Engineering Mutation-Tolerant Genes

by

Prince Kwame Ampofo

## A Thesis Presented in Partial Fulfillment of the Requirements for the Degree Master of Science

Approved April 2019 by the Graduate Supervisory Committee:

Xiaojun Tian, Chair Samira Kiani Yang Kuang

ARIZONA STATE UNIVERSITY

May 2019

#### ABSTRACT

Ideas from coding theory are employed to theoretically demonstrate the engineering of mutation-tolerant genes, genes that can sustain up to some arbitrarily chosen number of mutations and still express the originally intended protein. Attention is restricted to tolerating substitution mutations. Future advances in genomic engineering will make possible the ability to synthesize entire genomes from scratch. This presents an opportunity to embed desirable capabilities like mutation-tolerance, which will be useful in preventing cell deaths in organisms intended for research or industrial applications in highly mutagenic environments. In the extreme case, mutation-tolerant genes (mutols) can make organisms resistant to retroviral infections.

An algebraic representation of the nucleotide bases is developed. This algebraic representation makes it possible to convert nucleotide sequences into algebraic sequences, apply mathematical ideas and convert results back into nucleotide terms. Using the algebra developed, a mapping is found from the naturally-occurring codons to an alternative set of codons which makes genes constructed from them mutationtolerant, provided no more than one substitution mutation occurs per codon. The ideas discussed naturally extend to finding codons that can tolerate t arbitrarily chosen number of mutations per codon. Finally, random substitution events are simulated in both a wild-type green fluorescent protein (GFP) gene and its mutol variant and the amino acid sequence expressed from each post-mutation is compared with the amino acid sequence pre-mutation.

This work assumes the existence of synthetic protein-assembling entities that function like tRNAs but can read k nucleotides at a time, with  $k \ge 5$ . The realization of this assumption is presented as a challenge to the research community.

# DEDICATION

 $This \ work \ is \ dedicated \ to \ my \ family \ for \ their \ love \ and \ care.$ 

## ACKNOWLEDGMENTS

I would like to express my gratitude to Dr. Xiaojun Tian, Dr. Samira Kiani and Dr. Yang Kuang for serving on my committee and for their constructive feedback on my work. Special thanks also to the Mastercard Foundation for the scholarship that made my education at Arizona State University, which culminated in this thesis,

possible.

		P	age
LIST	OF 7	TABLES	v
CHAI	PTEF	ł	
1	ALC	GEBRAIC REPRESENTATION OF NUCLEOTIDE BASES	1
	1.1	Introduction	1
	1.2	Field	1
	1.3	Constructing Field of order 4	3
	1.4	Assignment of Field elements to nucleotide Bases	6
	1.5	Doing Algebra with the Base Representations	7
	1.6	Algebraic representation of Codons	10
2	COI	DONS AS VECTORS	11
	2.1	Introduction	11
	2.2	Codon Vector Space	11
	2.3	Subspaces of $\mathbb{B}^n$	17
	2.4	Linear Combination, Span, and Linear Independence	18
	2.5	Bases and Dimension of $\mathbb{B}^n$	20
	2.6	Orthogonal Complements	22
3	TOI	LERATING SUBSTITUTION MUTATIONS	25
	3.1	Introduction	25
	3.2	Hamming Distance and Weight	26
	3.3	Tolerating $t$ Substitution Mutations per Codon	28
	3.4	Example: Tolerating 1 Substitution Mutation per Codon	31
REFE	EREN	CES	45
APPI	ENDE	X	
А	3-LH	ETTER CODONS TO 5-LETTER CODONS	46

# TABLE OF CONTENTS

# LIST OF TABLES

Table	Η	Page
1.1	Algebraic representation of Nucleotide Bases	7
1.2	Arithmetic table for $GF(2^2)$	9
3.1	Mapping to Alternative Codons	34
3.1	Mapping to Alternative Codons	35
3.1	Mapping to Alternative Codons	36
3.3	GFP DNA Sequence and Its Mutol Variant	41
3.3	GFP DNA Sequence and Its Mutol Variant	42
3.3	GFP DNA Sequence and Its Mutol Variant	43
3.4	Mutation Simulation Results	44
A.1	Degeneracy Mapping	47
A.1	Degeneracy Mapping	48
A.1	Degeneracy Mapping	49
A.1	Degeneracy Mapping	50
A.1	Degeneracy Mapping	51
A.1	Degeneracy Mapping	52
A.1	Degeneracy Mapping	53
A.1	Degeneracy Mapping	54
A.1	Degeneracy Mapping	55
A.1	Degeneracy Mapping	56

#### Chapter 1

## ALGEBRAIC REPRESENTATION OF NUCLEOTIDE BASES

#### 1.1 Introduction

In our effort to develop genes that can tolerate mutations, we will be employing some mathematical ideas. To make this possible, we need to represent nucleotide bases, the building blocks from which DNA and RNA are constructed, with suitable mathematical entities. This way, we can easily apply mathematical ideas to algebraic representations of codons or nucleotide sequences and convert results back into nucleotide terms. This chapter is devoted to finding suitable algebraic representations for the bases: adenine (A), guanine (G), cytosine (C), thymine (T), and uracil (U).

In section 1.2, an interesting algebraic structure called a field is introduced. In section 1.3, we construct a field whose elements will be used as algebraic representations of the bases. In section 1.4, we assign elements of the field constructed in section 1.3 to the bases. Finally, in section 1.5, we discuss how computation is done in a field.

Most concepts utilized throughout this work are from the areas of linear algebra and coding theory. The uninitiated will find the following references helpful: Axler (2007), Ling (2004), and Lin and Costello (1983).

## 1.2 Field

A field F is a set together with two binary operations addition (denoted by +) and multiplication (denoted by \*) defined such that the following axioms are satisfied:

1. Closure.  $\forall x, y \in F, x + y \in F$  and  $x * y \in F$ 

- 2. Commutativity.  $\forall x, y \in F$ , x + y = y + x and x \* y = y \* x
- 3. Associativity.  $\forall x, y, z \in F$ , x + (y + z) = (x + y) + z and x \* (y \* z) = (x \* y) \* z
- 4. Existence of additive and multiplicative identities. (a) There exists an additive identity 0 ∈ F such that x+0 = x for all x ∈ F. (b) There exists a multiplicative identity 1 ∈ F such that x \* 1 = x for all x ∈ F
- 5. Existence of additive and multiplicative inverses. (a) For every x ∈ F, there exists an additive inverse denoted -x ∈ F such that x+(-x) = 0. (b) For every x ∈ F except 0, there exists an inverse denoted x<sup>-1</sup> ∈ F such that x \* x<sup>-1</sup> = 1.
- 6. Distributivity of multiplication over addition. For all  $x, y, z \in F$ , x \* (y + z) = x \* y + x \* z
- 7. Distinctness of additive and multiplicative identities.  $1 \neq 0$

The axioms outlined above demonstrate why a field is ideal for use in representing nucleotide bases algebraically. The algebraic properties a field possesses, like closure, commutativity, associativity, distributivity, etcetra, as outlined in axioms 1 through 6, are essential for ease of algebraic manipulations. By using elements from a field for the representation of nucleotides, we will be able to easily perform algebraic manipulations and apply mathematical ideas to algebraic equivalents of nucleotide sequences, toward the goal of developing mutation-tolerant genes.

In the next section, we will construct a field with 4 elements. Adenine, Guanine, and Cytosine will each be assigned one of these algebraic elements. Thymine and Uracil will both be assigned the same fourth element since Uracil is the RNA "equivalent" of Thymine.

#### 1.3 Constructing Field of order 4

In general, if p is prime and addition and multiplication are done in modulo-p, the set  $\{0, 1, 2, \ldots, p-1\}$  satisfies all the field axioms outlined in Section 1.2 above and thus forms a field of p elements. The number of elements in a field is called its order. When the order of a field is finite, we call it a Galois field and denote it by  $GF(number \ of \ elements)$ . Also, when the number of elements in a Galois field is prime, we call the field a prime field.

We desire to construct a field of 4 elements: GF(4). However, 4 is not prime and so  $\{0, 1, 2, \ldots, 4-1\} = \{0, 1, 2, 3\}$  is not a field. We can however employ a well-known idea from field theory which allows the construction of any field  $GF(p^q)$  from GF(p), where p is prime and q is any positive integer. Thus, we can construct our desired field  $GF(2^2)$  from  $GF(2) = \{0, 1\}$ , which is the field formed from  $\{0, 1, 2, \ldots, p-1\}$ when p = 2. When GF(2) is extended to get  $GF(2^2)$ , GF(2) is called the base field and  $GF(2^2)$  is called an extension field. To utilize this strategy, we first introduce two concepts:

Polynomials over fields. If we construct an n-degree polynomial f(α) = f<sub>0</sub> + f<sub>1</sub>α+f<sub>2</sub>α<sup>2</sup>+···+f<sub>n-1</sub>α<sup>n-1</sup>+f<sub>n</sub>α<sup>n</sup> such that its coefficients f<sub>0</sub>, f<sub>1</sub>, f<sub>2</sub>,..., f<sub>n-1</sub>, f<sub>n</sub> are all taken from a field called F, we say that f(α) is a "polynomial over field F". As an example, polynomials over GF(2) are polynomials with all coefficients either 0 or 1. In general, there are 2<sup>n</sup> polynomials of degree n that can be constructed over GF(2). To see this, notice that for any polynomial of degree n over GF(2), its coefficients f<sub>0</sub>, f<sub>1</sub>, f<sub>2</sub>,..., f<sub>n-1</sub> can each have two possible values (0 or 1) whereas coefficient f<sub>n</sub> must be 1 since if it is 0, the polynomial will no longer be of degree n. Thus, the total number of n-degree polynomials over GF(2) is 2 \* 2 \* ··· \* 2 / n times

over GF(2) that are of degree 2. They are:

$$\alpha^2$$
,  $\alpha + \alpha^2$ ,  $1 + \alpha^2$ , and  $1 + \alpha + \alpha^2$  (1.1)

2. Primitive polynomials. A q-degree polynomial  $f(\alpha)$  over a prime field GF(p)is called primitive if the smallest positive integer n for which  $f(\alpha)$  divides  $\alpha^n + 1$ without a remainder is  $n = p^q - 1$ . Primitive polynomials are important because they are used in constructing extension fields from base fields. It is difficult to test all polynomials over a given field to find ones that are primitive. As such, tables of primitive polynomials exist from which desired primitive polynomials can be obtained. Using one such table, we find that of the four polynomials of degree 2 over GF(2) (listed in 1.1), only one is primitive:

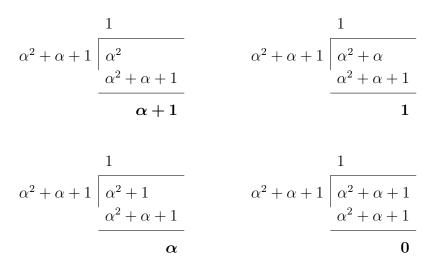
$$p(\alpha) = 1 + \alpha + \alpha^2 \tag{1.2}$$

With these two concepts introduced, we state without proof a theorem with which we can construct  $GF(2^2)$  from GF(2).

**Theorem 1.1.** If  $f(\alpha)$  is a q-degree primitive polynomial over GF(p), the set of all q-degree polynomials over GF(p) taken as modulo  $f(\alpha)$  forms  $GF(p^q)$ .

*Remark.* Just as an integer like 9 can be expressed in modulo 3 by dividing 9 by 3 and keeping the remainder, a polynomial  $f_1(\alpha)$  can be expressed as modulo of another polynomial  $f_2(\alpha)$  by dividing  $f_1(\alpha)$  by  $f_2(\alpha)$  and keeping the remainder.

Using Theorem 1.1 above, we will divide each of the 4 polynomials in Equation (1.1) by the primitive polynomial in Equation (1.2), keeping the remainder in each case. The set of these four remainders, under modulo 2 addition and multiplication, will satisfy all the field axioms in Section 1.2 and thus form the desired field of four elements,  $GF(2^2)$ .



The remainders are 0, 1,  $\alpha$ , and  $\alpha + 1$ . Using Theorem 1.1, our desired field<sup>1</sup> is:

 $GF(2^2) = \{ 0, 1, \alpha, \alpha + 1 \}, \text{ with } + \text{ and } * \text{ done in mod } 2$  (1.3)

**Compact form of**  $\alpha + 1$ . For every extension field, all the non-zero elements of the set are expressed in terms of one special element called the **primitive element**. By noticing that 1 can be written as  $\alpha^0$ , it is easy to spot that  $\alpha$  is the primitive element in the field we just constructed (Equation 1.3 above). A well-known result in field theory is that a primitive element is a root of its corresponding primitive polynomial. Thus,  $p(\alpha) = 1 + \alpha + \alpha^2 = 0$ , where p is the primitive polynomial used in constructing  $GF(2^2)$ . Using this relation, we can find another useful form of the  $\alpha + 1$  element in  $GF(2^2)$ :

<sup>&</sup>lt;sup>1</sup>The phrase "with + and \* done in mod 2" in (1.3) is vital because a set by itself is not a field. When we mention the set without the operations, it must be taken as implied.

$$1 + \alpha + \alpha^{2} = 0$$

$$1 + \alpha + (\alpha^{2} + \alpha^{2}) = 0 + \alpha^{2} \qquad (added \ \alpha^{2} \text{ to both sides})$$

$$\alpha^{2} + \alpha^{2} = (1+1)\alpha^{2} = (0)\alpha^{2} = 0 \qquad (1+1=0 \text{ due to modulo 2 addition})$$

$$\therefore \quad 1 + \alpha = \alpha^{2} \qquad (1.4)$$

According to (1.4) above, the element  $\alpha + 1$  in  $GF(2^2)$  is equal to  $\alpha^2$ . We can thus write  $GF(2^2) = \{0, 1, \alpha, \alpha+1\}$  as  $\{0, 1, \alpha, \alpha^2\}$ . Because  $\alpha^2$  is more compact, we will use it more often. Being in exponent form, it will also be used primarily in multiplication operations, whereas the  $\alpha + 1$  form will be used mainly in addition operations.

## 1.4 Assignment of Field elements to nucleotide Bases

Our goal at the beginning of this chapter was to find suitable algebraic entities to represent the nucleotide bases. We constructed a field of four elements,  $GF(2^2)$ , in the preceding section. Here, we assign these elements to the bases.

As outlined in Table 1.1 below, we will represent Adenine with 0, Guanine with  $\alpha$ , and Cytosine with  $\alpha^2$ . Because the RNA "equivalent" of Thymine is Uracil, they are both assigned 1. These assignments are not arbitrary. They have been done such that adding 1 to the algebraic representation of any base gives the algebraic representation of its Watson-Crick base pair. Examples<sup>2</sup>:

- 1) 0 (Adenine) + 1 = 1 (Thymine)
- 2) 1 (Thymine) + 1 = 0 (Adenine)
- 3)  $\alpha$  (Guanine) + 1 =  $\alpha^2$  (Cytosine)
- 4)  $\alpha^2$  (Cytosine) + 1 = ( $\alpha$  + 1) + 1 =  $\alpha$  (Guanine)

<sup>&</sup>lt;sup>2</sup>All operations are done in modulo 2. Also, note that  $\alpha^2 = \alpha + 1$  (Equation 1.4)

Base	Algebraic Representation
Adenine (A)	0
Guanine (G)	lpha
Cytosine (C)	$lpha^2$
Thymine (T) & Uracil (U)	1

**Table 1.1:** Algebraic representation of Nucleotide Bases. Note that  $\alpha^2 = \alpha + 1$ 

We will refer to the set of elements used to represent the bases as  $\boldsymbol{B}$  (as in Bases), where the subscripts used are simply to indicate which base the element represents and may be left out for the sake of compactness:

$$\boldsymbol{B} = \{0_A, \ 1_T \text{ or } U, \ \alpha_G, \ \alpha^2_C\}$$
(1.5)

As stated earlier, notice that 1 is used to represent both Thymine (T) and Uracil (U) since each can be thought of as the DNA or RNA "equivalent" of the other. As such, when 1 appears in the algebraic representation of a nucleotide sequence, it should be seen as representing Thymine or Uracil based on the context. Obviously, if 1 appears in a DNA sequence, it represents Thymine. In an RNA sequence, it represents Uracil.

#### 1.5 Doing Algebra with the Base Representations

Our goal for this chapter was to find suitable algebraic representations for the nucleotide bases so that we can write any nucleotide sequence as a sequence of algebraic entities, apply mathematical ideas to it, and map results back into nucleotide terms, with the aim of developing mutation-tolerant genes. Having found such suitable algebraic entities (elements of set B), we now discuss how addition, multiplication, subtraction, and division can be done with them.  $GF(2) = \{0, 1\}$ , the base field from which our extension field of 4 elements,  $\mathbf{B} = \{0, 1, \alpha, \alpha^2\}$ , was constructed, is only a field under mod 2 addition and multiplication. By implication,  $\mathbf{B}$  is also a field under mod 2 addition and multiplication. As such all operations will be done in mod 2. Note that because we used elements from a field to represent the bases, all the field axioms outlined in Section 1.2 hold. As such, for all operations, the commutative, associative, and distributive properties can be used. It's also vital to remember that  $\alpha^2 = \alpha + 1$  (Equation 1.4).

- Addition. Adding any elements of B is just like normal addition, keeping in mind that in mod 2, 0+0=0, 0+1=1, 1+0=1 and 1+1=0. An addition table is shown in Table 1.2a. Adding any element to itself gives 0. Also, as indicated earlier, adding α<sup>2</sup> (which equals α + 1) to any element gives its Watson-Crick base pair.
- Multiplication. Multiplication can also be done by keeping in mind that in mod 2, 0 \* 0 = 0, 0 \* 1 = 0, 1 \* 0 = 0 and 1 \* 1 = 1. Note that for any element b,

$$b^{i} * b^{k} = \underbrace{b * b * \cdots * b}_{i \text{ times}} * \underbrace{b * b * \cdots * b}_{k \text{ times}} = b^{i+k}$$

One relation that will be useful in simplifying multiplication operations is  $\alpha^3 = 1$ . To see this, notice that  $\alpha^3 = \alpha^2 * \alpha$ , and from the multiplication table (Table 1.2b),  $\alpha^2 * \alpha = 1$ . For example,  $\alpha * \alpha^2 * \alpha^2 = \alpha^5$  can be simplified as  $\alpha^5 = \alpha^3 * \alpha^2 = 1 * \alpha^2 = \alpha^2$ .

*Remark.* Though only the two operations discussed above are explicitly defined on a field, we can still do subtraction and division because subtraction and division can be written in terms of addition and multiplication, respectively.

+	0	1	$\alpha$	$\alpha^2$		*	0	1	$\alpha$	$\alpha^2$
0	0	1	$egin{array}{c} lpha \\ lpha^2 \\ 0 \\ 1 \end{array}$	$\alpha^2$	_	0	0	0	$0$ $\alpha$ $\alpha^{2}$ $1$	0
1	1	0	$\alpha^2$	$\alpha$		1	0	1	$\alpha$	$\alpha^2$
$\alpha$	α	$\alpha^2$	0	1		$\alpha$	0	$\alpha$	$\alpha^2$	1
$\alpha^2$	$\alpha^2$	$\alpha$	1	0		$\alpha^2$	0	$\alpha^2$	1	$\alpha$
(a) Addition Table									tion T	

**Table 1.2:** Arithmetic table for  $GF(2^2)$ . Note that + and \* are done mod 2.

- Subtraction. Note that a − b = a + (−b), where −b is the additive inverse of b. Thus, to subtract b from a, we first find the additive inverse of b (which is denoted by −b) and add it to a. Note that an additive inverse of an element b is the element which when added to b, gives 0. Looking at the addition table in Table 1.2a, we notice that adding any element to itself gives 0. This means that for set B, the additive inverse of an element from another is the same as adding them. For example, a − b = a + (−b) = a + b since −b, which denotes the inverse of b, equals b itself. 2) Because of implication 1, the addition table shown in Table 1.2a also acts as a subtraction table for the elements in B.
- Division. Note that <sup>a</sup>/<sub>b</sub> = a \* b<sup>-1</sup>, where b<sup>-1</sup> is the multiplicative inverse of b. Thus, to divide a by b, we first find the multiplicative inverse of b then multiply it by a. Remember that the multiplicative inverse of an element b is an element which when multiplied by b, gives 1. For example, to find <sup>1</sup>/<sub>α<sup>2</sup></sub>, we first find the multiplicative inverse of α<sup>2</sup> by looking at the multiplicative table (Table 1.2b) to find the element which when multiplied by α<sup>2</sup>, gives 1. Since α<sup>2</sup> \* α = 1, α is the multiplicative inverse of α<sup>2</sup>. Thus, <sup>1</sup>/<sub>α<sup>2</sup></sub> = 1 \* (α<sup>2</sup>)<sup>-1</sup> = 1 \* α = α.

## 1.6 Algebraic representation of Codons

Having found suitable algebraic representation for the nucleotide bases, we can now write each of the 64 DNA sense  $(5' \rightarrow 3')$  codons in terms of the algebraic representation of its nucleotide bases.

