

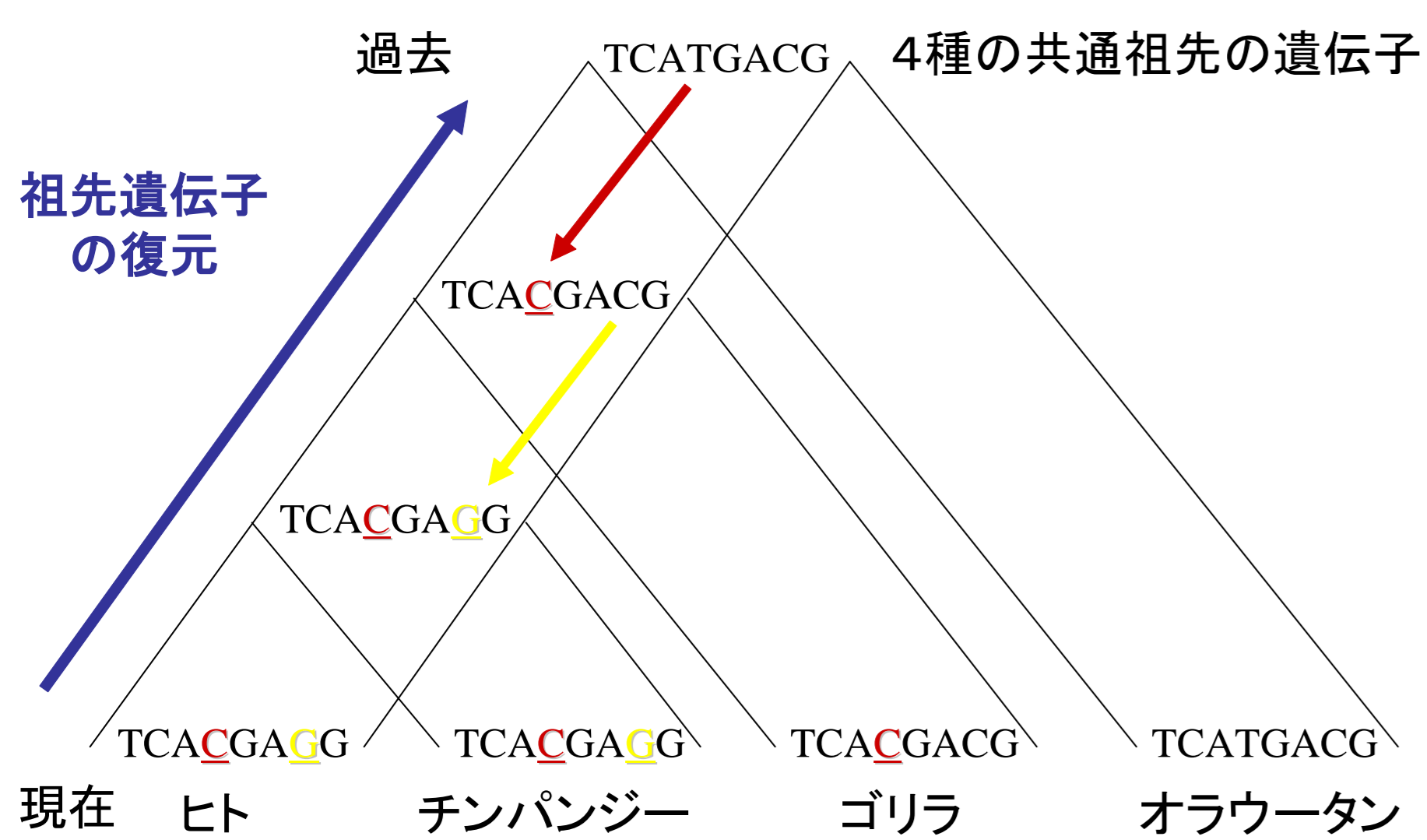
分子進化のモデリングと分子系統樹の推定

足立 淳 データ科学研究系 准教授

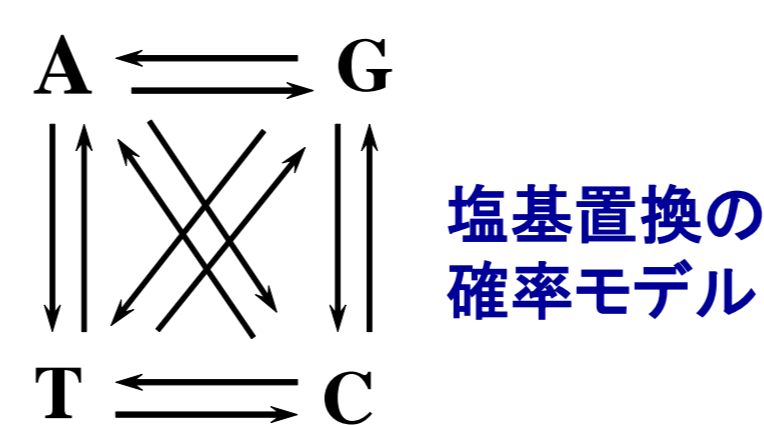
2014年6月13日 統計数理研究所 オープンハウス

分子系統学とは

遺伝子の変異差から種の系統を推定

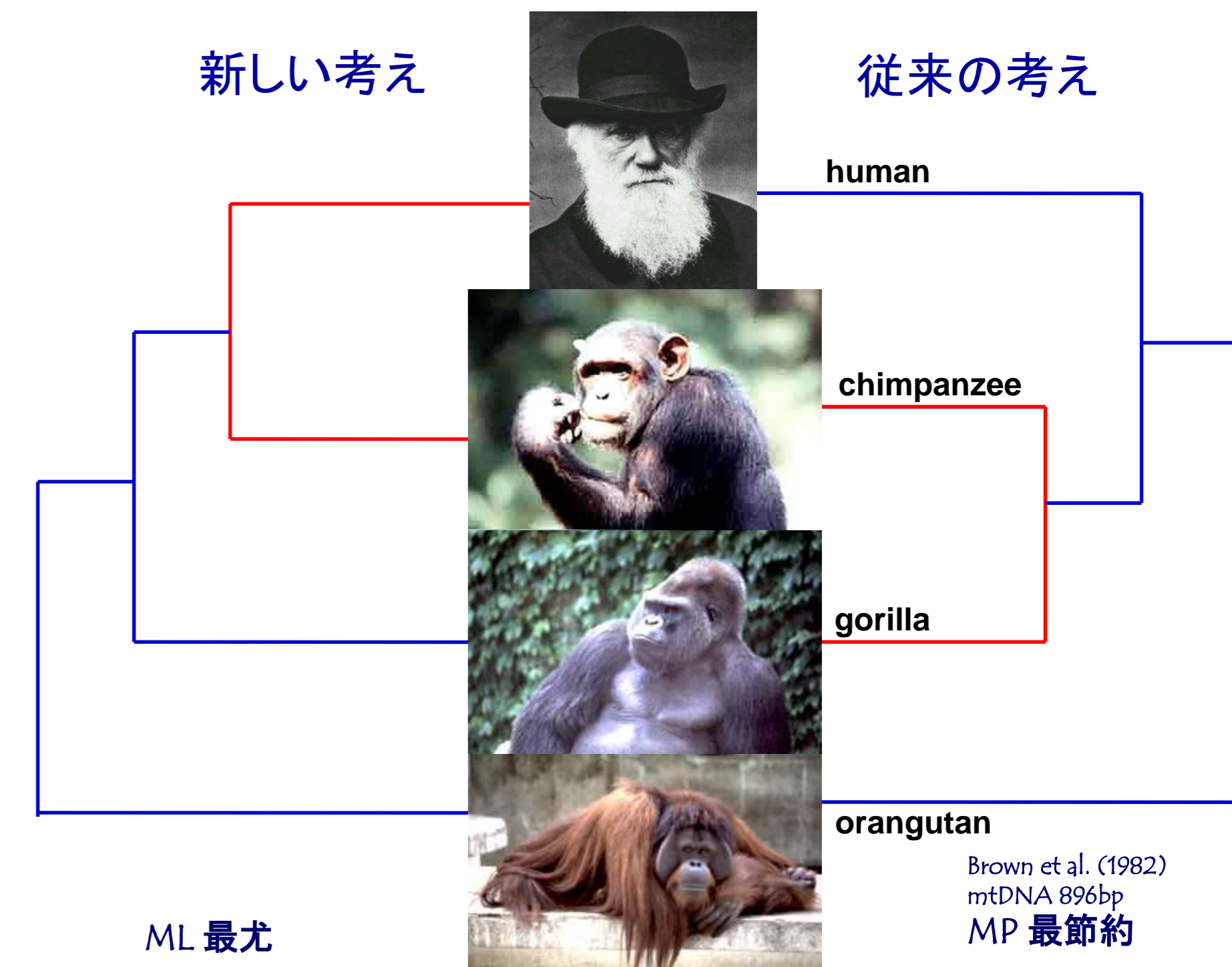


