#### A Compact Mathematical Programming Formulation for DNA Motif Finding

#### **Carl Kingsford**

Center for Bioinformatics & Computational Biology University of Maryland, College Park

#### Elena Zaslavsky and Mona Singh

Computer Science Dept. and Lewis-Sigler Institute for Integrative Genomics, Princeton University





# DNA -> mRNA -> Protein



Finding transcription factor binding sites can tell us about the cell's regulatory network.

## Motif Finding

Transcription factor

- 1. ttgccacaaaataatccgccttcgcaaattgaccTACCTCAATAGCGGTAgaaaaacgcaccactgcctgacag
- 2. gtaagtacctgaaagttacggtctgcgaacgctattccacTGCTCCTTTATAGGTAcaacagtatagtctgatgga
- 3. ccacacggcaaataaggagTAACTCTTTCCGGGTAtgggtatacttcagccaatagccgagaatactgccattccag
- 4. ccatacccggaaagagttactccttatttgccgtgtggttagtcgcttTACATCGGTAAGGGTAggggattttacagca
- 5. aaactattaagatttttatgcagatgggtattaaggaGTATTCCCCATGGGTAacatattaatggctctta
- 6. ttacagtctgttatgtggtggctgttaaTTATCCTAAAGGGGTAtcttaggaatttactt

# Given p sequences, find the most mutually similar length-k subsequences, one from each sequence:

$$\underset{s_1, \dots, s_p}{\operatorname{argmin}} \sum_{i < j} \operatorname{dist}(s_i, s_j)$$

 $dist(s_i, s_j) = Hamming distance between s_i and s_j.$ Hundreds of papers, many formulations (Tompa05)

### **Graph Formulation**

- For p sequences, complete p-partite graph.
- Node for each sliding window of length k.
- Weight on edge (u,v) = dist(u,v) = # of differences between subsequences u and v.

gctgttaaTTATCCGGGGGTAtcttagga

Goal: Choose one node from each part to minimize weight of the induced subgraph.



### Hardness

- NP-hard (Wang+94, Akutsu+00).
- General distance measure ⇒ inapproximatible within O(|∨|) = # nodes in the graph (Chazelle+04).
- Triangle inequality  $\Rightarrow$  constant-factor approximation (Bafna+97).
- Interested in provably optimal solutions.

#### Integer Programming Formulation

- Binary variables x<sub>u</sub> for each node
- Binary variables x<sub>uv</sub> for each edge



$$\begin{array}{ll} \text{Minimize} & \sum_{\{u,v\}\in E} w_{uv} x_{uv} \\ \text{I.} & \sum_{u\in V_j} x_u = 1 & \text{for every part j} \\ \text{2.} & \sum_{u\in V_j} x_{uv} = x_v & \text{for every part j, node v} \end{array}$$
(IP1)

(Zaslavsky & Singh, 05)

#### Problem: Instances are Huge

#### 205,146 to 16,637,889 variables

(For the reasonably sized instances in our test set)

Goal: Can we exploit features of the instances to reduce sizes (or get better formulation)

## New Formulation

Ø Only small number of possible edge weights, ≤ window length.

Don't care which edge is chosen, only that correct cost is paid.

→ Merge edges that have same cost; vastly reduce # of variables.





**Binary** variables:

 $\oslash$  X<sub>u</sub> on the nodes.

•  $Y_{ujc}$  for each node u, position j not containing u, and weight c. Idea:  $Y_{ujc} = 1$  if we choose an edge of weight c between node u and position j.

## Reduced # Variables

• N = # nodes per position

• k = length of window = # of possible edge weights - 1



N<sup>2</sup> edge variables

Need to ensure compatible edge-groups are chosen. 2N(k+1) edge-group variables k << N



 $X_u, Y_{ujc} \in \{0,1\}$ 

#### IP1 vs. IP2

Optimum IP2 = optimum IP1.

IP2 has a factor of
 O(k/N) fewer variables, and
 O(k) more constraints than IP1.

k = motif length N = sequnce length

k is fixed by transcription factor geometry; N will grow as longer sequences are considered.

▲ LP relaxation of IP2 is weaker than that of IP1.
 ➡ Add constraints to make LP2 as tight as LP1.

# Constraining Y Variables



Neighbors of a set of Y vars:  $N(ujc) = \{(v,i,c) : cost(u,v) = c\}$  $N(Q) = \cup N(ujc)$ 

Neighbors are compatible



# Tightening LP2

**Thm.** If all constraints of the form (\*) are included, then the resulting LP has the same optimum as LP1.

Thm (separation algorithm). There is a polytime algorithm to find a violated constraint of the form (\*).

⇒ Despite exponential # of constraints, new LP can be solved in polytime.

# Testing Data Set

- 39 families of E. coli transcription factors from (Robison+98).
- Binding sites found via various experimental techniques.
- 3 20 sequences of length ≥ 300 in each family.
  Length of motif ranges: 11 to 48
  Up to 5,960 nodes in resulting graphs.

#### Formulation is Accurate

Test on real data.

Finds biologically relevant solutions.

Compare performance to state-of-the-art probabilistic approach based on Gibbs sampling.

## Accuracy Comparisons

 Real motif

 False positives (FP)

 Fillennn Million

 Prediction Million positives (TP)

Nucleotide level measure of accuracy (Pevzner&Sze,00):

nPC = nTP / (nTP + nFN + nFP)

Compares favorably with Gibbs sampling strategy: plot our nPC minus Gibbs nPC.



# Timing Comparison

Speed up for LP2 + cutting planes to reach same objective function as LP1.

10x faster is not uncommon.



(green indicates LP1 did not finish in ≤ 5 hours)

#### Conclusion

- Able to find provably optimal solutions to real transcription factor binding site discovery problems.
- Large speed up by using bounded number of objective function costs.
- Finds as good or better motifs than other motif-finding approaches.
- Open: how to use triangle inequality and overlapping windows to further shrink IP.

### Acknowledgments

- Co-authors: Mona Singh, Elena Zaslavsky
- Additional funds provided by:
  - Center for Bioinformatics and Computational Biology, University of Maryland, College Park <u>http://www.cbcb.umd.edu</u>
  - Program in Integrative Information, Computer and Application Sciences (PICASso), Princeton University <u>http://www.cs.princeton.edu/picasso/</u>

#### Times & Sizes

