

- photosystem I
- photosystem II
- cytochrome b/f complex
- ATP synthase
- NADH dehydrogenase
- RubisCO large subunit
- RNA polymerase
- ribosomal proteins (SSU)
- ribosomal proteins (LSU)
- clpP, matK
- other genes
- hypothetical chloroplast reading frames (ycf)
- transfer RNAs
- ribosomal RNAs

# SAM files

nkane@Serval-Professional: /media/Documents/Genomics2013/testing

5:23 PM

```
@SQ      SN:mt      LN:366924
SRR307232.1 4 * 0 0 * * 0 0 CAAAGGAACNCAGATATAATGCTGTGCAAAATAAT DCDDDC=C#C:@>- :769<DCCBDD:BB?CC/A
SRR307232.2 4 * 0 0 * * 0 0 CGAATCANCACAGCGGAACGGTAAGGCCGTGA GGGGFEEEE#DD:CBDDCC=:CBB@BBCBDD:C@
SRR307232.3 4 * 0 0 * * 0 0 TGAACCAANGATTATGAGTCCCCTGCTAACC D?:55CC>#CC-@=-:B:AAB?=D=B:B?:B5CDB
SRR307232.4 4 * 0 0 * * 0 0 CCGACCTGANTGCCTGCATCCAAAGGGCAACAGG GGGGGECEC#ADDCFCGGGEGGD:GFGBBGGGG
SRR307232.5 4 * 0 0 * * 0 0 ATGGGTCACTCCGTCGGCGATCTGTTACGCGG EFFFDBDB#C>CCC=CB:EEBEDFBF94:BCDC
SRR307232.6 16 mt 53886 37 35M * 0 0 CAGGCTCTTGACCTCGATCGAGCNACGAGCAG <6A0A?AB-EEBEE=??EECC?C#?=-CCEEDE
SRR307232.7 4 * 0 0 * * 0 0 CGTTCATGCNCCGATCCTCGTCCACACGGTTCGA DDD-DAAAC#DBA:CDDADBDD?DBD==:CBDD
SRR307232.8 4 * 0 0 * * 0 0 TGAAGCCANAAACAGCACCGTGCACAGCGACC DDBBCC-#9.-B=:DADBBABBD-B?:D??AD
SRR307232.9 4 * 0 0 * * 0 0 CGTAAATTTNGCGACTAAATAAATAGACGGTAGA BADD=C=#C.;6@=DDD??78@DAD=?C-CA
SRR307232.10 4 * 0 0 * * 0 0 TGATTCAGGTCGCGCAAAACACTGAACGCTGA DCDD55CC#DDCDDDD?BDA->DD=DDDDDD
SRR307232.11 4 * 0 0 * * 0 0 CATTGTCTCNGACTTCCCCGGTTCCTAATAGTCG FFBEDBD#CACACFFFEFEFDFCDFFFFFE
SRR307232.12 4 * 0 0 * * 0 0 ACAATATTCNACCTACAGCATCCACATTGGCTGC DD55-C:5C?C?C=>-A?5DD5B=?A5D5-D
SRR307232.13 4 * 0 0 * * 0 0 TTGACGACTNGATACCGCAAGATATGTTGGCAT FFFDFBDD#DADD=EDEEEEDDEEED:DEDEED
SRR307232.14 4 * 0 0 * * 0 0 GTTGTCTGTCGATGCTGCTGATCGGTCGCAAAAC DEEEDDD#C5AAB=?BAEEEEEEDEEEED?
SRR307232.15 4 * 0 0 * * 0 0 GTTGTGTANTATGAACTAAGAAACATAATGCTC EEEEACC=@#A-AA?BDD?:E=C3CBC?BEEBE
SRR307232.16 4 * 0 0 * * 0 0 TTGTCAATTNGACTGTGGTCAATTCCGGTTGAG GGGDDDD#DBCBDFEFDF?BCEFFDAFCFFAF
SRR307232.17 4 * 0 0 * * 0 0 AACCGCGATNTGGTCAGCGCCAGCAACCGCA EEE=EDDB@#>:A<C?BBB?5DB?5-@,BA?DA5
SRR307232.18 4 * 0 0 * * 0 0 TCGCCGGCCGACGCCATGGAAATATCCAGCCCC FFBDFDBD#DB:AD:@;BDD?CDBD?DD?D
SRR307232.19 16 mt 126424 37 35M * 0 0 TTAACCACTACCAATTAAGAGTTAAAGCTCAC ??@DCA@C?<CCCA?BDEFAC@>=#@CCCE@BE
SRR307232.20 4 * 0 0 * * 0 0 TAGCTGGTGNCAATTTGTGCAATGAAATCCAAG GEGGBDD#DDDCGGGGGGGGGGGGGGGGGGG
SRR307232.21 4 * 0 0 * * 0 0 CTGGGCTGNTTCCCGGGGGCGGACGCTCGGCNN EACBB@A#@A@A@E@EB@9B:B?@E@#####
SRR307232.22 4 * 0 0 * * 0 0 CTTGGAAAANTAGGAAAGTGGGAGATAGGGGATA DCBDA?CBA#CD?BEEBEDD?B8->CCAC;@>
SRR307232.23 4 * 0 0 * * 0 0 AACCTTACNGATAAAACAGACCTGGTACAAAGCG GFGGEDE#DCBADDEEEEEEDEEEEEEDEE
SRR307232.24 4 * 0 0 * * 0 0 CGCTCGCCNTGCACGAAATCCAGGTGGAAGCGG GFGDDDD#DDCDFEEFFEECEDDFE?CEEDE
SRR307232.25 4 * 0 0 * * 0 0 CGCCATATNGAAAGACCTACAAGCCCTTCA FFEFEBE#DA-DDDDAD5DDDD=D?B?D=DDC
SRR307232.26 4 * 0 0 * * 0 0 TCTTACTCAGTCAACACTGACATCAAGTCAT GGGGGGGGGGEGGGGGGGGGGFGGGGEGGGG
SRR307232.27 4 * 0 0 * * 0 0 TTCCATGGCNTGGCGCCAACACGAGCTGTTGTA EGGG?DEEE#B@BDDDDDD=BCDDBDDCCDBB
SRR307232.28 4 * 0 0 * * 0 0 GCAATTAATGATTTTGGTAATTAACAGCGCCA D?:BDEEDDE??D?=-DEEB5BE5E:EEDA=EBEA?
SRR307232.29 4 * 0 0 * * 0 0 CTCAAAACACTAATTAAGCTTCTTAGCTTCTC A:BB?DD=:>C5CA5C:CCDD=-?:@=BD:A#
SRR307232.30 4 * 0 0 * * 0 0 CCAGGTGATCCGCTTGTGCACGAGGACACCGGG @@@@:@@??:DDDDADDDDD=DD5?DDADBD
SRR307232.31 4 * 0 0 * * 0 0 AGTGGTAGAGTGGCTTGTGTCACGATCCACTGA ?AA<>'?:21)??3=5@@C=D5DADD5C5DC;,B5
SRR307232.32 4 * 0 0 * * 0 0 GCAGTATATCTCAGTGTGCTCTTTGCTTGT :@:7=*787=DDDDDDDDDDDDDD=DDDDDDDB
SRR307232.33 4 * 0 0 * * 0 0 TCAATGCGAGCTGATGACTCGCCCTACTAGGAA @@Q=*>@;EEEE?DFFFFFFFBBFFDFFFD
SRR307232.34 4 * 0 0 * * 0 0 TCTCGGTCGTTGCTGCTGCTGCTGCTGCTGCTG GFGDFDFD=EEC=EBDE=C@DDA=EB?CDFEFD
SRR307232.35 4 * 0 0 * * 0 0 GGGTGAGCTGTAATGCGAATGTGGTGGGCGG EEEAA@C>CAC>>@CBBDB=:BD@?EE2EACC=EA
SRR307232.36 4 * 0 0 * * 0 0 CATACGGATGATGATGTTGAAAGTCAAGAAAG GDGDFFFFFFFFDFDFDFDFDFDFDFDFDFDF
SRR307232.37 4 * 0 0 * * 0 0 GCTACACAGGATGATGCTCGCCGCTCAGCGAT FGFBDGGGDF;D55CDA@EEDEEECEEEEE5;
SRR307232.38 4 * 0 0 * * 0 0 GTCATGCAGTCCGACAAAGGCTGATTCGCGTCG B=B=B->?B:B=BEFFBDBADFFFAAFFFFFBF
SRR307232.39 4 * 0 0 * * 0 0 GGATGGTTCGGTGTCTGTTTCGCCGAAAGCTC FG?GGFAAFBEBBEBEBAEDCE?EEBFB?F:==5E
SRR307232.40 4 * 0 0 * * 0 0 TATGACCCCTACCGCTCTCGCCGCGGCTTCGAT GFFGGFGDGGGFGG=?EAEEBEE5BDD=EEA5B
SRR307232.41 4 * 0 0 * * 0 0 GACCATCTGCTCGCCGCTGGCCTTCGGCGCAGC GEGGFGDF=BEEDDAD?EBBD?:C?:B)A;C5
SRR307232.42 4 * 0 0 * * 0 0 CTATGGATTGACTCGATGACTGACGCAAGGAGATC DC:AD?C??AC?DB:=-C->?<?>=???5-?
SRR307232.43 4 * 0 0 * * 0 0 ATAGAAATACAGGCTTTGGCAAGACGTAGACTTA AEDEBEAE;E>@=C)@;277DD:D?D?BBD-DD?
SRR307232.44 4 * 0 0 * * 0 0 TGCGCGCCACACCGTGTGGTGGTTCCTTGG E=EBDEEBEFFEFA?C?C?ACA=@SA@A?CB?@B
SRR307232.45 4 * 0 0 * * 0 0 ACCATAAAATATCTCCAGTTGATGTGCTTTTG DDDDDDD;BCCDDDDDDCDDCDBDDDDDDDD
SRR307232.46 4 * 0 0 * * 0 0 TGCCAGCAATACCTGGCCACAATACGACGCCA FFFEDDFD?:EEEE=BDDDBDD:D;D?CDBD?D
SRR307232.47 4 * 0 0 * * 0 0 TCACCTGCCACTGTGGGACAGCCCTGGTAAA :C?>5;@AAEE;=5BDBDB-?B?75A-CCC=?
SRR307232.48 4 * 0 0 * * 0 0 ACCGCCAATCGATCGGCGACCGGATGCTCG GGGGGGFGBDGGEFGGEE9FFBE?5B:DBFF?
SRR307232.49 4 * 0 0 * * 0 0 CTCAAACCTTCGCAAGCTCAACAGCTTCA D@BD?@BBBBD-ADB?:@?CCA?*>;8;B?B#
SRR307232.50 4 * 0 0 * * 0 0 AACCGCCGAGGAGTTCGGCCCTCGCTGAGAG EEDDDDE;BDD:5B:DA=AAA/C,=,=CC:-?
SRR307232.51 4 * 0 0 * * 0 0 GCTCAGGGCATCTGATTAAATCGCTACCAAT FFFAFEFDF=FBEBEBAEBEBEBEBEEFEFEA
SRR307232.52 4 * 0 0 * * 0 0 CGGCTTTCGGGTGTCAGGAAAGACTCGCGCTCG E:E:?:EBED?564,>>@=BDDBCB?5DB-C
SRR307232.53 4 * 0 0 * * 0 0 AATTATCAGTCACTGCTCAGTCAAAATTTATGC EDEEECEBE@E:CCCEEB?DA:5=>7=ABB?BC
SRR307232.54 4 * 0 0 * * 0 0 CCAGCAGCAGTGTTCGACGCTCATGAATTCATGA EAAEEDDE;:E@E@EDA=D?<CA?=?>?B?@B
SRR307232.55 0 mt 329525 37 35M * 0 0 TTCAACAAACAGGTTTCATCTGCTGATGATGAT GGGGGE?BFGDGGFGDGGGDFDFBFAFBF
SRR307232.56 4 * 0 0 * * 0 0 CGCGAAATGGATAGAGCCGAGTACAGATACGGC D5DDAD=-D=D-BDD=:DB5DA=?-?:5C->
```

