

Bioinformatics and Advanced Programming

Jan T. Kim



BCS Advanced Programming SG, 11 Dec 2014

Abstract

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

Bioinformatics can be defined strictly as the science of information in biological systems, or more broadly as developing and applying computational tools for analysing biological data. The biological processes that generate this information, particularly evolution, are highly complex, and therefore analysis of biological information is often computationally challenging. I will present the following selected topics and highlight the advanced computing challenges they involve, and also outline advances in the biosciences that have been enabled by tackling these challenges.

Many bioinformatics analyses are based on DNA sequences which today can be determined at very high volume through "Next Generation Sequencing" (NGS) techniques. As a result, the volume of publicly available sequence data has reached the range of petabytes. Searching this body of data requires highly optimised computational approaches, such as BLAST ("Basic Local Alignment Search Tool").

More recently, NGS methods that generate very large numbers of "short reads", i.e. strings of sequence . Central computational challenges resulting from these new technologies are "de novo assembly" of the original long sequence(s) from short reads, and mapping very large numbers of short reads to a known reference sequence.

Phylogeny analysis, i.e. reconstruction of ancestry relationships among species, is a classical field of bioinformatics which typically involves two steps, first a multiple alignment of the sequences is computed which subsequently is used to compute a tree. Computing multiple alignments is an optimisation problem that can only be approximately solved.

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

The Pirbright Institute



Preventing and controlling viral diseases



Core funding:



Project funding by BBSRC
and many others.



Outline

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

① Introduction

Molecular Biology Basics

Resolving the Phylogeny of Land Plants

Reconstructing Foot and Mouth Disease Transmission Trees

② Sequence Analysis

Pairwise Alignment

BLAST

③ “Next Generation” Sequencing Challenges

④ Multiple Sequence Alignment (MSA)

Bioinformatics: Definition(s)

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

- Scientific inquiry into **information in biological systems**.
- Computational analysis of **biological data**.
- Computer assisted mining of **biological literature**.

Outline

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

1 Introduction

Molecular Biology Basics

Resolving the Phylogeneny of Land Plants

Reconstructing Foot and Mouth Disease Transmission Trees

2 Sequence Analysis

Pairwise Alignment

BLAST

3 “Next Generation” Sequencing Challenges

4 Multiple Sequence Alignment (MSA)

DNA: Structure

Introduction

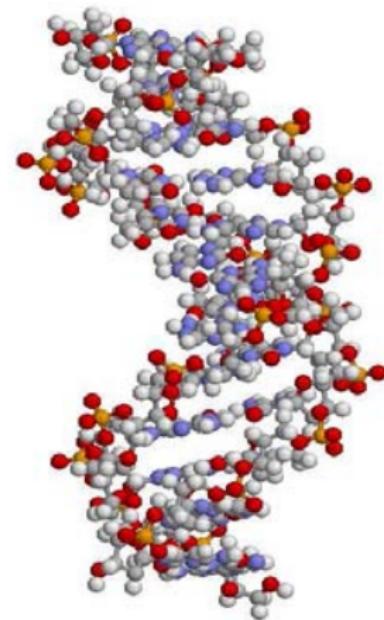
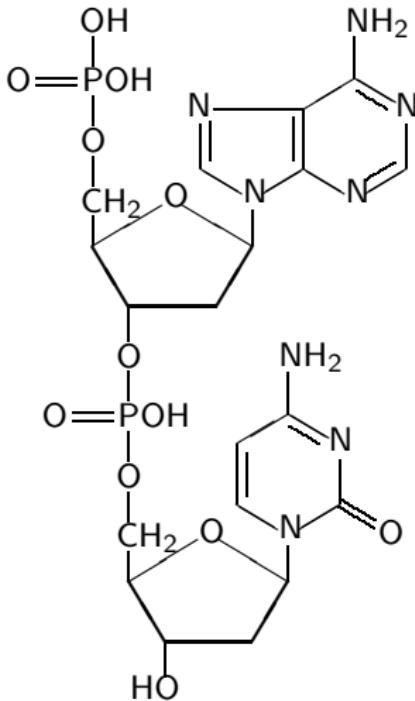
Mol. Bio. Basics

Plant Phylogeny

FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA



Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

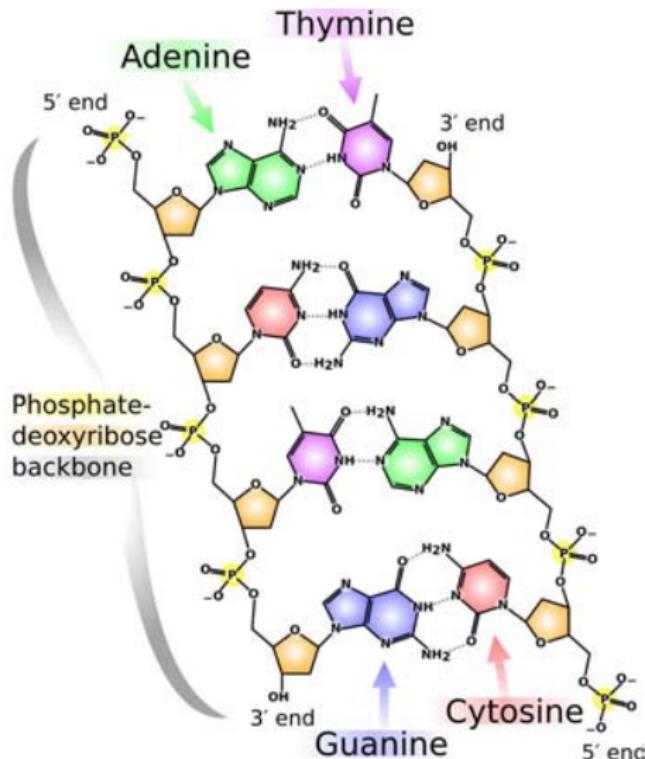
Transmission

Sequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Base Complementarity



http://commons.wikimedia.org/wiki/File:DNA_chemical_structure.svg

The “Central Dogma”

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

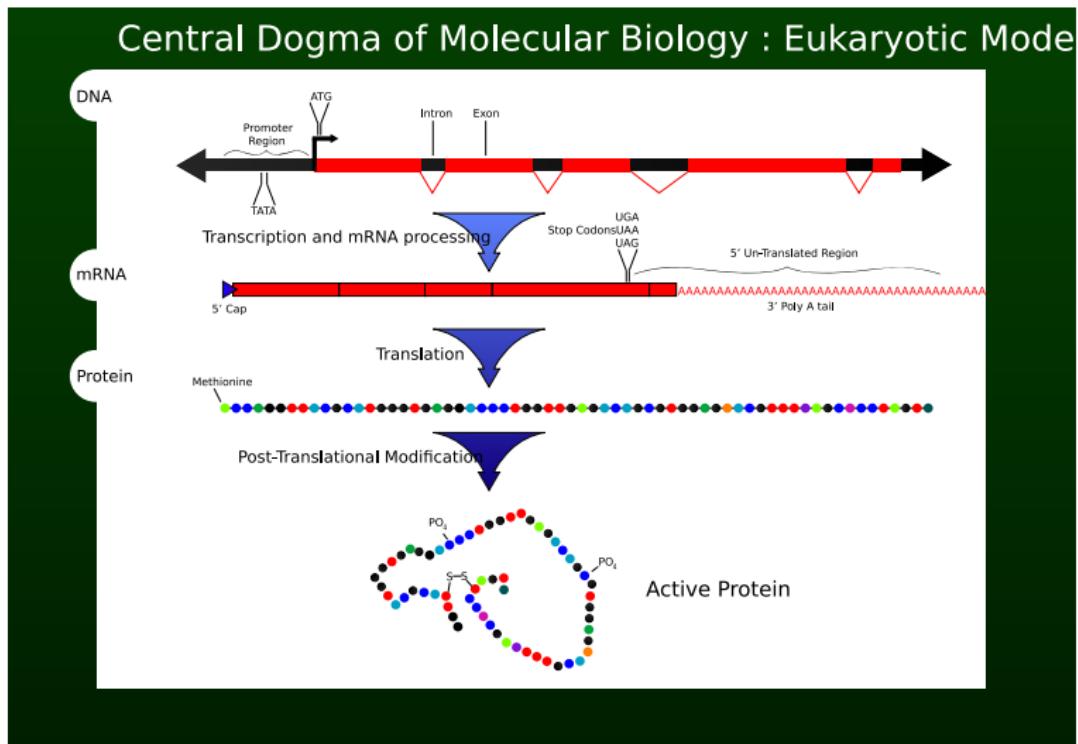
Pairwise

Alignment

BLAST

NGS

MSA



http://en.wikipedia.org/wiki/Central_dogma_of_molecular_biology

<http://en.wikipedia.org/wiki/File:Cdmb.svg>

The Success of Bioinformatics

- **The Object: Information in biological systems:**

In living systems, a dynamics of information has gained control over the dynamics of energy, which determines the behavior of most non-living systems.

[Langton, 1992]

- Genetic information is **digital**.

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

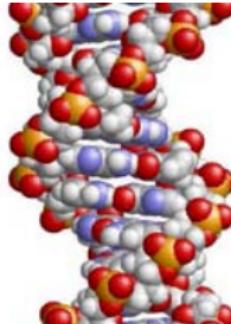
The Success of Bioinformatics

- **The Object: Information in biological systems:**

In living systems, a dynamics of information has gained control over the dynamics of energy, which determines the behavior of most non-living systems.

[Langton, 1992]

- Genetic information is **digital**.



Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

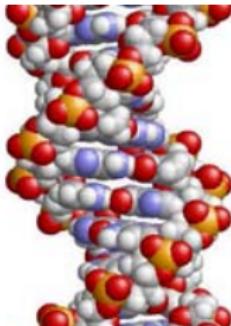
The Success of Bioinformatics

- **The Object: Information in biological systems:**

In living systems, a dynamics of information has gained control over the dynamics of energy, which determines the behavior of most non-living systems.

[Langton, 1992]

- Genetic information is **digital**.



```
TACCGTCAC  
CTACACCAT  
ACCTACATG  
TTCACATTAA
```

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Sequence Data Is Big Data

- NCBI-GenBank Flat File Release 204.0 (15 Oct 2014):

<ftp://ftp.ncbi.nih.gov/genbank/gbrel.txt>

- 178,322,253 loci,
- 181,563,676,918 bases.

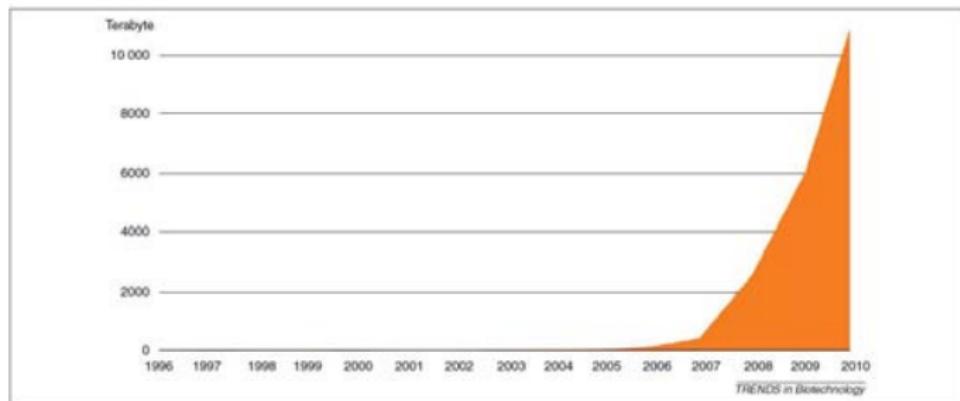


Figure 1. The scale of data growth. The chart shows the trend in storage capacity needed to store biological data at EMBL-EBI (a terabyte is a million million bytes).

[Crosswell and Thornton, 2012]

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

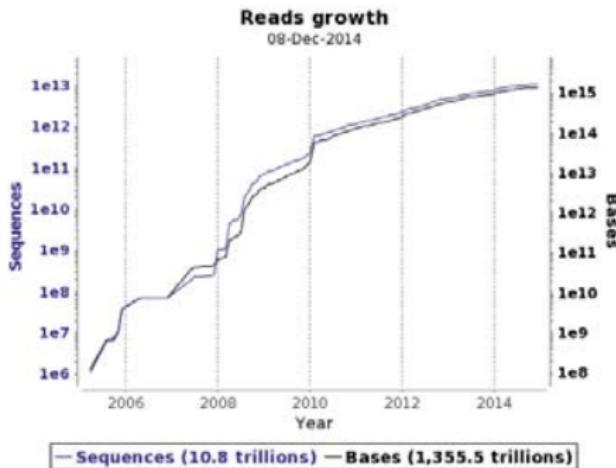
MSA

Sequence Data Is Big Data

- NCBI-GenBank Flat File Release 204.0 (15 Oct 2014):

<ftp://ftp.ncbi.nih.gov/genbank/gbrel.txt>

- 178,322,253 loci,
- 181,563,676,918 bases.



<http://www.ebi.ac.uk/ena/about/statistics>

Outline

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

1 Introduction

Molecular Biology Basics

Resolving the Phylogeneny of Land Plants

Reconstructing Foot and Mouth Disease Transmission Trees

2 Sequence Analysis

Pairwise Alignment

BLAST

3 “Next Generation” Sequencing Challenges

4 Multiple Sequence Alignment (MSA)

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA

Land Plant Phylogeny



Angiosperms
(flowering plants)



Gnetales



Gymnosperms

Introduction

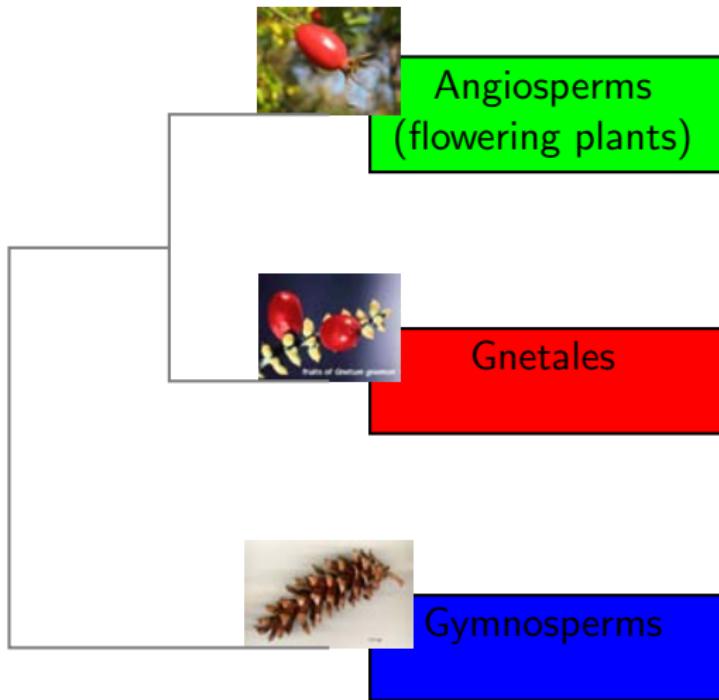
Mol. Bio. Basics

Plant PhylogenyFMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Land Plant Phylogeny



Introduction

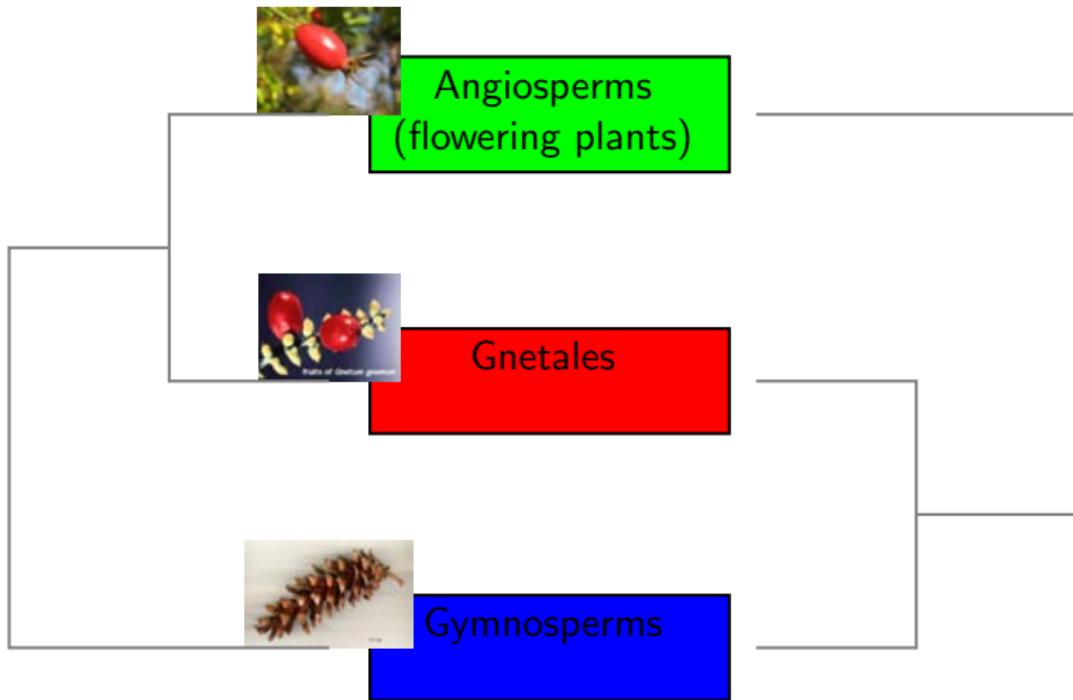
Mol. Bio. Basics

Plant PhylogenyFMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Land Plant Phylogeny



