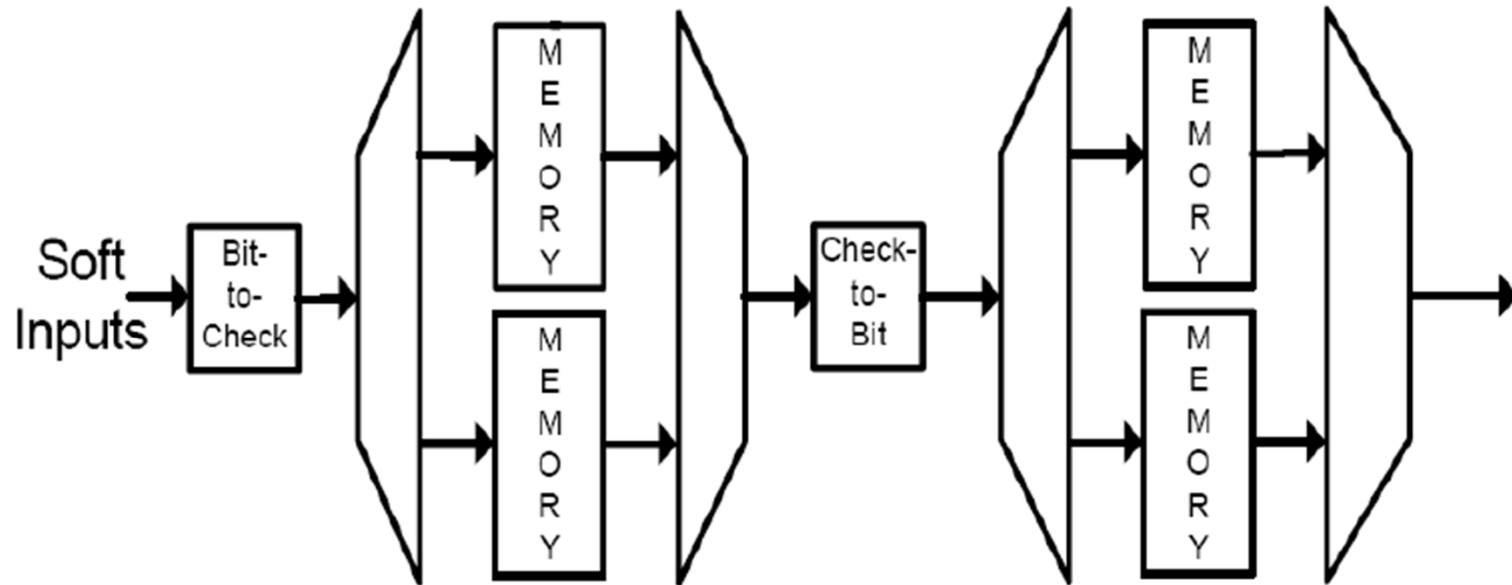
- Check updates and bit updates in a serial fashion.
- Huge Memory requirement. Memory in critical path.



Serial Architecture

[1] Levine, .etal Implementation of near Shannon limit error-correcting codes using reconfigurable hardware
IEEE Field-Programmable Custom Computing Machines, 2000

Decoder Architectures, Serial



Serialized and fully pipelined implementation requires two memory buffers per stage, alternating between read/write.

Serial Architecture.

[2] E. Yeo, "VLSI architectures for iterative decoders in magnetic recording channels," *IEEE Trans. Magnetics*, vol.37, no.2, pp. 748-55, March 2001.

Semi-parallel Architectures

- Check updates and bit updates using several units.
- Partitioned memory by imposing structure on H matrix.
- There are several semi-parallel architectures proposed.
- Initial semi-parallel architectures were complex and very low throughput. Needed huge amounts of memory.

On-the-fly Computation

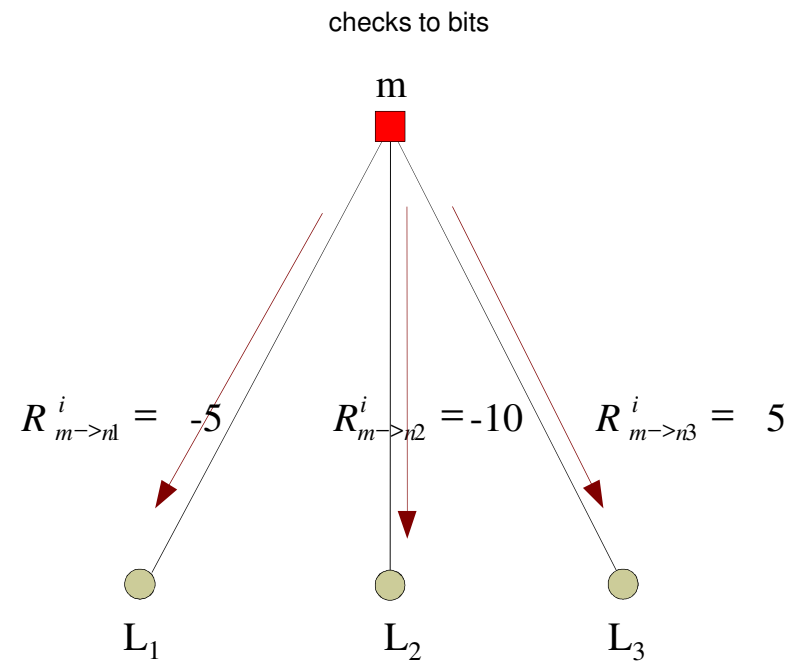
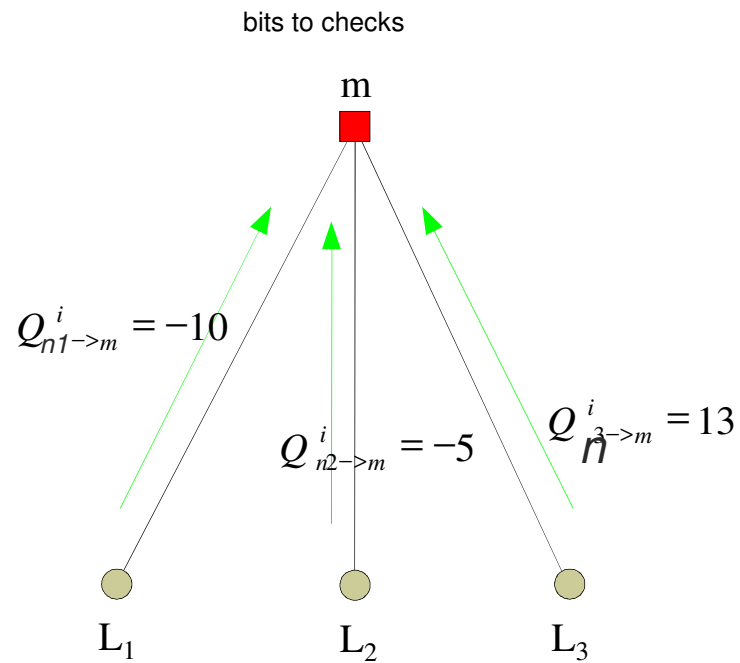
Our previous research ([1-13]) introduced the following concepts to LDPC decoder implementation. References P1-P4 are more comprehensive and are the basis for this presentation.

1. Block serial scheduling
2. Value-reuse,
3. Scheduling of layered processing,
4. Out-of-order block processing,
5. Master-slave router,
6. Dynamic state,
7. Speculative Computation
8. Run-time Application Compiler [support for different LDPC codes with in a class of codes. Class:802.11n,802.16e,Array, etc. Off-line re-configurable for several regular and irregular LDPC codes]

All these concepts are termed as On-the-fly computation as the core of these concepts are based on minimizing memory and re-computations by employing just in-time scheduling. For this presentation, we will focus on concept 4.

Check Node Unit (CNU) Design

$$\kappa_{m \rightarrow l}^{(i)} = |R_{m \rightarrow l}^{(i)}| = \min_{l' \in N(m) \setminus l} |Q_{l' \rightarrow m}^{(i-1)}|$$



Check Node Unit (CNU) Design

$$\kappa_{mn}^{(i)} = |R_{mn}^{(i)}| = \min_{n' \in N(m) \setminus n} |Q_{n'm}^{(i-1)}| \quad (2)$$

The above equation (2) can be reformulated as the following set of equations.

$$M1_m^{(i)} = \min_{n' \in N(m)} |Q_{mn'}^{(i-1)}| \quad (6)$$

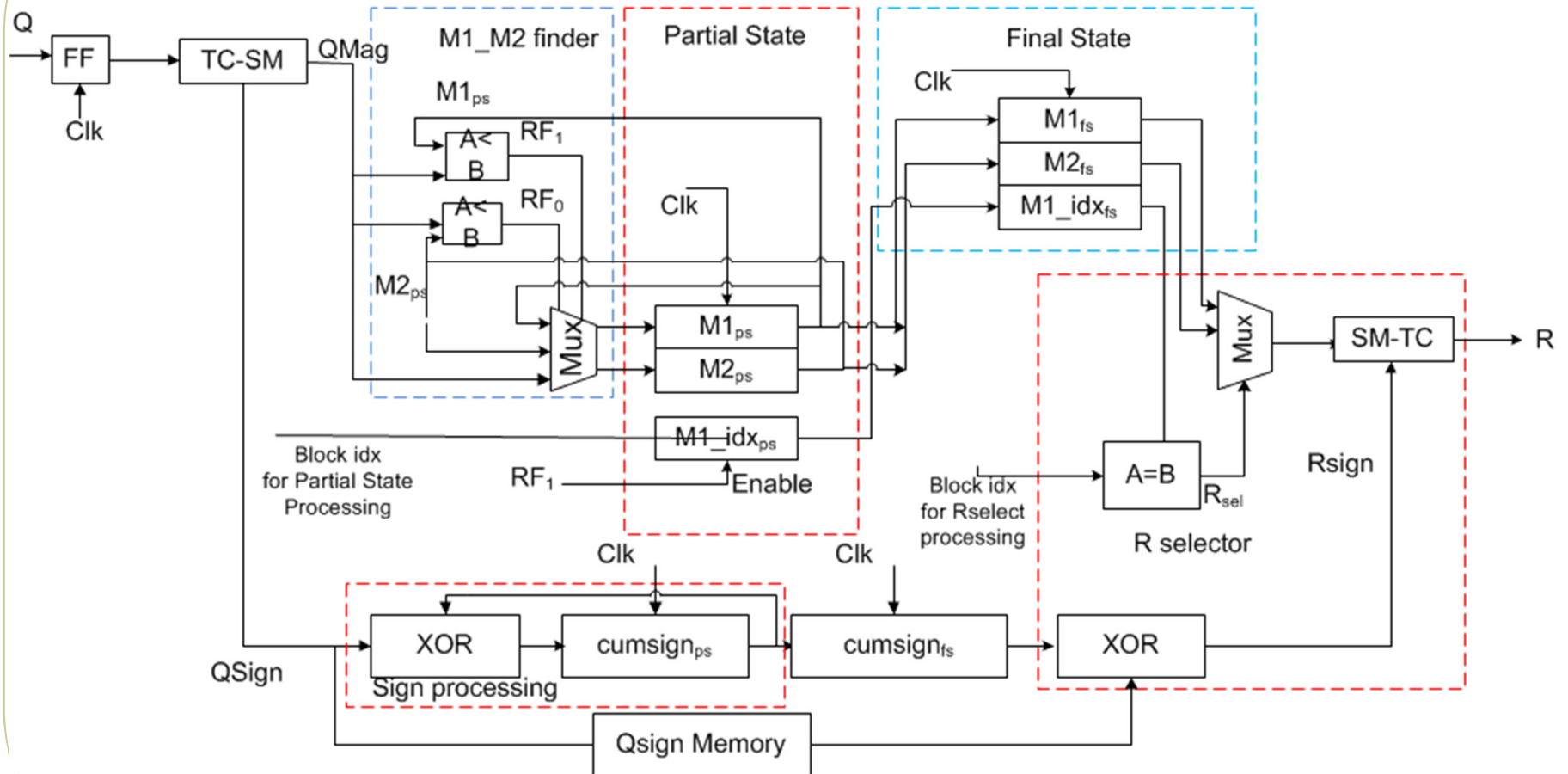
$$M2_m^{(i)} = 2nd \min_{n' \in N(m)} |Q_{mn'}^{(i-1)}| \quad (7)$$