$ \begin{array}{c} \left(\alpha \ \alpha^{2} \ 1\right)_{GCT} \\ \left(\alpha \ \alpha^{2} \ 0\right)_{GCA} \\ \left(\alpha \ \alpha^{2} \ \alpha^{2}\right)_{GCC} \\ \left(\alpha \ \alpha^{2} \ \alpha\right)_{GCG} \end{array} \right\} \text{Ala}  \begin{array}{c} \left(0 \ 0 \ 1\right)_{AAT} \\ \left(0 \ 0 \ \alpha^{2}\right)_{AAC} \end{array} \right\} \text{Asn} $	$ \begin{pmatrix} \alpha & 0 & 1 \end{pmatrix}_{GAT} \\ (\alpha & 0 & \alpha^2)_{GAC} \end{pmatrix} \operatorname{Asp} \qquad \begin{pmatrix} 1 & \alpha & 1 \end{pmatrix}_{TGT} \\ (1 & \alpha & \alpha^2)_{TGC} \end{array} \right\} \operatorname{Cys} $
$ \begin{pmatrix} (\alpha \ 0 \ 0) \\ (\alpha \ 0 \ \alpha) \\ GAG \end{pmatrix} Glu  \begin{pmatrix} (\alpha^2 \ 0 \ 0) \\ (\alpha^2 \ 0 \ \alpha) \\ CAG \end{pmatrix} Gln $ $ (1 \ 1 \ 0) \\ TTA $	$ \begin{pmatrix} (\alpha \ \alpha \ 1) \ _{GGT} \\ (\alpha \ \alpha \ 0) \ _{GGA} \\ (\alpha \ \alpha \ \alpha^2) \ _{GGC} \\ (\alpha \ \alpha \ \alpha) \ _{GGG} \end{pmatrix} Gly \qquad \begin{pmatrix} (\alpha^2 \ 0 \ 1) \ _{CAT} \\ (\alpha^2 \ 0 \ \alpha^2) \ _{CAC} \end{pmatrix} His $
$ \begin{array}{c} (0 \ 1 \ 1)_{AAT} \\ (0 \ 1 \ \alpha^2)_{ATC} \\ (0 \ 1 \ 0)_{ATA} \end{array} \right\} IIe \qquad \begin{array}{c} (1 \ 1 \ \alpha)_{TTG} \\ (\alpha^2 \ 1 \ 1)_{CTT} \\ (\alpha^2 \ 1 \ \alpha)_{CTA} \\ (\alpha^2 \ 1 \ \alpha^2)_{CTC} \end{array} \right\} Le $	$ \begin{array}{c} \text{eu} & (0 \ 0 \ 0)_{AAA} \\ (0 \ 0 \ \alpha)_{AAG} \end{array} \right\} \text{Lys} \qquad (0 \ 1 \ \alpha)_{ATG} \text{ Met} $
$(1 \alpha^{-} 1) TCT$	$ \begin{array}{c} \begin{array}{c} \left( 0 \ \alpha^{2} \ 1 \right)_{ACT} \\ \left( 0 \ \alpha^{2} \ \alpha^{2} \right)_{ACC} \\ \left( 0 \ \alpha^{2} \ 0 \right)_{ACA} \\ \left( 0 \ \alpha^{2} \ \alpha \right)_{ACG} \end{array} \right\} \operatorname{Thr} \left( \begin{array}{c} \left( 1 \ 0 \ 1 \right)_{TAT} \\ \left( 1 \ 0 \ \alpha^{2} \right)_{TAC} \end{array} \right\} \operatorname{Tyr} \\ \left( \begin{array}{c} \left( \alpha^{2} \ \alpha \ 1 \right)_{CGT} \end{array} \right) \end{array} \right\} $
$ \begin{array}{c} (\alpha \ 1 \ 1) \ _{GTT} \\ (\alpha \ 1 \ 0) \ _{GTA} \\ (\alpha \ 1 \ \alpha^2) \ _{GTC} \\ (\alpha \ 1 \ \alpha^2) \ _{GTC} \\ (\alpha \ 1 \ \alpha^2) \ _{GTG} \end{array} \right\} Val  \begin{array}{c} (1 \ \alpha^2 \ \alpha^2) \ _{TCC} \\ (1 \ \alpha^2 \ \alpha) \ _{TCG} \\ (1 \ \alpha^2 \ \alpha) \ _{TCG} \\ (0 \ \alpha \ 1) \ _{AGT} \\ (0 \ \alpha \ \alpha^2) \ _{AGC} \end{array} \right\} Ser \\ \begin{array}{c} (1 \ 0 \ \alpha) \ _{TAA} \\ (1 \ 0 \ \alpha) \ _{TAG} \\ (1 \ \alpha \ 0) \ _{TGA} \end{array} \right\} Stop $	$(1 \alpha \alpha)_{TGG} \operatorname{Trp} \left( \begin{array}{c} (\alpha^2 \alpha \alpha^2)_{CGC} \\ (\alpha^2 \alpha 0)_{CGA} \\ (\alpha^2 \alpha \alpha)_{CGG} \\ (0 \alpha 0)_{AGA} \\ (0 \alpha \alpha)_{AGG} \end{array} \right) \operatorname{Arg}$

#### Chapter 2

## CODONS AS VECTORS

#### 2.1 Introduction

In the previous chapter, a suitable algebraic representation was found for the nucleotide bases:  $\mathbf{B} = \{0_A, 1_T \text{ or } U, \alpha_G, \alpha^2_C\}$ . In section 1.6, each DNA codon was then written in terms of the algebraic representation of its nucleotides. For example, codon ATG in terms of the algebraic representation of its nucleotide bases is  $01\alpha$  since A, T, and G are represented by 0, 1, and  $\alpha$ , respectively.

$$ATG \xrightarrow{\text{algebraic form}} 0 \ 1 \ \alpha$$

To distinguish between codons and their algebraic forms, we will call the algebraic forms <u>**p-codons**</u>. Thus, the p-codon of ATG is  $01\alpha$ . In section 2.2, we consider the set of all 64 p-codons corresponding to the 64 DNA codons and show that this set forms a vector space where the p-codons are vectors and 0, 1,  $\alpha$ , and  $\alpha^2$  are scalars. From section 2.3 through 2.5, relevant vector space concepts are discussed in the context of this codon vector space.

### 2.2 Codon Vector Space

Consider set C below whose elements are the 64 p-codons. We wish to show that this set is a vector space. This is important because by confirming that C is indeed a vector space, we can make use of all theory that has been developed about vector spaces, toward our goal of developing mutation-tolerant genes. Let's first define a vector space.

**Definition 1. (Vector Space)** A vector space V is a set of objects (called vectors) that can be added (+) together and multiplied ( $\cdot$ ) by scalars such that the following axioms are satisfied for any choice of vectors  $\boldsymbol{u}$ ,  $\boldsymbol{v}$ , and  $\boldsymbol{w}$  and scalars a and b :

- 1. Closure.  $\boldsymbol{u} + \boldsymbol{w} \in V$  and  $\boldsymbol{a} \cdot \boldsymbol{v} \in V$
- 2. Commutativity.  $\boldsymbol{u} + \boldsymbol{w} = \boldsymbol{w} + \boldsymbol{u}$
- 3. Additive identity. There exists a zero vector, denoted  $\mathbf{0}$ , in C such that  $\mathbf{u} + \mathbf{0} = \mathbf{u}$
- 4. Multiplicative identity.  $1 \cdot \boldsymbol{v} = \boldsymbol{v}$
- 5. Additive inverses. For any vector  $\boldsymbol{u} \in V$ , there exists  $-\boldsymbol{u} \in V$  such that  $\boldsymbol{u} + (-\boldsymbol{u}) = \boldsymbol{0}$
- 6. Associativity.  $\boldsymbol{u} + (\boldsymbol{v} + \boldsymbol{w}) = (\boldsymbol{u} + \boldsymbol{v}) + \boldsymbol{w}$  and  $a \cdot (b \cdot \boldsymbol{v}) = (a \cdot b) \cdot \boldsymbol{v}$
- 7. Distributivity.  $a \cdot (\boldsymbol{u} + \boldsymbol{v}) = (a \cdot \boldsymbol{u}) + (a \cdot \boldsymbol{v})$  and  $(a + b) \cdot \boldsymbol{u} = (a \cdot \boldsymbol{u}) + (b \cdot \boldsymbol{u})$

What the definition says is that if we form a set and define a way of adding the elements of the set and also define multiplication of the elements by some scalars such that the addition and multiplication operations obey all the 7 axioms listed, then we have a vector space. Thus, to demonstrate that the set of p-codons **C** (with  $\mathbf{B} = \{0, 1, \alpha, \alpha^2\}$  as scalars) is a vector space, we need to do two things

as per the definition: 1) define two operations: the addition of two p-codons and multiplication of a p-codon by a scalar and 2) show that the 7 axioms are satisfied under those operations. We do these in sections 2.2.1 through 2.2.3 below. Addition of p-codons and multiplication of p-codons by scalars are defined in terms of tuples of length n so as to generalize to our later needs. In the context of C, we simply set n = 3.

## 2.2.1 Addition of p-codons.

We define the addition of p-codons  $\boldsymbol{a} = (a_1, a_2, ..., a_n)$  and  $\boldsymbol{b} = (b_1, b_2, ..., b_n)$ as  $\boldsymbol{a} + \boldsymbol{b} = (a_1 + b_1, a_2 + b_2, ..., a_n + b_n)$  where the + in  $a_i + b_i$  for i = 1, 2, ..., n is modulo 2 addition as given in addition table 1.2a.

#### 2.2.2 Multiplication of p-codons by scalars.

We define the multiplication of a p-codon  $\boldsymbol{a} = (a_1, a_2, ..., a_n)$  by a scalar  $\lambda$  in **B** as  $\lambda \cdot \boldsymbol{a} = (\lambda a_1, \lambda a_2, ..., \lambda a_n)$ . We will sometimes write  $\lambda \cdot \boldsymbol{a}$  simply as  $\lambda \boldsymbol{a}$ .

## 2.2.3 Proof that axioms 1 through 7 are satisfied.

Having defined addition of p-codons and scalar multiplication, we are now in a position to show that set **C** satisfies axioms 1 through 7 and is thus a vector space. For the proofs below, it is important to keep two things in mind. First, the set of scalars  $\mathbf{B} = \{0, 1, \alpha, \alpha^2\}$  is a field and so addition or multiplication of any members of the set always results in another member of the set (see Table 1.2). Second, **C** is the set of all 3-tuples whose 3 elements are chosen from **B** and so to show that some arbitrary 3-tuple ( $\lambda_1 \ \lambda_2 \ \lambda_3$ ) is a member of **C**, we only need to show that  $\lambda_1, \lambda_2$ , and  $\lambda_3$  are elements of **B**.

1. Closure.

- (a) For any p-codons u and w in C, u + w is another p-codon in C
  Let u = (u<sub>1</sub> u<sub>2</sub> u<sub>3</sub>) and w = (w<sub>1</sub> w<sub>2</sub> w<sub>3</sub>). Then u + w = (u<sub>1</sub> + w<sub>1</sub> u<sub>2</sub> + w<sub>2</sub> u<sub>3</sub> + w<sub>3</sub>). Since the sum of any elements in B yields another element in B, u<sub>1</sub> + w<sub>1</sub>, u<sub>2</sub> + w<sub>2</sub>, and u<sub>3</sub> + w<sub>3</sub> are in B. We conclude that u + w is in C, as desired.
- (b) For any p-codon v in C and scalar a in B, a ⋅ v is a p-codon in C
  Let v = (v<sub>1</sub> v<sub>2</sub> v<sub>3</sub>). Then a ⋅ v = (av<sub>1</sub> av<sub>2</sub> av<sub>3</sub>). Since the product of any scalars in B yields another in B, av<sub>1</sub>, av<sub>2</sub>, and av<sub>3</sub> are in B. We conclude that a ⋅ v is in C, as desired.
- 2. Commutativity. For any p-codons  $\boldsymbol{u}$  and  $\boldsymbol{w}$  in C,  $\boldsymbol{u} + \boldsymbol{w} = \boldsymbol{w} + \boldsymbol{u}$ Let  $\boldsymbol{u} = (u_1 \ u_2 \ u_3)$  and  $\boldsymbol{w} = (w_1 \ w_2 \ w_3)$ . Then  $\boldsymbol{u} + \boldsymbol{w} = (u_1 + w_1 \ u_2 + w_2 \ u_3 + w_3)$  and  $\boldsymbol{w} + \boldsymbol{u} = (w_1 + u_1 \ w_2 + u_2 \ w_3 + u_3)$ . Addition of scalars from B is commutative and so  $u_i + w_i = w_i + u_i$  for i = 1, 2, 3. Thus,  $\boldsymbol{u} + \boldsymbol{w} = \boldsymbol{w} + \boldsymbol{u}$ .
- 3. Additive identity. There exists a zero p-codon, denoted **0**, in C such that for any p-codon  $\boldsymbol{u}$  in C,  $\boldsymbol{u} + \boldsymbol{0} = \boldsymbol{u}$

Note that 0 is a scalar in *B*. When a = 0 in 1(*b*) above, we find that there exists a vector  $\mathbf{0} = (0, 0, 0)$  in *C*. Let  $\mathbf{u} = (u_1 \ u_2 \ u_3)$ . Then  $\mathbf{u} + \mathbf{0} = (u_1 \ u_2 \ u_3) + (0 \ 0 \ 0) = (u_1 + 0 \ u_2 + 0 \ u_3 + 0) = (u_1 \ u_2 \ u_3) = \mathbf{u}$ , as desired.

- 4. Multiplicative Identity. For any p-codon  $\boldsymbol{v}$  in C,  $1 \cdot \boldsymbol{v} = \boldsymbol{v}$ Note that 1 is the multiplicative identity element in B. Let  $\boldsymbol{v} = (v_1 \ v_2 \ v_3)$ . Then  $1 \cdot \boldsymbol{v} = (1 \cdot v_1 \ 1 \cdot v_2 \ 1 \cdot v_3) = (v_1 \ v_2 \ v_3) = \boldsymbol{v}$ , as desired.
- 5. Additive inverses. For any p-codon  $u \in C$ , there exists an additive inverse  $-u \in C$  such that u + (-u) = 0. Notice that each element in B has an additive inverse which is itself because addition is done in modulo 2 : 0 + 0 = 0,

1 + 1 = 0,  $\alpha + \alpha = 0$ , and  $\alpha^2 + \alpha^2 = 0$ . Because addition of p-codons is defined componentwise, for any p-codon  $\boldsymbol{u} = (u_1 \ u_2 \ u_3)$ ,  $\boldsymbol{u} + \boldsymbol{u} = \boldsymbol{0}$  and so  $\boldsymbol{u}$ is its own additive inverse. Thus,  $-\boldsymbol{u} = -1 \cdot \boldsymbol{u} = 1 \cdot \boldsymbol{u} = \boldsymbol{u}$ . We conclude that each p-codon in *C* has an additive inverse, as required.

6. Associativity. For any  $\boldsymbol{u}, \boldsymbol{v}, \boldsymbol{w} \in C$  and  $a, b \in B$ ,

(a) 
$$\underline{u + (v + w)} = (u + v) + w$$
  
 $u + (v + w) = (u_1 \ u_2 \ u_3) + (v_1 + w_1, \ v_2 + w_2, \ v_3 + w_3)$   
 $= (u_1 + v_1 + w_1, \ u_2 + v_2 + w_2, \ u_3 + v_3 + w_3)$   
 $= (u_1 + v_1, \ u_2 + v_2, \ u_3 + v_3) + (w_1 \ w_2 \ w_3)$   
 $= (u + v) + w$ 

(b)  $\underline{a \cdot (b \cdot \boldsymbol{v})} = (a \cdot b) \cdot \boldsymbol{v}$ 

$$a \cdot (b \cdot \boldsymbol{v}) = a \cdot (b \cdot v_1, \ b \cdot v_2, \ b \cdot v_3)$$
$$= (a \cdot b \cdot v_1, \ a \cdot b \cdot v_2, \ a \cdot b \cdot v_3)$$
$$= (a \cdot b) \cdot (v_1 \ v_2 \ v_3)$$
$$= (a \cdot b) \cdot \boldsymbol{v}$$

7. Distributivity. For any  $\boldsymbol{u}, \boldsymbol{v} \in C$  and  $a, b \in B$ ,

(a) 
$$\underline{a \cdot (u + v)} = (a \cdot u) + (a \cdot v)$$
$$a \cdot (u + v) = a \cdot (u_1 + v_1, u_2 + v_2, u_3 + v_3)$$
$$= (a \cdot (u_1 + v_1), a \cdot (u_2 + v_2), a \cdot (u_3 + v_3))$$
$$= (a \cdot u_1 + a \cdot v_1, a \cdot u_2 + a \cdot v_2, a \cdot u_3 + a \cdot v_3) \quad (\text{step } 3)$$
$$= (a \cdot u_1, a \cdot u_2, a \cdot u_3) + (a \cdot v_1, a \cdot v_2, a \cdot v_3)$$
$$= a \cdot (u_1 \ u_2 \ u_3) + a \cdot (v_1 \ v_2 \ v_3)$$
$$= (a \cdot u) + (a \cdot v)$$

(b)  $(a+b) \cdot \boldsymbol{u} = (a \cdot \boldsymbol{u}) + (b \cdot \boldsymbol{u})$ 

$$(a+b) \cdot \boldsymbol{u} = (a+b) \cdot (u_1 \ u_2 \ u_3)$$
  
=  $((a+b) \cdot u_1, \ (a+b) \cdot u_2, \ (a+b) \cdot u_3)$   
=  $(a \cdot u_1 + b \cdot u_1, \ a \cdot u_2 + b \cdot u_2, \ a \cdot u_3 + b \cdot u_3)$  (step 3)  
=  $(a \cdot u_1, \ a \cdot u_2, \ a \cdot u_3) + (b \cdot u_1, \ b \cdot u_2, \ b \cdot u_3)$   
=  $a \cdot (u_1 \ u_2 \ u_3) + b \cdot (u_1 \ u_2 \ u_3)$   
=  $(a \cdot \boldsymbol{u}) + (b \cdot \boldsymbol{u})$ 

Step 3 in (a) and (b) above both follow from their respective step 2 due to field axiom 6 in section 1.2. All other steps follow from the definition of p-codon addition and multiplication of p-codons by scalars given earlier in this section.

Having shown that C, the set of all 64 p-codons is a vector space, we can refer to the p-codons as vectors. We next discuss vector space concepts that are essential to our work. In doing this, we adopt some conventions. We will write  $\mathbf{B}^3$  to mean the set of all 3-tuples whose elements are chosen from the set of scalars B. Since the vector space C fits this description, we will henceforth refer to it as  $\mathbf{B}^3$ . We also note that the arguments given to show that C satisfies the vector space axioms extend naturally to  $\mathbf{B}^n$ , where n is a positive integer. Thus,  $\mathbf{B}^n$  is a vector space under the same addition and scalar multiplication defined in sections 2.2.1 and 2.2.2 and is defined as:

$$\mathbf{B}^{n} = \{(u_{1}, u_{2}, ..., u_{n}) : u_{i} \in B \text{ for } i = 1, 2, ..., n\}$$
(2.1)

The vector space  $\mathbf{B}^n$ , where n > 3 will be needed in the next chapter and so the discussions below are done in the general context of  $\mathbf{B}^n$ . Obviously, when n = 3 we obtain the vector space C.

#### 2.3 Subspaces of $B^n$

A subspace of  $\mathbf{B}^n$  is defined as any nonempty subset of  $\mathbf{B}^n$  which is itself a vector space under the same vector addition and scalar multiplication defined on  $\mathbf{B}^n$  (given in sections 2.2.1 and 2.2.2). Let C be a subset of  $\mathbf{B}^n$ . To verify that C is a subspace, we only need to check that the closure axiom is satisfied on C. Thus, we only need to verify (2.2) below:

For any 
$$\boldsymbol{u}, \boldsymbol{v}, \boldsymbol{w} \in C$$
 and  $\lambda \in B$ , (a)  $\boldsymbol{u} + \boldsymbol{v} \in C$  (b)  $\lambda \cdot \boldsymbol{w} \in C$  (2.2)

To see why the closure axiom is sufficient to verify if C is a subspace, notice that because C is a subset of a vector space, axioms<sup>1</sup> 2, 4, 6, and 7 are automatically satisfied. Also, because every vector in C is its own additive inverse due to modulo 2 addition, axiom 5 is satisfied on C, leaving the closure (axiom 1) and additive identity (axiom 3) axioms to be checked. The additive identity axiom requires that there be an identity vector which when added to any vector, leaves the vector unchanged. Because we have defined addition as componentwise, we only need show that  $\mathbf{0} = (0, 0, ..., 0)$ is in C to satisfy the additive identity axiom. When  $\lambda = 0$  in (2.2) above, the closure

<sup>&</sup>lt;sup>1</sup>Given in the vector space definition (Definition 1)

axiom verifies the membership of **0** in *C*. Thus, any nonempty subset of  $\mathbf{B}^n$  that satisfies the closure axiom satisfies all the vector space axioms and is a subspace. Examples: {(0,0,0), (0,0,1),  $(0,0,\alpha)$ ,  $(0,0,\alpha^2)$ } and {(0,0,0),  $(0,1,\alpha)$ ,  $(0,\alpha,\alpha^2)$ ,  $(0,\alpha^2,1)$ } are two subspaces of  $\mathbf{B}^3$  since the closure axiom (stated in (2.2)) is satisfied on them.

## 2.4 Linear Combination, Span, and Linear Independence

## 2.4.1 Linear Combination

Let  $\mathbf{v_1}, \mathbf{v_2}, \ldots, \mathbf{v_k}$  be vectors in  $\mathbf{B}^n$ . A linear combination of  $\mathbf{v_1}, \mathbf{v_2}, \ldots, \mathbf{v_k}$  is another vector in  $\mathbf{B}^n$  of the form  $\lambda_1 \mathbf{v_1} + \lambda_2 \mathbf{v_2} + \cdots + \lambda_k \mathbf{v_k}$ , where  $\lambda_1, \lambda_2, \ldots, \lambda_k$  are scalars chosen from B. For example, one linear combination of the vectors  $(0, 1, \alpha)$ , (0, 0, 1), and  $(0, 1, \alpha^2)$  in  $\mathbf{B}^3$  is  $0 \cdot (0, 1, \alpha) + \alpha \cdot (0, 0, 1) + \alpha^2 \cdot (0, 1, \alpha^2) = (0, \alpha^2, 0)$ , which is another vector in  $\mathbf{B}^3$ . Another linear combination is obtained by choosing different scalars:  $1 \cdot (0, 1, \alpha) + 0 \cdot (0, 0, 1) + 0 \cdot (0, 1, \alpha^2) = (0, 1, \alpha)$ .

#### 2.4.2 Span

Let  $U = \{u_1, u_2, \dots, u_k\}$  be a nonempty subset of  $\mathbf{B}^n$ . The span of U, denoted span(U), is the set of all linear combinations of the vectors in U.

$$\operatorname{span}(U) = \{\lambda_1 \boldsymbol{u_1} + \lambda_2 \boldsymbol{u_2} + \dots + \lambda_k \boldsymbol{u_k} : \lambda_i \in B\}$$
(2.3)

For example, let  $U = \{(1,0,1)\}$  be a subset of  $\mathbf{B}^3$ . Then  $\operatorname{span}(U) = \{0 \cdot (1,0,1), 1 \cdot (1,0,1), \alpha \cdot (1,0,1), \alpha^2 \cdot (1,0,1)\} = \{(0,0,0), (1,0,1), (\alpha,0,\alpha), (\alpha^2,0,\alpha^2)\}$ . We note that  $\operatorname{span}(U) = \{(0,0,0), (1,0,1), (\alpha,0,\alpha), (\alpha^2,0,\alpha^2)\}$  is a subspace of  $\mathbf{B}^3$ . This claim is verified by observing that the addition of any two vectors in  $\operatorname{span}(U)$  yields a vector in  $\operatorname{span}(U)$  and multiplying any vector in  $\operatorname{span}(U)$  by a scalar in B yields a vector in  $\operatorname{span}(U)$ .  $\operatorname{Span}(U)$  is not unique in being a subspace. In general, the span

of any nonempty subset of  $\mathbf{B}^n$  is a subspace. This claim is proved below.

#### Proof: If U is a nonempty subset of $\mathbf{B}^n$ , span(U) is a subspace.

Let  $U = \{\boldsymbol{u}_1, \boldsymbol{u}_2, \dots, \boldsymbol{u}_k\}$  be a nonempty subset of  $\mathbf{B}^n$ . Then  $\operatorname{span}(U)$  is given in (2.3) and vectors in  $\operatorname{span}(U)$  are of the form  $\lambda_1 \boldsymbol{u}_1 + \lambda_2 \boldsymbol{u}_2 + \dots + \lambda_k \boldsymbol{u}_k$  where  $\lambda_i \in B$ . Let  $\boldsymbol{x} = \eta_1 \boldsymbol{u}_1 + \eta_2 \boldsymbol{u}_2 + \dots + \eta_k \boldsymbol{u}_k$  and  $\boldsymbol{y} = \gamma_1 \boldsymbol{u}_1 + \gamma_2 \boldsymbol{u}_2 + \dots + \gamma_k \boldsymbol{u}_k$  be any two vectors in  $\operatorname{span}(U)$ .

Then  $\mathbf{x} + \mathbf{y} = (\eta_1 + \gamma_1)\mathbf{u}_1 + (\eta_2 + \gamma_2)\mathbf{u}_2 + \dots + (\eta_k + \gamma_k)\mathbf{u}_k$ . Because  $\eta_i + \gamma_i \in B$ ,  $\mathbf{x} + \mathbf{y} \in \operatorname{span}(U)$ . Let  $\lambda$  be any scalar in B. Then  $\lambda \mathbf{x} = (\lambda \eta_1)\mathbf{u}_1 + (\lambda \eta_2)\mathbf{u}_2 + \dots + (\lambda \eta_k)\mathbf{u}_k$ . Because  $\lambda \eta_i \in B$ ,  $\lambda \cdot \mathbf{x} \in \operatorname{span}(U)$ . It has been shown that the closure axiom is satisfied on  $\operatorname{span}(U)$ . As discussed in section 2.3, this is sufficient to conclude that  $\operatorname{span}(U)$  is a subspace.

