



Finding Signals from the
Candidates of Interaction Sites
for Photosystem I Reduction

SHONDALYN A. SMITH

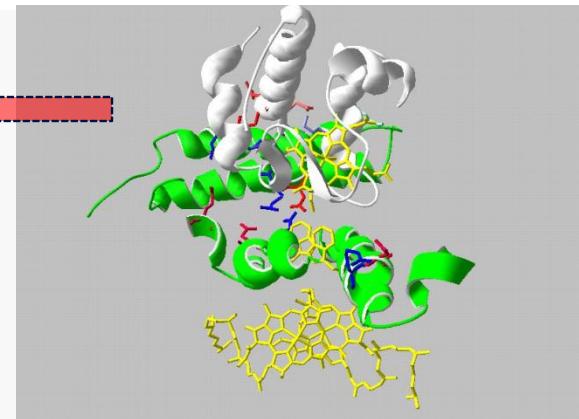
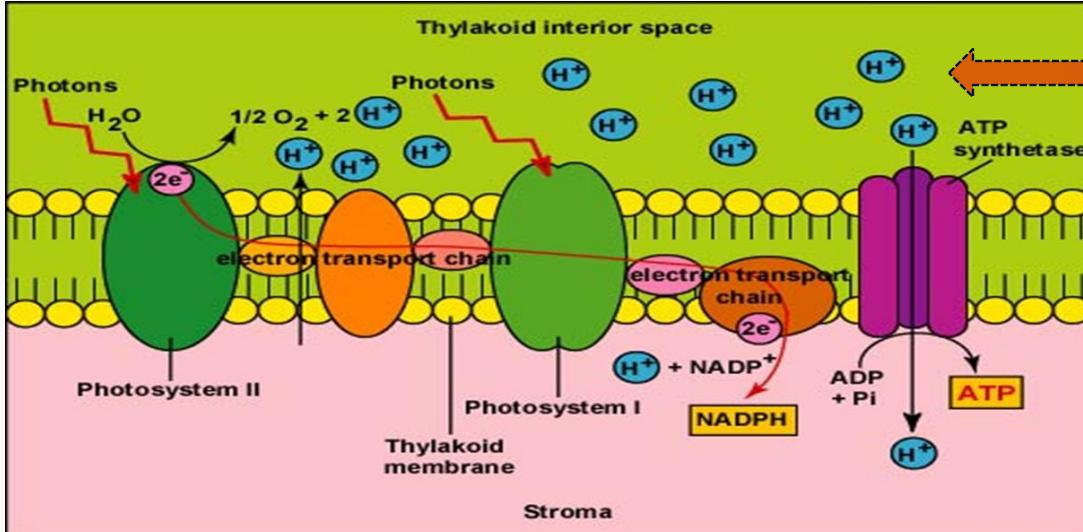
STEVYN SAWYER

DR. WEI CHEN, ADVISOR

DEPARTMENT OF COMPUTER SCIENCE

TENNESSEE STATE UNIVERSITY

Research Background



Artificially adding c₆ and psaF in Photosystem I speed up the process and produce more hydrogen

- Natural system produces hydrogen but in small amounts.
- Artificially adding c₆ and psaF proteins produces more amounts quickly.
- Previous senior project group predicted the candidates for the interaction.

In order to make the prediction more precise, we want to find more common signals from these candidates !!!

Problem Statement



Find common signal (motif) from the candidates of interaction sites in psaF and c6 protein families.

GGGCT	A T c C A g C T	GGGT CGT CAC AT TCC C C T T C G A
TGAGGGTGCCCAATAA	g g G C A A C T	CCAA AGCG GGACA
	AT G g A t C T	GATGCC GTT GAC GAC CT AA AT CA AC GG
	GG A a G C A A C c	CCAGGAGCGC CTT GCTGG TT C T ACC
TTTTCTAAAAAGATTATAATGTCGGTCC	t T G g A A C T	
GCTGTACACACTGGATCATGCTGC	AT G C c A t T	TTCA
CATGATCTT TG	AT G g c A C T	TGGATGAGGGAATGAT

Challenge – Then same motif in different proteins can be different after have had mutation. Finding motif is computationally very expensive.

It is NP(No-deterministic Polynomial(-hard))!!

Goal and Objectives



Goal:

Find further information (motif) from the identified candidates of interaction sites in c6 and psaF protein families.

Research Objectives:

- 1) Design new algorithm for finding motif effectively and efficiently.
- 2) Implement and evaluate the algorithm and analyze the results with existing algorithms.

Approach



- Investigate known algorithms.
- Design new algorithm.
- Compare the proposed algorithm with existing algorithms.