$$k = Min1Index \quad (8)$$

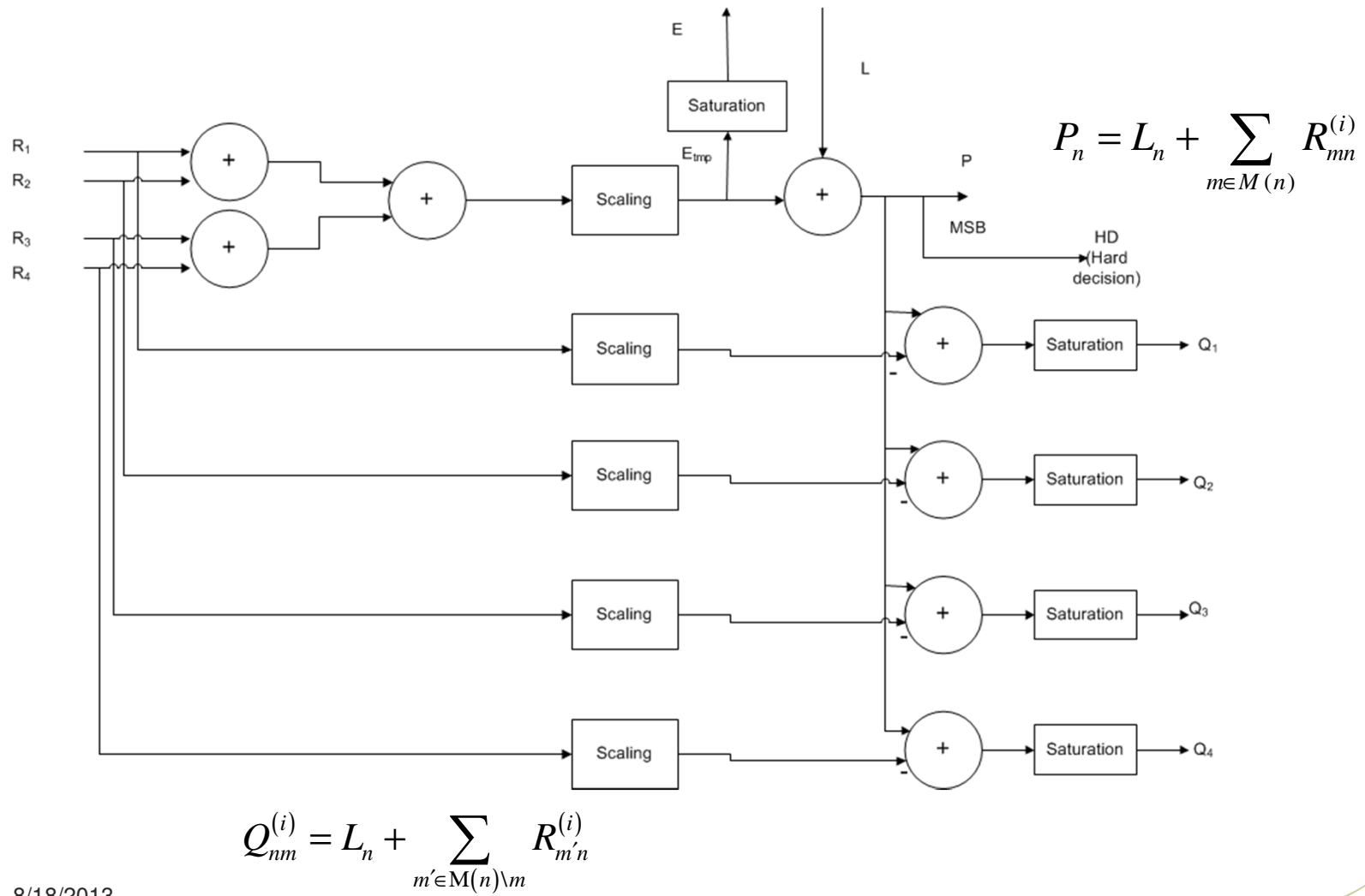
$$\begin{aligned} \kappa_{mn}^{(i)} &= M1_m^{(i)}, \forall n \in N(m) \setminus k \\ &= M2_m^{(i)}, n = k \end{aligned} \quad (9)$$

- Simplifies the number of comparisons required as well as the memory needed to store CNU outputs.
- Additional design choices: The correction has to be applied to only two values instead of distinct values. Need to apply 2's complement to only 1 or 2 values instead of values at the output of CNU.

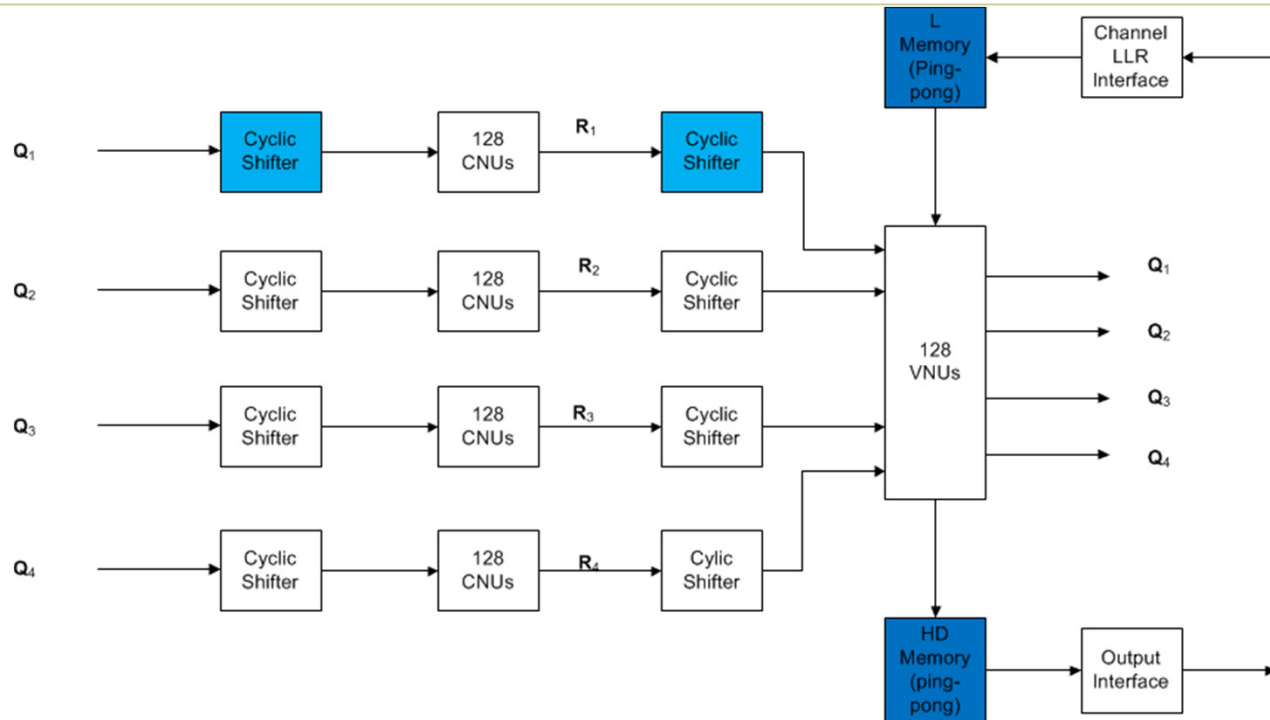
CNU Micro Architecture for min-sum



VNU Micro Architecture



Non-Layered Decoder Architecture



Supported H matrix parameters

r row degree=36

c column degree =4

Sc circulant size=128

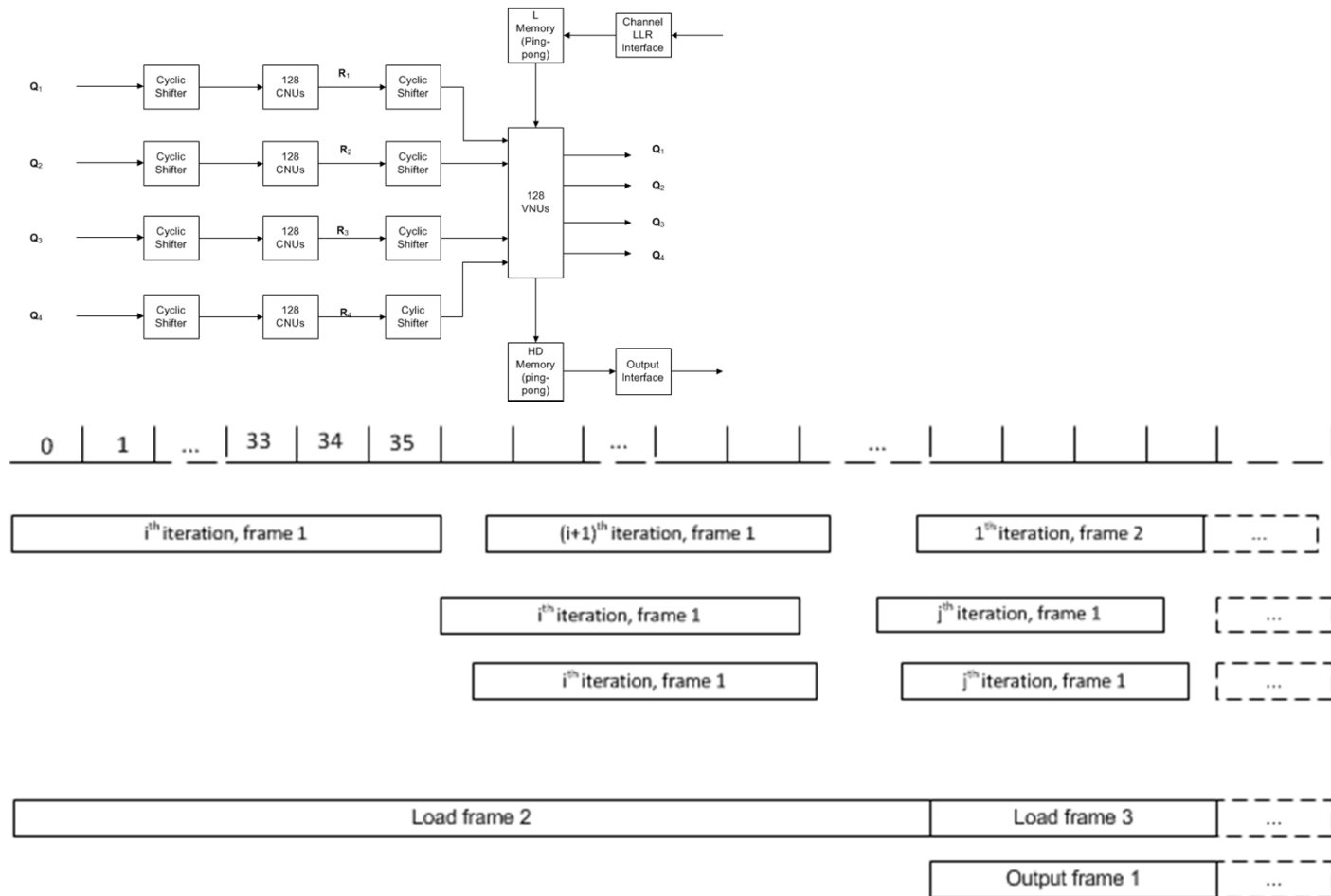
$N=Sc*r=4608$

L Memory=> Depth 36, Width=128*5

HD Memory=> Depth 36, Width=128

Possible to remove the shifters (light-blue) by re-arranging H matrix's first layer to have zero shift coefficients.