Important terminology: Let U be a nonempty subset of  $\mathbf{B}^n$  and C a subspace of  $\mathbf{B}^n$ . If  $\operatorname{span}(U) = C$ , we say that U spans C and call U a generating or spanning set of C. This terminology is appropriate because we can have only U and still obtain all vectors in the subspace C simply by forming linear combinations of the vectors in U. Using this terminology in the context of the example given earlier, the subset  $U = \{(1,0,1)\}$  of  $\mathbf{B}^3$  spans or generates the subspace C = $\{(0,0,0), (1,0,1), (\alpha,0,\alpha), (\alpha^2,0,\alpha^2)\}$  since  $\operatorname{span}(U) = C$ .

### 2.4.3 Linear Independence

A set of vectors from  $\mathbf{B}^n \{ \boldsymbol{u}_1, \boldsymbol{u}_2, \dots, \boldsymbol{u}_k \}$  is linearly independent if the only choice of scalars  $\lambda_1, \lambda_2, \dots, \lambda_k$  from B that makes  $\lambda_1 \boldsymbol{u}_1 + \lambda_2 \boldsymbol{u}_2 + \dots + \lambda_k \boldsymbol{u}_k$  equal  $\mathbf{0}$  is  $\lambda_1 = \lambda_2 = \dots = \lambda_k = 0$ . For example,  $\{(1, 0, 0), (0, 0, \alpha)\}$  is a linearly independent set of vectors from  $\mathbf{B}^3$  because  $\lambda_1$  and  $\lambda_2$  must both be 0 in order for  $\lambda_1 \cdot (1, 0, 0) + \lambda_2 \cdot (0, 0, \alpha) = (0, 0, 0) = \mathbf{0}$  to be true. In contrast,  $\{(\alpha, 0, 1), (1, 0, \alpha^2)\}$  is not linearly independent since we can choose values for  $\lambda_1$  and  $\lambda_2$  other than  $\lambda_1 = \lambda_2 = 0$  such that  $\lambda_1 \cdot (1,0,0) + \lambda_2 \cdot (0,0,\alpha) = \mathbf{0}$ . For example,  $\alpha^2 \cdot (1,0,0) + 1 \cdot (0,0,\alpha) = \mathbf{0}$ . Such a set that is not linearly independent is called linearly dependent.

Consider the subset  $U = \{(1, 0, 1)\}$  of  $\mathbf{B}^3$ . As discussed in section 2.4.2, U spans the subspace  $C = \{(0, 0, 0), (1, 0, 1), (\alpha, 0, \alpha), (\alpha^2, 0, \alpha^2)\}$ . Thus,  $\operatorname{span}(U) = C$ . Note that U is linearly independent since the only  $\lambda$  that makes  $\lambda \cdot (1, 0, 1) = \mathbf{0}$  is  $\lambda = 0$ . Now consider the set  $V = \{(1, 0, 1), (\alpha, 0, \alpha), (\alpha^2, 0, \alpha^2)\}$  which was formed by adding two vectors,  $(\alpha, 0, \alpha)$  and  $(\alpha^2, 0, \alpha^2)$ , to U. Note that V is linearly dependent since  $1 \cdot (1, 0, 1) + 1 \cdot (\alpha, 0, \alpha) + 1 \cdot (\alpha^2, 0, \alpha^2) = \mathbf{0}$ . It can be checked that  $\operatorname{span}(V) = C$ . Thus, despite having two extra vectors, V still has the same span as U. This is explained by the fact that V is linearly dependent. In any linearly dependent set, there is at least one vector which is already in the span of the other vectors in the set such that its removal does not change the span of the set<sup>2</sup>.

This illustrates an important fact which we will not formally prove: given a linearly independent set and a linearly dependent set that both span a subspace of  $\mathbf{B}^n$ , the linearly independent set will have fewer number of vectors.

## 2.5 Bases and Dimension of $B^n$

#### 2.5.1 Bases

A basis of the vector space  $\mathbf{B}^n$  is defined as a linearly independent set that spans  $\mathbf{B}^n$ . Let W be such a basis of  $\mathbf{B}^n$ . Then by definition, W spans  $\mathbf{B}^n$  and so every vector in  $\mathbf{B}^n$  can be obtained from a linear combination of the vectors in W. Thus, though  $\mathbf{B}^n$  has  $4^n$  vectors, we need only keep track of the |W| vectors in the basis W. Also, the linear independence requirement in the definition implies that |W| is

<sup>&</sup>lt;sup>2</sup>This does not hold when the linearly dependent set in question is  $\{(0,0,0)\}$ . We could however define span of the empty set to be  $\{(0,0,0)\}$  for this case.

the smallest number of vectors needed to span  $\mathbf{B}^n$ .

Several sets of vectors that are linearly independent and span  $\mathbf{B}^n$  exist. Thus, there are several possible sets to use as basis of  $\mathbf{B}^n$ , but we will use the simplest:

$$W = \{\underbrace{(1,0,\ldots,0)}_{1\text{st}}, \underbrace{(0,1,0,\ldots,0)}_{2\text{nd}}, \ldots, \underbrace{(0,\ldots,0,1)}_{n\text{th}}\}$$
(2.4)

|W| = n and each vector in W is an n-tuple. The first vector has 1 in its first position and 0 in the other positions, the second vector has 1 in its second position and 0 in the other positions, and so forth. Finally, the nth vector has 1 in its nth position and 0 in other positions. Next, we prove that W is indeed linearly independent and spans  $\mathbf{B}^{n}$ .

$$\lambda_1 \cdot (1, 0, \dots, 0) + \lambda_2 \cdot (0, 1, 0, \dots, 0) + \dots + \lambda_n \cdot (0, \dots, 0, 1) = (\lambda_1, \lambda_2, \dots, \lambda_n) \quad (2.5)$$

The linear combination of the basis vectors is given in (2.5), where  $\lambda_i \in B$ . Because the only scalars that make  $(\lambda_1, \lambda_2, \ldots, \lambda_n) = (0, 0, \ldots, 0)$  true is  $\lambda_1 = \cdots = \lambda_n = 0$ , we conclude that (2.4) is linearly independent. With the linear combination of the vectors in W given in (2.5) as  $(\lambda_1, \lambda_2, \ldots, \lambda_n)$ , it is clear that by choosing  $\lambda_1 = v_1, \lambda_2 = v_2, \ldots, \lambda_n = v_n$ , any vector  $(v_1, v_2, \ldots, v_n)$  in  $\mathbf{B}^n$  can be written as a linear combination of the vectors in W. We conclude that W spans  $\mathbf{B}^n$ . As such, we have verified that (2.4) is indeed a basis of  $\mathbf{B}^n$ .

A basis is by definition required to be linearly independent. This linear independence requirement leads to the fact that there is a unique way of expressing any vector  $\boldsymbol{v}$  in  $\mathbf{B}^n$  as a linear combination of the basis vectors. For example, in  $\mathbf{B}^3$ , (2.4) reduces to the basis  $W = \{(1,0,0), (0,1,0), (0,0,1)\}$  and the only way of expressing the vector  $(1,\alpha,\alpha^2)$  in  $\mathbf{B}^3$  as a linear combination of the basis is  $1 \cdot (1,0,0) + \alpha \cdot (0,1,0) + \alpha^2 \cdot (0,0,1)$ . Though we prove this fact in the context of our chosen basis (given in (2.4)), it holds for every other basis of  $\mathbf{B}^n$ . Let  $\boldsymbol{v} \in \mathbf{B}^n$ . Assume that  $\boldsymbol{v}$  can be written as two different linear combinations of the basis vectors. Thus  $\boldsymbol{v} = \lambda_1(1, 0, \dots, 0) + \dots + \lambda_n(0, \dots, 0, 1)$  and  $\boldsymbol{v} = \eta_1 \cdot (1, 0, \dots, 0) + \dots + \eta_n \cdot (0, \dots, 0, 1)$ . Subtracting these two expressions results in

$$\mathbf{0} = (\lambda_1 - \eta_1) \cdot (1, 0, \dots, 0) + \dots + (\lambda_n - \eta_n) \cdot (0, \dots, 0, 1)$$
(2.6)

Recalling that  $\{(1, 0, ..., 0), (0, 1, 0, ..., 0), ..., (0, ..., 0, 1)\}$  is a linearly independent set, (2.6) holds only if  $(\lambda_1 - \eta_1) = (\lambda_2 - \eta_2) = ... = (\lambda_n - \eta_n) = 0$ . This implies that  $\lambda_1 = \eta_1$ ,  $\lambda_2 = \eta_2$ , ...,  $\lambda_n = \eta_n$ . We thus conclude that the two supposedly different linear combinations are actually the same and so the expression of  $\boldsymbol{v}$  as a linear combination of the basis vectors is unique.

#### 2.5.2 Dimension

Dimension of a vector space is defined as the number of vectors in its basis.  $W = \{(1, 0, ..., 0), (0, 1, 0, ..., 0), ..., (0, ..., 0, 1)\}$  is a basis of  $\mathbf{B}^n$ . |W| = n and so the dimension of  $\mathbf{B}^n$  is n. As mentioned earlier, several other basis of  $\mathbf{B}^n$  exist. It can however be shown that each such basis has n vectors<sup>3</sup>. As such, the dimension of  $\mathbf{B}^n$  is independent of which basis of  $\mathbf{B}^n$  is being considered. We will write dim(V)to mean the dimension of some vector space V.

#### 2.6 Orthogonal Complements

Let  $\boldsymbol{u} = (u_1, u_2, \dots, u_n)$  and  $\boldsymbol{w} = (w_1, w_2, \dots, w_n)$  be vectors in  $\mathbf{B}^n$ . The inner product of  $\boldsymbol{u}$  and  $\boldsymbol{w}$ , denoted  $\boldsymbol{u} \cdot \boldsymbol{w}$ , is defined as  $\boldsymbol{u} \cdot \boldsymbol{w} = u_1 w_1 + u_1 w_1 + \dots + u_n w_n$ . When  $\boldsymbol{u} \cdot \boldsymbol{w} = 0$  we say  $\boldsymbol{u}$  and  $\boldsymbol{v}$  are orthogonal to each other.

Let A be a nonempty subset of  $\mathbf{B}^n$ . The orthogonal complement of A, denoted  $A^{\perp}$ ,

<sup>&</sup>lt;sup>3</sup>Proof of this statement is not given as it's not directly relevant to our work.

is the set of all vectors in  $\mathbf{B}^n$  that are orthogonal to every vector in A. Mathematically,

$$A^{\perp} = \{ \boldsymbol{x} \in \mathbf{B}^n : \boldsymbol{x} \cdot \boldsymbol{a} = 0 \text{ for all } \boldsymbol{a} \in A \}$$

$$(2.7)$$

The orthogonal complement of any nonempty subset of  $\mathbf{B}^n$  is a subspace. We offer a proof below. As discussed in section 2.3, the closure axiom alone is sufficient to prove that a subset of a vector space is a subspace.

# Proof: Given any nonempty subset $A \subseteq \mathbf{B}^n$ , $A^{\perp}$ is a subspace.

Let  $\boldsymbol{u}, \boldsymbol{v}, \boldsymbol{w} \in A^{\perp}, \boldsymbol{a} \in A$ , and  $\lambda \in \mathbf{B}$ . Then  $(\boldsymbol{u} + \boldsymbol{v}) \cdot \boldsymbol{a} = \boldsymbol{u} \cdot \boldsymbol{a} + \boldsymbol{v} \cdot \boldsymbol{a} = 0$ , which implies that  $(\boldsymbol{u} + \boldsymbol{v}) \in A^{\perp}$ . Also,  $(\lambda \boldsymbol{w}) \cdot \boldsymbol{a} = \lambda(\boldsymbol{w} \cdot \boldsymbol{a}) = 0$ , which implies  $(\lambda \boldsymbol{w}) \in A^{\perp}$ . Because  $(\boldsymbol{u} + \boldsymbol{v}) \in A^{\perp}$  and  $(\lambda \boldsymbol{w}) \in A^{\perp}$ , the closure axiom is satisfied and we conclude that  $A^{\perp}$  is a subspace.

The fact proved above will be useful to our work for the case where A is not only a subset but also a subspace. As such, we take A as a subspace in the discussions that follow. Being subspaces, A and  $A^{\perp}$  each have a set of basis vectors and hence a dimension. A well-known linear algebra result is that the dimensions of a subspace plus the dimension of its orthogonal complement equals the dimension of the vector space they are subsets of. In our case, since  $dim(\mathbf{B}^n) = n$  and  $A, A^{\perp} \subseteq \mathbf{B}^n$ ,

$$\dim(A) + \dim(A^{\perp}) = n \tag{2.8}$$

Let G and H be matrices whose rows are the basis vectors of A and  $A^{\perp}$ , respectively. Such matrices are called generator matrices. If dim(A) = k, G is a  $k \times n$ matrix and H is an  $(n - k) \times n$  matrix. Note that G and H both have n columns because vectors in  $\mathbf{B}^n$  are of length n. Because every vector in A is orthogonal to every other vector in  $A^{\perp}$ , the following equations hold and will be used in chapter 3.

$$\boldsymbol{a}H^T = \boldsymbol{0} \text{ for any } \boldsymbol{a} \in A \tag{2.9}$$

$$GH^T = \mathbf{0} \tag{2.10}$$

#### Chapter 3

## TOLERATING SUBSTITUTION MUTATIONS

#### 3.1 Introduction

In the previous chapter, we noted that the algebraic representations of naturallyoccurring codons form a vector space, which we denoted as  $\mathbf{B}^3$ . One vector in  $\mathbf{B}^3$  is  $(0\ 1\ \alpha)$ , which is the algebraic representation of ATG, which in a DNA sense strand, represents Methionine. Consider a substitution mutation event which changes the p-codon  $(0\ 1\ \alpha)$  to  $(0\ \alpha^2\ \alpha)$ , which is the algebraic representation of ACG and will ultimately be translated as Threonine, instead of the originally intended Methionine.

Due to the degeneracy of the genetic code, some substitution mutations result in another codon that still codes for the same amino acid. Take Cysteine as an example. In a DNA sense strand, Cysteine is coded for by both TGT (p-codon  $(1 \ \alpha \ 1)$ ) and TGC (p-codon  $(1 \ \alpha \ \alpha^2)$ ). If the last nucleotide is mutated from C to T or T to C, there will be no subsequent change in the resulting amino acid after translation. Any other substitution mutation will result in a codon which codes for an unintended amino acid.

Our goal in this chapter is to find a set of 64 'unnatural' codons that are each able to tolerate an arbitrarily chosen number of substitution mutations to any nucleotide. That is, codons that can sustain some chosen number of mutation errors and still guarantee that upon translation, the amino acid sequence will not differ from that originally intended. To do this, we will employ ideas from coding theory. These ideas sum up to building some sort of redundancy into the naturally occuring codons. As such, the codons we develop will necessarily each have more than 3 nucleotides and will be chosen from the vector space  $\mathbf{B}^n$ , where n > 3. Exactly which n we will use will be noted later. We introduce a few concepts that will be needed.

## 3.2 Hamming Distance and Weight

## 3.2.1 Hamming Distance

The Hamming distance (Hamming, 1950) between two vectors  $\boldsymbol{u} = (u_1, u_2, ..., u_n)$ and  $\boldsymbol{v} = (v_1, v_2, ..., v_n)$  in  $\mathbf{B}^n$  is the number of places where  $u_i \neq v_i$ . For example, the Hamming distance between  $\boldsymbol{u} = (0 \ 1 \ \alpha)$  and  $\boldsymbol{v} = (1 \ 0 \ \alpha)$ , is 2 since they differ in two places. We will denote the Hamming distance between  $\boldsymbol{u}$  and  $\boldsymbol{v}$  by  $d(\boldsymbol{u}, \boldsymbol{v})$ . Note that  $d(\boldsymbol{u}, \boldsymbol{v})$  can only take the values 0, 1, 2, ..., n. Also note that  $d(\boldsymbol{u}, \boldsymbol{v}) = d(\boldsymbol{v}, \boldsymbol{u})$ . Another important property of the Hamming distance is that it satisfies the triangle inequality<sup>1</sup>. Thus, given vectors  $\boldsymbol{u}, \boldsymbol{v}, \boldsymbol{w}, \quad d(\boldsymbol{u}, \boldsymbol{w}) \leq d(\boldsymbol{u}, \boldsymbol{v}) + d(\boldsymbol{v}, \boldsymbol{w})$ .

## 3.2.2 Hamming Weight

Let  $\boldsymbol{u} = (u_1, u_2, ..., u_n)$  be some vector in  $\mathbf{B}^n$ . Its Hamming weight, which we will denote as  $wt(\boldsymbol{u})$ , is the number of places where  $u_i \neq 0$ . In other words,  $wt(\boldsymbol{u}) = d(\boldsymbol{u}, 0)$ . For example, the weight of  $\boldsymbol{u} = (0 \ 1 \ \alpha^2)$  is  $wt(\boldsymbol{u}) = 2$ .

#### 3.2.3 Important Relationship between Hamming Distance and Weight

Consider vectors  $\boldsymbol{u}$  and  $\boldsymbol{v}$  in  $\mathbf{B}^n$  and let  $\boldsymbol{w} = \boldsymbol{u} + \boldsymbol{v}$ . Then, for i = 1, 2, ..., n,  $\boldsymbol{w}_i = 0$  if and only if  $\boldsymbol{u}_i = \boldsymbol{v}_i$  since we use base 2 addition. In other words, the coordinates of  $\boldsymbol{w}$  will be non-zero only at coordinates where  $\boldsymbol{u}$  and  $\boldsymbol{v}$  differ. As such, we can find the Hamming distance between two vectors by first adding them and then

<sup>&</sup>lt;sup>1</sup>A proof of this will not be directly relevant to our work. As such, it is not given.

taking the Hamming weight of the result. Mathematically, we write

$$d(\boldsymbol{u}, \boldsymbol{v}) = wt(\boldsymbol{u} + \boldsymbol{v}) \tag{3.1}$$

The above relationship becomes very useful when we desire to find the Hamming distance between all unique pairs of vectors in a vector space (or subspace) and report the minimum, called the minimum Hamming distance and denoted  $d_{min}$ . Take  $\mathbf{B}^n$ as the vector space of interest. Since  $\mathbf{B}^n$  has  $4^n$  vectors, if we naively tried to find  $d_{min}(\mathbf{B}^n)$  without the relationship given in the previous paragraph, we will have to find the Hamming distance between  $(4^n * (4^n - 1))/2$  pairs of vectors and report the minimum as  $d_{min}(\mathbf{B}^n)$ . By making use of the relationship  $d(\mathbf{u}, \mathbf{v}) = wt(\mathbf{u} + \mathbf{v})$ , we need only find the weight of each of the  $4^n - 1$  non-zero vectors and report the minimum as  $d_{min}(\mathbf{B}^n)$ . Thus, the claim is that given some vector space V,

$$d_{min}(V) = wt_{min}(V) \tag{3.2}$$

where  $wt_{min}(V)$  is the minimum of the weights of the non-zero vectors in V.

To see why (3.2) is true, first note that if  $v_1$  and  $v_2$  are different vectors in the vector space V, then their sum,  $\boldsymbol{w} = \boldsymbol{v_1} + \boldsymbol{v_2}$ , is a non-zero vector also in V. Thus, if we put all the vectors in V into all possible pairs of two, say  $(\boldsymbol{v_i}, \boldsymbol{v_j})$ , where  $i \neq j$ , then the sum of each pair will yield some  $\boldsymbol{w}$  such that the set of all the  $\boldsymbol{w}$ s will be equal to the non-zero vectors in V. Combining this with (3.1), it follows that if we want to find  $d_{min}(V)$ , instead of finding the Hamming distance between all possible pairs of unique vectors and checking the minimum, we can simply check the weight of each non-zero vector in V and record the minimum as our  $d_{min}(V)$  because each non-zero vector in V is the result of summing some two different vectors in V.

#### 3.3 Tolerating t Substitution Mutations per Codon

Our goal has been to find some 64 'unnatural' codons to used in place of the naturally-occurring codons in the construction of genes such that the genes are capable of tolerating an arbitrarily chosen number of substitution mutations per codon. We will put the p-codons of these 64 unnatural codons in a set  $M = \{m_1, m_2, \ldots, m_{64}\}$ . To make use of equation (3.2), we will require that M be a subspace of  $\mathbf{B}^n$ , where n > 3. This requirement also simplifies our work greatly because as a subspace, M is completely defined by its basis. So, we need only find the basis vectors of M. All the other vectors can be obtained via linear combinations of the basis vectors.

#### 3.3.1 Substitution Mutation as Vector Addition

Consider the p-codon  $\boldsymbol{u} = (0 \ 1 \ \alpha) \in \mathbf{B}^3$  and assume that some substitution mutation modified  $\boldsymbol{u}$  into  $\tilde{\boldsymbol{u}} = (1 \ 0 \ \alpha)$ . Since  $\tilde{\boldsymbol{u}} = \boldsymbol{u} + (1 \ 1 \ 0)$ , we can refer to  $(1 \ 1 \ 0)$ as the <u>mutation vector</u>. In general, for any  $\boldsymbol{u} \in \mathbf{B}^n$ , if a mutation event results in  $\boldsymbol{u}$  becoming  $\tilde{\boldsymbol{u}}$  such that  $\tilde{\boldsymbol{u}} = \boldsymbol{u} + \boldsymbol{w}$ , where  $\boldsymbol{w} \in \mathbf{B}^3$ , we shall call  $\boldsymbol{w}$  the mutation vector and write  $\tilde{m}(\boldsymbol{u}, \boldsymbol{w})$  to mean the vector  $\tilde{\boldsymbol{u}}$  that results from mutating  $\boldsymbol{u}$  with the mutation vector  $\boldsymbol{w}$ . Thus,

$$\widetilde{m}(\boldsymbol{u}, \boldsymbol{w}) = \widetilde{\boldsymbol{u}} = \boldsymbol{u} + \boldsymbol{w}$$
(3.3)

We now present a theorem that is essential to finding M. It gives the restriction in terms of Hamming distance that must be placed on the vectors in M to ensure that they have properties (also given in the theorem) necessary to tolerate mutations.