Convert to fastq

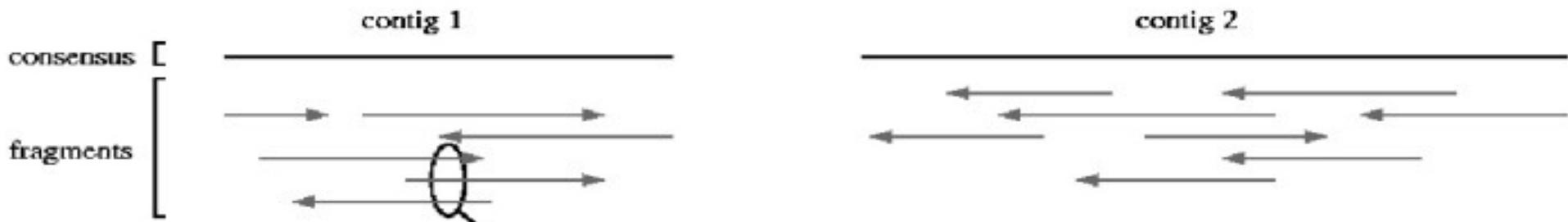
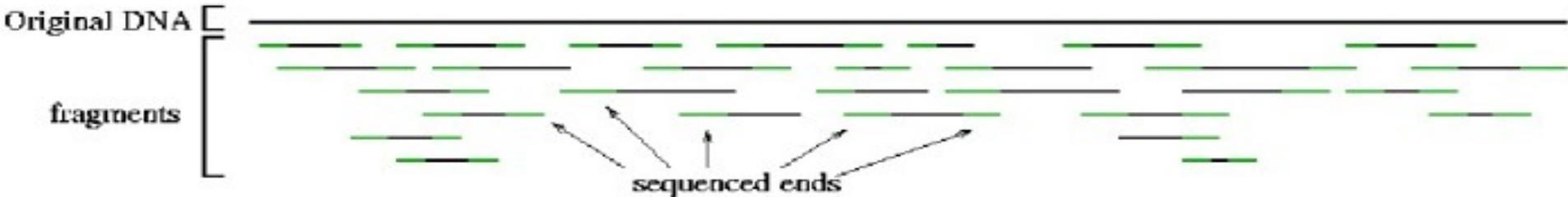
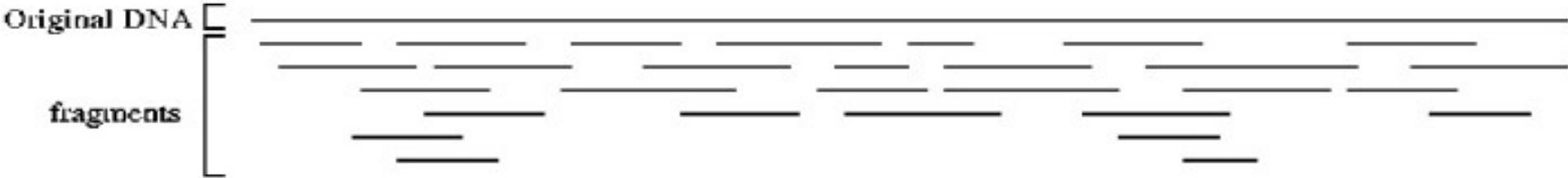
Convert to fastq

