

TITLE

Subtitle

Seminar Datenmanagement

Prof. Dr. Lena Wiese

Semester 7



Institute of Computer Science

Goethe-Universität Frankfurt a. M.

Author: NAKO NACHEV

6464679

s6327063@stud.uni-frankfurt.de

Degree: Master Informatik, Semester 7

Module ID: DM-MS, M-DS-S, DM-BS, B-ATAI-S

Date: January 21, 2024

Work based on:

- Paper 1: *Put Original Author Names here* Put Original Title(s) here
- Paper 2: *Put Original Author Names here* Put Original Title(s) here
- ...

Abstract

Contents

1	Introduction	4
1.1	Summary Overview Paper	4
1.2	Important Terms and Definitions	5
2	In-Depth Explanation of Error-correcting Code Chosen	6
3	Implementation of Error-correcting Code Chosen	9
3.1	DNA Fountain implementation	9
3.2	Benchmark Results	11
4	Conclusion	11

1 Introduction

1.1 Summary Overview Paper

The paper in [1] provides a detailed survey about different error correcting codes, their implementation and associated obstacles with each of the codes. The paper begins by examining the historical significance of error synchronisation problems in communication systems dating back to Samuel Morse, which remain an important topic even in today's technological systems. Although often treated as different problems, synchronisation and additive noise tend to have the same effect on communication channels by reducing their capacity. According to the article, designing error-correcting codes to tackle these problems proves to be quite challenging, with many techniques implemented not based on coding or even failing to mention synchronisation-correcting codes.

This article continues with a brief introduction on the error-correcting codes and explaining the different preliminaries associated with the survey. Such preliminaries include the mathematical notations used, the different types of synchronisation errors and their definitions. These range from deletion errors, insertion errors, duplication or repetition errors to bitshift errors. Last but not least, obstacles associated with the error correcting codes are described. Examples of such are huge bursts of substitution errors, splitting long messages of data into blocks of with block boundaries sometimes being unknown to receiver or codes being unable to keep the same correction capability when encoding several blocks. The authors describe in section III different synchronisation-correcting codes, how they are implemented, the historic research done by various authors and specific limitations for the different types. The section starts with binary algebraic block codes using $(0,1)$ alphabets, nonbinary and perfect codes, codes capable of correcting bursts of synchronisation errors, synchronizable, marker codes, codes for weak synchronisation errors, convolutional codes, expurgated and codes for random synchronisation channels.

The last section of the paper discusses future directions of synchronisation error-correcting codes. Binary block codes are currently unable to detect more than one synchronisation errors per block. Further challenges arise in estimating the capacity of channels corrupted by synchronisation errors and such where random deletion and/or insertion errors can occur. According to the authors, none of the codes described in the paper have all the properties

required to be used in practical communication systems. One solution could be concatenating different coding schemes in order to use the strengths of the various codes for their specific implementation use cases.

Deriving from the examples in the article from above, one could divide the error-correcting codes in multiple different categories. One separation used could be based on the number of symbols used in the alphabet - either binary alphabet consisting of $(0,1)^*$ symbols or non-binary with 2 or more symbols. Another categorization could follow based on other characteristics - such as block codes vs convolutional codes. Block codes break the data down to fixed-sized blocks and append protective or unique bits for error handling. Convolutional codes on the other side process the data continuously as a stream and do not split the data into blocks. Error correction here is done on the fly.

As far as characteristics of the error-correcting codes are concerned, performance, redundancy, error detection and correction seem to be the most important ones. Redundancy finds its application in appending the extra bits for the receiver to make use of and identify errors that may have happened during the transmission. Performance concerns the rate at which algorithms are able to correct the data, where similar to most systems a good balance between correctness achieved and computation should be achieved. Error detection and correction seems to be closely related and are the general characteristics that prove how valuable a given error-correcting code is.

There are various implementations of error-correcting codes. One such example is by using error-correcting output coding for improving text classification and thus reducing text classification errors by 66 percent. [2] Another example is by using error-correcting codes for real time communication in wireless networks. [3] The authors simulate using an adaptive scheme for reducing bandwidth overhead by comparing it to a single code scheme.