Phylogeny of MADS Proteins

Introduction

Mol. Bio. Basics

Plant Phylogeny

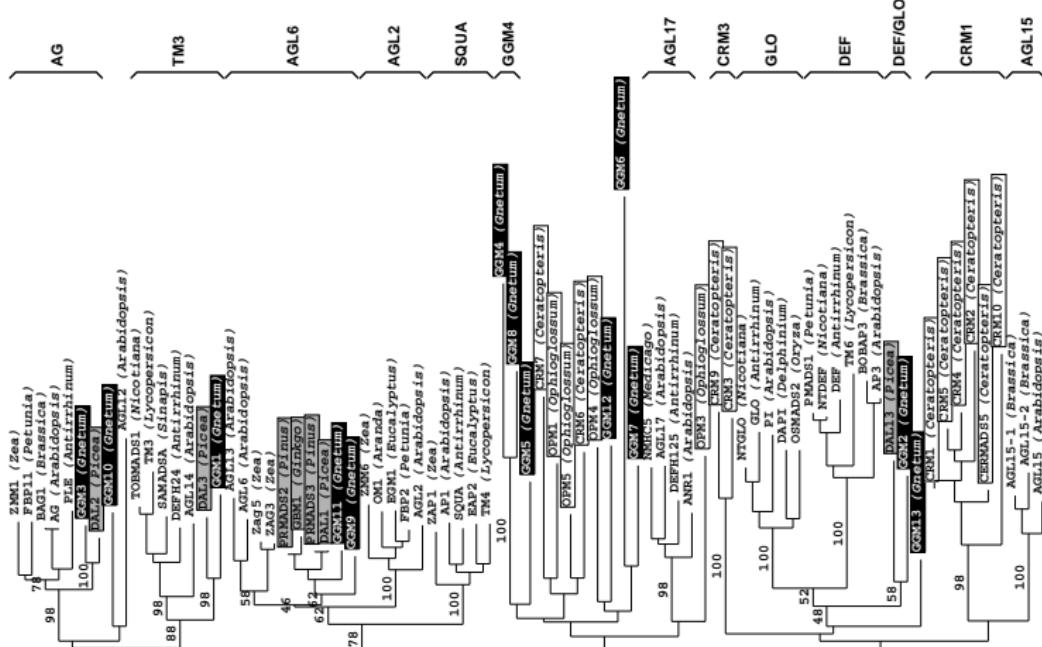
FMD
Transmission

Sequence Analysis

Pairwise Alignment BLAST

NGS

MSA



Angiosperms

Gnetales

Gymnosperms

The AG and AGL6 Subfamilies

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

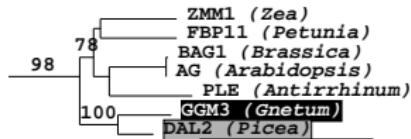
Pairwise

Alignment

BLAST

NGS

MSA



AG subfamily

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

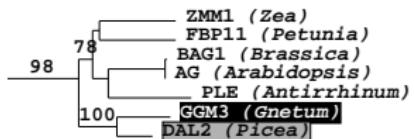
Alignment

BLAST

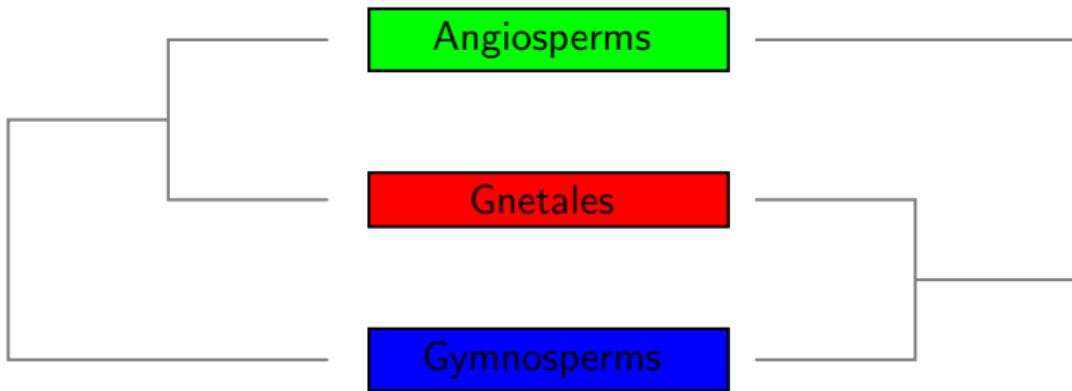
NGS

MSA

The AG and AGL6 Subfamilies



AG subfamily



Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

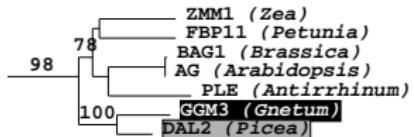
Alignment

BLAST

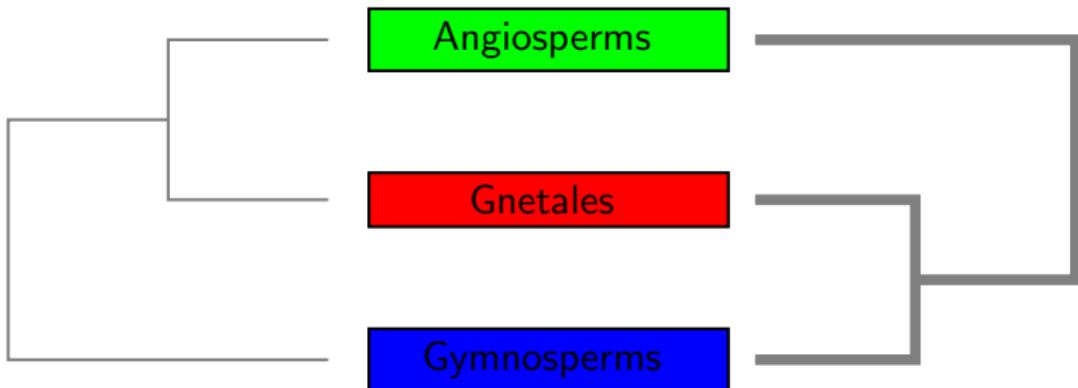
NGS

MSA

The AG and AGL6 Subfamilies



AG subfamily



Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

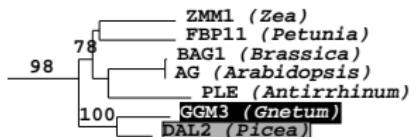
Alignment

BLAST

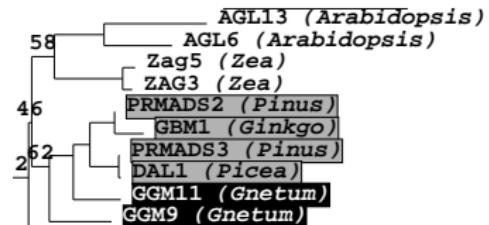
NGS

MSA

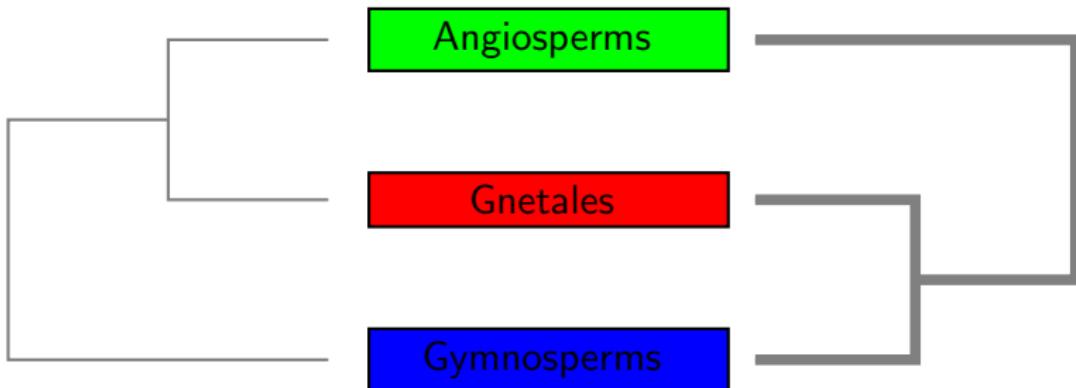
The AG and AGL6 Subfamilies



AG subfamily



AGL6 subfamily



Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

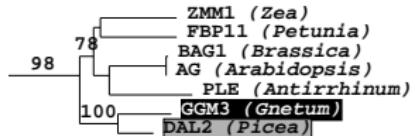
Alignment

BLAST

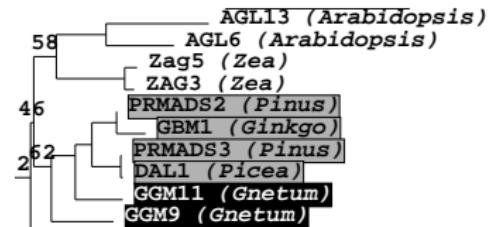
NGS

MSA

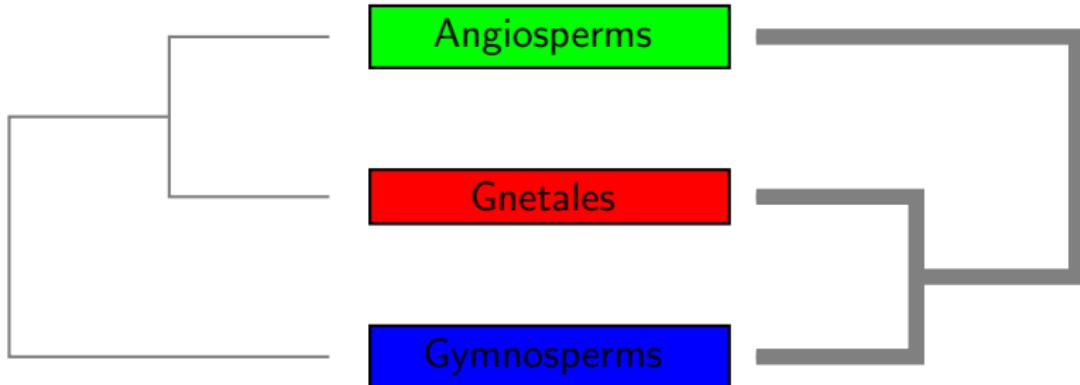
The AG and AGL6 Subfamilies



AG subfamily



AGL6 subfamily



Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

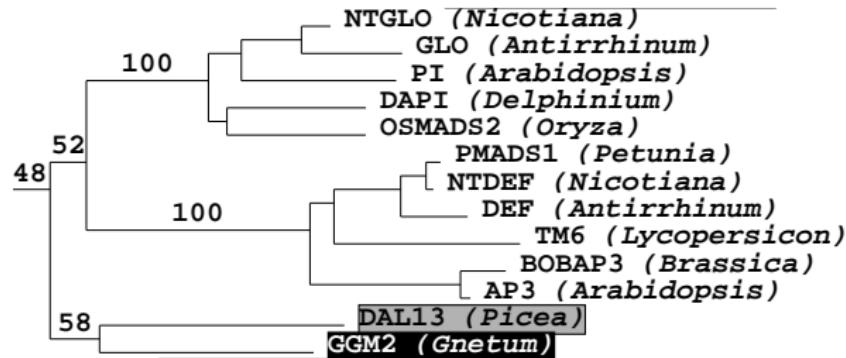
Sequence
AnalysisPairwise
Alignment

BLAST

NGS

MSA

The DEF and GLO Subfamilies



Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

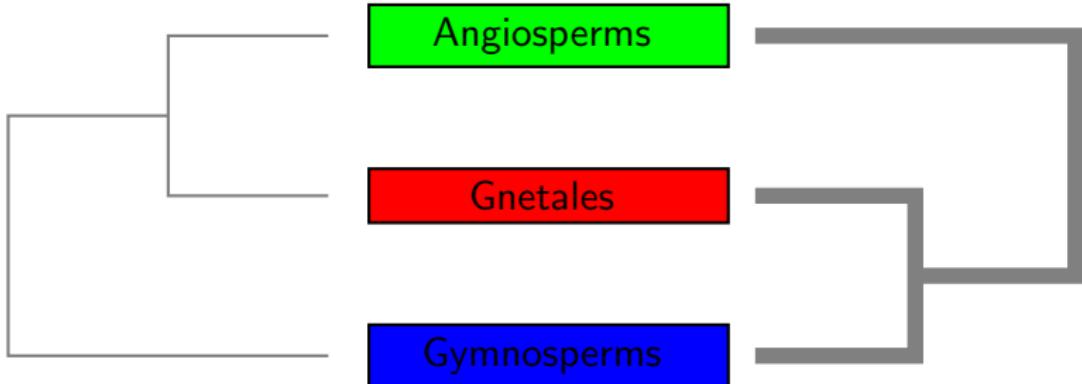
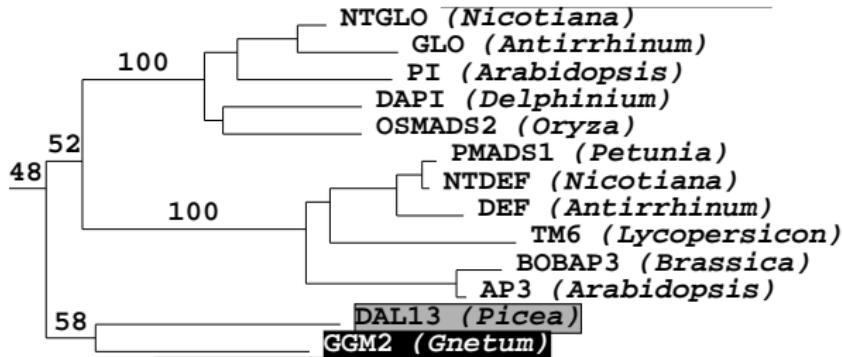
Sequence
AnalysisPairwise
Alignment

BLAST

NGS

MSA

The DEF and GLO Subfamilies



Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

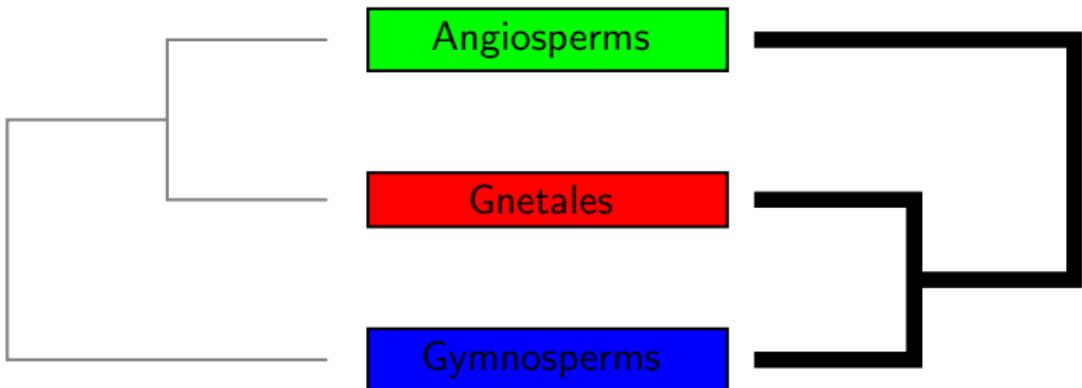
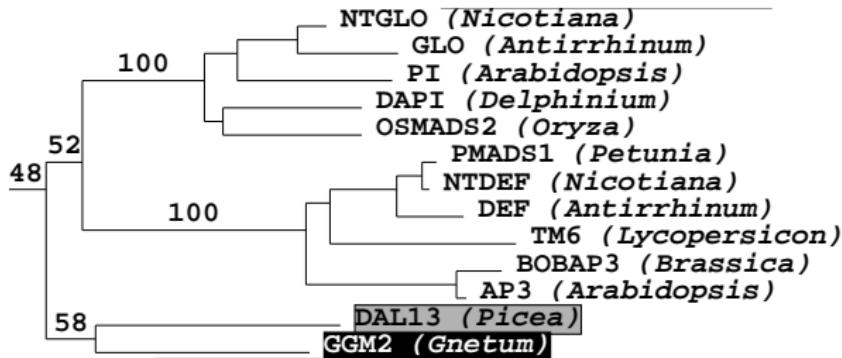
Alignment

BLAST

NGS

MSA

The DEF and GLO Subfamilies



Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

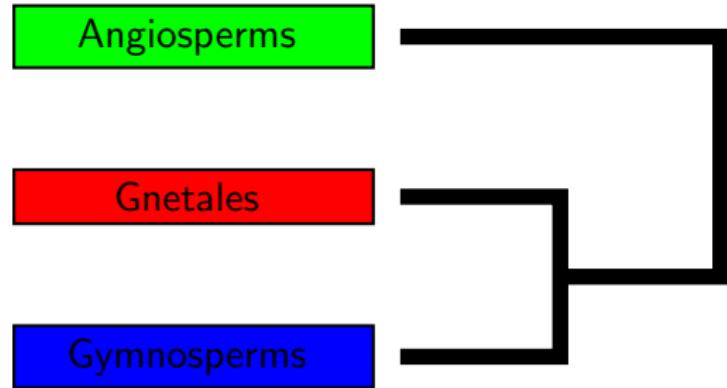
BLAST

NGS

MSA

Conclusion: Land Plant Phylogeny

MADS-Box Genes Reveal That Gnetales Are More Closely Related to Conifers than to Flowering Plants
[Winter et al., 1999].