配列の違いは共通祖先から進化してきた結果生じたものであり、進化の歴史を反映している。
→ 最尤法による系統樹推定

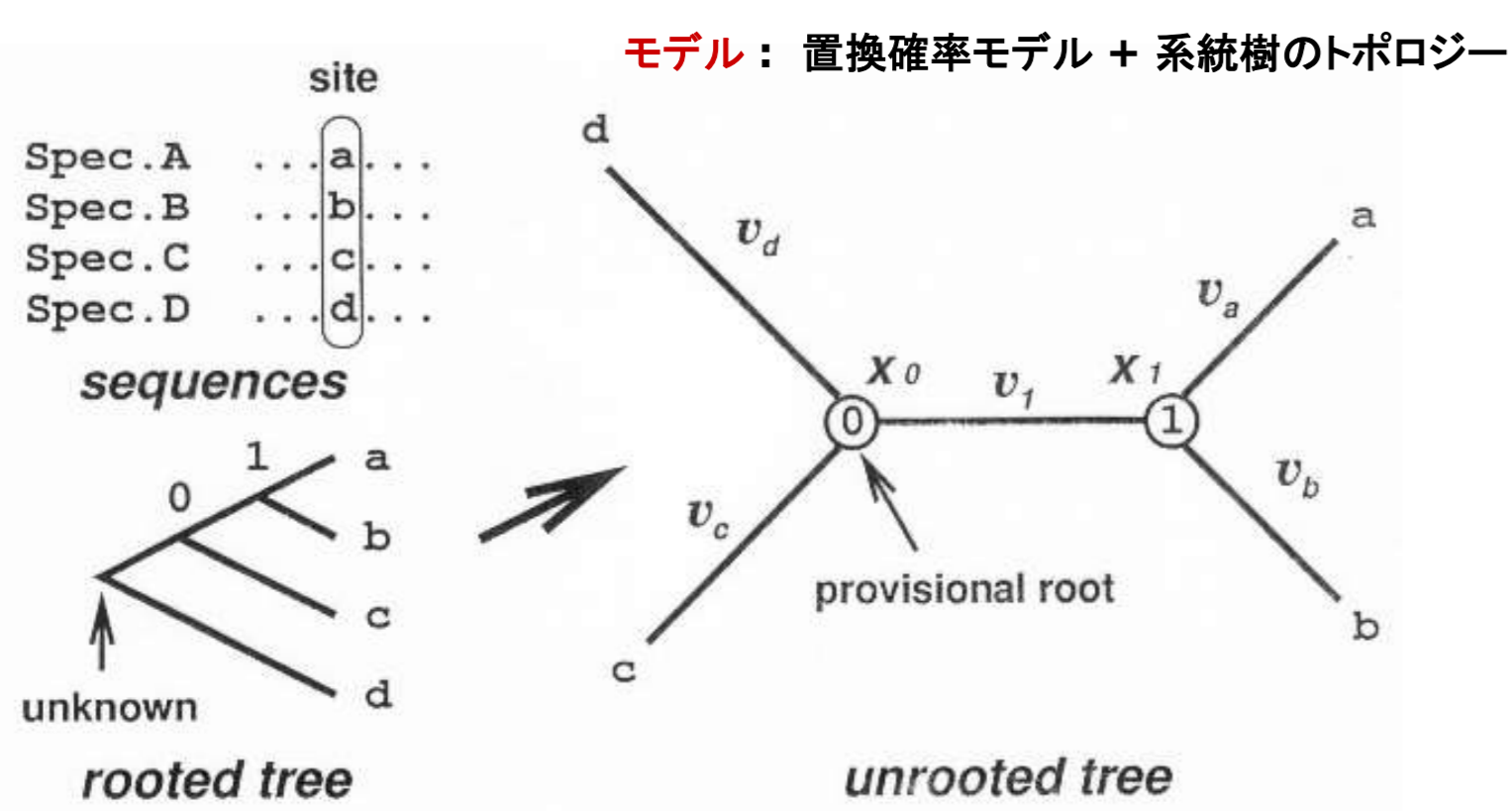


1. Human ヒト
2. Chimpanzee チンパンジー
3. Gorilla ゴリラ
4. Orangutan オランウータン

1 CTAGGCTATATACAACACTACGCAAAGGCCCAACGTTGTAGGCCCTAC
