

1º Simpósio  **LAVIREO**

Análise Filogenética: uma introdução

Edson Delatorre

Lab. de Genética Molecular de Microrganismos
Instituto Oswaldo Cruz/FIOCRUZ

delatorre.ioc@gmail.com

Rio de Janeiro – RJ, janeiro de 2019

Filogenética

Estudo da evolução e relacionamentos de uma coleção de “coisas” (genes, proteínas, organismos) que derivam de um ancestral comum.

Dados utilizados para análises filogenéticas

- Caracteres morfológicos
- Dados genéticos
 - RFLP
 - Frequências alélicas
 - Sequências de AA ou NT
- Uma combinação destes

Dados utilizados para análises filogenéticas

- Caracteres morfológicos
- Dados genéticos
 - RFLP
 - Frequências alélicas
 - **Sequências de AA ou NT**
- Uma combinação destes

Homologia x Analogia x Similaridade

- **Homologia:**

- Termo absoluto, significa que os OTU **compartilham um ancestral comum** recente o suficiente para que a variação observada entre as sequências carregue informação suficiente para as análises filogenéticas

- **Similaridade:**

- Termo relativo, comparação entre duas sequências expressa comumente como uma porcentagem.

- **Analogia:**

- Sequências que podem até ter alta similaridade, porém apresentam origens evolutivas separadas. São o resultado de evolução convergente.

Homologia x Analogia x Similaridade

- **Homologia:**

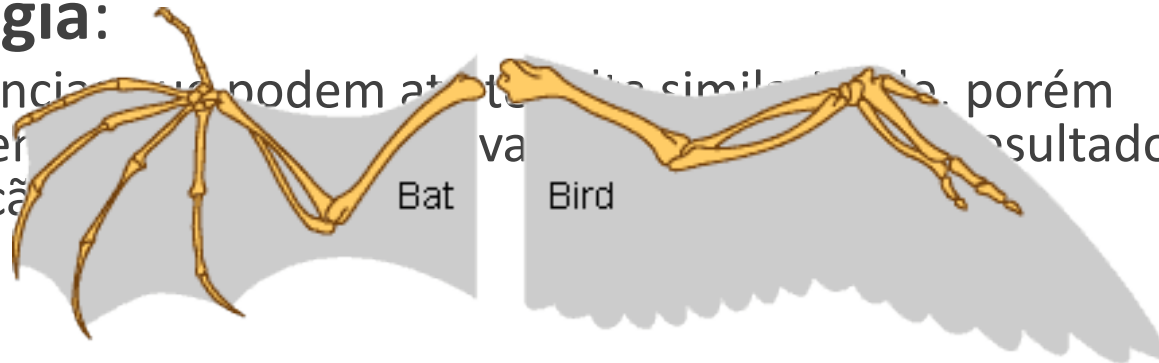
- Termo absoluto, significa que os OTU **compartilham um ancestral comum** recente o suficiente para que a variação observada entre as sequências carregue informação suficiente para as análises filogenéticas

- **Similaridade:**

- Termo relativo, comparação entre duas sequências expressa comumente como uma porcentagem.

- **Analogia:**

- Sequências podem apresentar similaridade, porém não são resultado de ancestral comum



Análises filogenéticas

- Somente sequências **homólogas** podem ser utilizadas para análises filogenéticas.
- **Recombinação:**
 - Perturba a comparação, uma vez que cada parte do genoma recombinado possui uma história evolutiva distinta.
- **Alinhamento de sequências**
 - Mutações pontuais
 - Indels: adição de “gaps” para conseguir a homologia posicional

Alinhamento

Human beta	-----VHLT PEEKSAVTALWGKVN -- VDEVGGEALGRLLVVYPWTQRFFESFGDLST
Horse beta	-----VQLS GEEKA AVLALWDKVN-- EEEVGGEALGRLLVVYPWTQRFFDSFGDLNS
Human alpha	-----VLSPADKTNVKA AWGKVG AHAGEYGA EALERMFLSFP TTKTYFPHF-DLS-
Horse alpha	-----VLSSAADKTNVKA AWSKVGGH AGEYGA EALERMFLGFP TTKTYFPHF-DLS-
Whale myoglobin	-----VLSEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSH PETLEK FDRFKHLKT
Lamprey globin	PIVDTG SVAPLSAAE TKIRSAWAPVYST YETSGVDILVKFFTSTPAAQE FFPKFKGLTT
Lupin globin	-----GALT ESQAALVKS SWEEFNANIPKH THRF FILVLEIAP AAKDL FSFLKGTSE
	* : : : * . : . : * : * : .

Human beta	PDAVMGN PKVKAH GKKVLGAFSDGLAHLDN-----LKGTFATLSEL HCD KLHVD PENFRL
Horse beta	PGAVMGN PKVKAH GKKVLHSFGEGVHHLDN-----LKGTF AA LSEL HCD KLHVD PENFRL
Human alpha	----HGSAQ VKG H GKKVAD ALTNAVAHVDD----MPNALSALS DLHA HKL RVD PVN FKL
Horse alpha	----HGSAQ VKAH GKKVGDALTLAVGHLD-----LPGALS NLSD L HA HKL RVD PVN FKL
Whale myoglobin	EAEMKAS EDLKKH GVTVLTALGAILKKKGH----HEAEL KPLAQ SHATKHKIP IKYLEF
Lamprey globin	ADQLKKS ADVRW HAERIINAVNDASMSDDT--EKMS MKLRD LSG KHAK SFQVDPQ YFKV
Lupin globin	VP--QNN PELQA HAGKVFKLVYEAAIQ LQVT GVVVT DATLKNLGS V HVSK GVAD- AHFPV
	. . : : * . : . : * . * . : .

Human beta	LGNVLVCVLA HHFGKEFTPPVQAA YQKV VAGVANALAHKYH-----
Horse beta	LGNVLVVVLAR HFGKDFTPELQAS YQKV VAGVANALAHKYH-----
Human alpha	LSHCLLVTLA AHLP AEFTPA VHAS LDKFLAS VSTVLTSKYR-----
Horse alpha	LSHCLLSTLA VHLPNDFTPAVHAS LDKFLSS VSTVLTSKYR-----
Whale myoglobin	ISEAIIHVL HSRHPGDFGADAQ GAMNKA LELFRKDIAAKYKELGYQG
Lamprey globin	LAAVIADTVAAG ---D-----AG FEKLMS MICILLRSAY-----
Lupin globin	VKEA ILKT IK EVVGAKWSEELNSAW TIAYDELA IVIKKEMNDAA---
	: : . : . . . :

NT ou AA?

- Para genes distantemente relacionados alinhamento de NTs podem se tornar ambíguos.
- Para genes proximamente relacionados, alinhamentos de AAs podem não conter uma quantidade suficiente de substituições para construção de uma árvore confiável.

