



Departament d'Enginyeria  
Telemàtica



UNIVERSITAT POLITÈCNICA DE CATALUNYA

## PhD Research Plan

---

# LOCALLY RECOVERABLE CODES

---

Petar Hlad Colic

Information Security Group  
Department of Network Engineering  
Universitat Politècnica de Catalunya

Advisor: Marcel Fernández Muñoz

Barcelona, July 2018



---

# CONTENTS

<b>1</b>	<b>Context and motivation</b>	<b>1</b>
<b>2</b>	<b>Coding Theory Preliminaries</b>	<b>3</b>
2.1	Linear codes . . . . .	3
<b>3</b>	<b>State of the Art</b>	<b>5</b>
3.1	Definition of LRC codes . . . . .	5
3.2	Bounds on parameters of LRC codes . . . . .	6
3.3	Algebraic Geometric Codes . . . . .	6
3.4	Cyclic and Binary LRC codes . . . . .	7
3.5	List decoding . . . . .	7
<b>4</b>	<b>Work Plan</b>	<b>9</b>
4.1	Goals . . . . .	9
4.2	Open Problems . . . . .	9
4.3	Time Plan and Methodology . . . . .	10
	<b>Bibliography</b>	<b>13</b>



---

---

# CHAPTER 1

---

## CONTEXT AND MOTIVATION

In recent years the explosion in the volumes of data being stored online has resulted in distributed storage systems transitioning to erasure coding based schemes in order to ensure reliability with low storage overheads. On such a massive scale, unreachable or failed servers are no longer an exception but a regular occurrence and recovery from such events has to be done efficiently.

Standard levels of RAID:

**RAID 1:** Consists of data mirroring. Data is written identically to  $N$  drives.

**Storage overhead:** With  $N$  times replication, the storage overhead is  $N$ .

**Failure Tolerance:** Can tolerate up to  $N - 1$  failed drives.

**Repair Procedure:** Need to read only a single drive to repair a failed one.

**RAID 5:** Consists of block-level striping with distributed parity. Parity information is distributed among the drive.

**Storage overhead:**  $\frac{N}{N-1}$

**Failure Tolerance:** Can tolerate only a single failed drive.

**Repair Procedure:** Need to read all  $N-1$  drives to repair failed one.

**RAID 6:** Consists of block-level striping with two distributed parities.

**Storage overhead:**  $\frac{N}{N-2}$

**Failure Tolerance:** Can tolerate any two failed drives.

**Repair Procedure:** Need to read all remaining  $N-2$  drives to repair two failed drives (or single one).

Other data structures use Maximum Distance Separable (MDS) codes, which are those with the greatest error correcting capability.

$[n, k]$  **MDS codes:**

**Storage overhead:**  $\frac{n}{k}$

**Failure Tolerance:** Can tolerate any  $n - k$  disk failures

**Repair Procedure:** Need to read  $k$  drives to repair a failed one.

All these classical erasure correcting codes guarantee that data can be recovered if a bounded number of codeword coordinates is erased, by accessing the remaining coordinates. However, all the described solutions handle very poorly the recovery of a single drive failure since it typically involves accessing large amount of coordinates. This is what is called the Repair Problem, to recover a failed drive with the minimum amount of resources.

In recent years Locally Recoverable Codes (LRC) emerged as the codes of choice for many such scenarios, as they solve very well the repair problem. They have been implemented in a number of large scale systems (see [9] and [14]). LRC codes have the property that a symbol of the codeword can be recovering accessing few other symbols of the codeword (called the *recovering set*).

Symbols can have more than one recovering set, and having more than one recovering set is beneficial in practice because it enables more users to access a given portion of data, thus enhancing data availability in the system.

Data storage applications require codes with small redundancy, low locality for information coordinates, large distance, and low locality for parity coordinates.

---

# CHAPTER 2

---

## CODING THEORY PRELIMINARIES

### 2.1 Linear codes

Let  $\mathbb{F}_q$  be the finite field with  $q$  elements, and consider  $\mathbb{F}_q^n$  the vector space of dimension  $n$  over  $\mathbb{F}_q$ .

