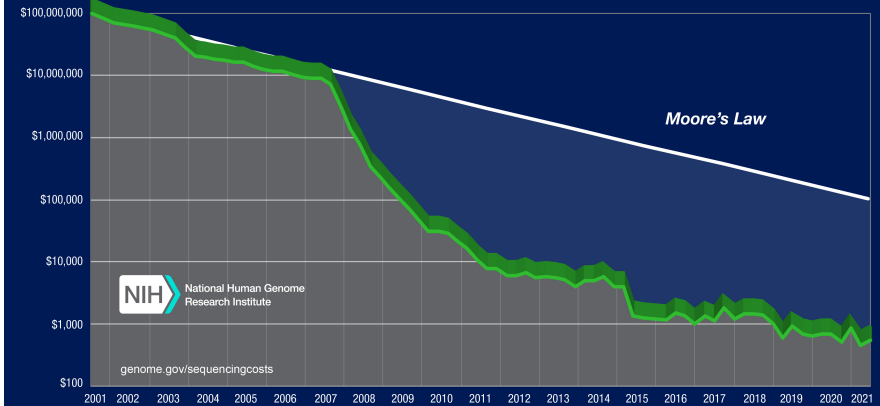




## Genome Assembly



### Cost per Human Genome



## High Throughput Sequencing



## Sequencing Output

Example applications:

- Sequencing a genome (DNA)
- Sequencing a transcriptome and gene expression studies (RNA)
- ChIP (chromatin immunoprecipitation)

Example platforms:

- 454
- [Illumina](#)
- SOLiD

- Hundreds of millions of sequencing reads, each ~200 nts in length
- We need to re-assemble the genome from these sequencing reads

