

ECE 551 Homework 2 (Spring 2016)
March 25th, 2016

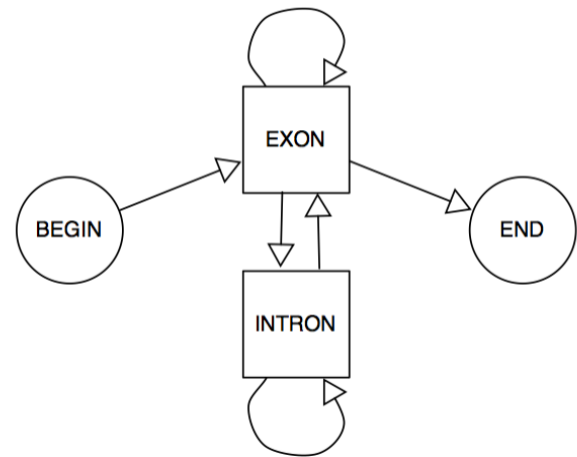
Q1: State Transition

For the small training set

ATGATACTTgtccgagTATATAG

ATGTTTTgtggcagAAAGA

ATGAATgtcgcgagTTTTATAG



	Exon	Intron
Exon	0,59	0,049
Intron	0.049	0.312

Emission Probabilities

	A	T	C	G
EXON	0.4	0.42	0.041	0.14
INTRON	0.123	0.123	0.213	0.541

Testing the set	gives PATH :	OUTPUT file
atgtaagtggccagttaatga	EEEEEEEEEEEEEEEEEEEE	ATgtaagtggccagTTAATGA
atgaaaagtggggccagtaatga	EEEEEEEEEEEEEEEEEEEE	ATGAAAAgtggggccagTAATGA
atggtcagtag	EEEEEEEE	ATGgtcagTAG

E denotes Exon state, and I to Intron state. The codes are here



Q2:

$$1) \text{ATG} = 0.7 * 0.1 * 0.7 * 0.2 * 0.7 * 0.4 * 0.9$$

Begin, M₁, M₂, M₃, end=4.3218E-3

$$2) \text{ATG_} = 0.1 * 0.25 * 1 * 0.5 * 0.1 * 0.25 * 1 * 0.1 * 0.2 * 1$$

Begin, I₁, M₁, I₂, M₂, D₃, end=

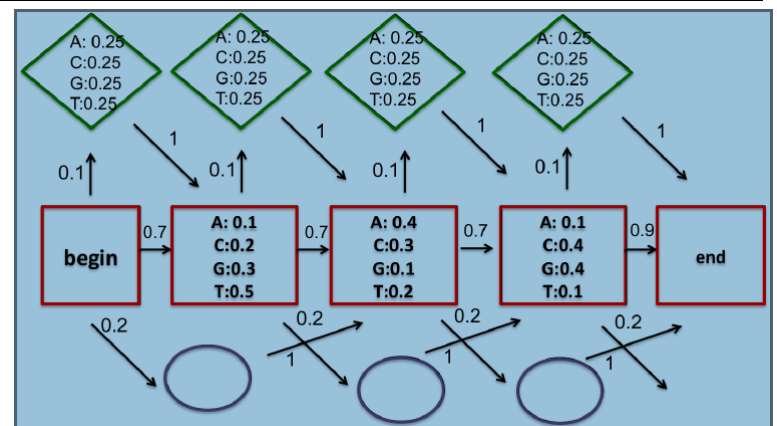
$$3) \text{_ATG_} = 0.2 * 1 * 0.4 * 0.1 * .25 * 1 * 0.4 * 0.9$$

Begin, D₁, M₂, I₃, M₃, end,

$$4) \text{AT_G} = 0.1 * 0.25 * 1 * 0.5 * 0.2 * 1 * 0.4 * 0.9$$

Begin, I₁, M₁, D₂, M₃, end

$$5) \text{A_TG} = 0.7 * 0.1 * 0.2 * 1 * 0.1 * 0.1 * 0.25 * 1$$

Begin, M₁, D₂, M₃, I₄, end.

Q3:

Motifs: GGAGG, GGAG or GAGG with probability 72%, 16% and 12%, comes after 4th or 6th position of a codon !