# De novo genome assembly

```
@HWI-ST765:7:1101:1318:2091#0/1
GGCCACCTATGACCGGCTCGCGCCGCTCGTCGGGGAGCGGCTGCTCGTCGTACCGGGGGCGCGCCCGCGGACGCCGTCCGCGGCCCGCTCCGCGCGCCCC
+
_____ccccggggghhhhh^b^c_UZFLZWacdBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
@HWI-ST765:7:1101:1628:2156#0/1
TCTTCGCGAGTATGTCTGTTGATGGCGCTGTGTCTATCTGCTCAAGGAAAGCAGCCCAACTCAATGTGTTACGCATTAGCGGCATTTGCTACATAATCCG
+
_____eeeefggggf_bddgeafgihdgehhgfeghhhifbgfhhhhhhhhhhdhiggfede`d`]bbdbcccccccccccccccccbcb`b`bdbcbcc
@HWI-ST765:7:1101:2627:2192#0/1
ATTATGAAGACTGGAGAAAGCCCTATATTTATTGTATTTCTTTTCTGGATCACAAAATCCTCCCCTCTGAAACAAAAGATGTAGTTGGAATAAATAAAAGG
+
bbbeeeeegfgegghffefghiiiiihhhhfghhicegihihhiihfiiiiihfihiihfihhhihfdggeceeee_bdddbccbdbdbccb_
@HWI-ST765:7:1101:3236:2246#0/1
GCGGAAAGAGGGCTTGAGGATGACTTCCCTCATAGACTGGGACCCCCACTTTGAGGTGGCTGACGTAGCCTTTAAACGGAGTCCCCGCATTCCCGGTATCT
+
bbbeeeeegfgggiiiiihhhiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiihghggfеееес`сdcccc`bcccc^bcccc]aacdccc[_ccd
@HWI-ST765:7:1101:3400:2241#0/1
GCGGACAGCTAATGCGTTCCTTATTTGAACAGGGTCTATGGTCCGTGACCCCGGATGCCGAAGGCGTCCTTGGGGTAATCTCGTAGTTCCTACG
+
_____cacc_eeaegfffZa`e]]de`egdfg[сgfсgZf]e^aX^G[Ze_agfffdgс`bXZ^[_]_aaa_GTTTTW_SX`aTX]`_bbaa_aacY`bbRO
@HWI-ST765:7:1101:4139:2060#0/1
NCTTCTCTTTCATCAGAGAGTAGAGTTGGGGCAATTGTGGGATCACGACGGGGACAGGGGCAGGTGCGGGCGGCGTCTCCGGTTGAGGAAGAGGCTGCC
+
BS\cceeeggggiiiiihifgiiiiiffhiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiggeccccccccccT__acX_c]][acc_cT[_`bcbaa``caa^`
@HWI-ST765:7:1101:4188:2089#0/1
ACAAGATATATTTGATATACTAAGATGATAGCTAGAGACTAGAGATGAGAGTGCAGGATCTAGATTTGTAACAAATATTCGACTTTGCTTATGCAAACCTGT
+
bbbeeeeeggggiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiihifghiiiiihiiiiifghiiiiiiiiiiiiiiiiiiiiihhhhhhhhhggggeeееееdcccccc
```

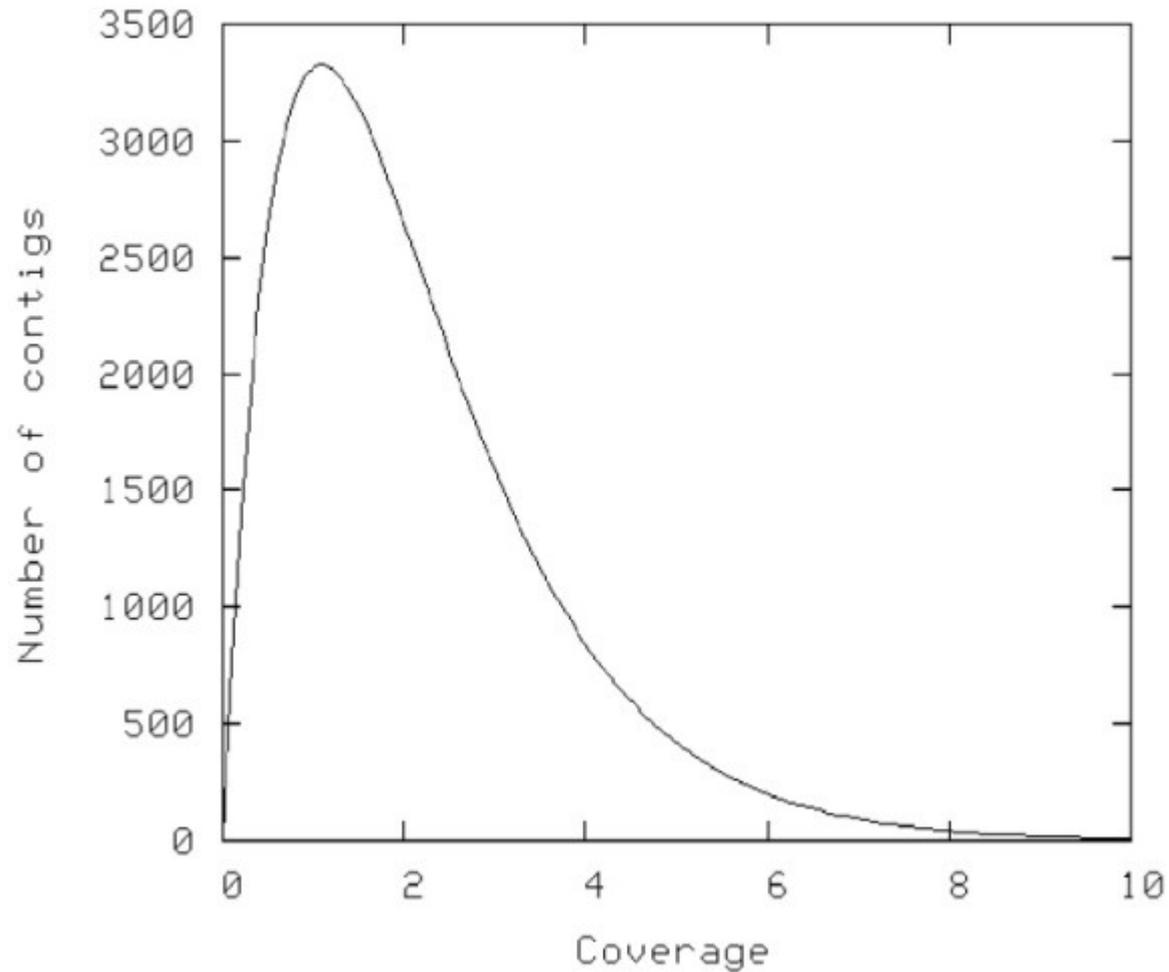


# DNA extraction, sequencing, assembly

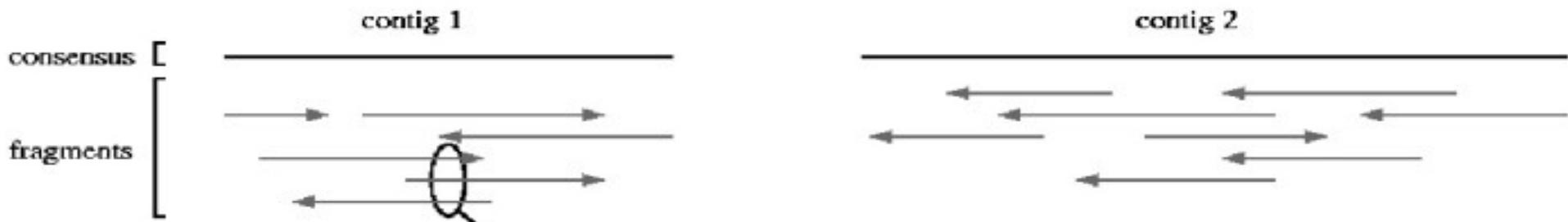
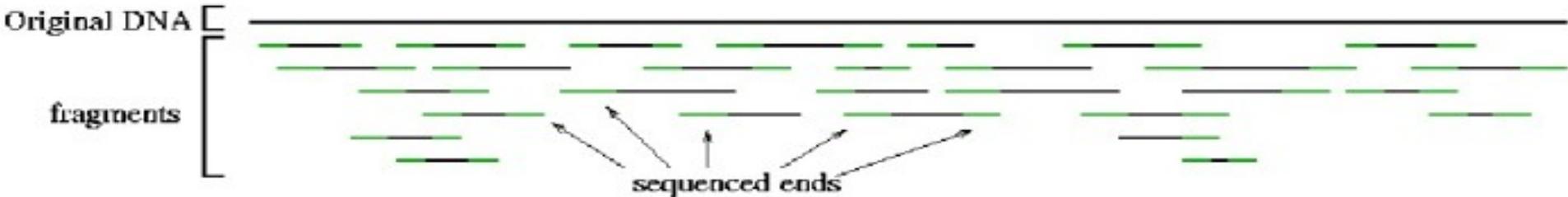
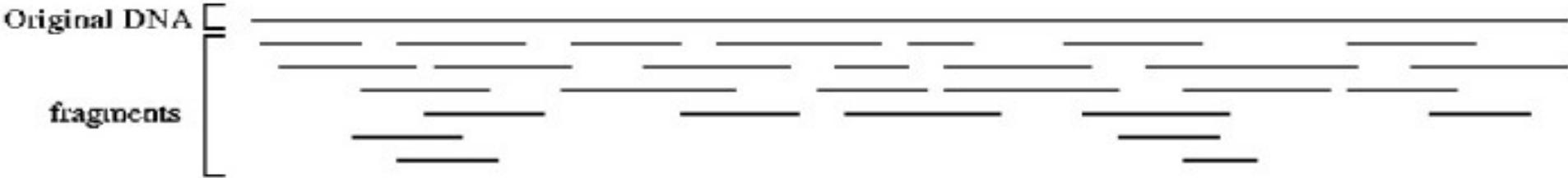


```
AAA A C T C G C C T G C T T A T C A A C C G A T C C C C C G C T A C C T T C T A C A G C C A T C A T T T  
AAA A C T C G C C T G C T T A T C A A C C G A T C C C C C G C T A C C T T C T A C A G C C A T C A T T T  
AAA A C T C G C C T G C T T A T C A A C C G A T C C C C C G C T A C C T T C T A C A G C C A T C A T T T
```

