

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

Segment Match Refinement and Applications

Aaron L. Halpern,
Daniel H. Huson*
& Knut Reinert

CELERA

* Now at: Center for Bioinformatics Tübingen

Daniel Huson, 2002

Copyright (c) 2008 Daniel Huson.

Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.2 or any later version published by the Free Software Foundation; with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts. A copy of the license can be found at <http://www.gnu.org/copyleft/fdl.html>

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

Overview

- Motivation: comparison of HGP and Celera assemblies of human
- Problem: one-to-one matching of sequences
- Solution 1: A greedy heuristic
- Solution 2: Optimal refinement algorithm
- Other applications
- Comparison of solutions

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

Two assemblies of human

Daniel Huson, 2002

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

Comparison of two assemblies

- What do they have in common?
 - Compute one-to-one matching of sequence
- How do they differ?
 - Determine how much of the remaining difference is unique,
 - unique-repetitive or
 - unique-under-collapsed sequence

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

Comparison of HGP and Celera assemblies

Daniel Huson, 2002

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

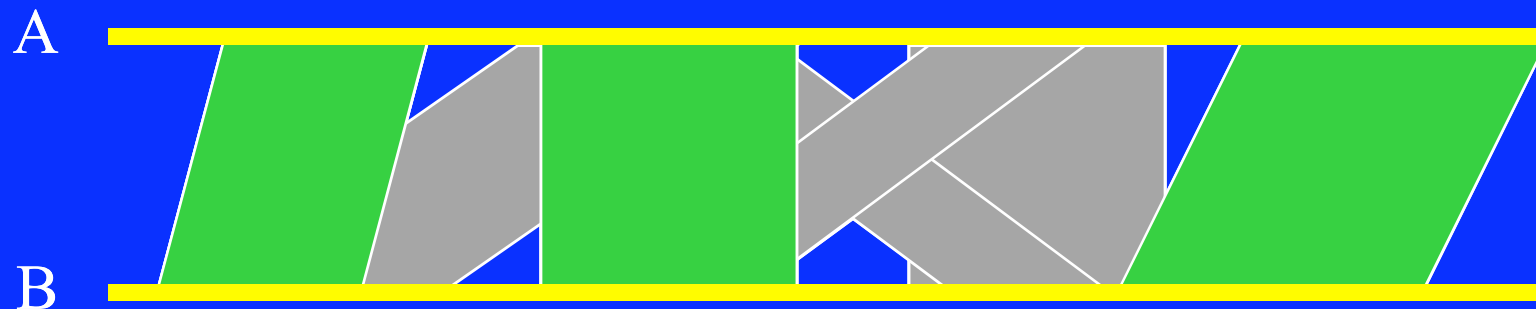
The Matched Sequence Problem

Given sequences A and B and a set Σ of segment matches between them. The **Matched Sequence Problem** is to compute a set of non-intersecting matches Σ' that are all submatches of Σ , such that the amount of sequence covered by the matched segments is maximized

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

Greedy one-to-one matching

- Given a set Σ of segment matches between two sequences A and B:

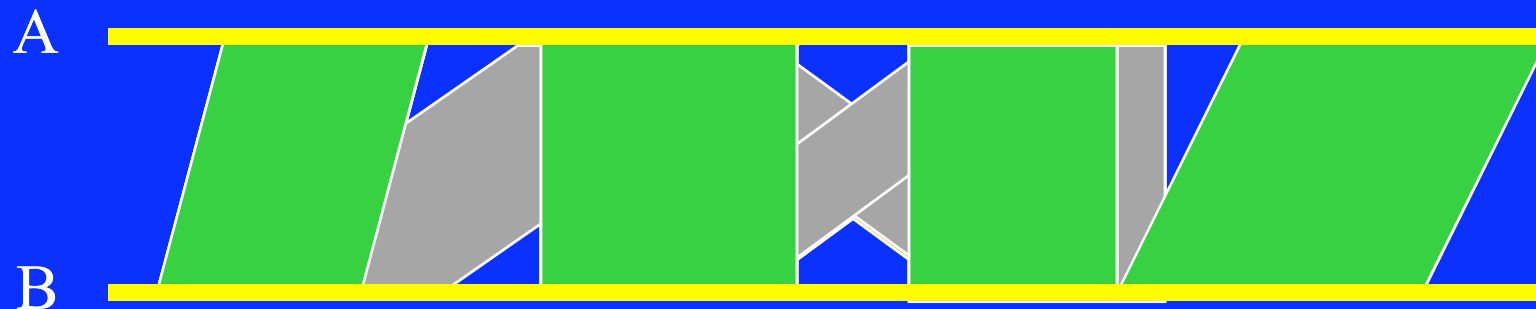


- A one-to-one matching of sequence can be obtained by greedily selecting a set of disjoint matches using a priority queue

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

Greedy matching with trimming

- Given a set Σ of segment matches between two sequences A and B:

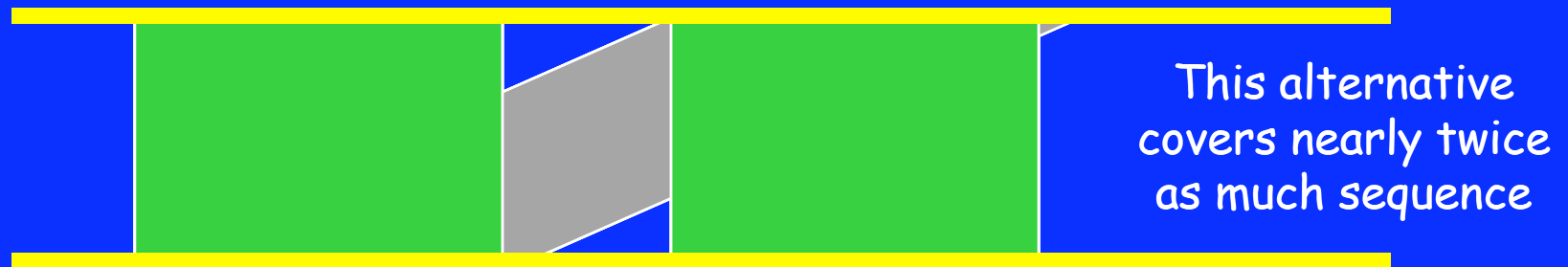


- Trim matches that partially overlap with already selected ones and re-insert into priority queue

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

How well does greedy selection with trimming work?

- Goal: obtain a set of non-overlapping matches that cover as much of both sequences as possible
- Greedy selection runs fast and obtains results that seem reasonable...
- BUT it could be missing up to half of the attainable coverage:



Daniel Huson, 2002

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

How about an optimal matching?