2 CTAGGCTACATACAACACTACGCAAAGGTCCCAACATTGTAGGTCCTTAC
3 TTAGGCTATATACAACACTACGTAAGGCCCAACGTCGTAGGCCCTAC
4 CTAGGCTATACACAACACTACGCAAGGACCTAACATCGTAGGCCCTAC



最尤法と尤度(Likelihood) $L = P(\text{data}|\text{model})$



Likelihood of a site

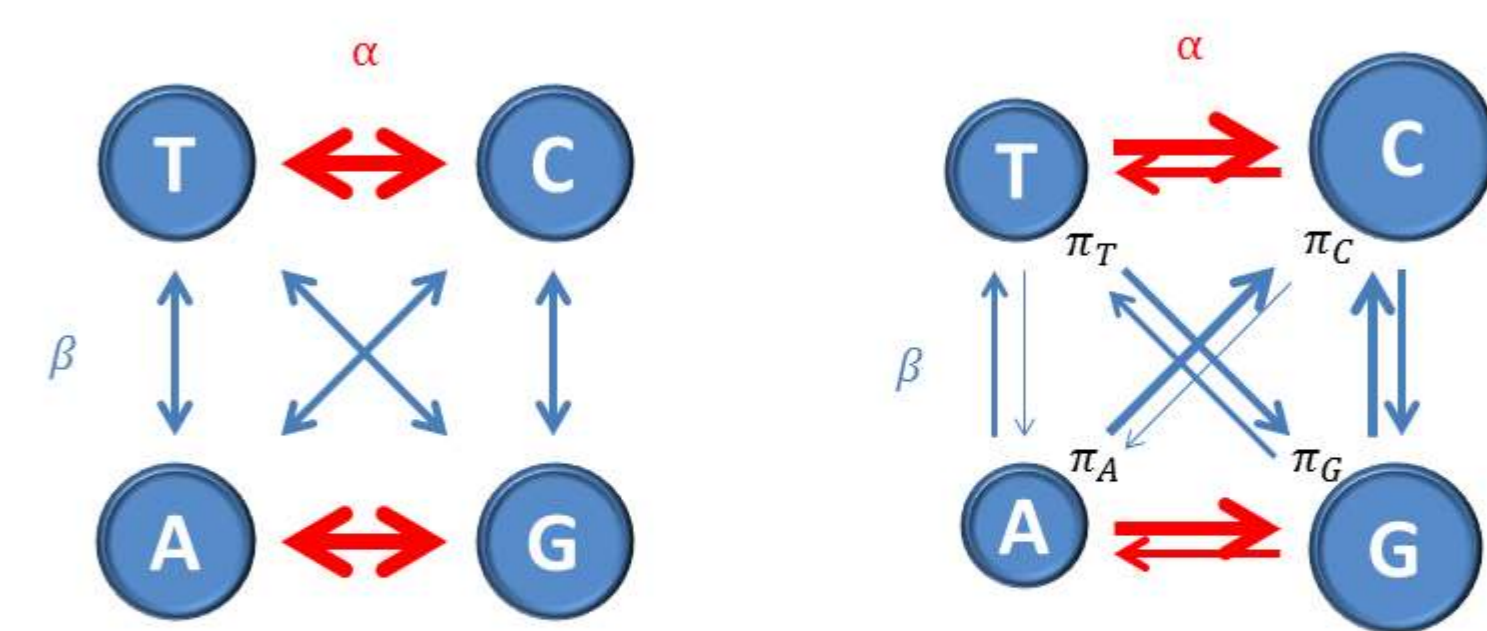
$$L = \sum_{x_0} \pi_{x_0} [P_{x_0c}(v_c) \times P_{x_0d}(v_d) \times \sum_{x_1} P_{x_0x_1}(v_1) P_{x_1a}(v_a) P_{x_1b}(v_b)]$$

