

The Big ORF Theory: The Intensive Enumeration & The Triplet Approximation

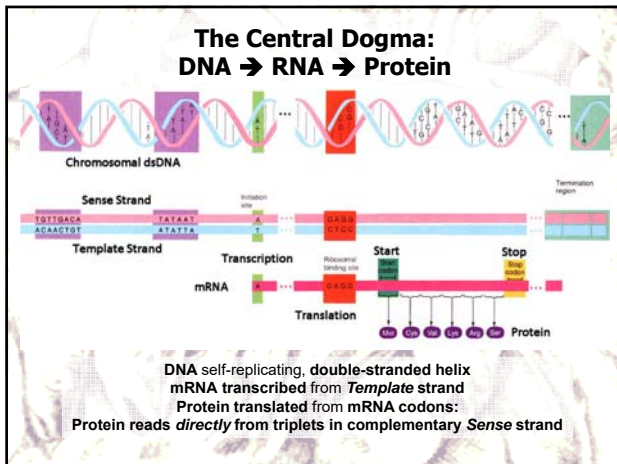
^{1,2} Steven M. Carr

(with contributions from ²Todd Wareham & ³Donald Craig)

¹Departments of Biology and ²Computer Science, and

³eHealth Research Unit (Faculty of Medicine)
Memorial University of Newfoundland

Biol4241 - 08 March 2016



The Genetic Code, 1965 (Nobel Prize, 1968)

TABLE 3. NUCLEOTIDE SEQUENCES OF RNA CODONS

1st Base	U	2nd Base C	A	G	3rd Base
U	PHE* PHE*	SER* SER*	TYR* TYR*	CYS* CYS*	U C
	leu* leu*	SER* SER*	TERM! TERM!	cyo! TRP*	A G
C	leu* leu*	pro* pro*	HIS* HIS*	ARG* ARG*	U C
	leu LEU	PRO* PRO	GLN* glu*	ARG* arg	A G
A	ILE* ILE*	THR* THR*	ASN* ASN*	SER* SER*	U C
	ile* MET*, F-MET	THR* THR	LVS* lys	arg* arg	A G
G	VAL* VAL*	ALA* ALA*	ASP* ASP*	GLY* GLY*	U C
	VAL* VAL	ALA* ALA	GLU* glu	GLY* GLY	A G

64 RNA codons: 61 coding & 3 "Termination"

	T	C	A	G	
T	F	S	Y	C	T
	F	S	Y	C	C
	L	*	*	*	A
	L	S	*	W	G
C	L	P	H	R	T
	L	P	H	R	C
	L	P	Q	R	A
	L	P	Q	R	G
A	I	T	N	S	T
	I	T	N	S	C
	I	T	K	S	A
	M	T	K	R	G
G	V	A	D	G	T
	V	A	D	G	C
	V	A	E	G	A
	V	A	E	G	G

E G A K K P M S A I I I N A P C A F V N
 G L S E Q S L E L M F L *
 1 V V V V V L D D D D D V V V V V V V V V V V V D D D V V D D V V D
 6 A T T T A C A G G G G C A T T A A T T C A A T G A T T G C T C A T G G C T A G C C T
 I Y T G A A L I L M I A H M A L
 L Q G H L F W L L M A L P

"5&1 condition" functions: enumerations, algorithms, approximations, simulations
Implications for ORF features; Alternative Genetic Codes