Pipeline for Non-layered Decoder



Layered Decoder Architecture

Optimized Layered Decoding with algorithm transformations for reduced memory and computations

$$\vec{R}_{l,n}^{(0)} = 0, \vec{P}_n = \vec{L}_n^{(0)} \quad [\text{Initialization for each new received data frame}], \quad (9)$$

$$\forall i = 1, 2, \dots, it_{\max} \quad [\text{Iteration loop}],$$

$$\forall l = 1, 2, \dots, j \quad [\text{Sub-iteration loop}],$$

$$\forall n = 1, 2, \dots, k \quad [\text{Block column loop}],$$

$$[\vec{Q}_{l,n}^{(i)}]^{s(l,n)} = [\vec{P}_n]^{s(l,n)} - \vec{R}_{l,n}^{(i-1)}, \quad (10)$$

$$\vec{R}_{l,n}^{(i)} = f\left([\vec{Q}_{l,n'}^{(i)}]^{s(l,n')}, \forall n' = 1, 2, \dots, k\right), \quad (11)$$

$$[\vec{P}_n]^{s(l,n)} = [\vec{Q}_{l,n}^{(i)}]^{s(l,n)} + \vec{R}_{l,n}^{(i)}, \quad (12)$$

where the vectors $\vec{R}_{l,n}^{(i)}$ and $\vec{Q}_{l,n}^{(i)}$ represent all the R and Q messages in each $p \times p$ block of the H matrix, $s(l,n)$ denotes the shift coefficient for the block in l^{th} block row and n^{th} block column of the H matrix.

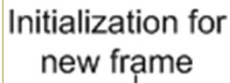
$[\vec{Q}_{l,n}^{(i)}]^{s(l,n)}$ denotes that the vector $\vec{Q}_{l,n}^{(i)}$ is cyclically shifted up by the amount $s(l,n)$

k is the check-node degree of the block row.

A negative sign on $s(l,n)$ indicates that it is a cyclic down shift (equivalent cyclic left shift).

$f(\cdot)$ denotes the check-node processing, which embodiments implement using, for example, a Bahl-Cocke-Jelinek-Raviv algorithm ("BCJR") or sum-of-products ("SP") or Min-Sum with scaling/offset.

11/11/2016



Compared to other work, this work has several advantages

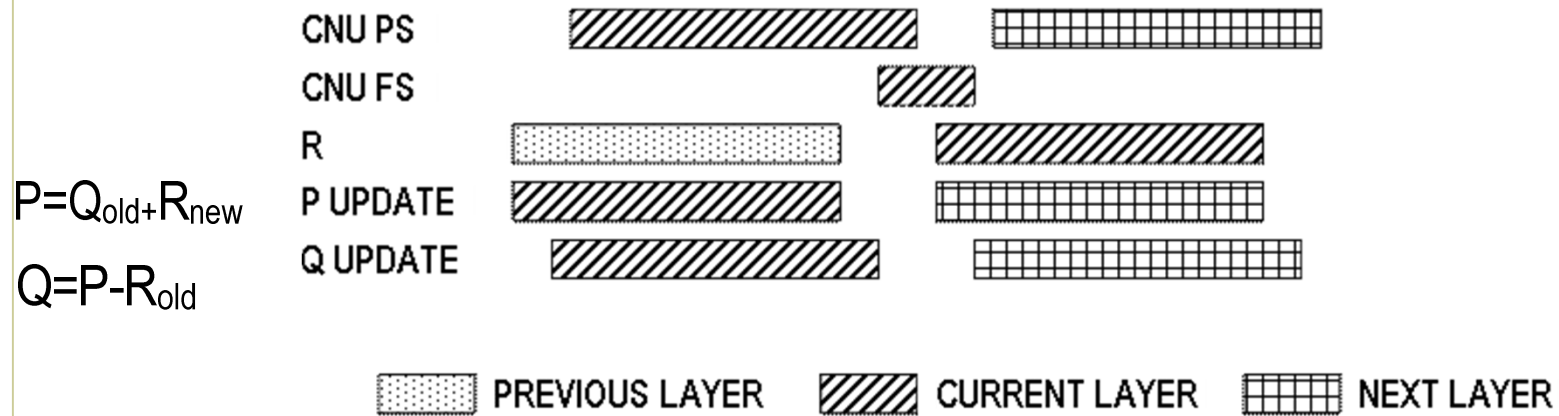
- 1) No need of separate memory for P.
- 2) Only one shifter instead of 2 shifters
- 3) Value-reuse is effectively used for both Rnew and Rold
- 4) Low complexity data path design-with no redundant data Path operations.
- 5) Low complexity CNU design.

$$\left[\vec{Q}_{l,n}^{(i)}\right]^{S(l,n)} = \left[\vec{P}_n\right]^{S(l,n)} - \vec{R}_{l,n}^{(i-1)}$$

$$Q = P - R_{old}$$

$$P = Q_{old} + R_{new}$$

Data Flow Diagram



Irregular QC-LDPC H Matrices

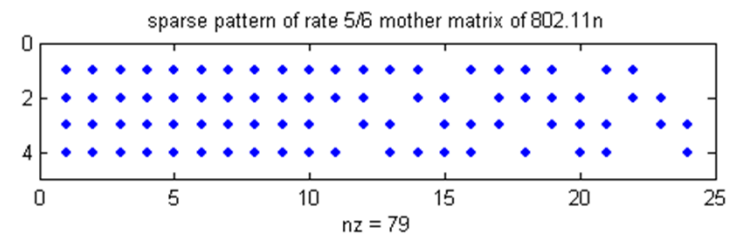
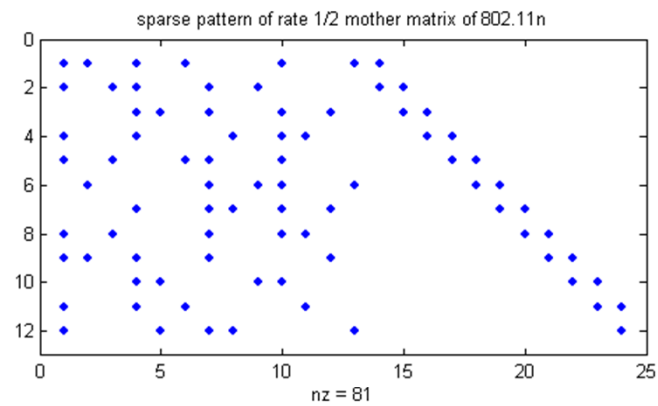
$$\mathbf{H} = \begin{bmatrix} \mathbf{P}_{0,0} & \mathbf{P}_{0,1} & \mathbf{P}_{0,2} & \cdots & \mathbf{P}_{0,n_b-2} & \mathbf{P}_{0,n_b-1} \\ \mathbf{P}_{1,0} & \mathbf{P}_{1,1} & \mathbf{P}_{1,2} & \cdots & \mathbf{P}_{1,n_b-2} & \mathbf{P}_{1,n_b-1} \\ \mathbf{P}_{2,0} & \mathbf{P}_{2,1} & \mathbf{P}_{2,2} & \cdots & \mathbf{P}_{2,n_b-2} & \mathbf{P}_{0,n_b-1} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \mathbf{P}_{m_b-1,0} & \mathbf{P}_{m_b-1,1} & \mathbf{P}_{m_b-1,2} & \cdots & \mathbf{P}_{m_b-1,n_b-2} & \mathbf{P}_{m_b-1,n_b-1} \end{bmatrix} = \mathbf{P}^{H_b}$$

Different base matrices to support different rates.

Different expansion factors (z) to support multiple lengths.

All the shift coefficients for different codes for a given rate are obtained from the same base matrix using modulo arithmetic

Irregular QC-LDPC H Matrices

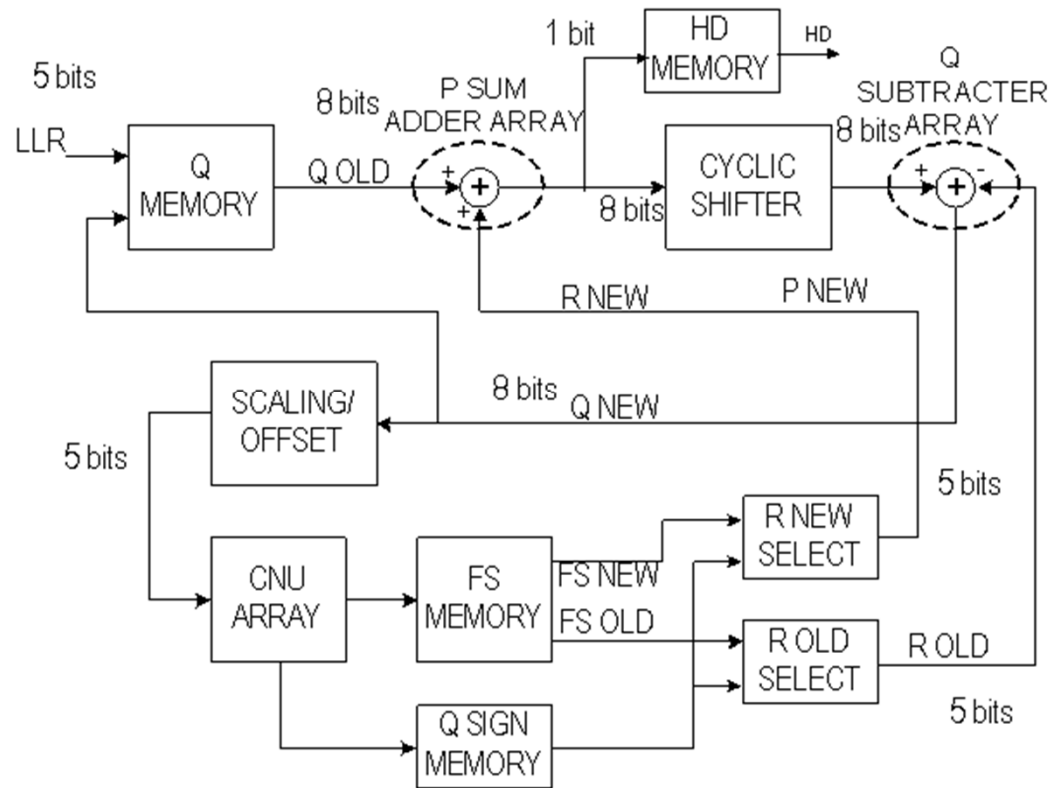


Irregular QC-LDPC H Matrices

- ❑ Existing implementations [*Hocevar*] for irregular QC-LDPC codes are still very complex to implement.
- ❑ These codes have the better BER performance and selected for IEEE 802.16e and IEEE 802.11n.
- ❑ It is anticipated that these type of codes will be the default choice for most of the standards.
- ❑ We show that with out-of-order processing and scheduling of layered processing, it is possible to design very efficient architectures.
- ❑ The same type of codes can be used in storage applications (holographic, flash and magnetic recording) if variable node degrees of 2 and 3 are avoided in the code construction for low error floor