Outline

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA

1 Introduction

Molecular Biology Basics

Resolving the Phylogeny of Land Plants

Reconstructing Foot and Mouth Disease Transmission Trees

2 Sequence Analysis

Pairwise Alignment

BLAST

3 “Next Generation” Sequencing Challenges

4 Multiple Sequence Alignment (MSA)

Transmission Trees

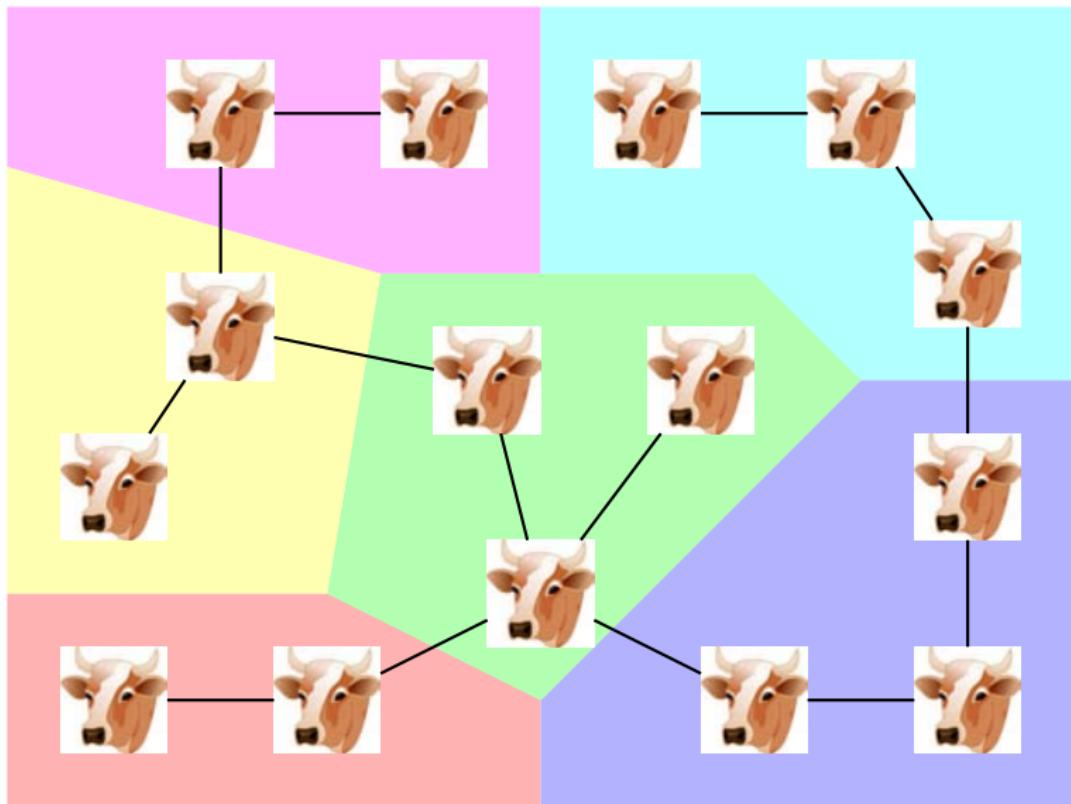
Introduction
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



Transmission Trees

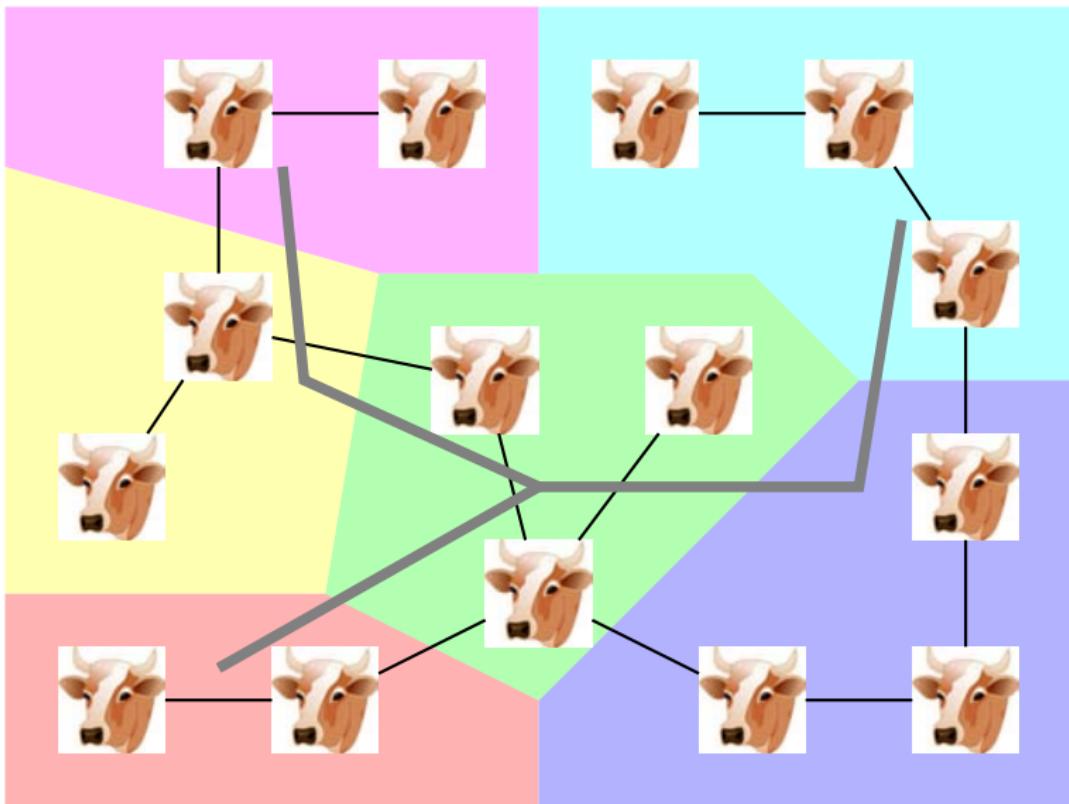
Introduction
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



Transmission Trees

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

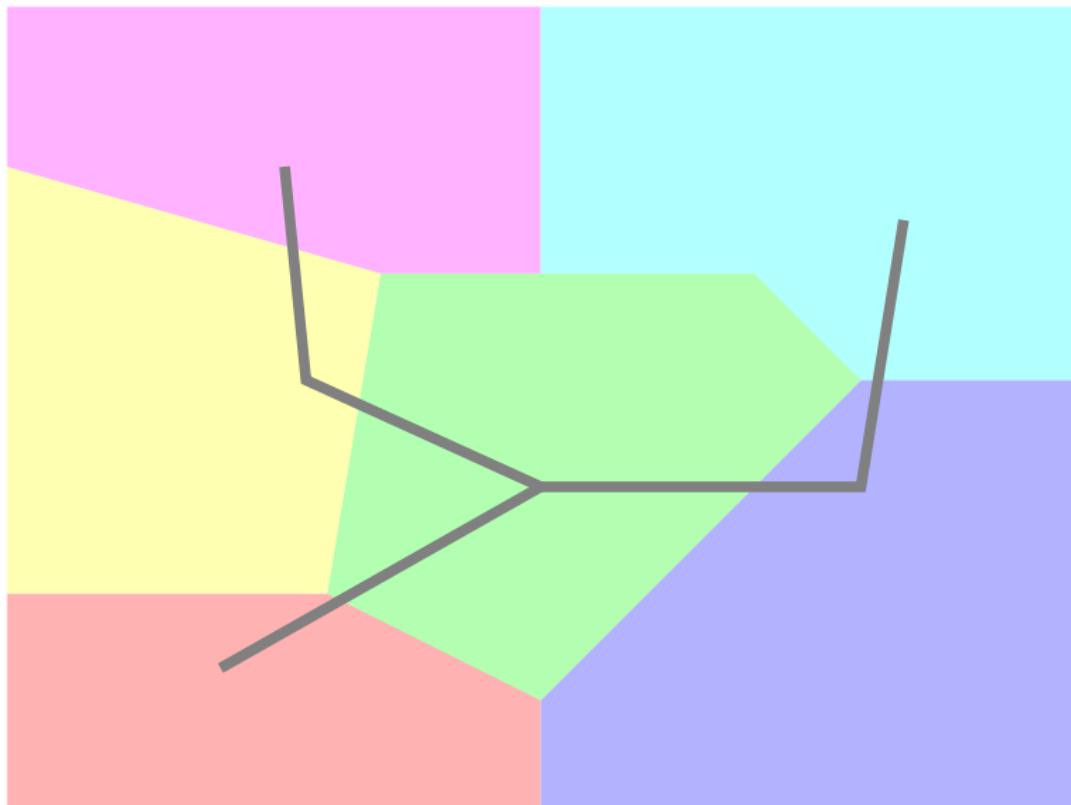
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



Transmission Trees

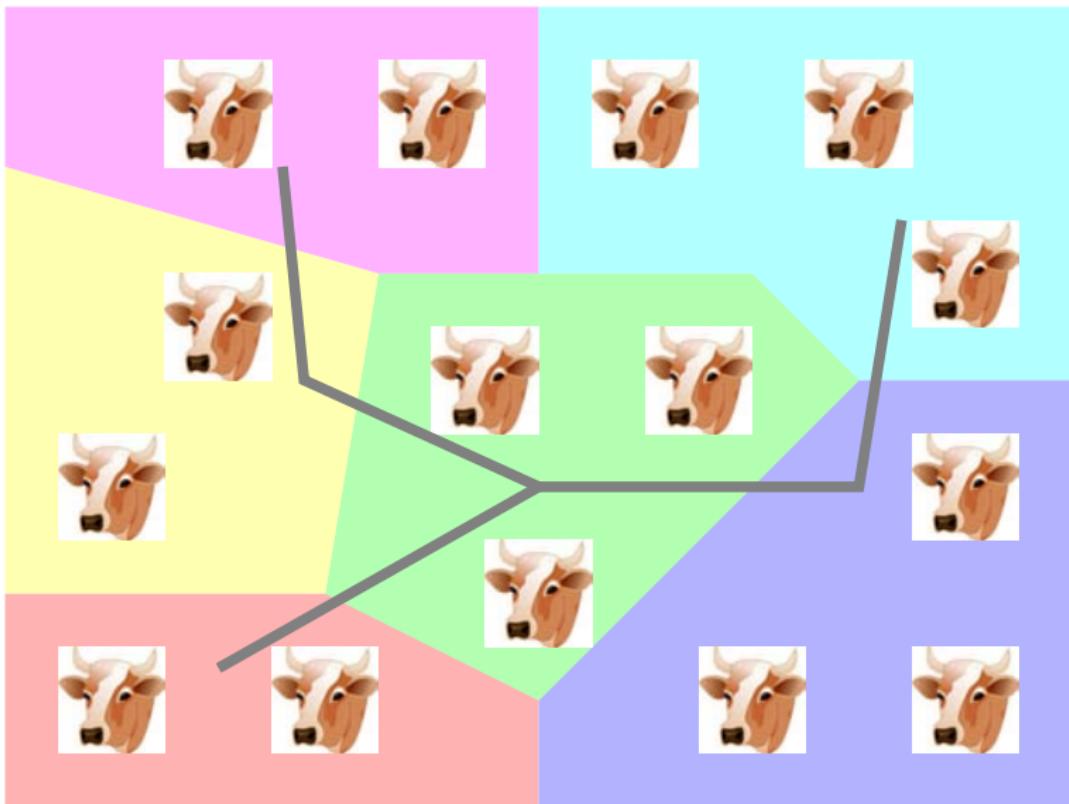
Introduction
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



Transmission Trees

Introduction

Mol. Bio. Basics

Plant Phylogeny

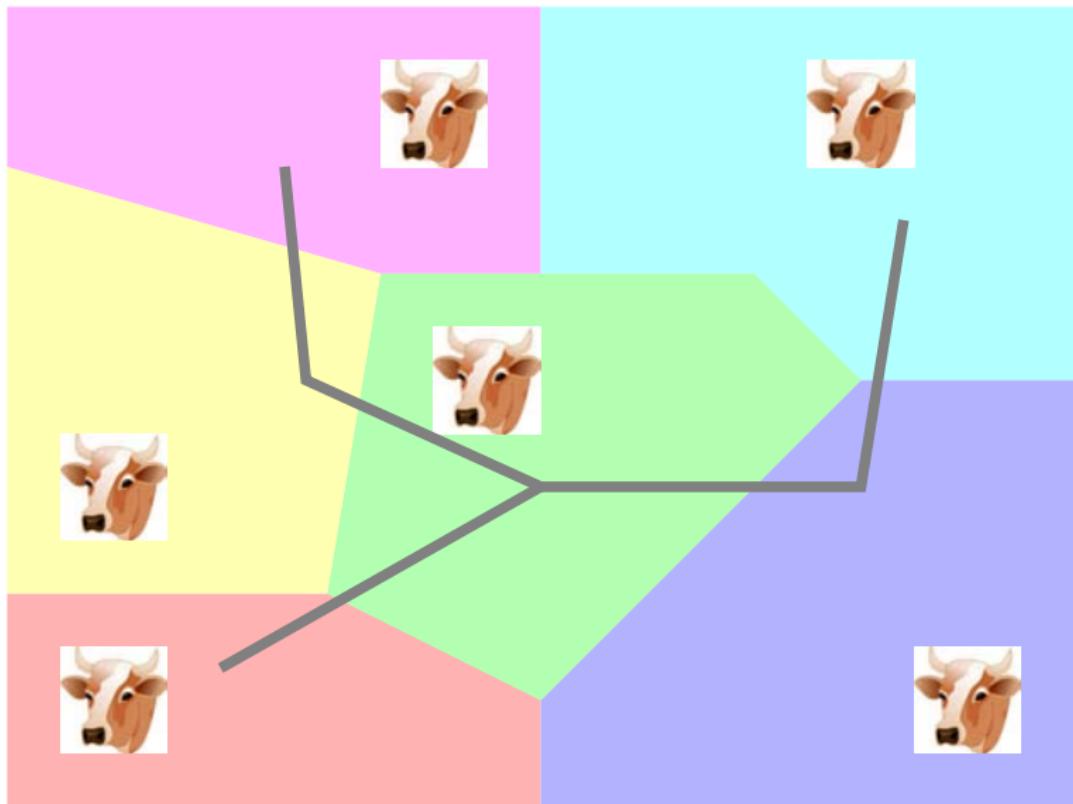
FMD

Transmission

Sequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA



Transmission Trees

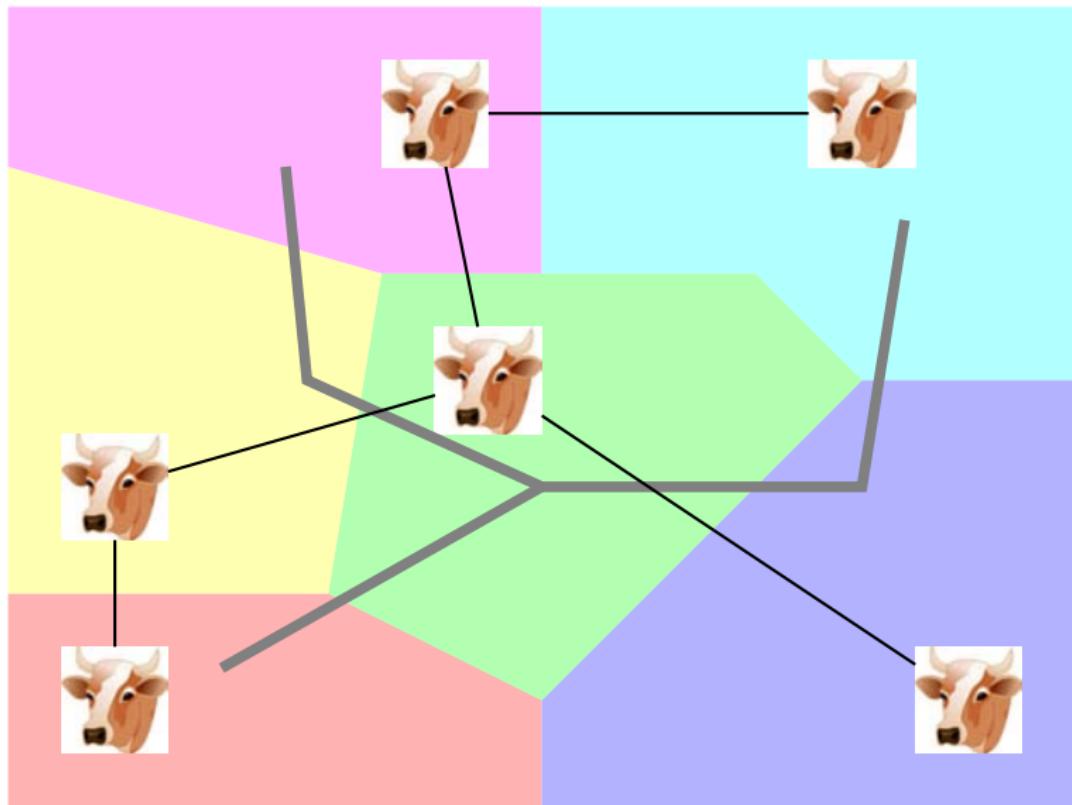
Introduction
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