# Number of contigs vs. genome coverage

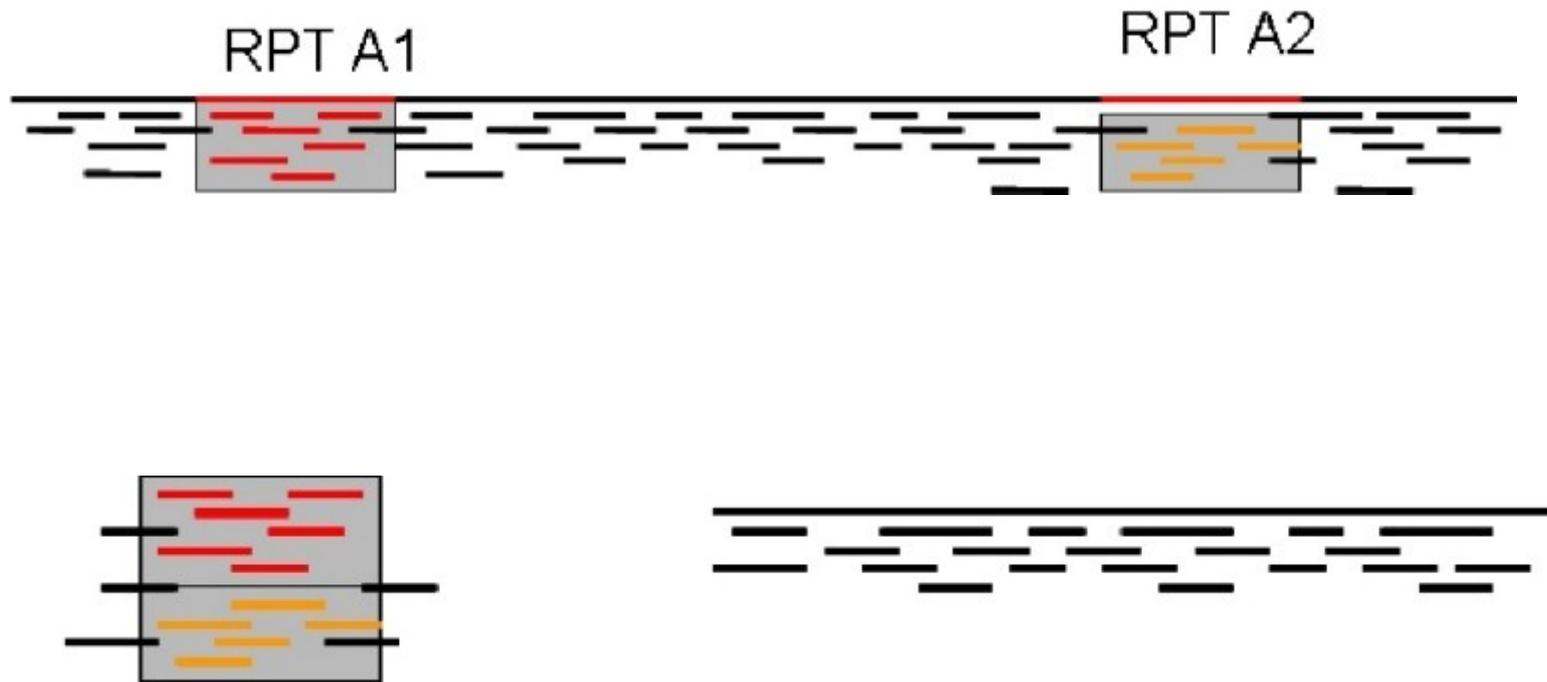


# DNA extraction, sequencing, assembly



```
AAA A C T C G C C T G C T T A T C A A C C G A T C C C C C G C T A C C T T C T A C A G C C A T C A T T T  
AAA A C T C G C C T G C T T A T C A A C C G A T C C C C C G C T A C C T T C T A C A G C C A T C A T T T  
AAA A C T C G C C T G C T T A T C A A C C G A T C C C C C G C T A C C T T C T A C A G C C A T C A T T T
```

