POLITECNICO DI TORINO Repository ISTITUZIONALE

Design and Optimization of Adaptable BCH Codecs for NAND Flash Memories

Original

Design and Optimization of Adaptable BCH Codecs for NAND Flash Memories / Fabiano, Michele; Indaco, Marco; DI CARLO, Stefano; Prinetto, Paolo Ernesto. - In: MICROPROCESSORS AND MICROSYSTEMS. - ISSN 0141-9331. - STAMPA. - 37:4-5(2013), pp. 407-419. [10.1016/j.micpro.2013.03.002]

Availability: This version is available at: 11583/2506420 since:

Publisher: Butterworth Heinemann Publishers:Linacre Editore attuale..ELSEVIER SCI LTD, THE BOULEVARD,

Published DOI:10.1016/j.micpro.2013.03.002

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)



Politecnico di Torino

Design and Optimization of Adaptable BCH Codecs for NAND Flash Memories

Authors: S. Di Carlo, M. Fabiano, M. Indago, and P. Prinetto

Published in the Microproce or and Microsystems Vol. 37 ,Issues. 4-5, 2013, pp. 407-419.

N.B. This is a copy of the ACCEPTED version of the manuscript. The final PUBLISHED manuscript is available on SienceDirect:

URL: http://www.sciencedirect.com/science/article/pii/S0141933113000471

DOI: 10.1016/i.micpro.2013.03.002

© 2013 Elsevier. Personal use of this material is permitted. Permission from Elsevier must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

Design and Optimization of Adaptable BCH Codecs for NAND Flash Memories

S. Di Carlo, M. Fabiano, M. Indaco, and P. Prinetto

Department of Control and Computer Engineering Politecnico di Torino, Corso Duca degli Abruzzi 24, I-10129 Torino, Italy E-mail: {stefano.dicarlo, michele.fabiano, marco.indaco, paolo.prinetto}@polito.it.

Abstract

NAND flash memories represent a key storage technology for solid-state storage systems. However, they suffer from serious reliability and endurance issues that must be mutigated by the use of proper error correction codes. This paper proposes the design and implementation of an optimized Bose-Chaudhuri Hocquenghem hardware codec core able to adapt its correction capability in a range of predefined values. Code adaptability makes it possible to efficiently trade-off, in-field reliability and code complexity. This feature is very important considering that the reliability of a NAND flash memory continuously decreases over time, meaning that the required correction capability is not fixed during the life of the device. Experimental results show that the proposed architecture enables to save resources when the device is in the early stages of its lifecycle, while introducing a limited overhead in terms of area.

Key words:

Flash memories, Error correcting codes, memory testing, BCH codes

Preprint submitted to Microprocessors and Microsystems

August 30, 2013

1. Introduction

NAND flash memories are a widespread technology for the development of compact, low-power, low-cost and high data throughput mass storage systems for consumer/industrial electronics and mission critical applications. Manufacturers are pushing flash technologies into smaller geometries to further reduce the cost per unit of storage. This includes moving from traditional single-level cell (SLC) technologies, able to store a single bit of information, to multi-level cell (MLC) technologies, storing more than one bit per cell.

1

The strong transistor miniaturization and the adoption of an increasing 10 number of levels per cell introduce serious issues related to yield, reliability, 11]. Error correction codes (ECCs) must therefore ? and endurance ?? 12 be systematically applied. ECCs are a cost-efficient technique to detect and 13 correct multiple errors]. Flash memories support ECCs by providing 14 spare storage cells dedicated to system management and parity bit storage, 15 while demanding the actual implementation to the application designer [? 16 Choosing the correction capability of an ECC is a trade-off between ? 1. 17 reliability and code complexity. It is therefore a strategic decision in the 18 deign of a flash-based storage system. A wrong choice may either overesti-19 mate or underestimate the required redundancy, with the risk of missing the 20 target failure rate. In fact, the reliability of a NAND flash memory continu-21 ously decreases over time, since program and erase operations are somehow 22 destructive. At the early stage of their life-time, devices have a reduced 23 error-rate compared to intensively used devices [?]. Therefore, designing an 24 ECC system whose correction capability can be modified in-field is an attrac-25 tive solution to adapt the correction schema to the reliability requirements ²⁶ the flash encounters during its life-time, thus maximizing performance and ²⁷ reliability. ²⁸

This paper proposes the hardware implementation of an optimized adapt-²⁹ able Bose - Chaudhuri - Hocquenghem (BCH) codec core for NAND flash³⁰ memories and a related framework for its automatic generation.³¹

Even though there is a considerable literature about efficient BCH en-32 coder/decoder software implementations ?], modern flash-based mem-? 33 ory systems (e.g., Solid State Drives (SSDs)) usually resort to specific high 34 speed hardware IP core [?,?] in order to minimize the memory latency. This 35 is motivated by the fact that contemporary high-density MLC flash mem-36 ories require a more powerful error correction capability, and, at the same 37 time, they have to meet more demanding requirements in terms of read/write 38 latency. 39

Given this premise we will tackle a BCH hardware implementation for 40 encoding and decoding tasks. In particular, the main contribution of the 41 proposed architecture is its adaptability. It enables in-field selection of the 42 desired correction capability, coupled with high optimization that minimizes 43 th required resources. Experimental results compare the proposed architecture with typical BCH codecs proposed in the literature. 45

The paper is organized as follows: Section ?? shortly introduces basic 46 notions and related works. Sections ?? and ?? present a solution to reduce 47 resources overhead, while Section ?? and ?? overview the proposed adaptable architecture. Section ?? provides experimental results and Section ?? 49 summarizes the main contributions of the work and concludes the paper. 50

2. Background and related works

Several hard- and soft-decision error correction codes have been proposed ⁵² in the literature, including Hamming based block codes [? ?], Reed-Solomon ⁵³ codes [?], Bose-Chaudhuri-Hocquenghem (BCH) codes [?], Goppa codes ⁵⁴ [?], Golay codes [?], etc. ⁵⁵

Even though selected classes of codes such as Goppa codes have been $_{56}$ demonstrated to provide high correction efficiency [?], when considering the $_{57}$ specific application domain of flash memories, the need to trade-off code efficiency, hardware complexity and performances have moved both the scientific $_{59}$ and industrial community toward a set of codes that enable very efficient and $_{60}$ optimized hardware implementations [??]

flash lesigns used very simple Hamming based block codes. Old SLC 62 Hamming eddes are relatively straightforward and simple to implement in 63 both software and hardware, but they offer very knited correction capability 64 [? ?]. As the error rate increased with successive generations of both SLC 65 and MLC NAND flash memories, designers moved to more complex and pow-66 erful codes including Reed-Solomon (RS) codes [?] and Bose-Chaudhuri-67 Hocquenghen (BCH) codes [?]. Both codes are similar and belong to the 68 larger class of cyclic codes which have efficient decoding algorithms due to 69 their strict algebraic architecture, and enable very optimized hardware im-70 plementations. RS codes perform correction over multi-bit symbols and are 71 better suited when errors are expected to occur in bursts, while BCH codes 72 perform correction over single-bit symbols and better perform when bit er-73 rors are not correlated, or randomly distributed. In fact, several studies have 74 reported that NAND flash memories manifest non-correlated or randomly 75

distributed bit errors over a page [?] making BCH codes more suitable for 76 their protection. 77

An exhaustive analysis of the mathematics governing BCH code is out of the scope of this paper. Only those concepts required to understand the proposed hardware implementation will be shortly discussed. It is worth to mention here that, since several publications proposed very efficient hardware implementations of Galois fields polynomial manipulations, such manipulation will be used in both encoding and decoding operations (???).