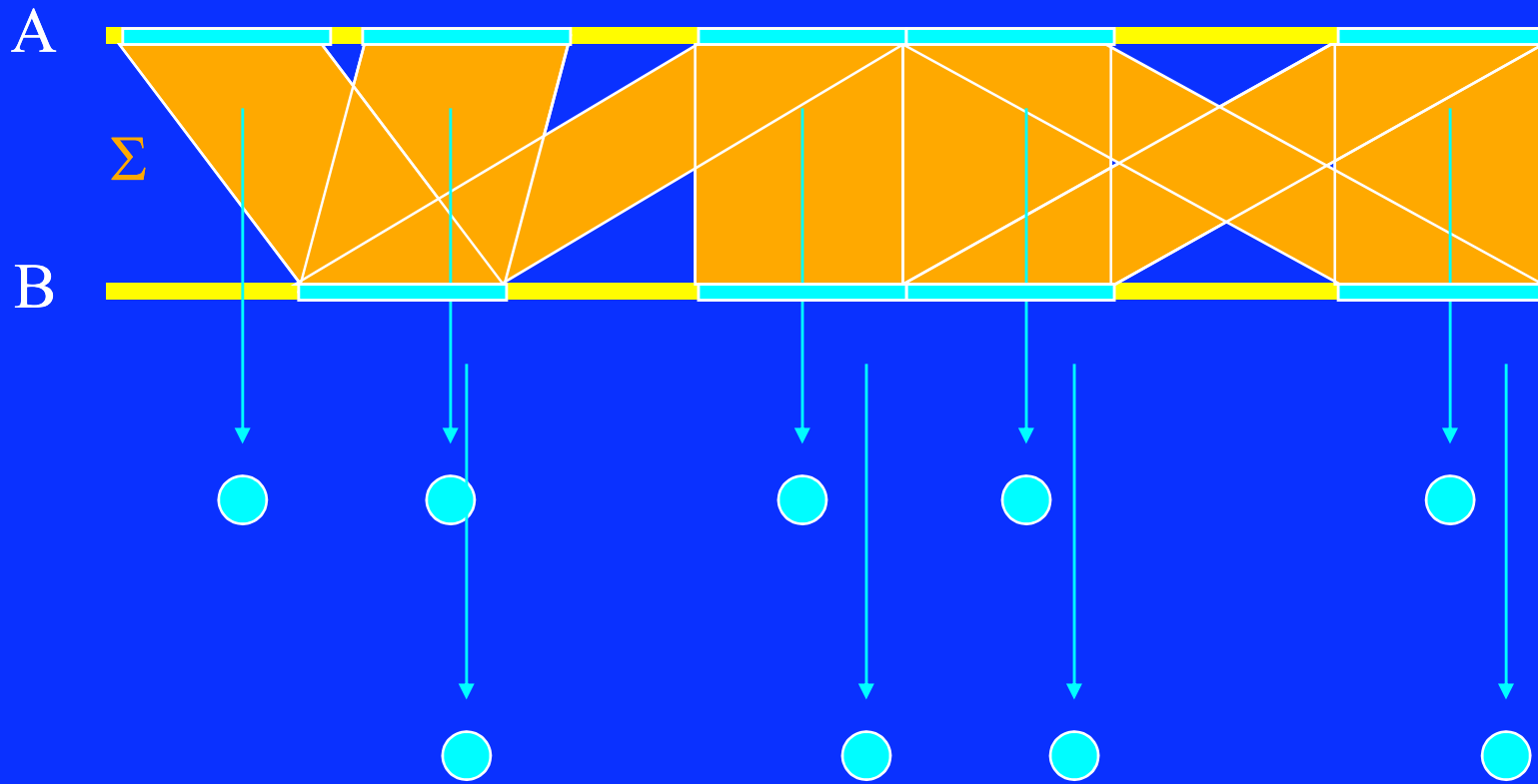
- Can an optimal set of non-overlapping matches be computed efficiently?
- How much better is an optimal solution than the one obtained by the greedy+trimming approach?