Homologia posicional

```

Human beta      -----VHLTPEEKSAVTALWGKVN--VDEVGGEALGRLLVVYPWTQRFFESFGDLST
Horse beta      -----VQLSGEEKAAVLALWDKVN--EEEVGGEALGRLLVVYPWTQRFFDSFGDLSN
Human alpha     -----VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-DLS-
Horse alpha     -----VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTKTYFPHF-DLS-
Whale myoglobin -----VLSEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKFDRFKHLKT
Lamprey globin  PIVDTGSVAPLSAAEKTKIRSAWAPVYSTYETSGVDILVKFFTSTPAAQEFFPKFKGLTT
Lupin globin    -----GALTESQAALVKSWEEFNANIPKHTHRFFILVLEIAPAAKDLFSFLKGTSE
                *  :  :  :  *  .  :  :  :  *  :  *  :  :

```

Human beta	PDAVMGNPKVKAH <u>GKKVLGAFSDGLAHL</u> DN----LKGTFATLS <u>SELHC</u> DKLHVDPENFRL
Horse beta	PGAVMGNPKVKAH <u>GKKVLHSFGEGVHH</u> LNDN----LKGTFAA <u>LSSELHC</u> DKLHVDPENFRL
Human alpha	----HGSAQVKGH <u>GKKVADALTNAVA</u> HVDD----MPNAL <u>SALSDLHA</u> HKLVRD PVNFKL
Horse alpha	----HGSAQVKAH <u>GKKVGDALTLA</u> VGHLDN----LP <u>GALSNSDLHA</u> HKLVRD PVNFKL
Whale myoglobin	EAEMKAS <u>EDLKKHGVTVLTALGA</u> ILKKKGH----HEAE <u>LKPLAQSH</u> ATKHKIP IKYLEF
Lamprey globin	ADQLKKSADV <u>RWHAERIINAV</u> NDAVASMDDT--EKMSMKLRDLS <u>GKHAKS</u> FQVD PQYFKV
Lupin globin	VP--QNNPELQA <u>HAGKVFLVYE</u> AIIQLQVTGVVVTDATLKNLG <u>SVHV</u> SKGVAD-AHFPV
	* * *

```
Human beta      LGNVLV CVLAHHF GKEFT PPVQAA YQKV VAGVANALAHKYH-----
Horse beta     LGNVLVV VLARHFG KDFTPELQAS YQKV VAGVANALAHKYH-----
Human alpha    LSHCLLV TLA AHLPAEFT PAVHAS LDKFLASVSTVLT SKYR-----
Horse alpha    LSHCLLV TLA VHLPNDFTP AVHAS LDKFLSSVSTVLT SKYR-----
Whale myoglobin ISEAIIH V LHSRHPGDFGADA QGAMNKALELFRKDIA AKYKELGYQG
Lamprey globin LAAVIADTVAAG---D-----AGFEKLMSMICILLRSAY-----
Lupin globin   VKEA ILKTIKEVVGAKWSEELNSAWT IAYDELAIVIKKEMNDAA---
```

Estratégias de alinhamento

- Algoritmos de alinhamento auxiliam quando se tem uma grande quantidade de sequências e quando estas são distantemente relacionadas.
- Existem diversos algoritmos de alinhamento
 - Diferem na penalidade de gaps, pontuação e outros detalhes
- Qual algoritmo escolher?

Estratégias de alinhamento

Table 3.2 Typical alignment tasks and recommended procedures

Input data	Recommendations
2–100 sequences of typical protein length (maximum around 10 000 residues) that are approximately globally alignable	Use PROBCONS , T-COFFEE , and MAFFT or MUSCLE , compare the results using ALTAVISIT . Regions of agreement are more likely to be correct. For sequences with low percent identity, PROBCONS is generally the most accurate, but incorporating structure information (where available) via 3DCOFFEE (a variant of T-COFFEE) can be extremely helpful
100–500 sequences that are approximately globally alignable	Use MUSCLE or one of the MAFFT scripts with default options. Comparison using ALTAVISIT is possible, but the results are hard to interpret with larger numbers of sequences unless they are highly similar
>500 sequences that are approximately globally alignable	Use MUSCLE with a faster option (we recommend maxiters-2) or one of the faster MAFFT scripts
Large numbers of alignments, high-throughput pipeline	Use MUSCLE with faster options (e.g. maxiters-1 or maxiters-2) or one of the faster MAFFT scripts
2–100 sequences with conserved core regions surrounded by variable regions that are not alignable	Use DIALIGN
2–100 sequences with one or more common domains that may be shuffled, repeated or absent	Use ProDA
A small number of unusually long sequences (say, >20 000 residues)	Use CLUSTALW . Other programs may run out of memory, causing an abort (e.g. a segmentation fault)

Estratégias de alinhamento

Table 3.2 Typical alignment tasks and recommended procedures

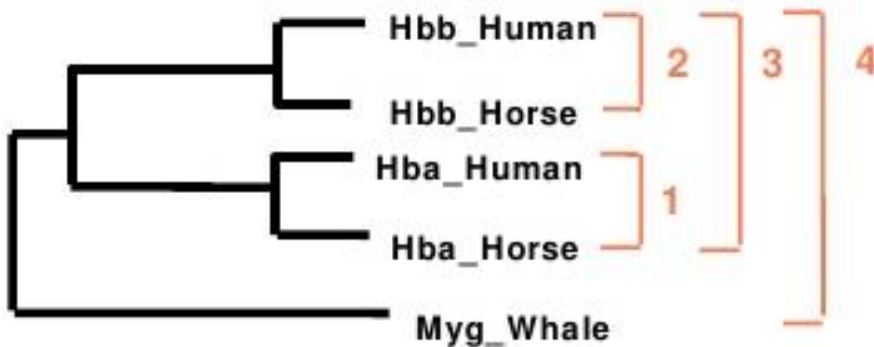
Input data	Recommendations
2–100 sequences of typical protein length (maximum around 10 000 residues) that are approximately globally alignable	Use PROBCONS , T-COFFEE , and MAFFT or MUSCLE ; compare the results using ALTAVISIT . Regions of agreement are more likely to be correct. For sequences with low percent identity, PROBCONS is generally the most accurate, but incorporating structure information (where available) via 3DCOFFEE (a variant of T-COFFEE) can be extremely helpful
100–500 sequences that are approximately globally alignable	Use MUSCLE or one of the MAFFT scripts with default options. Comparison using ALTAVISIT is possible, but the results are hard to interpret with larger numbers of sequences unless they are highly similar
>500 sequences that are approximately globally alignable	Use MUSCLE with a faster option (we recommend maxiters-2) or one of the faster MAFFT scripts
Large numbers of alignments, high-throughput pipeline	Use MUSCLE with faster options (e.g. maxiters-1 or maxiters-2) or one of the faster MAFFT scripts
2–100 sequences with conserved core regions surrounded by variable regions that are not alignable	Use DIALIGN
2–100 sequences with one or more common domains that may be shuffled, repeated or absent	Use ProDA
A small number of unusually long sequences (say, >20 000 residues)	Use CLUSTALW . Other programs may run out of memory, causing an abort (e.g. a segmentation fault)

CLUSTAL W

Nucl. Acids Res. 1994, 22:4673-4680

Hbb_Human	1	-			
Hbb_Horse	2	.17	-		
Hba_Human	3	.59	.60	-	
Hba_Horse	4	.59	.59	.13	-
Myg_Whale	5	.77	.77	.75	.75

Alinhamento de todos os pares.
Cálculo de uma matrix de distâncias.



A partir da matrix de distâncias
se constrói uma árvore filogenética
que servirá como guia.