Here I made 10 upstream start of codons respecting the probabilities of the motifs:

CODON1: _____ GGAGG

CODON2: _____ GGAGG __

CODON3: _____ GGAGG _

CODON4: _____ GGAGG

CODON5: _____ GGAGG _

CODON6: _____ GGAGG

CODON7: _____ GGAGG _

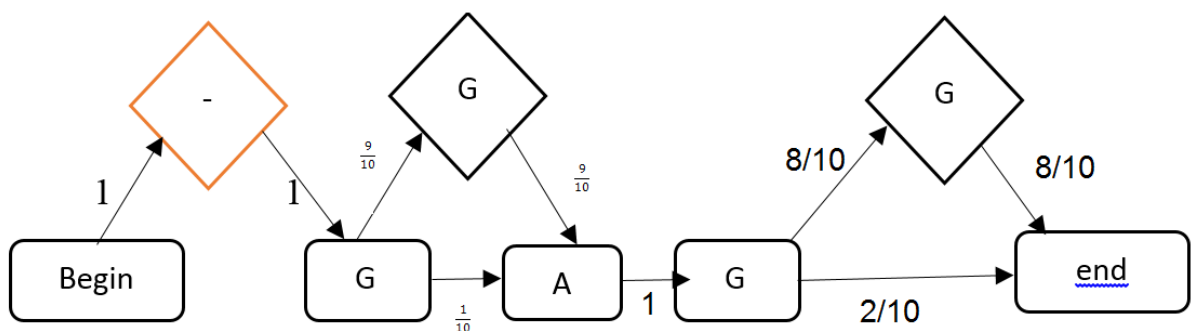
CODON8: _____ GGAG __

CODON9: _____ GGAG ____

CODON10: _____ GAGG _

The probabilities won't be accurate unless if we fix the conservative columns, if we consider that alignment we can do this: since they are randomly generated, but the emission observations are almost constant (G,A) starts with G and ends with G and has an A, Gs are duplicate 72% in both sides and 16% in first and 12% in the last side so we can have only 3 columns that is one for matching As and other two for matching Gs, and calculate the looping probabilities for insertions and deletion. We have number of max length of string is 11.