# Repeats can cause challenges



# Assembly algorithms

- Overlap-layout-consensus

# Assembly algorithms

- De Bruin graph

# What is a k-mer?

- A k-mer is a string (sequence of letters) of length k
- ATGTAATAATG
- ATGT

TGTA

GTAA

TAAT

AATA

ATAA

TAAT

AATG

# Assembly algorithms

- it was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness

# Assembly algorithms

- it was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness

it was the best

it was the age

age of foolishness

it was the worst

times, it was

was the age of

it was the

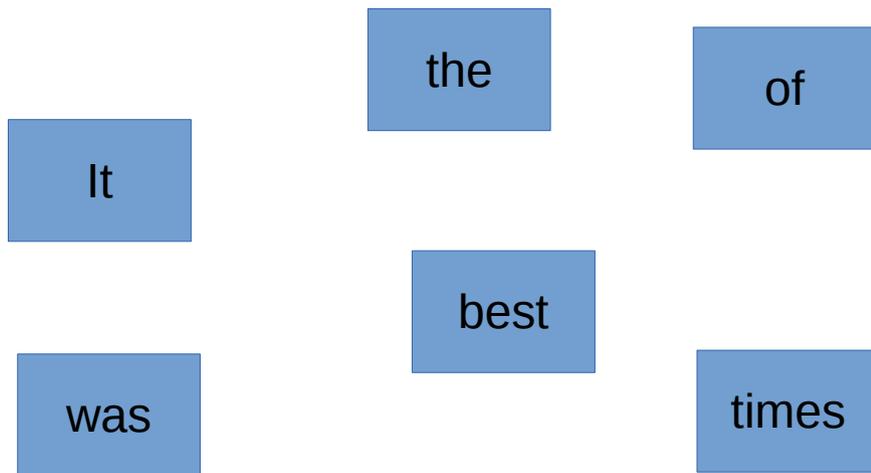
wisdom, it was

was the best of

the best of times

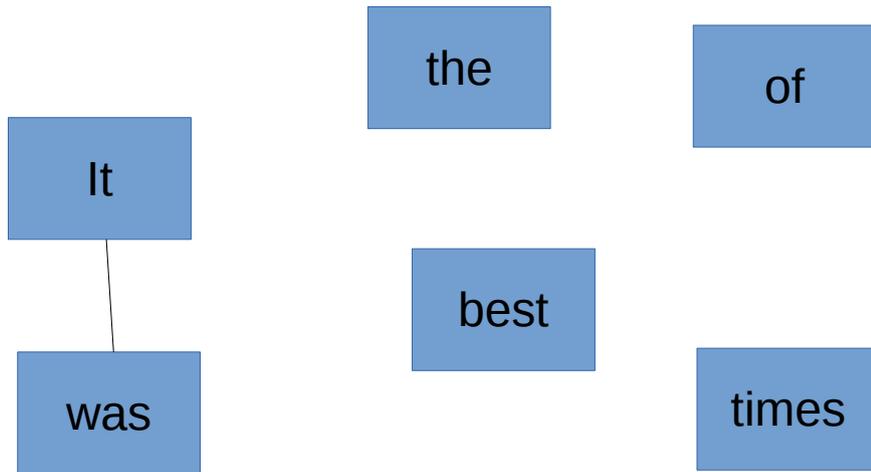
# Joining reads with k-mers

- It was the best
- was the best of
- the best of times



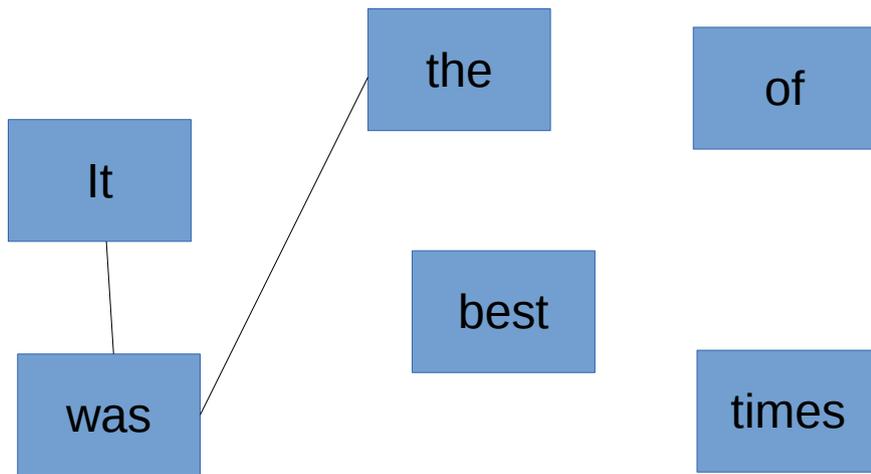
# Joining reads with k-mers

- It was the best
- was the best of
- the best of times



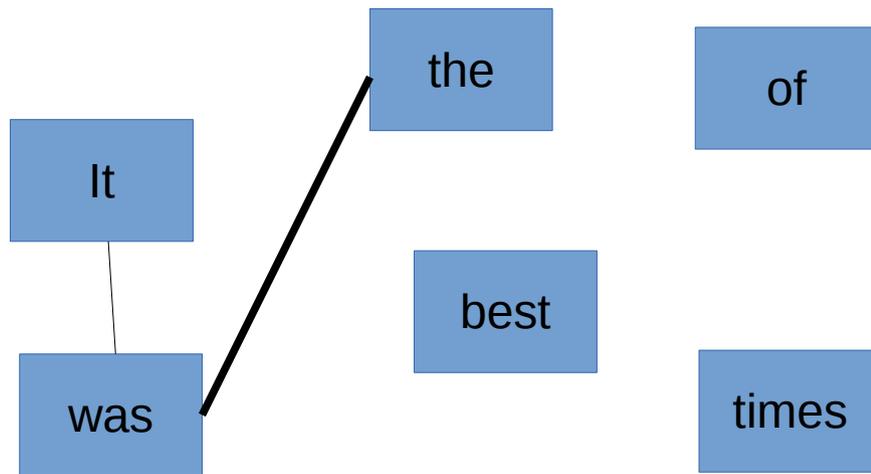
# Joining reads with k-mers

- It was the best
- was the best of