Given a finite Galois field $GF(2^m)$ (with $m \ge 3$), a *t*-error-correcting BCH 84 code, denoted as BCH[n,k,t], encodes a k-bit message $b_{k-1}b_{k-2}...b_0$ ($b_i \in b_{k-1}b_{k-2}...b_0$) 85 GF(2) to a n-bit codeword $b_{k-1}b_{k-2}\dots b_0 p_{r-1}p_{r-2}\dots p_0$ $(b_i, p_i \in GF(2))$ by 86 adding r parity bits to the original message. The number r of parity bits 87 required to correct t errors in the n-bit codeword is computed by finding 88 the minimum m that solves the inequality k + p 2^m $\langle \langle \rangle$ - 1, where r =89 k∕∤ $< 2^m - 1$, the **BCH** code is called *shortened* $m \cdot t$. Whenever $m \geq 1$ r90 or *polynomial*. In a shortened BCH code the codeword includes less binary 91 symbols than the ones the selected Galois field would allow. The missing 92 information symbols are imagined to be at the beginning of the codeword 93 \checkmark are considered to be 0. Let α be a primitive element of $GF(2^m)$ and 94 $\psi_{1}(x)$ a primitive polynomial with α as a root. Starting from $\psi_{1}(x)$ a set of 95 minimal polynomials $\psi_i(x)$ having α^i as root can be always constructed [? 96]. For the same $GF(2^m)$, different valid $\psi_1(x)$ may exist [?]. The generator 97 polynomial q(x) of a t-error-correcting BCH code is computed as the Least 98 Common Multiple (LCM) among 2t minimal polynomials $\psi_i(x)$ $(1 \le i \le 2t)$. 99 Given that $\psi_i(x) = \psi_{2i}(x) \ (\forall i \in [1, t])$ [?], only t minimal polynomials must 100 be considered and g(x) can therefore be computed as:

$$g(x) = LCM[\psi_1(x), \psi_3(x)..., \psi_{2t-1}(x)]$$
(1)

101

When working with BCH codes, the message and the codeword can be 102 represented as two polynomials: (1) b(x) of degree k-1 and (2) c(x) of degree 103 n-1. Given this representation, both the encoding and the decoding process 104 can be defined by algebraic operations among polynomials in $GF(2^m)$. The 105 encoding process can be expressed as: 106

$$c(x) = m(x) \cdot x^{r} + Rem(m(x) \cdot x^{r})_{g(x)}$$
(2)

where $Rem(m(x) \cdot x^r)_{g(x)}$ denotes the remainder of the division between the 107 message left shifter of r positions and the generator polynomial g(x). This 108 remainder represents the r parity bits to append to the original message. 109

The BCH decoding process searches for the position of erroneous bits 110 in the codeword. This operation requires three main computational steps: 111 1) syndrome computation, 2) error locator polynomial computation, and 3) 112 error position computation. 113

Given the selected correction capability t, the decoding process requires 114 first the computation of 2t syndromes of the codeword c(x), each associated with one of the 2t minimal polynomials $\psi_i(x)$ generating the code. 116 Syndromes are calculated by first computing the remainders $R_i(x)$ of the 117 division between c(x) and each minimal polynomial $\psi_i(x)$. If all remainders 118 are null, c(x) does not contain any error and the decoding stops. Otherwise, 119 the 2t syndromes are computed by evaluating each remainder $R_i(x)$ in α^i : 120 $S_i = R_i(\alpha^i)$. Practically, according to (??), given that $\psi_i(x) = \psi_{2i}(x)$, only 121 t remainders must be computed and evaluated in 2t elements of $GF(2^m)$.

The most used algebraic method to compute the coefficients of the error 123 locator polynomial from the syndromes is the Berlekamp-Massey algorithm 124 [?]. Since the complexity of this algorithm grows linearly with the correction 125 capability of the code, it enables efficient hardware implementations. The 126 equations that link syndromes and error locator polynomial can be expressed 127 as:

$$\begin{pmatrix} S_{t+1} \\ S_{t+2} \\ \vdots \\ S_{2t} \end{pmatrix} = \begin{pmatrix} S_1 & S_2 & \dots & S_t \\ S_2 & S_3 & \dots & S_{t+1} \\ \vdots & \vdots & \vdots \\ S_t & S_{t+1} & \dots & S_{2t-1} \end{pmatrix} \begin{pmatrix} \lambda_t \\ \lambda_{t-2} \\ \vdots \\ \lambda_0 \end{pmatrix}$$
(3)

129

The Berlekamp-Massey algorithm iteratively solves the system of equa-130 tions defined in (??) using consecutive approximations 131 Finally, the Chien Machine searches for the roots of the error locator 132 polynomial $\chi(x)$ computed by the Berlekamp-Massey algorithm [?]. It 133 basically evaluates the polynomial $\lambda(x)$ in each element α^{i} of $GF(2^{m})$. If α^{i} 134 satisfies the equation $1 + \lambda_1 \alpha^i + \lambda_2 \alpha^{2i} + ... + \lambda_t (\alpha^i)^t = 0$, α^i is a root of the 135 er r locator polynomial $\lambda(x)$, and its reciprocal $2^m - 1 - i$ reveals the error 136 position. In practice, this computation is performed exploiting the iterative 137 relation: 138

$$\lambda\left(\alpha^{j+1}\right) = \lambda_0 + \sum_{k=1}^{t-1} \left[\lambda_k \left(\alpha^j\right)^k\right] \alpha^k \tag{4}$$

Several publications proposed optimized hardware implementations of ¹³⁹ BCH codecs with fixed correction capability [???????]. However, ¹⁴⁰ to the best of our knowledge, only Chen et al. proposed a solution allowing ¹⁴¹ limited adaptation by extending a standard BCH codec implementation [? ¹⁴²]. One of the main contributions of Chen et al. is a Programmable Parallel ¹⁴³ Linear Feedback Shift Register (PPLFSR), whose generic architecture is reported in Fig. ??. It enables to dynamically change the generator polynomial ¹⁴⁵ of the LFSR. This is a key feature in the implementation of an adaptable ¹⁴⁶ BCH encoder.



The gray box of Fig. ?? highlights the basic adaptable block of this ¹⁴⁸ circuit. It exploits a multiplexer, controlled by one of the coefficients of the ¹⁴⁹ desired divisor polynomial, to dynamically insert an XOR gate at the output ¹⁵⁰ of one of the related D-type flip-flops composing the register. The *s* vertical ¹⁵¹ stages of the circuit implement the parallelism of the PPLFSR computing ¹⁵² the state at clock cycle i + s, based on the state at cycle *i*. However, this ¹⁵³

solution has high overhead. In fact such PPLFSR is able to divide by all $_{154}$ possible *r*-bit polynomials, while just well selected divisor polynomials are $_{155}$ required.

Although Chen at al. deeply analyze the encoding process and the is-157 sues related to the storage of parity bits, the decoding process is scarcely 158 analyzed, without providing details on how adaptability is achieved. Four 159 different correction modes, namely t = (9, 14, 19, 24) are considered in ?? 160 for a BCH code defined on $GF(2^{13})$ with a block size of 5/2B (every 2KB 161 page of the flash is split in four procks). The selection of the 4 modes is based 162 on considerations about the number of parity bits to stope. However, there 163 is no provision to understand whether additional modes can be easily imple-164 mented. As an example, when selecting correction modes in which the size 165 of the codeword is not a multiple of the parallelism of the decoder, alignment 166 problems arise, which are completely neglected in the paper 167

3. Optimized Architectures of Programmable Parallel LFSRs

In this section, we will introduce an optimized block to perform an adaptable remainder computation. In fact, one of the most recurring operations in CH encoding/decoding is the remainder computation between a polynomial representing a message to encode/decode and a generator/minimal polynomial of the code, that depends on the selected correction capability. The PPLFSR of Fig. **??** can perform this operation [**?**].

