

A Study of Error Floor Behavior in QC-MDPC Codes

Sarah Arpin, Tyler Billingsley, Daniel Hast, Jun Lau, Ray Perlner,
Angela Robinson

PQCrypto - 2022

September 28-30, 2022

Overview

- 1 Background
 - A primer on QC-MDPC codes
 - BIKE at a high level
 - A primer on error floors
- 2 Our approach
- 3 Analysis of DFR on Special Sets
- 4 Conclusions and Next Steps

Section 1

Background

Motivation

BIKE (**Bit-Flipping Key Encapsulation**) is a code-based KEM (key encapsulation mechanism) based on QC-MDPC (Quasi-Cyclic Moderate-Density Parity-Check) codes.

- Relevance to post-quantum cryptography
 - One of four remaining KEMs in the 4th round of the NIST PQC Standardization process
 - One of three code-based KEMs still under consideration
- IND-CCA security
 - The GJS key-recovery attack exploits decoding failures in an IND-CCA security model
 - Decoding failure rate (DFR) of a code-based KEM that claims IND-CCA security must be sufficiently low to prevent GJS attack

BIKE: Bit flipping key encapsulation - <https://bikesuite.org>

Guo, Johansson, and Stankovski. A Key Recovery Attack on MDPC with CCA Security Using Decoding Errors (2016)

Background

- A **binary linear code** $C = C(n, k)$ is a linear subspace of \mathbb{F}_2^n of dimension k . Vectors in C are called **codewords**.
- A **parity check matrix** of C is a $(n - k) \times n$ matrix H such that for all $v \in V$, we have $Hv^T = 0$ if and only if $v \in C$.
- That is, the rows of H give linear relations satisfied by codewords. Note that H determines C .
- For any $x \in \mathbb{F}_2^n$, Hx^T is called the **syndrome** of x .

QC-MDPC codes

- A **Moderate-Density Parity-Check** code (MDPC code) is a binary linear code $C(n, k)$ that has a parity check matrix H such that each row has weight $w \approx \sqrt{n}$.
- A **circulant matrix** is a matrix in which each row is obtained by shifting the previous row one element to the right.

Definition

A QC-MDPC code (QC = quasicyclic) is a MDPC code with a parity check matrix composed of circulant blocks.

BIKE at a high level

- Based on binary linear codes with
 - Quasi-cyclic structure: private key composed of two circulant blocks H_0, H_1
 - Moderately-dense parity check matrices
- Let r denote circulant block length. Let t denote maximum error weight.
 - Secret key $H \in \mathbb{F}_2^{r \times 2r}$ is of the form $H = [H_0 | H_1]$
 - Public key $H' = H_0^{-1}(H)$
 - Message encoded as error vector $e \in \mathbb{F}_2^{2r}$ of weight t
 - Ciphertext is syndrome $s = He^T \in \mathbb{F}_2^r$. Decrypt using Black-Grey-Flip syndrome decoder

BIKE at a high level

Parameters

r : block length

w : row weight of secret key

t : maximum error weight

λ : security parameter

Design principles

r prime

$x^r - 1$ has only two irreducible factors

$w \in \mathcal{O}(\sqrt{n})$

$w = 2d, d$ odd

$\lambda \approx t - \frac{1}{2} \log_2 r \approx w - \frac{1}{2} \log_2 r$

What is an error floor?

Graphs of DFRs on a log scale for low- to moderate-density parity check codes with iterative decoders display a phenomenon:

- Initial, rapid decrease of decoding failures (**waterfall region**)
- Eventual plateau, more linear decrease (**error floor region**)

To accurately predict the DFR for higher code length (signal-to-noise ratio), one must account for the error floor region.

