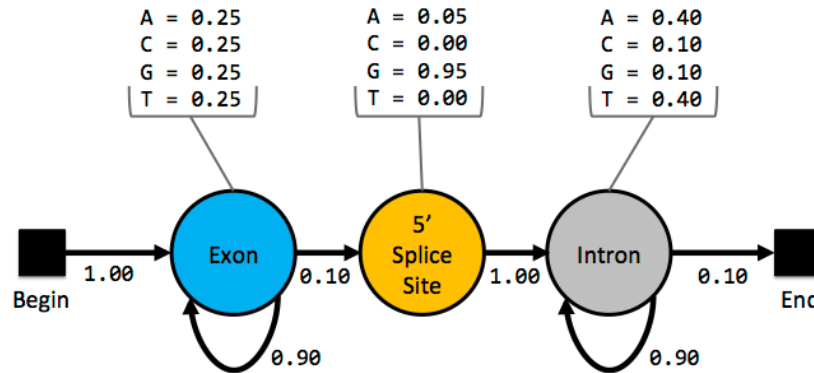


# Hidden Markov Model Basics

Written by Zane Goodwin and adapted from work by Anton E. Weisstein, Truman State University, 2013

1. Consider the following simple state machine:



a. Using the state machine above, manually calculate the probability of each of the following state paths:

Sequence:		A	C	G	C	G	
Path 1:	Begin	Exon	Exon	Exon	5'SS	Intron	End
Path 2:	Begin	Exon	Exon	5'SS	Intron	Intron	End
Path 3:	Begin	Exon	5'SS	Intron	Intron	Intron	End

(Note: 5'SS = 5' Splice Site)

- b. Based on your calculation of the probability of the three state paths above, which is the best state path and what is the *likelihood* of the state path?
- c. Each of the three state paths listed above has the 5' splice site at a different position in the sequence. Which position in the sequence has the highest probability of being the 5' splice site? How confident are you that this position is the correct choice for the 5' splice site?





- c. If your model predicts multiple positions in the sequence have the same likelihood to be a splice site, how could you use RNA-Seq data to identify the best splice site candidate?
- 
- 5. **QUESTION FOR THOUGHT:** Many gene predictors use a collection of known genes (i.e. training set) to estimate the transition and emission probabilities in a HMM. For example, one can estimate the transition probabilities for the exon and intron states using the length distributions of exons and introns in the training set.
    - a. If the training set used to train a gene predictor contains many short genes and a few long genes, would you expect the HMM to predict more long genes or more short genes? Why?

*(Note: Because the Excel workbook HMM\_intron.xls models only part of a single gene, you cannot use the workbook to address this question.)*

Last Update: 07/11/2013