**Theorem 3.1.** Let  $M = \{m_1, m_2, ..., m_{64}\}$  be a subspace of  $B^n$ , where n > 3. If  $d_{min}(M) \ge 2t+1$ , then if t or less substitution mutations occur in any codon  $m_i \in M$  so that it becomes  $\widetilde{m_i}$ , then

- 1.  $d(\boldsymbol{m_i}, \widetilde{\boldsymbol{m_i}}) \leq t$
- 2.  $\widetilde{\boldsymbol{m}_i} \notin M$
- 3.  $d(\widetilde{m_i}, m_j) > t$  for any  $m_j \in M$ , where  $m_j \neq m_i$

#### <u>Proof</u>

- 1. Since t or less substitution mutations occurred in  $m_i$ , it follows that  $\widetilde{m}_i$  and  $m_i$  differ in at most t places and hence  $d(m_i, \widetilde{m}_i) \leq t$ .
- 2. It takes a mutation vector of weight  $\geq 2t + 1$  to turn any codon in M into a different codon in M since  $d_{min}(M) \geq 2t + 1$ . Since  $\leq t$  substitution mutations occurred in  $\boldsymbol{m_i}, \widetilde{\boldsymbol{m_i}}$  cannot possibly be in M.
- 3. We prove this by contradiction so assume that d(m̃<sub>i</sub>, m<sub>j</sub>) ≤ t. We know from
  1) that d(m<sub>i</sub>, m̃<sub>i</sub>) ≤ t. Then, d(m<sub>i</sub>, m̃<sub>i</sub>) + d(m̃<sub>i</sub>, m<sub>j</sub>) ≤ 2t. Using the triangle inequality, we have that d(m<sub>i</sub>, m<sub>j</sub>) ≤ d(m<sub>i</sub>, m̃<sub>i</sub>) + d(m̃<sub>i</sub>, m<sub>j</sub>) ≤ 2t. This is a contradiction, since we know that d(m<sub>i</sub>, m<sub>j</sub>) ≥ 2t+1 because d<sub>min</sub>(M) ≥ 2t+1. We conclude that d(m̃<sub>i</sub>, m<sub>j</sub>) > t, as desired. ■

Using theorem 3.1, our task of developing codons capable of tolerating t mutations is simplified and can be stated as: find some set of vectors  $M = \{m_1, m_2, \ldots, m_{64}\}$ such that  $d_{min}(M) \ge 2t + 1$ . Upon finding such an M, we are guaranteed that any mutation event with a mutation vector of weight t or less that occurs in any codon in M will result in a new codon that is none of the vectors in M but still closest in Hamming distance to the original codon and can thus be decoded as the original codon. We next use the concept of orthogonal complements previously discussed in section 2.6.

## 3.3.2 Finding Generator Matrices for $M^{\perp}$ and M

Let  $H = [h_1, h_2, \ldots, h_n]$ , where  $h_i$  is the *i*th column of H, be the generator matrix for  $M^{\perp}$ . Then from equation (2.9) we know that  $\mathbf{a}H^T = \mathbf{0}$  for  $\mathbf{a} \in M$ . Thus, if  $\mathbf{a} = (a_1, a_2, \ldots, a_n), \quad a_1h_1 + a_2h_2 + \ldots a_nh_n = \mathbf{0}.$  Let k = 2t + 1. According to equation (3.2),  $d_{min}(M) = k$  implies that  $wt_{min}(M) = k$ , which has two implications that combine with the fact that  $a_1h_1 + a_2h_2 + \ldots a_nh_n = \mathbf{0}$  to tell us how to construct H. These are discussed below:

- 1.  $wt_{min}(M) = k$  implies that at least one vector in M has weight k. Let  $m = (m_1, m_2, ..., m_k) \in M$  be such a vector with weight k and let  $m_{i1}, m_{i2}, ..., m_{ik}$  be its k non-zero coordinates. Then  $m_1h_1 + m_2h_2 + ... + m_kh_k = m_{i1}h_{i1} + m_{i2}h_{i2} + ... + m_{ik}h_{ik} = \mathbf{0}$ . This implies that the generator matrix H has some k columns that are linearly dependent.
- 2.  $wt_{min}(M) = k$  implies that there is no vector in M with weight k 1. This means there is no  $m \in M$  with the k - 1 non-zero coordinates  $m_{i1}, m_{i2}$  $\dots, m_{i(k-1)}$  such that  $m_{i1}h_{i1} + m_{i2}h_{i2} + \dots + m_{i(k-1)}h_{i(k-1)} = 0$ . This implies that any k - 1 columns of H are linearly independent.

The two discussions above tell us how the requirement of  $d_{min}(M) = k$  on M translates to restrictions on the generator matrix H of its orthogonal complement  $M^{\perp}$ . These restrictions tell us how to construct H: construct H such that 1) some k columns are linearly dependent and 2) any k - 1 columns are linearly independent, where k = 2t + 1 and t is the arbitrary number of substitution mutations we want to tolerate per codon. Once H is constructed, we can then find the generator matrix of M, call it G, from the equation  $GH^T = \mathbf{0}$ , which was given in section 2.6. Note that G completely tells us the vectors in  $M = \{m_1, m_2, \ldots, m_{64}\}$  since its rows are a basis of M. The approach is demonstrated in the next section.

#### 3.4 Example: Tolerating 1 Substitution Mutation per Codon

Here, we use the ideas discussed in the previous sections, particularly section 3.3.2, to find an example set of 64 alternative codons  $M = \{m_1, m_2, \ldots, m_{64}\}$ , where  $M \subseteq \mathbf{B}^n$  and n > 3, such that a gene constructed from these codons is guaranteed to tolerate any 1 substitution mutation per codon<sup>2</sup>.

We know that if  $d_{min}(M) = wt_{min}(M) \ge 2t + 1$ , then M will have the structure necessary for tolerating any mutation with mutation vector of weight less than or equal to t. For our example, t = 1 and so we will form an M with  $d_{min}(M) = wt_{min}(M) = 3$ . Using ideas discussed in the prior section, this requires that we find G, the generator matrix of M, from  $GH^T = \mathbf{0}$  where H, the generator matrix of  $M^{\perp}$ , is a matrix with two restrictions: 1) some 3 columns are linearly dependent and 2) any 2 columns of H are linearly independent.

As noted in section 2.6, if G has size  $k \times n$ , H will have size  $(n - k) \times n$ . To find H and then G, we need to know what n and k are. M is a subspace of 64 vectors, one to 'replace' each naturally-occurring codon. Given that we have 4 scalars,  $0, 1, \alpha$ , and  $\alpha^2$ , M must have 3 basis vectors so that all possible linear combinations of the 3 basis vectors yields 64 vectors. By definition, the rows of G are the basis vectors of M. Thus, G has k = 3 rows. We will choose the subspace M from the vector space  $\mathbf{B}^5$  and so n = 5. This choice will be explained in the next section. Given n = 5 and k = 3, G and H have sizes  $3 \times 5$  and  $2 \times 5$ , respectively. Let's construct H.

Consider the set of vectors  $A = \{(1, \alpha^2), (1, \alpha), (\alpha, \alpha), (1, 0), (0, 1)\}$ . Since  $1 \cdot (\alpha, \alpha) + \alpha \cdot (1, 0) + \alpha \cdot (0, 1) = \mathbf{0}, (\alpha, \alpha), (1, 0)$  and (0, 1) are linearly dependent. Also, one can verify that any two vectors from A are linearly independent. The vectors in

 $<sup>^{2}</sup>$ Thus, any codon can tolerate a mutation with mutation vector of weight 1.

A thus satisfy the two requirements on the the column vectors of H. As such, we can use these vectors as the column vectors of H, as shown below.

$$H = \begin{bmatrix} 1 & 1 & \alpha & 1 & 0 \\ \alpha^2 & \alpha & \alpha & 0 & 1 \end{bmatrix}$$
(3.4)

Note that H is not unique. Any matrix of size  $2 \times 5$  with column vectors that satisfy the two requirements mentioned earlier can be used. Also note that the submatrix formed from the last two columns is the  $2 \times 2$  identity matrix  $I_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ . This choice was made so we could make use of the fact that if  $H = [P \mid I_{n-k}]$ , where  $I_{n-k}$  is the  $n - k \times n - k$  identity matrix, then  $G = [I_k \mid P^T]$  satisfies  $GH^T = \mathbf{0}$ . Using this fact and H given in equation (3.4), we have G as:

$$G = \begin{bmatrix} 1 & 0 & 0 & 1 & \alpha^2 \\ 0 & 1 & 0 & 1 & \alpha \\ 0 & 0 & 1 & \alpha & \alpha \end{bmatrix}$$
(3.5)

The rows of G form a basis of the set  $M = \{m_1, m_2, \dots, m_{64}\}$ . Thus, every  $m_i \in M$ can be obtained from a linear combination of the form  $m_i = \lambda_1 \cdot (1, 0, 0, 1, \alpha^2) + \lambda_2 \cdot (0, 1, 0, 1, \alpha) + \lambda_3 \cdot (0, 0, 1, \alpha, \alpha)$ , where  $\lambda_i$  are scalars. This linear combination can be written as the matrix multiplication given below:

$$\boldsymbol{m_i} = (\lambda_1, \lambda_2, \lambda_3) \cdot \boldsymbol{G} \tag{3.6}$$

Thus, G maps each vector  $(\lambda_1, \lambda_2, \lambda_3) \in \mathbf{B}^3$  uniquely to a vector  $(\omega_1, \omega_2, \omega_3, \omega_4, \omega_5) \in M$ . In essence, G tells us which unnatural codon of length 5 to use in place of each naturally-occuring codon of length 3 such that genes constructed from the unnatural

codons can tolerate up to one substitution mutation per codon. For example, the pcodon of ATG is  $(0 \ 1 \ \alpha)$ .  $(0 \ 1 \ \alpha)$ .  $G = (0 \ 1 \ \alpha \ \alpha \ 1)$ , which is the p-codon for ATGGT. We will thus use the synthetic codon ATGGT in place of ATG. Table 3.1 shows results from evaluating equation (3.6) for all p-codons of the naturally occuring codons. It shows how all 64 naturally-occuring codons should be mapped to the mutationtolerant (mutol) codons of length 5 in set M such that genes constructed from the mutol codons will be tolerant of up to one substitution mutation per codon. Since we represented both thymine and uracil with scalar 1 as outlined in table 1.1, T in the codons column of table 3.1 can simply be replaced with U to get the corresponding RNA codons.

Three observations can be made about the 5-letter codons in table 3.1 we claim will confer mutation-tolerance upon genes constructed from them: 1) each of the 5-letter unnatural codons can have any 1 of its nucleotides replaced with any other nucleotide without turning into any of the other 63 codons. 2) upon substituting any 1 nucleotide in any of the 5-letter codons, the resulting 5-letter codon is still most similar to the originally intended codon in terms of Hamming Distance. 3) any 1 nucleotide substitution in any of the 5-letter codons results in a 5-letter codon that is not in Table 3.1. These observations are desirable and are simply consequences of theorem 3.1 for the case t = 1. Observation 3 however poses a challenge, which is addressed in the next section.

p-e	codons	DNA Codons
$(0 \ 0 \ 0)$	$(0 \ 0 \ 0 \ 0 \ 0)$	AAA AAAAA
$(0 \ 0 \ 1)$	$(0 \ 0 \ 1 \ a \ a)$	AAT AATGG
$(0 \ 0 \ a)$	$(0 \ 0 \ a \ a^2 \ a^2)$	AAG AAGCC
$(0 \ 0 \ a^2)$	$(0 \ 0 \ a^2 \ 1 \ 1)$	AAC AACTT
$(0\ 1\ 0)$	$(0\ 1\ 0\ 1\ a)$	ATA ATATG
$(0\ 1\ 1)$	$(0 \ 1 \ 1 \ a^2 \ 0)$	ATT ATTCA
$(0 \ 1 \ a)$	$(0\ 1\ a\ a\ 1)$	ATG ATGGT
$(0 \ 1 \ a^2)$	$(0 \ 1 \ a^2 \ 0 \ a^2)$	ATC ATCAC
$(0 \ a \ 0)$	$(0 \ a \ 0 \ a \ a^2)$	AGA AGAGC
$(0 \ a \ 1)$	$(0\ a\ 1\ 0\ 1)$	AGT AGTAT
$(0 \ a \ a)$	$(0 \ a \ a \ 1 \ 0)$	AGG AGGTA
$(0 \ a \ a^2)$	$(0  a  a^2  a^2  a)$	AGC AGCCG
$(0 \ a^2 \ 0)$	$(0 \ a^2 \ 0 \ a^2 \ 1)$	ACA ACACT
$(0 \ a^2 \ 1)$	$(0 \ a^2 \ 1 \ 1 \ a^2)$	ACT ACTTC
$(0 \ a^2 \ a)$	$(0  a^2  a  0  a)$	ACG ACGAG
$(0 \ a^2 \ a^2)$	$(0 \ a^2 \ a^2 \ a \ 0)$	ACC ACCGA
$(1 \ 0 \ 0)$	$(1 \ 0 \ 0 \ 1 \ a^2)$	TAA TAATC
$(1 \ 0 \ 1)$	$(1 \ 0 \ 1 \ a^2 \ 1)$	TAT TATCT
$(1 \ 0 \ a)$	$(1 \ 0 \ a \ a \ 0)$	TAG TAGGA
$(1 \ 0 \ a^2)$	$(1 \ 0 \ a^2 \ 0 \ a)$	TAC TACAG
$(1 \ 1 \ 0)$	$(1\ 1\ 0\ 0\ 1)$	ΤΤΑ ΤΤΑΑΤ
$(1\ 1\ 1)$	$(1\ 1\ 1\ a\ a^2)$	TTT TTTGC

 Table 3.1: Mapping to Alternative Codons

р	-codons	DNA Codons
$(1 \ 1 \ a)$	$(1 \ 1 \ a \ a^2 \ a)$	TTG TTGCG
$(1 \ 1 \ a^2)$	$(1 \ 1 \ a^2 \ 1 \ 0)$	TTC TTCTA
$(1 \ a \ 0)$	$(1 \ a \ 0 \ a^2 \ 0)$	TGA TGACA
$(1 \ a \ 1)$	$(1 \ a \ 1 \ 1 \ a)$	TGT TGTTG
$(1 \ a \ a)$	$(1 \ a \ a \ 0 \ a^2)$	TGG TGGAC
$(1 \ a \ a^2)$	$(1 \ a \ a^2 \ a \ 1)$	TGC TGCGT
$(1 \ a^2 \ 0)$	$(1 \ a^2 \ 0 \ a \ a)$	TCA TCAGG
$(1 \ a^2 \ 1)$	$(1 \ a^2 \ 1 \ 0 \ 0)$	TCT TCTAA
$(1 \ a^2 \ a)$	$(1 \ a^2 \ a \ 1 \ 1)$	TCG TCGTT
$(1 \ a^2 \ a^2)$	$(1 \ a^2 \ a^2 \ a^2 \ a^2)$	TCC TCCCC
$(a \ 0 \ 0)$	$(a \ 0 \ 0 \ a \ 1)$	GAA GAAGT
$(a \ 0 \ 1)$	$(a \ 0 \ 1 \ 0 \ a^2)$	GAT GATAC
$(a \ 0 \ a)$	$(a \ 0 \ a \ 1 \ a)$	GAG GAGTG
$(a \ 0 \ a^2)$	$(a \ 0 \ a^2 \ a^2 \ 0)$	GAC GACCA
$(a \ 1 \ 0)$	$(a \ 1 \ 0 \ a^2 \ a^2)$	GTA GTACC
$(a \ 1 \ 1)$	$(a \ 1 \ 1 \ 1 \ 1)$	GTT GTTTT
$(a \ 1 \ a)$	$(a \ 1 \ a \ 0 \ 0)$	GTG GTGAA
$(a \ 1 \ a^2)$	$(a\ 1\ a^2\ a\ a)$	GTC GTCGG
$(a \ a \ 0)$	$(a \ a \ 0 \ 0 \ a)$	GGA GGAAG
$(a \ a \ 1)$	$(a \ a \ 1 \ a \ 0)$	GGT GGTGA
$(a \ a \ a)$	$(a \ a \ a \ a^2 \ 1)$	GGG GGGCT
$(a \ a \ a^2)$	$(a \ a \ a^2 \ 1 \ a^2)$	GGC GGCTC
$(a \ a^2 \ 0)$	$(a \ a^2 \ 0 \ 1 \ 0)$	GCA GCATA

 Table 3.1: Mapping to Alternative Codons

p-o	codons	DNA	Codons
$(a \ a^2 \ 1)$	$(a \ a^2 \ 1 \ a^2 \ a)$	GCT	GCTCG
$(a \ a^2 \ a)$	$(a \ a^2 \ a \ a \ a^2)$	GCG	GCGGC
$(a \ a^2 \ a^2)$	$(a \ a^2 \ a^2 \ 0 \ 1)$	GCC	GCCAT
$(a^2 \ 0 \ 0)$	$(a^2 \ 0 \ 0 \ a^2 \ a)$	CAA	CAACG
$(a^2 \ 0 \ 1)$	$(a^2 \ 0 \ 1 \ 1 \ 0)$	CAT	CATTA
$(a^2 \ 0 \ a)$	$(a^2 \ 0 \ a \ 0 \ 1)$	CAG	CAGAT
$(a^2 \ 0 \ a^2)$	$(a^2 \ 0 \ a^2 \ a \ a^2)$	CAC	CACGC
$(a^2 \ 1 \ 0)$	$(a^2 \ 1 \ 0 \ a \ 0)$	CTA	CTAGA
$(a^2 \ 1 \ 1)$	$(a^2 \ 1 \ 1 \ 0 \ a)$	CTT	CTTAG
$(a^2 \ 1 \ a)$	$(a^2 \ 1 \ a \ 1 \ a^2)$	CTG	CTGTC
$(a^2 \ 1 \ a^2)$	$(a^2 \ 1 \ a^2 \ a^2 \ 1)$	CTC	CTCCT
$(a^2 \ a \ 0)$	$(a^2 \ a \ 0 \ 1 \ 1)$	CGA	CGATT
$(a^2 \ a \ 1)$	$(a^2 \ a \ 1 \ a^2 \ a^2)$	CGT	CGTCC
$(a^2 \ a \ a)$	$(a^2 \ a \ a \ a \ a)$	CGG	CGGGG
$(a^2 \ a \ a^2)$	$(a^2 \ a \ a^2 \ 0 \ 0)$	$\operatorname{CGC}$	CGCAA
$(a^2 \ a^2 \ 0)$	$(a^2 \ a^2 \ 0 \ 0 \ a^2)$	CCA	CCAAC
$(a^2 \ a^2 \ 1)$	$(a^2 \ a^2 \ 1 \ a \ 1)$	CCT	CCTGT
$(a^2 a^2 a)$	$(a^2 \ a^2 \ a \ a^2 \ 0)$	CCG	CCGCA
$(a^2 \ a^2 \ a^2)$	$(a^2 \ a^2 \ a^2 \ 1 \ a)$	CCC	CCCTG

 Table 3.1: Mapping to Alternative Codons

#### 3.4.1 Codon Degeneracy

We first demonstrate observation 3 made in the previous section. Let  $\boldsymbol{x} =$  ATGGT, which is the 5-letter mutation-tolerant variant of the naturally-occuring codon ATG. Assume a mutation event that substitutes the third nucleotide of  $\boldsymbol{x}$ , guanine, with adenine. The resulting codon will be  $\tilde{\boldsymbol{x}}$ =ATAGT. We can confirm that ATAGT is indeed not one of the 64 5-letter codons given in table 3.1. This is also true of any 5-letter codon which results from a 1 nucleotide substitution mutation of any of the other 63 codons.

This poses the following challenge. Assume we have a gene, call it K, constructed using the mutation-tolerant 5-letter codons given in table 3.1 and artificial tRNAs that recognize the 64 corresponding 5-letter RNA codons. Further assume that some 1 nucleotide substitution has occurred in a codon  $\boldsymbol{x}$  of K to become  $\tilde{\boldsymbol{x}}$ . We know that  $\tilde{\boldsymbol{x}}$  is not in table 3.1. Thus, the translation of gene K cannot progress beyond the RNA equivalent of codon  $\tilde{\boldsymbol{x}}$  since its pattern is not one of the 64 patterns recognized by the the artificial tRNAs. In fact, there will be  $4^5 - 64$  RNA codons not recognized by the artificial tRNAs.

The solution is obvious. We ought to have artificial tRNAs that recognize RNA codons beyond those for the 64 5-letter codons given in 3.1. For example, instead of having an artificial tRNA decode only AUGGU as methionine (the amino acid coded for by the naturally-occurring AUG) we will have several artificial tRNAs that decode several 5-letter RNA codons, including the 'original' AUGGU, as methionine. Our approach to doing that is as follows.

As discussed in section, nucleotide substitutions can be described as the addition of a mutation vector to a p-codon. As such, being tolerant of single nucleotide substitutions means that the 5-letter codons in table 3.1 tolerate any mutation event whose

mutation vector	is	one	of	the rows	3	in	Z	below.
-----------------	----	-----	----	----------	---	----	---	--------

	г				
	0	0	0	0	0
	1	0	0	0	0
	0	1	0	0	0
	0	0	1	0	0
	0	0	0	1	0
	0	0	0	0	1
	$\alpha$	0	0	0	0
Z =	0	$\alpha$	0	0	0
-	0	0	$\alpha$	0	0
	0	0	0	$\alpha$	0
	0	0	0	0	$\alpha$
	$\alpha^2$	0	0	0	0
	0	$\alpha^2$	0	0	0
	0	0	$\alpha^2$	0	0
	0	0	0	$\alpha^2$	0
	0	0	0	0	$\alpha^2$
					-

Using this fact, given any 5-letter codon  $\boldsymbol{x}$  from table 3.1, our approach to deciding which 5-letter codons to decode in the same manner as  $\boldsymbol{x}$  is as follows: add the p-codon of  $\boldsymbol{x}$  to each row of Z, then have artificial tRNAs that decode all resulting codons as you would  $\boldsymbol{x}$ . Combining this approach with equation (3.6), we give equation (3.8) which takes as input the p-codon  $(\lambda_1, \lambda_2, \lambda_3)$  of a naturally-occurring 3-letter codon and yields K, a 16 × 1 matrix whose rows are p-codons of 16–5-letter codons where the first codon is the mutation-tolerant codon to be used in place of the naturallyoccurring 3-letter codons are

<sup>&</sup>lt;sup>3</sup>The zero mutation vector leads to a silent mutation but is added to facilitate use of Z in a latter equation.

to be decoded by artificial tRNAs during translation as the 3-letter codon would be in the natural genetic code.

$$K = \varphi_{16} [ (\lambda_1 \lambda_2 \lambda_3) \cdot G ] + Z$$
  
where  $\varphi_n[i] = [\underbrace{i, i, \dots, i}_{n \text{ times}}]^T$  (3.8)

Take codon ATG as an example. In a DNA sense strand, it codes for methionine. Its p-codon is  $(01\alpha)$ . Using equation (3.8),  $K = [(01\alpha\alpha1), (11\alpha\alpha1), (00\alpha\alpha1), (01\alpha^2\alpha1), (01\alpha\alpha^21), (01\alpha\alpha0), (\alpha1\alpha\alpha1), (0\alpha^2\alpha\alpha1), (010\alpha1), (01\alpha01), (01\alpha\alpha\alpha^2), (\alpha^21\alpha\alpha1), (0\alpha\alpha\alpha1), (011\alpha1), (01\alpha11), (01\alpha\alpha\alpha)]^T$ . The first p-codon in K is  $(01\alpha\alpha1)$ . In nucleotide terms, it is ATGGT. Thus, in 'converting' a naturally-occurring gene to its mutol variant, codon ATGGT should be used in place of ATG. Also, we ought to have artificial tRNAs that recognize and decode all 16 5-letter codons corresponding to the p-codons in K (given below) as ATG would naturally be decoded — methionine.