**Definition 2.1** (Linear Code). An  $(n, k)$  linear code over  $\mathbb{F}_q$  is a subspace  $\mathcal{C} \subset \mathbb{F}_q^n$  with  $\dim(\mathcal{C}) = k$ .  $\mathcal{C}$  is said to have length  $n$  and dimension  $k$ . Every element  $c \in \mathcal{C}$  is called a *codeword*.

Consider an  $(n, k)$  linear code  $\mathcal{C}$  over  $\mathbb{F}_q$ , and let  $\mathbf{x} = x_1 \dots x_n$ ,  $\mathbf{y} = y_1 \dots y_n$  two codewords of  $\mathcal{C}$ .

**Definition 2.2** (Hamming Distance). The Hamming distance between two vectors  $x = x_1 \dots x_n$  and  $y = y_1 \dots y_n$  is the number of coordinates where they differ, and is denoted by  $\text{dist}(x, y)$ .

**Definition 2.3** (Hamming Weight). The Hamming Weight of a vector  $x = x_1 \dots x_n$  is the number of nonzero coordinates of  $x$  and is denoted by  $\text{wt}(x)$ .

**Remark 2.4.** For a linear code  $\mathcal{C}$  and any two codewords  $x, y \in \mathcal{C}$ :

$$\text{dist}(x, y) = \text{wt}(x - y) \quad (2.1)$$

**Definition 2.5** (Minimum distance). A code  $\mathcal{C}$  has minimum distance  $d$  if any two codewords differ in at least  $d$  coordinates.

$$d = \min \text{dist}(x, y) = \min \text{wt}(x - y), \quad x, y \in \mathcal{C}, x \neq y \quad (2.2)$$

A linear code of length  $n$ , dimension  $k$ , and minimum distance  $d$  will be called an  $[n, k, d]$  linear code.

There are several known and well studied bounds on the size of a code considering its length and distance. One that will be very important in this work is the Singleton Bound.

**Theorem 2.6** (Singleton Bound). *If  $\mathcal{C}$  is an  $[n, k, d]$  code, then*

$$d \leq n - k + 1 \quad (2.3)$$

**Definition 2.7** (MDS Codes). A Maximum Distance Separable code  $\mathcal{C}$  is a code that attains the Singleton bound with equality. That is, a code s.t.

$$d = n - k + 1 \quad (2.4)$$





---

# CHAPTER 3

---

## STATE OF THE ART

### 3.1 Definition of LRC codes

Consider a linear  $[n, k, d]_q$  code  $\mathcal{C} \subset \mathbb{F}_q^n$ , where  $q$  is a prime power. We say that the  $i$ -th coordinate of  $\mathcal{C}$  has locality  $r$ , if the value at this coordinate can be recovered from accessing some other  $r$  coordinates of  $\mathcal{C}$ . We say that the code  $\mathcal{C}$  has locality  $r$  if every symbol of the codeword  $x \in \mathcal{C}$  can be recovered from a subset of  $r$  other symbols of  $x$ .

**Definition 3.1** (LRC Codes). A code  $\mathcal{C} \subset \mathbb{F}_q^n$  is a *locally recoverable code* (LRC) with locality  $r$  if for every  $i \in [n]$  there exists a subset  $\mathcal{R}_i \subset [n] \setminus \{i\}$ ,  $|\mathcal{R}_i| \leq r$  and a map  $\phi_i$  such that for every codeword  $\mathbf{x} \in \mathcal{C}$  we have

$$\mathbf{x}_i = \phi_i(\{\mathbf{x}_j, j \in \mathcal{R}_i\}) \quad (3.1)$$

This definition can be also rephrased as follows. Given  $a \in \mathbb{F}_q$  consider the sets of codewords

$$\mathcal{C}(i, a) = \{x \in \mathcal{C} : x_i = a\}, \quad i \in [n]$$

The code  $\mathcal{C}$  is said to have locality  $r$  if for every  $i \in [n]$  there exists a subset  $\mathcal{R}_i \subset [n] \setminus \{i\}$ ,  $|\mathcal{R}_i| \leq r$  such that the restrictions of the sets  $\mathcal{C}(i, a)$  to the coordinates in  $\mathcal{R}_i$  for different  $a$  are disjoint:

$$\mathcal{C}_{I_i}(i, a) \cap \mathcal{C}_{I_i}(i, a') = \emptyset, \quad a \neq a' \quad (3.2)$$

The subset  $I_i$  is called a *recovering set* for the symbol  $x_i$ .