1.2 Important Terms and Definitions

Upon reading the overview paper, following important terms and definitions could be listed based on the information in [1]:

- synchronisation error - it is either a deletion or an insertion error and excludes substitution errors. For example transmitting 0011 instead of 00011 is a deletion error and transmitting 0011 instead of 011 - insertion error.

- duplication or repetition error - it is a special insertion that replaces a bit by two copies of the same bit.
- bitshift error - it a special kind of error that might reverse a given string - for example 01 being transmitted as 10.
- additive noise - disturbances or interference when transmitting signal over a communication network
- redundancy - additional information introduced by error-correcting codes to correction synchronisation errors
- run - a substring of identical symbols of maximum length
- blocks - used to form chunks of words that are used to divide long messages

2 In-Depth Explanation of Error-correcting Code Chosen

Current storage medias do not have the capacity to meet the demand for the total amount of worldwide data. On top of that, the costs for maintaining and transferring data tend to increase as well, therefore better solutions are needed [4]. One such promising solution are the innovate DNA storage systems. They enable us to store digital data into DNA molecules, which are syntetically syntesized and far exceed the lifespan of general storage mediums [4]. But transferring and storing data comes with its complications and certain errors could appear, leading to corrupted or incomplete data. Thus, solutions are needed for correcting these errors. The storing strategy provided in [5] makes use of the so called DNA fountains. Fountain codes are typically used for channels with erasures, where files are split into packets and each of these packets is either received without errors or dropped [6]. Standard transfer protocols keep transmitting packets until they are successfully received from the other side, where another channel is also required to keep track of the state of transmitted packets or those that require retransmission. Fountain codes on the other side create packets as a result of functions of the whole file and the receiver only has to collect any N packets, with $N > K$ (original size) so that the whole file can be recovered. In the methaphorical sense, a

dna fountain generates droplets/drops (packets) und everyone who wants to collect these drops has to put a bucket under the fountain and wait until the amount collected is larger than original size K of the data [6]. Due to the nature of the fountain, such DNA fountains can be classified as rateless, since the number of packets generated could be limitless.

The strategy for data storage using DNA fountains proposed in [5] splits the encoding process into three steps. The first step is the preprocessing stage, where the binary file provided as input is split into segmets of identical length. A simple example would be splitting the code 001001100101 into segments of length 4, which would result in:

- Segment 1: 0010
- Segment 2: 0110
- Segment 3: 0101

The second step of the stragy is the so-called **luby-transform stage**. The segments generated in the first step are sampled randomly using a robust soliton distribution [7] and than added bitwise together. In the same time, a random seed is generated //TODO: probably add an explanation for the lfsr and after that attached to the result of the bitwise addition, thus forming the final droplet. One simple example might be:

- Choose two random segments: 0010 and 0110
- Apply bitwise addition:

$$\begin{array}{r} 0010 \\ \oplus 0110 \\ \hline 0100 \end{array}$$

- Generate random seed: 11
- Attach seed to the result of the bitwise addition to form the droplet: 110101

The luby transform stage concludes with the screening stage. This stage first converts the achieved droplet to DNA using the following conversion rules: $\{00,01,10,11\}$ to $\{A,C,G,T\}$. After that, the result dna sequence is screened against invalid sequences, such a homopolymer runs or GC content.

Long homopolymer runs and or sequences with high GC concentration are undesirable [5] due to being more difficult to synthesize or are more prone to sequencing errors. DNA sequences containing one of the above described properties are rejected and those that pass added to the oligo design file. Below is the screening stage described with continuation of the example above:

- Convert the droplet to a DNA sequence: $110101 \rightarrow TCC$
- Check against screening rules: TCC does not contain long homopolymer runs or high % GC content, therefore does not get rejected
- The oligo file now contains the droplet TCC