- the best of times



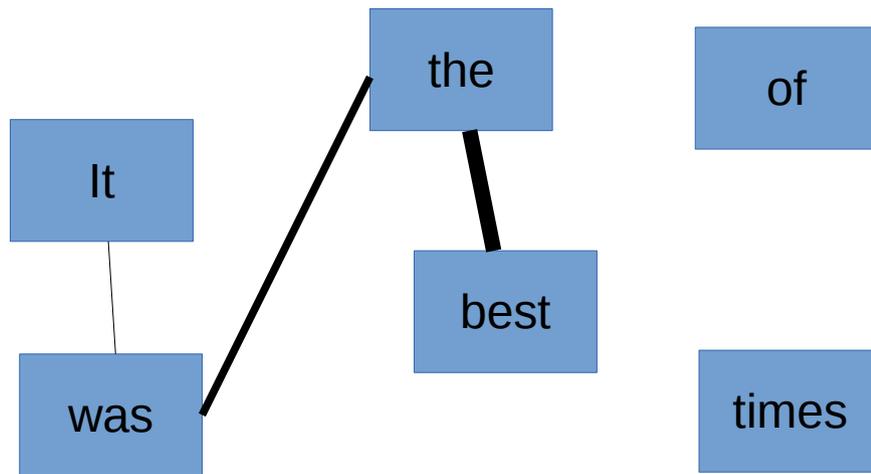
# Joining reads with k-mers

- It was the best
- was the best of
- the best of times



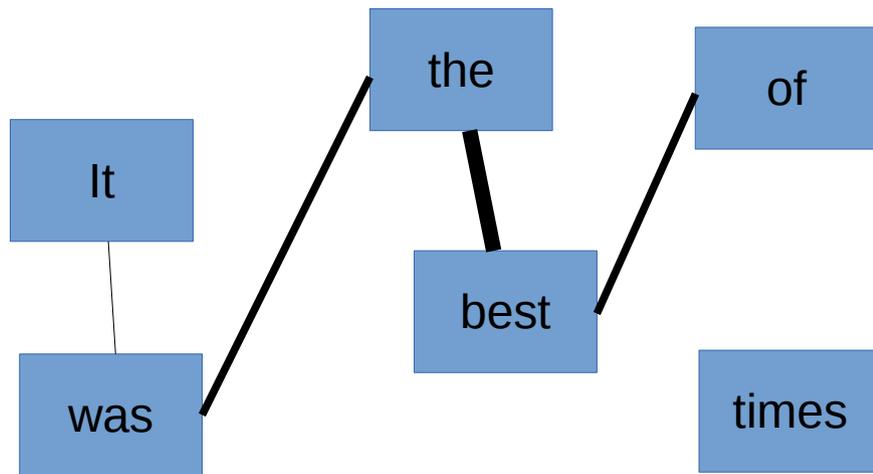
# Joining reads with k-mers

- It was the best
- was the best of
- the best of times



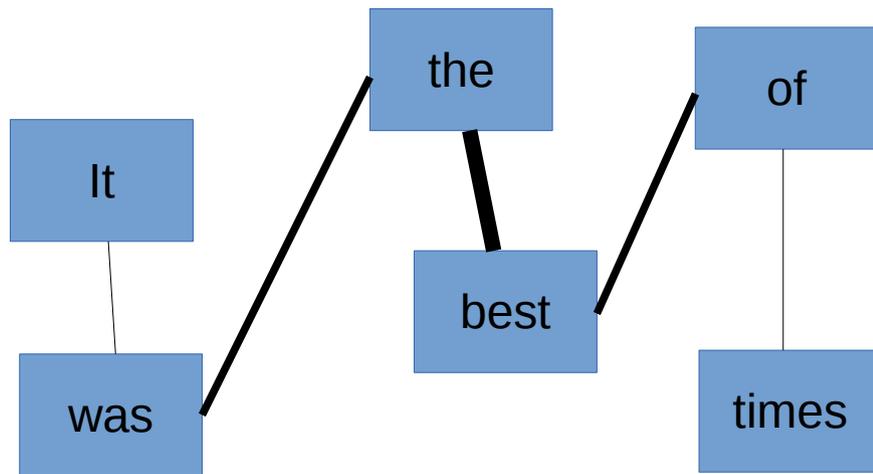
# Joining reads with k-mers

- It was the best
- was the best of
- the best of times



# Joining reads with k-mers

- It was the best
- was the best of
- the best of times



# Joining reads with k-mers

- It was the best
- was the best of
- the best of times



# Assembly algorithms

- it was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness

it was the best

it was the age

age of foolishness

it was the worst

times, it was

was the age of

it was the

wisdom, it was

was the best of

the best of times

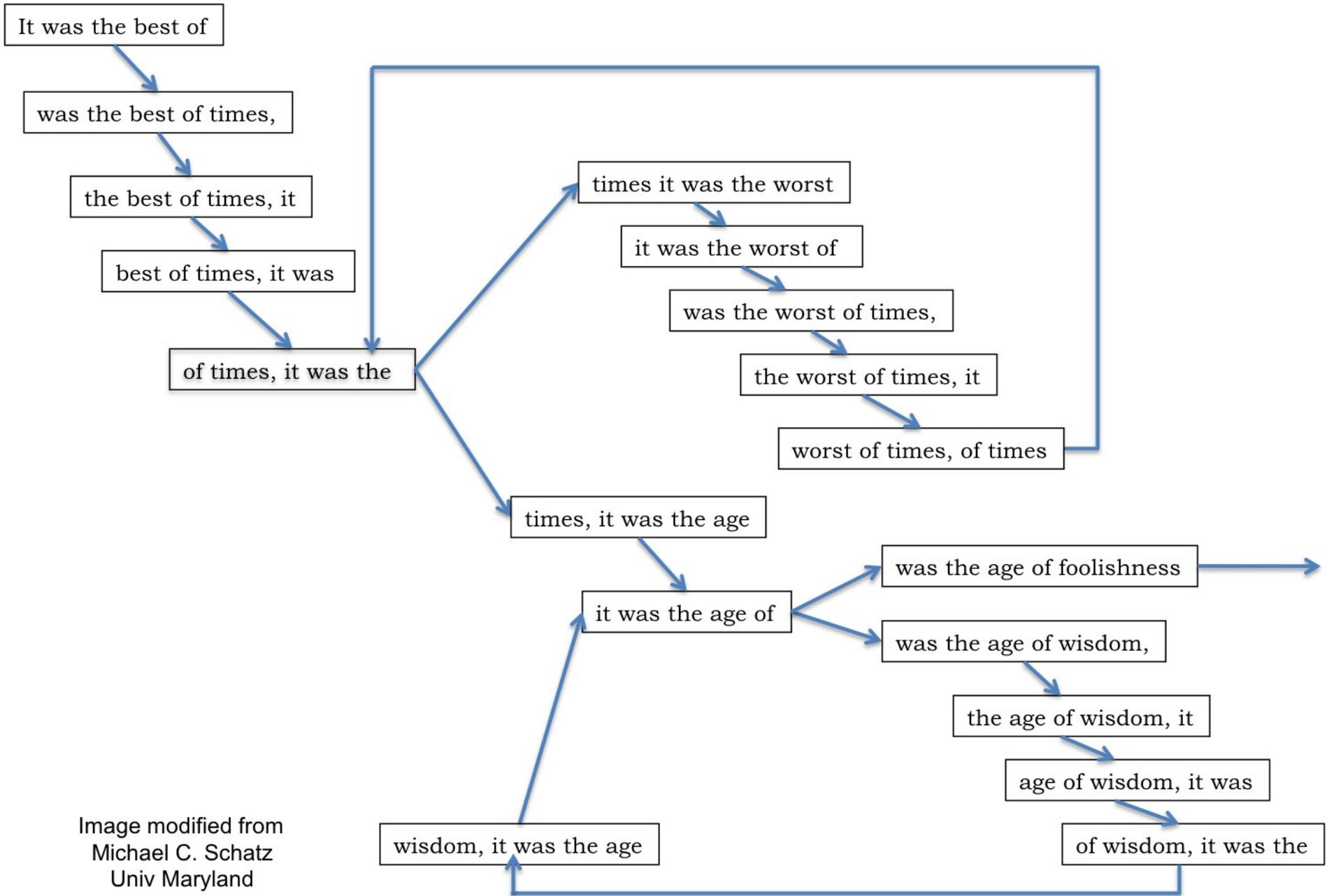
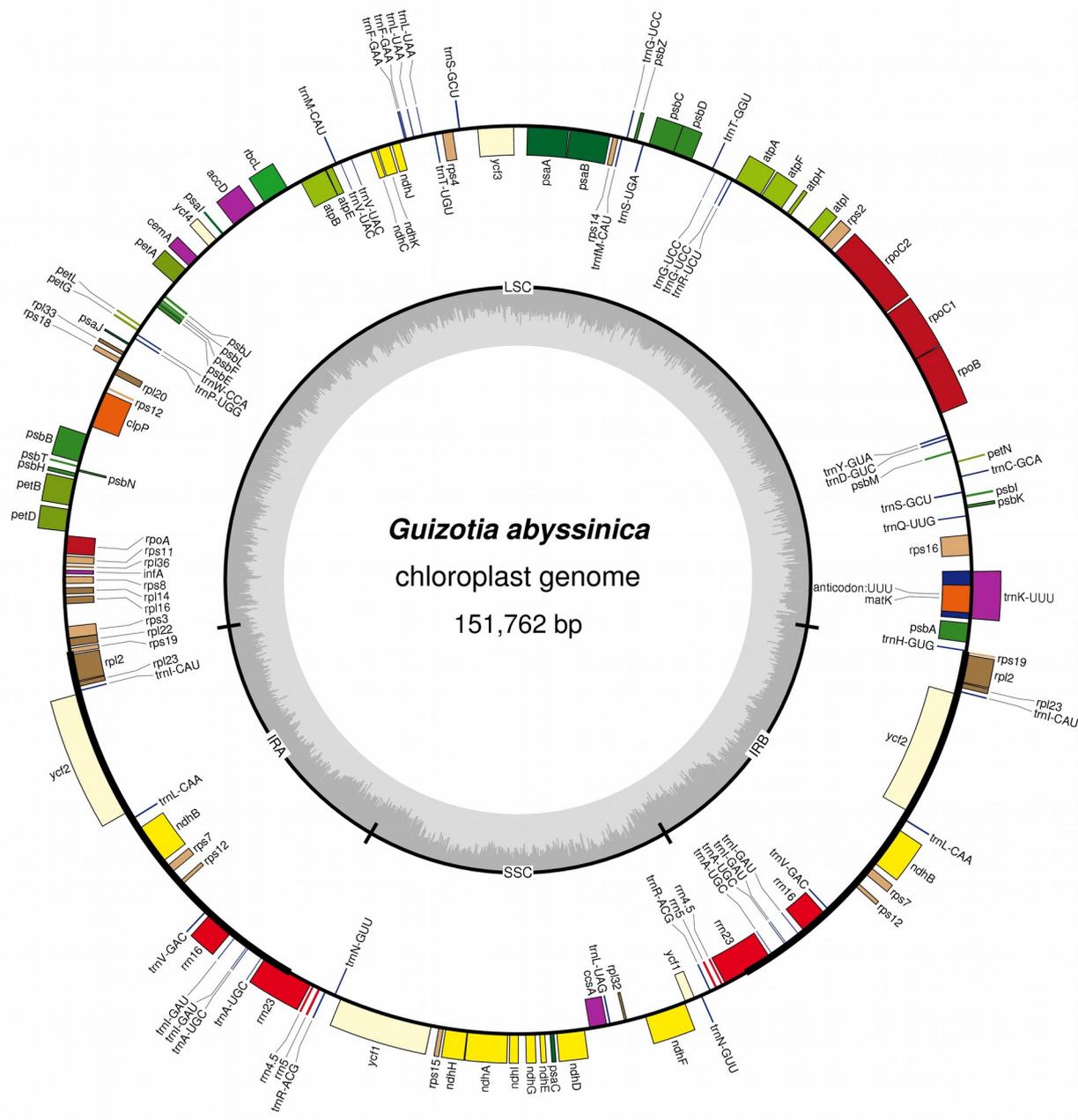


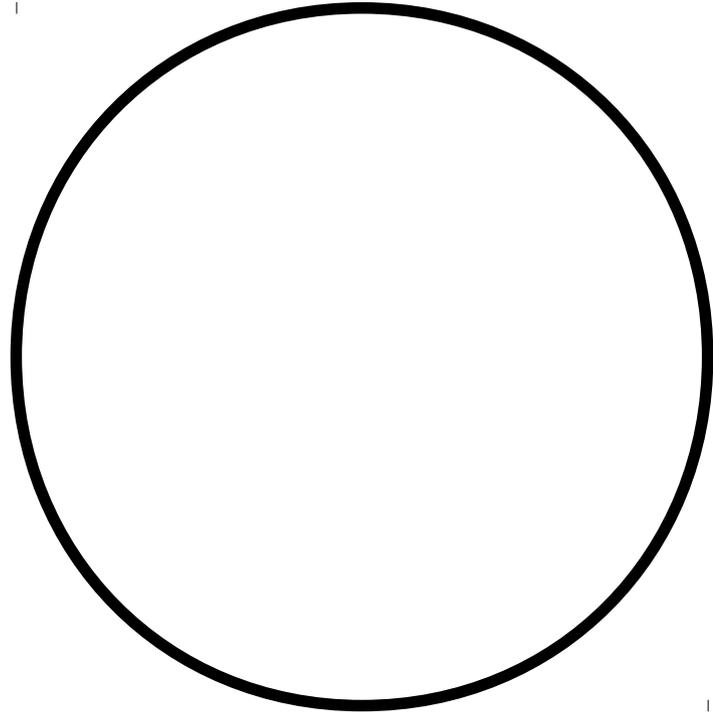
Image modified from  
 Michael C. Schatz  
 Univ Maryland