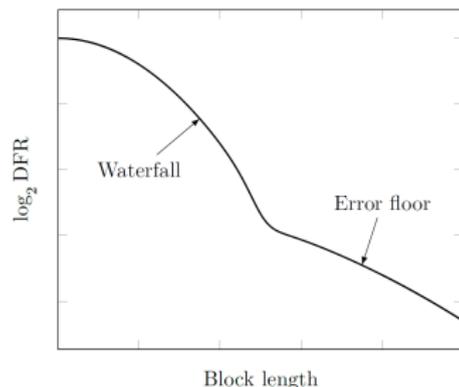


Image credit: Vasseur, V. (2021). Post-quantum cryptography: a study of the decoding of QC-MDPC codes. PhD thesis, Université de Paris

LDPC code approach

Represent code in Tanner graph form:

- Sparse bipartite graph
- Results on minimum distance based on girth (length of shortest cycle)
- Prevalence of small, closed loops increase probability of decoding failure

Definition

Let H be a parity-check matrix describing a code C . A (u, v) -**near codeword** is an error vector e of weight u whose syndrome $s = He^T$ has weight v .

McKay, Postol (2003): near codewords with small u, v and low-weight codewords cause high error floor for certain LDPC codes.

Marco Baldi. QC-LDPC Code-Based Cryptography (2014)

David J.C. MacKay, Michael S. Postol. Weaknesses of Margulis & Ramanujan-Margulis Low-Density Parity-Check Codes (2003)

Tom Richardson. Error floors of LDPC codes (2003)

Gerd Richter. Finding small stopping sets in the Tanner graphs of LDPC codes (2006)

MDPC code approach

Tanner graph is less sparse for MDPC codes. Too expensive to directly extend LDPC code results and techniques.

Approaches towards studying error floors of MDPC codes:

- Baldi et al. (2021) rigorously prove existence of error floor for QC-MDPC codes as a function of code length.
- Vasseur (2021) defines three sets of near-codewords and low-weight codewords, and analyzes their impact on the BIKE DFR.

Marco Baldi, Alessandro Barenghi, Franco Chiaraluce, Gerardo Pelosi, and Paolo Santini, Performance bounds for QC-MDPC codes decoders (2021)

Valentin Vasseur. Post-quantum cryptography: a study of the decoding of QC-MDPC codes (2021)

Valentin Vasseur. QC-MDPC codes DFR and the IND-CCA security of BIKE (2021)

Section 2

Our approach

20-bit DFR simulations

To better understand the error floor behavior of BIKE DFR curves, we experimentally consider BIKE at the 20-bit security level.

- 1 Use BIKE design parameters to generate parameter sets for $\lambda = 20$.
- 2 Use Boston University Shared Computing Cluster to run highly parallelizable experiments.
- 3 Examine factors that increase decoding failures, affecting the error floor.

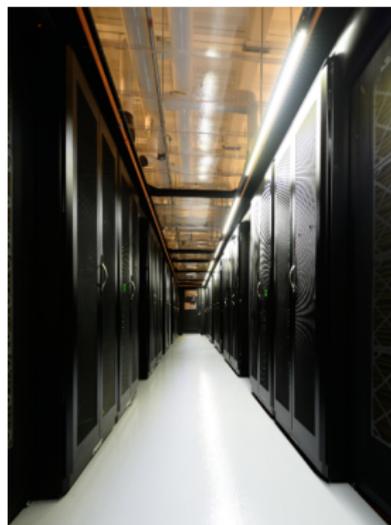


Image credit: <https://www.bu.edu/tech/services/research/computation/scc/>

Methods

Parameter selection:

- $(r, w, t, \lambda) = (523, 30, 18, 20)$
- Extend set of block sizes to “find” error floor

Weak key considerations:

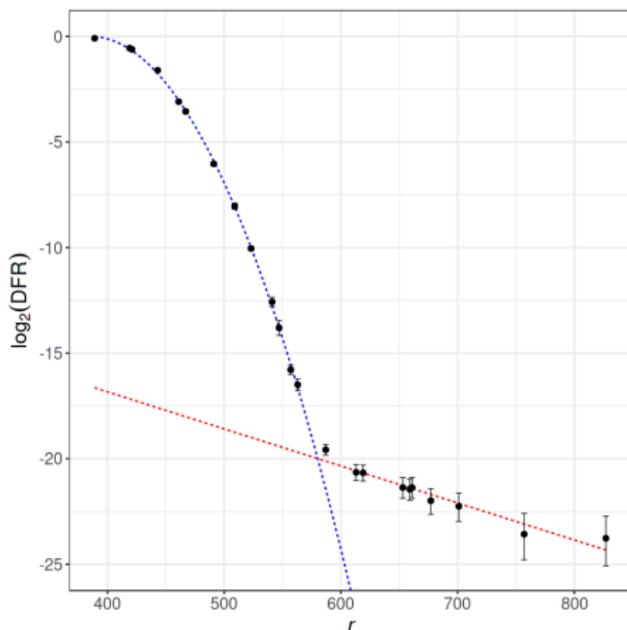
- Vasseur identifies 3 classes of BIKE weak keys based on threshold T .
- We determine $T = 3$ for $\lambda = 20$.