$\operatorname{ATGGT}_{(01\alpha\alpha1)}$	$\operatorname{ATGCT}_{(01\alpha\alpha^2 1)}$	$ATAGT_{(010\alpha 1)}$	$AGGGT_{(0\alpha\alpha\alpha1)}$
$\mathrm{TTGGT}_{(11\alpha\alpha1)}$	$\mathrm{ATGGA}_{(01lpha lpha 0)}$	$\operatorname{ATGAT}_{(01\alpha01)}$	$\text{ATTGT}_{(011\alpha 1)}$
$AAGGT_{(00\alpha\alpha1)}$	$\mathrm{GTGGT}_{(\alpha 1 \alpha \alpha 1)}$	$\mathrm{ATGGC}_{(01lphalpha^2)}$	$ATGTT_{(01\alpha 11)}$
$\text{ATCGT}_{(01\alpha^2\alpha 1)}$	$ACGGT_{(0\alpha^2\alpha\alpha 1)}$	$CTGGT_{(\alpha^{2}1\alpha\alpha1)}$	$\mathrm{ATGGG}_{(01lphalphalpha)}$

Just as we did above for codon ATG, equation (3.8) has been evaluated for all 64 3-letter codons and the results are given in Appendix A in table A.1.

#### 3.4.2 Demonstration of Mutation-Tolerance Capability

Here, we simulate random substitution mutations in both a wild type GFP gene (NCBI Reference Sequence: NC\_011521.1) and its mutol variant to demonstrate the superior mutation-tolerance of the mutol variant. By mutol variant<sup>4</sup>, we mean the DNA sequence obtained by replacing the naturally occurring 3-letter codons with the 5-letter codons we constructed. In Table 3.3 below, the top sequence in each row is the wild type GFP DNA sequence. The middle sequence on each row is the DNA sequence of the mutol variant with each 5-letter codon right below its corresponding 3-letter codon. As previously mentioned, the underlying assumption is that there exists some synthetic tRNAs capable of 'reading' these 5-letter codons. The bottom sequence in each row also indicates the amino acid coded for by the codons. The first and last codon in each row are numbered for ease of reference.

We randomly selected 20 same-numbered codons from the wild type GFP DNA and its mutol variant. For each of these codons, we randomly substituted one of its nucleotides. We then noted the amino acids coded for by the post-mutation codons in both the wild type and the mutol variant. The results for each of 20 mutated codons in both sequences are given in table 3.4. The columns named "Mutation" have a syntax of the form  $X \rightarrow Y$ , which simply means the original codon was X and became Y after the mutation. The same syntax is used in the columns named "Translation", where the amino acid on the left of the arrow is the originally intended amino acid and the amino acid on the right is the amino acid coded for by the mutated codon. From the table, we see that only codon number 25 and codon number 167 in the wild type sequence still coded for the original amino acid post-mutation. In contrast, all 20 codons that were mutated in the mutol variant still coded for the originally intended amino acid post-mutation.

 $<sup>^{4}</sup>$ Mutol variant as in **mu**tation-**tol**erant variant

5'ATG 1AGT AGTATAAA AAAAAGGA GGAAGGAA GAAGTCTT 7 CTTAG5'ATGCTSKGGA GGAA GCAACT CCAA CATT 14 A CTTCA FACT TGGA GTTGTT GA TGTT GAA CGTT GTT GTT CGTA GAA CCCAA CATT 14 A CCTT 15 CTTAG FGTT CTTAGGAT CTTAG CGAT CTTA GAAA TGAT GGTA GGTA CGGT GATAC CGAT CGAT CGAT CCTT 15 CTTAG CTTAG CGTT CTTA CGAT CAA CTTA CGAA CAA CAAAA CTTT CT CAA CGAT CAA CAAAA CTTT CTTA CAA CGAT CAA CAAAA CTTT CTTA CAAAA CGAT CTT CAA CGAA CAAAA CAA CTTT CTTA CAA CGAA CAAAA CAA CGAA CAAAA CAA CTTT CAAAA CGAA CAAAA CTTT CAAAA CGAA CAAAA CCTT 7 CAAA CGAA CAAAA CCTT 4 CCTT 4 CCTT 4 CCTT 4 CACC 43 CCTA CTTC TCTT CAA CAAAA CAA CAA CTTT CAA CCTT 4 CCAA CAAAA CACA 4 CAAAA CAAA CAA CCTT 4 CCAA CAAAA CCTT 4 CCTT 4 CCTT 4 CCAA CCTT 4 CCAA CCTT 4 CCAA CCTT 4 CCAA CCTT 4 CCAA C </th <th></th> <th></th> <th></th> <th></th> <th></th> <th></th> <th></th>							
$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	5' ATG $^1$	AGT	AAA	GGA	GAA	GAA	CTT $^7$
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	5' ATGCT						
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	M						
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	TTC <sup>8</sup>	ACT	GGA	GTT	GTC	CCA	ATT $^{14}$
$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	TTCTA	ACTTC	GGAAG	GTTTT	GTCGG	CCAAC	ATTCA
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	F	Т	G	V	V	Р	Ι
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	CTT <sup>15</sup>	GTT	GAA	TTA	GAT	GGT	GAT $^{21}$
$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	CTTAG	GTTTT	GAAGT	TTAAT	GATAC	GGTGA	GATAC
$ \begin{array}{ccccc} \operatorname{GTTTT} & \operatorname{AATGG} & \operatorname{GGGCT} & \operatorname{CACGC} & \operatorname{AAAAA} & \operatorname{TTTGC} & \operatorname{TCTAA} \\ \operatorname{V} & \operatorname{N} & \operatorname{G} & \operatorname{G} & \operatorname{GAG} & \operatorname{GGT} & \operatorname{GAA} & \operatorname{GGT} & \operatorname{S} & \operatorname{S} \\ \operatorname{GTCGG} & \operatorname{AGTAT} & \operatorname{GGAA} & \operatorname{GAGTG} & \operatorname{GAGTG} & \operatorname{GGTGA} & \operatorname{GAAGT} & \operatorname{GGTGA} \\ \operatorname{V} & \operatorname{S} & \operatorname{G} & \operatorname{CACA} & \operatorname{TAC} & \operatorname{GGA} & \operatorname{AAAA} & \operatorname{CTT} & \operatorname{CTAG} \\ \operatorname{GATAC} & \operatorname{GCATA} & \operatorname{ACA} & \operatorname{TAC} & \operatorname{GGA} & \operatorname{AAAA} & \operatorname{CTT} & \operatorname{CTTAG} \\ \operatorname{D} & \operatorname{A} & \operatorname{T} & \operatorname{Y} & \operatorname{G} & \operatorname{GAAG} & \operatorname{AAAAA} & \operatorname{CTTAG} \\ \operatorname{D} & \operatorname{A} & \operatorname{T} & \operatorname{Y} & \operatorname{G} & \operatorname{GAA} & \operatorname{CAC} & \operatorname{AAAAA} \\ \operatorname{CTTAG} & \operatorname{ACC} & \operatorname{ACA} & \operatorname{TAC} & \operatorname{GGA} & \operatorname{AAAAA} & \operatorname{CTT} & \operatorname{GGT} \\ \operatorname{ACC} & \operatorname{CTTAG} & \operatorname{AAAA} & \operatorname{TTT} & \operatorname{ATT} & \operatorname{TGC} & \operatorname{ACT} & \operatorname{C} & \operatorname{C} \\ \operatorname{T} & \operatorname{L} & \operatorname{K} & \operatorname{F} & \operatorname{I} & \operatorname{C} & \operatorname{C} & \operatorname{C} \\ \operatorname{T} & \operatorname{C} \\ \operatorname{GGAAG} & \operatorname{AAAAA} & \operatorname{CTA} & \operatorname{CCT} & \operatorname{CT} & \operatorname{CCA} & \operatorname{C} & \operatorname{C} & \operatorname{C} & \operatorname{C} \\ \operatorname{ACT} & \operatorname{C} & \operatorname{GGAA} & \operatorname{AAAAA} & \operatorname{CTAGA} & \operatorname{CTGT} & \operatorname{CC} & \operatorname{C} & $	L	V	$\mathbf{E}$	L	D	G	D
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	GTT <sup>22</sup>	AAT	GGG	CAC	AAA	TTT	TCT $^{28}$
$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	GTTTT	AATGG	GGGCT	CACGC	AAAAA	TTTGC	TCTAA
$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	V	Ν	G	Η	Κ	$\mathbf{F}$	$\mathbf{S}$
$\begin{array}{c ccccc} V & S & G & E & G & E & G \\ \hline GAT \ ^{36} & GCA & ACA & TAC & GGA & AAA & CTT \ ^{42} & GATAC & GCATA & ACACT & TACAG & GGAAG & AAAA & CTTAG & D & A & T & Y & G & K & L \\ \hline ACC \ ^{43} & CTT & AAA & TTT & ATT & TGC & ACT \ ^{49} & ACCGA & CTTAG & AAAA & TTTGC & ATTCA & TGCGT & ACTTC & T & L & K & F & I & C & T \\ \hline ACC \ ^{50} & GGA & AAA & CTA & CTA & CCT & GTT & CCA \ ^{56} & ACTTC & GGAAG & AAAAA & CTAGA & CTGT & GTTT & CCAACC & T & G & CCTGT & GTTTT & CCAACC & T & G & CCTAC & GTTTT & CCAACC & CTAC & CTTG & ACT & ACTTC & CCAACC & CTAC & ACACT & CTTAG & GTCGG & ACTTC & ACTTC & CCAACC & CCAAC & ACACT & CTTAG & GTCGG & ACTTC & ACTT & CCAACC & CCAAC & ACACT & CTTAG & GTCGG & ACTTC & ACTTC & CCAACC & CCAACC & ACACT & CTTAG & GTCGG & ACTTC & ACTTC & CCACC & G & CCAACC & CCAACC & CCAACC & CCAAC & CTT & GTC & ACT & ACTTC & CCAACC & CCATTA & ACACG & TTTC & TTC \ TTTGC \ TTT \ TTTC \ ACT \ TTC \ T$	GTC <sup>29</sup>	AGT	GGA	GAG	GGT	GAA	GGT $^{35}$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	GTCGG	AGTAT	GGAAG	GAGTG	GGTGA	GAAGT	GGTGA
$ \begin{array}{c cccc} GATAC & GCATA & ACACT & TACAG & GGAAG & AAAAA & CTTAG \\ D & A & T & Y & G & K & L \\ \hline \\ ACC ^{43} & CTT & AAA & TTT & ATT & TGC & ACT ^{49} \\ ACCGA & CTTAG & AAAAA & TTTGC & ATTCA & TGCGT & ACTTC \\ T & L & K & F & I & C & T \\ \hline \\ ACT ^{50} & GGA & AAA & CTA & CCT & GTT & CCA ^{56} \\ ACTTC & GGAAG & AAAA & CTAA & CTGT & GTTT & CCAAC \\ T & G & K & L & P & V & P \\ \hline \\ TGG ^{57} & CCA & ACA & ACA & CTT & GTC & ACT & ACT ^{63} \\ TGGAC & CCAAC & ACA & CTT & GTC & ACT & ACT C \\ W & P & T & L & V & T & T \\ \hline \\ TTC ^{64} & GGT & TAT & GGT & GTTTT & CAA & TGC ^{70} \\ TTCTA & GGTGA & TATCT & GGTGA & GTTT & CAA & TGC ^{70} \\ TTCTA & GGTGA & TATCT & GGTGA & TTT & CAAC & TGCGT \\ F & G & Y & G & V & Q & C \\ \hline \\ TTT ^{71} & GCG & AGA & TAC & CCA & GAT & CAT ^{77} \\ TTTGC & GCGGC & AGAG & TAC & CAA & CAAC & CAT & TTT \\ TTC ^{78} & AAA & CAG & CAT & GAC & TTT & TTC ^{84} \\ ATGGT & AAAAA & CAGAT & CATA & GACCA & TTTGC & TTCT \\ \end{array}$	V	$\mathbf{S}$	G	$\mathbf{E}$	G	$\mathbf{E}$	G
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	GAT $^{36}$	GCA	ACA	TAC	GGA	AAA	CTT $^{42}$
$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	GATAC	GCATA	ACACT	TACAG	GGAAG	AAAAA	CTTAG
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	D	А	Т	Υ	G	Κ	$\mathbf{L}$
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	ACC $^{43}$	CTT	AAA	TTT	ATT	TGC	ACT $^{49}$
$\begin{array}{cccccccccccccccccccccccccccccccccccc$		CTTAG	AAAAA	TTTGC			
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	T	L	К	F	Ι	С	Т
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	ACT $^{50}$	GGA	AAA	CTA	$\mathbf{CCT}$	GTT	CCA $^{56}$
$\begin{array}{cccccccccccccccccccccccccccccccccccc$		GGAAG	AAAAA	CTAGA	CCTGT	GTTTT	CCAAC
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	Т	G	К	L	Р	V	Р
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	TGG $^{57}$	CCA	ACA	CTT	GTC	ACT	ACT $^{63}$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		CCAAC	ACACT	CTTAG	GTCGG		
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	W	Р	Т	L	V	Т	Т
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	TTC <sup>64</sup>	GGT	TAT	GGT	GTT	CAA	TGC $^{70}$
$\begin{array}{c cccccc} TTT & TT & GCG & AGA & TAC & CCA & GAT & CAT ^{77} \\ TTTGC & GCGGC & AGAGC & TACAG & CCAAC & GATAC & CATTA \\ F & A & R & Y & P & D & H \end{array}$							
$\begin{array}{cccc} TTTGC \\ F \end{array} & \begin{array}{c} GCGGC \\ A \end{array} & \begin{array}{c} AGAGC \\ R \end{array} & \begin{array}{c} TACAG \\ Y \end{array} & \begin{array}{c} CCAAC \\ P \end{array} & \begin{array}{c} GATAC \\ D \end{array} & \begin{array}{c} CATTA \\ H \end{array} \\ \end{array} \\ \begin{array}{c} ATG^{78} \\ ATGGT \end{array} & \begin{array}{c} AAA \\ AAAAA \end{array} & \begin{array}{c} CAG \\ CAGAT \end{array} & \begin{array}{c} CAT \\ CATTA \end{array} & \begin{array}{c} GAC \\ GACC \end{array} & \begin{array}{c} TTT \\ TTC \end{array} & \begin{array}{c} TTC ^{84} \\ TTTGC \end{array} \\ \end{array}$	F	G	Y	G	V	Q	С
FARYPDHATG 78AAACAGCATGACTTTTTC 84ATGGTAAAAACAGATCATTAGACCATTTGCTTCTA							
ATG <sup>78</sup> AAA CAG CAT GAC TTT TTC <sup>84</sup> ATGGT AAAAA CAGAT CATTA GACCA TTTGC TTCTA							
ATGGT AAAAA CAGAT CATTA GACCA TTTGC TTCTA	F	А	R	Y	Р	D	Н
M K Q H D F F							
	M	K	Q	H	D	F	F

 Table 3.3: GFP DNA Sequence and Its Mutol Variant

$ \begin{array}{cccccccccccccccccccccccccccccccccccc$							
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	AAG <sup>85</sup>	AGT	GCC	ATG	CCT	GAA	$GGT$ $^{91}$
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$							
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$					Р		
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	TAT <sup>92</sup>	GTA	CAG	GAA	AGA	ACT	ATA 98
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	TATCT						
$\begin{array}{cccccccccccccccccccccccccccccccccccc$			$\mathbf{Q}$			Т	
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	TTT <sup>99</sup>	TTC	AAA	GAT	GAC	GGG	AAC $^{105}$
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	TTTGC	TTCTA	AAAAA	GATAC	GACCA	GGGCT	AACTT
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	F	$\mathbf{F}$	Κ	D	D	G	Ν
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	TAC <sup>106</sup>	AAG	ACA	CGT	GCT	GAA	GTC $^{112}$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	TACAG	AAGCC	ACACT	CGTCC	GCTCG	GAAGT	GTCGG
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	Y	Κ	Т	R	А	$\mathbf{E}$	V
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	AAG <sup>113</sup>	TTT	GAA	GGT	GAT	ACC	CTT $^{119}$
$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	AAGCC	TTTGC	GAAGT	GGTGA	GATAC	ACCGA	CTTAG
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	K	F	Ε	G	D	Т	L
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	-	AAT	AGA	ATC	GAG	TTA	AAA $^{126}$
$ \begin{array}{cccccccccccccccccccccccccccccccccccc$		AATGG	AGAGC	ATCAC	GAGTG	TTAAT	AAAAA
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	V	Ν	R	Ι	Ε	L	K
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	GGT $^{127}$	ATT	GAT		AAA	GAA	GAT $^{133}$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$			GATAC				GATAC
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	G	Ι	D	F	Κ	Ε	D
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	GGA $^{134}$						AAA $^{140}$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$							
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	G	Ν	Ι	L	G	Н	K
$\begin{array}{c ccccc} L & E & Y & N & Y & N & S \\ \hline CAC ^{148} & AAT & GTA & TAC & ATC & ATG & GCA ^{154} \\ CACGC & AATGG & GTACC & TACAG & ATCAC & ATGGT & GCATA \\ H & N & V & Y & I & M & A \\ \hline GAC ^{155} & AAA & CAA & AAG & AAT & GGA & ATC ^{161} \\ GACCA & AAAAA & CAACG & AAGCC & AATGG & GGAAG & ATCAC \\ D & K & Q & K & N & G & I \\ \hline AAA ^{162} & GTT & AAC & TTC & AAA & ATT & AGA ^{168} \\ AAAAA & GTTTT & AACTT & TTCTA & AAAAA & ATTCA & AGAGC \\ \hline \end{array}$	TTG $^{141}$	GAA	TAC	AAC	TAT	AAC	TCA $^{147}$
$ \begin{array}{c ccccccc} CAC & ^{148} & AAT & GTA & TAC & ATC & ATG & GCA & ^{154} \\ CACGC & AATGG & GTACC & TACAG & ATCAC & ATGGT & GCATA \\ H & N & V & Y & I & M & A \\ \hline GAC & ^{155} & AAA & CAA & AAG & AAT & GGA & ATC & ^{161} \\ GACCA & AAAAA & CAACG & AAGCC & AATGG & GGAAG & ATCAC \\ D & K & Q & K & N & G & I \\ \hline AAA & ^{162} & GTT & AAC & TTC & AAA & ATT & AGA & ^{168} \\ AAAAA & GTTTT & AACTT & TTCTA & AAAAA & ATTCA & AGAGC \\ \hline \end{array} $	TTGCG	GAAGT	TACAG	AACTT	TATCT	AACTT	TCAGG
$ \begin{array}{c cccccc} CACGC & AATGG & GTACC & TACAG & ATCAC & ATGGT & GCATA \\ H & N & V & Y & I & M & A \\ \hline GAC ^{155} & AAA & CAA & AAG & AAT & GGA & ATC ^{161} \\ GACCA & AAAAA & CAACG & AAGCC & AATGG & GGAAG & ATCAC \\ D & K & Q & K & N & G & I \\ \hline AAA ^{162} & GTT & AAC & TTC & AAA & ATT & AGA ^{168} \\ AAAAA & GTTTT & AACTT & TTCTA & AAAAA & ATTCA & AGAGC \\ \hline \end{array} $	L	Ε	Υ	Ν	Υ	Ν	S
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	CAC $^{148}$	AAT	GTA	TAC	ATC	ATG	GCA $^{154}$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	CACGC	AATGG	GTACC	TACAG	ATCAC	ATGGT	GCATA
$ \begin{array}{c ccccc} GACCA & AAAAA & CAACG & AAGCC & AATGG & GGAAG & ATCAC \\ D & K & Q & K & N & G & I \\ \hline AAA^{162} & GTT & AAC & TTC & AAA & ATT & AGA^{168} \\ AAAAA & GTTTT & AACTT & TTCTA & AAAAA & ATTCA & AGAGC \\ \hline \end{array} $	Н	Ν	V	Υ	Ι	Μ	А
DKQKNGIAAAGTTAACTTCAAAATTAGAAGAAAAAAGTTTTAACTTTTCTAAAAAAATTCAAGAGC	GAC $^{155}$	AAA	CAA	AAG	AAT		ATC $^{161}$
AAA <sup>162</sup> GTT AAC TTC AAA ATT AGA <sup>168</sup> AAAAA GTTTT AACTT TTCTA AAAAA ATTCA AGAGC	GACCA		CAACG				
AAAAA GTTTT AACTT TTCTA AAAAA ATTCA AGAGC	D	Κ	Q	Κ	Ν	G	Ι
	AAA <sup>162</sup>	GTT	AAC	TTC	AAA	ATT	AGA $^{168}$
K V N F K I R							
	К	V	Ν	F	Κ	Ι	R

 Table 3.3: GFP DNA Sequence and Its Mutol Variant

$\begin{array}{c} {\rm CAC} \ ^{169} \\ {\rm CACGC} \\ {\rm H} \end{array}$	AAC AACTT N	ATT ATTCA I	GAA GAAGT E	GAT GATAC D	GGA GGAAG G	$\begin{array}{c} \mathrm{AGC} \ ^{175} \\ \mathrm{AGCCG} \\ \mathrm{S} \end{array}$
$\begin{array}{c} {\rm GTT} \ ^{176} \\ {\rm GTTTT} \\ V \end{array}$	CAA CAACG Q	CTA CTAGA L	GCA GCATA A	GAC GACCA D	CAT CATTA H	$\begin{array}{c} {\rm TAT} \ ^{182} \\ {\rm TATCT} \\ {\rm Y} \end{array}$
$\begin{array}{c} \mathrm{CAA} \ ^{183} \\ \mathrm{CAACG} \\ \mathrm{Q} \end{array}$	CAA CAACG Q	AAT AATGG N	ACT ACTTC T	CCA CCAAC P	ATT ATTCA I	$\begin{array}{c} \mathrm{GGC} \ ^{189} \\ \mathrm{GGCTC} \\ \mathrm{G} \end{array}$
$\begin{array}{c} {\rm GAT} \ ^{190} \\ {\rm GATAC} \\ {\rm D} \end{array}$	GGC GGCTC G	CCT CCTGT P	GTC GTCGG V	CTT CTTAG L	TTA TTAAT L	$\begin{array}{c} {\rm CCA} \ ^{196} \\ {\rm CCAAC} \\ {\rm P} \end{array}$
GAC <sup>197</sup>	AAC	CAT	TAC	CTG	TCC	ACA <sup>203</sup>
GACCA	AACTT	CATTA	TACAG	CTGTC	TCCCC	ACACT
D	N	H	Y	L	S	T
$\begin{array}{c} \mathrm{CAA} \ ^{204} \\ \mathrm{CAACG} \\ \mathrm{Q} \end{array}$	TCT	GCC	CTT	TCG	AAA	GAT <sup>210</sup>
	TCTAA	GCCAT	CTTAG	TCGTT	AAAAA	GATAC
	S	A	L	S	K	D
CCC <sup>211</sup>	AAC	GAA	AAG	AGA	GAC	$\begin{array}{c} {\rm CAC} \ ^{217} \\ {\rm CACGC} \\ {\rm H} \end{array}$
CCCTG	AACTT	GAAGT	AAGCC	AGAGC	GACCA	
P	N	E	K	R	D	
ATG <sup>218</sup>	GTC	CTT	CTT	GAG	TTT	$\begin{array}{c} {\rm GTA} \ ^{224} \\ {\rm GTACC} \\ V \end{array}$
ATGGT	GTCGG	CTTAG	CTTAG	GAGTG	TTTGC	
M	V	L	L	E	F	
$\begin{array}{c} \mathrm{ACA} \ ^{225} \\ \mathrm{ACACT} \\ \mathrm{T} \end{array}$	GCT	GCT	GGG	ATT	ACA	CAT <sup>231</sup>
	GCTCG	GCTCG	GGGCT	ATTCA	ACACT	CATTA
	A	A	G	I	T	H
$\begin{array}{c} \mathrm{GGC} \ ^{232} \\ \mathrm{GGCTC} \\ \mathrm{G} \end{array}$	ATG	GAT	GAA	CTA	TAC	AAA <sup>238</sup>
	ATGGT	GATAC	GAAGT	CTAGA	TACAG	AAAAA
	M	D	E	L	Y	K
TAA <sup>239</sup> 3' TAATC 3' Stop						