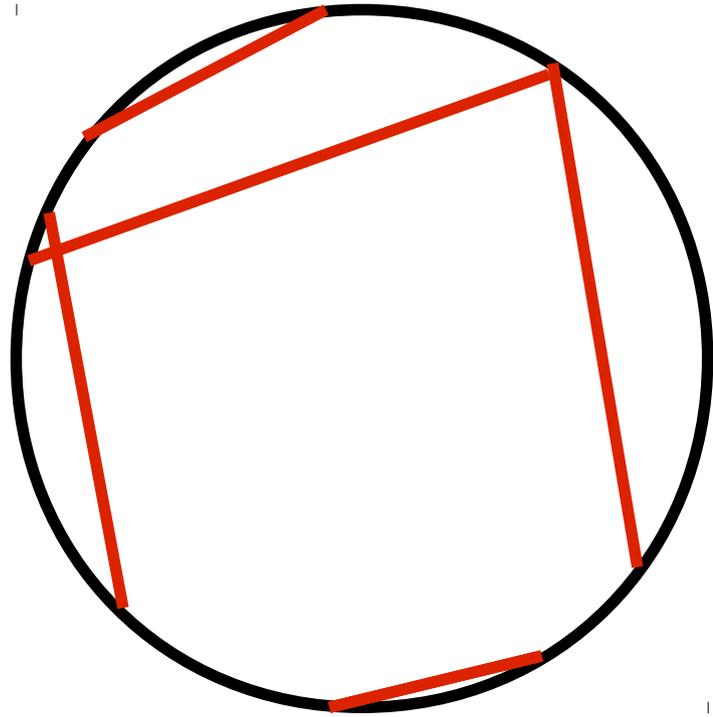
# Assembly algorithms

- It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of Light, it was the season of Darkness, it was the spring of hope, it was the winter of despair, we had everything before us, we had nothing before us, we were all going direct to Heaven, we were all going direct the other way

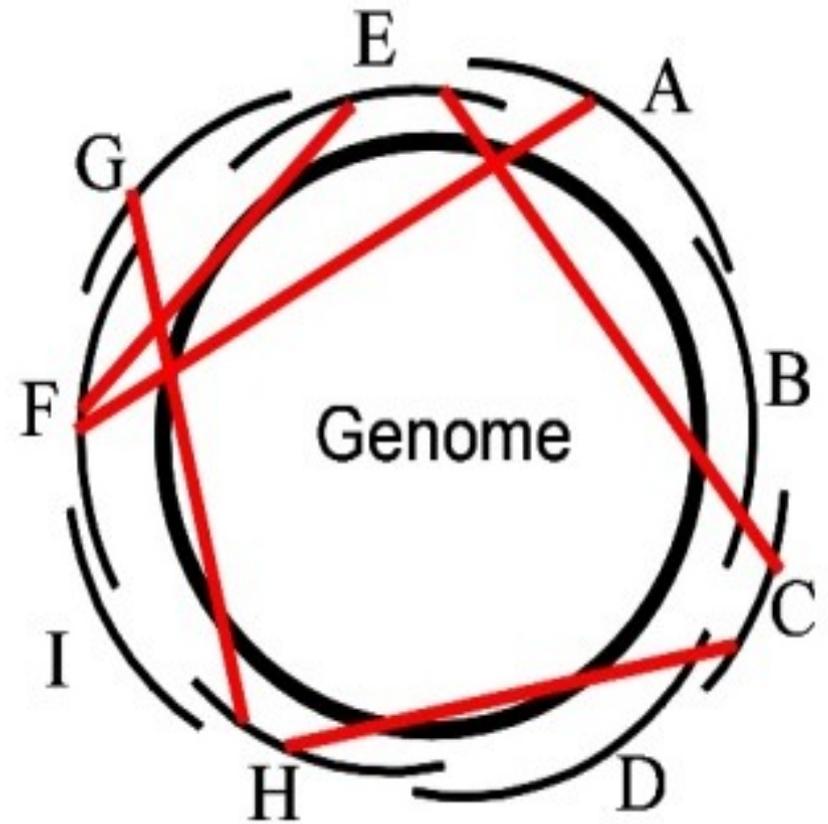


- photosystem I
- photosystem II
- cytochrome b/f complex
- ATP synthase
- NADH dehydrogenase
- RubisCO large subunit
- RNA polymerase
- ribosomal proteins (SSU)
- ribosomal proteins (LSU)
- clpP, matK
- other genes
- hypothetical chloroplast reading frames (ycf)
- transfer RNAs
- ribosomal RNAs