Hocevar, D.E., "A reduced complexity decoder architecture via layered decoding of LDPC codes,"
IEEE Workshop on Signal Processing Systems, 2004. SIPS 2004. .pp. 107- 112, 13-15 Oct. 2004

Layered Decoder Architecture



Advantages

- 1) Q memory (some times we call this as LPQ memory) can be used to store L/Q/P instead of 3 separate memories- memory is managed at circulant level as at any time for a given circulant we need only L or Q or P.
 - 2) Only one shifter instead of 2 shifters
 - 3) Value-reuse is effectively used for both Rnew and Rold
 - 4) Low complexity data path design-with no redundant data
- Path operations.
- 5) Low complexity CNU design.
 - 6) Out-of-order processing at both layer and circulant level for all the processing steps such as Rnew and PS processing to eliminate the pipeline and memory access stall cycles.

Out-of-order layer processing for R Selection

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
0	1	2			3	4		5	6		7	8					9	10						
1			11 ¹		12 ⁶			13 ⁷	14 ⁸			15 ⁹	16 ²		17 ³	18 ⁴		19 ¹⁰	20 ⁵					
2			21	22		23		24		25		26		27	28				29	30				
3			31	32		33	34		35		36				37	38				39	40			
4	41		42			43	44			45		46		47			48				49	50		
5			51		52	53			54		55		56		57	58						59	60	
6	61	62		63		64			65			66		67	68								69	70
7		71	72				73		74	75		76	77		78		79							80

Normal practice is to compute R new messages for each layer after CNU PS processing.

However, here we decoupled the execution of R new messages of each layer with the execution of corresponding layer's CNU PS processing. Rather than simply generating Rnew messages per layer, we compute them on basis of circulant dependencies.

R selection is out-of-order so that it can feed the data required for the PS processing of the second layer. For instance Rnew messages for circulant 29 which belong to layer 3 are not generated immediately after layer 3 CNU PS processing .

Rather, Rnew for circulant 29 is computed when PS processing of circulant 20 is done as circulant 29 is a dependent circulant of circulant of 20.

Similarly, Rnew for circulant 72 is computed when PS processing of circulant 11 is done as circulant 72 is a dependent circulant of circulant of 11.

Here we execute the instruction/computation at precise moment when the result is needed!

Out-of-order block processing for Partial State

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
0	1	2			3	4		5	6		7	8					9	10						
1			11 ¹		12 ⁶			13 ⁷	14 ⁸			15 ⁹	16 ²		17 ³	18 ⁴		19 ¹⁰	20 ⁵					
2			21	22		23		24		25		26		27	28				29	30				
3			31	32		33	34		35		36				37	38				39	40			
4	41		42			43	44			45		46		47			48				49	50		
5			51		52	53			54		55		56		57	58						59	60	
6	61	62		63		64			65			66		67	68								69	70
7		71	72				73		74	75		76	77		78		79							80

Re-ordering of block processing . While processing the layer 2,

the blocks which depend on layer 1 will be processed last to allow for the pipeline latency.

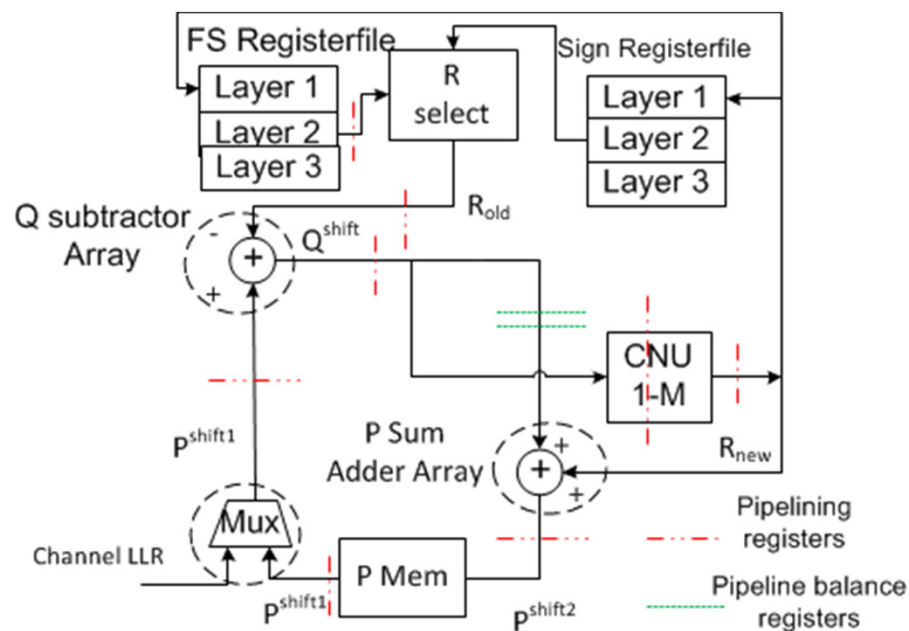
In the above example, the pipeline latency can be 5.

The vector pipeline depth is 5.so no stall cycles are needed while processing the layer 2 due to the pipelining. [In other implementations, the stall cycles are introduced – which will effectively reduce the throughput by a huge margin.]

Also we will sequence the operations in layer such that we process the block first that has dependent data available for the longest time.

This naturally leads us to true out-of-order processing across several layers. In practice we wont do out-of-order partial state processing involving more than 2 layers.

Block Parallel Layered Decoder



Compared to other work, this work has several advantages

- 1) Only one memory for holding the P values.
- 2) Shifting is achieved through memory reads. Only one memory multiplexer network is needed instead of 2 to achieve delta shifts
- 3) Value-reuse is effectively used for both R_{new} and R_{old}
- 4) Low complexity data path design-with no redundant data Path operations.
- 5) Low complexity CNU design with high parallelism.
- 6) Smaller pipeline depth
- 7) Out-of-order row processing to hide the pipeline latencies.

Here M is the row parallelization (i.e. number of rows in H matrix Processed per clock).

Cyclic Shifter

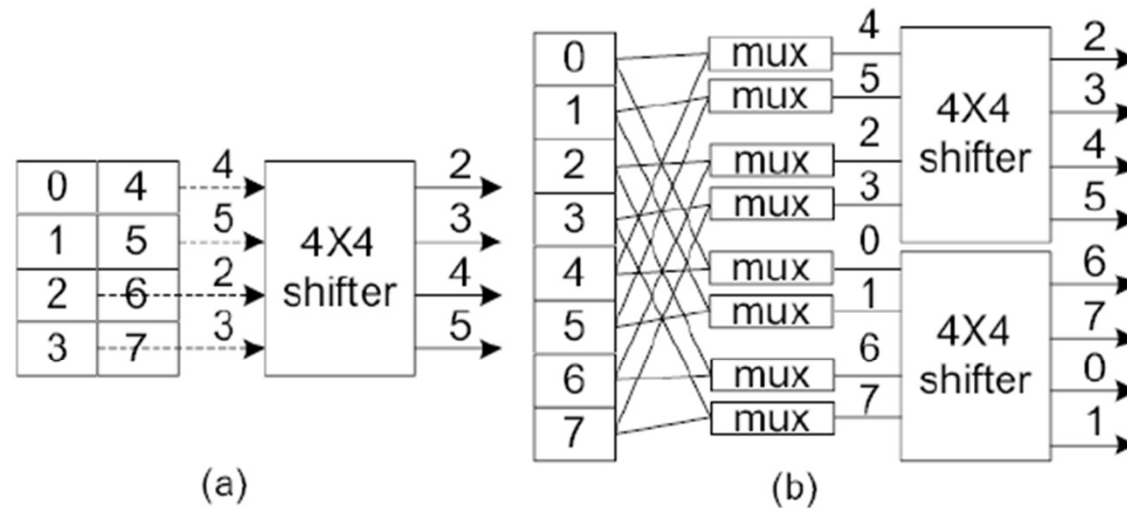


Fig. 5. (a) 8×8 shifting with 4×4 cyclic shifter. (b) 8×8 shifting with two 4×4 cyclic shifters.

This arrangement can support the base matrices having the expansion factors multiples of z by using $z \times z$ cyclic shifters.

Works for 802.11n in which the expansion factors are 27,54,81

Works for limited configurations of 802.16e in which the expansion factors are 24,28,...,96.[24 x 24 shifter works for 24,48 and 96]

Benes Network

Omega network is blocking network which can do cyclic permutations and other limited permutations.

Benes network is a non-blocking network which can do any permutation on the input vector. It is essentially combination of two omega n/w.

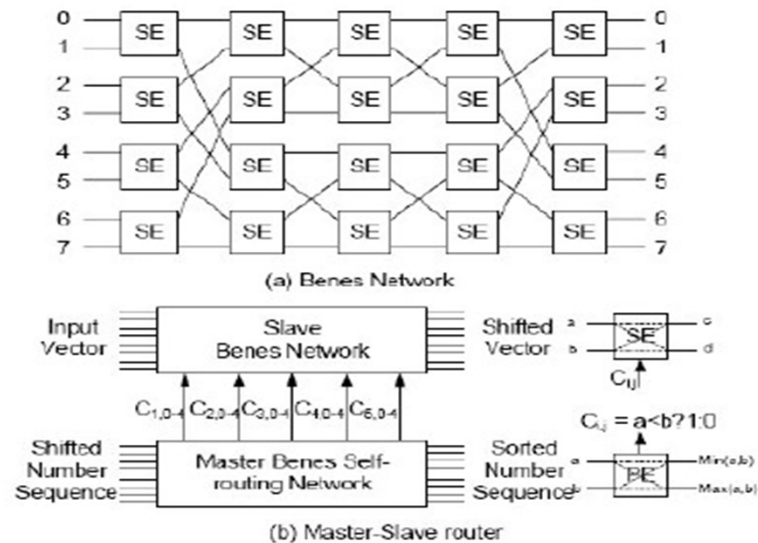
However the control complexity is more for both the networks.