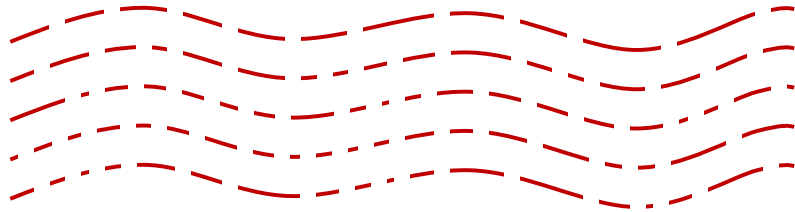
```

- ACGTAGTCGTAGTCGGTTACGATTGCGTACGTCAGTCTACGTTA
- CCAGCGTACTCTCGATGACGTGTACGTTACGATGACTGTAGTCAGTGT
- TTTGATCGTAGTGTCACTGAGCAACACCATTACTACTACTCTGGACATC
- TGGGGGATCGAGGATTCTAGTTATCGAGTGTCCGGGATTAATCGGATCGAA
- GGCATCATAGTCGATGCGATTACGATTAGCAGTGTCCGGGCTATACGTCGCGAT
- AGCCGGGTCGACGAGCGTACGGTCACTCGGATCGGATCGGATCGGATGAA
- TCGGTCGATCGAGTGTCTCCGGCTCTCGGAGCGCTAGGTAGAGAGCTG
- CTCCTCCAGCGTACTGCGATACGTTACGATCGGATCGGATCGGATCGGAT
- CCATTCGCTAGTCCGATCGATCGGATGATTGAGTCGCTAGTTACGATGT
- TACGGCGAGTTCGGCTAGTCACTGATCGGATCGGATCGGATCGGATCGG
- GAGCGTACGTCAGTCTACGTTCCGATCGGATCGGATCGGATCGGATCGG
- ACGCACGATCGATCTATGCATCGATGTCGATGTCGATGTCGATGTCGAT
- GTGCGTAGTCGTATATGCATAGCATGTTAGTCTAGCGTAGTCAGTCAGTAC
- ATCATCGGCGATAGTCTACGATGTTATATCTACGCGGGCCCACTTCGCAA
- CGAGTAGGAGTCGTAGTCGATGCGATGCGATCGGATCGGATCGGATCGGAT
- GACAGTCGCGAGTATGATAGCAGAGTTCGATGCGATGTCAGTCAGTGTAGCAGTATGTA
- CAACTTCCGCACTCTCCGGTCTCTCTCTCTAGATAGAGACTTACGATCG
- TCCGCGCATGCTAGTCCGGGCACTAGCGGACATCGGCGACACACGTCGATACAC
- AGACTGATATATCGGGCGGATGCGACTGTAGCTATATACGGCATCGGTC
- GGATCGATCATCGTGTAGTCGATGCGATGCGATGCGATGCGATGCGATGCGAT
- CGGATCGATCATCGTGTAGTCGATGCGATGCGATGCGATGCGATGCGATGCGAT
- TTATTATCGGCGAGTGTGCTAGTGTGATGCGATGCGATGCGATGCGATGCGAT
- AGTACGTAGTATCTGAGCGTCTCTCTACGGACATCGATGCGATGCGATGCGAT
- TTTATTACGACGATAGTGGCGATTCGCTATGCGATGCGATGCGATGCGATGCGAT
- GTGCGTAGTCGTATATGCATAGCATGTTAGTCTAGCGTAGTCAGTCAGTAC
- ATCATCGGCGATAGTCTACGATGTTATATCTACGCGGGCCCACTTCGCAA
- CGAGTAGGAGTCGTAGTCGATGCGATGCGATGCGATGCGATGCGATGCGATGCGAT
- GACAGTCGCGAGTATGATAGCAGAGTTCGATGCGATGTCAGTCAGTGTAGCAGTATGTA
- CAACTTCCGCACTCTCCGGTCTCTCTCTCTAGATAGAGACTTACGATCG
- TCCGCGCATGCTAGTCCGGGCACTAGCGGACATCGGCGACACACGTCGATACAC
- AGACTGATATATCGGGCGGATGCGACTGTAGCTATATACGGCATCGGTC
- GGATCGATCATCGTGTAGTCGATGCGATGCGATGCGATGCGATGCGATGCGAT
- CGGATCGATCATCGTGTAGTCGATGCGATGCGATGCGATGCGATGCGATGCGAT
- TTATTATCGGCGAGTGTGCTAGTGTGATGCGATGCGATGCGATGCGATGCGAT
- AGTACGTAGTATCTGAGCGTCTCTCTACGGACATCGATGCGATGCGATGCGAT
- TTTATTACGACGATAGTGGCGATTCGCTATGCGATGCGATGCGATGCGATGCGAT
- GTGCGTAGTCGTATATGCATAGCATGTTAGTCTAGCGTAGTCAGTCAGTAC
- ATCATCGGCGATAGTCTACGATGTTATATCTACGCGGGCCCACTTCGCAA
- CGAGTAGGAGTCGTAGTCGATGCGATGCGATGCGATGCGATGCGATGCGATGCGAT
- GACAGTCGCGAGTATGATAGCAGAGTTCGATGCGATGTCAGTCAGTGTAGCAGTATGTA
- CAACTTCCGCACTCTCCGGTCTCTCTCTCTAGATAGAGACTTACGATCG
- TCCGCGCATGCTAGTCCGGGCACTAGCGGACATCGGCGACACACGTCGATACAC
- AGACTGATATATCGGGCGGATGCGACTGTAGCTATATACGGCATCGGTC
- GGATCGATCATCGTGTAGTCGATGCGATGCGATGCGATGCGATGCGATGCGAT
- CGGATCGATCATCGTGTAGTCGATGCGATGCGATGCGATGCGATGCGATGCGAT
- TTATTATCGGCGAGTGTGCTAGTGTGATGCGATGCGATGCGATGCGATGCGAT
- AGTACGTAGTATCTGAGCGTCTCTCTACGGACATCGATGCGATGCGATGCGAT
- TTTATTACGACGATAGTGGCGATTCGCTATGCGATGCGATGCGATGCGATGCGAT