 Table 3.3: GFP DNA Sequence and Its Mutol Variant

	Wildtype	e GFP	Mutol Variant		
Codon No.	Mutation	Translation	Mutation	Translation	
10	$GGA \rightarrow GCA$	$\mathbf{G} \to \mathbf{A}$	$GGAAG \rightarrow GCAAG$	$\mathbf{G} \to \mathbf{G}$	
13	$CCA \rightarrow GCA$	$\mathbf{P} \to \mathbf{A}$	$CCAAC \rightarrow CCAAG$	$\mathbf{P} \to \mathbf{P}$	
15	$CTT \rightarrow GTT$	$\mathrm{L} \to \mathrm{V}$	$\mathrm{CTTAG} \to \mathrm{CCTAG}$	$\mathrm{L} \to \mathrm{L}$	
18	$TTA \rightarrow GTA$	$\mathrm{L} \to \mathrm{V}$	$\mathrm{TTAAT} \rightarrow \mathrm{TTAGT}$	$\mathrm{L} \to \mathrm{L}$	
19	$GAT \rightarrow CAT$	$\mathrm{D} \to \mathrm{H}$	$GATAC \rightarrow AATAC$	$\mathrm{D}\to\mathrm{D}$	
23	$AAT \rightarrow AAA$	$\mathbf{N} \to \mathbf{K}$	$\mathbf{AATGG} \rightarrow \mathbf{AATGT}$	$\mathrm{N} \rightarrow \mathrm{N}$	
25	$CAC \rightarrow CAT$	$\mathrm{H} \to \mathrm{H}$	$CACGC \rightarrow CACAC$	$\mathrm{H} \to \mathrm{H}$	
39	$TAC \rightarrow TTC$	$\mathbf{Y} \to \mathbf{F}$	$\mathrm{TACAG} \to \mathrm{TCCAG}$	$\mathbf{Y} \to \mathbf{Y}$	
55	$GTT \rightarrow ATT$	$\mathbf{V} \to \mathbf{I}$	$GTTTT \rightarrow GCTTT$	$V \rightarrow V$	
83	$TTT \rightarrow TAT$	$\mathbf{F} \to \mathbf{Y}$	$\mathrm{TTTGC} \to \mathrm{TTGGC}$	$F \rightarrow F$	
94	$CAG \rightarrow CAC$	$\mathbf{Q} \to \mathbf{H}$	$\mathrm{CAGAT} \to \mathrm{CAGCT}$	$\mathbf{Q} \to \mathbf{Q}$	
102	$GAT \rightarrow GTT$	$\mathrm{D} \to \mathrm{V}$	$GATAC \rightarrow GCTAC$	$\mathrm{D}\to\mathrm{D}$	
108	$ACA \rightarrow ATA$	$\mathrm{T} \to \mathrm{I}$	$\mathbf{ACACT} \to \mathbf{ACAGT}$	$\mathrm{T} \to \mathrm{T}$	
112	$GTC \rightarrow GCC$	$\mathbf{V} \to \mathbf{A}$	$\mathrm{GTCGG} \to \mathrm{TTCGG}$	$\mathbf{V} \to \mathbf{V}$	
130	$TTT \rightarrow ATT$	$\mathbf{F} \to \mathbf{I}$	$\mathrm{TTTGC} \to \mathrm{TATGC}$	$\mathbf{F} \to \mathbf{F}$	
138	$GGA \rightarrow TGA$	$\mathbf{G} \to \mathbf{Stop}$	$GGAAG \rightarrow GGGAG$	$\mathbf{G} \to \mathbf{G}$	
150	$GTA \rightarrow GGA$	$\mathbf{V} \to \mathbf{G}$	$GTACC \rightarrow GTACA$	$\mathbf{V} \to \mathbf{V}$	
167	$ATT \rightarrow ATA$	$\mathrm{I} \to \mathrm{I}$	$\mathrm{ATTCA} \to \mathrm{AGTCA}$	$\mathrm{I} \to \mathrm{I}$	
211	$CCC \rightarrow ACC$	$\mathbf{P} \to \mathbf{T}$	$\mathrm{CCCTG} \to \mathrm{CTCTG}$	$\mathbf{P} \to \mathbf{P}$	
233	$ATG \rightarrow ATA$	$\mathbf{M} \to \mathbf{I}$	$\mathrm{ATGGT} \to \mathrm{ATTGT}$	$M \rightarrow M$	

 Table 3.4:
 Mutation Simulation Results

### REFERENCES

Axler, S., Linear Algebra Done Right (Springer, 2007).

- Hamming, R. W., "Error detecting and error correcting codes", Bell System Technical Journal 29, 2, 147160 (1950).
- Lin, S. and D. J. Costello, Error control coding: fundamentals and applications (Prentice-Hall, 1983).
- Ling, S., *Coding theory : a first course* (Cambridge University Press, Cambridge, UK New York, 2004).

# APPENDIX A

3-LETTER CODONS TO 5-LETTER CODONS

	$AAAAA_{(00000)}$	$AAATA_{(00010)}$	$AAGAA_{(00a00)}$	$ACAAA_{(0a^2000)}$
AAA	$TAAAA_{(10000)}$	$AAAAT_{(00001)}$	$AAAGA_{(000a0)}$	$AACAA_{(00a^200)}$
(000)	$ATAAA_{(01000)}$	$\mathrm{GAAAA}_{(a0000)}$	$AAAAG_{(0000a)}$	$AAACA_{(000a^20)}$
	$AATAA_{(00100)}$	$AGAAA_{(0a000)}$	$CAAAA_{(a^20000)}$	$AAAAC_{(0000a^2)}$
	$AATGG_{(001aa)}$	$AATCG_{(001a^2a)}$	$AACGG_{(00a^2aa)}$	$ACTGG_{(0a^21aa)}$
AAT	$TATGG_{(101aa)}$	$AATGC_{(001aa^2)}$	$AATAG_{(0010a)}$	$AAGGG_{(00aaa)}$
(001)	$ATTGG_{(011aa)}$	$GATGG_{(a01aa)}$	$AATGA_{(001a0)}$	$AATTG_{(0011a)}$
	$AAAGG_{(000aa)}$	$AGTGG_{(0a1aa)}$	$CATGG_{(a^201aa)}$	$AATGT_{(001a1)}$
	$AAGCC_{(00aa^2a^2)}$	$AAGGC_{(00aaa^2)}$	$AAACC_{(000a^2a^2)}$	$ACGCC_{(0a^2aa^2a^2)}$
AAG	$\mathrm{TAGCC}_{(10aa^2a^2)}$	$AAGCG_{(00aa^2a)}$	$AAGTC_{(00a1a^2)}$	$AATCC_{(001a^2a^2)}$
(00a)	$\operatorname{ATGCC}_{(01aa^2a^2)}$	$GAGCC_{(a0aa^2a^2)}$	$AAGCT_{(00aa^21)}$	$AAGAC_{(00a0a^2)}$
	$AACCC_{(00a^2a^2a^2)}$	$AGGCC_{(0aaa^2a^2)}$	$CAGCC_{(a^20aa^2a^2)}$	$AAGCA_{(00aa^20)}$
	$AACTT_{(00a^211)}$	$AACAT_{(00a^201)}$	$AATTT_{(00111)}$	$ACCTT_{(0a^2a^211)}$
AAC	$TACTT_{(10a^211)}$	$AACTA_{(00a^210)}$	$AACCT_{(00a^2a^21)}$	$AAATT_{(00011)}$
$(00a^2)$	$ATCTT_{(01a^211)}$	$GACTT_{(a0a^211)}$	$AACTC_{(00a^21a^2)}$	$AACGT_{(00a^2a1)}$
	$AAGTT_{(00a11)}$	$\mathrm{AGCTT}_{(0aa^211)}$	$CACTT_{(a^20a^211)}$	$AACTG_{(00a^21a)}$
	$ATATG_{(0101a)}$	$ATAAG_{(0100a)}$	$\operatorname{ATGTG}_{(01a1a)}$	$AGATG_{(0a01a)}$
ATA	$\mathrm{TTATG}_{(1101a)}$	$\operatorname{ATATC}_{(0101a^2)}$	$\text{ATACG}_{(010a^2a)}$	$\text{ATCTG}_{(01a^21a)}$
(010)	$AAATG_{(0001a)}$	$\operatorname{GTATG}_{(a101a)}$	$ATATA_{(01010)}$	$ATAGG_{(010aa)}$
	$\text{ATTTG}_{(0111a)}$	$ACATG_{(0a^201a)}$	$CTATG_{(a^2101a)}$	$\mathrm{ATATT}_{(01011)}$
	$\text{ATTCA}_{(011a^20)}$	$ATTGA_{(011a0)}$	$\text{ATCCA}_{(01a^2a^20)}$	$AGTCA_{(0a1a^20)}$
ATT	$\text{TTTCA}_{(111a^20)}$	$\text{ATTCT}_{(011a^21)}$	$\text{ATTTA}_{(01110)}$	$\text{ATGCA}_{(01aa^20)}$
(011)	$AATCA_{(001a^20)}$	$\text{GTTCA}_{(a11a^20)}$	$\text{ATTCG}_{(011a^2a)}$	$ATTAA_{(01100)}$
	$ATACA_{(010a^20)}$	$ACTCA_{(0a^21a^20)}$	$CTTCA_{(a^2 11a^2 0)}$	$\operatorname{ATTCC}_{(011a^2a^2)}$
	ATGGT <sub>(01aa1)</sub>	$\text{ATGCT}_{(01aa^21)}$	ATAGT <sub>(010a1)</sub>	$AGGGT_{(0aaa1)}$
ATG	$\mathrm{TTGGT}_{(11aa1)}$	$ATGGA_{(01aa0)}$	$\operatorname{ATGAT}_{(01a01)}$	$\operatorname{ATTGT}_{(011a1)}$
(01a)	$AAGGT_{(00aa1)}$	$GTGGT_{(a1aa1)}$	$\operatorname{ATGGC}_{(01aaa^2)}$	$\operatorname{ATGTT}_{(01a11)}$
	$\text{ATCGT}_{(01a^2a1)}$	$ACGGT_{(0a^2aa1)}$	$CTGGT_{(a^21aa1)}$	$ATGGG_{(01aaa)}$
			× /	

 Table A.1: Degeneracy Mapping

$\begin{array}{c} \text{ATC} \\ (01a^2) \end{array}$	$\begin{array}{l} \operatorname{ATCAC}_{(01a^20a^2)} \\ \operatorname{TTCAC}_{(11a^20a^2)} \\ \operatorname{AACAC}_{(00a^20a^2)} \\ \operatorname{ATGAC}_{(01a0a^2)} \end{array}$	$\begin{array}{l} \text{ATCTC}_{(01a^21a^2)} \\ \text{ATCAG}_{(01a^20a)} \\ \text{GTCAC}_{(a1a^20a^2)} \\ \text{ACCAC}_{(0a^2a^20a^2)} \end{array}$	$\begin{array}{l} \operatorname{ATTAC}_{(0110a^2)} \\ \operatorname{ATCGC}_{(01a^2aa^2)} \\ \operatorname{ATCAT}_{(01a^201)} \\ \operatorname{CTCAC}_{(a^21a^20a^2)} \end{array}$	$\begin{array}{l} \mathrm{AGCAC}_{(0aa^20a^2)} \\ \mathrm{ATAAC}_{(0100a^2)} \\ \mathrm{ATCCC}_{(01a^2a^2a^2)} \\ \mathrm{ATCCAA}_{(01a^200)} \end{array}$
$\begin{array}{c} AGA\\ (0a0) \end{array}$	$\begin{array}{l} \operatorname{AGAGC}_{(0a0aa^2)} \\ \operatorname{TGAGC}_{(1a0aa^2)} \\ \operatorname{ACAGC}_{(0a^20aa^2)} \\ \operatorname{AGTGC}_{(0a1aa^2)} \end{array}$	$\begin{array}{l} \mathrm{AGACC}_{(0a0a^2a^2)}\\ \mathrm{AGAGG}_{(0a0aa)}\\ \mathrm{GGAGC}_{(aa0aa^2)}\\ \mathrm{AAAGC}_{(000aa^2)} \end{array}$	$\begin{array}{l} \operatorname{AGGGGC}_{(0aaaa^2)} \\ \operatorname{AGAAC}_{(0a00a^2)} \\ \operatorname{AGAGT}_{(0a0a1)} \\ \operatorname{CGAGC}_{(a^2a0aa^2)} \end{array}$	$\begin{array}{l} \mathrm{ATAGC}_{(010aa^2)} \\ \mathrm{AGCGC}_{(0aa^2aa^2)} \\ \mathrm{AGATC}_{(0a01a^2)} \\ \mathrm{AGAGA}_{(0a0a0)} \end{array}$
$\begin{array}{c} \text{AGT} \\ (0a1) \end{array}$	$\begin{array}{l} \operatorname{AGTAT}_{(0a101)} \\ \operatorname{TGTAT}_{(1a101)} \\ \operatorname{ACTAT}_{(0a^2101)} \\ \operatorname{AGAAT}_{(0a001)} \end{array}$	$\begin{array}{l} \mathrm{AGTTT}_{(0a111)} \\ \mathrm{AGTAA}_{(0a100)} \\ \mathrm{GGTAT}_{(aa101)} \\ \mathrm{AATAT}_{(00101)} \end{array}$	$\begin{array}{l} \mathrm{AGCAT}_{(0aa^201)} \\ \mathrm{AGTGT}_{(0a1a1)} \\ \mathrm{AGTAC}_{(0a10a^2)} \\ \mathrm{CGTAT}_{(a^2a101)} \end{array}$	$\begin{array}{l} \operatorname{ATTAT}_{(01101)} \\ \operatorname{AGGAT}_{(0aa01)} \\ \operatorname{AGTCT}_{(0a1a^21)} \\ \operatorname{AGTAG}_{(0a10a)} \end{array}$
AGG (0aa)	$\begin{array}{l} \operatorname{AGGTA}_{(0aa10)} \\ \operatorname{TGGTA}_{(1aa10)} \\ \operatorname{ACGTA}_{(0a^2a10)} \\ \operatorname{AGCTA}_{(0aa^210)} \end{array}$	$\begin{array}{l} \operatorname{AGGAA}_{(0aa00)} \\ \operatorname{AGGTT}_{(0aa11)} \\ \operatorname{GGGTA}_{(aaa10)} \\ \operatorname{AAGTA}_{(00a10)} \end{array}$	$\begin{array}{l} \text{AGATA}_{(0a010)} \\ \text{AGGCA}_{(0aaa^20)} \\ \text{AGGTG}_{(0aa1a)} \\ \text{CGGTA}_{(a^2aa10)} \end{array}$	$\begin{array}{l} \text{ATGTA}_{(01a10)} \\ \text{AGTTA}_{(0a110)} \\ \text{AGGGA}_{(0aaa0)} \\ \text{AGGTC}_{(0aa1a^2)} \end{array}$
$\begin{array}{c} \text{AGC} \\ (0aa^2) \end{array}$	$\begin{array}{l} \operatorname{AGCCG}_{(0aa^2a^2a)} \\ \operatorname{TGCCG}_{(1aa^2a^2a)} \\ \operatorname{ACCCG}_{(0a^2a^2a^2a)} \\ \operatorname{AGGCG}_{(0aaa^2a)} \end{array}$	$\begin{array}{l} \operatorname{AGCGG}_{(0aa^2aa)} \\ \operatorname{AGCCC}_{(0aa^2a^2a^2)} \\ \operatorname{GGCCG}_{(aaa^2a^2a)} \\ \operatorname{AACCG}_{(00a^2a^2a)} \end{array}$	$\begin{array}{l} \operatorname{AGTCG}_{(0a1a^2a)} \\ \operatorname{AGCTG}_{(0aa^21a)} \\ \operatorname{AGCCA}_{(0aa^2a^20)} \\ \operatorname{CGCCG}_{(a^2aa^2a^2a)} \end{array}$	$\begin{array}{l} \operatorname{ATCCG}_{(01a^2a^2a)} \\ \operatorname{AGACG}_{(0a0a^2a)} \\ \operatorname{AGCAG}_{(0aa^20a)} \\ \operatorname{AGCCT}_{(0aa^2a^21)} \end{array}$
$\begin{array}{c} ACA\\ (0a^20) \end{array}$	$\begin{array}{c} \operatorname{ACACT}_{(0a^20a^21)} \\ \operatorname{TCACT}_{(1a^20a^21)} \\ \operatorname{AGACT}_{(0a0a^21)} \\ \operatorname{ACTCT}_{(0a^21a^21)} \end{array}$	$\begin{array}{l} \text{ACAGT}_{(0a^20a1)} \\ \text{ACACA}_{(0a^20a^20)} \\ \text{GCACT}_{(aa^20a^21)} \\ \text{ATACT}_{(010a^21)} \end{array}$	$\begin{array}{l} \operatorname{ACGCT}_{(0a^2aa^21)} \\ \operatorname{ACATT}_{(0a^2011)} \\ \operatorname{ACACC}_{(0a^20a^2a^2)} \\ \operatorname{CCACT}_{(a^2a^20a^21)} \end{array}$	$\begin{array}{l} \text{AAACT}_{(000a^21)} \\ \text{ACCCT}_{(0a^2a^2a^21)} \\ \text{ACAAT}_{(0a^2001)} \\ \text{ACAACG}_{(0a^20a^2a)} \end{array}$
$\begin{array}{c} \text{ACT} \\ (0a^21) \end{array}$	$\begin{array}{l} \operatorname{ACTTC}_{(0a^211a^2)} \\ \operatorname{TCTTC}_{(1a^211a^2)} \\ \operatorname{AGTTC}_{(0a11a^2)} \\ \operatorname{ACATC}_{(0a^201a^2)} \end{array}$	$\begin{array}{l} \operatorname{ACTAC}_{(0a^210a^2)} \\ \operatorname{ACTTG}_{(0a^211a)} \\ \operatorname{GCTTC}_{(aa^211a^2)} \\ \operatorname{ATTTC}_{(0111a^2)} \end{array}$	$\begin{array}{l} \operatorname{ACCTC}_{(0a^2a^21a^2)} \\ \operatorname{ACTCC}_{(0a^21a^2a^2)} \\ \operatorname{ACTTT}_{(0a^2111)} \\ \operatorname{CCTTC}_{(a^2a^211a^2)} \end{array}$	$\begin{array}{l} \text{AATTC}_{(0011a^2)} \\ \text{ACGTC}_{(0a^2a1a^2)} \\ \text{ACTGC}_{(0a^21aa^2)} \\ \text{ACTTA}_{(0a^2110)} \end{array}$

 Table A.1: Degeneracy Mapping

	$ACGAG_{(0a^2a0a)}$	$ACGTG_{(0a^2a1a)}$	$ACAAG_{(0a^200a)}$	$AAGAG_{(00a0a)}$
ACG	$\mathrm{TCGAG}_{(1a^2a0a)}$	$ACGAC_{(0a^2a0a^2)}$	$ACGGG_{(0a^2aaa)}$	$ACTAG_{(0a^210a)}$
$(0a^2a)$	$\mathrm{AGGAG}_{(0aa0a)}$	$\operatorname{GCGAG}_{(aa^2a0a)}$	$ACGAA_{(0a^2a00)}$	$ACGCG_{(0a^2aa^2a)}$
	$ACCAG_{(0a^2a^20a)}$	$\operatorname{ATGAG}_{(01a0a)}$	$\operatorname{CCGAG}_{(a^2a^2a0a)}$	$ACGAT_{(0a^2a01)}$
	$ACCGA_{(0a^2a^2a0)}$	$ACCCA_{(0a^2a^2a^20)}$	$ACTGA_{(0a^21a0)}$	$AACGA_{(00a^2a0)}$
ACC	$\mathrm{TCCGA}_{(1a^2a^2a0)}$	$ACCGT_{(0a^2a^2a1)}$	$ACCAA_{(0a^2a^200)}$	$ACAGA_{(0a^20a0)}$
$(0a^2a^2)$	$AGCGA_{(0aa^2a0)}$	$\mathrm{GCCGA}_{(aa^2a^2a0)}$	$ACCGG_{(0a^2a^2aa)}$	$\text{ACCTA}_{(0a^2a^210)}$
	$ACGGA_{(0a^2aa0)}$	$\text{ATCGA}_{(01a^2a0)}$	$CCCGA_{(a^2a^2a^2a0)}$	$ACCGC_{(0a^2a^2aa^2)}$
	$\mathrm{TAATC}_{(1001a^2)}$	$\mathrm{TAAAC}_{(1000a^2)}$	$\mathrm{TAGTC}_{(10a1a^2)}$	$\mathrm{TCATC}_{(1a^201a^2)}$
TAA	$AAATC_{(0001a^2)}$	$TAATG_{(1001a)}$	$\mathrm{TAACC}_{(100a^2a^2)}$	$\mathrm{TACTC}_{(10a^21a^2)}$
(100)	$\mathrm{TTATC}_{(1101a^2)}$	$CAATC_{(a^2001a^2)}$	$\mathrm{TAATT}_{(10011)}$	$\mathrm{TAAGC}_{(100aa^2)}$
	$\mathrm{TATTC}_{(1011a^2)}$	$\mathrm{TGATC}_{(1a01a^2)}$	$\mathrm{GAATC}_{(a001a^2)}$	$\mathrm{TAATA}_{(10010)}$
	$TATCT_{(101a^21)}$	$TATGT_{(101a1)}$	$\text{TACCT}_{(10a^2a^21)}$	$\mathrm{TCTCT}_{(1a^21a^21)}$
TAT	$AATCT_{(001a^21)}$	$TATCA_{(101a^20)}$	$TATTT_{(10111)}$	$\text{TAGCT}_{(10aa^21)}$
(101)	$\mathrm{TTTCT}_{(111a^21)}$	$CATCT_{(a^201a^21)}$	$\mathrm{TATCC}_{(101a^2a^2)}$	$TATAT_{(10101)}$
	$\mathrm{TAACT}_{(100a^21)}$	$\mathrm{TGTCT}_{(1a1a^21)}$	$GATCT_{(a01a^21)}$	$\mathrm{TATCG}_{(101a^2a)}$
	$TAGGA_{(10aa0)}$	$TAGCA_{(10aa^20)}$	$TAAGA_{(100a0)}$	$TCGGA_{(1a^2aa0)}$
TAG	$AAGGA_{(00aa0)}$	$TAGGT_{(10aa1)}$	$TAGAA_{(10a00)}$	$TATGA_{(101a0)}$
(10a)	$TTGGA_{(11aa0)}$	$CAGGA_{(a^20aa0)}$	$TAGGG_{(10aaa)}$	$TAGTA_{(10a10)}$
	$TACGA_{(10a^2a0)}$	$\mathrm{TGGGA}_{(1aaa0)}$	$GAGGA_{(a0aa0)}$	$\mathrm{TAGGC}_{(10aaa^2)}$
	$TACAG_{(10a^20a)}$	$TACTG_{(10a^21a)}$	$TATAG_{(1010a)}$	$\mathrm{TCCAG}_{(1a^2a^20a)}$
TAC	$AACAG_{(00a^20a)}$	$\mathrm{TACAC}_{(10a^20a^2)}$	$TACGG_{(10a^2aa)}$	$TAAAG_{(1000a)}$
$(10a^2)$	$\mathrm{TTCAG}_{(11a^20a)}$	$CACAG_{(a^20a^20a)}$	$TACAA_{(10a^200)}$	$\mathrm{TACCG}_{(10a^2a^2a)}$
	$\mathrm{TAGAG}_{(10a0a)}$	$\mathrm{TGCAG}_{(1aa^20a)}$	$GACAG_{(a0a^20a)}$	$\mathrm{TACAT}_{(10a^201)}$
	$TTAAT_{(11001)}$	$TTATT_{(11011)}$	$\mathrm{TTGAT}_{(11a01)}$	$\mathrm{TGAAT}_{(1a001)}$
TTA	$ATAAT_{(01001)}$	$TTAAA_{(11000)}$	$\mathrm{TTAGT}_{(110a1)}$	$\text{TTCAT}_{(11a^201)}$
(110)	$TAAAT_{(10001)}$	$\operatorname{CTAAT}_{(a^2 1001)}$	$\mathrm{TTAAC}_{(1100a^2)}$	$\mathrm{TTACT}_{(110a^21)}$
	$\mathrm{TTTAT}_{(11101)}$	$\mathrm{TCAAT}_{(1a^2001)}$	$\operatorname{GTAAT}_{(a1001)}$	$TTAAG_{(1100a)}$