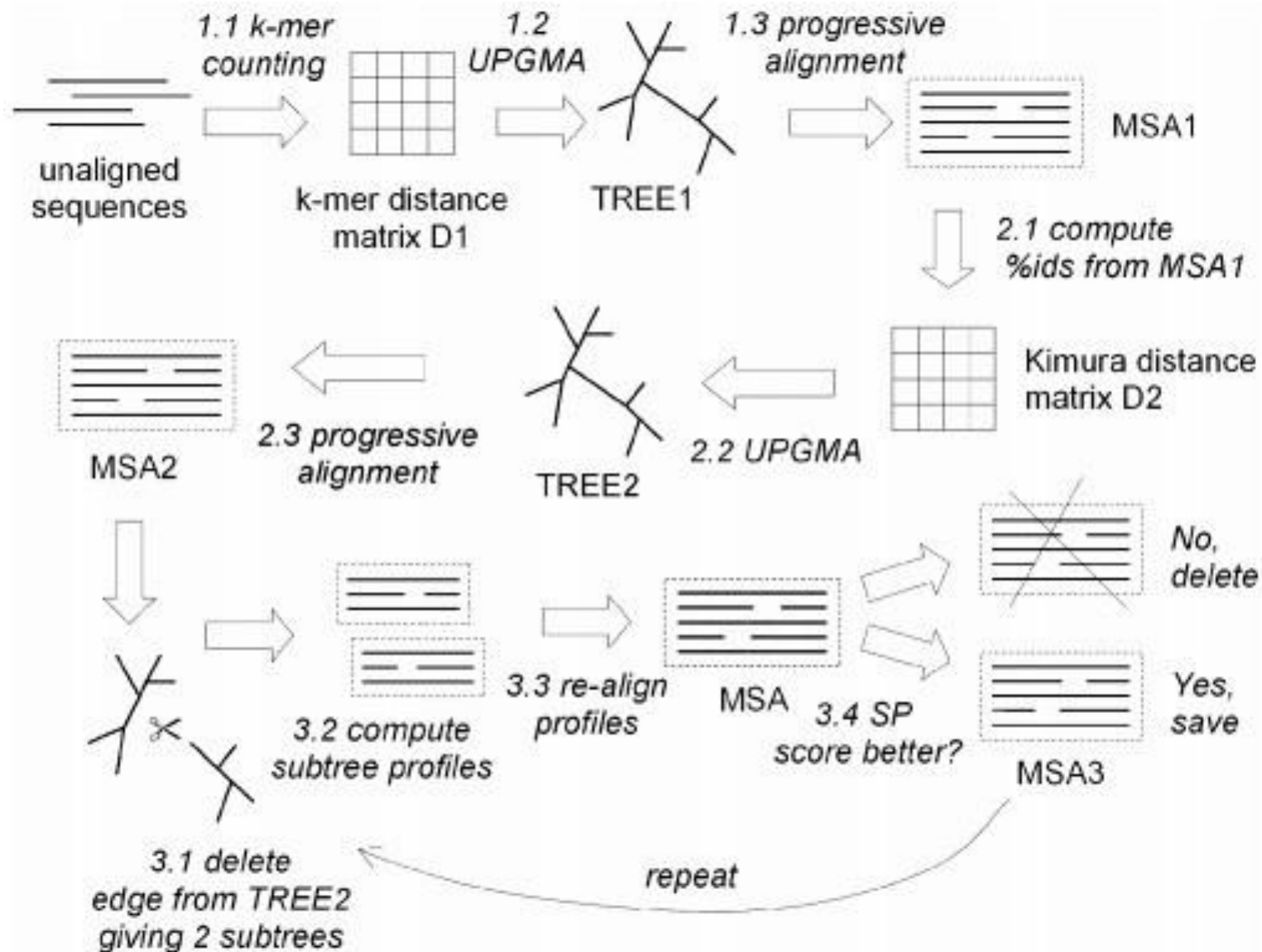
alpha-helices

1	PEEKSAVTALWGKVN--VDEVGG			
2	GEEKAAVLALWDKVN--EEEVGG			
3	PADKTNVKAAWGKVGAGHAGEYGA			
4	AADKTNVKAAWSKVGGHAGEYGA			
5	EHEWQLVLHVWAKVEADVAGHGQ			

Alinhamento global progressivo
das folhas até a raiz.