168

A *r*-bit PPLFSR can potentially divide by any *r*-bit polynomial by properly controlling its configuration signals $(g_0 \dots g_{r-1})$. However, in BCH encoding/decoding, even considering an adaptable codec, just well selected divisor polynomials are required (e.g., the generators polynomials $g_9(x)$, $g_{14}(x)$, $_{178}$ $g_{19}(x)$, $g_{24}(x)$ of the four implemented correction modes of [?]). This computational block is therefore highly inefficient. Moreover, the set of divisor polynomials required in a BCH codec usually share common terms among each other. Such terms can be exploited to generate an optimized PPLFSR (OPPLFSR) architecture.

Let us consider, as an example, the design of a r=15-bit programmable 184 LFSR able to divide by two polynomials $p_1(x) = x^{15} + x^{13} + x^{10} + x^5 + x^3 + x + 1$ 185 and $p_2(x) = x^{13} + x^{12} + x^{10} + x^5 + x^4 + x^3 + x^2 + x + 1$ using a s=8-bit parallelism. 186

A traditional PPFLSR implementation would require $15 \times 8 = 120$ gray boxes (i.e., 120 XORs-MUXs) According to this implementation, this PP-LFSR could divide by any $2^{15} = 32,768$ possible 15-bit polynomials, even if just 2 polynomials (i.e., the 0.006% of its full potential) are required.

An analysis of the target divisor polynomials can be exploited to optimize 191 the PPLFSR architecture. Table ?? reports the binary representation of the 192 two polynomials.

Looking at Table ??, three categories of polynomial terms can be identified:

1. Common terms (represented in bold), i.e., terms defined in all considered polynomials $(x^{13}, x^{10}, x^5, x^3, x, \text{ and } 1 \text{ in Table } ??)$. For these terms, 197 an XOR will be always required in the PPLFSR, thus saving the area 198 dedicated to the MUX and the related control logic. 199

2. Missing terms (represented in underlined italic zeros), i.e., terms not 200 defined in any of the considered polynomials, $(x^{14}, x^{11}, x^9, x^8, x^7 \text{ and } 201$ x^6 in Table ??). For these terms both the XOR and the related MUX 202



can be avoided.

3. Specific terms, i.e., terms that are specific of a subset of the considered 204 polynomials $(x^{15}, x^{12}, x^4, x^2 \text{ in Table ??})$. These terms are the only 205 ones actually required. 206

We can therefore implement an optimized programmable LFSR (OP 207 PLFSR) with three main building blocks:

- 1. each common present term (i.e., columns of all "1" of Table ??) needs 209 an XOR, only; 210
- 2. each common absent term (i.e., columns of all "0" of Table ??) needs ²¹¹ neither XOR nor MUX; ²¹²
- 3. each specific term has a gray box, as Fig. ??; 213



Figure 2: Example of the resulting PPLFSR (a) and OPPLFSR (b) with 8-bit parallelism for x^{15} , x^{14} and x^{13} of $p_1(x)$ and $p_2(x)$

This optimization also applies on polynomials with very different lengths. ²¹⁵ As an example, an OPPLFSR with single bit parallelism and able to divide ²¹⁶ by $p_1(x) = x^{225} + x + 1$ and $p_2(x) = x + 1$, would only require a single ²¹⁷ adaptable block, compared to the 226 blocks required by a normal PPLFSR. ²¹⁸

203

214

Furthermore, the advantage of the OPPLFSR increases with the parallelism $_{219}$ of the block. In fact, with the same 2 polynomials, a 8-bit OPPLFSR would $_{220}$ require 8 adaptable blocks compared to $226 \times 8 = 1,808$ adaptable blocks of $_{221}$ a traditional PPLFSR. $_{222}$

For sake of generality, Fig. ?? shows the high-level architecture of a 223 generic OPPLFSR. Such a block is able to divide by a set $p_1(x), ..., p_M(x)$ 234 of polynomials. We denote with q the number of required gray poxes. 225



The OPPLFSR interface includes: a s-bit input port (b) used to feed 226 the data, a $\lfloor \log_2(M) \rfloor$ -bit input port (sel) used to select the polynomial of 227 th division, and a s-bit port (o) providing the result of the division. Two 228 blocks compose the OPPLFSR: $OPPLFSR_{net}$ and ROM. The OPPLFSR_{net} 229 represents the complete network, partially shown in the example of Fig. ??. 230 Given the output of the ROM, the q-bit signal g controls the MUXs of the 231 q gray boxes (Fig. ??) according to the selected polynomial. The ROM is 232 optimized accordingly with the design of the OPPLFSR, which leads to a 233 reduced ROM and to a lower area overhead w.r.t. a full PPLFSR. 234

4. BCH Code Design Optimization

In this section, we address first the issue of choosing the most suitable ²³⁶ set of polynomials for an optimized adaptable BCH code. Then, we propose ²³⁷ a novel block, shared between the adaptable BCH encoder and the decoder, ²³⁸ which reduces the area overhead of the resulting codec core.

4.1. The choice of the set of polynomials

The optimization offered by the OPPLFSR introduced in Section ??, may ²⁴¹ become ineffective if not properly exploited. It depends on the number and ²⁴² on the terms of the shared divisor polynomials implemented in the block. As ²⁴³ an example, an excessive number of shared polynomials may make it difficult ²⁴⁴ to find common terms, leading to an unwilled increase of the area overhead. ²⁴⁵ Therefore, the choice of the polynomials to share is critical and must be ²⁴⁶ properly tailored to the overall design. ²⁴⁷

Let us denote by Ω the set of t generators $g_{t}(x)$ and t minimal polynomials 248 ψ_i which fully characterize an adaptable BCH code (see Section ??). Since 249 for $GF(\mathbb{Z}^n)$ several primitive polynomials $\psi_i(x)$ can be used to define the 250 code, several set Ω_i can be constructed. Choosing the most suitable set Ω_i is 251 drived to obtain an effective design of the OPPLFSR. On the one hand, it 252 can be shown that the complexity of Ω_i increases with m [???]. On the 253 other hand, the current trend is to adopt BCH codes with high values of m254 (e.g., $GF(2^{15})$) because current flash devices features a worse bit error rate [? 255]. Therefore, a simple visual inspection of each set Ω_i is not feasible to find 256 the most suitable set of polynomials. An algorithmic approach is therefore 257 mandatory. 258

240

Each set Ω_i can be classified resorting to a *Maximum Correlation Index* ²⁵⁹ (MCI). We define as $MCI(p_1, p_2, ..., p_N)$ the maximum number of common ²⁶⁰ terms shared by a generic set of polynomials $p_1, p_2, ..., p_N$. As an example, ²⁶¹ the polynomials of Table **??** have $MCI(p_1, p_2) = 12$.

In the sequel, we introduce an algorithm to assess each set Ω_i according to its MCI. Given $i = \{1, ..., Y\}$, for each set Ω_i .

