

Introduction to Assembly Algorithms

Mihai Pop

Recap...

- Sequencing technologies only "read" small chunks of DNA, yet genomes are substantially larger
- The shotgun sequencing approach generates many random fragments from the original DNA
- The task of the assembly program is to stitch together the many small pieces into a reconstruction of the genome
- Essentially..... a huge jigsaw puzzle
- Think: shred a collection of Harry Potter books at random then try to rebuild the original without any additional information.

Assembling two cities

it was the best

was the *age of*

best *of times* it

wisdom it was the *it was the age of times* it was

it was the best

was the best *of*

the worst *of times*

was the worst *of*

was the best *of*

times it was the

it was the *age*

times it was the

was the *age of*

the best *of times*

worst *of times* it

age of wisdom it

it was the *age*

of wisdom it was

it was the worst

the *age of* wisdom *of times* it was

the *age of* foolishness

Shortest common superstring problem

What are we looking for? (mathematically)

Given a set of strings, $\Sigma = (s_1, \dots, s_n)$, determine the shortest string S such that every s_i is a sub-string of S .

NP-hard

approximations: 4, 3, 2.89, ...

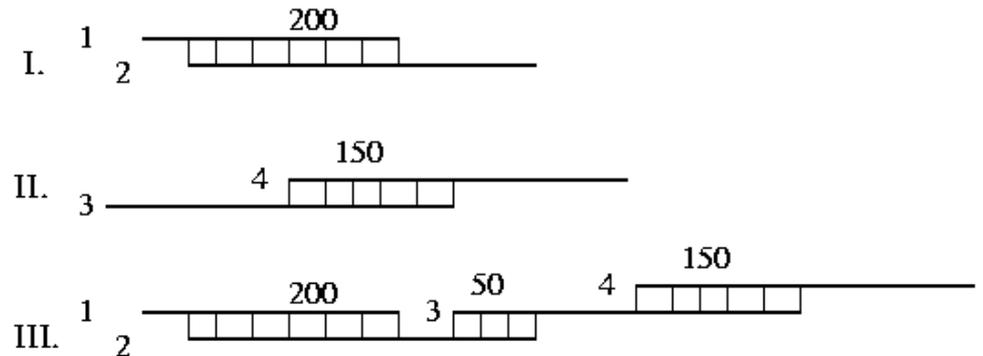
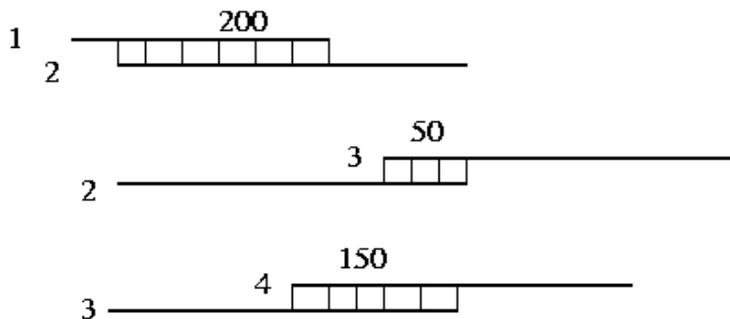
...ACAGGACTGCACAGATTGATAG

ACTGCACAGATTGATAGCTGA...

Greedy algorithm details

- Compute all pairwise overlaps
- Pick best (e.g. in terms of alignment score) overlap
- Join corresponding reads
- Repeat from * until no more joins possible

Greedy algorithm (4-approximation)



Greedy approach gets 'stuck'

it was the best
was the best *of*
the best *of times*
best *of times*
of times
times

