

ChE/BE 163 Fall 2020
Problem Set #1
 Due 1pm Thursday, October 22

Problem 1 (Equilibrium analysis of a complex of interacting nucleic acid strands, 15 pts). Consider the hybridization chain reaction (HCR) mechanism of Figure 1 (*Proc Natl Acad Sci USA* **101**, 15275–15278, 2004), in which metastable DNA hairpins H1 and H2 polymerize upon exposure to initiator sequence I1. Note that each hairpin is intended to have a 6-nt toehold, an 18-bp stem, and a 6-nt loop.

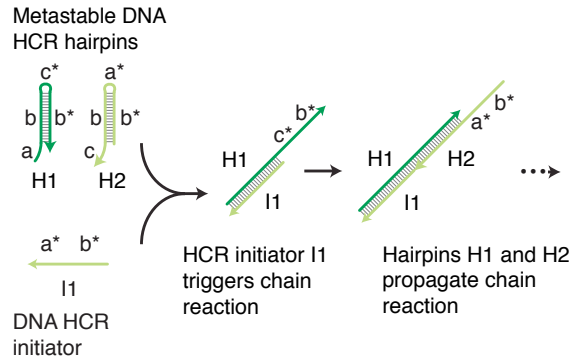


Figure 1: Conditional self-assembly via hybridization chain reaction (HCR).

Consider two different sequence designs for HCR that are intended to detect initiator:

I1: CCACACCACAACAACCACATCTCG

HCR Design 1

H1: CGAGATGTGGTTGTTGTGGTGTGGATACAACCACACCACAACAACCAC

H2: CCACACCACAACAACCACATCTCGGTGGTTGTTGTGGTGTGGTTGTAT

HCR Design 2

H1: CGAGATGTGGTTGTTGTGGTGTGGGTGGTCCACACCACAACAACCAC

H2: CCACACCACAACAACCACATCTCGGTGGTTGTTGTGGTGTGGAACCAC

Use the Utilities page of the NUPACK web application (nupack.org) to analyze these two different designs at 23 °C by examining the equilibrium base-pairing properties of the reactant and product complexes (each complex comprising one or more interacting strands) for two elementary steps in the HCR reaction pathway:

Step 1: I1 + H1 → I1·H1

Step 2: I1·H1 + H2 → I1·H1·H2

Use MFE structures and equilibrium base-pairing probabilities to explain which design you prefer.

Problem 2 (Ensemble size via dynamic programming, 20 pts).

In this problem, you will adapt the dynamic programming algorithm for computing the partition function of a single strand over the ensemble of unpsuedoknotted secondary structures, Γ , to the simpler problem of computing the number of possible unpsuedoknotted secondary structures, $|\Gamma|$. The primary reference for this problem is [Dirks and Pierce, *J. Comp. Chem.*, **24**, 1664 \(2003\)](#), with an emphasis on the $O(N^4)$ Algorithm (equations [7] and [8], Figures 3, 4 and 5, and the pseudocode in Figure 6).

- a) (3 pts) The partition function algorithm enables you to efficiently calculate

$$Q(\phi) = \sum_{s \in \Gamma} e^{-\Delta G(\phi, s)/k_B T}$$

for a strand with sequence ϕ . Here, we wish to calculate the size of the ensemble

$$|\Gamma(\phi)| = \sum_{s \in \Gamma(\phi)} 1 \tag{1}$$

where each structure $s \in \Gamma(\phi)$ contains only Watson-Crick pairs (i.e., no wobble pairs) and we note that two bases cannot pair if $j - i < 4$ (due to steric constraints). How can you make a simple change to the energy $\Delta G(\phi, s)$ to use the partition function algorithm to calculate $|\Gamma(\phi)|$?

- b) (10 pts) Consider the DNA sequence **CGAGATACCTCGATCACGCG**. Write a Python script to execute the dynamic program and calculate $|\Gamma(\phi)|$. Include your script and show the final state of the matrices Q , Q^b and Q^m . Hint: *What is the value of $Q_{i,j}^b$ if i and j cannot base-pair (either because they are not Watson-Crick complements or because they are sterically prevented from pairing)?* Hint: *Do not use NUPACK's `ensemble_size` command to check your script. It will give you different answers because it includes wobble pairs. You can test your program with the sequence **ATAGTTTTCTCGAAAACGAT**, which has 2349 unpsuedoknotted secondary structures.*
- c) (7 pts) Generate 5 random sequences for each length $N \in \{10, 20, 30, 40, 50, 60, 70, 80, 90, 100\}$ nt (select the sequence of each nucleotide from a uniform distribution over $\{\mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}\}$). Use your script to calculate $|\Gamma(\phi)|$ for each sequence; display $|\Gamma(\phi)|$ vs N as a scatter plot to observe the blowup in ensemble size as N increases. Suppose $|\Gamma(\phi)| = ae^{bN}$ so that taking log of both sides yields $\log |\Gamma(\phi)| = \log a + bN$. Display $\log |\Gamma(\phi)|$ vs N as a scatter plot and use a built-in least squares fitting procedure (e.g., `polyfit` in NumPy) to estimate a and b ; display the line of best fit in your plot.

Problem 3 (Objective function and nucleic acid sequence design, 50 pts).

In this problem you will design sequences intended to form an “RNA stick figure” at equilibrium (Figure 2).

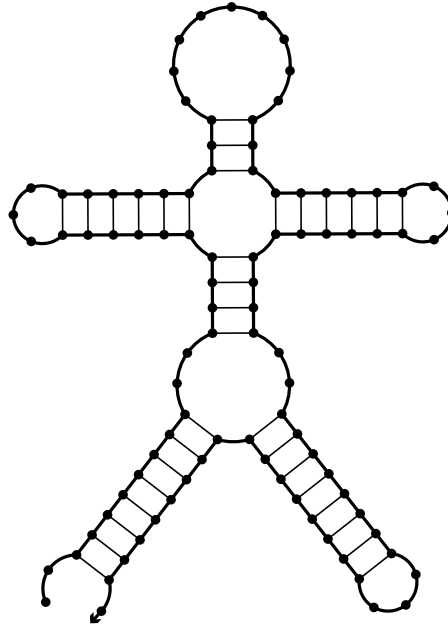


Figure 2: RNA stick figure.

We specify the secondary structure using *dot-parens* notation. Each unpaired base is represented by a dot and each base pair by matching parentheses. If a given base has a dot, it is unpaired. If it has an open parenthesis, it is paired with the parenthesis that closes it (e.g., “.(((...)))” represents a secondary structure where the second base is paired with the last, the third with the ninth, and the fourth with the eighth). The stick figure is denoted below:

..(((((((..((((((((((...))))))((.....))((((((...)))))))))..(((((((...)))))))))

For unspseudoknotted structures, the dot-parens representation of the secondary structure has the same information as the drawing above, the corresponding polymer graph (not shown), and the secondary structure matrix S . The dot-parens representation is a convenient way to represent the structure in a computer program because it is a simple string, and you will use it in your script.

In doing the design calculations, you will make use of Utilities commands in the NUPACK Python module (but not the Design commands). The [NUPACK User Guide](#) provides a useful online reference describing the NUPACK python module. Additionally, you will need to use other Python-based software, such as NumPy. You can use a computing resource on the NUPACK server cluster that we have set up for you that already has NUPACK 4.0 installed, along with Python wrappers to aid in its use, and other Python packages you will need. For instructions on how to launch and use a NUPACK server, see the Computing section of the [Course Info page](#).

Unless specifically instructed to use the NUPACK web application, you will write your own Python script that uses commands from the NUPACK Python module (again, with the obvious exception of