Average DFR per r for all messages:

- 1 Sample a random key H , reject if H is a weak key.
- 2 Sample a random message e .
- 3 Compute $s = He^T$.
- 4 Run BGF decoder on input (H, s) . For output e' , decoder is said to have failed if $e' \neq e$.
- 5 Repeat N times.

20-bit DFR plot

- 95% confidence intervals shown on plot.
- Tested 10^8 keys for $r \in [587, 827]$.
- Fewer keys tested for smaller r because higher DFR means fewer trials needed to narrow confidence intervals.
- Fit lines are quadratic (blue) for waterfall region and linear (red) for error floor region.



Section 3

Analysis of DFR on Special Sets

Special sets of problematic error vectors

Near codewords of small weight and syndrome weight are expected to cause decoding failures.

Vasseur identified three sets:

- \mathcal{C} : Weight w codewords in $\text{null}(H)$.
- \mathcal{N} : Weight d , syndrome-weight d near-codewords.
- $2\mathcal{N}$: Sums of two elements of \mathcal{N} .

Vasseur also introduced sets of vectors which are near $\mathcal{C}, \mathcal{N}, 2\mathcal{N}$:

For $\mathcal{S} \in \{\mathcal{C}, \mathcal{N}, 2\mathcal{N}\}$ and a general vector $e \in \mathbb{F}_2^{2r}$, $\ell := \max_{s \in \mathcal{S}} |e \star s|$.

Vectors with ℓ close to $|e|$ are near \mathcal{S} .

Let δ denote the distance from \mathcal{S} . Vectors with small δ are near \mathcal{S} .

Vectors in/near $\mathcal{C}, \mathcal{N}, 2\mathcal{N}$ are difficult to distinguish in syndrome decoding
 \Rightarrow **decoding failures.**

The sets \mathcal{C} , \mathcal{N} , and $2\mathcal{N}$

$$H = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

\mathcal{C} : nullspace of H .

\mathcal{N} : half-rows of H

$2\mathcal{N}$: $n_1 + n_2$, $n_i \in \mathcal{N}$

$$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$$

Our Contribution: experimentally identify which vectors cause decoding failures at the 20-bit security level.

DFR for special sets

For $r = 523, 587, 659$, we computed the DFR on vectors of weight $t = 18$ of varying distances δ from the sets \mathcal{C} , \mathcal{N} , and $2\mathcal{N}$.

\mathcal{C} :
Weight $w = 30$
 $\Rightarrow \delta$ in $\{12, 14, \dots, 48\}$

\mathcal{N} :
Weight $d = 15$
 $\Rightarrow \delta$ in $\{3, 5, \dots, 33\}$

$2\mathcal{N}$:
Weight at most 30
 $\Rightarrow \delta \leq 48$, even

How does the DFR of weight- t vectors close to \mathcal{C} , \mathcal{N} , $2\mathcal{N}$ compare to the DFR on generic weight- t vectors?