In the existing work, the control signals are pre-computed and stored assuming only one H matrix need to be supported.

[3] Mansour et al "A 640-Mb/s 2048-Bit Programmable LDPC Decoder Chip"-IEEE Journal of Solid-State Circuits, March 2006

[12] Malema, G.; Liebelt, M., "Interconnection Network for Structured Low-Density Parity-Check Decoders," Communications, 2005 Asia-Pacific Conference on , vol., no.pp. 537- 540, 03-05 Oct. 2005

Proposed Master-slave Router



In 802.16e in which the expansion factors are 24,28,...,96.

Assume that we have a monolithic 96 x 96 cyclic shifter. It can not do cyclic shift on vector size less than 96.

To accommodate all different vector sizes, we need to use a generic router such as Bene's n/w.

Control signals can be generated by using another self routing n/w

Gunnam, KK; Choi, G. S.; Yeary, M. B.; Atiquzzaman, M.; "VLSI Architectures for Layered Decoding for Irregular LDPC Codes of WiMax," Communications, 2007. ICC '07. IEEE International Conference on 24-28 June 2007 Page(s):4542 - 4547

Master-slave router

Note that this memory for providing control signals to this network is equal to $\frac{M}{2}(2^{\log_2(M)} - 1)$ bits for every shift value that needs to be supported.

This will be a very huge requirement for supporting all the WiMax codes. Note that, the memory needed for storing control signals for Omega network is around 1.22 mm² in out of the decoder chip area of 14.1 mm². [3].

To support 19 different expansion factors and 6 types of base H matrices *in run time*, the control signal memory needs approximately 139.08 mm².

Assume that we need to perform a cyclic shift of 2 on a message vector of length 4 using a 8 x8 Slave Benes network.

Supply the integers(2,3,0,1,4,5,6,7) to the Master Benes network which is always configured to sort the inputs and output (0,1,2,...7).

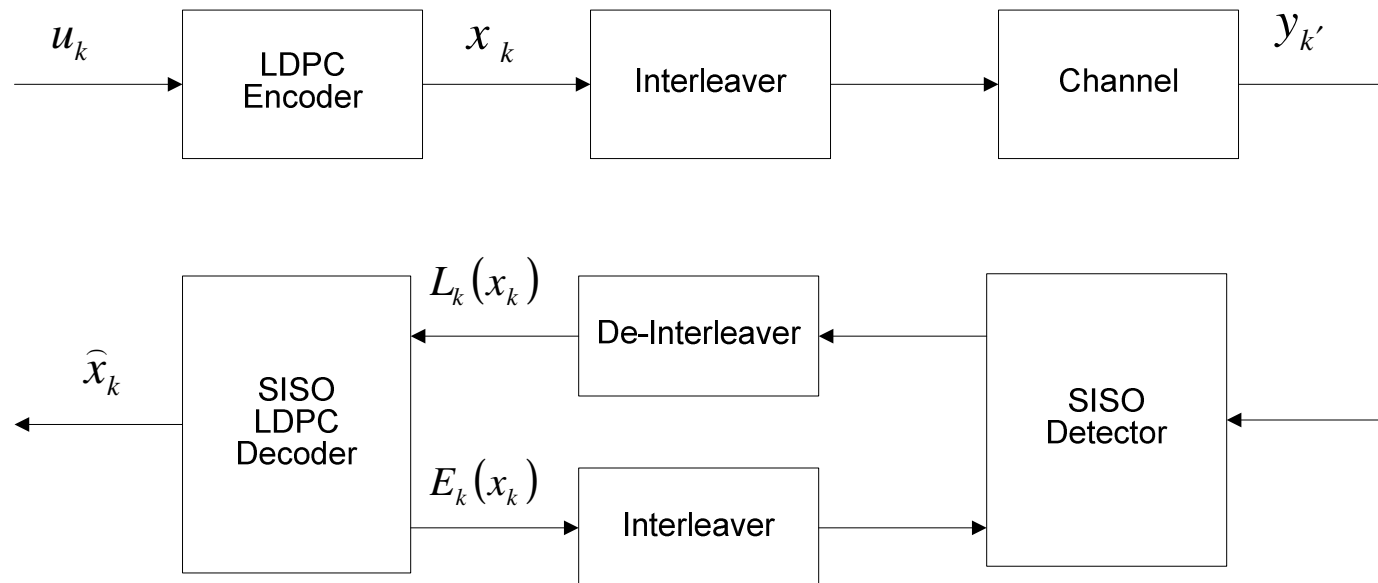
Though the sorting is based on [13], the new proposal is the conception of master-slave router.

[13] Gocal, B. 1996. *Bitonic Sorting on Bene Networks*. In *Proceedings of the 10th international Parallel Processing Symposium (April 15 - 19, 1996)*. IPPS. IEEE Computer Society, Washington, DC, 749-753.

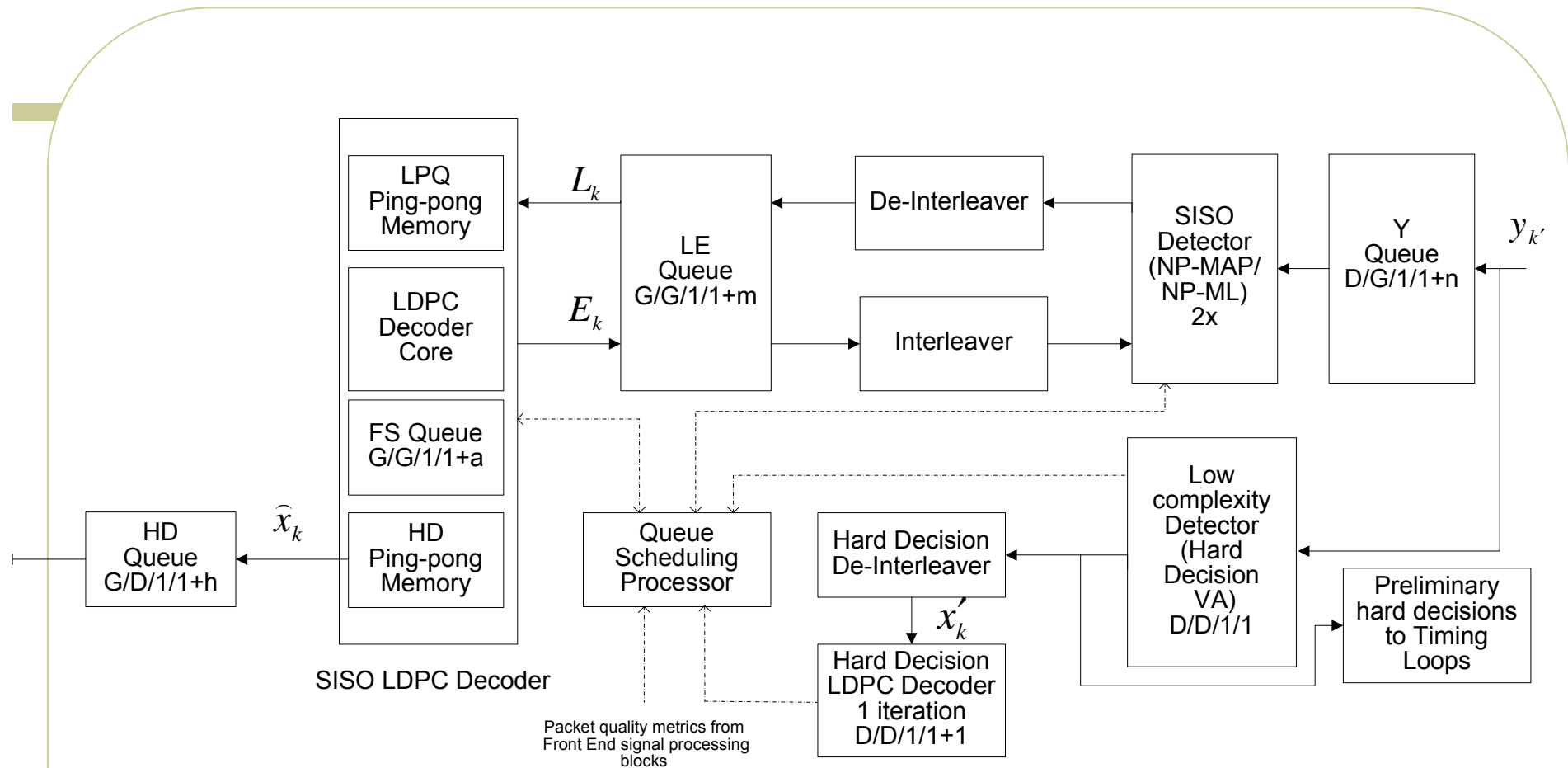
TURBO EQUALIZATION

8/18/2013

System Model for Turbo Equalization



Proposed System Level Architecture for Turbo Equalization



Gunnam et. al, "Next generation iterative LDPC solutions for magnetic recording storage," Signals, Systems and Computers, 2008 42nd Asilomar Conference on, Publication Year: 2008 , Page(s): 1148 – 1152

Local-Global Interleaver

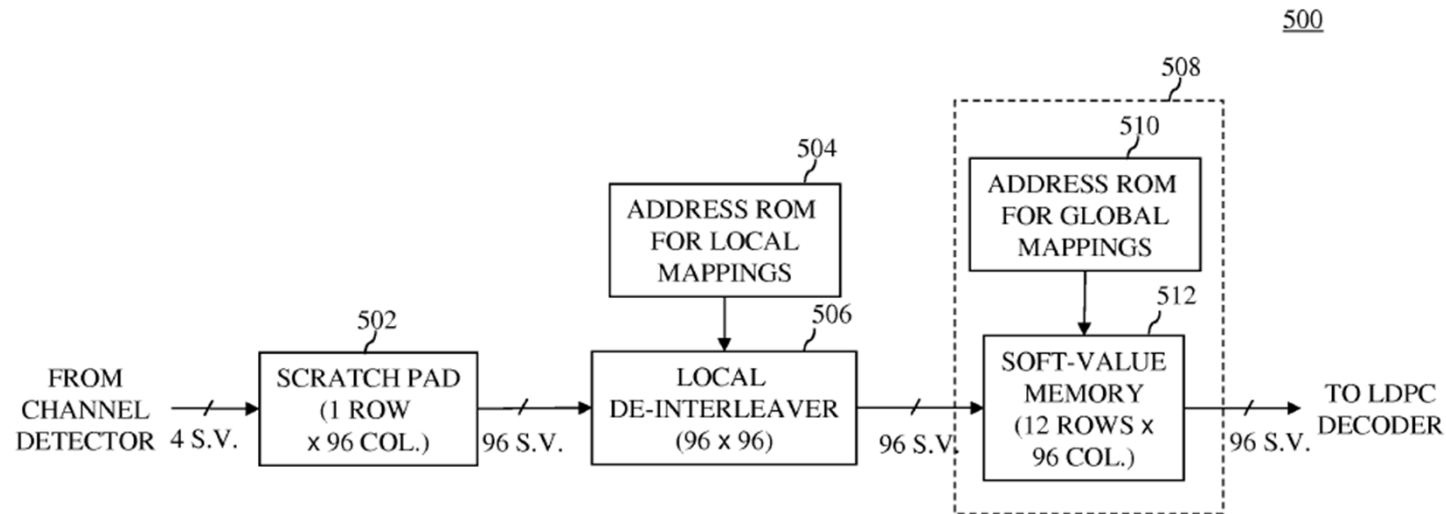


FIG. 5

Row-Column interleavers need to have memory organized such that it can supply the data samples for both row and column access.