1. consider
$$\Omega_i = \{p_1, ..., p_N\}$$
 and $v_0 = p_i$;

- 2. determine the polynomial p_h such that the partition $S_1 = (v_0, p_h)$ has the maximum $MCI(v_0, p_h)$ where $h = \{1, ..., N\}$ and $p_h \neq v_0$; 267
- 3. determine the polynomial p_k such that the partition $S_{i,1} = ((v_0, p_h), p_k)$ has the maximum $MCI(v_0, p_h, p_k)$, where $k = \{1, ..., N\}$ and $p_k \neq p_h \neq$ 269 v_0 ;
- 4. repeat step 3 until all polynomials have been considered in the partition $_{271}$ $S_{i,1}$; $_{272}$
- 5. change the starting polynomial to the next one, e.g., $v_0 = p_2$, considering z_{73} $S_{i,2}$ and repeat steps 2-4; z_{74}
- 6. when $v_0 = p_N$, consider the next set Ω_{i+1} ; 275

The algorithm ends when all sets Ω_i have been analyzed. For each Ω_i , 276 the output is a set of partitions: 277

$$S_{i,j} = \{S_{i,1}, S_{i,2}, \dots, S_{i,N}\}$$
(5)

Fig. ?? graphically shows the MCI of two partitions generated from two $_{278}$ different starting points, for an hypothetical set Ω_i . $_{279}$

Fig. ?? shows that MCI always has a decreasing trend with the size of $_{280}$ the partition S. This is straightforward since adding a polynomial may only $_{281}$



decrease or keep constant the current value of MCI. The curves, reported ²⁸² in ? are critical in the choice of the most suitable set of polynomials for ²⁸³ an optimized BCH code. For each partition $S_{i,j}$ with $j = \{1...N\}$, we can ²⁸⁴ compute the average MCI (MCI_{avg}) as: ²⁸⁵

$$MCI_{avg}(S_{i,j}) = \frac{1}{N} \sum_{l=1}^{N-1} MCI_l$$
 (6)

Eq. ?? applies to each set Ω_i where $i = \{1...Y\}$. The best partition of the set Ω_i is then computed selecting the one with 287

16

maximum MCI_{avg} :

$$S_{best_i} = \underset{i}{argmax} \left[MCI_{avg} \left(S_{i,j} \right) \right] \tag{7}$$

Finally, Eq. ?? compares the best partition of each set Ω_i to find the best 289 set of polynomials:



Let us provide an example to support the understanding of the algorithm. ²⁹³ Suppose to consider a single set Ω_i composed of the polynomials of Table ??. ²⁹⁴ The steps of the algorithm are: ²⁹⁵

1. Let us start with $v_0 = p_1$

17

288

296

- 2. We first evaluates $MCI(p_1, p_2) = 3$, $MCI(p_1, p_3) = 4$, $MCI(p_1, p_4) = {}^{297}$ 3. Since $MCI(p_1, p_3) = 4$ is the maximum, the resulting partition is ${}^{298}S_{i,1} = \{p_1, p_3\}$
- 3. The next step considers $MCI((p_1, p_3), p_2) = 3$ and $MCI((p_1, p_3), p_4) = 300$ 3. It is straightforward that the choice of either p_2 or p_4 does not affect 301 the final value of the MCI_{avg} .

Given Ω_i with starting point p_1 , it can be shown that the final partition 303 is $S_{i,1} = \{((p_1, p_3), p_4), p_2\}$ with a $MCI_{ang} = (4+3+3)/4 = 2.5$ from Eq. ??. 304

The complete algorithm iterates this computation for all possible starting ³⁰⁵ points. Fig. ?? graphically shows the output of the MCI associated with each ³⁰⁶ partition $S_{i,j}$ calculated for the following starting point $j = \{1, 2, 3, 4\}$. ³⁰⁷



Figure 5: The MCI Trend of Table ??

According to Eq. ??, $S_{i,2}$ (the bold line) is the S_{best_i} of the example of 308

Table ??, with a $MCI_{avg}(S_{i,j}) = 4$.

4.2. Shared Optimized Programmable Parallel LFSRs

Let us assume to design an adaptable BCH code with correction capability $_{311}$ from 1 up to t_M . Such a code needs to compute remainders of the division $_{312}$ of:

- the message m(x) by (potentially) all generator polynomials from g_1 314 up to g_{t_M} , for the encoding (??); 315
- the codeword c(x) by (potentially) all minimal polynomials from $\psi_1(x)$ ³¹⁶ up to $\psi_{2t_M-1}(x)$, to compute the set of syndromes required during the ³¹⁷ decoding phase. ³¹⁸

In a traditional implementation, these computations are performed by ³¹⁹ two separate set of LFSRs. In this paper, we propose to devise a shared ³²⁰ set of LFSRs able to: (i) perform all these computations, and (ii) reduce the ³²¹ overall cost in terms of resources overhead. Therefore, we can adopt the same ³²² shared set of LFSRs both in the encoding and decoding processes. This is ³²³ possible since in a flash memory these operations are, in general, not required ³²⁴ at the same time. ³²⁵

The OPPLFSR, introduced in Section ??, is the main building block of $_{326}$ the set of shared LFSRs. Therefore, we will refer hereafter to such set of $_{327}$ LFSRs as shared OPPLFSR (shOPPLFSR). Fig. ?? shows the high-level $_{328}$ architecture of the shOPPLFSR. Its interface includes: a *s*-bit input port $_{329}$ (IN) used to input the data to be divided, a $\lceil \log_2(N) \rceil$ -bit input port (en) $_{330}$ used to enable each OPPLFSR, an input port (sel) used to select the proper $_{331}$

310

polynomial by which each OPPLFSR has to divide, and a N \times s-bit port (p) $_{332}$ providing the result of the division. $_{333}$



Figure 6: The shOPPLFSR architecture is composed by multiple OPPLFSRs

Given N OPPLFSRs and a maximum correction capability t_M , each ³³⁴ OPPLFSR_i performs the division by a set of generator polynomials g(x) and ³³⁵ minimal polynomials $\psi(x)$. Such shOPPLFSR can be seen as an optimized ³³⁶ programmable LESR able to: ³³⁷

• divide by all generator polynomials from $g_1(x)$ to $g_{t_M}(x)$;

• divide by specific subsets of minimal polynomials from Eq. ??, as well. 339

338

An improper choice of the shared polynomials g(x) and $\psi(x)$ can dramatically reduce the performance of the overall BCH codec. Also the partitioning strategy adopted is critical to maximize the optimization in terms of area, minimizing the impact on the latency of encoding/decoding operations. 341

The algorithm presented in Section ?? provides a valuable support for the exploration of this huge design space. In fact, the proposed method can be sploited to properly partition polynomials into the different OPPFLSRs of 346 Fig. ??, in order to maximize the optimization of the resulting shOPPFLSR. ³⁴⁷ Such optimization should not be obtained following blindly the outcomes of ³⁴⁸ the algorithm, but always tailoring them to the specific design. Regarding ³⁴⁹ this topic, Section ?? provides more details about our experimental setup ³⁵⁰ and the related experimental results.