DFR for special sets

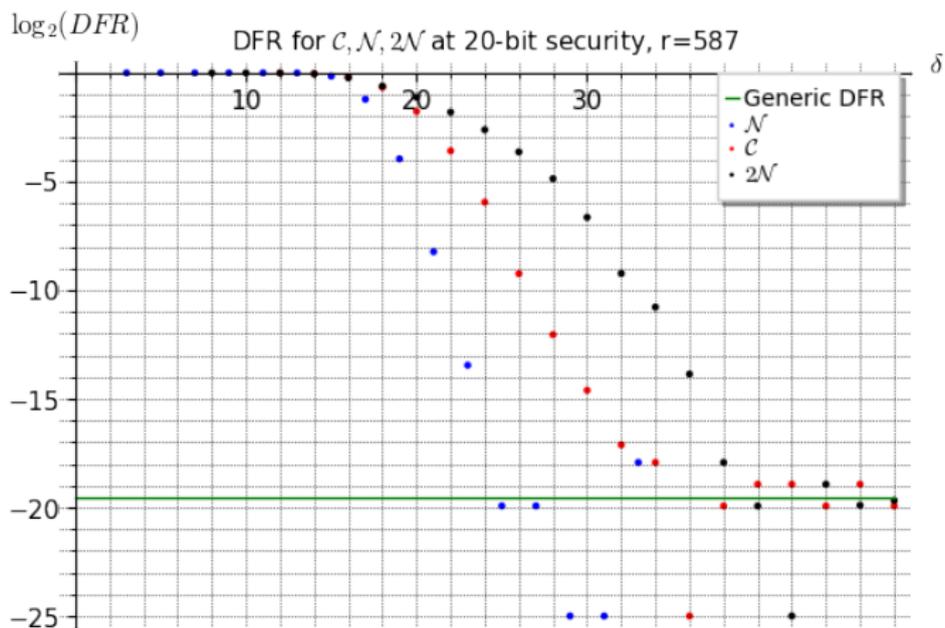


Figure: DFR versus δ for $r = 587$

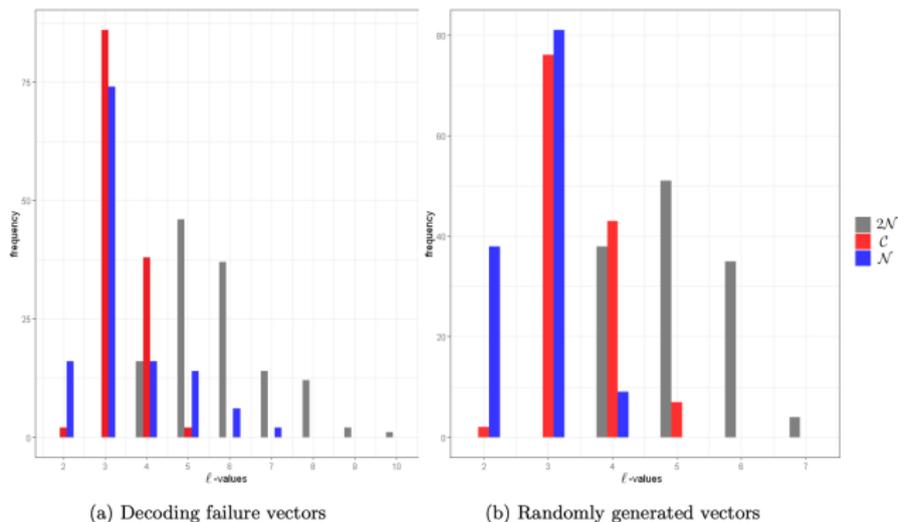
Vectors with small δ are near \mathcal{S} .

Special sets vs. general error vectors

In our generic DFR computation, we recorded decoding failures.

How many overlaps do decoding failure vectors have with $\mathcal{C}, \mathcal{N}, 2\mathcal{N}$?
 $r = 587$. For each decoding failure vector, we found the maximum number of overlaps with an element of \mathcal{S} for each $\mathcal{S} \in \{\mathcal{C}, \mathcal{N}, 2\mathcal{N}\}$.