Transmission Trees

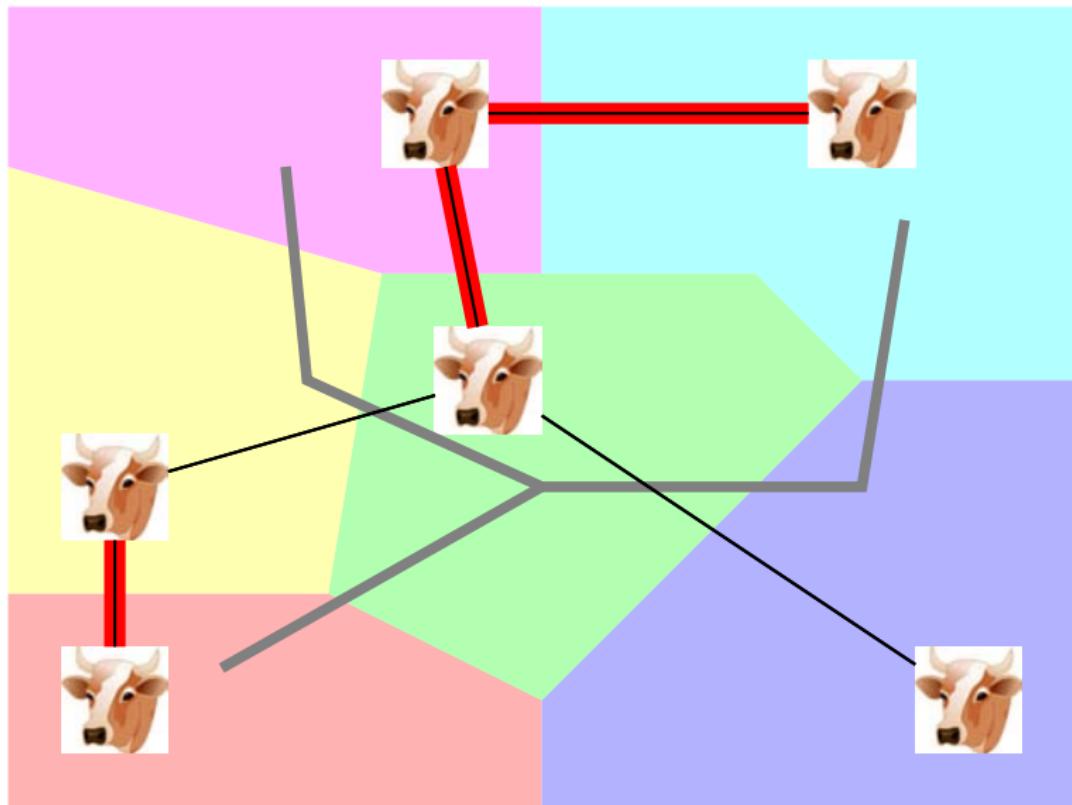
Introduction
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



Transmission Trees

Introduction

Mol. Bio. Basics

Plant Phylogeny

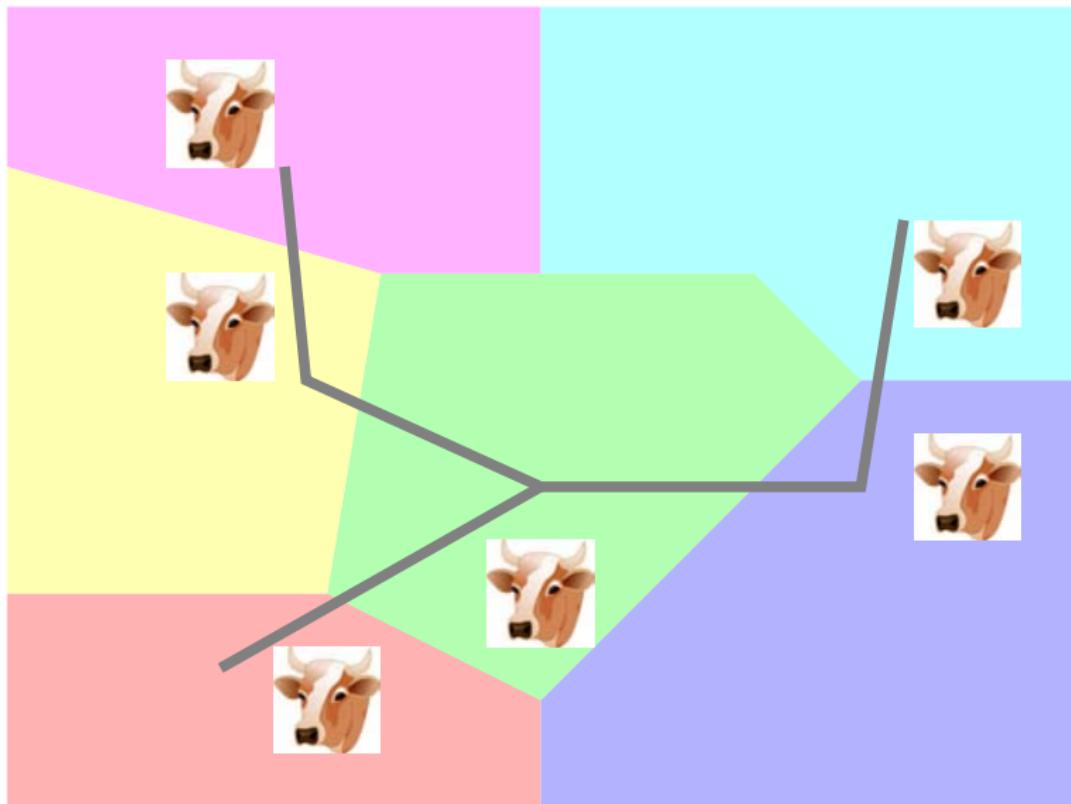
FMD

Transmission

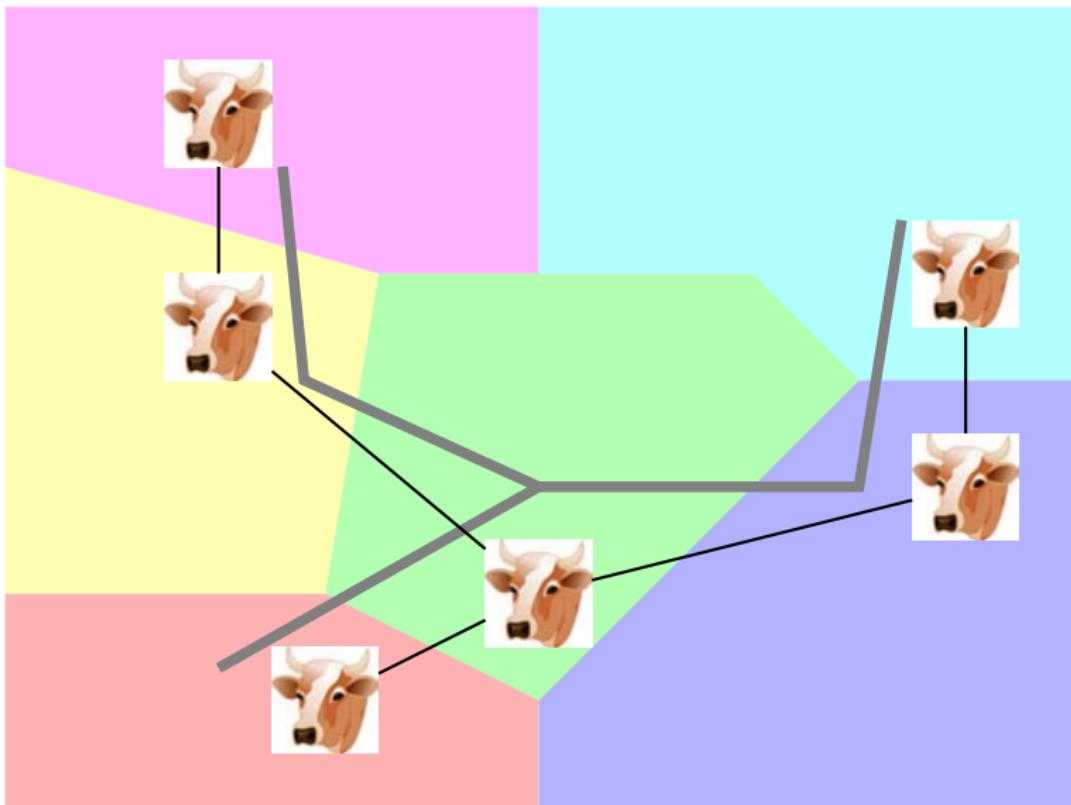
Sequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA



Transmission Trees



Consensus Sequence Data

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

site	sequences	
p0 (ancestor)	2	<ul style="list-style-type: none">• 2007 FMD outbreak in UK
p1b	6	
p2b	7	
p2c	3	
p3b	8	
p3c	2	
p4b	2	
p5	5	
p6b	3	
p7	8	
p8	1	
sum	47	
combinations	967680	

Introduction

Mol. Bio. Basics

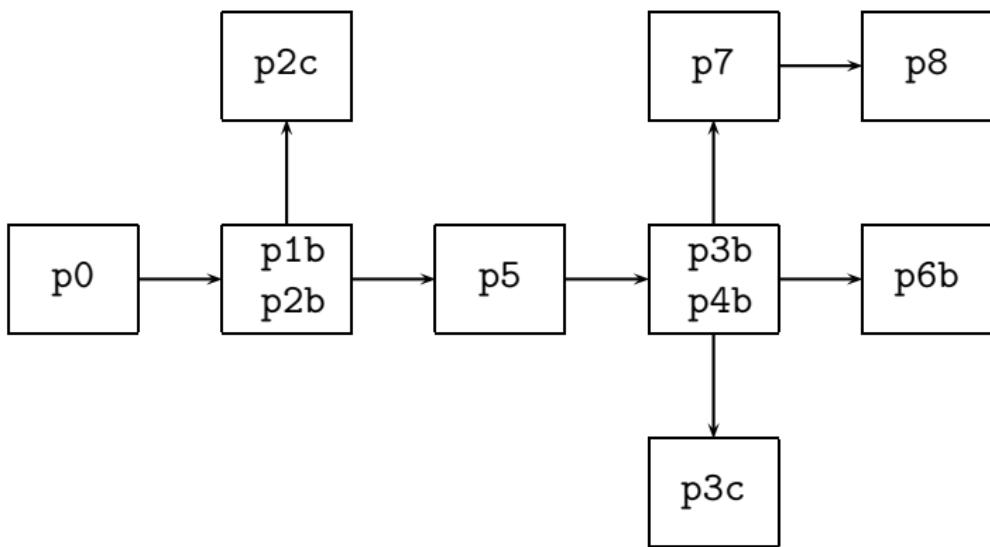
Plant Phylogeny

FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Reference Transmission Tree

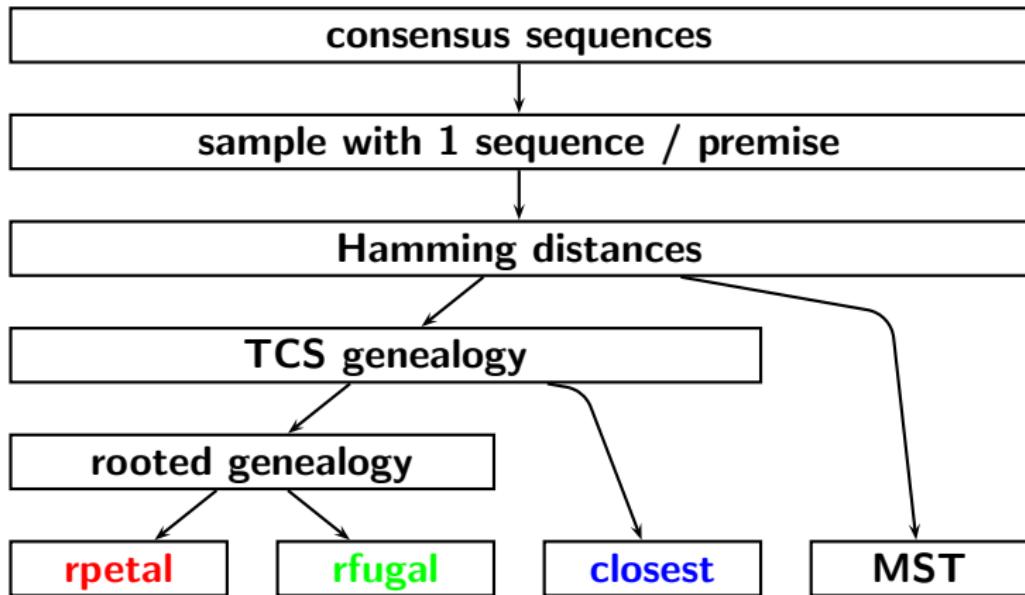


Based on a TCS genealogy of all 47 samples, and additional background knowledge.

Flowchart

Introduction
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence Analysis
Pairwise Alignment
BLAST
NGS
MSA



Analysis carried out for **1000 random samples** containing one consensus sequence from each site.