Functional Requirements



1. Protein pairs of c6 and psaF family from the same organism need to be provided.
2. The candidates of interaction sites of the given protein pairs need to be provided.
3. Identify and compare the existing software/tools/databases for motif/signal search of proteins for c6 and psaF.
4. Efficient algorithms and software that fits motifs/signals search from the candidates of interaction sites of c6 and psaF family are required.

Non-functional Requirements



1. Utilize standard database to find signals from the candidates of the interaction sites.
2. Compare the results from standard database and propose algorithm.
3. Algorithm implementation will utilize C# in Visual Studios.

System Design



- Dataset preparation to find signal from candidates of interaction sites.
- Utilize existing database/tools/software.
- Algorithm design.

System Design: Dataset preparation



- *Example 1: one pair of psaF and c6 proteins.*
- **psaF:**
 - DIAGLTPCSESKAYAKLEKKELKTLEKRLKQYEADSAPAVALKAT
MERTKARFANYAKAGLLCGNDGLPHLIADPGLALKYGHAGEVFI
PTFGFLYVAGYIGYVGRQYLIAVKGEEKPTDKEIIDVPLATKLAW
QGAGWPLAAVQELQRGTLLEKEENITVSPR
- **c6:**
 - ADLALGAQVFNGNCAACHMGGRNSVMPEKTLDKAALEQYLDGG
FKVESIIYQVENKGAMPAWADRLSEEEIQAVAЕYVFKQATDAA
WKY

Algorithm Design

Calculate the score of a motif

(All algorithms will use this method)

Example: Seven 36- nucleotide DNA sequences with a hid pattern TGCAACT mutated in two positions.

GGGCTAtcCAgCTGGGTCGTACATTCCCCTTCGA
TGAGGGTGCCAATAAggGCAACTCCAAAGCGGACA
ATGgAtCTGATGCCGTTGACGACCTAAATCAACGG
GGAaGCAACcCCAGGAGGCCCTTGCTGGTTCTACC
TTTTCTAAAAAGATTATAATGTCGGTCCtTGgAACT
GCTGTACACACTGGATCATGCTGCATGCcAtTTCA
CATGATCTTTGATGgcACTTGGATGAGGGAATGAT

GGGCT	<u>ATcCAgCT</u>	GGGTCGTACATTCCCCTTCGA
TGAGGGTGCCAATAA	<u>g g G C A A C T</u>	CCAAAGCGGACA
	<u>ATGgA t C T</u>	GATGCCGTTGACGACCTAAATCAACGG
GG	<u>A a G C A A C c</u>	CCAGGAGGCCCTTGCTGGTTCTACC
TTTTCTAAAAAGATTATAATGTCGGTCC	<u>t T G g A A C T</u>	
GCTGTACACACTGGATCATGCTGC	<u>ATG C c A t T</u>	TTCA
CATGATCTTTG	<u>ATG g c A C T</u>	TGGATGAGGGAATGAT

A	5	1	0	5	5	0
T	1	5	0	0	1	6
G	1	1	6	3	0	1
C	0	0	1	4	2	0

Consensus **ATGCAACT**

position = (6,17,1,3,29,25,13)
Score = 5+5+6+4+5+5+6+6 = 42

Existing Algorithms for Finding Motifs

BruteForce

BruteForce(DNA, t, n, l)

```
bestScore := 0;  
for i1 := 1 to n-l+1  
    for i2 := 1 to n-l+1  
        .....  
        for it := 1 to n-l+1  
            S = (i1, i2, ..., it)  
            if (Score(S DNA) > bestScore)  
                bestScore := Score(S, DNA)  
                bestMotifPosition = S  
return bestScore & bestMotifPosition;
```

DNA: DNA sequences

t: number of DNA sequences

n: length of DNA sequences

l: length of the motif

Time Complexity : $O(n^t lt)$

Example: t=7, n=36, l=8,

GGGCT	A T c C A g C T	GGGT CGT CAC AT TCC CTT TCGA
TGAGGGTGC CCAATAA	g g G C A A C T	CCAA AGCGGACA
	AT G g A t C T	GATGCC GTTTGACGACCTAAATCAACGG
	GG A a G C A A C c	CCAGGAGCGCCTTTGCTGGTTCTACC
TTTTCTAAAAAGATTATAATGTCGGTCC	t T G g A A C T	
GCTGTACACACTGGATCATGCTGC	AT G C c A t T	TTCA
CATGATCTTTTG	AT G g c A C T	TGGATGAGGGAATGAT