**Definition 3.2** (t-LRC Codes). A code  $\mathcal{C}$  is said to have  $t$  disjoint recovering sets if for every  $i \in [n]$  there are pairwise disjoint subsets  $R_i^1, \dots, R_i^t \subset [n] \setminus \{i\}$  such that for all  $j = 1, \dots, t$  and every pair of symbols  $a, a' \in \mathbb{F}_q$ ,  $a \neq a'$

$$\mathcal{C}(i, a)_{R_i^j} \cap \mathcal{C}(i, a')_{R_i^j} = \emptyset \quad (3.3)$$

For linear LRC codes, the relation between a symbol  $i$  and its recovering set  $I_i$  is linear. Thus, any symbol in  $I_i \cup \{i\}$  can be recovered from the remaining symbols. We then call  $I_i \cup \{i\}$  a *repair group*.

## 3.2 Bounds on parameters of LRC codes

Gopalan et al. proved in [7] the following bounds:

**Theorem 3.3.** *Let  $\mathcal{C}$  be an  $(n, k, r)$  LRC code. The rate of  $\mathcal{C}$  satisfies*

$$\frac{k}{n} \leq \frac{r}{r+1} \quad (3.4)$$

*The minimum distance of  $\mathcal{C}$  satisfies:*

$$d \leq n - k - \left\lceil \frac{k}{r} \right\rceil + 2 \quad (3.5)$$

**Theorem 3.4** ([13, 21]). *For  $(n, k, r, t)$  LRC codes with  $t \geq 2$  disjoint recovering sets:*

$$d \leq n - k + 2 - \left\lceil \frac{t(k-1) + 1}{t(r-1) + 1} \right\rceil \quad (3.6)$$

We will refer to codes attaining the bound 3.5 (the bound 3.6 in case  $t \geq 2$ ) as optimal LRC codes.

In [17], Tamo, Barg, and Frolov found new bounds on the distance and rate of LRC codes as well as asymptotic bounds.

**Theorem 3.5.** *Let  $\mathcal{C}$  be an  $(n, k, r, t)$  LRC code with  $t$  disjoint recovering sets of size  $r$ . Then the rate of  $\mathcal{C}$  satisfies*

$$\frac{k}{n} \leq \frac{1}{\prod_{j=1}^t (1 + \frac{1}{jr})} \quad (3.7)$$

*The minimum distance of  $\mathcal{C}$  is bounded above as follows:*

$$d \leq n - \sum_{i=0}^t \left\lfloor \frac{k-1}{r^i} \right\rfloor \quad (3.8)$$

$$R_q(r, \delta) \geq 1 - \min_{0 < s \leq 1} \left\{ \frac{1}{r+1} \log_q((1 + (q-1)s)^{r+1} + (q-1)(1-s)^{r+1}) - \delta \log_q s \right\} \quad (3.9)$$

## 3.3 Algebraic Geometric Codes

A family of optimal LRC codes was described by Tamo and Barg in [16] which are subcodes of Reed-Solomon codes, considering a polynomial that is constant on each part of a partition of the evaluation points set. The length of these codes is upper

bounded by  $q$  (same as RS codes). Then an algebraic geometric approach of an optimal LRC code is obtained from maps of degree  $r + 1$  from  $\mathbb{P}_{\mathbb{F}_q}^1$  to  $\mathbb{P}_{\mathbb{F}_q}^1$ .

In [2], the authors generalized this idea and constructed LRC codes from morphisms on algebraic curves (Hermitian curves and Garcia-Stichtenoth curves). In [3] the authors expand on the constructions of [2] to produce families of LRC codes coming from a larger variety of curves, as well as from higher-dimensional varieties. They construct optimal LRC codes with code length larger than the size of the alphabet (e.g.  $n = q^2 + 2$  and  $n = q^2 - 1$ ).