MUSCLE

Nucl. Acids Res. 2004 32:1792-1797

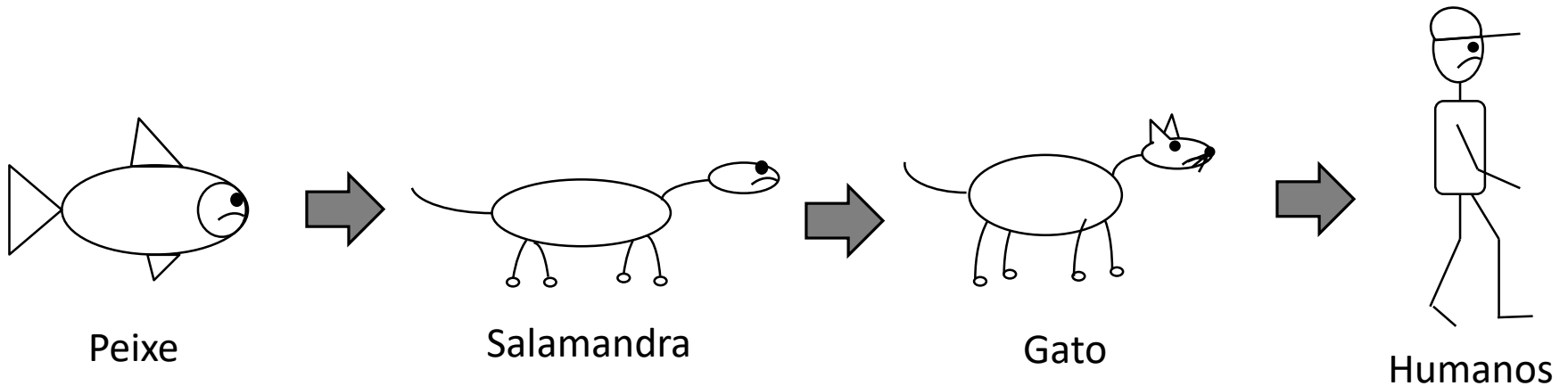


Rascunho

Progressivo

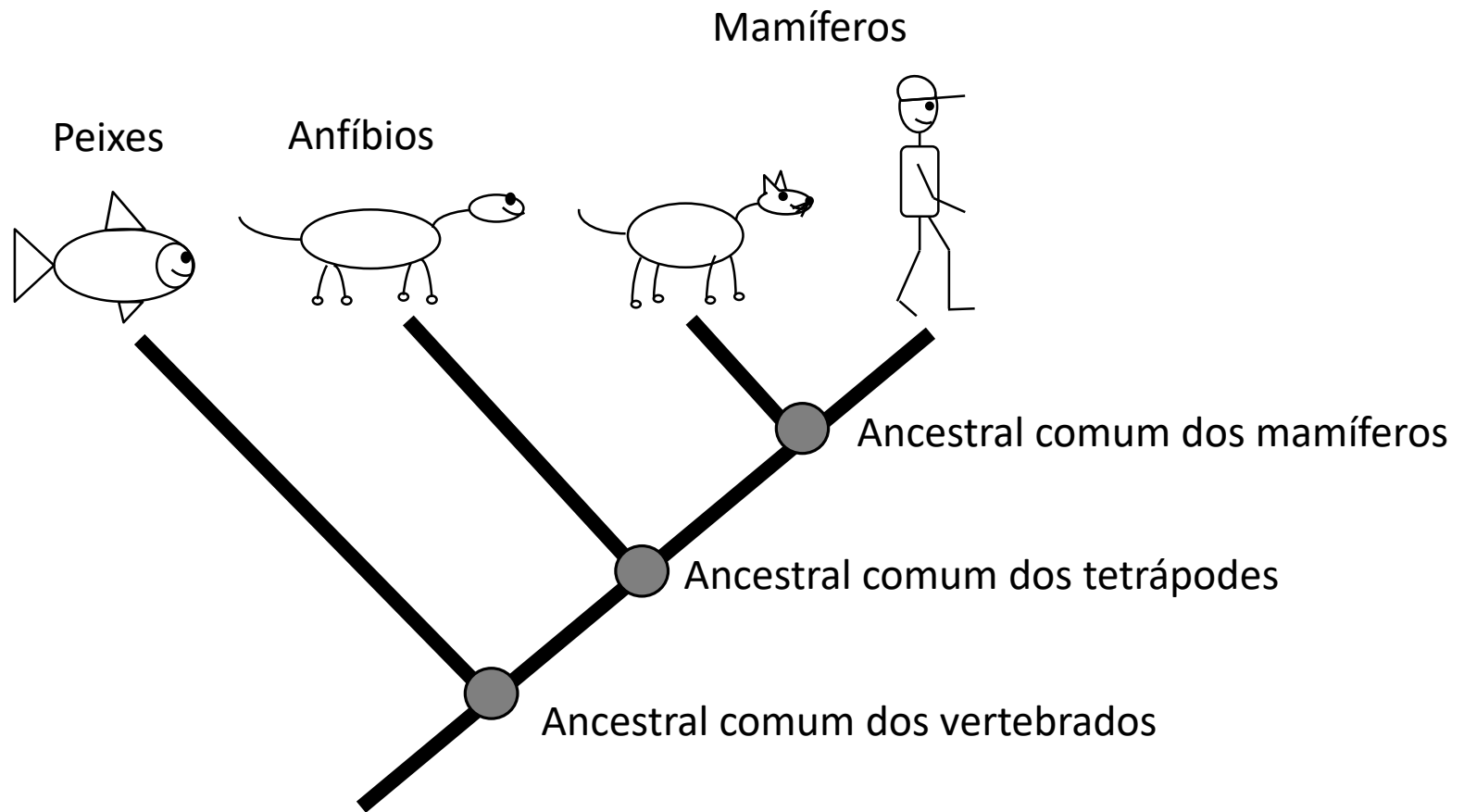
Refinamento

Árvores filogenéticas e evolução



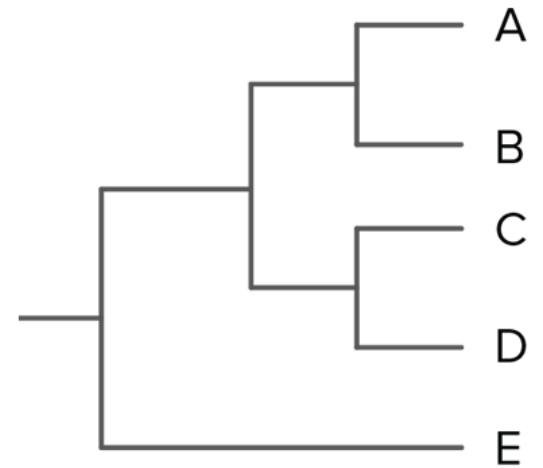
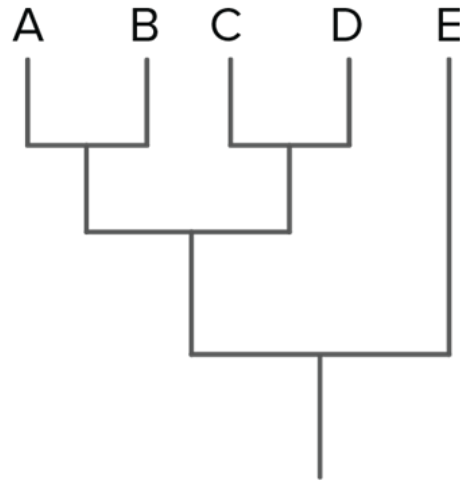
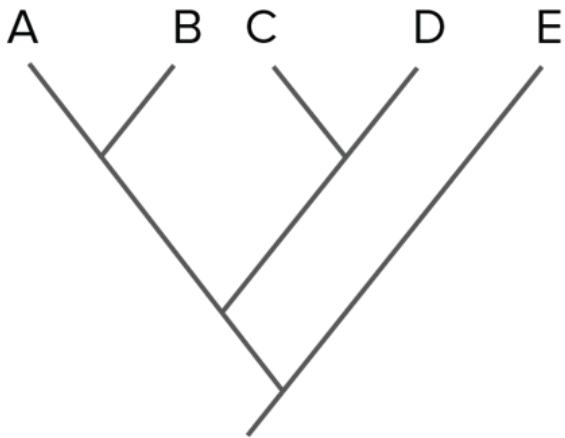
Isto **NÃO** é evolução