Low latency memory efficient interleaver compared to traditional row-column interleaver. Only one type of access (i.e. row access) is needed for both detector and decoder.

Data flow in Local-Global Interleaver

600

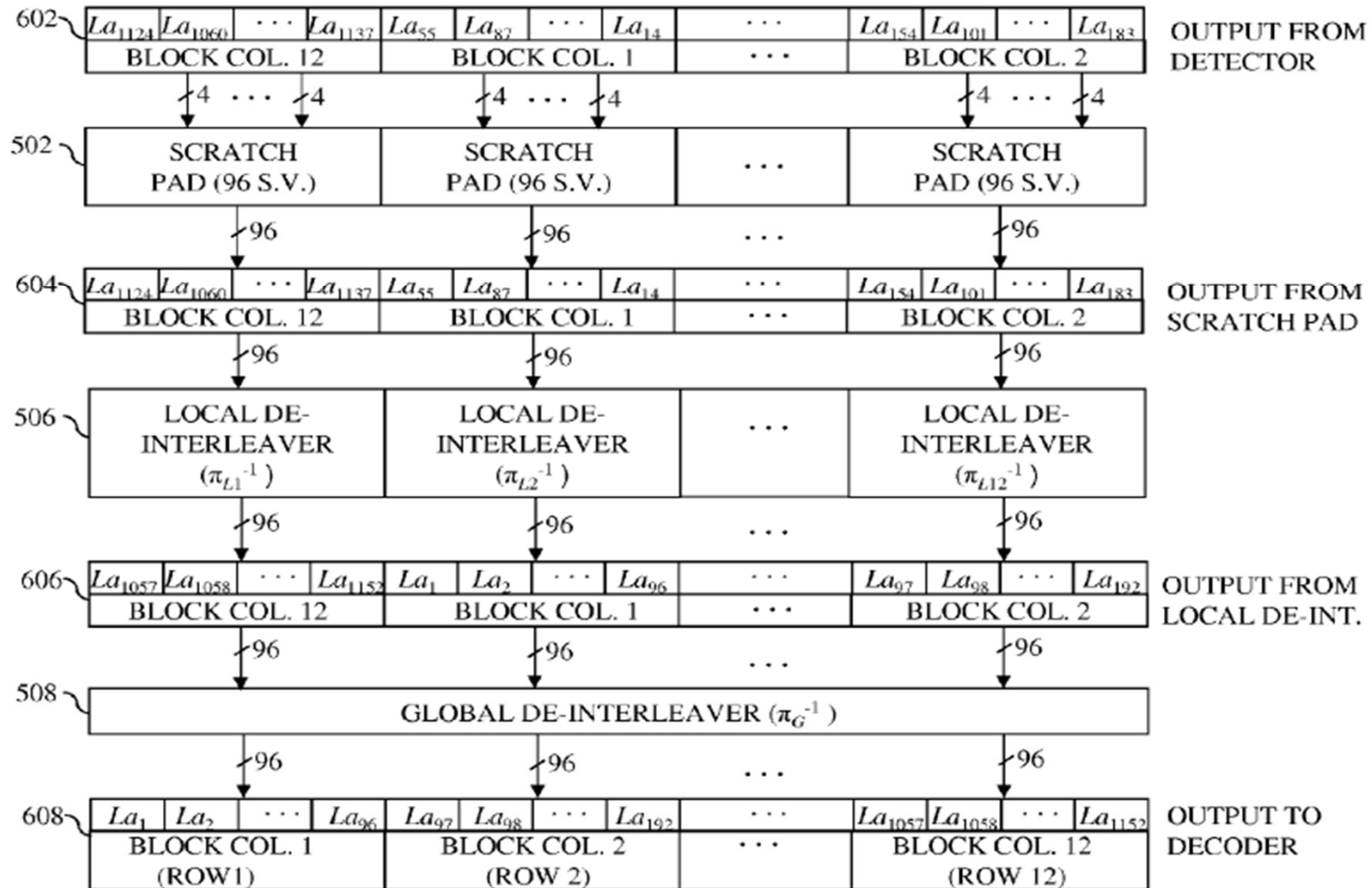
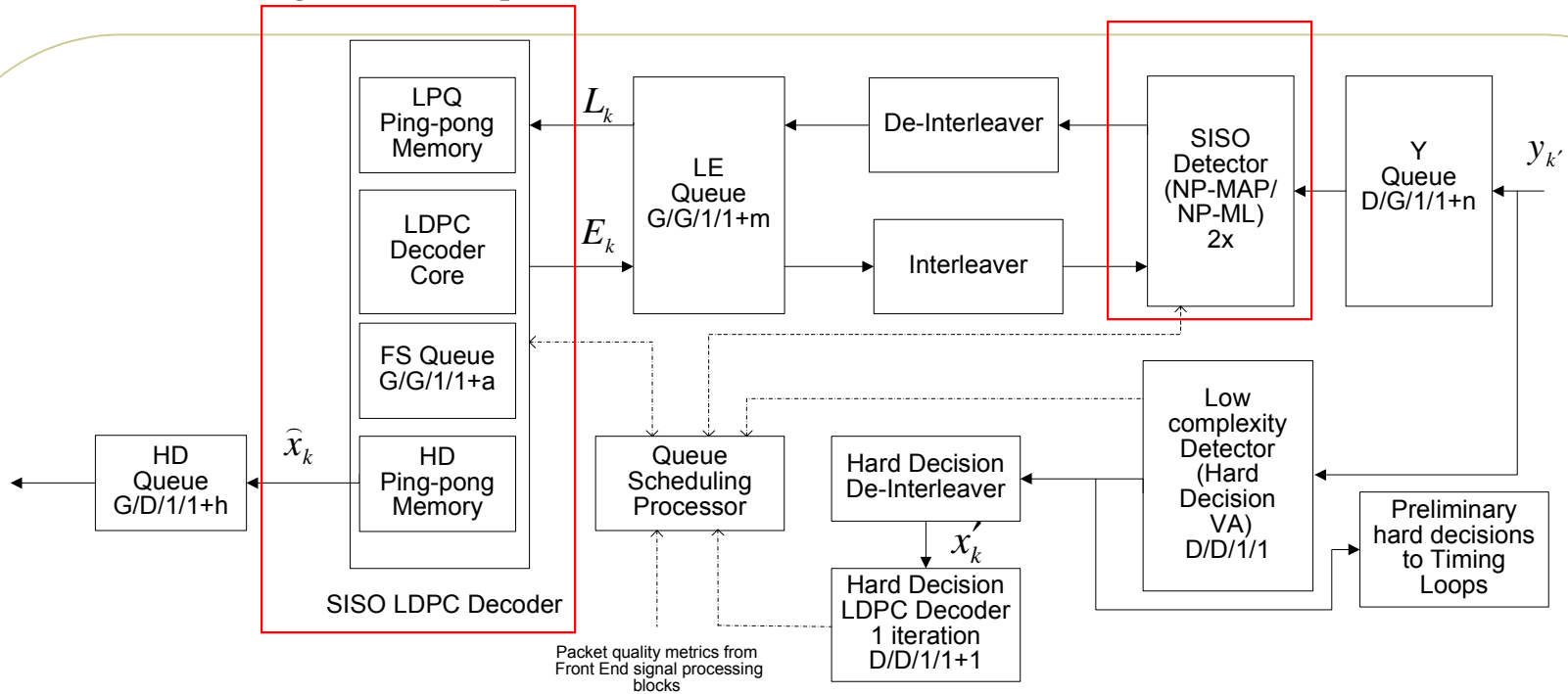


FIG. 6

Why Statistical Buffering?

- The innovation here is the novel and efficient arrangement of queue structures such that we would get the performance of a hardware system that is configured to run h (which is set to 20 in the example configuration) maximum global iterations while the system complexity is proportional to the hardware system that can 2 maximum global iterations.
- $D/G/1/1+n$ is Kendall's notation of a queuing model. The first part represents the input process, the second the service distribution, and the third the number of servers.
D- Deterministic
G- General

Primary Data path



- The primary data-path contains one SISO LDPC decoder and one SISO detector.
- The LDPC decoder is designed such that it can handle the total amount of average iterations in two global iterations for each packet.
- The SISO detector is in fact two detector modules that operate on the same packet but different halves of the packet thus ensuring one packet can be processed in 50% of the inter-arrival time T . Each detector processes 4 samples per clock cycle.
- Thus both the detector and the LDPC decoder can sustain maximum of two global iterations per each packet if no statistical buffering is employed.