Introduction

Mol. Bio. Basics

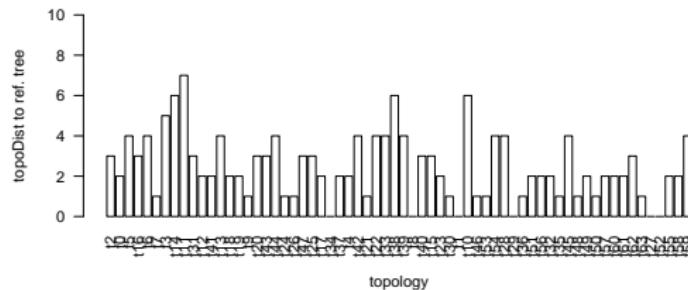
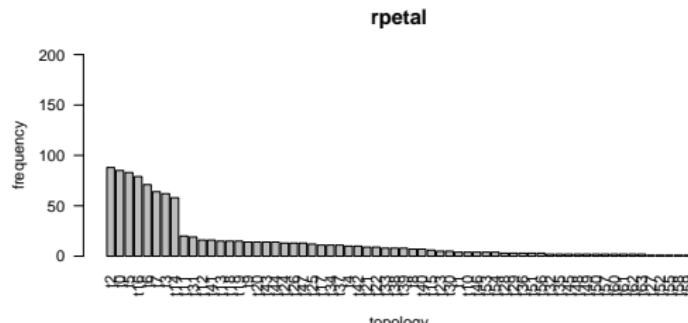
Plant Phylogeny

FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Results: Tree Topology



radipetal: branch nodes merged towards root (p_0)

Introduction

Mol. Bio. Basics
Plant Phylogeny

FMD
Transmission

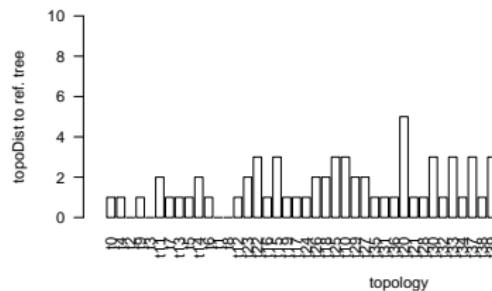
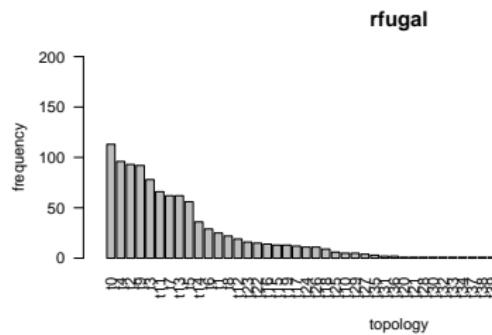
Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA

Results: Tree Topology



radifugal: branch nodes merged away from root (p0)

Introduction

Mol. Bio. Basics

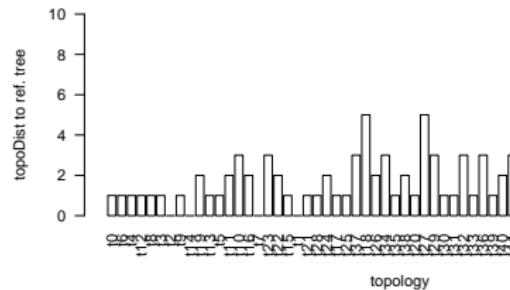
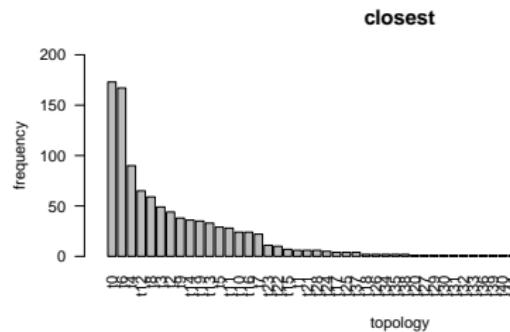
Plant Phylogeny

FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Results: Tree Topology



closest: branch nodes merged towards closest premise

Introduction

Mol. Bio. Basics
Plant Phylogeny

FMD
Transmission

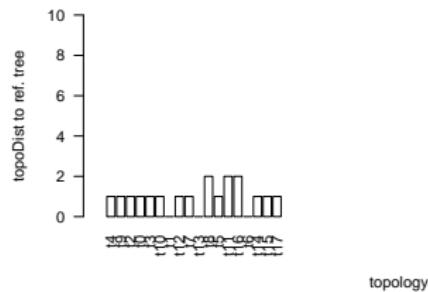
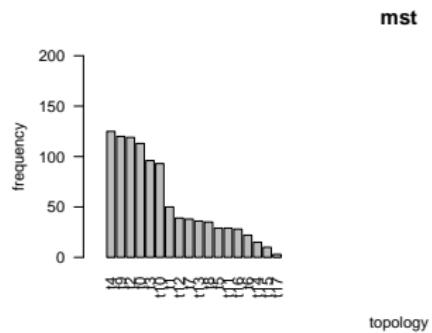
Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA

Results: Tree Topology



MST

Summary: FMD Transmission Trees

- Algorithms for constructing transmission trees
 - from TCS genealogies: **radipetal**, **radifugal**, **closest**,
 - minimum spanning tree (MST).
- Comparison based on the 2007 outbreak.
 - **closest** provides TCS based transmission trees best precision.
 - MST provides even marginally better precision.
- **Outlook:**
 - Try more sophisticated distance measures.
 - Include further transmission tree reconstruction methods.
 - Larger data sets, NGS “beyond the consensus”
[Wright et al., 2011]

Outline

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
Analysis

Pairwise
Alignment

BLAST

NGS

MSA

1 Introduction

Molecular Biology Basics

Resolving the Phylogeneny of Land Plants

Reconstructing Foot and Mouth Disease Transmission Trees

2 Sequence Analysis

Pairwise Alignment

BLAST

3 “Next Generation” Sequencing Challenges

4 Multiple Sequence Alignment (MSA)

Outline

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
Analysis

**Pairwise
Alignment**

BLAST

NGS

MSA

① Introduction

Molecular Biology Basics

Resolving the Phylogeny of Land Plants

Reconstructing Foot and Mouth Disease Transmission Trees

② Sequence Analysis

Pairwise Alignment

BLAST

③ “Next Generation” Sequencing Challenges

④ Multiple Sequence Alignment (MSA)

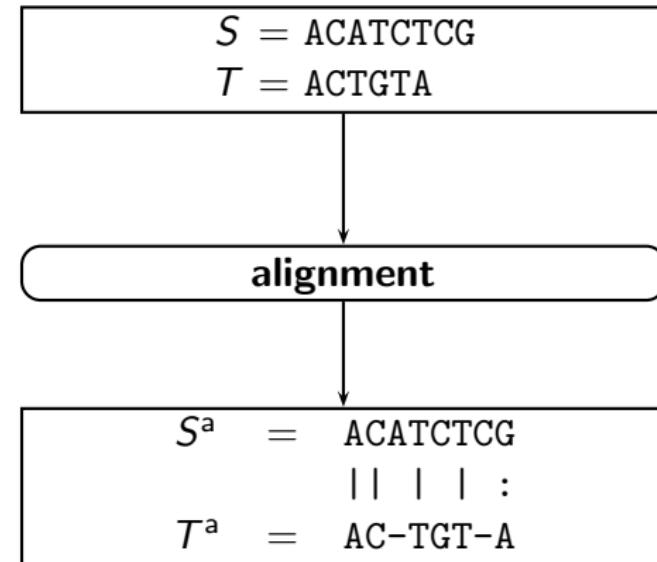
Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Pairwise Alignment: Idea



Formal Definition

- **Extend** sequences S and T by inserting **gaps** to S^a and T^a .
 - aligned sequences have **equal length**: $|S^a| = |T^a|$
 - gaps cannot be paired with gaps
- Biological background: **homology**, symbols in a column should derive from same **common ancestor**.
- **Match**: column with **equal** symbols in S^a and T^a .
- **Indel**: column with a **gap symbol** in S^a or T^a .
- **Mismatch**: column with **different** symbols (non-gap) in S^a and T^a .

ACATCTCG
AC-TGT-A

Formal Definition

- **Extend** sequences S and T by inserting **gaps** to S^a and T^a .
 - aligned sequences have **equal length**: $|S^a| = |T^a|$
 - gaps cannot be paired with gaps
- Biological background: **homology**, symbols in a column should derive from same **common ancestor**.
- **Match**: column with **equal** symbols in S^a and T^a .
- **Indel**: column with a gap symbol in S^a or T^a .
- **Mismatch**: column with **different** symbols (non-gap) in S^a and T^a .

ACATCTCG
AC-TGT-A

Formal Definition

- **Extend** sequences S and T by inserting **gaps** to S^a and T^a .
 - aligned sequences have **equal length**: $|S^a| = |T^a|$
 - gaps cannot be paired with gaps
- Biological background: **homology**, symbols in a column should derive from same **common ancestor**.
- **Match**: column with **equal** symbols in S^a and T^a .
- **Indel**: column with a **gap symbol** in S^a or T^a .
- **Mismatch**: column with **different** symbols (non-gap) in S^a and T^a .

ACATCT**CG**
AC-TGT-A

Formal Definition

- **Extend** sequences S and T by inserting **gaps** to S^a and T^a .
 - aligned sequences have **equal length**: $|S^a| = |T^a|$
 - gaps cannot be paired with gaps
- Biological background: **homology**, symbols in a column should derive from same **common ancestor**.
- **Match**: column with **equal** symbols in S^a and T^a .
- **Indel**: column with a **gap symbol** in S^a or T^a .
- **Mismatch**: column with **different** symbols (non-gap) in S^a and T^a .

ACAT**CTCG**
AC-TGT-A

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
Analysis

Pairwise

Alignment

BLAST

NGS

MSA

Scoring of Alignments

The score $m(k)$ of column k is

- the **space penalty** $m(k) = -g$, if one symbol is the gap symbol, here: $g = 2$),
- otherwise the **pair score** $m(k) = \mu(s^a(k), t^a(k))$, here

$$\mu(x, y) = \begin{cases} 1, & \text{if } x = y, \\ -1, & \text{otherwise} \end{cases}$$

$S^a =$	A	C	A	T	C	T	C	G	
$T^a =$	A	C	-	T	G	T	-	A	
<i>score:</i>	+1	+1	-2	+1	-1	+1	-2	-1	= -2

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Optimal Alignments

- **Objective:** Find the alignment with **maximal** score.
- **Problem:** The number of alignments is

$$\binom{|S|+|T|}{|S|} \cdot \binom{|S|+|T|}{|T|}$$

- Trying out all alignments is **impossible**.
 - Recursion results in trying out all alignments.

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

Optimal Alignments

- **Observation:** A prefix alignment of an optimal alignment is optimal (as well).
 - Otherwise, a **contradiction** results: The optimal alignment could be improved by changing the prefix.
- **Dynamic programming:** Tabulate optimal scores of prefix alignments

ACATCTCG

AC-TGT-A

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

Optimal Alignments

- **Observation:** A prefix alignment of an optimal alignment is optimal (as well).
 - Otherwise, a **contradiction** results: The optimal alignment could be improved by changing the prefix.
- **Dynamic programming:** Tabulate optimal scores of prefix alignments

ACATCTCG

AC-TGT-A

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

Optimal Alignments

- **Observation:** A prefix alignment of an optimal alignment is optimal (as well).
 - Otherwise, a **contradiction** results: The optimal alignment could be improved by changing the prefix.
- **Dynamic programming:** Tabulate optimal scores of prefix alignments

ACATCTCG

AC-TGT-A

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

Optimal Alignments

- **Observation:** A prefix alignment of an optimal alignment is optimal (as well).
 - Otherwise, a **contradiction** results: The optimal alignment could be improved by changing the prefix.
- **Dynamic programming:** Tabulate optimal scores of prefix alignments

ACATCTCG

AC-TGT-A

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

Optimal Alignments

- **Observation:** A prefix alignment of an optimal alignment is optimal (as well).
 - Otherwise, a **contradiction** results: The optimal alignment could be improved by changing the prefix.
- **Dynamic programming:** Tabulate optimal scores of prefix alignments

ACATCTCG

AC-TGT-A

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
AnalysisPairwise
Alignment

BLAST

NGS

MSA

Optimal Alignments

- **Observation:** A prefix alignment of an optimal alignment is optimal (as well).
 - Otherwise, a **contradiction** results: The optimal alignment could be improved by changing the prefix.
- **Dynamic programming:** Tabulate optimal scores of prefix alignments

ACATCTCG

AC-TGT-A

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

Optimal Alignments

- **Observation:** A prefix alignment of an optimal alignment is optimal (as well).
 - Otherwise, a **contradiction** results: The optimal alignment could be improved by changing the prefix.
- **Dynamic programming:** Tabulate optimal scores of prefix alignments

ACATCTCG

AC-TGT-A

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

Optimal Alignments

- **Observation:** A prefix alignment of an optimal alignment is optimal (as well).
 - Otherwise, a **contradiction** results: The optimal alignment could be improved by changing the prefix.
- **Dynamic programming:** Tabulate optimal scores of prefix alignments

ACATCTCG

AC-TGT-A

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

Optimal Alignments

- **Observation:** A prefix alignment of an optimal alignment is optimal (as well).
 - Otherwise, a **contradiction** results: The optimal alignment could be improved by changing the prefix.
- **Dynamic programming:** Tabulate optimal scores of prefix alignments

ACATCTCG

AC-TGT-A

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Table of Prefix-Alignment Scores

	-	A	G	A	C
-	0.0	-2.0	-4.0	-6.0	-8.0
A	-2.0	1.0	-1.0	-3.0	-5.0
G	-4.0	-1.0	2.0	0.0	-2.0
C	-6.0	-3.0	0.0	1.0	1.0

The **optimal** alignment score is 1.0.

☞ Notice $O(n^2)$ complexity.

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

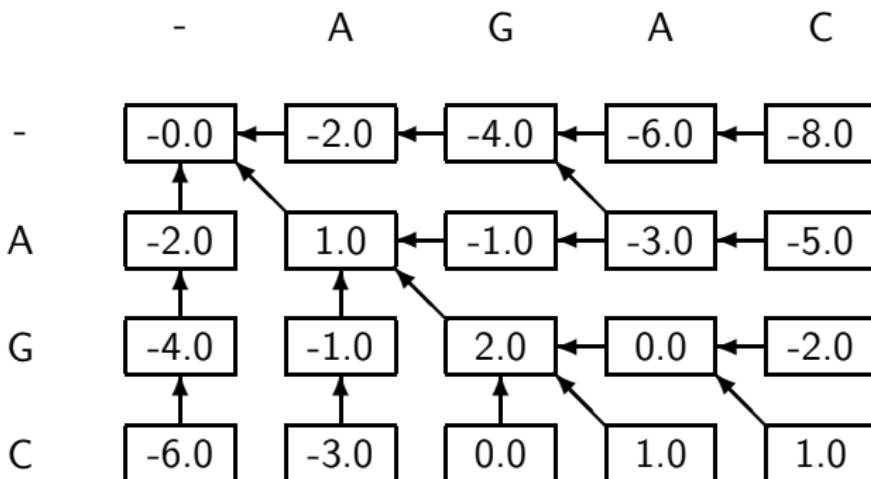
Sequence
AnalysisPairwise
Alignment

BLAST

NGS

MSA

Backtracking the Alignment



AG-C

AGAC

Outline

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

① Introduction

Molecular Biology Basics

Resolving the Phylogeneny of Land Plants

Reconstructing Foot and Mouth Disease Transmission Trees

② Sequence Analysis

Pairwise Alignment

BLAST

③ “Next Generation” Sequencing Challenges

④ Multiple Sequence Alignment (MSA)

BLAST: Basic Local Alignment Search Tool

- **Objective:** Given a **query sequence**, find **similar sequences** in a database.
- Size of database prohibits pairwise alignment of query to all entries.
- Algorithm outline [Altschul et al., 1997]:
 - ① **Scan** for **hits**, i.e. gapless short word alignments exceeding a threshold score.
 - ② **Extend** hits maximally to obtain **HSPs** (high scoring pairs).
 - ③ **Combine** HSPs to (gapped) alignments.
- E-values indicate **expected number** of HSPs with given score.
 - Interesting HSPs have $E \ll 1$.