Árvores filogenéticas e evolução

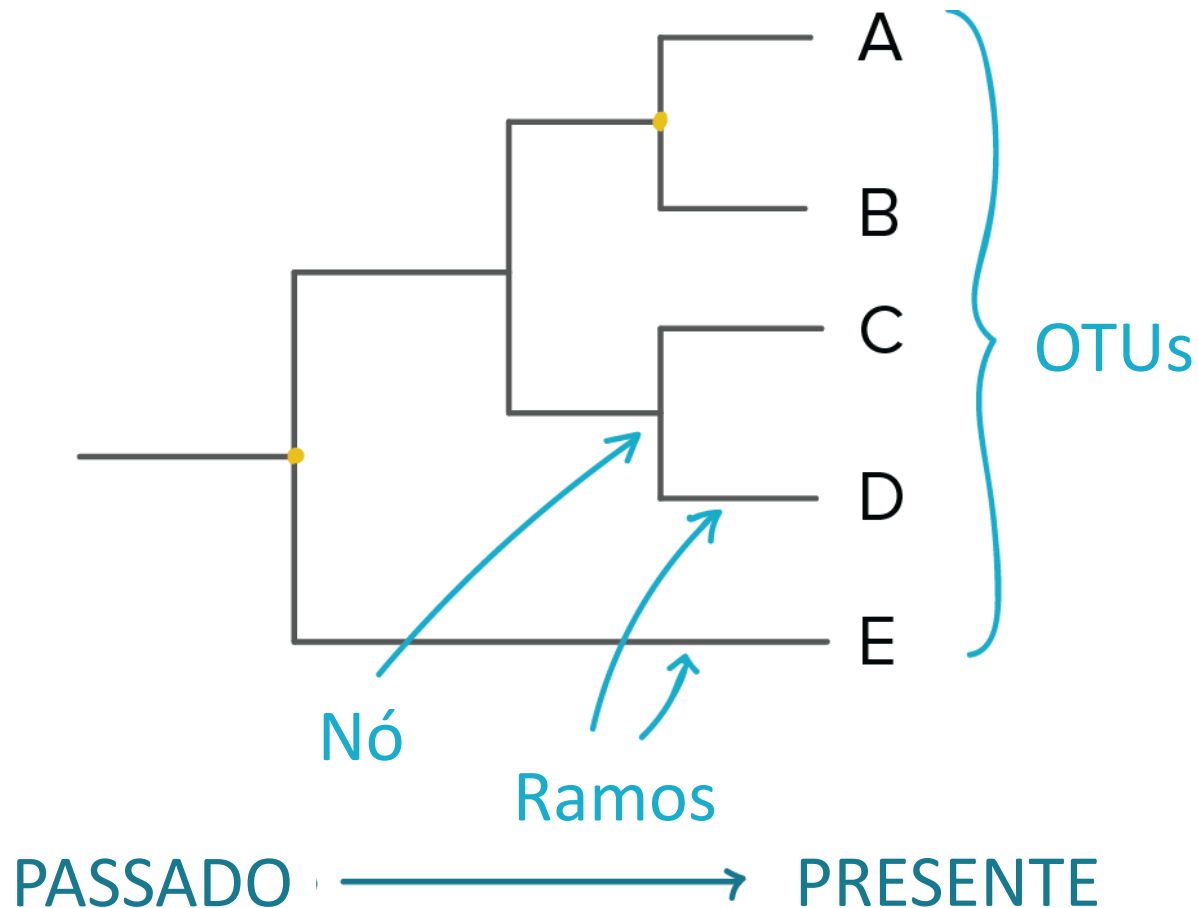


Isto **É** evolução

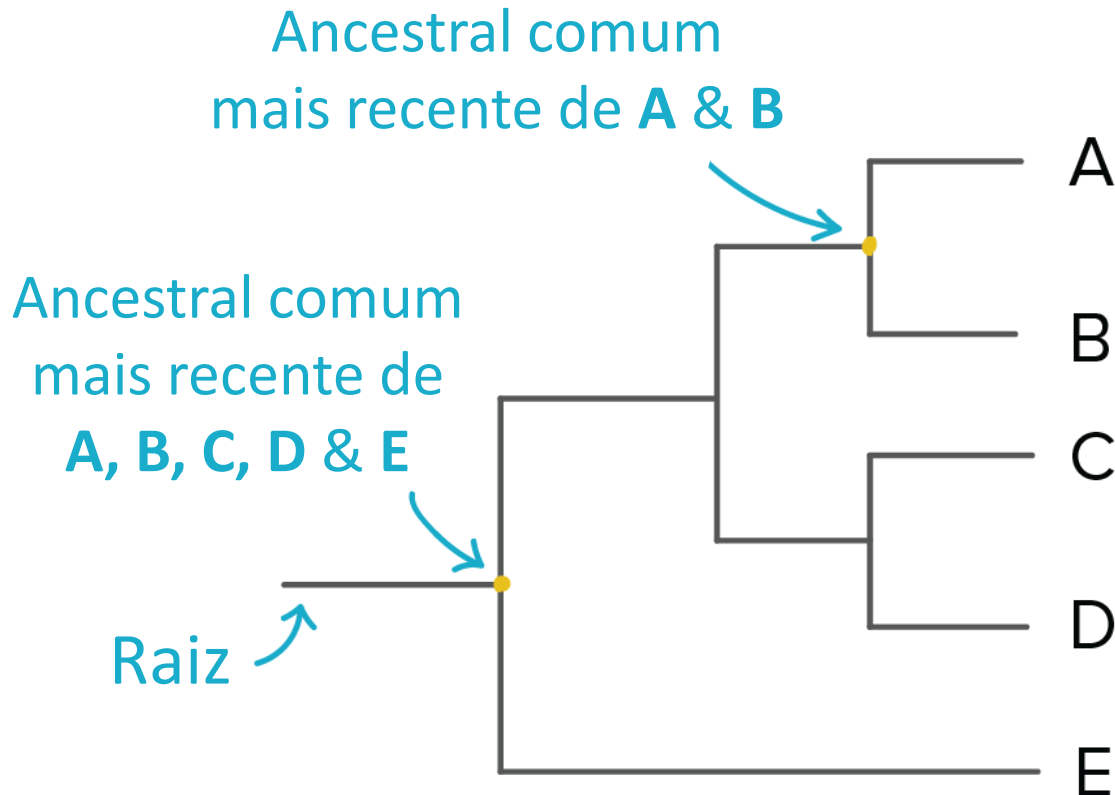
Anatomia de uma árvore filogenética



Anatomia de uma árvore filogenética

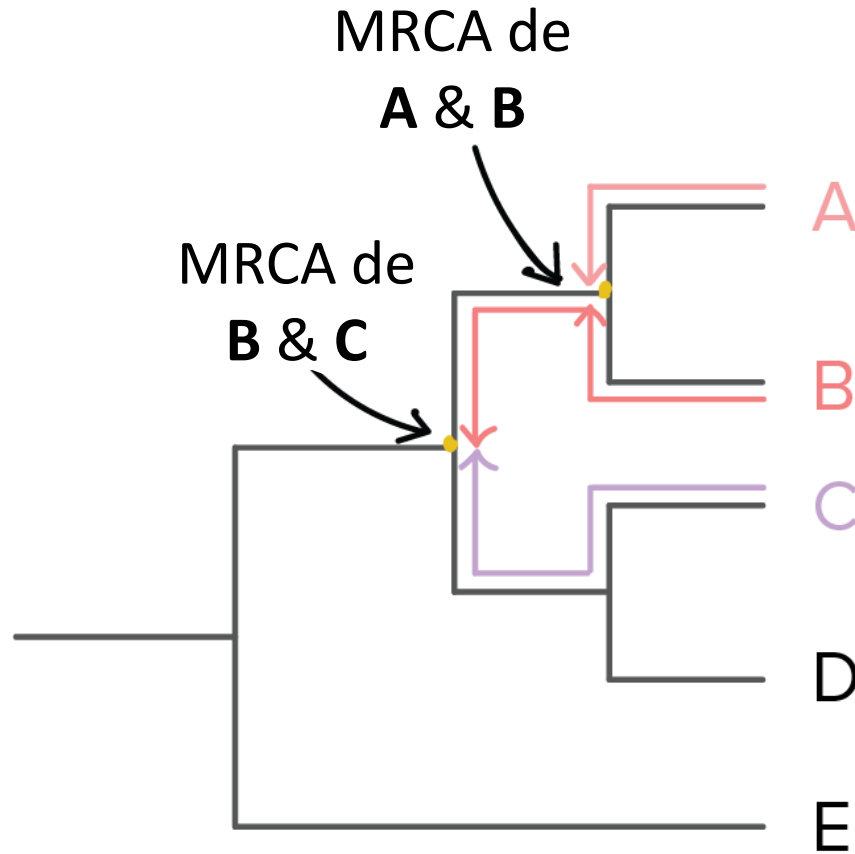


Anatomia de uma árvore filogenética

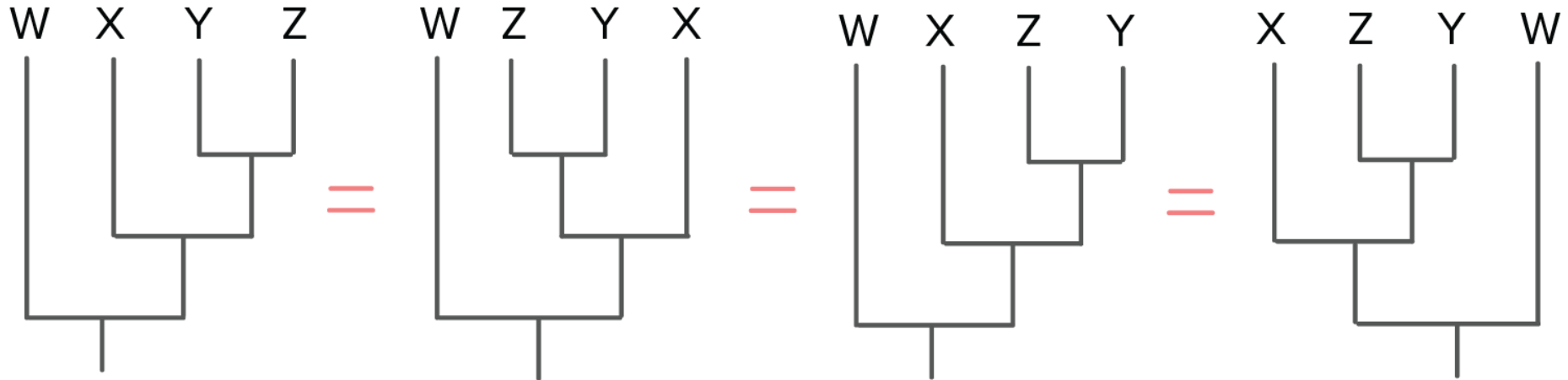
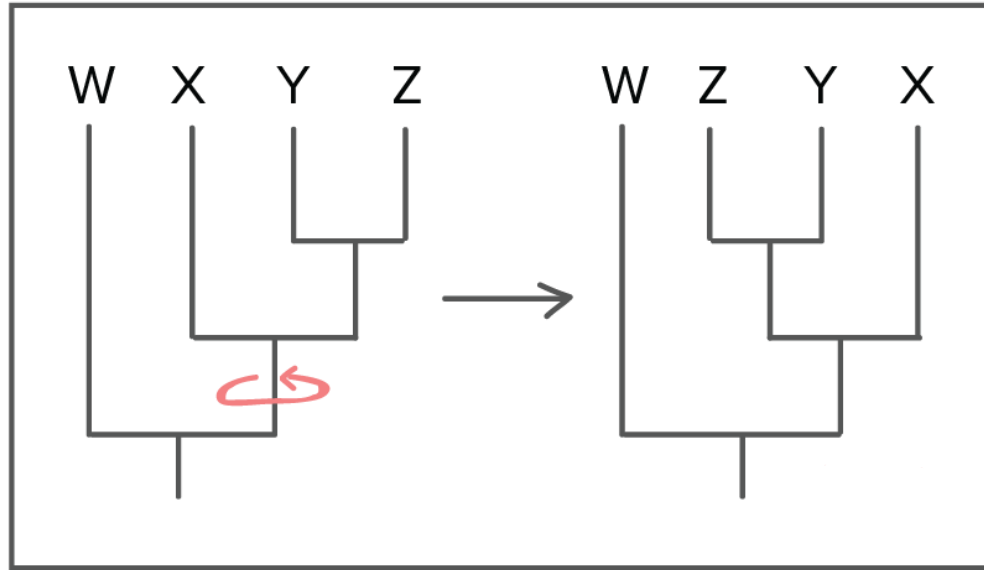