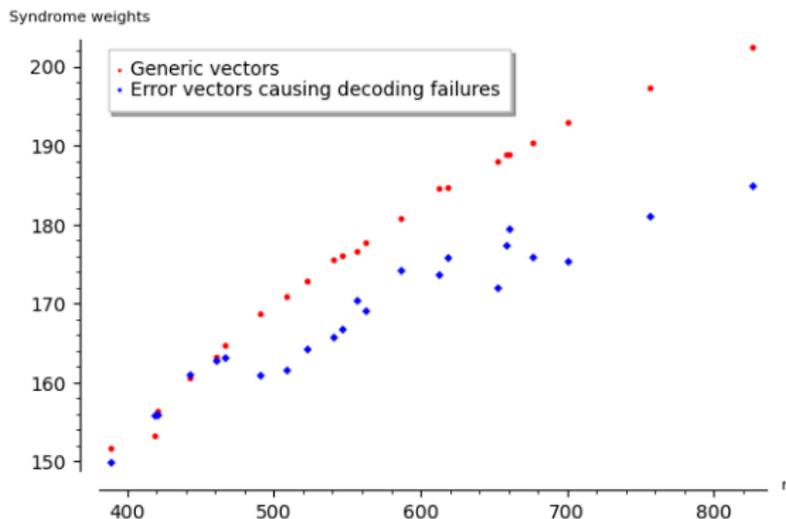
We repeat for the same number of random vector of weight $t = 18$.



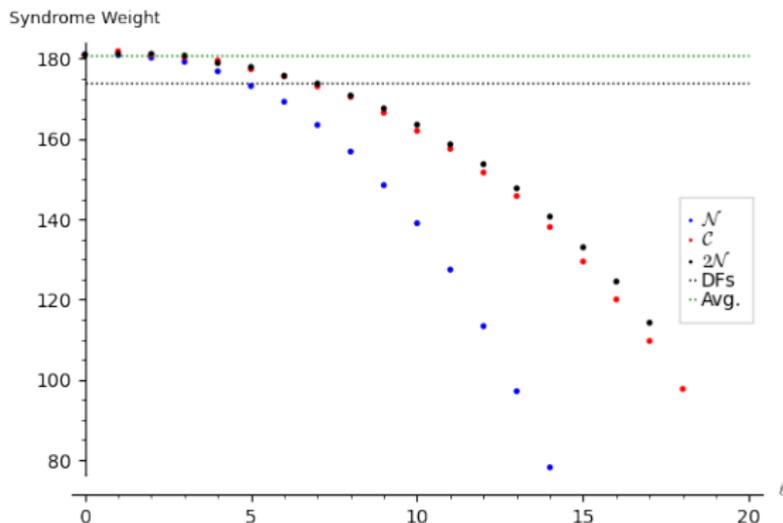
Decoding failure vectors are not unusually close to $\mathcal{C}, \mathcal{N}, 2\mathcal{N}$.

Syndrome weights of decoding failure vectors

Given a parity check matrix H , a vector v , the **syndrome** of v is $s := Hv^T$. The **syndrome weight**, $|s|$, is the number of nonzero entries of s . Vectors causing decoding failures have smaller-than-average syndrome weights.



$r = 587$, syndrome weights vs. overlaps with $\mathcal{C}, \mathcal{N}, 2\mathcal{N}$



For decoding failure vectors, we had the following average l numbers of overlaps with \mathcal{C} , \mathcal{N} , and $2\mathcal{N}$:

- \mathcal{C} : mean $l \approx 3.31$
- \mathcal{N} : mean $l \approx 3.42$
- $2\mathcal{N}$: mean $l \approx 5.77$

Section 4

Conclusions and Next Steps

Conclusions

- The waterfall/error-floor DFR picture remains at the 20-bit security level for BIKE.
- \mathcal{C} , \mathcal{N} , $2\mathcal{N}$, and the corresponding $\mathcal{A}_{t,\ell}$ sets are not overly represented among decoding failures at the 20-bit level.
- Decoding failure vectors do have lower-than-average syndrome weight.

Next Steps

- Where does the error floor begin at higher security levels? Can we make a conjecture based on our 20-bit data?
- How can we classify the error vectors which heavily contribute to decoding failures? Are there new categories, in addition to \mathcal{C} , \mathcal{N} , $2\mathcal{N}$, and the $\mathcal{A}_{t,\ell}$ sets?
- Is it possible for the error vectors to be brought closer to \mathcal{C} , \mathcal{N} , $2\mathcal{N}$ during the iterative decoding process?
- Can Tanner graph techniques from LDPC codes be useful in identifying new classes of vectors which contribute to decoding failure?

Thank you.

<https://eprint.iacr.org/2022/1043>