### 3.4 Cyclic and Binary LRC codes

Binary error correcting codes are of special interest due to practical reasons.

In [11] authors show how the presence of locality within a binary cyclic code can be exploited to improve decoding performance and reduce decoding complexity. They approach the problem with ordered statistics decoding (OSD) method and with trellis decoding.

In [19], authors consider linear cyclic codes with the locality property. They focus on optimal cyclic codes that arise from the construction of RS like LRC codes in [16], and give a characterization of these codes in terms of their zeros, and observe that there are many equivalent ways of constructing optimal cyclic LRC codes over a given field.

### 3.5 List decoding

A code of length  $n$  is called  $(\tau, \ell)$ -list decodable if the Hamming sphere of radius  $\tau$  centered at any vector  $v$  of length  $n$  always contains at most  $\ell$  codewords  $c \in \mathcal{C}$ . It was shown by Johnson in [10] that any code of length  $n$  and distance  $d$  is  $(\tau_J, \ell)$ -list decodable where  $\tau_J = n - \sqrt{n(n-d)}$  is the Johnson radius and  $\ell \in \text{poly}(n)$ .

It was recently shown by Holzbaur and Wachter-Zeh in [8] that the list decoding radius of certain LRC codes exceed the Johnson radius and give a general list decoding algorithm. The complexity of the algorithm is polynomial in  $n$  when the number of repairing groups is constant, otherwise it grows exponentially.



---

# CHAPTER 4

---

## WORK PLAN

### 4.1 Goals

The goal of this PhD program is to make contributions in the field of Coding Theory, specifically in Locally Recoverable Codes. To do it, a number of open problems are detailed in section 4.2, that will be studied.

Alexander Barg is one of the main contributors in the topic of LRC codes. When discussing with him about open problems in the field, he said that the main problems have already been worked on, and the remaining problems are non-trivial and hard. Therefore, specific problems will be studied in this PhD thesis.

### 4.2 Open Problems

**Improvement of binary LRC codes decoding:** Following the work in [11], decoding of binary LRC codes that are not cyclic will be studied.

**New constructions of LRC codes on algebraic curves:** Following the work in [2] and [3], new constructions of LRC codes over algebraic varieties will be searched, looking for optimal LRC codes with small field size.

**Improvement of bounds for LRC- $t$  codes:** The lower bound 3.7 appears to be far from tight. Tamo and Barg in [18] said they believe that the rate  $\left(\frac{r}{r+1}\right)^t$  is the largest possible for a LRC- $t$  code as long as  $t$  is not too large (e.g.  $t \in O(\log n)$ ). This rate can be achieved constructing a  $t$ -fold power of the binary  $(r+1, r)$  single-parity-check code.

Theorem 3.5 is proved applying probabilistic method techniques on the properties of a graph. The problem of optimizing the bound of the rate of LRC- $t$  codes will be studied, and one of the ways could be following a similar proof considering some restrictions that were not considered in [17].

**List decoding of LRC codes:** the problem will be studied to search for new families of LRC codes that could be list decoded beyond the Johnson radius.

## 4.3 Time Plan and Methodology

The work and research on the open problems is planned for a temporal span of 3 years.

**First year:** Initial research and open problems statement

**Second year:** Research on stated problems

**Third year:** Writing

### First Year

During the first year several activities have been done to determine the state of the art of Locally Recoverable Codes, and state the problems that will be the object of study in this PhD thesis.

**Master's Degree** Obtention of Master's Degree in Advanced Mathematics and Mathematical Engineering, attending specific courses related to the problems that will be worked on. Courses: Coding Theory, Commutative Algebra, Algebraic Geometry, Combinatorics and Graph Theory.

**Stays** Two week stay in University of Maryland as a visitor student invited by Prof. Alexander Barg. In Those two weeks collaborating with Prof. Barg, a better understanding on the state of the art and which are the interesting problems to solve.