Ancestral comum mais recente = **MRCA**

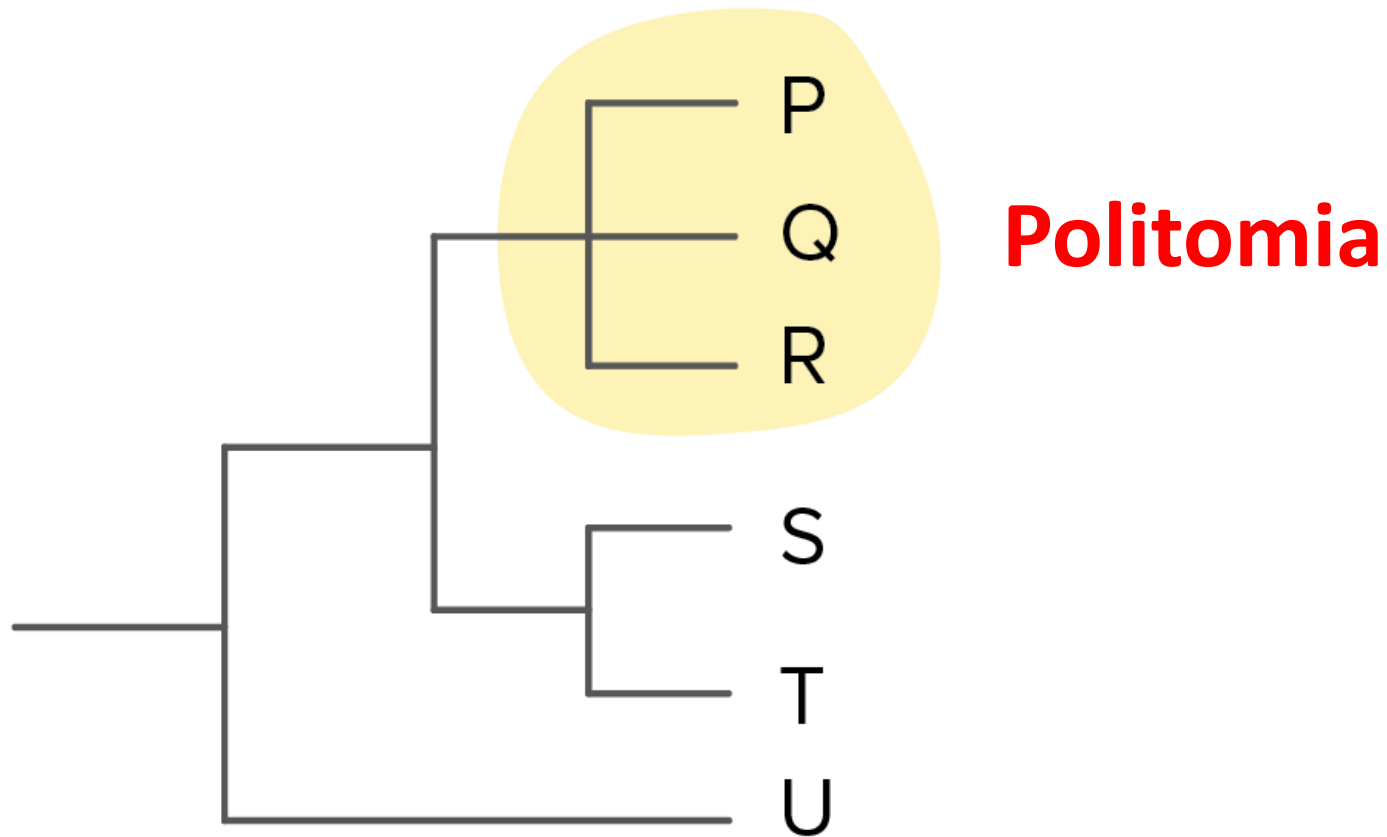
Anatomia de uma árvore filogenética



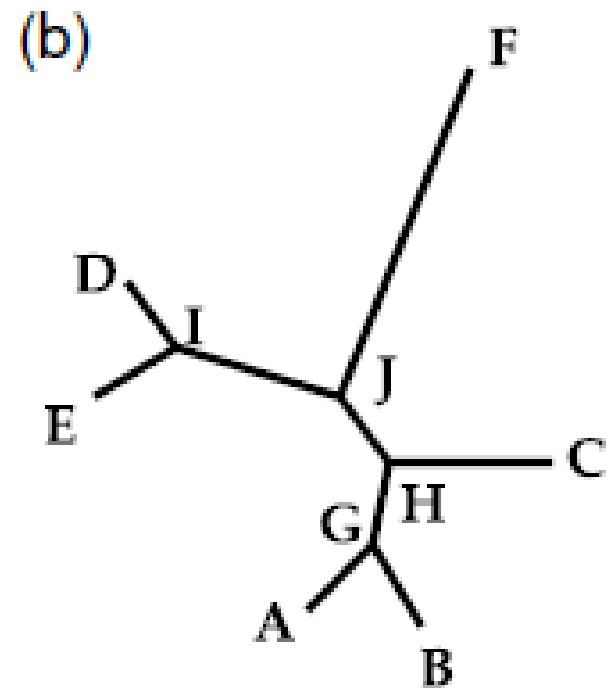
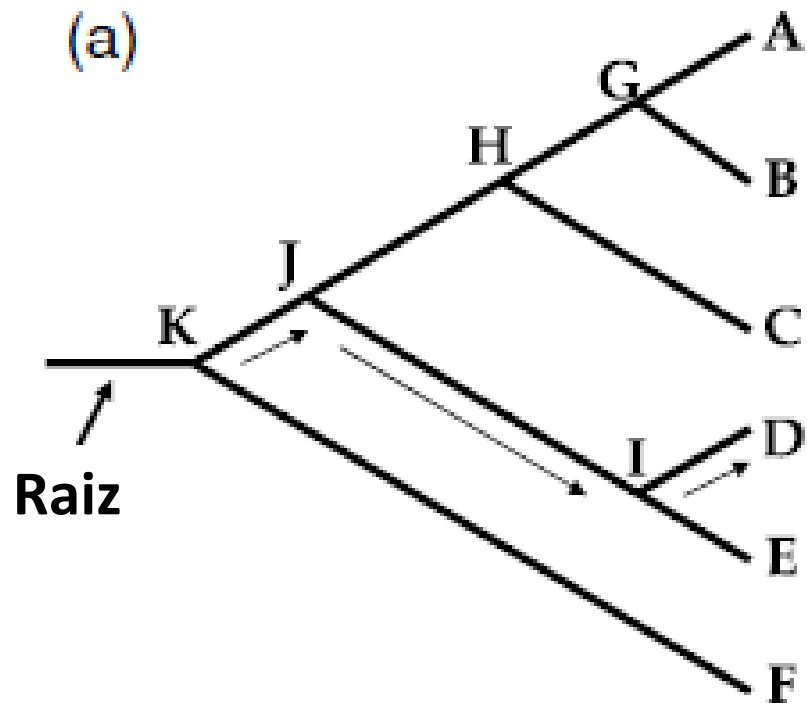
Anatomia de uma árvore filogenética