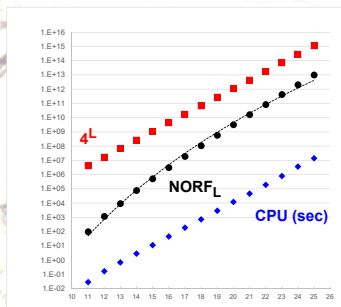
```
function CompleteSequence(OidNum, rh)  
    if rh = 7 then  
        return CompleteSequence(OidNum, 5)  
    elif rh <= OidNum  
        return CompleteSequence(OidNum, rh + 1)  
    else  
        res = Null  
        for i = StartOfSequence(OidNum) to r in randomisation of 1  
            if i is possible pair position for reading frame of  
                3 do  
                    St = place stop codon OidNum at position pos relative  
                        to i  
                    if stop placement is possible then  
                        res = GenerateSequence(OidNum, rh + 1)  
                    if res is equal to Null then  
                        exit for loop  
        return res  
  
function CompleteRandomize(OidNum, 5)  
    if 5 has more unrefined bases  
        res = Null  
        pos = position of unfilled base i  
        for base in randomisation of list ("A", "G", "C", "T") do  
            St = 5  
            Sipped = 0  
            while number of stops in open reading frame of St equals 0 then  
                res = CompleteSequence(OidNum, 5)  
                Sipped = Sipped + 1  
            exit for loop  
  
function generateRandomOid(sequen)  
    Let st be a sequence of lengths equal to unfilled numbers  
    Random = random selection from reading frame numbers 1..6  
    return concatenate(st, Random)
```

(I) Enumeration of "5&1" solutions

	NORF(L)	4 ^L	ratio
10	0	1,048,576	-
11	96	4,194,304	43,691
12	1,152	16,777,216	14,564
13	9,216	67,108,864	7,282
14	76,320	268,435,456	3,517
15	511,104	1,073,741,824	2,101
16	3,122,688	4,294,967,296	1,375
17	19,286,112	17,179,869,184	891
18	108,498,048	68,719,476,736	633
19	588,598,272	274,877,906,944	467
20	3,204,880,608	1,099,511,627,776	343
21	16,526,184,576	4,398,046,511,104	266
22	83,667,290,112	17,592,186,044,416	210
23	424,013,102,496	70,368,744,177,664	166
24	2,072,865,313,536	281,474,976,710,656	136
25	10,026,040,699,392	1,125,899,906,842,620	112

There are no "5&1" solutions for $L \leq 10$ bp : QED
 96 for $L=11$ bp ; NORF_L increases exponentially

Solutions are rare wrt 4^L :
 Random search for exemplars is inefficient

Enumerated "5&1" solutions length L (NORF_L)
 versus total sequences (4^L) & required CPU time


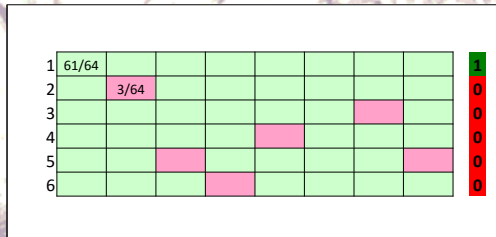
Enumeration succumbs to storage and (or) CPU limitations @ L=25
 How does NORF_L behave wrt 4^L if L >> 25?

(II) Triplet Approximation: Theory

Genetic Code comprises 61 coding & 3 stop triplets:

Disregard: triplet composition, reading frame overlap, 5' → 3', etc.

$$p(\text{coding}) = 61/64 \quad p(\text{stop}) = 3/64$$



(II) Triplet Approximation: Computation

Consider a dsDNA length L bp with $T = L / 3$ triplets
 Genetic Code comprises 61 coding & 3 stop triplets:
 Disregard triplet composition, reading frame overlap, etc.
 Prob that random triplet is coding is $C = 61/64$

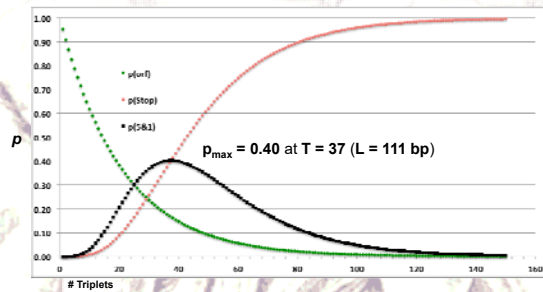
For any set (string) of T random triplets:

- 1) Prob of an ORF: $p(\text{ORF}) = C^T$
- 2) Prob of ≥ 1 stops: $p(\text{Stop}) = 1 - C^T$
- 3) Prob RF1 Open & RFs2-6 closed: $= (C^T)(1 - C^T)^5$

Then: Joint Prob for any of six RFs $p(5\&1) = (6)(C^T)(1 - C^T)^5$

(II.i) Triplet approximation:

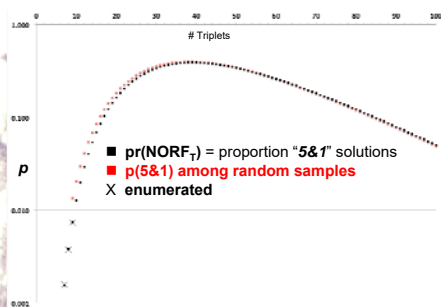
$$p(5\&1) = (6)(p(\text{ORF}))(p(\text{Stop}))^5$$



0.1 Kbp dsDNA has optimal (40%) chance of including single ORF
 Solutions limited first by "multiple open", later by "all closed" frames
 Optimum size for Natural Selection to acquire / maintain / extend function?

(II.ii) Triplet simulation:

Take 10^6 random samples of dsDNA length T



Approximation provides good upper bound for Monte Carlo Simulation
 [Simulation is CPU-intensive]

(III) Alternative Genetic Codes

	T	C	A	G	
T	F	S	Y	C	T
	F	S	Y	C	C
	L	S	Y	W	A
	L	S	*	W	G
C	L	P	H	R	T
	L	P	H	R	C
	L	P	Q	R	A
	L	P	Q	R	G
A	I	T	N	S	T
	I	T	N	S	C
	I	T	K	R	A
	M	T	K	R	G
G	V	A	D	G	T
	V	A	D	G	C
	V	A	E	G	A
	V	A	E	G	G

There are 3 stops: what if there were 4, 2, or 1 ?

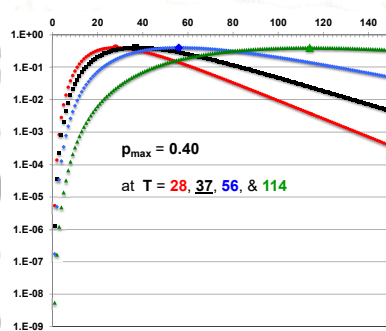
Evolution of animal mtDNA Genetic Codes: Triplet re-assignments

Code	Phylogenetic occurrence	Triplet Position	1	2	3
1	Universal		T	T	T
4	Ophisthokonta (Fungi + Animalia)		T	T	C
3	Fungi (inc. Yeast)		T	T	C
5	Protostomia, Urochordata I. (Appendicularia)		T	T	C
9	Acoelomata (Platyhelminthes), Echinodermata, Hemichordata II (Enteropneusta)		T	T	C
14	Trichobalhorzia, Radolphus		T	T	C
24	Hemichordata II (Pterobranchia)		T	T	C
13	Urochordata II (Ascidia + Thaliacea)		T	T	C
2	Chordata		T	T	C

NCBI GenBank codes, October 2014

Which (Code) came first, the **Acoelomates** or the **Protostomes** ?

p(5&1) for genetic codes with S = 4, 3, 2, or 1 stop triplets



Why *three* stops? Does $S = 3$ optimize something ?

Summary:

- **Carr's Conjecture** is answered: *no* ORF exemplars < 11 bp
- Enumeration of "5&1" DNA exemplars is exponential & CPU limited
- Triplet approximation shows
 - $NORF_1$ function has unexpected shape & implications
 - Initially limited by improbability of simultaneous closed frames
 - Ultimately limited by improbability of ORF
 - Alternative Genetic Codes with $n = 3$ vs 1, 2, or 4 stop codons
 - **Random ORFs** are not rare: $P_{max} = 0.40$
 - Optimum size varies $84 \sim (111) \sim 352$ bp
 - Confirmed by random sampling of **dsDNA** space
- "There's an app for that:" **RandomORF**
 - Carr, Craig, & Wareham (2014) *CBE Life Sci Educ* 13,68
 - <http://www.ucs.mun.ca/~donald/orf/biocomp/>

Acknowledgements

Todd Wareham, Dept of Computer Science &
Donald Craig, eHealth Research Unit (Faculty of
Medicine)
Memorial University of Newfoundland
NSERC Discovery Grants to Drs Carr & Wareham

Acknowledgements



Students in Biol2250 – Principles of Genetics