**Seminars** Algebraic Geometry Seminar, two sessions per week. Following:

- Book on algebraic curves [4]
- Lecture notes on algebraic geometry and algebraic curves [5, 6]

**Self-Learning** Followed:

- Book on Probabilistic Method [1]

### Second Year

During the second year, the main focus will be on the stated problems.

### Seminars

- Algebraic Geometric Codes, two sessions per week. Following: Books on Algebraic Geometric Codes [15, 20]
- Probabilistic method, two sessions per week. Following: Book on Probabilistic Method [1]

### Self-Learning Will follow:

- Book on Coding Theory [12]

### Third Year

During the third year, the goals will be to conclude with the research on the problems and to write the PhD thesis.





---

# BIBLIOGRAPHY

- [1] Noga Alon and Joel H. Spencer. *The Probabilistic Method*. 3rd. Wiley Publishing, 2008 (cit. on pp. 10, 11).
- [2] A. Barg, I. Tamo, and S. Vladuts. “Locally Recoverable Codes on Algebraic Curves”. In: *IEEE Transactions on Information Theory* 63.8 (2017), pp. 4928–4939. ISSN: 0018-9448. DOI: [10.1109/TIT.2017.2700859](https://doi.org/10.1109/TIT.2017.2700859) (cit. on pp. 7, 9).
- [3] Alexander Barg, Kathryn Haymaker, Everett W. Howe, Gretchen L. Matthews, and Anthony Várilly-Alvarado. “Locally Recoverable Codes from Algebraic Curves and Surfaces”. In: *Algebraic Geometry for Coding Theory and Cryptography*. Ed. by Everett W. Howe, Kristin E. Lauter, and Judy L. Walker. Cham: Springer International Publishing, 2017, pp. 95–127. ISBN: 978-3-319-63931-4 (cit. on pp. 7, 9).
- [4] William Fulton. *Algebraic Curves*. 1969 (cit. on p. 10).
- [5] Andreas Gathmann. *Algebraic Geometry*. 2014. URL: <http://www.mathematik.uni-kl.de/~gathmann/class/alggeom-2014/alggeom-2014.pdf> (cit. on p. 10).
- [6] Andreas Gathmann. *Plane Algebraic Curves*. 2018. URL: <http://www.mathematik.uni-kl.de/~gathmann/class/curves-2018/curves-2018.pdf> (cit. on p. 10).
- [7] P. Gopalan, C. Huang, H. Simitci, and S. Yekhanin. “On the Locality of Codeword Symbols”. In: *IEEE Transactions on Information Theory* 58.11 (2012), pp. 6925–6934. ISSN: 0018-9448. DOI: [10.1109/TIT.2012.2208937](https://doi.org/10.1109/TIT.2012.2208937) (cit. on p. 6).
- [8] Lukas Holzbaur and Antonia Wachter-Zeh. “List Decoding of Locally Repairable Codes”. In: *CoRR* abs/1801.04229 (2018). arXiv: [1801.04229](https://arxiv.org/abs/1801.04229). URL: <http://arxiv.org/abs/1801.04229> (cit. on p. 7).