Apart from the seed and payload, the final droplet also needs redundancy in order to be able to correct various errors that may appear during the synthesizing process. This is where Reed-Solomon codes (R-S) come into place. Reed-Solomon codes are nonbinary cyclic codes with symbols up to m-bit, with m being any positive integer having a value greater than 2. R-S (n,k) codes of m-bit symbols comply to the following rules:

$$0 < k < n < 2^m + 2$$

with k being the number of symbols encoded and n the total number of symbols in the encoded block [8]. Furthermore, R-S(n,k) codes also include parity symbols:

$$(n, k) = (2^m - 1, 2^m - 1 - 2t)$$

with $2t = n - k$ being the number of parity symbols [8]. For example a R-S(5,3) 3-bit code contains 5 code word bytes, of which 3 are data and 2 parity (n-k). Therefore, this code is capable of correcting up to $2t = 2 \Rightarrow t = 1$ errors.

- Read papers on your chosen Error-correcting Code.
- Do a search on the Scholar page for more information on the Error-correcting Code.
- What approach is proposed? What is the formal definition of the chosen Error-correcting Code? What types of errors can be corrected/detected by the code?
- What (in the opinion of the article authors) makes the presented Error-correcting Code better than other approaches?

3 Implementation of Error-correcting Code Chosen

- Upload your own implementation of the Error-correcting Code into Moodle.
- Run your implementation on the provided benchmark dataset.
- Describe your implementation and the benchmark results

3.1 DNA Fountain implementation

The suggested in [5] dna fountain method was implemented in python by making use of standard libraries and a couple of extra libraries for ready implementations of the LFSR and RS-Codes. The project contains the following files:

- `dna_fountain.py`
- `decoder.py`
- `helper.py`
- `image_creator.py`
- `models.py`
- `soliton.py`
- `input.py`

The starting point of the project is the `dna_fountain.py` file. The file begins by initializing the parameters needed for the soliton distribution and the encoding process. The first step involves choosing the input data to use. For the purpose of the seminar, two different inputs were predefined (`picture_short` and `picture_long`), both of which are saved in the `input.py` file. Both inputs are Python lists with 0s and 1s, corresponding to the different bits saved from the pictures. Based on the chosen input, different globals and defaults are created. Examples of such are `segment_size`, `initial_seed`, polynomial and `segments_multiplier`. The function continues by first splitting the input list into multiple segments based on the initialized parameters.

If the input file is not divisible by the number of bits per segment due to the chosen segment size, the last segment is filled with 0s on the left side to match the size of the other segments. After the list of segments is created, it is passed to the `robust_soliton_distribution` function from the `soliton.py` file, which returns the distribution of the values based on the total number of segments generated. `Soliton.py` also contains a `ideal_soliton_distribution` function, which is not used, but remains an option.

After all defaults, parameters and values are initialized, the oligo creation process via the `oligo_creation` functions begins. It first defines the total number of oligos to be generated and starts iterating until the desired number is reached. The `luby_transform` function is responsible for the creation of oligos. It first creates a droplet (with the corresponding seed, data payload and RS-Code (if exists)) and checks it against homopolymer runs or high GC content. If any of the conditions are met, the function will not return anything, meaning the iterator has to try again. Otherwise, the function returns the created oligo, and it is passed to the final list. The core of the `luby_transform` function is the `create_droplet_bin` function. It creates the seed by using the LFSR with the predefined parameters (seed, polynomial) and chooses a number N of segments to combine, with the help of the generated distribution from the robust soliton implementation. The seed is then used to generate N indices for random segments of the list, which are XOR-ed together to form the final payload of the droplet. Along the way, different helper functions and models are used. Typical examples of helper functions include binary arrays to bytes, XOR-ing integers or strings or different utility functions, all of which provided in `helper.py` file. `Models.py` contains classes representing different types, which, although not necessary in python due to missing static checking, might improve readability and debugging (if necessary).