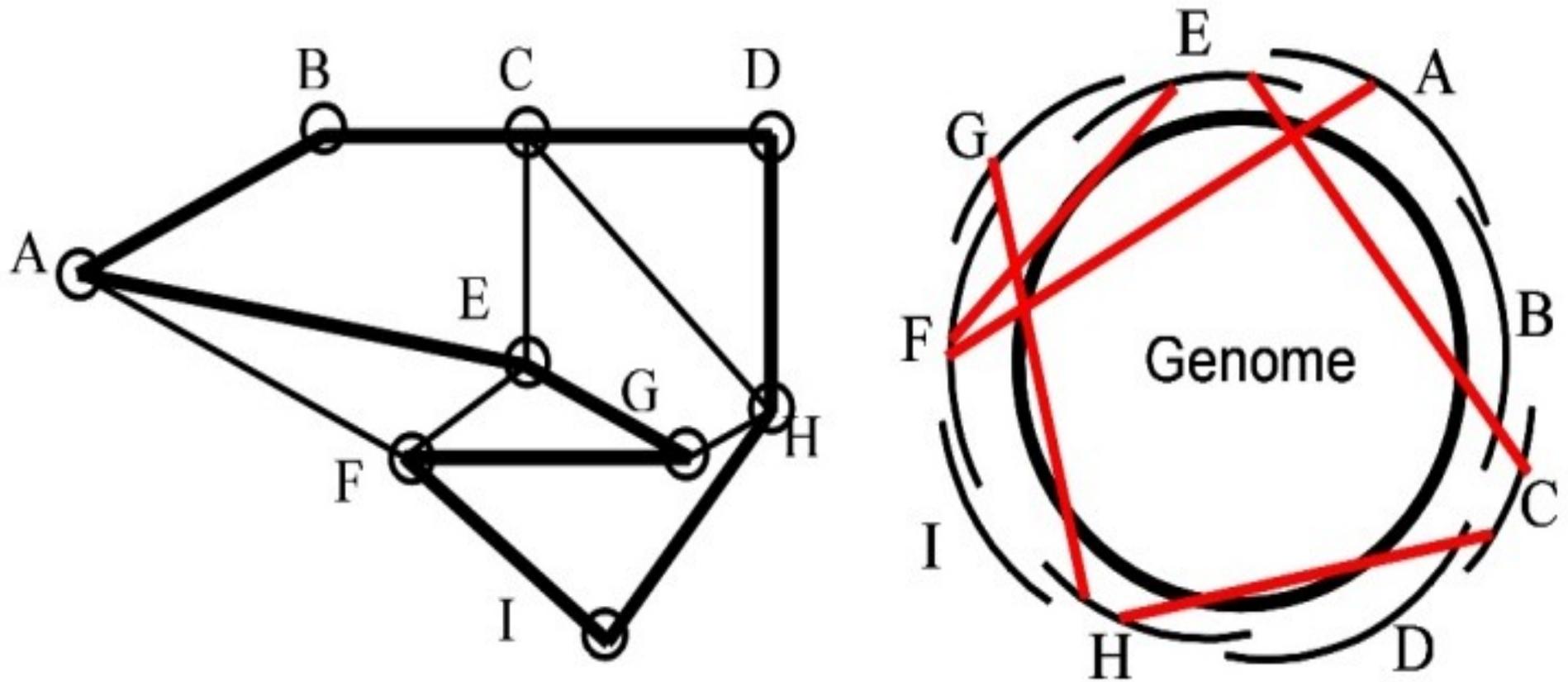




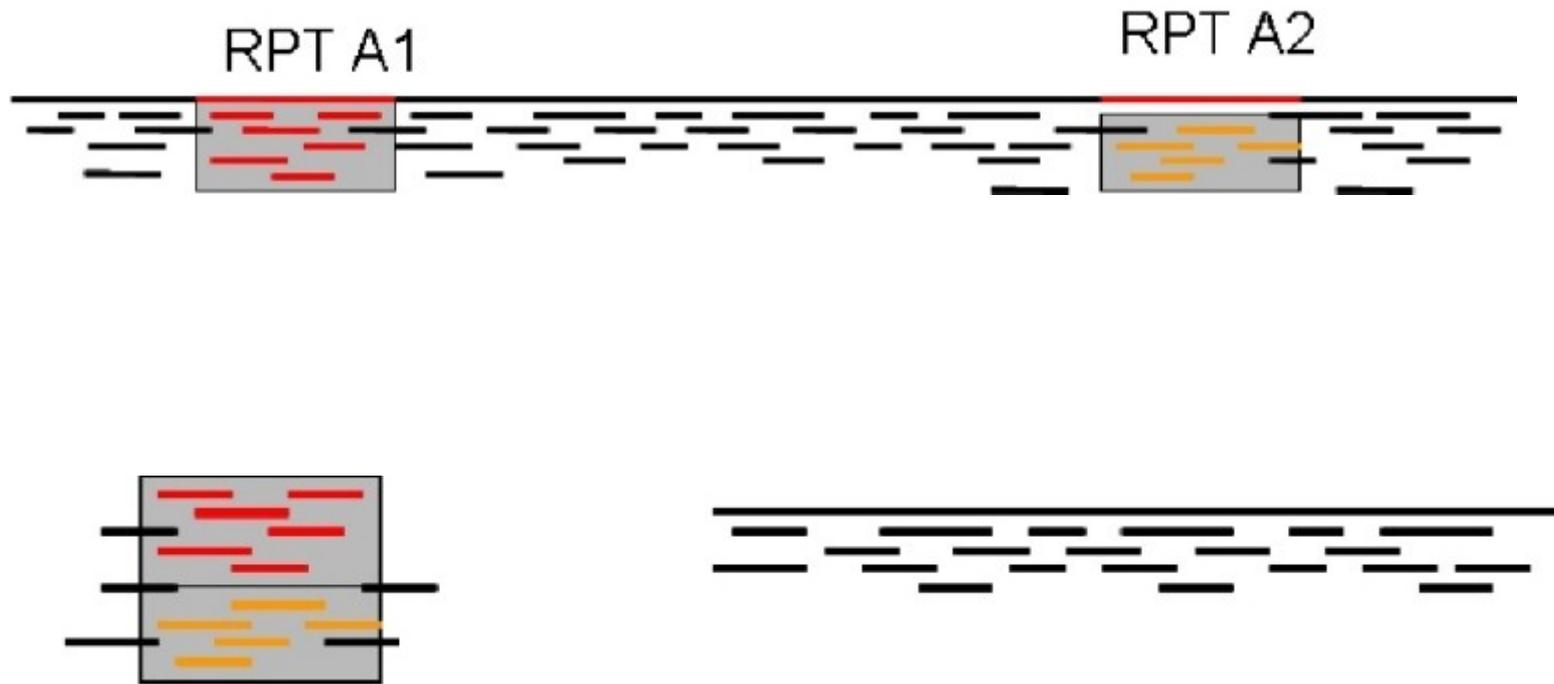
K-mers connect reads -> assembly



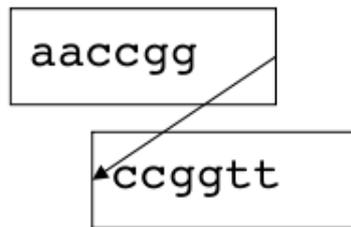
K-mers connect reads -> assembly



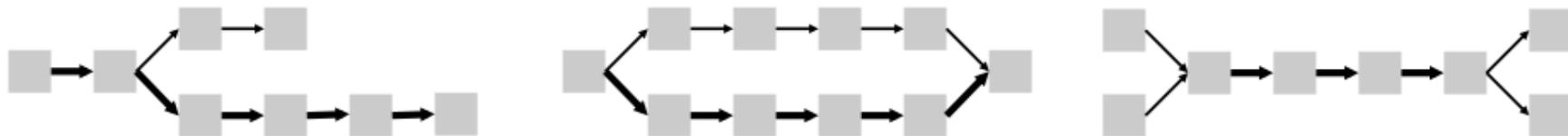
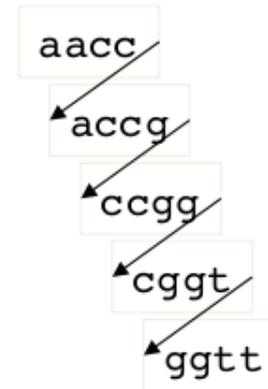
# Repeats can cause errors



# Graph Algorithms



Assembly algorithms construct graphs, which are nodes connected by edges. In traditional assemblers (left) nodes are reads, edges are overlaps. In de Bruijn assemblers (right) nodes are K-mers, edges are exact matches of K-1.



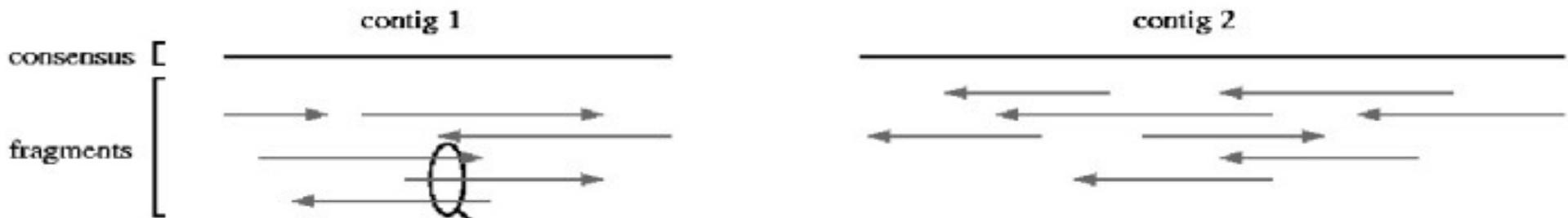
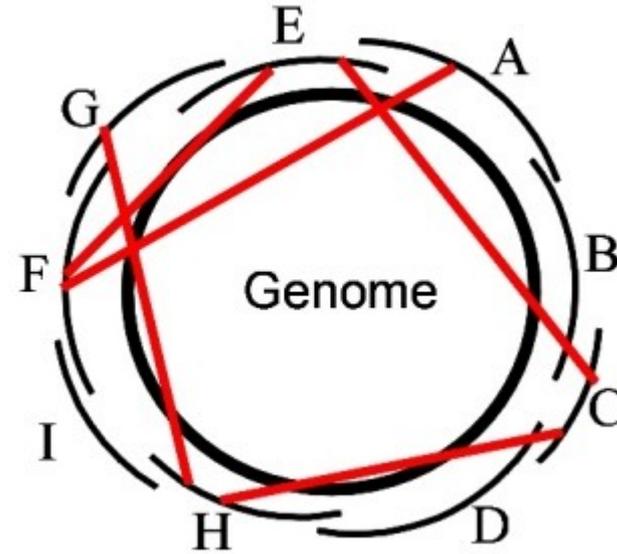
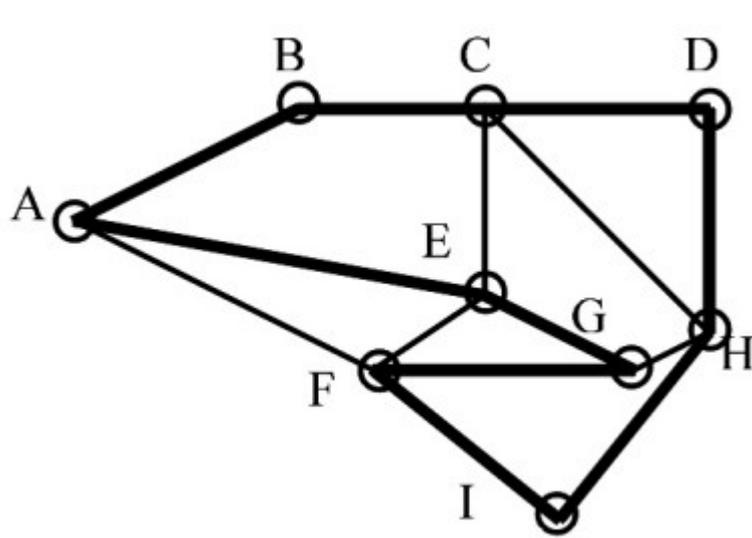
Assembly algorithms reduce graph complexity. Boxes are reads or K-mers. Edges are overlaps or matches. Edge thickness can indicate amount of support in reads.

Left: A spur is induced by bad sequence at a read end or low coverage and polymorphism.

Middle: A bubble is induced by polymorphism or sequencing error.

Right: A collapsed repeat might get teased apart or isolated or multiply placed.

# Evaluating the assembly – is it right? Which assembly is better?



```
AAA A C T C G C C T G C T T A T C A A C C G A T C C C C C G C T A C C T T C T A C A G C C A T C A T T T
AAA A C T C G C C T G C T T A T C A A C C G A T C C C C C G C T A C C T T C T A C A G C C A T C A T T T
AAA A C T C G C C T G C T T A T C A A C C G A T C C C C C G C T A C C T T C T A C A G C C A T C A T T T
```

# Assembly: varying kmer size