Profile	A	51005500	
	T	15000116	$S = (6, 17, 1, 3, 29, 25, 13)$
	G	11630100	$\text{Score}(S) = 5+5+6+4+5+5+6+6 = 42$
	C	00142061	

Running Time: 7 hours

Existing Algorithms for Finding Motifs

Greedy Algorithm

Greedy-MotifFinding(DNA, t, n, l)

bestMotif := (1,1,...,1);

s := (1,1,...,1)

for s1 := 1 to $n-l+1$

 for s2 := 1 to $n-l+1$

 S := (s1, s2, 1, ..., 1)

 if (Score(S, Seq) > bestScore)

 bestScore := Score(S, DNA);

 bestMotif Position:= S

for i := 3 to t

 for si := 1 to $n-l+1$

 S:= (s1, s2, ..., si, 1, ..., 1)

 if (Score(S, DNA) > bestScore)

 bestScore := Score(S, Seq);

 bestMotif Pos:= S;

return bestScore & bestMotifPosition

Time Complexity : $O(n^2tl + nt^2l)$

Local
optimal
solution for
first 2 strings

Local
optimal
solution for
3rd - tth
strings

DNA[0] = "GCTGTACACA**CTGGATCA**TGCTGCATGCCATTTC" (highlighted in yellow)

DNA[1] = "CATGATCTTTGATGGCACTTGGATGA**GGGAATGAT**" (highlighted in red)

DNA[2] = "A TGG A TCT**GATGCCGTTGACGACCTAAATCAACGG**" (highlighted in green)

DNA[3] = "GGAAGCAA**CCCCCAGGAGCGCCTTGCTGGTTCTACC**" (highlighted in red)

DNA[4] = "T TT TCTAA**AAAGATTATAATGTCGGTCCTTGGAACT**" (highlighted in red)

.....

DNA[13] = "TGAGGGTG**CCCAATAAGGGCAACTCCAAAGCGGACA**" (highlighted in red)



DNA[0] = "GCTGTACACA**CTGGATCA**TGCTGCATGCCATTTC" (highlighted in yellow)

DNA[1] = "CATGATCTTTGATGGCACTTGGATGA**GGGAATGAT**" (highlighted in red)

DNA[2] = "**ATGGATCTGATGCCGTTGACGACCTAAATCAACGG**" (highlighted in green)

DNA[3] = "GGAAGCAA**CCCCCAGGAGCGCCTTGCTGGTTCTACC**" (highlighted in red)

DNA[4] = "T TT TCTAA**AAAGATTATAATGTCGGTCCTTGGAACT**" (highlighted in red)

.....

DNA[13] = "TGAGGGTG**CCCAATAAGGGCAACTCCAAAGCGGACA**" (highlighted in red)



DNA[0] = "GCTGTACACA**CTGGATCA**TGCTGCATGCCATTTC" (highlighted in yellow)

DNA[1] = "CATGATCTTTGATGGCACTTGGATGA**GGGAATGAT**" (highlighted in red)

DNA[2] = "**ATGGATCTGATGCCGTTGACGACCTAAATCAACGG**" (highlighted in green)

DNA[3] = "**GGAAGCAA**CCCCCAGGAGCGCCTTGCTGGTTCTACC" (highlighted in red)

DNA[4] = "T TT TCTAA**AAAGATTATAATGTCGGTCCTTGGAACT**" (highlighted in red)

.....

DNA[13] = "TGAGGGTG**CCCAATAAGGGCAACTCCAAAGCGGACA**" (highlighted in red)



Running Time: instant

Proposed Algorithms for Finding Motifs

Improved Greedy Algorithm (using Divide-and-Conquer)



ImprovedMotifFinding(DNA[i..j], t, n, l)

```
if (j-i) < 4  
    return Greedy(DNA[i..j], t, n, l)  
else  
    k = (i+j-1)/2
```

x = **ImprovedMotifFinding(DNA[i..k], t, n, l)**

y = **ImprovedMotifFinding(DNA[k+1..j], t, n, l)**

```
if x.score > y.score  
    improve DNA[k+1..j] by the motifs in DNA[i..k]  
else  
    improve DNA[i..j] by the motifs in DNA[K_1..j]  
return bestScore and bestMotifPosition
```