it was
it was the
it was the age
it was the best
it was the worst

wisdom it was the

the worst *of times*
was the best *of*

it was the *age*

was the *age of*

worst *of times* it *of times* it was

age of wisdom it

it was the *age*

of wisdom it was

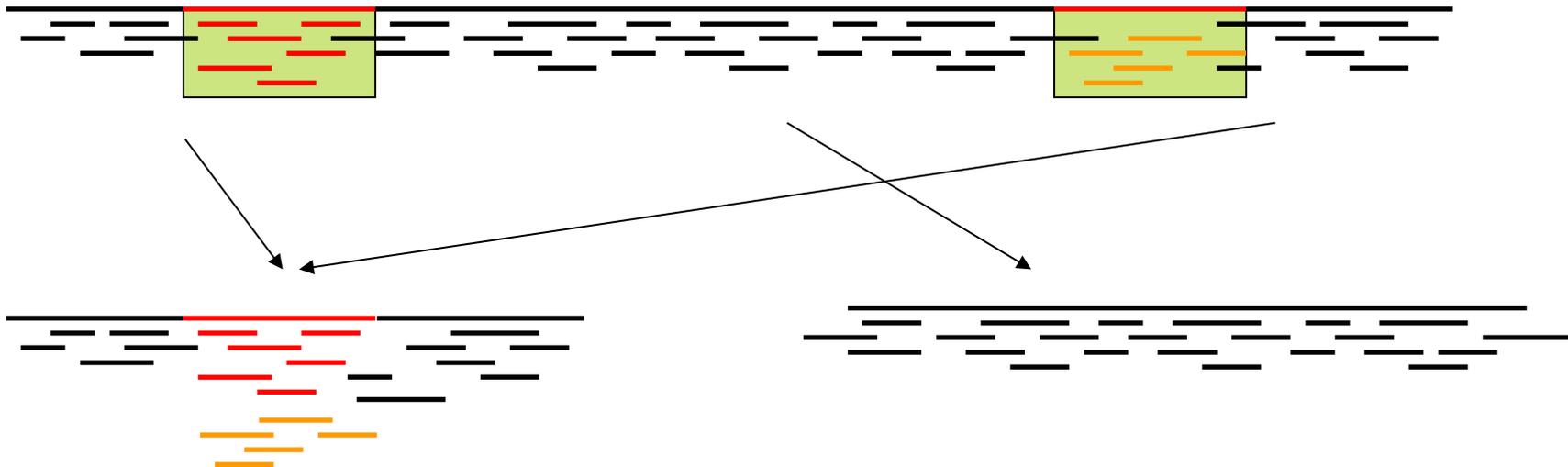
the *age of* wisdom *of times* it was

the *age of* foolishness

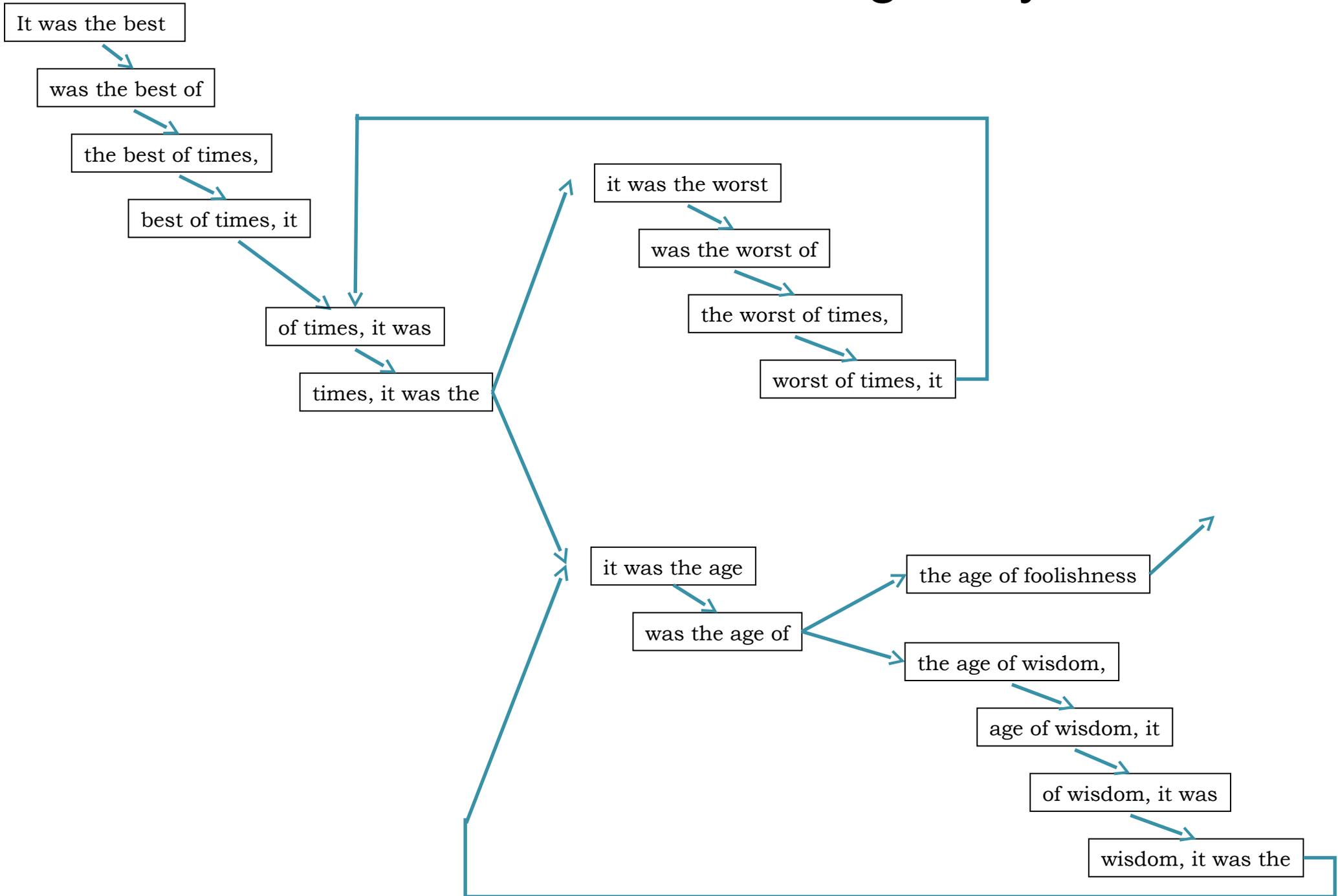
Repeats (where greedy fails)

AAAAAAAAAAAAAAAAAAAAAAAA
AAAAAA AAAAAA AAAAAA
 AAAAAA AAAAAA
 AAAAAA AAAAAA

AAAAAA
AAAAAA
AAAAAA
AAAAAA
AAAAAA
AAAAAA
AAAAAA
AAAAAA



Can we do better than greedy?



Graph-based assembly

- Better than Greedy (can see better what is going on)
- Repeats still a problem

Next: Assembly paradigms