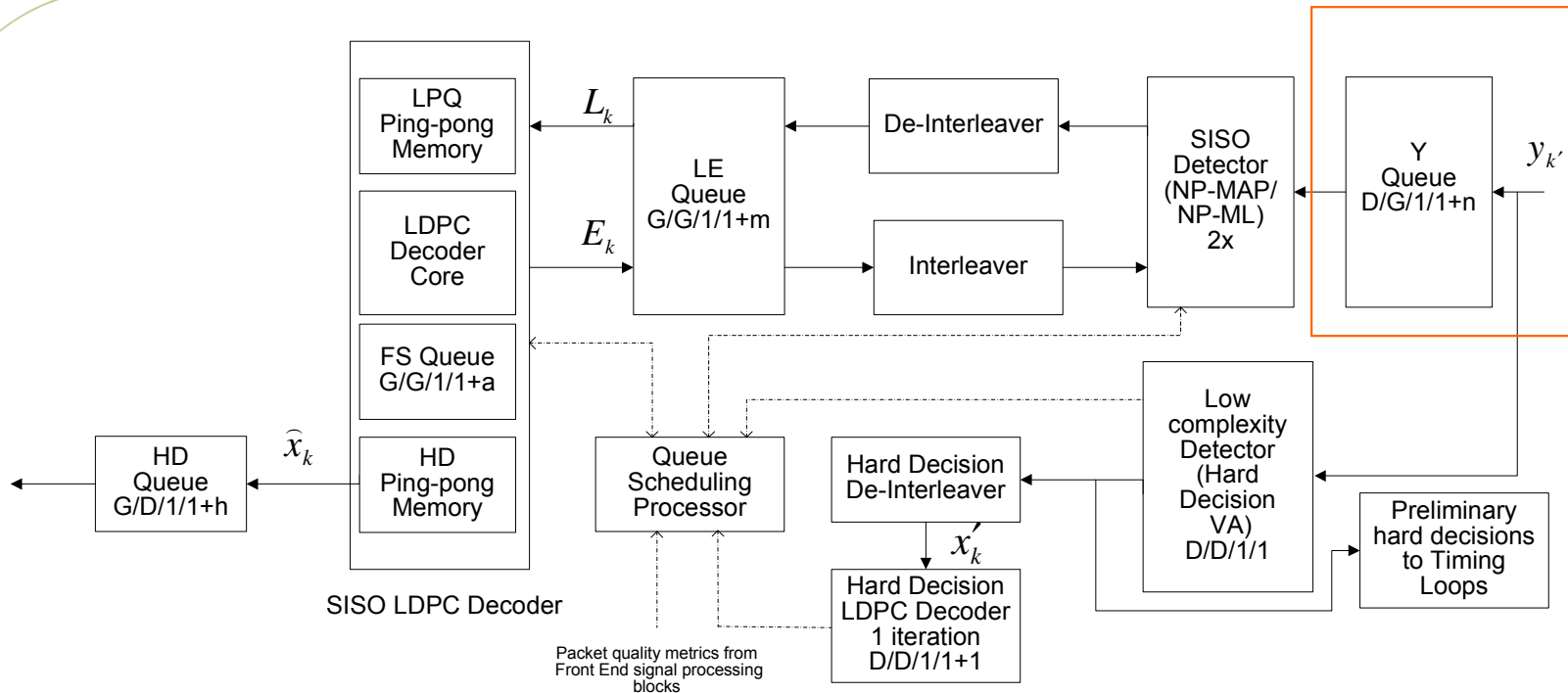
Secondary Data path

- The secondary data path contains the low complexity detector based on hard decision Viterbi algorithm, a hard decision interleaver followed by hard decision LDPC decoder that is sized for doing only one iteration.
- The secondary path thus does one reduced complexity global iteration and operates on the incoming packets immediately.
 - 1) it can generate preliminary decisions immediately (with a latency equal to T) to drive the front end timing loops thus making the front end processing immune from the variable time processing in the primary data path
 - 2) it can generate quality metrics to the queue scheduling processor.
- The low complexity detector is in the arrangement D/D/1/1 according to Kendall Notation[8]: the arrival times are deterministic, the processing time/service times are deterministic, one processor and one memory associated with the processor.
- The low complexity decoder is in the arrangement D/D/1/1+1 – this is similar to the low complexity detector except that there is one additional input buffer to enable the simultaneous filling of the input buffer while the hard decision iteration is being performed on a previous packet. Note that LDPC decoder needs the complete codeword before it can start the processing

Variations in number of global and local iterations

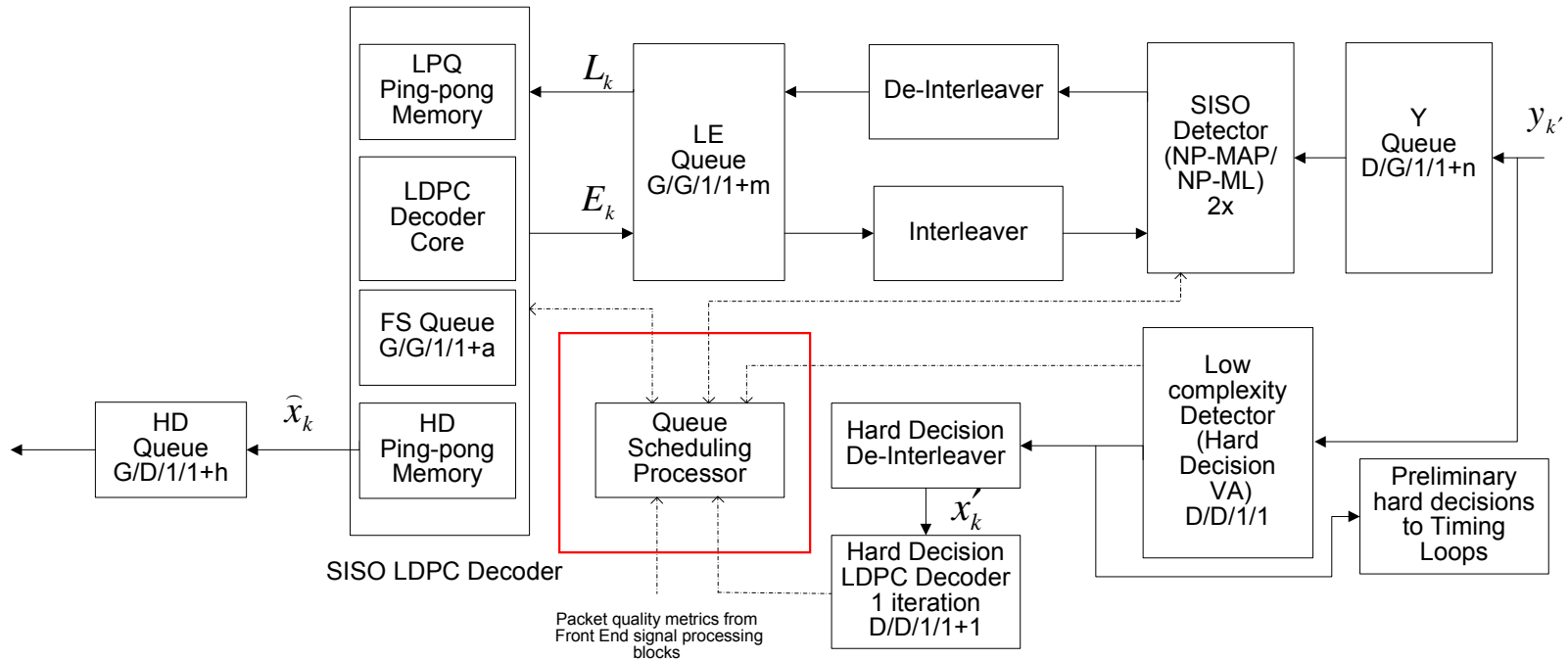
- In the last successful global iteration the LDPC decoder does the variable number of local iterations.
- The left-over LDPC decoder processing time is shared to increase the number of local iterations in following global iterations for the next packet.
- For each packet, at least one global iteration is performed and the distribution of required global iterations follows a general distribution that heavily depends on the signal to noise ratio.

Y Queue



- The y-sample data for each arriving packet is buffered in Y queue. Since the data comes at deterministic time intervals in a real-time application, the arrival process is D and the inter-arrival time is T.
- The overall processing/service time for each packet is variable and is a general distribution G. Assume that 4 y-samples per clock in each packet are arriving and the packets are coming continuously. In real-time applications, we need to be able to process 4-samples per clock though some latency is permitted.

Queue scheduling processor

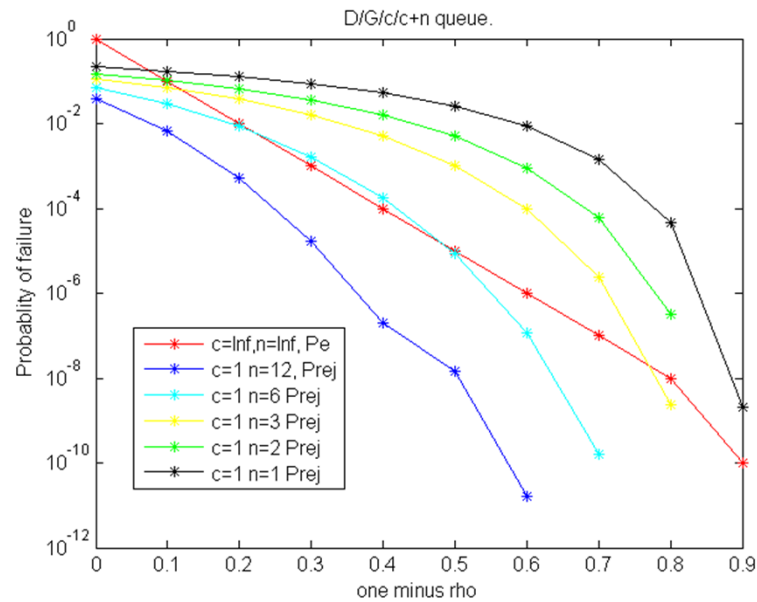


- The queue scheduling processor takes the various quality metrics from the secondary reduced complexity data path as well as the intermediate processing from the primary data path.
- One example of a quality metric is the number of unsatisfied checks from LDPC decoder. All the queues in the primary data path are preemptive such that the packets are processed according to the quality metric obtained through preprocessing.

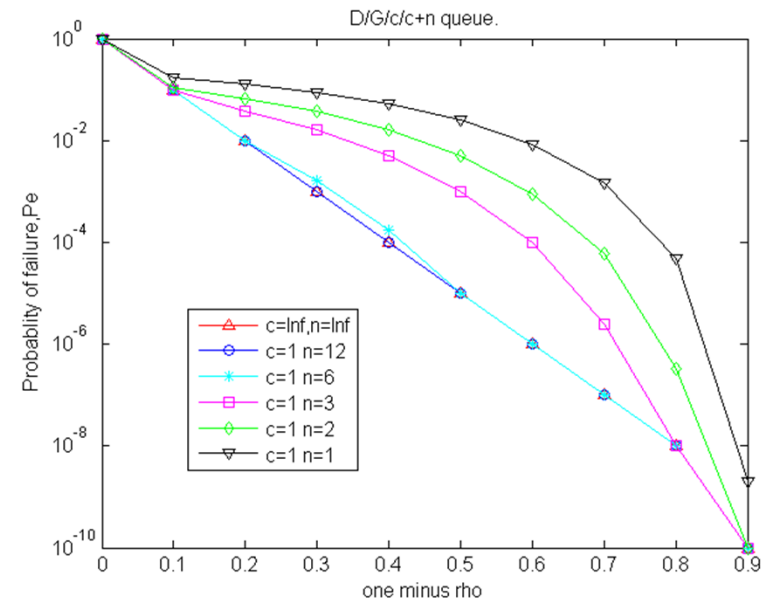
One example configuration

- In the example configuration of y -queues $D/G/c/1+n$, $c=1$ as we have one LDPC processor that can complete the processing of packet, the number of additional y -buffers, n .
 - Assume that all the other queues are optimized and have the values $m=4, a=3, h=20$.
 - Here $\rho = \lambda \cdot E(S)$ where λ is the average arrival rate and $E(S)$ is the average service time.
 - The performance measures are calculated under the assumption that λ is $1/T$ (i.e. 1 packet is coming every T time units) and constant and the average service time $E(S)$ (is less than or equal to T time units) varying based on the SNR.
 - Thus the value of ρ is between 0 and 1. $1 - \rho$ represents $1 - \rho$ and is indicator of the system's average availability.
 - The main requirement is that rejection probability should be kept low.
 - A) Should be less than $1e-6$ at ρ of 0.5
 - B) Should be less than $1e-9$ at ρ of 0.9
 - C) Asymptotically should reach 0 as ρ increases beyond 0.9
- The above requirements are based on the magnetic recording channel.

Different queue configurations



Probability of packet rejection
for different queue
configurations



Probability of overall packet
failure (Pe) for different queue
configurations

Queuing systems summary

- By comparing the previous two figures, we can see that we can either increase the processing power by 3 times (which is more expensive) or increase the number of y-buffers in the system n to 12 to achieve identical results. While it is not shown, we can get more benefits by doing both of them!
- Note that both the configurations in the previous two figures are still statistically buffered systems with $m=4, a=3, h=20$.
- If statistical buffering is disabled for other buffers in the system, then we need much higher number of processors up to 10 to gain the performance of a system that has no statistical buffering.
- As the average number of global iteration varies from 1 to 2 based on the SNR and the required number of global iterations vary from 1 to 20, the system with 10 processors with no statistical buffering would be idle for most of the time the proposed system with statistical buffering needs to have only one processor and can do the global iterations from 1 to 20.
- In conclusion, we show that statistical buffering if carefully done brings significant performance gains to the system while maintaining the low system complexity.