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

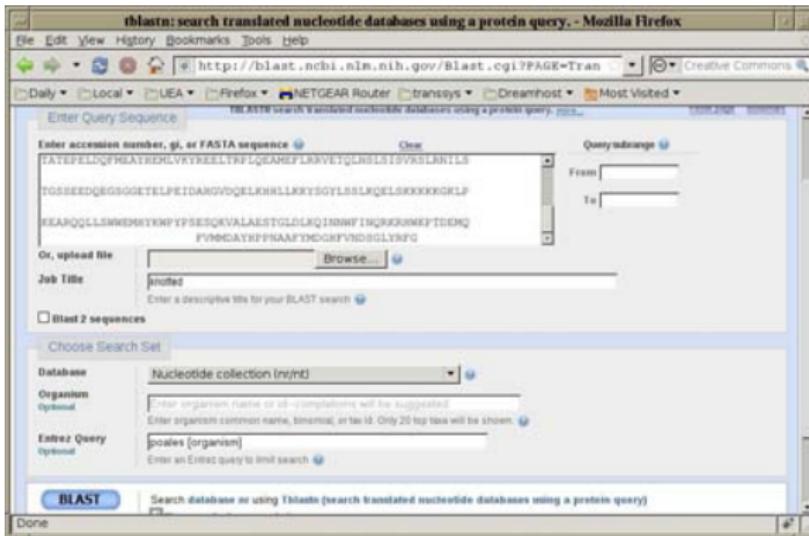
Sequence
AnalysisPairwise
Alignment

BLAST

NGS

MSA

BLAST: Search Engine for Sequences



<http://blast.ncbi.nlm.nih.gov/>

BLAST: Search Engine for Sequences

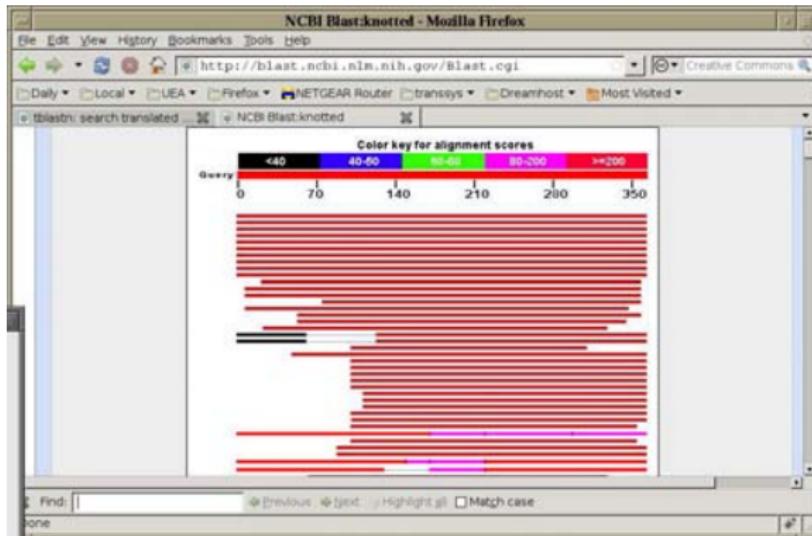
Introduction
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



<http://blast.ncbi.nlm.nih.gov/>

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
AnalysisPairwise
Alignment

BLAST

NGS

MSA

BLAST:

Search Engine for Sequences

NCBI Blast:knotted - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://blast.ncbi.nlm.nih.gov/Blast.cgi

Daily Local UEA Firefox NETGEAR Router transsys Dreamhost Most Visited

blastn: search translated NCBI Blast:knotted

Descriptions

Sequences producing significant alignments:

	Score	E	Value
obj AF022390.1 AF022390	Hordeum vulgare knotted class 1 homeod... 751	0.0	U
obj AF224498.1 AF224498	Triticum aestivum KNOTTED-1-like homeo... 672	0.0	U
obj AF224500.1 AF224500	Triticum aestivum KNOTTED-1-like homeo... 651	0.0	U
obj AF224499.1 AF224499	Triticum aestivum KNOTTED-1-like homeo... 651	0.0	U
obj 1015001..1 P1CQH01	Oryza sativa Japonica Group OSH1 mRNA fo... 513	2e-159	U*
obj EU981052.1	Zea mays clone 329250 homeobox protein OSH1 mR... 516	4e-145	U
obj AY107752.1	Zea mays full-length cDNA clone ZM_BFc0029823 ... 513	8e-145	U
obj XK61308.1 ZMKN1	Z.mays Knotted-1 (Kn-1) gene 514	1e-144	U
obj EF798860.1	Zea mays clone 322776 homeobox protein OSH1 mR... 514	1e-144	U
obj BT924231.1	Zea mays full-length cDNA clone ZM_BFc0029823 ... 514	1e-144	U
obj DQ0337422.1	Leersia virginica KNOTTED1-like homeodomain pr... 505	3e-142	U
obj DQ0337421.1	Chasmantium latifolium KNOTTED1 homeodomain p... 505	3e-142	U
obj DQ0337420.1	Setaria italica KNOTTED1 homeodomain protein (... 505	9e-142	U

Find: Previous Next Highlight Match case

none

<http://blast.ncbi.nlm.nih.gov/>

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
AnalysisPairwise
Alignment

BLAST

NGS

MSA

NCBI Blast:knotted - Mozilla Firefox

File Edit View History Bookmarks Tools Help

Daily Local UEA Firefox NETGEAR Router transsys Dreamhost Most Visited

blastn: search translated ... NCB Blast:knotted

>Oryza sativa (japonica cultivar-group) genomic DNA, chromosome 3 Length=36192742

Features in this part of subject sequence:
Dae03g0727000

Score = 172 bits (437), Expect(2) = 6e-64, Method: Compositional matrix adjust.
Identities = 123/222 (55%), Positives = 137/222 (61%), Gaps = 60/222 (27%)
Frame = -2

Query	Subject	Score
1	MEIIGHIFGL-GATA-----HQNNHSQIWCSSPLISAVIISPPFPQQQQQHQQQQGAGYLARSP	55
29403579	MEEI I HIFG+ GA+ H R PWGS LSA+ +PPPF Q Q Q Q AG +AB+P	
56	L5LN TAPPGVSNHGGGSGCSNPV1QLQLANGSIL EACAKAAKEPSSSSYAADVEAIKAKIISH	115
29403402	L+LNTA V NEVLQLANGSLL+AC KA + +S-SYAADVEAIKAKIISH	
116	FHSSLLAAYLCQK-----	130
29403246	PH+SSLLAAYLCQKASPTHMN*PS*LPPCM8IALAAARIYSIMSCQARSQLT*RNCRQGL	
131	-VGAPPEVEARLTAVAQDLELRQTLAGLGATPELELDQFM	171
	VGADREW+ARLTAVAGGLELRQTLAGLGATPELELDQFM	

Find: Previous Next Highlight Match case

None

<http://blast.ncbi.nlm.nih.gov/>

Outline

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

1 Introduction

Molecular Biology Basics

Resolving the Phylogeneny of Land Plants

Reconstructing Foot and Mouth Disease Transmission Trees

2 Sequence Analysis

Pairwise Alignment

BLAST

3 “Next Generation” Sequencing Challenges

4 Multiple Sequence Alignment (MSA)

"Next Generation" Sequencing

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

- Long DNA sequences cannot be read like a tape.
- Short fragments from random genomic locations can be sequenced.
- NGS generates **very large** number of (very) **short** sequencing reads.



<http://www.illumina.com/systems/miseq.ilmn>

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

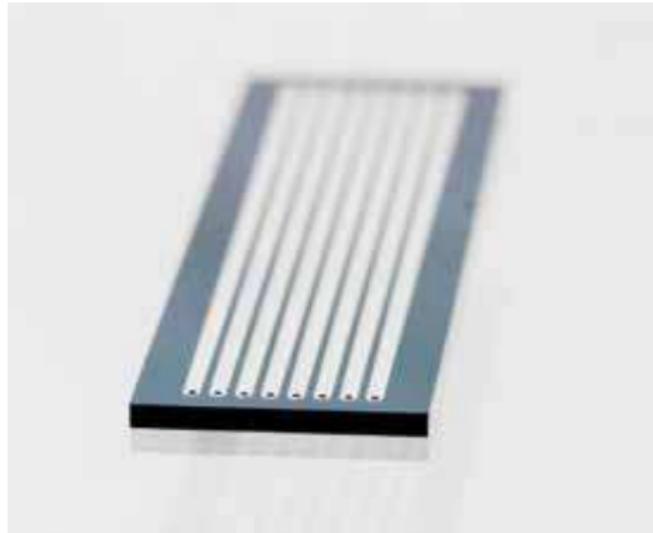
Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA

Illumina NGS Sequencing



Massive numbers of sequencing reactions take place in one
flow cell.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

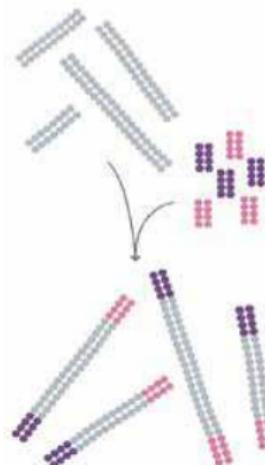
Illumina NGS Sequencing

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA



DNA is **fragmented** and **adapters** are ligated.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Illumina NGS Sequencing

Introduction

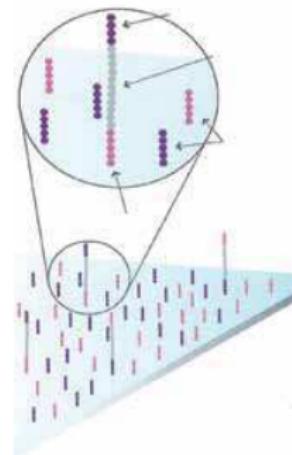
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



Fragments (with adapters) are **attached** to the slide in a flow cell.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Illumina NGS Sequencing

Introduction

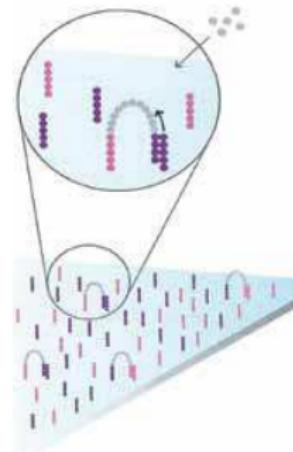
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



The slide is studded with **primers**, facilitating **bridge amplification** . . .

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Illumina NGS Sequencing

Introduction

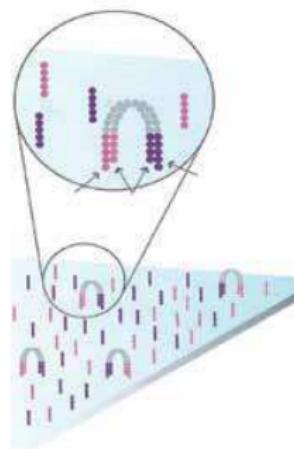
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



... resulting in **double stranded fragments** ...

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Illumina NGS Sequencing

Introduction

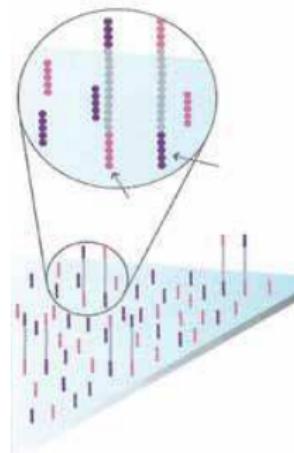
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



... which are then **denatured**.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Illumina NGS Sequencing

Introduction

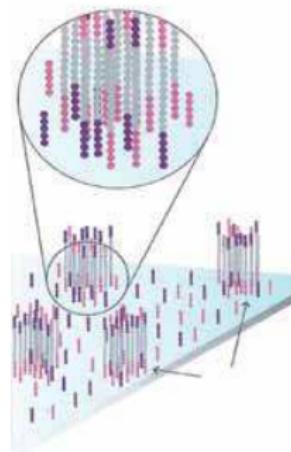
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



Multiple rounds of amplification result in a **cluster** from each initial fragment.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

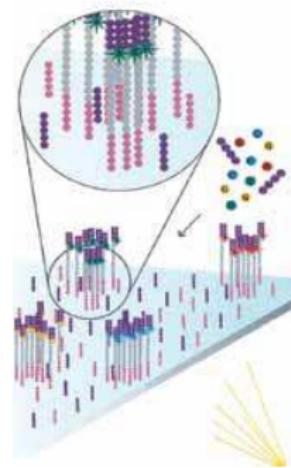
Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

Illumina NGS Sequencing



Reversible terminator nucleotides are added.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA

Illumina NGS Sequencing



Incorporated nucleotide fluoresce at different wave lengths.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Illumina NGS Sequencing

Introduction

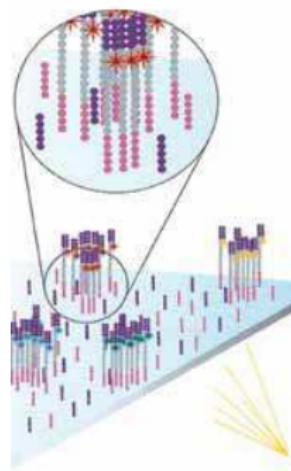
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



After removal of the terminator, the next nucleotide is added

...

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Illumina NGS Sequencing

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



... and the fluorescent light is imaged.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Illumina NGS Sequencing

Introduction

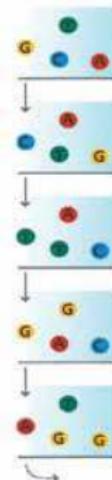
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



Each image yields one base **for each cluster**.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

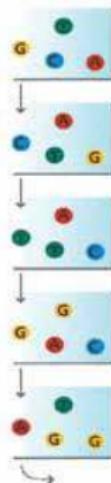
Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

Illumina NGS Sequencing



```
read0 =  
read1 =  
read2 =  
read3 =
```

Each image yields one base **for each cluster**.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

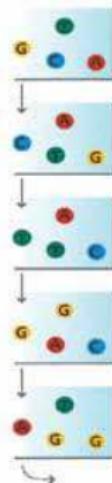
Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

Illumina NGS Sequencing



read0 = G
read1 = T
read2 = C
read3 = A

Each image yields one base **for each cluster**.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

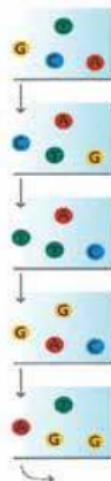
Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

Illumina NGS Sequencing



read0 = GC
read1 = TA
read2 = CT
read3 = AG

Each image yields one base **for each cluster**.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

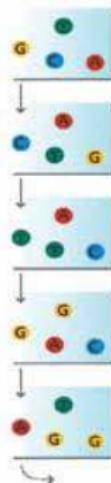
Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

Illumina NGS Sequencing



read0 = GCT^T
read1 = TAA^A
read2 = CTT^T
read3 = AGC^C

Each image yields one base **for each cluster**.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

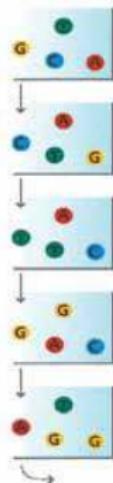
Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

Illumina NGS Sequencing



read0 = GCT**G**
read1 = TAAG**G**
read2 = CTT**A**
read3 = AGCC**C**

Each image yields one base **for each cluster**.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

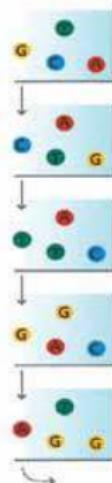
Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Illumina NGS Sequencing



read0 = GCTGA
read1 = TAAGT
read2 = CTTAG
read3 = AGCCG

Each image yields one base **for each cluster.**

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

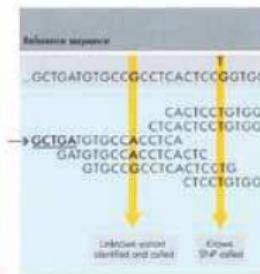
Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

Illumina NGS Sequencing



Reads are further processed, e.g. in sequence assembly.

http://www.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Applications of “Next Generation” Sequencing

- Novel genome sequencing
- Re-sequencing to discover **genomic variation**
 - Single nucleotide polymorphisms (SNPs), and their association to phenotypic traits,
 - Evolution of genomic variation patterns.
- Metagenomics
- *-Seq techniques
 - gene expression measurement: RNA-Seq
 - binding sites: ChIP-Seq
 - microRNA-Seq

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA

Mapping NGS Reads

- Task: align billions of to a known reference genome.
- Not feasible using dynamic programming.
- Feasible using advanced **indexing** of the reference.
 - e.g. Burrows-Wheeler transform
 - Pigeonhole principle

GA~~TAGAGT~~AGACGATGAGACCCATGACA

GGC GAGTAGACGAT GACCCATGATA
GGCT GAGTAGCCGATG CCCATGACA
GGCT AGTAGCCGATGAG CTCATGACA
GGCTAG GTAGACGATGAGA CATGACA
GGCTAGA AGACGATGAGA ATGACA
GGCTAGAG AGCCGATGAGACC ATGACA
GGCTAGAGT GACGATGAGACCC
CCGATGAGACCCAT

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
AnalysisPairwise
Alignment

BLAST

NGS

MSA

Finding Single Nucleotide Polymorphisms (SNPs)

GGCTAGAGT GACGATGAGACCCATGACA

GGC GAGTAGACGAT GACCCATGATA

GGCT GAGTAGCCGATG CCCATGACA

GGCT AGTAGCCGATGAG CTCATGACA

GGCTAG GTAGACGATGAGA CATGACA

GGCTAGA AGACGATGAGA ATGACA

GGCTAGAG AGCCGATGAGACC ATGACA

GGCTAGAGT GACGATGAGACCC

CCGATGAGACCCAT

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
Analysis

Pairwise

Alignment

BLAST

NGS

MSA

Finding Single Nucleotide Polymorphisms (SNPs)

GA~~T~~AGACTAGACGATGAGACCCATGACA

GGC GAGTAGACGAT GACCCATGATA

GGCT GAGTAGCCGATG CCCATGACA

GGCT AGTAGCCGATGAG CT~~T~~CATGACA

GGCTAG GTAGACGATGAGA CATGACA

GGCTAGA AGACGATGAGA ATGACA

GGCTAGAG AGCCGATGAGACC ATGACA

GGCTAGAGT GACGATGAGACCC

CCGATGAGACCCAT

Aligning Reads to Reference Sequences

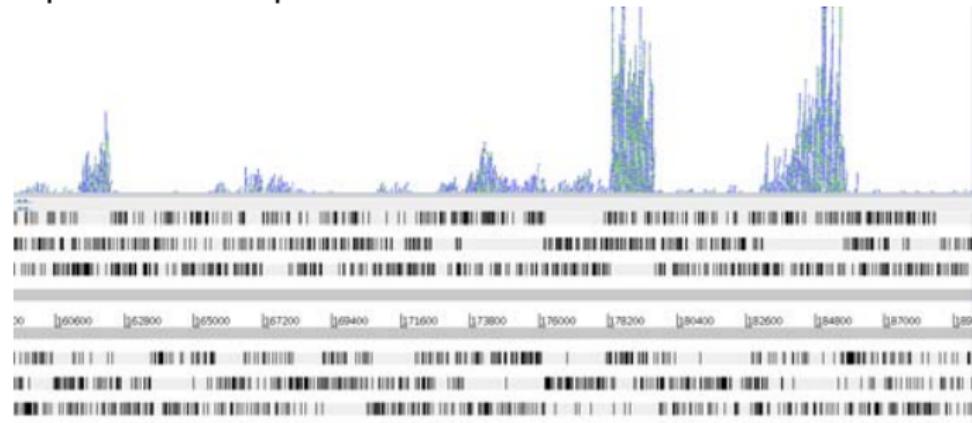
Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Example: RNA-Seq



Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
AnalysisPairwise
Alignment

BLAST

NGS

MSA

Assembling NGS Reads

GCTGATGTGCCGCCTCACTCCGGTGG

CACTCCGGTGG

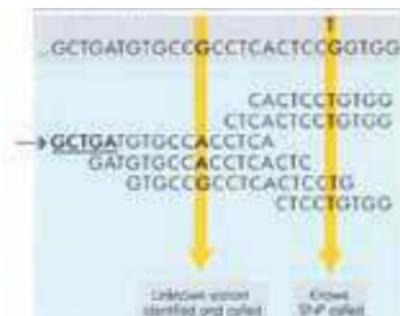
CTCACTCCTGTGG

GCTGATGTGCCACCTCA

GATGTGCCGCCTCACTC

GTGCCACCTCACTCCGG

CTCCGGTGG



- **Many copies** of a genome are **fragmented**
- Each base has **quality**, giving its **probability** of being correct.