Once the oligos are created, they are passed, along with other parameters to the `decoder.py` file. This file starts by preprocessing the oligos. The first two stages of preprocessing, as suggested by [9], involve retaining only sequences of a certain length and collapsing identical sequences and sorting them so that sequences that occur more often are at the front of the list. Since the created oligos are still in string format (due to nucleotides being the letter A,C,G,T), we need to convert the data back to binary. This is achieved via the `create_oligo` function, which also extracts the seed, data payload and potential error correcting codes that were appended at the end of the oligo. The described in [5] process does not attempt to make error corrections

and for that purpose oligos with errors are dropped and do not fall into the decoding process. Last but not least, the `recursively_infer_segments` function is called, which is the final stage of the decoding process. It iterates through the preprocessed oligos and attempts to infer the segments until no more updates can be made. It finally returns the inferred segments and they are passed back to the `dna_fountain.py` file.

The final part of the `dna_fountain.py` file is the image creation. The `image_creator.py` file is called, which contains a function for generating a .PNG picture based on the returned data (list of 0s and 1s). That way we can compare the output picture to the initial picture and compare whether we have successfully decoded the initial file. Since some of the bits may be different, which most of the time is might not be noticable with a human eye, the `calc_similarity` function from the `helper.py` file can compare both input and output lists and return a value between $[0,1]$ indicating how similar the decoded file is to the initial. Value of 1 means the initial file was decoded without any errors or missing information.

3.2 Benchmark Results

4 Conclusion

- Summarize the main points and achievements
- Add your own assessment/criticism on the topic

Some Latex-specific hints:

- You can use abbreviations, like Active Directory (AD), Microsoft (MS) or Compact Disc (CD).
- There are also symbols like for example π , φ and λ .
- Last but not least, use glossary entries like Active Directory and File.
- Do not forget to cite related work, like [10] or [11].

References

- [1] Hugues Mercier, Vijay K Bhargava, and Vahid Tarokh. A survey of error-correcting codes for channels with symbol synchronization errors. *IEEE Communications Surveys & Tutorials*, 12(1):87–96, 2010.
- [2] Rayid Ghani. Using error-correcting codes for text classification. In *ICML*, pages 303–310. Citeseer, 2000.
- [3] Moncef Elaoud and Parameswaran Ramanathan. Adaptive use of error-correcting codes for real-time communication in wireless networks. In *Proceedings. IEEE INFOCOM’98, the Conference on Computer Communications. Seventeenth Annual Joint Conference of the IEEE Computer and Communications Societies. Gateway to the 21st Century (Cat. No. 98*, volume 2, pages 548–555. IEEE, 1998.
- [4] Yiming Dong, Fajia Sun, Zhi Ping, Qi Ouyang, and Long Qian. Dna storage: research landscape and future prospects. *National Science Review*, 7(6):1092–1107, 2020.
- [5] Yaniv Erlich and Dina Zielinski. Dna fountain enables a robust and efficient storage architecture. *science*, 355(6328):950–954, 2017.
- [6] David JC MacKay. Fountain codes. *IEE Proceedings-Communications*, 152(6):1062–1068, 2005.
- [7] Weiqing Yao, Benshun Yi, Taiqi Huang, and Weizhong Li. Poisson robust soliton distribution for lt codes. *IEEE Communications Letters*, 20(8):1499–1502, 2016.
- [8] Bernard Sklar. Reed-solomon codes. *Downloaded from URL <http://www.informit.com/content/images/art.sub.-sklar7.sub.-reed-solomo-n/elementLinks/art.sub.-sklar7.sub.-reed-solomon.pdf>*, pages 1–33, 2001.
- [9] Yaniv Erlich and Dina Zielinski. Supplementary materials for dna fountain enables a robust and efficient storage architecturedna fountain enables a robust and efficient storage architecture. *science*, 355(6328):950–954, 2017.

- [10] Lior Okman, Nurit Gal-Oz, Yaron Gonen, Ehud Gudes, and Jenny Abramov. Security issues in nosql databases. In *10th International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom)*, pages 541–547. IEEE, 2011.
- [11] Dhruba Borthakur, Jonathan Gray, Joydeep Sen Sarma, Kannan Muthukkaruppan, Nicolas Spiegelberg, Hairong Kuang, Karthik Ranganathan, Dmytro Molkov, Aravind Menon, Samuel Rash, et al. Apache hadoop goes realtime at facebook. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management of data*, pages 1071–1080. ACM, 2011.