 Table A.1: Degeneracy Mapping

TTT (111)	$TTTGC_{(111aa^2)}$ $ATTGC_{(011aa^2)}$ $TATGC_{(101aa^2)}$ $TTAGC_{(110aa^2)}$	$TTTCC_{(111a^2a^2)}$ $TTTGG_{(111aa)}$ $CTTGC_{(a^211aa^2)}$ $TCTGC_{(1a^21aa^2)}$	$TTCGC_{(11a^2aa^2)}$ $TTTAC_{(1110a^2)}$ $TTTGT_{(111a1)}$ $GTTGC_{(a11aa^2)}$	$TGTGC_{(1a1aa^2)}$ $TTGGC_{(11aaa^2)}$ $TTTTC_{(1111a^2)}$ $TTTGA_{(111a0)}$
$\begin{array}{c} {\rm TTG} \\ (11a) \end{array}$	$TTGCG_{(11aa^2a)}$ $ATGCG_{(01aa^2a)}$ $TAGCG_{(10aa^2a)}$ $TTCCG_{(11a^2a^2a)}$	$TTGGG_{(11aaa)}$ $TTGCC_{(11aa^2a^2)}$ $CTGCG_{(a^21aa^2a)}$ $TCGCG_{(1a^2aa^2a)}$	$TTACG_{(110a^2a)}$ $TTGTG_{(11a1a)}$ $TTGCA_{(11aa^20)}$ $GTGCG_{(a1aa^2a)}$	$TGGCG_{(1aaa^2a)}$ $TTTCG_{(111a^2a)}$ $TTGAG_{(11a0a)}$ $TTGCT_{(11aa^21)}$
$\begin{array}{c} \text{TTC} \\ (11a^2) \end{array}$	$\begin{array}{c} {\rm TTCTA}_{(11a^210)} \\ {\rm ATCTA}_{(01a^210)} \\ {\rm TACTA}_{(10a^210)} \\ {\rm TTGTA}_{(11a10)} \end{array}$	$\begin{array}{l} {\rm TTCAA}_{(11a^200)} \\ {\rm TTCTT}_{(11a^211)} \\ {\rm CTCTA}_{(a^21a^210)} \\ {\rm TCCTA}_{(1a^2a^210)} \end{array}$	$\begin{array}{l} {\rm TTTTA}_{(11110)} \\ {\rm TTCCA}_{(11a^2a^20)} \\ {\rm TTCTG}_{(11a^21a)} \\ {\rm GTCTA}_{(a1a^210)} \end{array}$	$\begin{array}{l} {\rm TGCTA}_{(1aa^{2}10)} \\ {\rm TTATA}_{(11010)} \\ {\rm TTCGA}_{(11a^{2}a0)} \\ {\rm TTCTC}_{(11a^{2}1a^{2})} \end{array}$
$TGA \\ (1a0)$	$TGACA_{(1a0a^20)}$ $AGACA_{(0a0a^20)}$ $TCACA_{(1a^20a^20)}$ $TGTCA_{(1a1a^20)}$	$TGAGA_{(1a0a0)}$ $TGACT_{(1a0a^21)}$ $CGACA_{(a^2a0a^20)}$ $TAACA_{(100a^20)}$	$TGGCA_{(1aaa^20)}$ $TGATA_{(1a010)}$ $TGACG_{(1a0a^2a)}$ $GGACA_{(aa0a^20)}$	$TTACA_{(110a^20)}$ $TGCCA_{(1aa^2a^20)}$ $TGAAA_{(1a000)}$ $TGACC_{(1a0a^2a^2)}$
$\begin{array}{c} \mathrm{TGT} \\ (1a1) \end{array}$	$TGTTG_{(1a11a)}$ $AGTTG_{(0a11a)}$ $TCTTG_{(1a^211a)}$ $TGATG_{(1a01a)}$	$\begin{array}{l} \mathrm{TGTAG}_{(1a10a)} \\ \mathrm{TGTTC}_{(1a11a^2)} \\ \mathrm{CGTTG}_{(a^2a11a)} \\ \mathrm{TATTG}_{(1011a)} \end{array}$	$\begin{array}{l} \mathrm{TGCTG}_{(1aa^21a)} \\ \mathrm{TGTCG}_{(1a1a^2a)} \\ \mathrm{TGTTA}_{(1a110)} \\ \mathrm{GGTTG}_{(aa11a)} \end{array}$	$\begin{array}{l} {\rm TTTTTG}_{(1111a)} \\ {\rm TGGTG}_{(1aa1a)} \\ {\rm TGTGG}_{(1a1aa)} \\ {\rm TGTTT}_{(1a111)} \end{array}$
$\begin{array}{c} {\rm TGG} \\ (1aa) \end{array}$	$TGGAC_{(1aa0a^2)}$ $AGGAC_{(0aa0a^2)}$ $TCGAC_{(1a^2a0a^2)}$ $TGCAC_{(1aa^20a^2)}$	$\begin{array}{l} \mathrm{TGGTC}_{(1aa1a^2)} \\ \mathrm{TGGAG}_{(1aa0a)} \\ \mathrm{CGGAC}_{(a^2aa0a^2)} \\ \mathrm{TAGAC}_{(10a0a^2)} \end{array}$	$TGAAC_{(1a00a^2)}$ $TGGGC_{(1aaaa^2)}$ $TGGAT_{(1aa01)}$ $GGGAC_{(aaa0a^2)}$	$\begin{array}{l} {\rm TTGAC}_{(11a0a^2)} \\ {\rm TGTAC}_{(1a10a^2)} \\ {\rm TGGCC}_{(1aaa^2a^2)} \\ {\rm TGGAA}_{(1aa00)} \end{array}$
$\begin{array}{c} \text{TGC} \\ (1aa^2) \end{array}$	$\begin{array}{c} \mathrm{TGCGT}_{(1aa^2a1)} \\ \mathrm{AGCGT}_{(0aa^2a1)} \\ \mathrm{TCCGT}_{(1a^2a^2a1)} \\ \mathrm{TGGGT}_{(1aaa1)} \end{array}$	$\begin{array}{l} \mathrm{TGCCT}_{(1aa^2a^21)}\\ \mathrm{TGCGA}_{(1aa^2a0)}\\ \mathrm{CGCGT}_{(a^2aa^2a1)}\\ \mathrm{TACGT}_{(10a^2a1)} \end{array}$	$TGTGT_{(1a1a1)}$ $TGCAT_{(1aa^201)}$ $TGCGC_{(1aa^2aa^2)}$ $GGCGT_{(aaa^2a1)}$	$TTCGT_{(11a^2a1)}$ $TGAGT_{(1a0a1)}$ $TGCTT_{(1aa^211)}$ $TGCGG_{(1aa^2aa)}$

 Table A.1: Degeneracy Mapping

TCA	$\mathrm{TCAGG}_{(1a^20aa)}$	$\mathrm{TCACG}_{(1a^20a^2a)}$	$\mathrm{TCGGG}_{(1a^2aaa)}$	$TAAGG_{(100aa)}$
	$ACAGG_{(0a^20aa)}$	$\mathrm{TCAGC}_{(1a^20aa^2)}$	$\mathrm{TCAAG}_{(1a^200a)}$	$\mathrm{TCCGG}_{(1a^2a^2aa)}$
$(1a^20)$	$\mathrm{TGAGG}_{(1a0aa)}$	$CCAGG_{(a^2a^20aa)}$	$TCAGA_{(1a^20a0)}$	$\mathrm{TCATG}_{(1a^201a)}$
	$\mathrm{TCTGG}_{(1a^21aa)}$	$\mathrm{TTAGG}_{(110aa)}$	$\mathrm{GCAGG}_{(aa^20aa)}$	$\mathrm{TCAGT}_{(1a^20a1)}$
	$\mathrm{TCTAA}_{(1a^2100)}$	$\mathrm{TCTTA}_{(1a^2110)}$	$\mathrm{TCCAA}_{(1a^2a^200)}$	$TATAA_{(10100)}$
TCT	$ACTAA_{(0a^2100)}$	$\mathrm{TCTAT}_{(1a^2101)}$	$\mathrm{TCTGA}_{(1a^21a0)}$	$TCGAA_{(1a^2a00)}$
$(1a^21)$	$\mathrm{TGTAA}_{(1a100)}$	$CCTAA_{(a^2a^2100)}$	$\mathrm{TCTAG}_{(1a^210a)}$	$\mathrm{TCTCA}_{(1a^21a^20)}$
	$TCAAA_{(1a^2000)}$	$TTTAA_{(11100)}$	$\operatorname{GCTAA}_{(aa^2100)}$	$\mathrm{TCTAC}_{(1a^210a^2)}$
	$\mathrm{TCGTT}_{(1a^2a11)}$	$\mathrm{TCGAT}_{(1a^2a01)}$	$\mathrm{TCATT}_{(1a^2011)}$	$TAGTT_{(10a11)}$
TCG	$ACGTT_{(0a^2a11)}$	$\mathrm{TCGTA}_{(1a^2a10)}$	$\mathrm{TCGCT}_{(1a^2aa^21)}$	$\mathrm{TCTTT}_{(1a^2111)}$
$(1a^2a)$	$\mathrm{TGGTT}_{(1aa11)}$	$\text{CCGTT}_{(a^2a^2a11)}$	$\mathrm{TCGTC}_{(1a^2a1a^2)}$	$\mathrm{TCGGT}_{(1a^2aa1)}$
	$\mathrm{TCCTT}_{(1a^2a^211)}$	$\mathrm{TTGTT}_{(11a11)}$	$\operatorname{GCGTT}_{(aa^2a11)}$	$\mathrm{TCGTG}_{(1a^2a1a)}$
	$\operatorname{TCCCC}_{(1a^2a^2a^2a^2)}$	$\mathrm{TCCGC}_{(1a^2a^2aa^2)}$	$\mathrm{TCTCC}_{(1a^21a^2a^2)}$	$\operatorname{TACCC}_{(10a^2a^2a^2)}$
TCC	$\operatorname{ACCCC}_{(0a^2a^2a^2a^2a^2)}$	$\mathrm{TCCCG}_{(1a^2a^2a^2a)}$	$\mathrm{TCCTC}_{(1a^2a^21a^2)}$	$\mathrm{TCACC}_{(1a^20a^2a^2)}$
$(1a^2a^2)$	$\mathrm{TGCCC}_{(1aa^2a^2a^2)}$	$\operatorname{CCCCC}_{(a^2a^2a^2a^2a^2a^2)}$	$\mathrm{TCCCT}_{(1a^2a^2a^21)}$	$\mathrm{TCCAC}_{(1a^2a^20a^2)}$
	$\mathrm{TCGCC}_{(1a^2aa^2a^2)}$	$\mathrm{TTCCC}_{(11a^2a^2a^2)}$	$\operatorname{GCCCC}_{(aa^2a^2a^2a^2)}$	$\mathrm{TCCCA}_{(1a^2a^2a^20)}$
	$GAAGT_{(a00a1)}$	$GAACT_{(a00a^21)}$	$GAGGT_{(a0aa1)}$	$\text{GCAGT}_{(aa^20a1)}$
GAA	$CAAGT_{(a^200a1)}$	$GAAGA_{(a00a0)}$	$GAAAT_{(a0001)}$	$GACGT_{(a0a^2a1)}$
(a00)	$\operatorname{GTAGT}_{(a10a1)}$	$AAAGT_{(000a1)}$	$GAAGC_{(a00aa^2)}$	$GAATT_{(a0011)}$
	$GATGT_{(a01a1)}$	$\mathrm{GGAGT}_{(aa0a1)}$	$TAAGT_{(100a1)}$	$GAAGG_{(a00aa)}$
	$GATAC_{(a010a^2)}$	$GATTC_{(a011a^2)}$	$GACAC_{(a0a^20a^2)}$	$\operatorname{GCTAC}_{(aa^210a^2)}$
GAT		$GATAG_{(a010a)}$	$GATGC_{(a01aa^2)}$	$GAGAC_{(a0a0a^2)}$
GAT	$CATAC_{(a^2010a^2)}$	GAIAG(a010a)	$(a01aa^2)$	$OnO(a0a0a^2)$
$\begin{array}{c} \text{GAT} \\ (a01) \end{array}$	$CATAC_{(a^2010a^2)}$ $GTTAC_{(a110a^2)}$			
I	$CATAC_{(a^2010a^2)}$ $GTTAC_{(a110a^2)}$ $GAAAC_{(a000a^2)}$	$\begin{array}{l} \text{GATAG}_{(a010a)} \\ \text{AATAC}_{(0010a^2)} \\ \text{GGTAC}_{(aa10a^2)} \end{array}$	$GATAT_{(a0101)}$ $TATAC_{(1010a^2)}$	$GATCC_{(a01a^2a^2)}$ $GATAA_{(a0100)}$
1	$\operatorname{GTTAC}_{(a110a^2)}$	$AATAC_{(0010a^2)}$	$GATAT_{(a0101)}$	$\operatorname{GATCC}_{(a01a^2a^2)}$
1	$\begin{array}{l} \operatorname{GTTAC}_{(a110a^2)} \\ \operatorname{GAAAC}_{(a000a^2)} \end{array}$	$\begin{array}{l} \text{AATAC}_{(0010a^2)} \\ \text{GGTAC}_{(aa10a^2)} \end{array}$	$\begin{array}{l} \text{GATAT}_{(a0101)} \\ \text{TATAC}_{(1010a^2)} \end{array}$	$GATCC_{(a01a^2a^2)}$ $GATAA_{(a0100)}$
(a01)	$\begin{array}{c} \operatorname{GTTAC}_{(a110a^2)} \\ \operatorname{GAAAC}_{(a000a^2)} \\ \\ \operatorname{GAGTG}_{(a0a1a)} \end{array}$	$\begin{array}{l} \text{AATAC}_{(0010a^2)} \\ \text{GGTAC}_{(aa10a^2)} \\ \\ \text{GAGAG}_{(a0a0a)} \end{array}$	$\begin{array}{c} \text{GATAT}_{(a0101)} \\ \text{TATAC}_{(1010a^2)} \\ \\ \text{GAATG}_{(a001a)} \end{array}$	$\begin{array}{c} \text{GATCC}_{(a01a^2a^2)}\\ \text{GATAA}_{(a0100)}\\ \\ \text{GCGTG}_{(aa^2a1a)} \end{array}$

 Table A.1: Degeneracy Mapping

$GAC (a0a^2)$	$\begin{array}{l} \text{GACCA}_{(a0a^2a^20)}\\ \text{CACCA}_{(a^20a^2a^20)}\\ \text{GTCCA}_{(a1a^2a^20)}\\ \text{GAGCA}_{(a0aa^20)} \end{array}$	$\begin{array}{l} \operatorname{GACGA}_{(a0a^2a0)}\\ \operatorname{GACCT}_{(a0a^2a^21)}\\ \operatorname{AACCA}_{(00a^2a^20)}\\ \operatorname{GGCCA}_{(aaa^2a^20)} \end{array}$	$\begin{array}{l} \operatorname{GATCA}_{(a01a^20)}\\ \operatorname{GACTA}_{(a0a^210)}\\ \operatorname{GACCG}_{(a0a^2a^2a)}\\ \operatorname{TACCA}_{(10a^2a^20)} \end{array}$	$\begin{array}{l} \operatorname{GCCCA}_{(aa^2a^2a^20)}\\ \operatorname{GAACA}_{(a00a^20)}\\ \operatorname{GACAA}_{(a0a^200)}\\ \operatorname{GACCC}_{(a0a^2a^2a^2)} \end{array}$
GTA ( <i>a</i> 10)	$\begin{array}{l} \operatorname{GTACC}_{(a10a^2a^2)}\\ \operatorname{CTACC}_{(a^210a^2a^2)}\\ \operatorname{GAACC}_{(a00a^2a^2)}\\ \operatorname{GTTCC}_{(a11a^2a^2)} \end{array}$	$\begin{array}{l} \operatorname{GTAGC}_{(a10aa^2)} \\ \operatorname{GTACG}_{(a10a^2a)} \\ \operatorname{ATACC}_{(010a^2a^2)} \\ \operatorname{GCACC}_{(aa^20a^2a^2)} \end{array}$	$\begin{array}{l} \operatorname{GTGCC}_{(a1aa^2a^2)} \\ \operatorname{GTATC}_{(a101a^2)} \\ \operatorname{GTACT}_{(a10a^21)} \\ \operatorname{TTACC}_{(110a^2a^2)} \end{array}$	$\begin{array}{l} \operatorname{GGACC}_{(aa0a^2a^2)} \\ \operatorname{GTCCC}_{(a1a^2a^2a^2)} \\ \operatorname{GTAAC}_{(a100a^2)} \\ \operatorname{GTACA}_{(a10a^20)} \end{array}$
GTT (a11)	$\begin{array}{l} \operatorname{GTTTT}_{(a1111)} \\ \operatorname{CTTTT}_{(a^2 1111)} \\ \operatorname{GATTT}_{(a0111)} \\ \operatorname{GTATT}_{(a1011)} \end{array}$	$\begin{array}{l} \operatorname{GTTAT}_{(a1101)} \\ \operatorname{GTTTA}_{(a1110)} \\ \operatorname{ATTTT}_{(01111)} \\ \operatorname{GCTTT}_{(aa^2111)} \end{array}$	$\begin{array}{l} \operatorname{GTCTT}_{(a1a^211)} \\ \operatorname{GTTCT}_{(a11a^21)} \\ \operatorname{GTTTC}_{(a111a^2)} \\ \operatorname{TTTTT}_{(11111)} \end{array}$	$\begin{array}{l} \text{GGTTT}_{(aa111)} \\ \text{GTGTT}_{(a1a11)} \\ \text{GTTGT}_{(a11a1)} \\ \text{GTTTG}_{(a111a)} \end{array}$
$\begin{array}{c} \text{GTG} \\ (a1a) \end{array}$	$\begin{array}{c} \operatorname{GTGAA}_{(a1a00)}\\ \operatorname{CTGAA}_{(a^21a00)}\\ \operatorname{GAGAA}_{(a0a00)}\\ \operatorname{GTCAA}_{(a1a^200)} \end{array}$	$\begin{array}{c} \mathrm{GTGTA}_{(a1a10)} \\ \mathrm{GTGAT}_{(a1a01)} \\ \mathrm{ATGAA}_{(01a00)} \\ \mathrm{GCGAA}_{(aa^2a00)} \end{array}$	$\begin{array}{c} \operatorname{GTAAA}_{(a1000)} \\ \operatorname{GTGGA}_{(a1aa0)} \\ \operatorname{GTGAG}_{(a1a0a)} \\ \operatorname{TTGAA}_{(11a00)} \end{array}$	$\begin{array}{l} \operatorname{GGGAA}_{(aaa00)} \\ \operatorname{GTTAA}_{(a1100)} \\ \operatorname{GTGCA}_{(a1aa^20)} \\ \operatorname{GTGAC}_{(a1a0a^2)} \end{array}$
$\begin{array}{c} \text{GTC} \\ (a1a^2) \end{array}$	$\begin{array}{c} \operatorname{GTCGG}_{(a1a^2aa)}\\ \operatorname{CTCGG}_{(a^21a^2aa)}\\ \operatorname{GACGG}_{(a0a^2aa)}\\ \operatorname{GTGGG}_{(a1aaa)} \end{array}$	$\begin{array}{l} \operatorname{GTCCG}_{(a1a^2a^2a)}\\ \operatorname{GTCGC}_{(a1a^2aa^2)}\\ \operatorname{ATCGG}_{(01a^2aa)}\\ \operatorname{GCCGG}_{(aa^2a^2aa)} \end{array}$	$\begin{array}{l} \operatorname{GTTGG}_{(a11aa)}\\ \operatorname{GTCAG}_{(a1a^20a)}\\ \operatorname{GTCGA}_{(a1a^2a0)}\\ \operatorname{TTCGG}_{(11a^2aa)} \end{array}$	$\begin{array}{l} \operatorname{GGCGG}_{(aaa^2aa)} \\ \operatorname{GTAGG}_{(a10aa)} \\ \operatorname{GTCTG}_{(a1a^21a)} \\ \operatorname{GTCGT}_{(a1a^2a1)} \end{array}$
GGA (aa0)	$\begin{array}{c} \mathrm{GGAAG}_{(aa00a)}\\ \mathrm{CGAAG}_{(a^2a00a)}\\ \mathrm{GCAAG}_{(aa^200a)}\\ \mathrm{GGTAG}_{(aa10a)} \end{array}$	$\begin{array}{l} \text{GGATG}_{(aa01a)} \\ \text{GGAAC}_{(aa00a^2)} \\ \text{AGAAG}_{(0a00a)} \\ \text{GAAAG}_{(a000a)} \end{array}$	$\begin{array}{l} \operatorname{GGGAG}_{(aaa0a)} \\ \operatorname{GGAGG}_{(aa0aa)} \\ \operatorname{GGAAA}_{(aa000)} \\ \operatorname{TGAAG}_{(1a00a)} \end{array}$	$\begin{array}{l} \operatorname{GTAAG}_{(a100a)} \\ \operatorname{GGCAG}_{(aaa^20a)} \\ \operatorname{GGACG}_{(aa0a^2a)} \\ \operatorname{GGAAT}_{(aa001)} \end{array}$
$\begin{array}{c} \text{GGT} \\ (aa1) \end{array}$	$\begin{array}{c} \mathrm{GGTGA}_{(aa1a0)}\\ \mathrm{CGTGA}_{(a^2a1a0)}\\ \mathrm{GCTGA}_{(aa^21a0)}\\ \mathrm{GGAGA}_{(aa0a0)} \end{array}$	$\begin{array}{l} \operatorname{GGTCA}_{(aa1a^20)} \\ \operatorname{GGTGT}_{(aa1a1)} \\ \operatorname{AGTGA}_{(0a1a0)} \\ \operatorname{GATGA}_{(a01a0)} \end{array}$	$\begin{array}{l} \operatorname{GGCGA}_{(aaa^2a0)} \\ \operatorname{GGTAA}_{(aa100)} \\ \operatorname{GGTGG}_{(aa1aa)} \\ \operatorname{TGTGA}_{(1a1a0)} \end{array}$	$\begin{array}{l} \operatorname{GTTGA}_{(a11a0)} \\ \operatorname{GGGGGA}_{(aaaa0)} \\ \operatorname{GGTTA}_{(aa110)} \\ \operatorname{GGTGC}_{(aa1aa^2)} \end{array}$