置換モデル

塩基置換やアミノ酸置換のマルコフモデル
Transition probability matrix $P(t)$ during time t
 $P(t) = e^{Qt}$

- DNAモデル(塩基置換)
 - $A \leftrightarrow G, T \leftrightarrow C$ トランジション
 - $(A,G) \leftrightarrow (T,C)$ トランスバージョン
- タンパクモデル(アミノ酸置換)
 - 経験的に得られた 20×20 の置換配列
- コドンモデル
 - アミノ酸置換 と 同義塩基置換 の組み合わせ
 - 深い分岐の系統と浅い分岐の系統を同時推定

塩基置換モデル



Transition / Transversion model K80 Kimura's 2 parameter (1980)
Ts / Tv + Proportional model HKY85 Hasegawa, Kishino & Yano 1985

コドンモデル CTA へ9通りの置換

code	Amino acid	code	Amino acid	code	Amino acid	code	Amino acid
TTT	F Phe	TCT	S Ser	TAT	Y Tyr	TGT	C Cys
TTC		TCC		TAC		TGC	
TTA		TCG		TAA	---	TGA	---
TTG		TCG		TAG	---	TGG	W Trp
CTT	L Leu	CCT		CAT	H His	CGT	
CTC		CCG		CAC		CGC	
CTA		CCA	P Pro	CAA	Q Gln	CGA	R Arg
CTG		CCG		CAG		CGG	
ATT		ACT		AAT	N Asn	AGT	S Ser
ATC	I Ile	ACC		AAC		AGC	
ATA		ACA	T Thr	AAA		AGA	
ATG	M Met	GCT		GAT	K Lys	AGG	R Arg
GTT		AGC		AAG		GGT	
GTC		GCA	A Ala	GAC	D Asp	GGC	
GTA	V Val	GAA		GAA		GGG	G Gly
GTG		GCG		GAG	E Glu	GGG	
T		C		A		G	