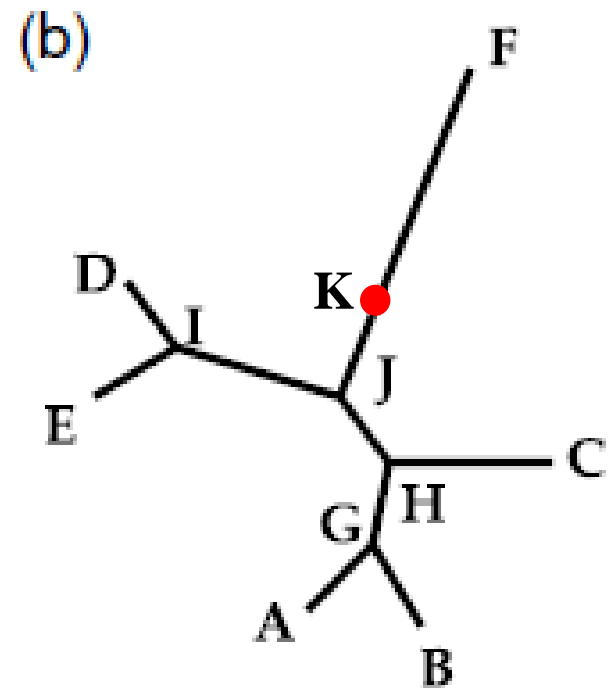
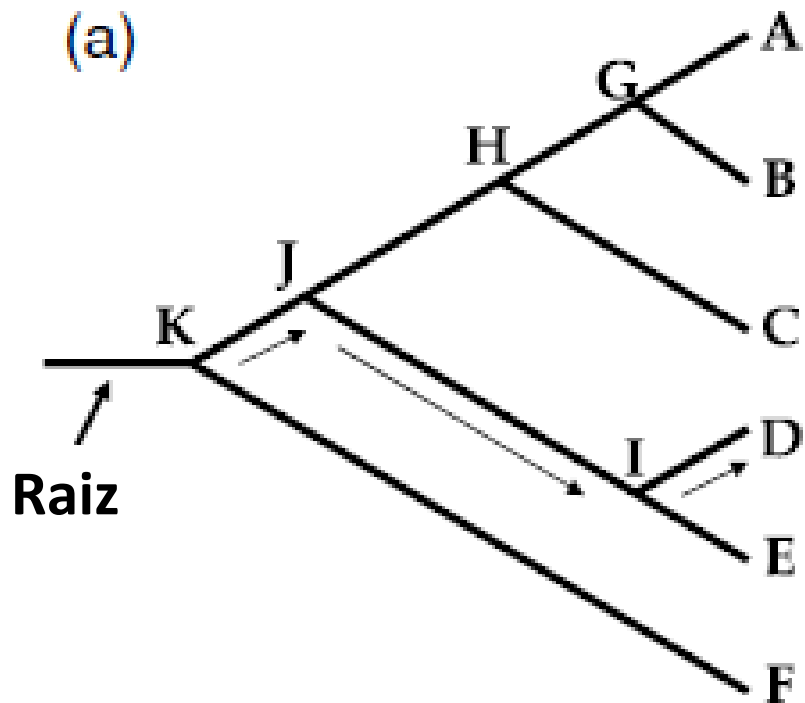
Anatomia de uma árvore filogenética