```



## DNA Sequencing

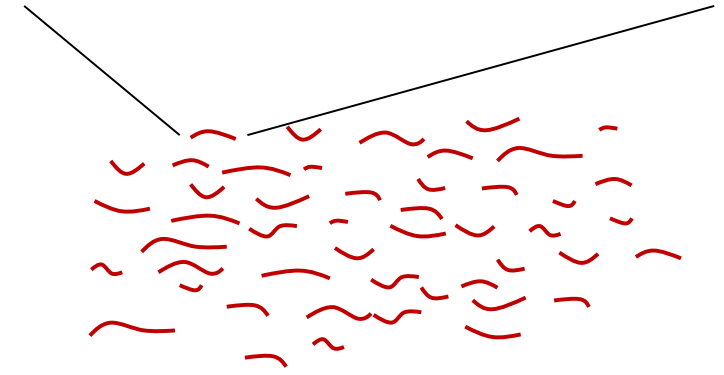


F-5



## DNA Sequencing

CGTAGTAGTCACAGTCTACGTATATGGGCTCAGCATATAGCGTATAGCGGACTTAGCCATCGTA



F-6



## DNA Sequencing

```

>CGTAGTAGTCACAGTCTACGTATATGGGCTCAGCATATAGCGTATAGCGGACTTAGCCATCG
>GCGTATAGTCTATATACGACTTATCGGCTCGGTCGCAGAGCAGATATATGCAGTTATATGCTAG
>CCTACGTTATATCGATACTACTAGTCTCGTCATGAGCGAGTAGATAGTATGACGAGCGACGATC
...
>CGATATTAGCCTAGCATCATTACGGCGAGACTCTCGGCTCGCTATATAGCGCTATAGCGAT
>CGGCTATAGCGCATATGCTCAGTAGCTATTAGCAGTATTACGATTATAGTCTCGGCGCATTAC
>TTTCGGGGATAAGTCTTCGTCTTATGCGACGATTATACGGCCGTATATTTGCATTTAGCATT
>GGCGTATGGCGGATATCGGCGGTCATAGCAGCCGATTAGGCTACGCCGATGCATCG
>CGCGATCGGCGGATCGCGTCAGTCGCGCAGTAGCGCGGCATAGTCGTATCGGCGCCG
>TGACAGAAGCTATAAGAGTCAGTAGATCTGAGTATTAGCATTATCGGCGCGATGCGCGATAACG
>GCGTATAGTCTATATACGACTTATCGGCTCGGTCGCAGAGCAGATATATGCAGTTATATGCTA
>CGCGATCGGCGGATCGCGTCAGTCGCGCAGTAGCGCGGCATAGTCGTATCGGCGCCGATCGC
>ATAGCAGCAGTATAGGATATGCTGCTCGTTCGACTATCATACTCGCTCGGCTAGCA
>TGACAGAAGCTATAAGAGTCAGTAGATCTGAGTATTAGCATTATCGGCGCGATGCGCGA
>CCTACGTTATATCGATACTACTAGTCTCGTCATGAGCGAGTAGATAGTATGACGAGCGACGATCC
>CGTAGTAGTCACAGTCTACGTATATGGGCTCAGCATATAGCGTATAGCGGACTTAGCCATCG
>TTTCGGGGATAAGTCTTCGTCTTATGCGACGATTATACGGCCGTATATTTGCATTTAGCATT
>GGCGTATGGCGGATATCGGCGGTCATAGCAGCCGATTAGGCTACGCCGATGCATCGTCGAGTA

```

F-7



## Assembly

```

>GGGCAGGC
>GGCTAGGG
>TAGGGCA
TAGGGCA
GGCTAGGG
GGGCAGGC
GGCTAGGGCAGGC

```

F-8



# Assembly

>GGGCAGGC  
 >GGCTAGGG  
 >TAGGGCA

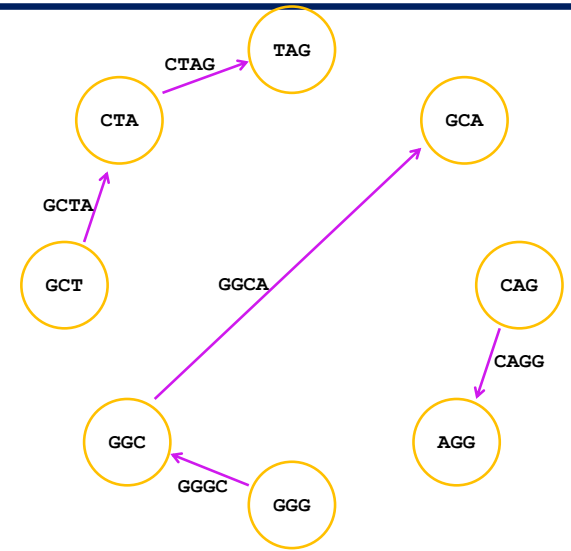
CTAG  
 GGGC  
 CAGG  
 GGCA  
 GCTA  
 AGGC  
 TAGG  
 AGGG  
 GCAG  
 GGCT  
 GGGC  
 TAGG  
 GGGC  
 AGGG



# Assembly

>GGGCAGGC  
 >GGCTAGGG  
 >TAGGGCA

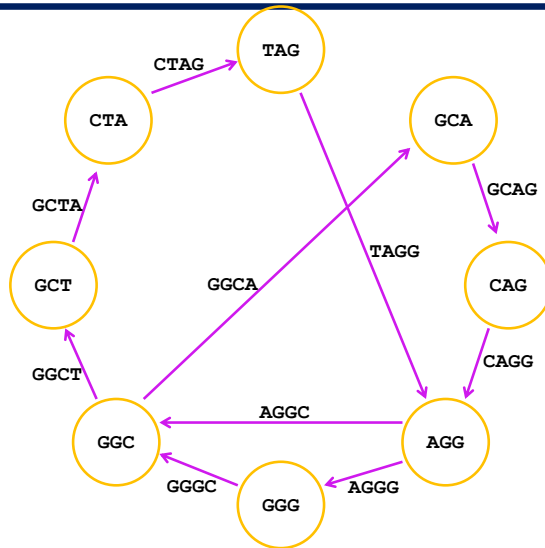
CTAG  
 GGGC  
 CAGG  
 GGCA  
 GCTA  
 AGGC  
 TAGG  
 AGGG  
 GCAG  
 GGCT  
 GGGC  
 TAGG  
 GGGC  
 AGGG