352

5. Adaptable BCH Encoder

In this section, we propose an adaptable BCH encoder which exploits the 353 shOPPLFSR of Section ??. According to the BCH theory the shOPPLFSR 354 of Fig. ?? is a very efficient circuit to perform the computation expressed in 355 Eq. ??. However, in the encoding phase, the message m(x) must be multi-356 plied by x^r before calculating the reminder of the division by g(x) (see Eq. 357 ??). This can be obtained without significant modifications of the architec-358 ture of shOPPFLSR. It is enough to input the bits of the message directly 359 in the most significant bit of the LFSR, instead than starting from least 360 significant bit. Fig. ?? shows the high level architecture of the adaptable 361 encoder 362

The encoder's interface includes: a s-bit input port (IN) used to input the k- is message to encode starting from the most significant bits, a $\lceil \log_2(t_M) \rceil$ - β_{64} bit input port (t) selecting the requested correction capability in a range between 1 and t_M , a start input signal used to start the encoding process and a s-bit output port (OUT) providing the r parity bits. Three blocks β_{67} compose the encoder: a shOPPLFSR, a flush logic and a controller. β_{68}

The shOPPLFSR performs the actual parity bits computation. According to the BCH theory, adaptation is achieved by supporting the computation ³⁷⁰



Figure 7: High-level architecture of the adaptable encoder highlighting the three main building blocks and their main connections.

of remainders with t_{M} generator polynomials, one for each value t may as-371 sume. The controller achieves this task in two steps: (i) enabling the proper 372 OPPLFSR through the len signal, and (ii) selecting the proper polynomial 373 according to the desired correction capability t. through the Isel signal. 374 Then, it manages the overall encoding process based on two internal param-375 eters: 1) the number of s-bit words composing the message (fixed at design 376 the number of produced s-bit parity words, that depends on time) and 2)377 the selected correction capability. The flush logic splits the r parity bits into 378 words, providing them in output, one per clock cycle. 379

To further optimize the encoding and the decoding process, since in a flash memory these operations are not required at the same time, the encoder's shOPPLFSR can be merged with the shOPPLFSRs that will be employed in the syndrome computation (see Section ??), thus allowing additional area saving.

6. Adaptable BCH Decoder

Fig. ?? presents the high-level architecture of the proposed adaptable 386 decoder. The decoder's interface includes: a s-bit input port (IN) used to 387 input the n-bit codeword to decode (starting from the most significant bits), 388 a $\lceil \log_2(t_M) \rceil$ -bit input port (t) to select the desired correction capability. a 389 start input signal to start the decoding and a set of output ports previding 390 information about detected errors. In particular:

- deterr is a $\lceil \log_2(t_M) \rceil$ -bip port providing the number of errors that have been detected in a codeword. In case of decoding failure it is set to 0; 394
- erradd and ermask provide information about the detected error positions. Assuming the codeword split into h-bit words, erradd is used 396 as a word address in the codeword and errmask is a h-bit mask whose 397 asserted bits indicate detected erroneous bits in the addressed word. 398 The parallelism h of the error mask depends on the parallelism of the 399 Chien machine, as explained later in this section; 400

• vmask is asserted whenever a valid error mask is available at the output 401 of the decoder; 402

- fail is asserted whenever an error occurred during the decoding process (e.g., the number of errors is greater than the selected correction 404 capability); 405
- end is asserted when the decoding process is completed. 406



Figure 8: High-level architecture of the adaptable decoder, highlighting the four main building blocks: the adaptable syndrome machine, the adaptable iBM machine, the adaptable Chien machine, and the controller in charge of managing the overall decoding process

The full decoder therefore includes four main blocks: (1) the Adapt *able Syndrome Machine*, computing the syndromes of the codeword, (2) the Adaptable inversion-less Berlekamp Massey (iBM) Machine, that elaborates the syndromes to produce the error locator polynomial, (3) the Adapt *able Chien Search Machine* in charge of searching for the error positions, and (4) the Controller coordinating the overall decoding process.

6.1. Adaptable Syndrome Machine

Fig. ?? shows the high-level architecture of the proposed adaptable syndrome machine with correction capability $1 \le t \le t_M$

413



Figure 9: Architecture of the adaptable Syndrome Machine

According to Section ??, remainders can be calculated by a set of Parallel ⁴¹⁶ LFSRs (PLFSRs) whose architecture is similar to the one of the PPLFSR ⁴¹⁷ of Fig. ??, with the only difference that the characteristic polynomial is ⁴¹⁸ fixed (XOR gates are inserted only where needed, without multiplexers). ⁴¹⁹

Each PLFSR computes the remainder of the division of the codeword by a 420 different minimal polynomial $\psi_i(x)$. Given two correction capabilities t_1 and 421 t_2 with $t_1 < t_2 \le t_M$, the set of $2t_1$ minimal polynomials generating the code 422 for t_1 is a subset of those generating the code for t_2 . To obtain adaptability 423 of the correction capability in a range between 1 and t_M , the syndrome 424 machine can therefore be designed to compute the maximum number $t_{\rm M}$ 425 of remainders required to obtain $2t_M$ syndromes. Based on the selected 426 correction capability t, only the first t PDFSRs out of the t_M available in the 427 circuit are actually enabled through the Enable div. network of Fig. ??. 428

A full parallel syndrome calculator, including the PLFSRs, requires a 429 considerable amount of resources that are underutilized in the early stages 430 of the flash lifetime when reduced correction capability is required. To opti-431 mize the adaptable syndrome machine and to trade-off between complexity 432 and performance, we exploit the shOPPLFSR introduced in Section ??. The 433 architecture proposed in Fig. ?? includes two sets of LFSRs for remainder 434 computation: (i) conventional PLFSRs, and (ii) shOPPLFSR. Conventional 435 PLFSRs are exploited for parallel fast computation of low order syndromes 436 required when the requested correction capability is below a given threshold. 437 shOPRDFSR is designed to divide for selected groups of minimal polynomials 438 not covered by the fixed PPLFSRs. It represents a shared resource utilized 439 when the requested correction capability increases. It enables area reduction 440 at the cost of a certain time overhead. The architectural design, chosen for 441 the fixed PLFSRs and the OPPLFSR, enables to trade-off hardware com-442 plexity and decoding time, as it will be discussed in Section ??. 443

It is worth to mention here that the parallel architecture of the PLFSR, 444



Figure 10: Example of the schema of a byte aligner for t = 2 and s = 3

coupled with the adaptability of the code, introduces a set of additional 445 word alignment problems that must be addressed to correctly adapt the 446 syndrome calculation to different values of t. The syndrome machine receives 447 the codeword in words of s bits, starting from the most significant word. 448 When the number of parity bits does not allow to align the codeword to the 449 parallelism s, the unused bits of the last word are filled with 0. To correctly 450 compute each syndrome, the parity bit r_0 of the codeword must enter the 451 least significant bit of each LFSR. The aligner block of Fig. ?? assures 452 this condition by properly right-shifting the codeword while it is input into 453 the syndrome machine. Let us consider the following example: k = 2KB, 454 t = 2, s = 8 and therefore $r = m \cdot t = 30$. Since 30 is not multiple of nn\ 455 $s \neq 8$, the codeword is filled with two zeros and p_0 is saved in position 2 of 456 the last byte of the codeword $(m_{2047} m_{2046} \dots m_1 m_0 p_{29} p_{28} \dots p_1 p_0 0 0)$. In this case 457 the PLFSRs require a 2-bit alignment, implemented by the network of Fig. 458 ??. It simply delays the last 2 input bits resorting to two flip-flops, whose 459 initial state has to be zero, and properly rotates the remaining input bits. 460 Changing the correction capability of the decoder changes the number of 461 parity bits of the codeword, and therefore the required alignment. Given the $_{462}$ parallelism s of the decoder, a maximum of s alignments must be provided $_{463}$ and implemented in the *Aligner* block of Fig. ??. $_{464}$

With the proper alignment, the PLFSRs can perform the correct division 465 and the evaluators can provide the required syndromes. The evaluators are 466 simple combinational networks involving XOR operations, according to the 457 Galois Fields theory (the reader may refer to [?] for specific implementation 468 details).