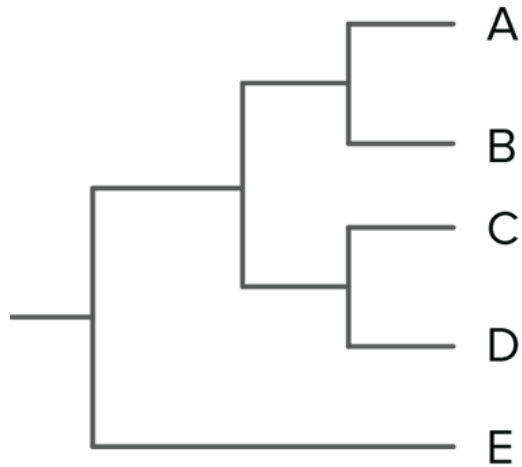
Anatomia de uma árvore filogenética



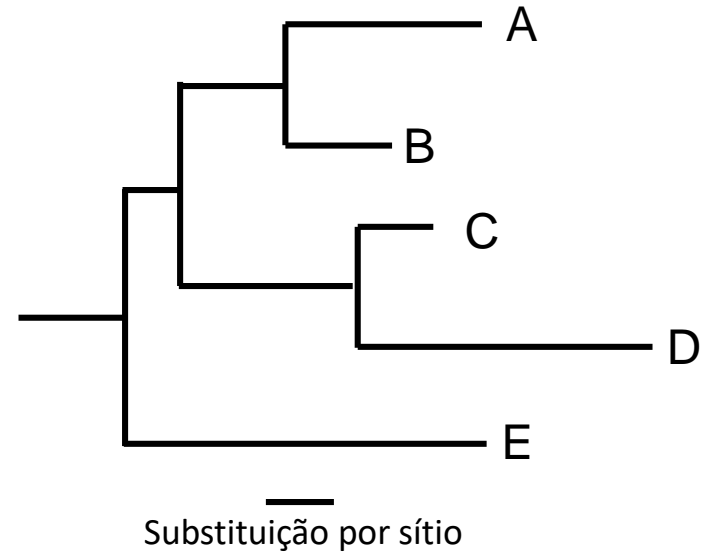
Anatomia de uma árvore filogenética



Anatomia de uma árvore filogenética

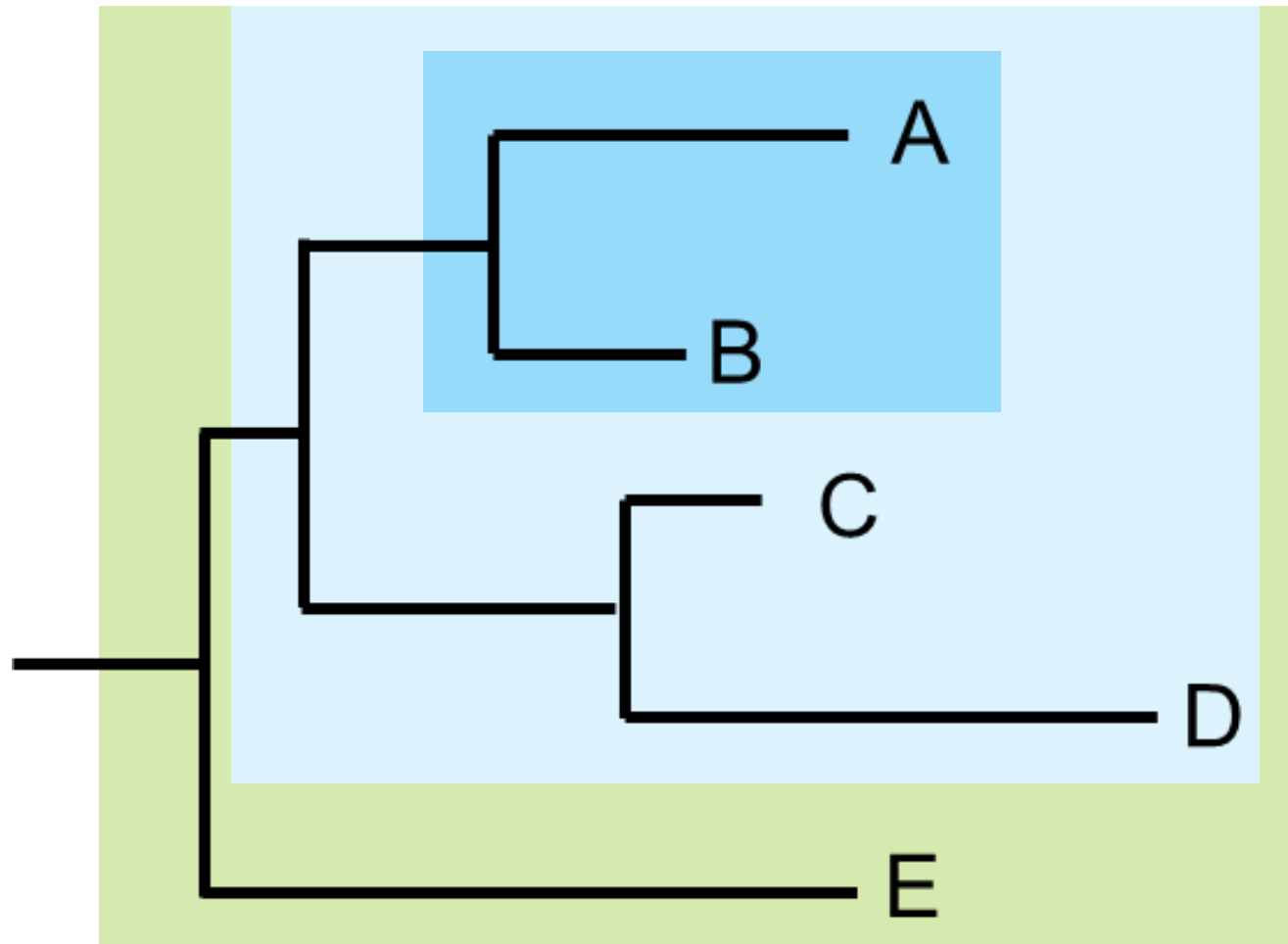


Cladograma



Filograma

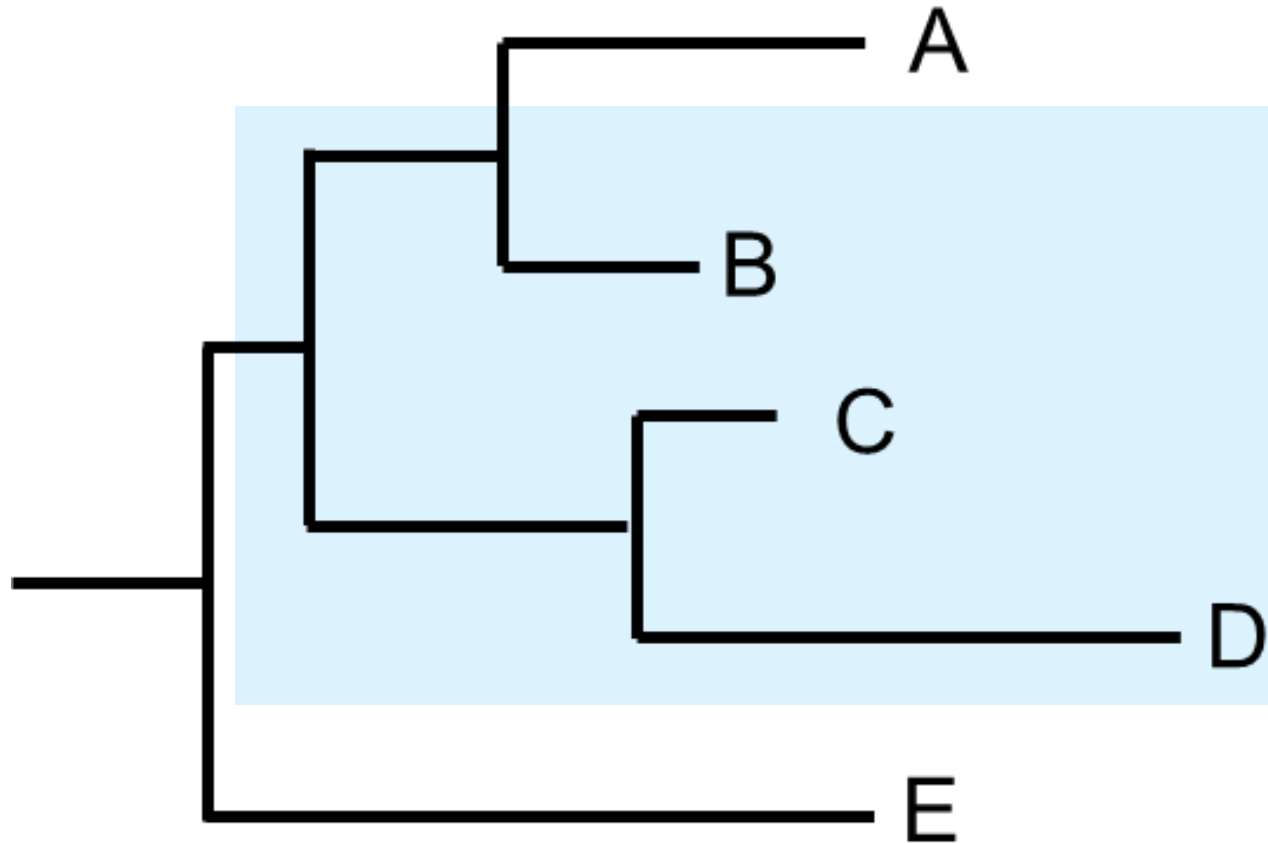
Entendendo filogenias



**São
CLADOS**

Monofilia

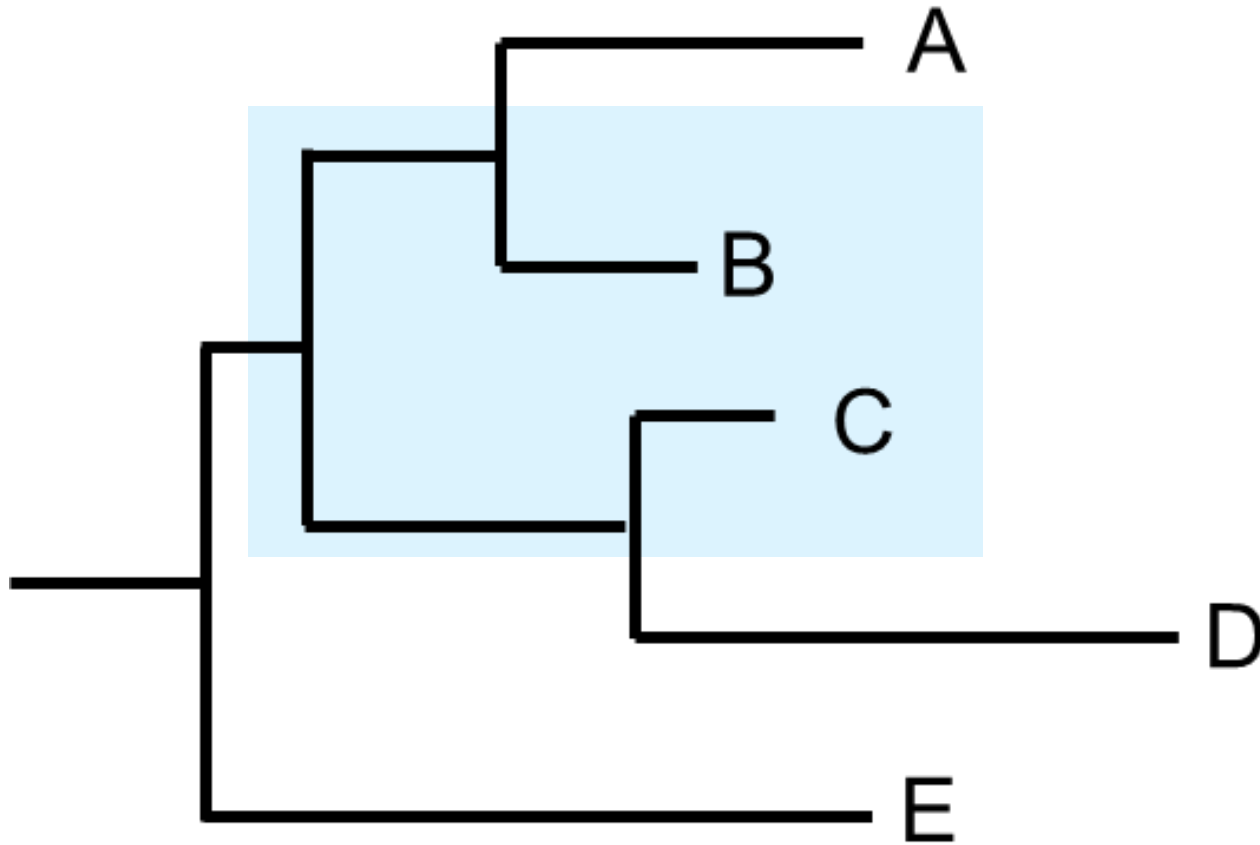
Entendendo filogenias



**NÃO
são
um
CLADO**

Parafilia

Entendendo filogenias



**NÃO
São
um
CLADO**

Polifilia