- [9] Cheng Huang, Huseyin Simitci, Yikang Xu, Aaron Ogus, Brad Calder, Parikshit Gopalan, Jin Li, and Sergey Yekhanin. “Erasure Coding in Windows Azure Storage”. In: *Presented as part of the 2012 USENIX Annual Technical Conference (USENIX ATC 12)*. Boston, MA: USENIX, 2012, pp. 15–26. ISBN: 978-931971-93-5. URL: <https://www.usenix.org/conference/atc12/technical-sessions/presentation/huang> (cit. on p. 2).
- [10] S. Johnson. “A new upper bound for error-correcting codes”. In: *IRE Transactions on Information Theory* 8.3 (1962), pp. 203–207. ISSN: 0096-1000. DOI: [10.1109/TIT.1962.1057714](https://doi.org/10.1109/TIT.1962.1057714) (cit. on p. 7).
- [11] M. Nikhil Krishnan, Bhagyashree Puranik, P. Vijay Kumar, Itzhak Tamo, and Alexander Barg. “Exploiting Locality for Improved Decoding of Binary Cyclic Codes”. In: *IEEE Trans. Communications* 66.6 (2018), pp. 2346–2358. DOI: [10.1109/TCOMM.2018.2797988](https://doi.org/10.1109/TCOMM.2018.2797988). URL: <https://doi.org/10.1109/TCOMM.2018.2797988> (cit. on pp. 7, 9).
- [12] F.J. MacWilliams and N.J.A. Sloane. *The Theory of Error-Correcting Codes*. Vol. 16. North Holland, 1983. ISBN: 9780444851932 (cit. on p. 11).
- [13] A. S. Rawat, D. S. Papailiopoulos, A. G. Dimakis, and S. Vishwanath. “Locality and Availability in Distributed Storage”. In: *IEEE Transactions on Information Theory* 62.8 (2016), pp. 4481–4493. ISSN: 0018-9448. DOI: [10.1109/TIT.2016.2524510](https://doi.org/10.1109/TIT.2016.2524510) (cit. on p. 6).
- [14] Maheswaran Sathiamoorthy, Megasthenis Asteris, Dimitris Papailiopoulos, Alexandros G. Dimakis, Ramkumar Vadali, Scott Chen, and Dhruba Borthakur. “XOR-ing Elephants: Novel Erasure Codes for Big Data”. In: *Proc. VLDB Endow.* 6.5 (Mar. 2013), pp. 325–336. ISSN: 2150-8097. DOI: [10.14778/2535573.2488339](https://doi.org/10.14778/2535573.2488339). URL: <http://dx.doi.org/10.14778/2535573.2488339> (cit. on p. 2).
- [15] Henning Stichtenoth. *Algebraic Function Fields and Codes*. Springer-Verlag Berlin Heidelberg, 2009. DOI: [10.1007/978-3-540-76878-4](https://doi.org/10.1007/978-3-540-76878-4) (cit. on p. 11).
- [16] I. Tamo and A. Barg. “A Family of Optimal Locally Recoverable Codes”. In: *IEEE Transactions on Information Theory* 60.8 (2014), pp. 4661–4676. ISSN: 0018-9448. DOI: [10.1109/TIT.2014.2321280](https://doi.org/10.1109/TIT.2014.2321280) (cit. on pp. 6, 7).
- [17] I. Tamo, A. Barg, and A. Frolov. “Bounds on the Parameters of Locally Recoverable Codes”. In: *IEEE Transactions on Information Theory* 62.6 (2016), pp. 3070–3083. ISSN: 0018-9448. DOI: [10.1109/TIT.2016.2518663](https://doi.org/10.1109/TIT.2016.2518663) (cit. on pp. 6, 9).
- [18] Itzhak Tamo and Alexander Barg. “Bounds on Locally Recoverable Codes with Multiple Recovering Sets”. In: *CoRR* abs/1402.0916 (2014). arXiv: [1402.0916](https://arxiv.org/abs/1402.0916). URL: <http://arxiv.org/abs/1402.0916> (cit. on p. 9).
- [19] Itzhak Tamo, Alexander Barg, Sreechakra Goparaju, and Robert Calderbank. “Cyclic LRC Codes, Binary LRC Codes, and Upper Bounds on the Distance of Cyclic Codes”. In: *Int. J. Inf. Coding Theory* 3.4 (Jan. 2016), pp. 345–364. ISSN: 1753-7703. DOI: [10.1504/IJICOT.2016.079496](https://doi.org/10.1504/IJICOT.2016.079496). URL: <https://doi.org/10.1504/IJICOT.2016.079496> (cit. on p. 7).

- [20] Serge Vladut, Dmitry Nogin, and Michael Tsfasman. *Algebraic Geometric Codes: Basic Notions*. Boston, MA, USA: American Mathematical Society, 2007. ISBN: 0821843060, 9780821843062 (cit. on p. 11).
- [21] A. Wang and Z. Zhang. “Repair Locality With Multiple Erasure Tolerance”. In: *IEEE Transactions on Information Theory* 60.11 (2014), pp. 6979–6987. ISSN: 0018-9448. DOI: [10.1109/TIT.2014.2351404](https://doi.org/10.1109/TIT.2014.2351404) (cit. on p. 6).