6.2. Adaptable Berlekamp Massey Machine

In our adaptable codec we implemented the inversion-less Berlekamp-Massey (iBM) algorithm proposed in [?] which is able to compute the error $_{472}$ locator polynomia () (x) in t iterations $_{473}$

470

The main steps of the computation are reported in Alg. **??**. At iteration 474 i (rows 2 to 12), the algorithm finds an error locator polynomial $\lambda(x)$ whose 475 coefficients solve the first i equations of (??) (now 4). It then tests if the 476 same polynomial solves also i + 1 equations (row 5). If not, it computes a 477 discrepancy term δ so that $\lambda(x) + \delta$ solves the first i + 1 equations (row 9). 478 This iterative process is repeated until all equations are solved. If, at the 479 end of the iterations, the computed polynomial has a degree lower than t, 480 it correctly represents the error locator polynomial and its degree represents 481 the number of detected errors; otherwise, the code is unable to correct the 482 given codeword. 483

The architecture of the iBM machine is intrinsically adaptive as long as 484 one guarantees that the internal buffers and the hardware structures are sized 485 to deal with the worst case design (i.e., $t = t_M$). The coefficients of $\lambda(x)$ are 486

Algorithm 1 Inversion-less Berlekamp-Massey alg.

1: $\lambda(x) = 1, k(x) = 1, \delta = 1$ 2: for i = 0 to t - 1 do $d = \sum_{j=1}^{t} \left(\lambda_j \cdot S_{2i-j} \right)$ 3: $\lambda(x) = \delta\lambda(x) + d \cdot x \cdot k(x)$ 4:if d = 0 OR $Deg(\lambda(x)) > i$ then 5: $k(x) = x^2 \cdot k(x)$ 6: 7: else 8: $k(x) = x \cdot k(x)$ $\delta = d$ 9: end if 10: 11: i=i+112: end for 13: if $Deg(\lambda(x)) < t$ then 14: output $\lambda(x)$; 15: else output FAILURE 16:17: end if

m-bit registers whose number depends on the correction capability. In the 487 worst case, up to t_M coefficients must be stored for each polynomial. 488 The adaptable iBM machine therefore includes two m-bit register files 489 with t_M registers to store these coefficients. Whenever the requested correc-490 tion capability is lower than t_M some of the registers will remain unused. The 491 number of multiplications performed during the computations also depends 492 on t. Row 3 requires t multiplications, while row 4 requires t multiplications 493 to compute $\delta \lambda_i(x)$ and t multiplications to compute $d \cdot x \cdot k(x)$. 494

We implemented a serial iBM Machine including 3 multipliers for $GF(2^m)_{495}$ to perform multiplications of rows 3 and 4. It can perform each iteration of $_{496}$ the iBM algorithm in 2t clock cycles (t cycles for row 3 and t cycles for ⁴⁹⁷ row 4) achieving a time complexity of $2t^2$ clock cycles. This implementation ⁴⁹⁸ is a good compromise between performance and hardware complexity. An ⁴⁹⁹ input t dynamically sets the number of iterations of the algorithm, thus ⁵⁰⁰

502

6.3. Adaptable Chien Machine

The overall architecture of the proposed adaptable Chien Machine is 503 shown in the Fig. ??. The machine first loads into t_M in-bit registers the 504 coefficients from λ_1 to λ_{t_M} of the error locator polynomial $\lambda(x)$ computed by 505 the iBM machine (1d = 0) The actual search is then started (1d = 1). At 506 each clock cycle, the block performs h parallel evaluations of $\lambda(x)$ in GF(2^m) 507 and outputs h-Wiword, denoted as errmask. Each bit of errmask corre-508 sponds to one of the h candidate error locations that have been evaluated. 509 Asserted bits denote detected evers. This mask can then be XORed (outside 510 the Chien Machine) with the related bits of the oddeword in order to correct 511 the detected erroneous bits. 512

The architecture of Fig. ?? provides an adaptable Chien machine with ⁵¹³ lower area consumption than other designs [?], having, at the same time, ⁵¹⁴ a marginal impact on performance. Four interesting features contribute to ⁵¹⁵ such optimization: (i) constant multipliers substructure sharing, (ii) adaptability to the correction capability, (iii) improved fast skipping to reduce the ⁵¹⁷ decoding time, and (iv) reduced full GF multipliers area. In the sequel, we ⁵¹⁸ briefly address each feature. ⁵¹⁹

The first feature is represented by the optimized GF Constant Multipliers $_{520}$ (optGFCM) networks of Fig. ??. The *h* parallel evaluations are based on $_{521}$



Figure 11: Architecture of the proposed parallel adaptable Chien Machine with parallelism equal to h

 (\mathbf{x}) , the parallel evaluation of equation equation (??? In the worst case (t)= 522 (??) requires a matrix of $t_M \times h$ constant Galois multipliers. They multi-523 ply the content of the t_M registers by $\alpha, \alpha^2, ..., \alpha^{4/4}$, respectively. However, 524 we can note that each column of constant GF multipliers shares the same 525 Therefore we can iteratively group their best-matching commultiplicand, 526 into the t_M optGFCM networks of Fig. ??. Such optGFCMs binations ? 527 provide up to 60% reduction of the hardware complexity of the machine with 528 no mpact on performance. 529

The second feature is the adaptability of the Chien machine. The rows of the matrix define the parallelism of the block (i.e., the number of evaluations per clock cycles), while the columns define the maximum correction capability of the block. Whenever the selected correction capability t is lower than t_M , the coefficients of the error locator polynomial of degree greater than t are equal to zero and do not contribute to equation (??), thus allowing us to t_{10}^{530} adapt the computation to the different correction capabilities.

The third feature stems from a simple observation. Depending on the 537 selected correction capability t, not all the elements of $GF(2^m)$ represent 538 realistic error locations. In fact, considering a codeword composed of k bits 539 of the original message and $r = m \cdot t$ parity bits, only $k + m \cdot t$ out of 2^{m} 540 elements of the Galois field represent realistic error locations. Given that an 541 error location L is the inverse of the related CF element $(L = 2^{\circ})$ thě 542 elements of $GF(2^m)$ in which the error locator polynomial must be evaluated 543 are in the following range: 544

$$\sum_{\text{cation L=0}}^{2m-1}, \sum_{\text{error location } L=k+m\cdot t-1}^{2^m-k-m\cdot t}$$
(9)

536

All elements between α^0 and $\alpha^{2^m-k-m'}$ can be skipped to reduce the 545 computation time. Differently from fixed correction capability fast skipping 546 Chien machines this interval is not constant here but depends on the se-547 lected t. The architecture of Fig. ?? implements an adaptable fast skipping 548 by initializing the internal registers to the coefficients of the error corrector 549 polynomial multiplied by a proper value $\beta_{ini}^t = \alpha^{2^m - k - m \cdot t - 1}$. For each value 550 t_M m bit constant values corresponding to β_{ini}^t , $(\beta_{ini}^t)^2$, ..., $(\beta_{ini}^t)^{t_M}$ 551 must be stored in an internal ROM (not shown in Fig. ??) and multiplied 552 by the coefficients λ_i using a full GF multiplier. 553

This is connected with the last feature, the reduced GF Full Multipliers ⁵⁵⁴ (redGFFM) network of Fig. ??. Each full GF multiplier has a high cost in ⁵⁵⁵ terms of area. Since they are used only during initialization of the Chien, the ⁵⁵⁶ redGFFM adopts only $z \leq t_M$ full GF multipliers. It also includes a (λ) input ⁵⁵⁷ port to input z coefficients, per clock cycles, of the error locator polynomial. ⁵⁵⁸ This network enables to reduce area consumption, at a reasonable cost in ⁵⁵⁹ terms of latency. ⁵⁶⁰

For the sake of brevity, a detailed description of the controller required 561 to fully coordinate the decoder's modules interaction is out of the scope of 562 this paper.