# Assembly

>GGGCAGGC  
 >GGCTAGGG  
 >TAGGGCA

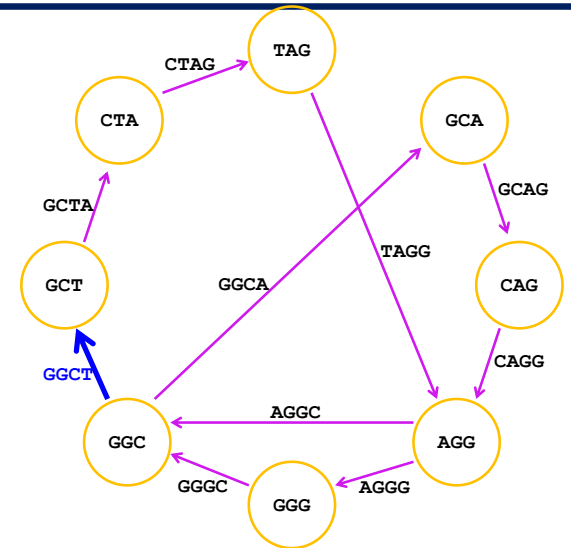
CTAG  
 GGGC  
 CAGG  
 GGCA  
 GCTA  
 AGGC  
 TAGG  
 AGGG  
 GCAG  
 GGCT  
 GGGC  
 TAGG  
 GGGC  
 AGGG