EF1α のアミノ酸配列: 動物、菌類、植物、原生動物、細菌で共通の配列が見られ

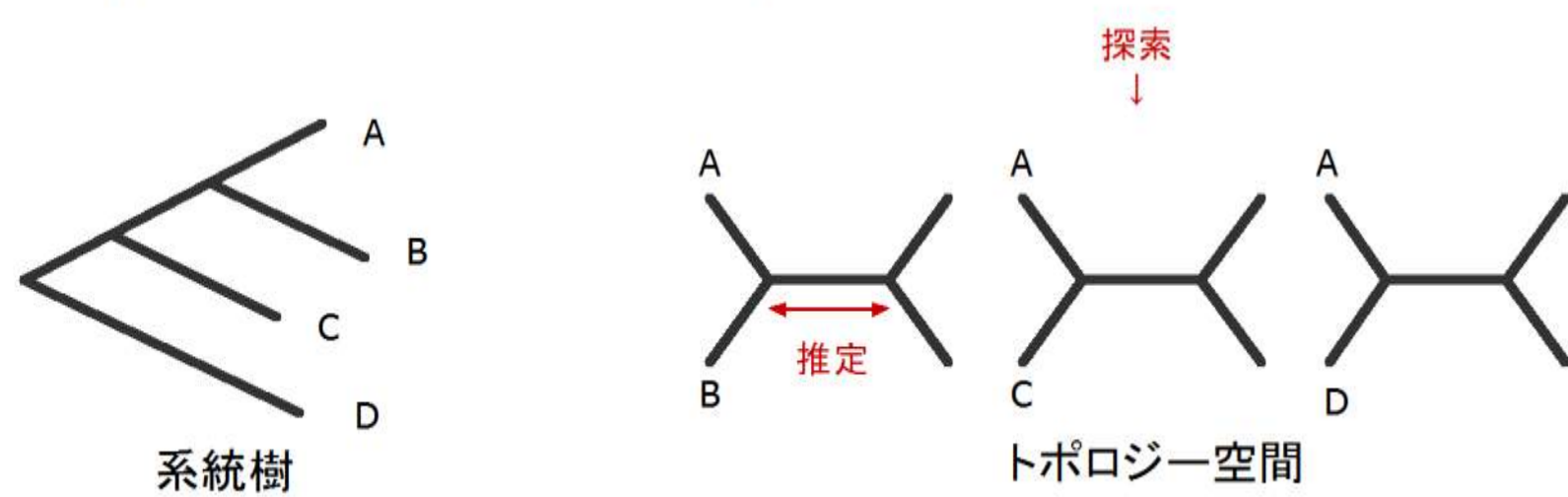
CONSENSUS	STTTTTHLYK	CGHDKRTIE	KFEKEAAE	GS	KGSFYAVWL	IKLKAEREER	ITIDIALWKF	ET	KY	VT	I	DPQGHDFDK
Homo sapiens
Gallus gallus
Xenopus laevis
Danio rerio
Apis mellifera
Bombus morio
Oncocercus
Saccharomyces
Ashbya gossypii
Candida albica
Trichoderma re
Podospora anae
Puccinia grami
Asbidia glauca
Arabidopsis th
Glycine max
Hordeum vulgare
Triticum aestiv
Trichomonas ta
Giardia lamblia
Hexamita infla
Glugea pleocoi
Sulfolobus sol
Halobacterium
Methanococcus

アミノ酸置換モデル mtREV

A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	31	135	69	328	41	63	727	84	595	165	1	772	35	321	2229	2078	1	41	818
R	135	135	1	1102	1	152	455	1	71	985	1	34	178	38	34	146	1	51	
N	69	1	4527	1	261	3495	334	828	88	4	175	1	35	1	319	98	67	893	
D	328	1	30	1	242	1	214	723	256	247	1	1	547	101	1788	1046	222	897	
Q	41	1102	731	261	242	1	1798	83	3208	72	281	2331	349	160	801	447	524	1	
E	63	1	503	3495	1	1798	191	306	1	5	1630	1	1	54	286	138	1	115	
G	727	152	314	334	214	83	191	1	40	10	89	5	1	1	651	73	49	1	
H	84	455	2656	828	723	3208	306	1	43	62	470	13	177	273	341	295	47	3171	
I	595	1	132	88	256	72	1	40	43	1875	65	2656	441	117	280	1805	1	200	
L	165	71	151	4	247	281	5	149	62	1875	37	2938	1112	184	388	717	190	231	
K	1	985	2712	175	1	2631	1630	89	470	65	37	440	52	293	765	808	183	233	
M	772	34	117	35	547	160	1	1	177	441	1112	52	546	93	310	267	57	2534	
F	321	178	508	1	101	801	54	1	273	117	184	293	89	93	720	732	33	102	
S	2229	38	2607	319	1788	447	286	651	341	280	388	765	714	310	720	3426	148	218	
T	2078	34	1271	98	1046	524	138	73	295	1805	717	808	2547	267	732	3426	125	158	
W	1	146	67	63	222	1	1	49	47	1	390	133	209	57	33	148	125	151	
Y	41	1	893	74	897	260	115	1	3171	200	231	233	297	2534	102	218	158	151	
V	818	51	76	1	1	78	137	20	1	6030	465	6	1745	34	48	1	1172	35	