7. Experimental Results

This section provides experimental data from the implementation of the 565 adaptable BCH codec proposed on a selected case study. 566

564

567

7.1. Automatic generation framework

To cope with the complexity of a manual design of these blocks, a semi-568 automatic generation tool named ADAGE (ADaptive ECC Automatic GEn-569 erator) [? Able to generate a fully synthesizable adaptable BCH codec core 570 following the proposed architecture has been designed and exploited in this 571 experimentation extending a preliminary framework previously introduced 572 The overall architecture of the framework is in Fig. ??. in |? (573 The code analyzer block represents the first computational step required 574 to select the desired code correction capability based on the Bit Error Rate 575 (BER) of a page of the selected flash [?]. The BER is the fraction of er-576 roneous bits of the flash. It is the key factor used to select the correction 577 capability. Two values of BER must be considered. The former is the raw 578 bit error rate (RBER), i.e., the BER before applying the error correction. 579 It is technology/environment dependent and increases with the aging of the 580 page ? ?]. The latter is the uncorrectable bit error rate (UBER), i.e., 581



Figure 12: BCH codec automatic generation framework.

the BER after the application of the ECC, which is application dependent. $_{582}$ It is computed as the probability of having more than t errors in the codeword (calculated as a binomial distribution of randomly occurred bit errors) $_{584}$ divided by the length of the codeword [?]: $_{585}$

$$UBER = \frac{P\left(E > t\right)}{n} = \frac{1}{n} \sum_{i=t+1}^{n} \binom{n}{i} \cdot RBER^{i} \cdot (1 - RBER)^{n-i}$$
(10)

Given the RBER of the flash and the target UBER, Eq. ?? can be exploited to compute the maximum required correction capability of the code and consequently the value of *m* that defines the target GF. Given these two parameters, the Galois Field manager exploits an internal polynomials database to generate the set of minimal polynomials and the related generator polynomials for the selected code.

Finally, the RTL VHDL code generator combines these parameters and 592 generates a RTL description of the BCH encoder and decoder implementing 593 the architecture illustrated in this paper. 594

The whole framework combines Matlab software modules with custom 595 C programs The full framework code is available for download at http: 596 //www.testgroup.polito.it in the Tools section of the website. 597

598

7.2. Experimental setup

Experiments have been performed, using as a case study a 2-bit per cell ⁵⁹⁹ MLC NAND Flash Memory featuring a 45nm manufacturing process designed for low-power applications, with page size of 2KB plus 64B of spare ⁶⁰¹ cells. The memory has an 8-bit I/O interface. Considering the design of ⁶⁰² the BCH code, the current trend is to enlarge the block size k over which ⁶⁰³ ECC operations are performed. In fact, longer blocks better handle higher 604 concentrations of errors, providing more protection while using fewer parity 605 bits [?]. For this reason, we adopted a block size k = 2KB, equal to the 606 page size of the selected memory. 607

Experiments performed on the flash provided that, in a range between 608 10 and 100,000 program/erase (P/E) cycles on a page, the estimated RBER 609 changes in a range $[9 \times 10^{-6} \div 3.5 \times 10^{-4}]$ With a target UBER of 610 10^{-13} , which is typical for commercial applications [? (, according to 611 equation (??) we need to design a code with correction capability in the 612 range $t_{min} = 5$ up to $t_{\mathcal{M}} = 24$. Since $k = 2^{1}$ and $t_{\mathcal{M}} = 24$, from the 613 expression $k + m \cdot t_M \leq 2^m$ Twe deduce m = 15, thus obtaining a maximum 614 of $r = m \cdot t_M \simeq 45$ B of parity information. Given the 8-bit I/O interface of 615 the memory both the encoder and the decoder have been designed with an 616 input parallelism of s = 8 bits. The values of h and z of the Chien Machine 617 are a trade-off between the complexity of the decoder and the decoding time. 618 Given the I/O parallelism of the flash and the area optimizations of Fig. ??, 619 we opted for a Chien machine with parallelism h = 8 and z = 1 full GF 620 multipliers. 621 622

in Table ??.

Arch. 1 is classic BCH architecture with fixed correction capability of 624 24 errors per page. It represents the reference to compare our adaptable 625 architectures. 626

Arch. 2 is an adaptable architecture with $t_{min} = 5 < t \leq 24$ using 627 a traditional PPLFSR for the encoder and 24 PLFSRs for the syndrome 628 calculation. It is worth mentioning here that, differently from what reported ⁶²⁹ in the previous sections, the minimum required correction capability of the ⁶³⁰ codec is higher than 1. This allows us to save space in the encoder PPLFSR ⁶³¹ since less polynomials must be stored, and in the Chien Machine's ROM ⁶³² since less β_{ini} terms must be stored.

Arch. 3 is an optimized version of Arch. 2 exploiting the use of a shop-PLFSR shared between the encoder and the decoder, to trade-off design 635 complexity and decoding time. In order to optimize the use of the shOP-636 PLFSR, we exploited the algorithm proposed in Section ??. Given our adapt-637 able BCH code, a set of ad-hoc Matlab simulation scripts implement this 638 preliminary analysis of 1,800¹ set Ω_i of polynomials. Each set Ω_i contains 639 $t_M - t_{min} - 1 = 20$ generator polynomials required in the encoder and $t_M = 24$ 640 minimal polynomials required in the decoder. This analysis aimed at finding 641 the most suitable set of shared generator and minimal polynomials to trade-642 A reasonable trade-off has been off between decodep's area and latency. 643 found using a shopples composed of $\mathbb{N} =$ 5 OPPLFSRs, each of which 644 dividing by the following set of polynomials: $\{g_5, \psi_{29}, \psi_{39}\}, \{g_6, \psi_{31}, \psi_{41}\}, \{g_6, \psi_{31}, \psi_{31}, \psi_{41}\}, \{g_6, \psi_{31}, \psi_{31}, \psi_{31}\}, \{g_6, \psi_{31}, \psi_{31}, \psi_{31}, \psi_{31}, \psi_{31}\}, \{g_6, \psi_{31}, \psi_{31}, \psi_{31}, \psi_{31}, \psi_{31}, \psi_{32}, \psi_{32}, \psi_{33}, \psi_{33$ 645 $\{g_7, \psi_{33}, \psi_{43}\}$, $\{g_8, \psi_{35}, \psi_{45}\}$, and $\{g_9, \dots, g_{24}, \psi_{37}, \psi_{47}\}$. The reader may refer 646 o the appendix of this paper for the full list of employed polynomials. All 647 other structures remain almost unchanged. The comparison between Arch.1 648 and Arch. 2 enables to highlight the benefits of using an adaptable codec, 649 while the comparison between Arch. 2 and Arch. 3 shows the advantages of 650 adding optimized shared blocks. 651

¹our BCH code has 1,800 primitive polynomials $\psi_1(x)$

Tab	Table 3: Characteristics of the analyzed architectures							
	Adaptable	OPPLFSRs	Chien Machine					
Arch. 1	No	-	h = 8, t = 24					
Arch. 2	Yes	-	$h = 8, t \in [5, 24]$					
Arch. 3	Yes	5	$h = 8, t \in [5, 24]$					
		~ / /						

652

7.3. Performance evaluations

Table ?? summarizes the main implementation details of the three selected architectures in terms of required parity bits and worst case encoding/decoding latency, expressed in terms of clock cycles. 655