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

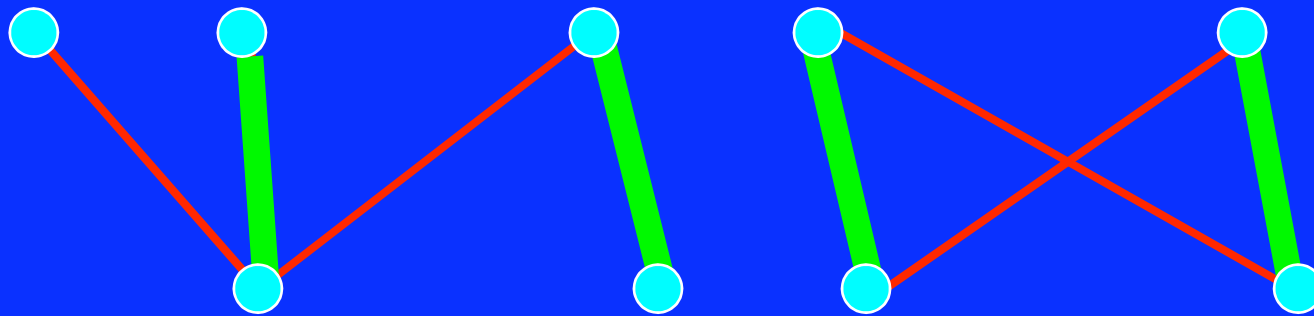
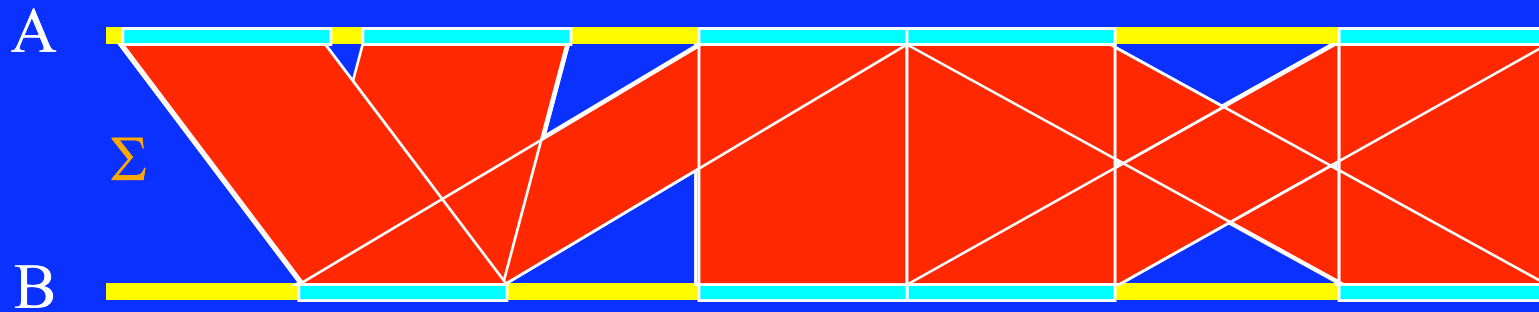
Sequence matching vs. graph matchings

- Given two sequences A and B and a set Σ of segment matches
- Assume that Σ is *resolved*, i.e. that the projections of any two matches onto each sequence are either disjoint or identical
- Then, the Sequence Matching Problem can be reduced to a weighted matching problem for a bipartite graph

The match graph



The match graph



Heaviest matching in bipartite graph gives optimal sequence matching

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

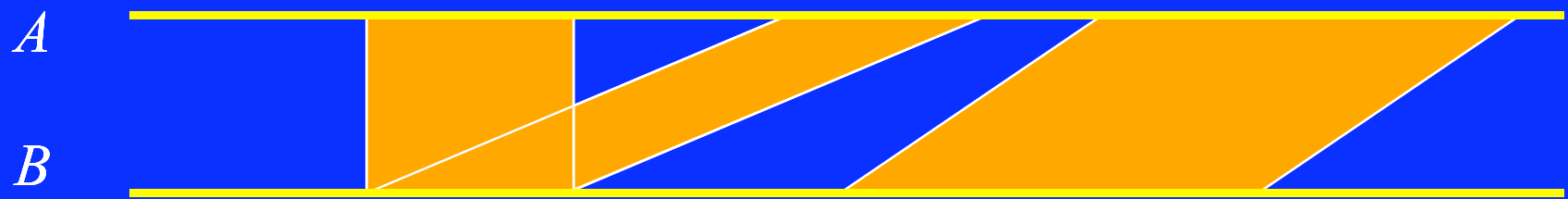
The resolved refinement

- Given a set Σ of matches between A and B.
- To be able to apply the bipartite graph maximal matching algorithm, we need to refine all overlapping matches so as to obtain a resolved refinement Σ' of Σ

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

Resolved segment matches

- Given a set Σ of *segment matches* between two sequences A and B



- Want to determine a *resolved* subset Σ' of Σ in which the projections of any two matches onto each sequence are either disjoint or identical.

A
C
G
T
T
C
G
A
C
T
A
A
C
G
C
T
C
G
C
G
T
A
C
G
T
T
A

In other words...

- Given many overlapping segment matches produced by programs such as MUMer or BLAST
- Determine a resolved set of matches that do not overlap improperly