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

NGS Assembly: Example

C A G A G C A

99 99 99 98 95 96 93

C A G A G C A G A C A

99 99 99 99 99 99 97 95 94 96 89

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

99 99 99 99 99 45 26 57 87 85 84 78

A A G T G C T A T C A

99 99 99 99 99 98 99 96 91 88 82

C T A T C A A C T

99 99 99 99 99 99 96 94 95

T A T C A A C T A G

99 99 99 99 99 94 97 95 91 88

A A C T A G

99 99 99 98 91 93

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

NGS Assembly: Example

C A G A G C A

99 99 99 98 95 96 93

C A G A G C A G A C A

99 99 99 99 99 99 97 95 94 96 89

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

99 99 99 99 99 45 26 57 87 85 84 78

A A G T G C T A T C A

99 99 99 99 99 98 99 96 91 88 82

C T A T C A A C T

99 99 99 99 99 99 96 94 95

T A T C A A C T A G

99 99 99 99 99 94 97 95 91 88

A A C T A G

99 99 99 98 91 93

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

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

NGS Assembly: Example

C A G A G C A

99 99 99 98 95 96 93

C A G A G C A G A C A

99 99 99 99 99 99 97 95 94 96 89

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

99 99 99 99 99 45 26 57 87 85 84 78

A A G T G C T A T C A

99 99 99 99 99 98 99 96 91 88 82

C T A T C A A C T

99 99 99 99 99 99 96 94 95

T A T C A A C T A G

99 99 99 99 99 94 97 95 91 88

A A C T A G

99 99 99 98 91 93

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

NGS Sequence Assembly

- Assembly depends on **overlaps** among reads.
- **Quality** of bases must be taken into account.
- Reads that are too short are **not informative**.
- **Repetitive sequences** make assembly difficult.
- Insufficient **depth** results in **multiple contigs**.
- Sufficient depth is a key success factor:
 - **Joining of contigs** depends on sufficient overlap (N_{50} value).
 - Resolving low quality bases depends on depth.
 - Depth does not help resolve repetitive sequences.

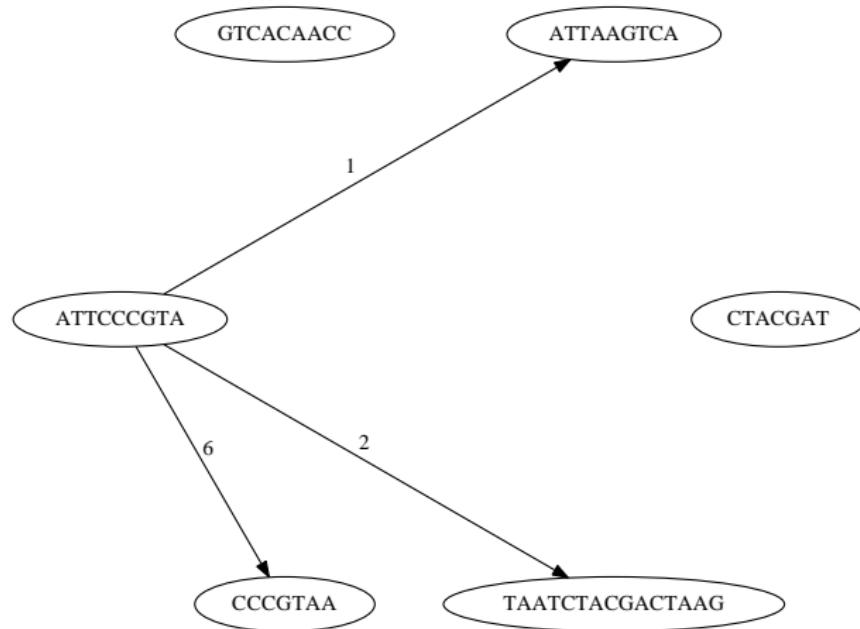
NGS Assembly: Overlap Approach

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA



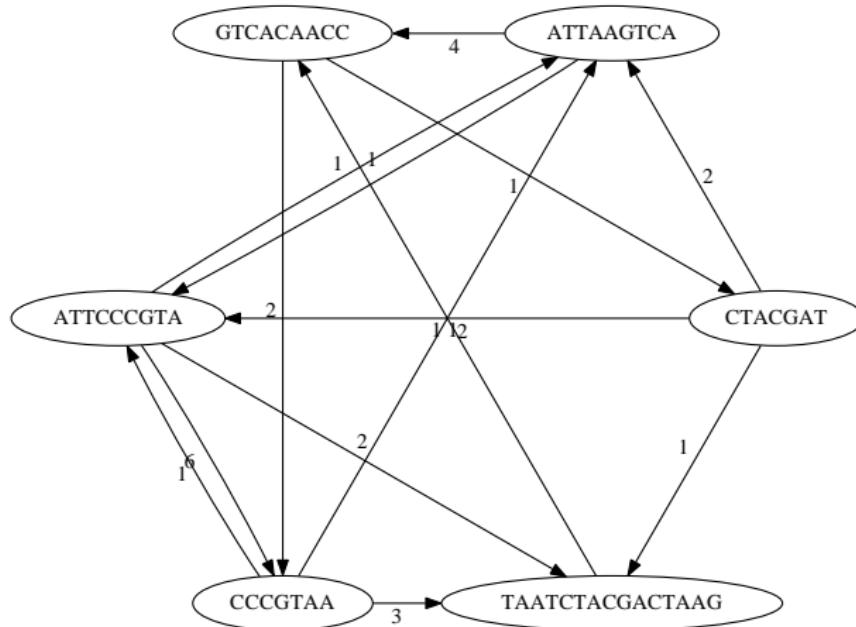
NGS Assembly: Overlap Approach

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA



Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

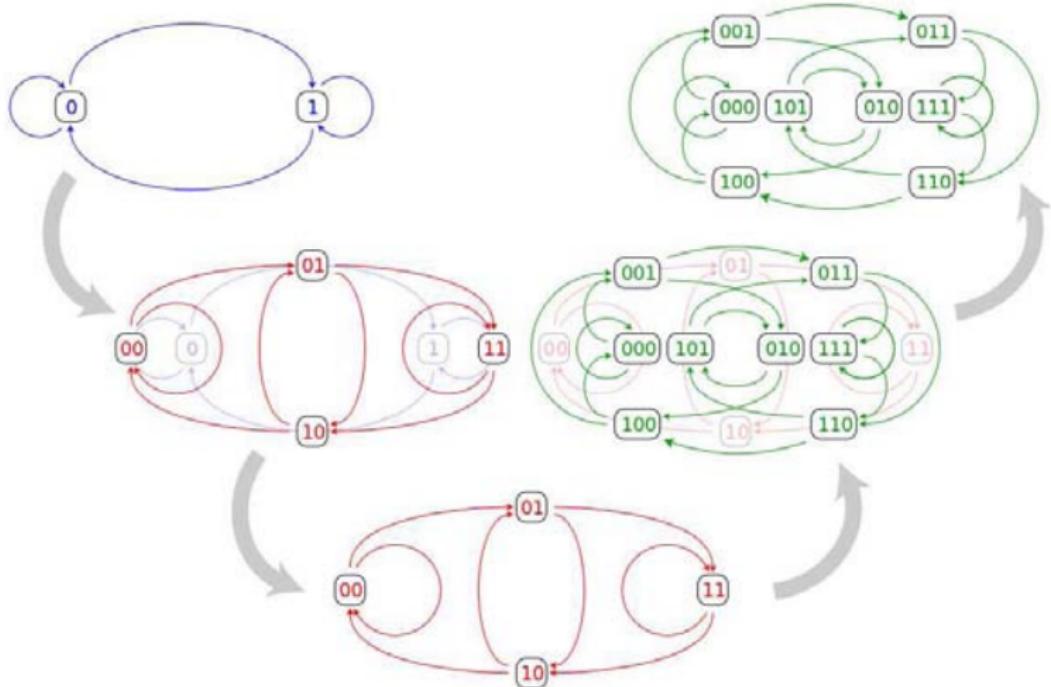
Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

De Bruijn Graph



<http://commons.wikimedia.org/wiki/File:DeBruijn-as-line-digraph.svg>

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

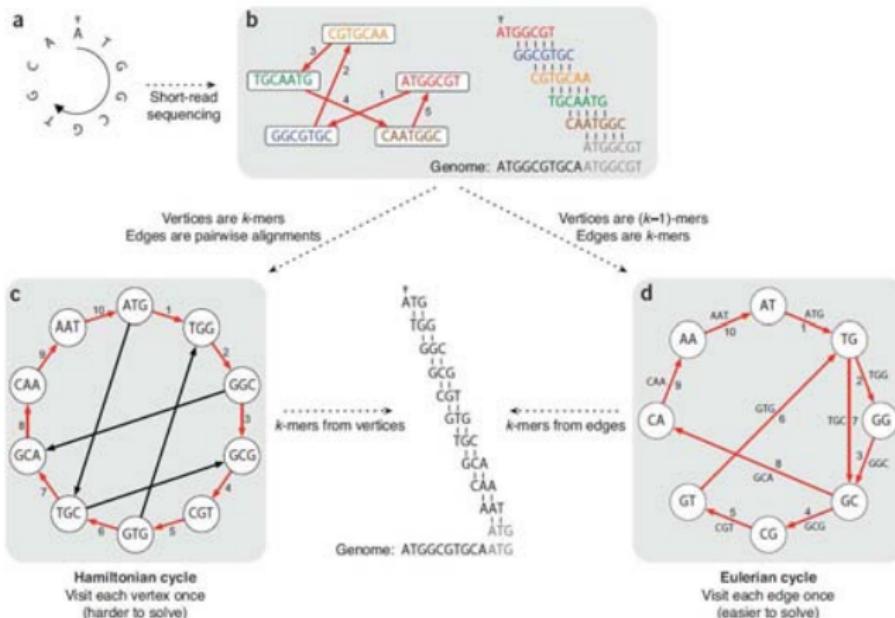
Transmission

Sequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

NGS Assembly: De Bruijn Approach



Compeau, Pevzner & Tesler, Nature Computational Biology 29 (2011): 987–991, Fig. 3

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

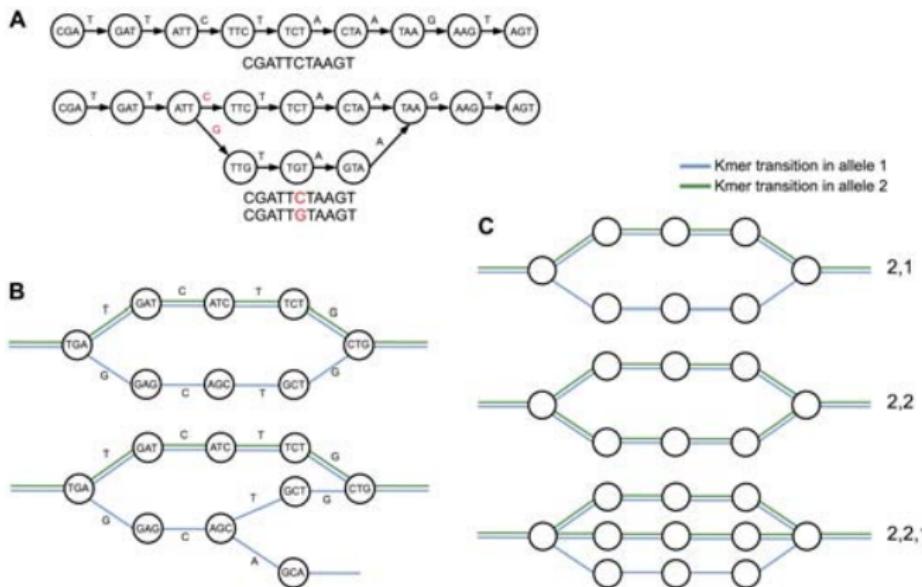
Transmission

Sequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Polymorphisms and de Bruijn Assembly



[Leggett et al., 2013, Fig. 1]

Summary 1: NGS Data Analysis

Introduction
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA

- **Mapping** to a reference sequence, using **indexing**
 - resequencing
 - detection of SNPs and other variants,
 - identification of genes (RNA-seq).
- **De novo assembly** of genomes or transcriptomes.
 - Resource intensive (particularly memory)
 - Overlap assembly: feasible with smaller sets
 - De Bruijn graph assembly of k -mers
- NGS metagenomics ...

⚠ Software has limitations and is evolving rapidly. ⚠

Outline

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence

Analysis

Pairwise

Alignment

BLAST

NGS

MSA

① Introduction

Molecular Biology Basics

Resolving the Phylogeneny of Land Plants

Reconstructing Foot and Mouth Disease Transmission Trees

② Sequence Analysis

Pairwise Alignment

BLAST

③ “Next Generation” Sequencing Challenges

④ Multiple Sequence Alignment (MSA)

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Multiple Alignment

- Extend pairwise approach?

- 2 sequences: table of n^2 prefix alignments
- 3 sequences: table of n^3 prefix alignments
- Warning: Very large numbers ahead

- Aligning 100 sequences of 300 symbols: about 10^{170} prefix alignments.
 - How much computing time does the universe have?

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

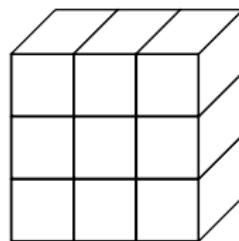
NGS

MSA

Multiple Alignment

- Extend pairwise approach?

- 2 sequences: table of n^2 prefix alignments
- 3 sequences: table of n^3 prefix alignments
- Warning: Very large numbers ahead



- Aligning 100 sequences of 300 symbols: about 10^{170} prefix alignments.
 - How much computing time does the universe have?

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

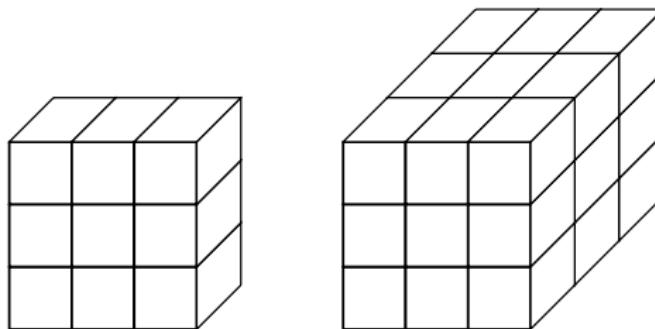
NGS

MSA

Multiple Alignment

- Extend pairwise approach?