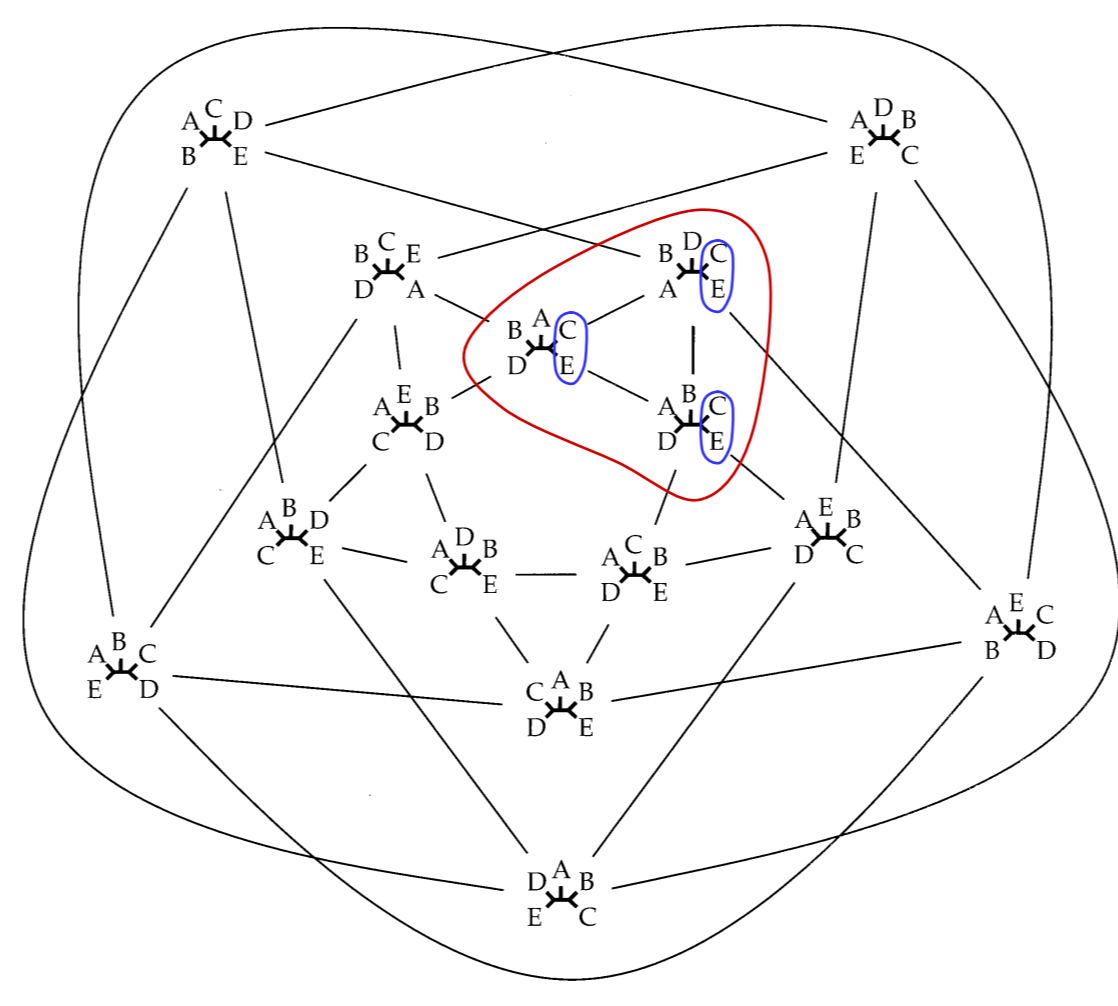
系統樹 = トポロジー + 距離

- トポロジー = 枝分かれの順番
トポロジー空間から探索