 Table A.1: Degeneracy Mapping

GGG (aaa)	$\begin{array}{l} \text{GGGCT}_{(aaaa^21)} \\ \text{CGGCT}_{(a^2aaa^21)} \\ \text{GCGCT}_{(aa^2aa^21)} \\ \text{GGCCT}_{(aaa^2a^21)} \end{array}$	$\begin{array}{l} \mathrm{GGGGGT}_{(aaaa1)} \\ \mathrm{GGGCA}_{(aaaa^20)} \\ \mathrm{AGGCT}_{(0aaa^21)} \\ \mathrm{GAGCT}_{(a0aa^21)} \end{array}$	$\begin{array}{l} \mathrm{GGACT}_{(aa0a^21)}\\ \mathrm{GGGTT}_{(aaa11)}\\ \mathrm{GGGCC}_{(aaaa^2a^2)}\\ \mathrm{TGGCT}_{(1aaa^21)} \end{array}$	$\begin{array}{l} \operatorname{GTGCT}_{(a1aa^21)}\\ \operatorname{GGTCT}_{(aa1a^21)}\\ \operatorname{GGGAT}_{(aaa01)}\\ \operatorname{GGGCG}_{(aaaa^2a)} \end{array}$
$\begin{array}{c} \text{GGC} \\ (aaa^2) \end{array}$	$\begin{array}{c} \text{GGCTC}_{(aaa^21a^2)} \\ \text{CGCTC}_{(a^2aa^21a^2)} \\ \text{GCCTC}_{(aa^2a^21a^2)} \\ \text{GGGTC}_{(aaa1a^2)} \end{array}$	$\begin{array}{c} \operatorname{GGCAC}_{(aaa^20a^2)} \\ \operatorname{GGCTG}_{(aaa^21a)} \\ \operatorname{AGCTC}_{(0aa^21a^2)} \\ \operatorname{GACTC}_{(a0a^21a^2)} \end{array}$	$\begin{array}{c} \text{GGTTC}_{(aa11a^2)} \\ \text{GGCCC}_{(aaa^2a^2a^2)} \\ \text{GGCTT}_{(aaa^211)} \\ \text{TGCTC}_{(1aa^21a^2)} \end{array}$	$\begin{array}{l} \operatorname{GTCTC}_{(a1a^21a^2)} \\ \operatorname{GGATC}_{(aa01a^2)} \\ \operatorname{GGCGC}_{(aaa^2aa^2)} \\ \operatorname{GGCTA}_{(aaa^210)} \end{array}$
$\begin{array}{c} \text{GCA} \\ (aa^20) \end{array}$	$\begin{array}{l} \operatorname{GCATA}_{(aa^2010)}\\ \operatorname{CCATA}_{(a^2a^2010)}\\ \operatorname{GGATA}_{(aa010)}\\ \operatorname{GCTTA}_{(aa^2110)} \end{array}$	$\begin{array}{l} \operatorname{GCAAA}_{(aa^2000)} \\ \operatorname{GCATT}_{(aa^2011)} \\ \operatorname{ACATA}_{(0a^2010)} \\ \operatorname{GTATA}_{(a1010)} \end{array}$	$\begin{array}{l} \operatorname{GCGTA}_{(aa^2a10)} \\ \operatorname{GCACA}_{(aa^20a^20)} \\ \operatorname{GCATG}_{(aa^201a)} \\ \operatorname{TCATA}_{(1a^2010)} \end{array}$	$\begin{array}{l} \text{GAATA}_{(a0010)} \\ \text{GCCTA}_{(aa^2a^210)} \\ \text{GCAGA}_{(aa^20a0)} \\ \text{GCATC}_{(aa^201a^2)} \end{array}$
$\begin{array}{c} \text{GCT} \\ (aa^21) \end{array}$	$\begin{array}{c} \operatorname{GCTCG}_{(aa^21a^2a)} \\ \operatorname{CCTCG}_{(a^2a^21a^2a)} \\ \operatorname{GGTCG}_{(aa1a^2a)} \\ \operatorname{GCACG}_{(aa^20a^2a)} \end{array}$	$\begin{array}{c} \operatorname{GCTGG}_{(aa^21aa)} \\ \operatorname{GCTCC}_{(aa^21a^2a^2)} \\ \operatorname{ACTCG}_{(0a^21a^2a)} \\ \operatorname{GTTCG}_{(a11a^2a)} \end{array}$	$\begin{array}{c} \operatorname{GCCCG}_{(aa^2a^2a^2a)} \\ \operatorname{GCTTG}_{(aa^211a)} \\ \operatorname{GCTCA}_{(aa^21a^20)} \\ \operatorname{TCTCG}_{(1a^21a^2a)} \end{array}$	$\begin{array}{l} \text{GATCG}_{(a01a^2a)} \\ \text{GCGCG}_{(aa^2aa^2a)} \\ \text{GCTAG}_{(aa^210a)} \\ \text{GCTCT}_{(aa^21a^21)} \end{array}$
$\begin{array}{c} \mathrm{GCG} \\ (aa^2a) \end{array}$	$\begin{array}{l} \operatorname{GCGGC}_{(aa^2aaa^2)} \\ \operatorname{CCGGC}_{(a^2a^2aaa^2)} \\ \operatorname{GGGGC}_{(aaaaa^2)} \\ \operatorname{GCCGC}_{(aa^2a^2aa^2)} \end{array}$	$\begin{array}{l} \operatorname{GCGCC}_{(aa^2aa^2a^2)} \\ \operatorname{GCGGG}_{(aa^2aaa)} \\ \operatorname{ACGGC}_{(0a^2aaa^2)} \\ \operatorname{GTGGC}_{(a1aaa^2)} \end{array}$	$\begin{array}{l} \operatorname{GCAGC}_{(aa^20aa^2)} \\ \operatorname{GCGAC}_{(aa^2a0a^2)} \\ \operatorname{GCGGT}_{(aa^2aa1)} \\ \operatorname{TCGGC}_{(1a^2aaa^2)} \end{array}$	$\begin{array}{l} \operatorname{GAGGC}_{(a0aaa^2)} \\ \operatorname{GCTGC}_{(aa^21aa^2)} \\ \operatorname{GCGTC}_{(aa^2a1a^2)} \\ \operatorname{GCGGA}_{(aa^2aa0)} \end{array}$
$\begin{array}{c} \text{GCC} \\ (aa^2a^2) \end{array}$	$\begin{array}{l} \operatorname{GCCAT}_{(aa^2a^201)}\\ \operatorname{CCCAT}_{(a^2a^2a^201)}\\ \operatorname{GGCAT}_{(aaa^201)}\\ \operatorname{GCGAT}_{(aa^2a01)} \end{array}$	$\begin{array}{l} \operatorname{GCCTT}_{(aa^2a^211)}\\ \operatorname{GCCAA}_{(aa^2a^200)}\\ \operatorname{ACCAT}_{(0a^2a^201)}\\ \operatorname{GTCAT}_{(a1a^201)} \end{array}$	$\begin{array}{l} \operatorname{GCTAT}_{(aa^2101)}\\ \operatorname{GCCGT}_{(aa^2a^2a1)}\\ \operatorname{GCCAC}_{(aa^2a^20a^2)}\\ \operatorname{TCCAT}_{(1a^2a^201)} \end{array}$	$\begin{array}{l} \operatorname{GACAT}_{(a0a^201)}\\ \operatorname{GCAAT}_{(aa^2001)}\\ \operatorname{GCCCT}_{(aa^2a^2a^21)}\\ \operatorname{GCCAG}_{(aa^2a^20a)} \end{array}$
$\begin{array}{c} \text{CAA} \\ (a^2 00) \end{array}$	$CAACG_{(a^200a^2a)}$ $GAACG_{(a00a^2a)}$ $CTACG_{(a^210a^2a)}$ $CATCG_{(a^201a^2a)}$	$CAAGG_{(a^200aa)}$ $CAACC_{(a^200a^2a^2)}$ $TAACG_{(100a^2a)}$ $CGACG_{(a^2a0a^2a)}$	$\begin{array}{c} \text{CAGCG}_{(a^20aa^2a)}\\ \text{CAATG}_{(a^2001a)}\\ \text{CAACA}_{(a^200a^20)}\\ \text{AAACG}_{(000a^2a)} \end{array}$	$CCACG_{(a^2a^20a^2a)}$ $CACCG_{(a^20a^2a^2a)}$ $CAAAG_{(a^2000a)}$ $CAAAG_{(a^200a^21)}$

 Table A.1: Degeneracy Mapping

$\begin{array}{c} \text{CAT} \\ (a^2 0 1) \end{array}$	CATTA <sub>(<math>a^{2}0110</math>)</sub> GATTA <sub>(<math>a0110</math>)</sub> CTTTA <sub>(<math>a^{2}1110</math>)</sub> CAATA <sub>(<math>a^{2}0010</math>)</sub>	$CATAA_{(a^20100)}$ $CATTT_{(a^20111)}$ $TATTA_{(10110)}$ $CGTTA_{(a^2a110)}$	$CACTA_{(a^20a^210)}$ $CATCA_{(a^201a^20)}$ $CATTG_{(a^2011a)}$ $AATTA_{(00110)}$	$CCTTA_{(a^2a^2110)}$ $CAGTA_{(a^20a10)}$ $CATGA_{(a^201a0)}$ $CATTC_{(a^2011a^2)}$
$\begin{array}{c} \text{CAG} \\ (a^20a) \end{array}$	$CAGAT_{(a^20a01)}$ $GAGAT_{(a0a01)}$ $CTGAT_{(a^21a01)}$ $CACAT_{(a^20a^201)}$	$CAGTT_{(a^{2}0a11)}$ $CAGAA_{(a^{2}0a00)}$ $TAGAT_{(10a01)}$ $CGGAT_{(a^{2}aa01)}$	$\begin{array}{l} {\rm CAAAT}_{(a^20001)}\\ {\rm CAGGT}_{(a^20aa1)}\\ {\rm CAGAC}_{(a^20a0a^2)}\\ {\rm AAGAT}_{(00a01)} \end{array}$	$CCGAT_{(a^2a^2a01)}$ $CATAT_{(a^20101)}$ $CAGCT_{(a^20aa^21)}$ $CAGAG_{(a^20a0a)}$
$CAC \\ (a^2 0 a^2)$	$CACGC_{(a^20a^2aa^2)}$ $GACGC_{(a0a^2aa^2)}$ $CTCGC_{(a^21a^2aa^2)}$ $CAGGC_{(a^20aaa^2)}$	$CACCC_{(a^20a^2a^2a^2)}$ $CACGG_{(a^20a^2aa)}$ $TACGC_{(10a^2aa^2)}$ $CGCGC_{(a^2aa^2aa^2)}$	$CATGC_{(a^201aa^2)}$ $CACAC_{(a^20a^20a^2)}$ $CACGT_{(a^20a^2a1)}$ $AACGC_{(00a^2aa^2)}$	$CCCGC_{(a^2a^2a^2a^2a^2)}$ $CAAGC_{(a^200aa^2)}$ $CACTC_{(a^20a^21a^2)}$ $CACGA_{(a^20a^2a0)}$
$CTA (a^2 10)$	$CTAGA_{(a^210a0)}$ $GTAGA_{(a10a0)}$ $CAAGA_{(a^200a0)}$ $CTTGA_{(a^211a0)}$	$\begin{array}{c} \text{CTACA}_{(a^210a^20)}\\ \text{CTAGT}_{(a^210a1)}\\ \text{TTAGA}_{(110a0)}\\ \text{CCAGA}_{(a^2a^20a0)} \end{array}$	$\begin{array}{c} \text{CTGGA}_{(a^21aa0)}\\ \text{CTAAA}_{(a^21000)}\\ \text{CTAGG}_{(a^210aa)}\\ \text{ATAGA}_{(010a0)} \end{array}$	$\begin{array}{c} \text{CGAGA}_{(a^2a0a0)} \\ \text{CTCGA}_{(a^21a^2a0)} \\ \text{CTATA}_{(a^21010)} \\ \text{CTAGC}_{(a^210aa^2)} \end{array}$
$\begin{array}{c} \text{CTT} \\ (a^2 11) \end{array}$	$CTTAG_{(a^2110a)}$ $GTTAG_{(a110a)}$ $CATAG_{(a^2010a)}$ $CTAAG_{(a^2100a)}$	$CTTTG_{(a^2111a)}$ $CTTAC_{(a^2110a^2)}$ $TTTAG_{(1110a)}$ $CCTAG_{(a^2a^210a)}$	$CTCAG_{(a^21a^20a)}$ $CTTGG_{(a^211aa)}$ $CTTAA_{(a^21100)}$ $ATTAG_{(0110a)}$	$CGTAG_{(a^2a10a)}$ $CTGAG_{(a^21a0a)}$ $CTTCG_{(a^211a^2a)}$ $CTTAT_{(a^21101)}$
$\begin{array}{c} \text{CTG} \\ (a^2 1 a) \end{array}$	$CTGTC_{(a^21a1a^2)}$ $GTGTC_{(a1a1a^2)}$ $CAGTC_{(a^20a1a^2)}$ $CTCTC_{(a^21a^21a^2)}$	$CTGAC_{(a^21a0a^2)}$ $CTGTG_{(a^21a1a)}$ $TTGTC_{(11a1a^2)}$ $CCGTC_{(a^2a^2a1a^2)}$	$\begin{array}{l} \text{CTATC}_{(a^2101a^2)} \\ \text{CTGCC}_{(a^21aa^2a^2)} \\ \text{CTGTT}_{(a^21a11)} \\ \text{ATGTC}_{(01a1a^2)} \end{array}$	$\begin{array}{c} \text{CGGTC}_{(a^2aa1a^2)} \\ \text{CTTTC}_{(a^2111a^2)} \\ \text{CTGGC}_{(a^21aaa^2)} \\ \text{CTGTA}_{(a^21a10)} \end{array}$
$\begin{array}{c} \text{CTC} \\ (a^2 1 a^2) \end{array}$	$CTCCT_{(a^21a^2a^21)}$ $GTCCT_{(a1a^2a^21)}$ $CACCT_{(a^20a^2a^21)}$ $CTGCT_{(a^21aa^21)}$	$CTCGT_{(a^21a^2a^1)}$ $CTCCA_{(a^21a^2a^20)}$ $TTCCT_{(11a^2a^21)}$ $CCCCT_{(a^2a^2a^2a^21)}$	$\begin{array}{l} \text{CTTCT}_{(a^211a^21)}\\ \text{CTCTT}_{(a^21a^211)}\\ \text{CTCCC}_{(a^21a^2a^2a^2)}\\ \text{ATCCT}_{(01a^2a^21)} \end{array}$	$CGCCT_{(a^2aa^2a^21)}$ $CTACT_{(a^210a^21)}$ $CTCAT_{(a^21a^201)}$ $CTCCG_{(a^21a^2a^2a)}$

 Table A.1: Degeneracy Mapping

[ ]				
	$CGATT_{(a^2a011)}$	$CGAAT_{(a^2a001)}$	$CGGTT_{(a^2aa11)}$	$CTATT_{(a^21011)}$
$\begin{array}{c} \text{CGA} \\ (a^2a0) \end{array}$	$GGATT_{(aa011)}$	$CGATA_{(a^2a010)}$	$\operatorname{CGACT}_{(a^2a0a^21)}$	$\operatorname{CGCTT}_{(a^2aa^211)}$
$(u \ u 0)$	$\operatorname{CCATT}_{(a^2a^2011)}$	$\mathrm{TGATT}_{(1a011)}$	$\mathrm{CGATC}_{(a^2a01a^2)}$	$\operatorname{CGAGT}_{(a^2a0a1)}$
	$\mathrm{CGTTT}_{(a^2a111)}$	$CAATT_{(a^20011)}$	$AGATT_{(0a011)}$	$\mathrm{CGATG}_{(a^2a01a)}$
	$\operatorname{CGTCC}_{(a^2a1a^2a^2)}$	$\mathrm{CGTGC}_{(a^2a1aa^2)}$	$\operatorname{CGCCC}_{(a^2aa^2a^2a^2)}$	$\operatorname{CTTCC}_{(a^2 1 1 a^2 a^2)}$
CGT	$\mathrm{GGTCC}_{(aa1a^2a^2)}$	$\mathrm{CGTCG}_{(a^2a1a^2a)}$	$\mathrm{CGTTC}_{(a^2a11a^2)}$	$\mathrm{CGGCC}_{(a^2aaa^2a^2)}$
$(a^2a1)$	$\operatorname{CCTCC}_{(a^2a^21a^2a^2)}$	$\mathrm{TGTCC}_{(1a1a^2a^2)}$	$\operatorname{CGTCT}_{(a^2a1a^21)}$	$\mathrm{CGTAC}_{(a^2a10a^2)}$
	$\mathrm{CGACC}_{(a^2a0a^2a^2)}$	$\operatorname{CATCC}_{(a^201a^2a^2)}$	$\operatorname{AGTCC}_{(0a1a^2a^2)}$	$\mathrm{CGTCA}_{(a^2a1a^20)}$
	$CGGGG(a^2aaaa)$	$CGGCG_{(a^2aaa^2a)}$	$CGAGG_{(a^2a0aa)}$	$CTGGG_{(a^21aaa)}$
CGG	$\mathrm{GGGGG}_{(aaaaa)}$	$\mathrm{CGGGC}_{(a^2aaaa^2)}$	$\mathrm{CGGAG}_{(a^2aa0a)}$	$\mathrm{CGTGG}_{(a^2a1aa)}$
$(a^2aa)$	$CCGGG_{(a^2a^2aaa)}$	$\mathrm{TGGGG}_{(1aaaa)}$	$CGGGA_{(a^2aaa0)}$	$\mathrm{CGGTG}_{(a^2aa1a)}$
	$\mathrm{CGCGG}_{(a^2aa^2aa)}$	$CAGGG_{(a^20aaa)}$	$\mathrm{AGGGG}_{(0aaaa)}$	$CGGGT_{(a^2aaa1)}$
	$CGCAA_{(a^2aa^200)}$	$CGCTA_{(a^2aa^210)}$	$CGTAA_{(a^2a100)}$	$CTCAA_{(a^21a^200)}$
CGC	$\mathrm{GGCAA}_{(aaa^200)}$	$\operatorname{CGCAT}_{(a^2aa^201)}$	$CGCGA_{(a^2aa^2a0)}$	$CGAAA_{(a^2a000)}$
$(a^2aa^2)$	$CCCAA_{(a^2a^2a^200)}$	$\mathrm{TGCAA}_{(1aa^200)}$	$\operatorname{CGCAG}_{(a^2aa^20a)}$	$\mathrm{CGCCA}_{(a^2aa^2a^20)}$
	$\mathrm{CGGAA}_{(a^2aa00)}$	$CACAA_{(a^20a^200)}$	$AGCAA_{(0aa^200)}$	$\mathrm{CGCAC}_{(a^2aa^20a^2)}$
	$\operatorname{CCAAC}_{(a^2a^200a^2)}$	$\operatorname{CCATC}_{(a^2a^201a^2)}$	$\operatorname{CCGAC}_{(a^2a^2a0a^2)}$	$CAAAC_{(a^2000a^2)}$
CCA	$\operatorname{GCAAC}_{(aa^200a^2)}$	$\operatorname{CCAAG}_{(a^2a^200a)}$	$\operatorname{CCAGC}_{(a^2a^20aa^2)}$	$\operatorname{CCCAC}_{(a^2a^2a^20a^2)}$
$(a^2a^20)$	$\mathrm{CGAAC}_{(a^2a00a^2)}$	$\mathrm{TCAAC}_{(1a^200a^2)}$	$\operatorname{CCAAT}_{(a^2a^2001)}$	$\operatorname{CCACC}_{(a^2a^20a^2a^2)}$
	$\operatorname{CCTAC}_{(a^2a^210a^2)}$	$\mathrm{CTAAC}_{(a^2100a^2)}$	$ACAAC_{(0a^200a^2)}$	$\mathrm{CCAAA}_{(a^2a^2000)}$
	$\operatorname{CCTGT}_{(a^2a^21a1)}$	$\operatorname{CCTCT}_{(a^2a^21a^21)}$	$\operatorname{CCCGT}_{(a^2a^2a^2a^1)}$	$CATGT_{(a^201a1)}$
CCT	$GCTGT_{(aa^21a1)}$	$CCTGA_{(a^2a^21a0)}$	$\operatorname{CCTAT}_{(a^2a^2101)}$	$\mathrm{CCGGT}_{(a^2a^2aa1)}$
$(a^2a^21)$	$\mathrm{CGTGT}_{(a^2a1a1)}$	$\mathrm{TCTGT}_{(1a^21a1)}$	$\mathrm{CCTGC}_{(a^2a^21aa^2)}$	$\operatorname{CCTTT}_{(a^2a^2111)}$
	$\operatorname{CCAGT}_{(a^2a^20a1)}$	$\mathrm{CTTGT}_{(a^211a1)}$	$ACTGT_{(0a^21a1)}$	$CCTGG_{(a^2a^21aa)}$
	$CCGCA_{(a^2a^2aa^20)}$	$CCGGA_{(a^2a^2aa0)}$	$CCACA_{(a^2a^20a^20)}$	$CAGCA_{(a^20aa^20)}$
CCG	$\mathrm{GCGCA}_{(aa^2aa^20)}$	$\mathrm{CCGCT}_{(a^2a^2aa^21)}$	$\mathrm{CCGTA}_{(a^2a^2a10)}$	$CCTCA_{(a^2a^21a^20)}$
$(a^2a^2a)$	$CGGCA_{(a^2aaa^20)}$	$\mathrm{TCGCA}_{(1a^2aa^20)}$	$\operatorname{CCGCG}_{(a^2a^2aa^2a)}$	$CCGAA_{(a^2a^2a00)}$
	$\mathrm{CCCCA}_{(a^2a^2a^2a^2)}$	$CTGCA_{(a^21aa^20)}$	$ACGCA_{(0a^2aa^20)}$	$\operatorname{CCGCC}_{(a^2a^2aa^2a^2)}$
L				

 Table A.1: Degeneracy Mapping

 Table A.1: Degeneracy Mapping

COO	$CCCTG_{(a^2a^2a^21a)}$	$CCCAG_{(a^2a^2a^20a)}$	$CCTTG_{(a^2a^211a)}$	$CACTG_{(a^20a^21a)}$
$\begin{array}{c} \text{CCC} \\ (a^2a^2a^2) \end{array}$	$\begin{array}{l} \operatorname{GCCTG}_{(aa^{2}a^{2}1a)}\\ \operatorname{CGCTG}_{(a^{2}aa^{2}1a)}\\ \operatorname{CCGTG}_{(a^{2}a^{2}a1a)}\end{array}$	$CCCTC_{(a^2a^2a^21a^2)}$ $TCCTG_{(1a^2a^21a)}$ $CTCTG_{(a^21a^21a)}$	$CCCCG_{(a^2a^2a^2a^2a)}$ $CCCTA_{(a^2a^2a^2a)}$ $ACCTG_{(0a^2a^2a)}$	$CCATG_{(a^2a^201a)}$ $CCCGG_{(a^2a^2a^2aa)}$ $CCCTT_{(a^2a^2a^211)}$