- 2 sequences: table of n^2 prefix alignments
- 3 sequences: table of n^3 prefix alignments
- Warning: Very large numbers ahead



- Aligning 100 sequences of 300 symbols: about 10^{170} prefix alignments.
 - How much computing time does the universe have?

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

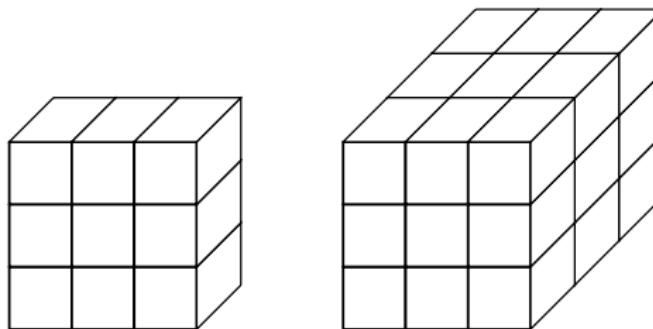
NGS

MSA

Multiple Alignment

- Extend pairwise approach?

- 2 sequences: table of n^2 prefix alignments
- 3 sequences: table of n^3 prefix alignments
- Warning: Very large numbers ahead



- Aligning 100 sequences of 300 symbols: about 10^{170} prefix alignments.
 - How much computing time does the universe have?

Progressive Multiple Alignment

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

- Compute **all** pairwise alignments
- Use **alignment dissimilarities** to produce a **guide tree**.
- Align most similar pair of sequences and **merge** them into a **profile**.
- **Progressively** align profiles.
- **Result:** All sequences aligned (and merged into one profile).
- Programs clustal, muscle

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
Analysis

Pairwise

Alignment

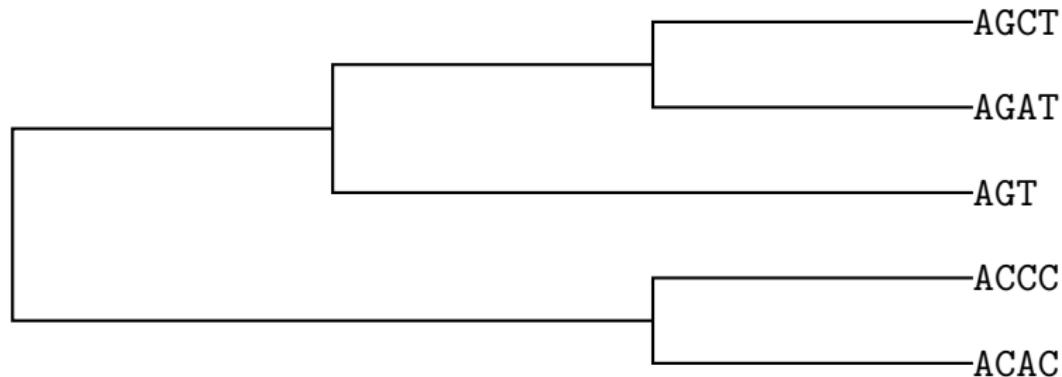
BLAST

NGS

MSA

More Uses of Continuous Sequences

- Profile searches (mostly superseded by HMMs)
- Progressive multiple alignment



Introduction

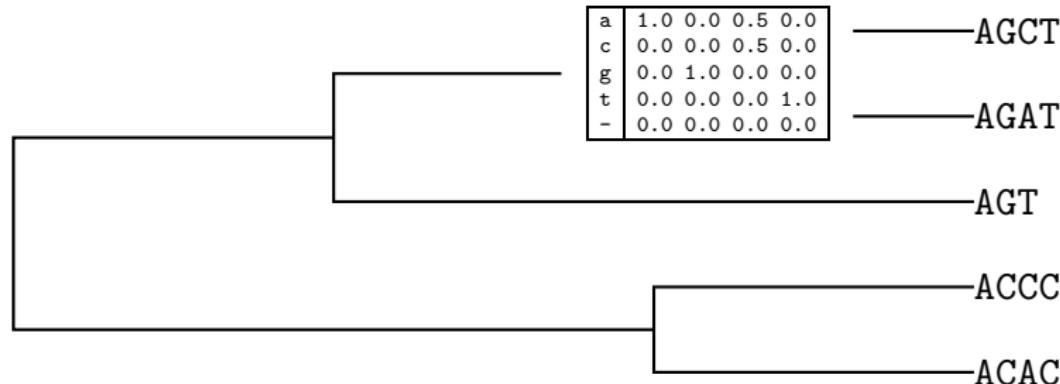
Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

More Uses of Continuous Sequences

- Profile searches (mostly superseded by HMMs)
- Progressive multiple alignment



Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
Analysis

Pairwise

Alignment

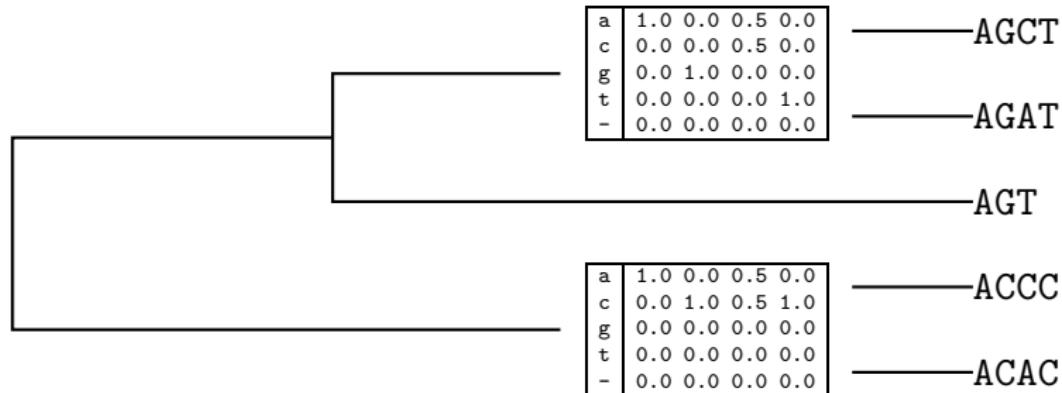
BLAST

NGS

MSA

More Uses of Continuous Sequences

- Profile searches (mostly superseded by HMMs)
- Progressive multiple alignment



Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

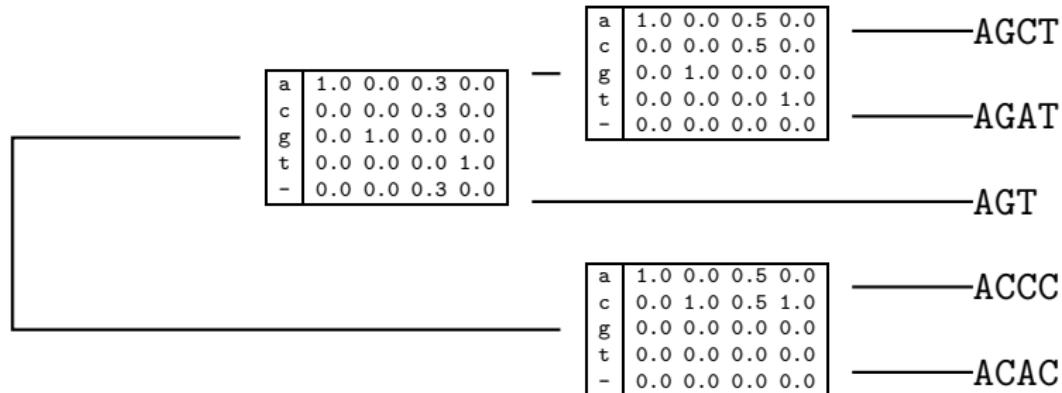
Sequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

More Uses of Continuous Sequences

- Profile searches (mostly superseded by HMMs)
- Progressive multiple alignment



Introduction

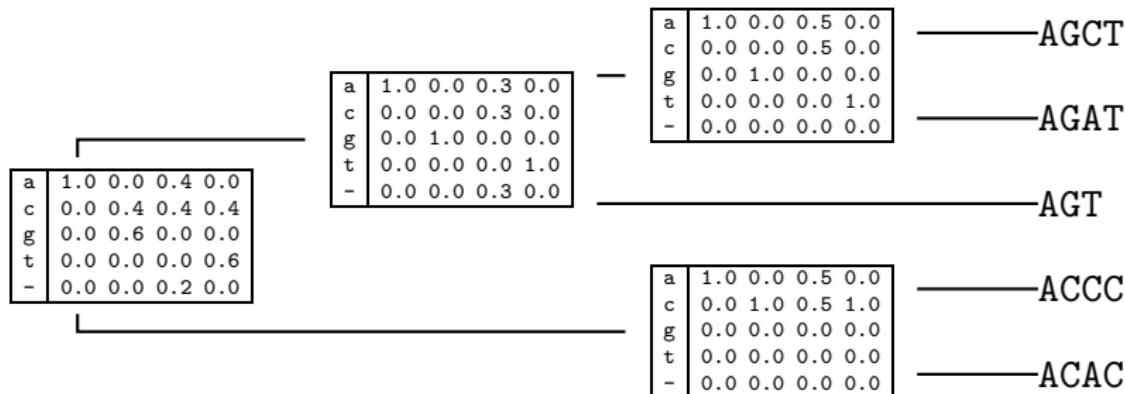
Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

More Uses of Continuous Sequences

- Profile searches (mostly superseded by HMMs)
- Progressive multiple alignment



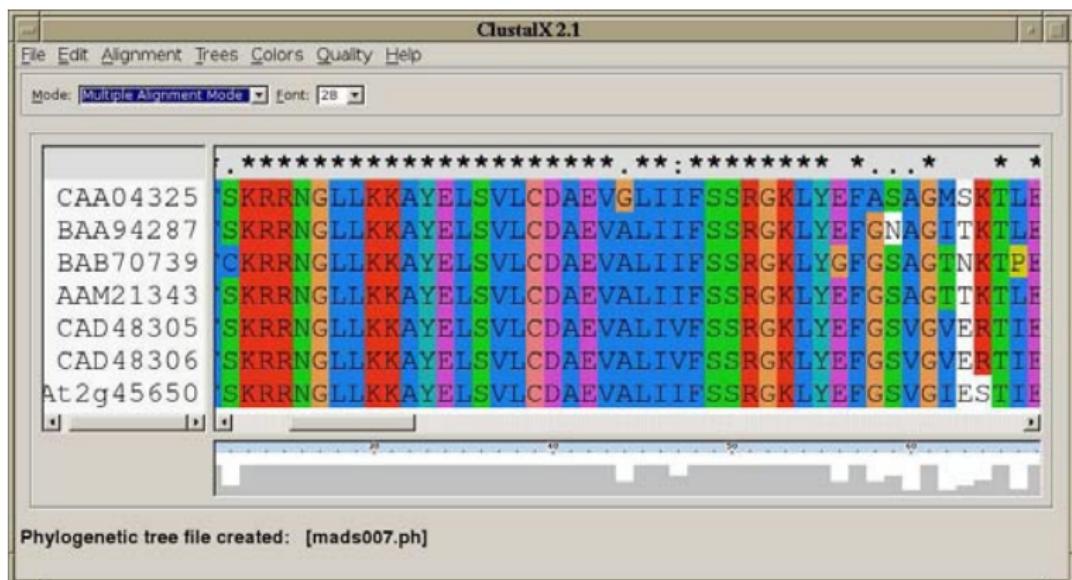
Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
TransmissionSequence
AnalysisPairwise
Alignment
BLAST

NGS

MSA

Multiple Alignment



(program: clustalx)

<http://www.clustal.org/>

Overview of Molecular Phylogeny

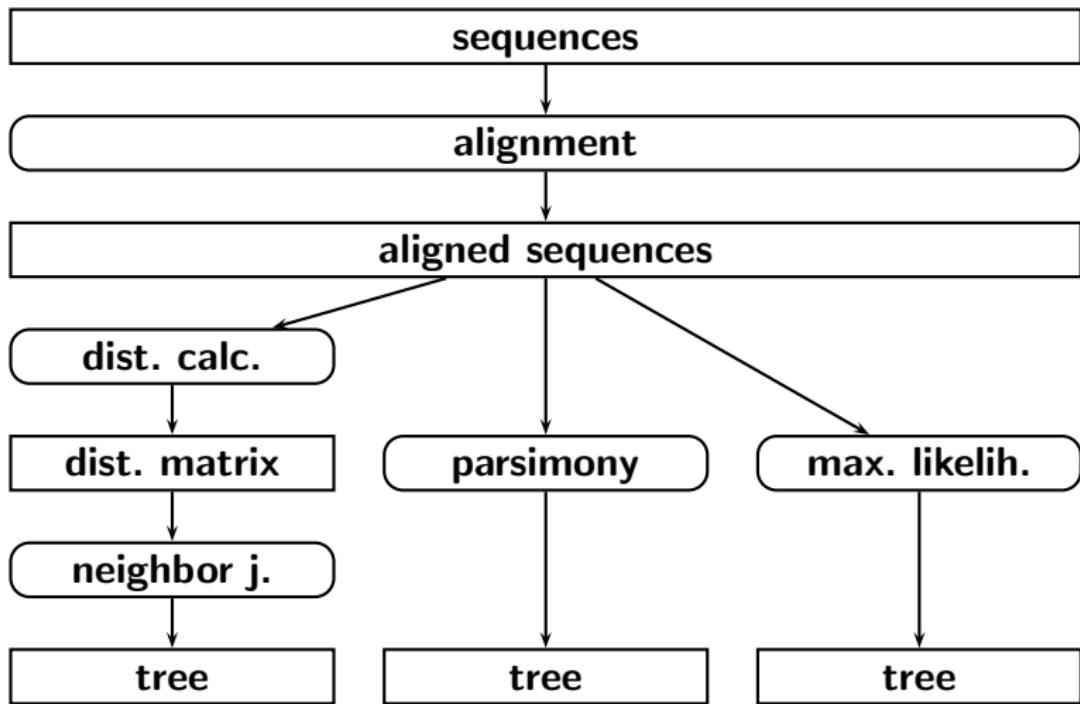
Introduction
Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence
Analysis

Pairwise
Alignment
BLAST

NGS

MSA



Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

Acknowledgements

- Kai-Uwe Winter, Thomas Münster, Luzie U. Wingen, Günter Theißen, Heinz Saedler
- Begoña Valdazo-Gonzalez, Nick Knowles, Don King
- Jan Gewehr, Thomas Martinetz, Daniel Polani, Simon Moxon, Vincent Moulton
- Anyela Camargo, Alessandra Devoto, John Turner

Introduction

Mol. Bio. Basics
Plant Phylogeny
FMD
Transmission

Sequence Analysis

Pairwise Alignment
BLAST

NGS

MSA

References

-  Altschul, S. F., Madden, T. L., Schäffer, A. A., Zhang, J., Zhang, Z., Miller, W., and Lipman, D. J. (1997). Gapped BLAST and PSI-BLAST: A new generation of protein database search programs. *Nucleic Acids Research*, 25:3389–3402. <http://nar.oupjournals.org/cgi/content/full/25/17/3389>.
-  Crosswell, L. C. and Thornton, J. M. (2012). ELIXIR: A distributed infrastructure for European biological data. *Trends in Biotechnology*, 30:241–242.
-  Langton, C. G. (1992). Preface. In Langton, C. G., Taylor, C., Farmer, J. D., and Rasmussen, S., editors, *Artificial Life II*, volume X of *Santa Fe Institute Studies in the Sciences of Complexity, Proceedings*, pages xiii–xviii, Redwood City, CA. Addison-Wesley.
-  Leggett, R. M., Ramirez-Gonzalez, R. H., Verweij, W., Kawashima, C. G., Iqbal, Z., Jones, J. D., Caccamo, M., and MacLean, D. (2013). Identifying and classifying trait linked polymorphisms in non-reference species by walking coloured de bruijn graphs. *PLoS One*, 8:e60058.
-  Winter, K.-U., Becker, A., Münster, T., Kim, J. T., Saedler, H., and Theißen, G. (1999). MADS-box genes reveal that gnetophytes are more closely related to conifers than to flowering plants. *Proceedings of the National Academy of Sciences, USA*, 96:7342–7347.
-  Wright, C. F., Morelli, M. J., Thébaud, G., Knowles, N. J., Merzky, P., Paton, D. J., Haydon, D. T., and King, D. P. (2011). Beyond the consensus: Dissecting within-host viral population diversity of foot-and-mouth disease virus by using next-generation genome sequencing. *Journal of Virology*, 85:2266–2275.

Introduction

Mol. Bio. Basics

Plant Phylogeny

FMD

Transmission

Sequence
Analysis

Pairwise
Alignment

BLAST

NGS

MSA

Thank You
for your attention and participation