Time Complexity :

$$\begin{aligned} T(n) &= 2T(n/2) + nt^2l/2 \quad \text{if } t > 4 \\ &= n^2tl \quad (\text{use greedy}) \quad \text{if } t \leq 4 \end{aligned}$$

$$T(n) = O(n^3tl)$$

If score of the motifs in first part is larger

CTGTACACACTGATCATGCTGCATGCCATTCA
CATGATCTTTGATGGCACTTGGATGAGGGAATGAT
ATGGATCTGATGCCGTTGACGACCTAAATCACCGG
GGAAGCAACCCCAGGAGCGCCTTGCTGGTTCTACC
TTTCTAAAAAGATTATAATGTCGGTCCTTGGAACT
TTTCTAAAAAGATTATAATGTCGGTCCTTGGAACT
GCTGTACACACTGGATCATGCTGCATGCCATTCA
CATGATCTTTGATGGCACTTGGATGAGGGAATGAT
GGGCTATCCAGCTGGTCGTACATTCCCCTTCGA
TGAGGGTGCCAATAAGGGCAACTCCAAAGCGGACA
TGGATCTGATGCCGTTGACGACCTAAATCACCGG
GGAAGCAACCCCAGGAGCGCCTTGCTGGTTCTACC
GGGCTATCCAGCTGGTCGTACATTCCCCTTCGA
TGAGGGTGCCAATAAGGGCAACTCCAAAGCGGACA



If score of the motifs in second part is larger

Test results and algorithm comparison

Dataset 1: Number of DNA sequences = 7,
Length of DNA sequence = 32, length of motif = 8



GGGCTATCCAGCTGGTCGTACATTCCCCTTCGA
TGAGGGTGCCCAATAAGGGCAACTCCAAAGCGGACA
ATGGATCTGATGCCGTTGACGACCTAAATCAACGG
Greedy/Improved-greedy
GGAAGCAACCCAGGAGCGCCTTGCTGGYTCTAAC
TTTTCTAAAAAGATTATAATGTCGGTC**CTTGGAACT**
GCTGTACACACTGGATCAT**GCTGCATGCCATTTC**A
CATGATCTTTGA**TGGCACTTGGATGAGGGAATGAT**

GGGCTATCCAGCTGGTCGTCA**CATTCCCCTTCGA**
TGAGGGTGCCCAATAAGGGCAACTCCAAAGCGGACA
ATGGATCTGATGCCGTTGACGACCTAAATCAACGG
GGAAGCAACCCAGGAGCGCCTTGCTGGTTCTACC
TTTTCTAAAAAGATTATAATGTC**GGTCCTTGGAACT**
GCTGTACACACTGGATCATGCTGCATGCCATTTC
CATGATCTTTGA**TGGCACTTGGATGAGGGAATGAT**



Brute Force

Algorithm	Score of Motif	Position of Motif	Run Time
Brute Force	42	5, 15, 0, 2, 27, 19, 12	7 hours
Greedy	37	14, 4, 2, 13, 23, 0, 0	4.2816 ms
Improved Greedy	37	14, 4, 2, 13, 23, 0, 0	7.9553 ms

Test results and algorithm comparison

Dataset 2: Number of DNA sequences = 14,
 Length of DNA sequence = 32, length of motif = 8

GCTGTACACACTGGATCATGCTGCATGCCATTTCAT
 CATGATCTTTGATGGCACTTGGATGAGGAATGAT
 ATGGATCTGATGCCGTTGACGACCTAAATCAACGG
 GGAAGCAACCCCAGGAGCGCCTTGCTGGTTCTACC
 TTTTCTAAAAAGATTATAATGTCGGTCTTGGAACT
 TTTTCTAAAAAGATTATAATGTCGGTCTTGGAACT
 GCTGTACACACTGGATCATGCTGCATGCCATTTCAT
 CATGATCTTTGATGGCACTTGGATGAGGAATGAT
 GGGCTATCCAGCTGGGTGTCACATTCCCCTTCGA
 TGAGGGTGCCAATAAGGGCAACTCCAAAGCGGACA
 ATGGATCTGATGCCGTTGACGACCTAAATCAACGG
 GGAAGCAACCCCAGGAGCGCCTTGCTGGTTCTACC
 GGGCTATCCAGCTGGGTGTCACATTCCCCTTCGA
 TGAGGGTGCCAATAAGGGCAACTCCAAAGCGGACA


 Improved-
 greedy

Algorithm	Score of Motif	Position of Motif	Run Time
Brute Force			Years
Greedy	68	10, 27, 0, 11, 8, 8, 10, 26, 0, 2, 0, 2, 1, 2	7.0717 ms
Improved Greedy	76	10, 26, 0, 11, 28, 28, 10, 26, 11, 2, 0, 11, 11, 2	21.0426 ms

Future Work



- Continue to find better algorithms.
- Apply proposed algorithm to the candidates of the interaction sites in c6 and psaF protein families.

Acknowledgements



Dr. Wei Chen

Dr. Ali Sekmen