ERROR FLOOR MITIGATION

8/18/2013

Error Floors of LDPC Codes

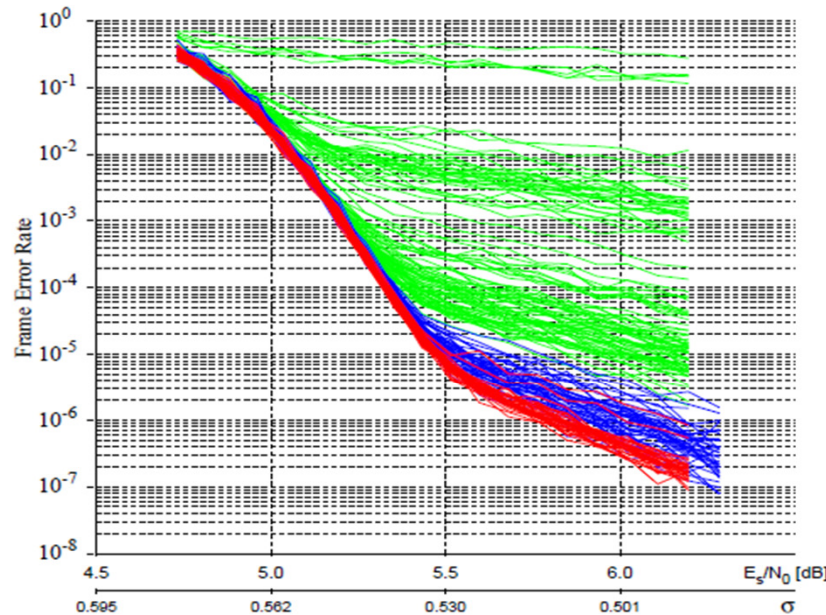


Figure 2: Frame error rates for random graphs (green), girth optimized (blue), and neighborhood optimized (red) graphs all of identical degree structure (nearly regular (3,15)).

When the BER/FER is plotted for conventional codes using classic decoding techniques, the BER steadily decreases in the form of a curve as the SNR condition becomes better. For LDPC codes and turbo codes that use iterative decoding, there is a point after which the curve does not fall as quickly as before, in other words, there is a region in which performance flattens. This region is called the *error floor region*. The region just before the sudden drop in performance is called the *waterfall region*. Error floors are usually attributed to low-weight codewords (in the case of Turbo codes) and trapping sets or near-codewords (in the case of LDPC codes).

Getting the curve steeper again: Error Floor Mitigation Schemes

- The effect of trapping sets is influenced by noise characteristics, H matrix structure, order of layers decoded, fixed point effects, quality of LLRs from detector. Several schemes are developed considering above factors. Some of them are

1) Changing the LLRs given to the decoder by using the knowledge of USCs, detector metrics and front end signal processing markers

2) With the knowledge of USCs, match the error pattern to a known trapping set. If the trapping set information is completely known, simply flip the bits and do the CRC.

If the trapping set information is partially known (i.e. only few bit error locations are stored due to storage issues), then do target bit adjustment using this information.

If no information on trapping set is stored, then identify the bits connected to USC based on H matrix information. Simply try TBE on each bit group.

Targetted bit adjustment on a bit/a bit group refers to the process of flipping the sign of these bits to opposite value and setting the magnitude of bit LLRs to maximum while keeping the other bit sign values unaltered but limiting their magnitude to around 5% of maximum LLR.

Couple of ways to reduce the number of experiments.

3) When multi-way interleaving is used, use of the separate interleavers on each component codeword.

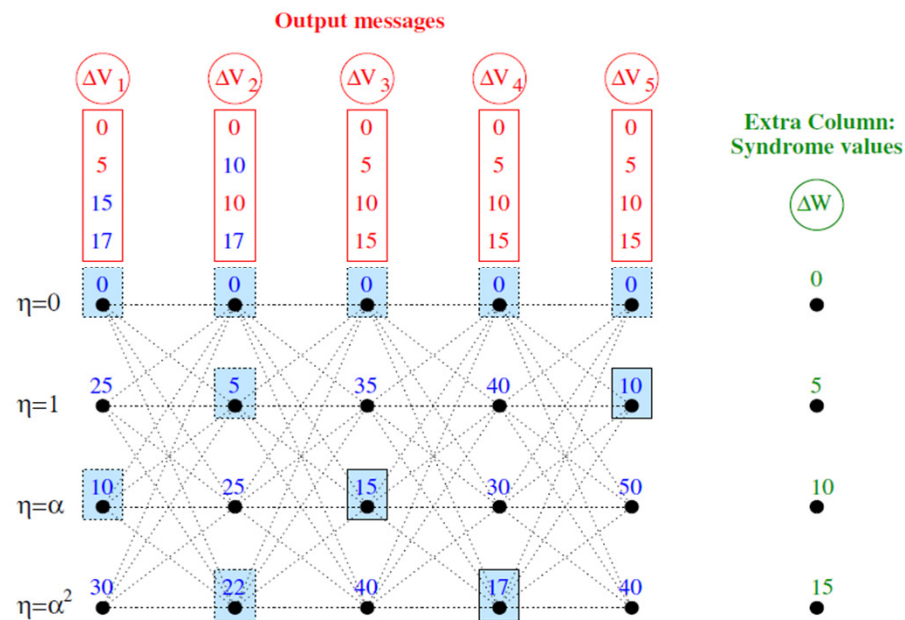
4) Skip-layer decoding: Conditionally decode a high row weight layer only when trapping set signature is present (USC < 32).

T-EMS, CHECK NODE UPDATE FOR NON-BINARY LDPC

8/18/2013

T-EMS

- Make use of the full trellis representation of messages in the **Delta-domain** messages.
- Select only a **small subset** of trellis nodes to build the most reliable configurations
 \Rightarrow minimization of the complexity.
- Add an extra-column to the trellis composed of **Syndrome** reliabilities for the parallel update of the output messages
 \Rightarrow improved decoding latency.



References

- [1] Gunnam, K.K.; Choi, G.S.; Yeary, M.B.; Shaohua Yang; Yuanxing Lee, "Next generation iterative LDPC solutions for magnetic recording storage," Signals, Systems and Computers, 2008 42nd Asilomar Conference on, Publication Year: 2008 , Page(s): 1148 – 1152
- [2] Kiran Gunnam, "LDPC Decoding: VLSI Architectures and Implementations", Invited presentation at Flash Memory Summit, Santa Clara, August 2012.
http://www.flashmemorysummit.com/English/Collaterals/Proceedings/2012/20120822_LDPC%20Tutorial_Module2.pdf
- [3] K. K. Gunnam, G. S. Choi, and M. B. Yeary, "Value-reuse properties of min-sum for GF(q)," Texas A&M University Technical Note, Oct.2006, published around August 2010.
Extended papers on [3]:
- [4] E. Li, K. Gunnam and D. Declercq, "Trellis based Extended Min-Sum for Decoding Nonbinary LDPC codes", in the proc. of ISWCS'11, Aachen, Germany, November 2011
http://perso-etis.ensea.fr/~declercq/PDF/ConferencePapers/Li_2011_ISWCS.pdf.gz
- [5] E. Li, D. Declercq and K. Gunnam, "Trellis based Extended Min-Sum Algorithm for Non-binary LDPC codes and its Hardware Structure", IEEE Trans. Communications., 2013

More references

6. Gunnam, KK; Choi, G. S.; Yeary, M. B.; Atiquzzaman, M.; "VLSI Architectures for Layered Decoding for Irregular LDPC Codes of WiMax," Communications, 2007. ICC '07. IEEE International Conference on 24-28 June 2007 Page(s):4542 - 4547
 7. Gunnam, K.; Gwan Choi; Weihuang Wang; Yeary, M.; "Multi-Rate Layered Decoder Architecture for Block LDPC Codes of the IEEE 802.11n Wireless Standard," Circuits and Systems, 2007. ISCAS 2007. IEEE International Symposium on 27-30 May 2007 Page(s):1645 – 1648
 8. Gunnam, K.; Weihuang Wang; Gwan Choi; Yeary, M.; "VLSI Architectures for Turbo Decoding Message Passing Using Min-Sum for Rate-Compatible Array LDPC Codes," Wireless Pervasive Computing, 2007. ISWPC '07. 2nd International Symposium on 5-7 Feb. 2007
 9. Gunnam, Kiran K.; Choi, Gwan S.; Wang, Weihuang; Kim, Euncheol; Yeary, Mark B.; "Decoding of Quasi-cyclic LDPC Codes Using an On-the-Fly Computation," Signals, Systems and Computers, 2006. ACSSC '06. Fortieth Asilomar Conference on Oct.-Nov. 2006 Page(s):1192 - 1199
 10. Gunnam, K.K.; Choi, G.S.; Yeary, M.B.; "A Parallel VLSI Architecture for Layered Decoding for Array LDPC Codes," VLSI Design, 2007. Held jointly with 6th International Conference on Embedded Systems., 20th International Conference on Jan. 2007 Page(s):738 – 743
 11. Gunnam, K.; Gwan Choi; Yeary, M.; "An LDPC decoding schedule for memory access reduction," Acoustics, Speech, and Signal Processing, 2004. Proceedings. (ICASSP '04). IEEE International Conference on Volume 5, 17-21 May 2004 Page(s):V - 173-6 vol.5
 12. GUNNAM, Kiran K., CHOI, Gwan S., and YEARY, Mark B., "Technical Note on Iterative LDPC Solutions for Turbo Equalization," Texas A&M Technical Note, Department of ECE, Texas A&M University, College Station, TX 77843, Report dated July 2006. Available online at <http://dropzone.tamu.edu> March 2010, Page(s): 1-5.
 13. K. Gunnam, G. Choi, W. Wang, and M. B. Yeary, "Parallel VLSI Architecture for Layered Decoding," Texas A&M Technical Report, May 2007. Available online at <http://dropzone.tamu.edu>
- Check <http://dropzone.tamu.edu> for technical reports.

Several features in this presentation and in references[1-13] are covered by the following 2 patents and other pending patent applications by Texas A&M University System (TAMUS).

- [P1] K. K. Gunnam and G. S. Choi, "Low Density Parity Check Decoder for Regular LDPC Codes," U.S. Patent 8,359,522
- [P2] K. K. Gunnam and G. S. Choi, "Low Density Parity Check Decoder for Irregular LDPC Codes," U.S. Patent 8,418,023
- [P3] K. K. Gunnam and G. S. Choi, "Low Density Parity Check Decoder for Irregular LDPC Codes," US Patent Application 20130151922
- [P4] K. K. Gunnam and G. S. Choi, "Low Density Parity Check Decoder for Regular LDPC Codes", US Patent Application 20130097469