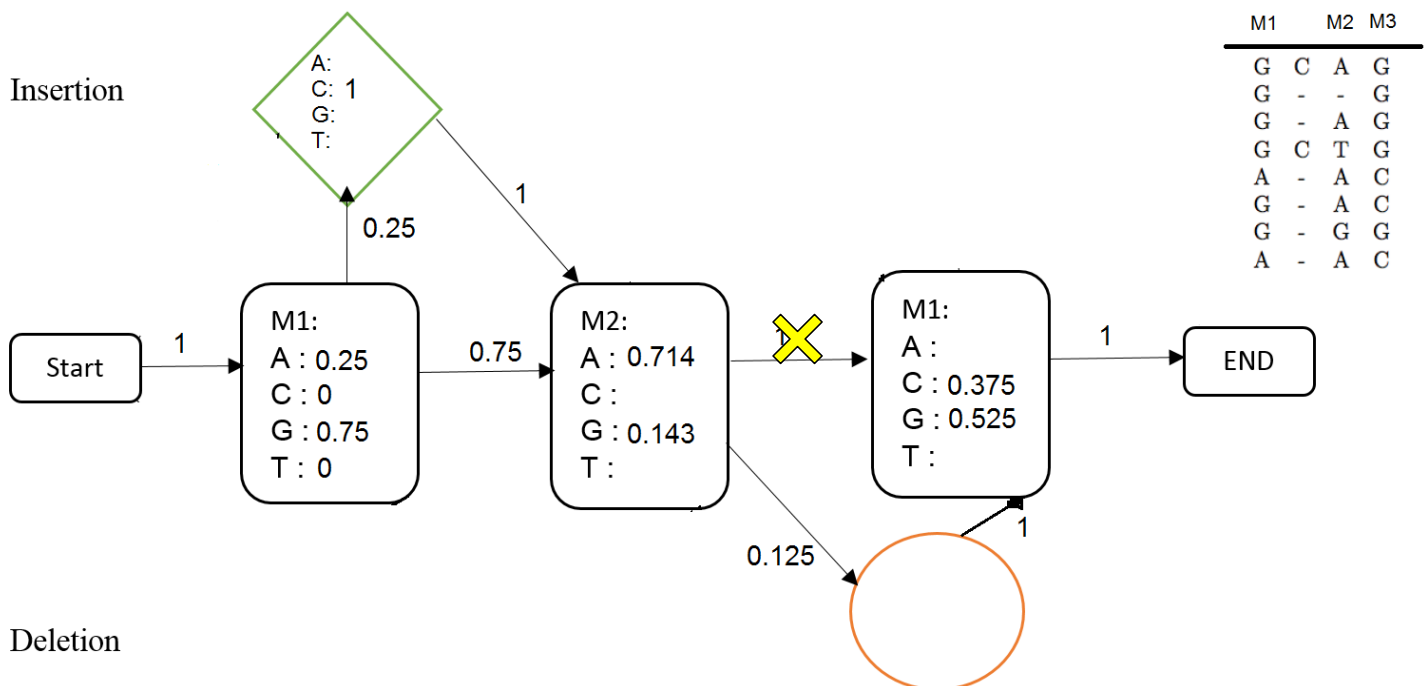
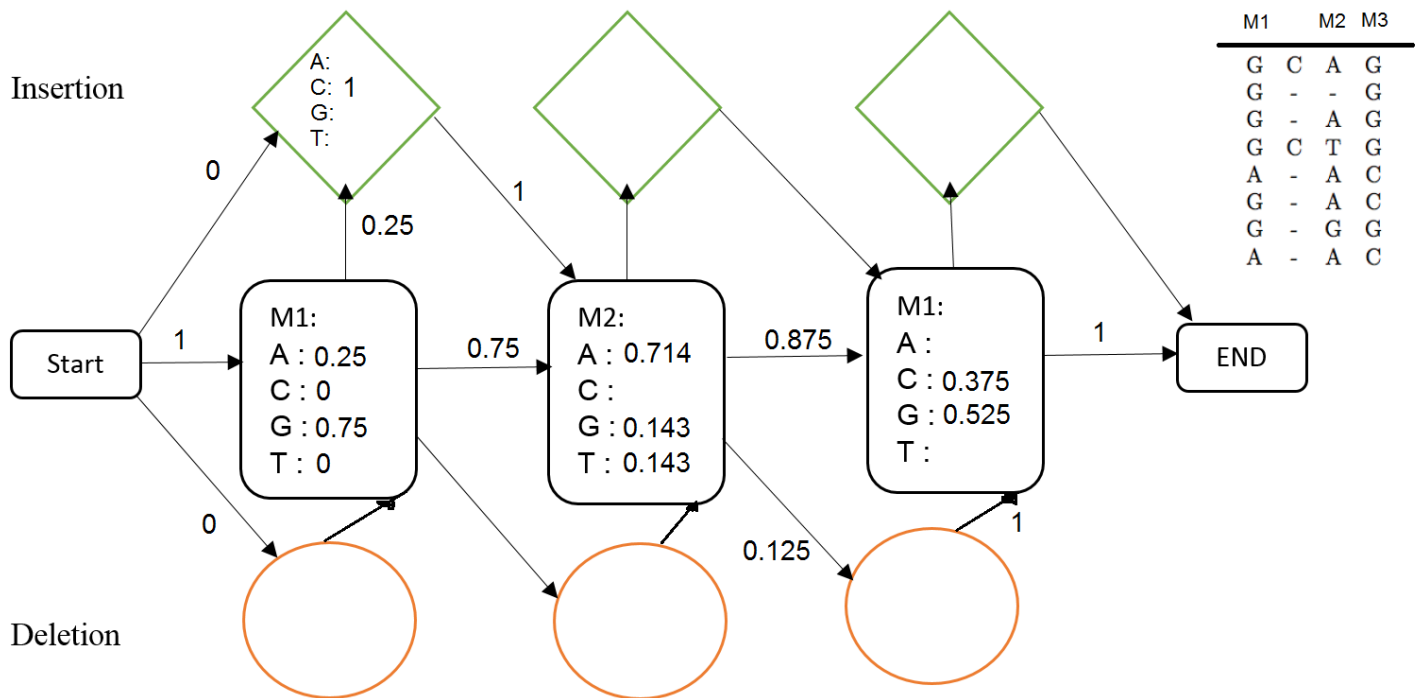
insertion



The orange insertion is a different HMM that need to be calculated with different emissions its length between 4 and 6



Q4: choosing 3 column out of 4 as most conservative columns, and one as insertion and one deletion occur in the second sequence in the transition from the second to the third column. (empty entries correspond to 0 probability)



B : since I have only 3 column and the first insertion doesn't loop that is it produce only once (c) , the sequence GCCAG can't be generated with the above model.