Let us start with the evaluation of the amount of redundancy introduced 656 which has a fixed correction capability by the two architectures. Arch. 657 to store $m \cdot t_M$ 360 parity bits of 24 errors per page, requires 15658 (about 45B) for each 2KB page of the flash. This accounts for about 70% of 659 the full spare area available for each page. Since the spare area cannot be 660 fully reserved for storing ECC information (high-level functions, such as file 661 system management and wear-leveling need to save considerable amount of 662 in smation in this area), this percentage represents a considerable overhead 663 for the selected device. Based on the results of Table ??, Fig. ?? shows how, 664 for the adaptable codecs of both Arch. 2 and Arch. 3, the percentage of spare 665 area dedicated for storing parity bits changes with the selected correction 666 capability. The total occupation ranges in this case from 15% to 70% of the 667 total spare area. This mitigates the overhead for storing parity bits whenever 668 the error rate enables to select low correction capabilities (e.g., for devices in 669





Figure 13: Percentage of spare area dedicated to parity bits while changing the correction capability of the adaptable codes of Arch. 2 and Arch. 3

For all implementations, the encoding latency depends on the size of the 671 incoming message and is therefore constant regardless the adaptability of the 672 (see Table???). The decoding latency is instead influenced by the encoder 673 correction capability, as reported in Table ??. Fig. ?? compares the decoding 674 lat acy of the three architectures for each considered correction capability. 675 Results are provided in number of clock cycles. It is worth mentioning here 676 that timing estimations of Table ?? and Fig. ?? depict the worst-case sce-677 nario in which the Chien Machine must search all possible positions prior to 678 find the detected number of errors. Fig. ?? highlights that, for the lowest 679 correction capability, both Arch. 2 and Arch. 3 enable 22% of decoding time 680 reduction when compared to the fixed decoding time of Arch. 1. The decod-681

ing time increases with the correction capability. For Arch. 2, it reaches the 682 same level of the fixed architecture when the correction capability reaches 683 t = 24. Arch. 3 deviates from this behavior for $t \ge 20$. This penalty is intro-684 duced by the use of the shOPPLFSR in the Syndrome Machine. In this case, 685 the codec includes 5 blocks to perform remainder computation with 10 min-686 imal polynomials $\{\psi_{29}, \psi_{39}, \psi_{31}, \psi_{41}, \psi_{33}, \psi_{43}, \psi_{35}, \psi_{45}, \psi_{37}, \psi_{47}\}$. This implies 687 doubling the syndrome computation time every time the required dorrection 688 capability reaches a level in which all these polynomials must be used. Nev-689 ertheless, we will show that this reduced performance is counterbalanced by 690 a reduced area overhead. 691



Figure 14: Worst case decoding latency for the three architectures considered.

7.4. Synthesis Results

Synopsys Design Vision and a CORE 45nm technology cell library have ⁶⁹³ been exploited to synthesize the designs. Table ?? shows the results of the ⁶⁹⁴

692

synthesis of the three architectures. The hardware structures required to 695 obtain the adaptability of the code introduce a certain area overhead. Con-696 sidering Arch. 2, the area of the encoder increases since 19 generator poly-697 nomials must be stored in its ROM, while the area of the decoder increases 698 due both to the aligners in the syndrome machine and to the ROM in the 699 Chien machine to adapt the fast skipping process. Nevertheless, the intro-700 duced overhead is about 14% which is still acceptable. Considering Arch) 3, 701 the introduced overhead is halved w. A. Arch. 2. The area of the encoder is 702 almost comparable with Arch. 2 However, it now includes the shOPPLFSR 703 and a smaller ROMs which contribute, with the LASR sharing, at decreasing 704 the area of the decoder. For both architectures we obtained a maximum clock 705 frequency of 100 MHz, which confirms that the adaptability does not impact 706 the maximum speed of the circuit. This area result is interesting if compared 707 with an estimation of the area for the adaptable architecture proposed in [? 708 designed a code working on blocks of data of 512B, smaller than |. |? 709 the 2KB used in this paper. Given the same maximum correction capability 710] uses a code defined on $CE(2^{13})$ instead of the code defined $(t_M = 24),$ 711 on $GF(2^{15})$ used in this paper. However, even if the code is simpler and the 712 hy ober of correction modes is smaller (only 4 correction modes), the area of 713 the codec accounts about 158.9K equivalent gates², which is higher than the 714 111.4K and the 105.2K equivalent gates of the Arch. 2 and Arch. 3 proposed. 715

Fig. ?? compares the decoder's dynamic power dissipation of the three 716 architectures computed using Synopsys PrimeTime. As for the decoding 717

²Equivalent gates for state-of-the-art architectures have been estimated from the information provided in the papers

		Table 5: Synthe	sis Results	
	Comp.	Max Clock	Equiv. Gates	Over-head
	Encoder	100 MHz	33.3 K	
Arch. 1	Decoder	$100 \mathrm{~MHz}$	64.1 K	
	Overall	$100 \mathrm{~MHz}$	97.4 K	(ref.)
	Encoder	$100 \mathrm{~MHz}$	40.8 K	\$
Arch. 2	Decoder	100 MHz	70.6 K	C
	Overall	100 MHz	111.4 K	14%
	Encoder	100 MHz	39.2 K	\diamond
Arch. 3	Decoder	160 MHz	66.0 K	
\square	Overall	100 MHz)) 105.2 K	7%
~ 1				

latency the analysis has been performed for a worst-case simulation in which 718 t errors are injected at the end of the codeword so that the Chien Machine 719 must search all possible positions prior to detect all errors. Considering Arch. 720 2, results show that the introduction of the adaptability enables up to 15% of 721 dynamic power saving when the lowest correction capability can be selected. 722 TT is due to the fact that the portions of the circuits not required for low 723 correction capabilities are disabled. The introduction of the optimizations 724 proposed in Arch. 3 has no significant impact on the dynamic power that 725 remains almost equal to the one of Arch. 2. 726



Figure 15: Worst case dynamic power consumption of the three decoders for the three considered architectures. Power is expressed in mW.

8. Conclusions

This paper proposed a BCH codec architectures and its related automatic generation framework which enables its code correction capability to be selected in a predefined range of values. Designing an ECC system whose correction capability can be modified in-field has the potentiality to adapt the correction schema to the reliability requirements the flash encounters during its life-time, thus maximizing performance and reliability. 728

727

Experimental results on a selected NAND flash memory architecture 734 proved that the proposed solution reduces spare area usage, decoding time, 735 and power dissipation whenever small correction capability can be selected. 736

Table 0. Willinia polynolials expressed with the corresponding nexace									
coefficients		5	\sqrt{C}	$\widetilde{\mathbf{D}}$	~	\otimes			
	ψ_1	0x F465	1217	0x B13D	-\$¢33	9x 8011			
	ψ_3	0x C209	ψ_{19}	9x B305	ψ_{35}	0x BA2B			
~ {	₩5	Ox B3B7	ψ_{21}	0x A495	ψ_{37}	$0 \ge 0 $ D95F			
$\langle h \rangle$	442	0x E6EB	ψ_{23}	0x 88C7	ψ_{39}	Qx BFF5			
\backslash	ψ_9	0x E647	\$25	0x C357	ψ ₄₁	0x BA87			
	ψ_{11}	0x D4E5	ψ_{27}	0x B2C1	<i>¥</i> 43	∫ox 9BEB			
\sim	ψ_{13}	0x 8371	ψ_{29}	0x 97DD	ψ_{45}	0x 93CB			
(15	0x EDD9	ψ_{31}	0x FA49	ψ_{47}	0x F385			
R	\mathcal{D}								
V									

Table 6: Minimal polynomials expressed with the corresponding hexadecimal string of

 \diamond

fficiente 4 -4 H C ŀ Tahle