# Assembly

GGCT

>GGGCAGGC  
 >GGCTAGGG  
 >TAGGGCA

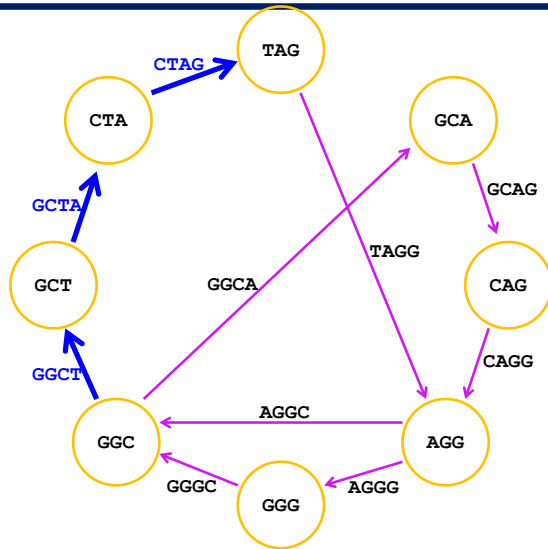




# Assembly

## GGCTAG

- >GGGCAGGC
- >GGCTAGGG
- >TAGGGCA



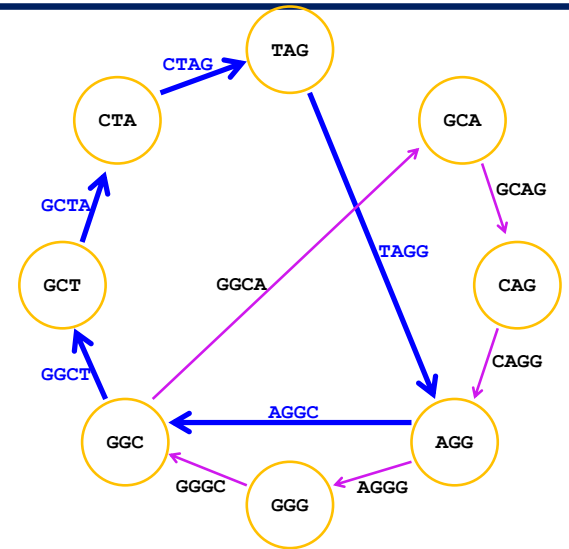
F-13



# Assembly

## GGCTAGGC

- >GGGCAGGC
- >GGCTAGGG
- >TAGGGCA



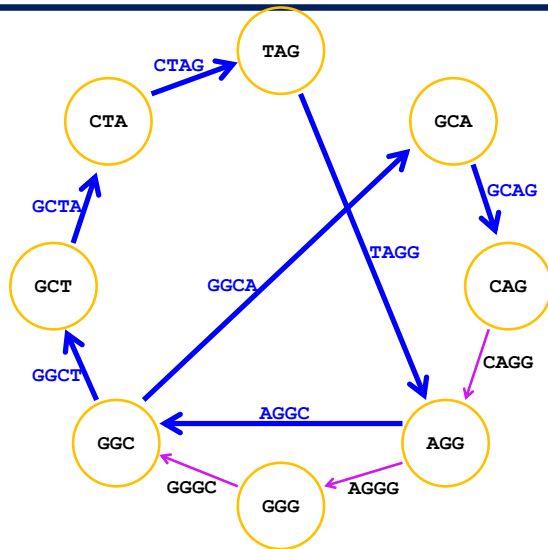
F-14



# Assembly

## GGCTAGGCAG

- >GGGCAGGC
- >GGCTAGGG
- >TAGGGCA



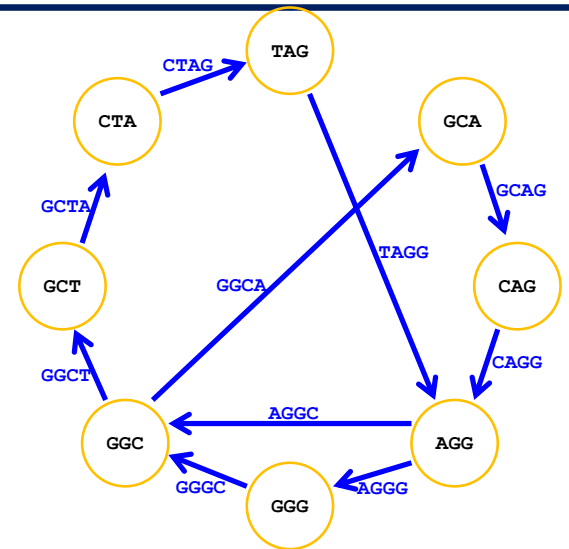
F-15



# Assembly

## GGCTAGGCAGGGC

- >GGGCAGGC
- >GGCTAGGG
- >TAGGGCA



F-16



## Challenges

- Repeats
- Sequencing errors

F-17



## Implementation

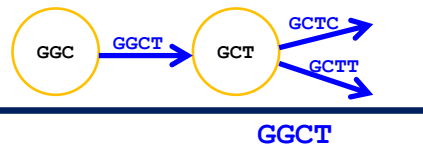
deBruijn graph can be implemented with a hash table

- Each entry in hash table corresponds to an edge in the graph (each key is a  $k$ -mer and each value is the number of occurrences of the  $k$ -mer).
- Nodes are stored implicitly.

F-18



## Implementation



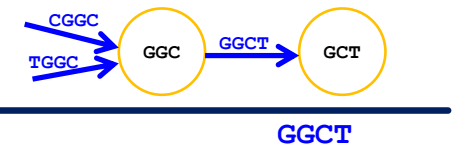
Assembly corresponds to Euler path through graph

- Genome sequence starts with any  $k$ -mer (edge in the graph)
- Repeatedly extend genome sequence **forward**, one nucleotide at a time, until no further extensions are possible
  - The genome sequence is extended and a nucleotide added to its end if there exists a nucleotide (A, C, G, or T) that can be added to the end of the  $k-1$  final nucleotides of the genome sequence to form a  $k$ -mer that is an edge in the graph.
  - If there are multiple individual nucleotides that can be added to the final  $k-1$  nucleotides in the genome sequence to form  $k$ -mer edges in the graph, then the nucleotide resulting in the  $k$ -mer edge with the largest number of occurrences is chosen.
  - Each time the genome sequence is extended by a nucleotide, the corresponding  $k$ -mer edge is removed from the graph.

F-19



## Implementation



Assembly corresponds to Euler path through graph

- Genome sequence starts with any  $k$ -mer (edge in the graph)
- Repeatedly extend genome sequence **backward**, one nucleotide at a time, until no further extensions are possible
  - The genome sequence is extended and a nucleotide added to its front if there exists a nucleotide (A, C, G, or T) that can be added to the front of the  $k-1$  first nucleotides of the genome sequence to form a  $k$ -mer that is an edge in the graph.
  - If there are multiple individual nucleotides that can be prepended to the first  $k-1$  nucleotides in the genome sequence to form  $k$ -mer edges in the graph, then the nucleotide resulting in the  $k$ -mer edge with the largest number of occurrences is chosen.
  - Each time the genome sequence is extended by a nucleotide, the corresponding  $k$ -mer edge is removed from the graph.

F-20