- 距離 = 枝の長さ (推定置換数)
置換確率モデルから計算



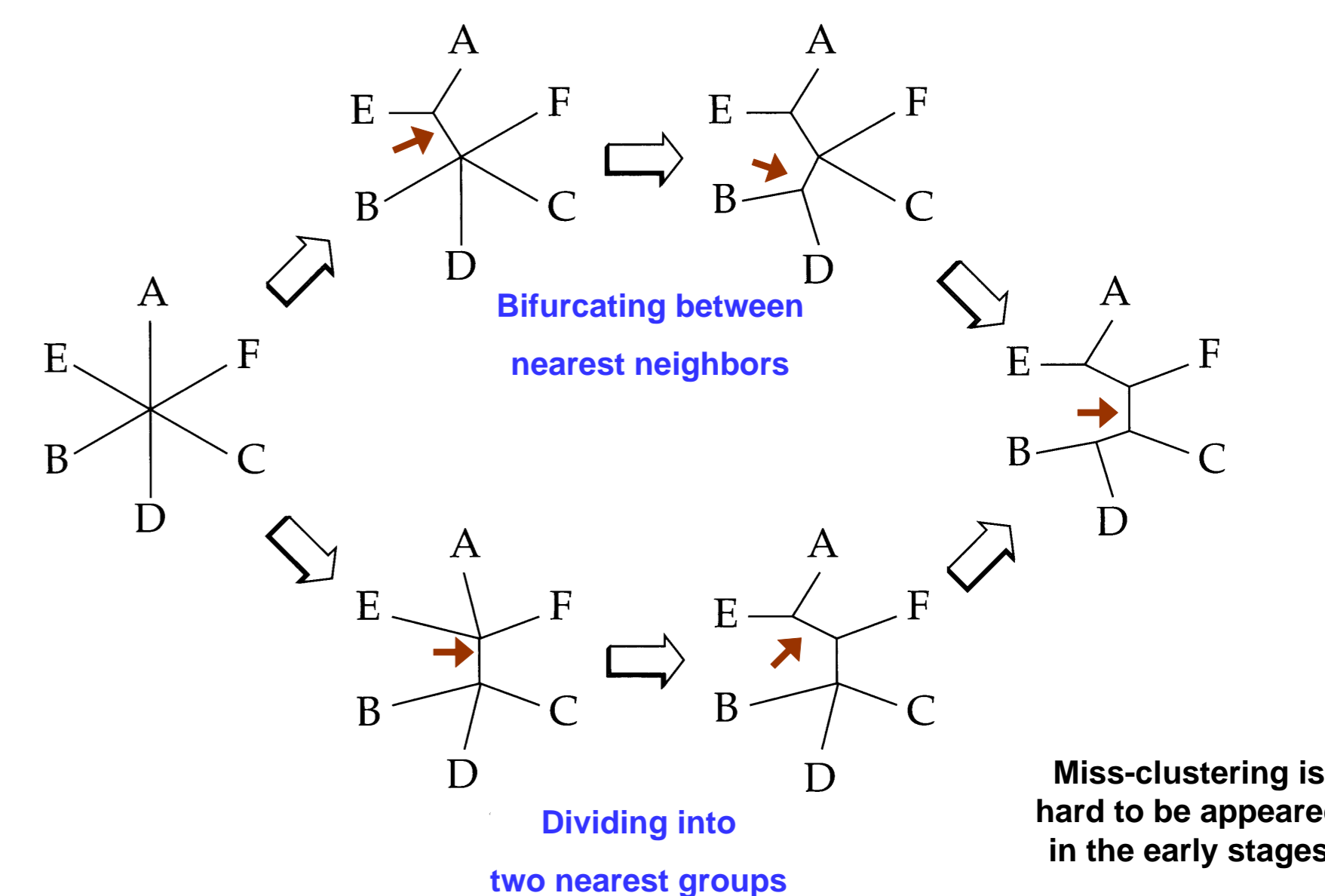
トポロジー空間と探索法

最近隣木間の関係



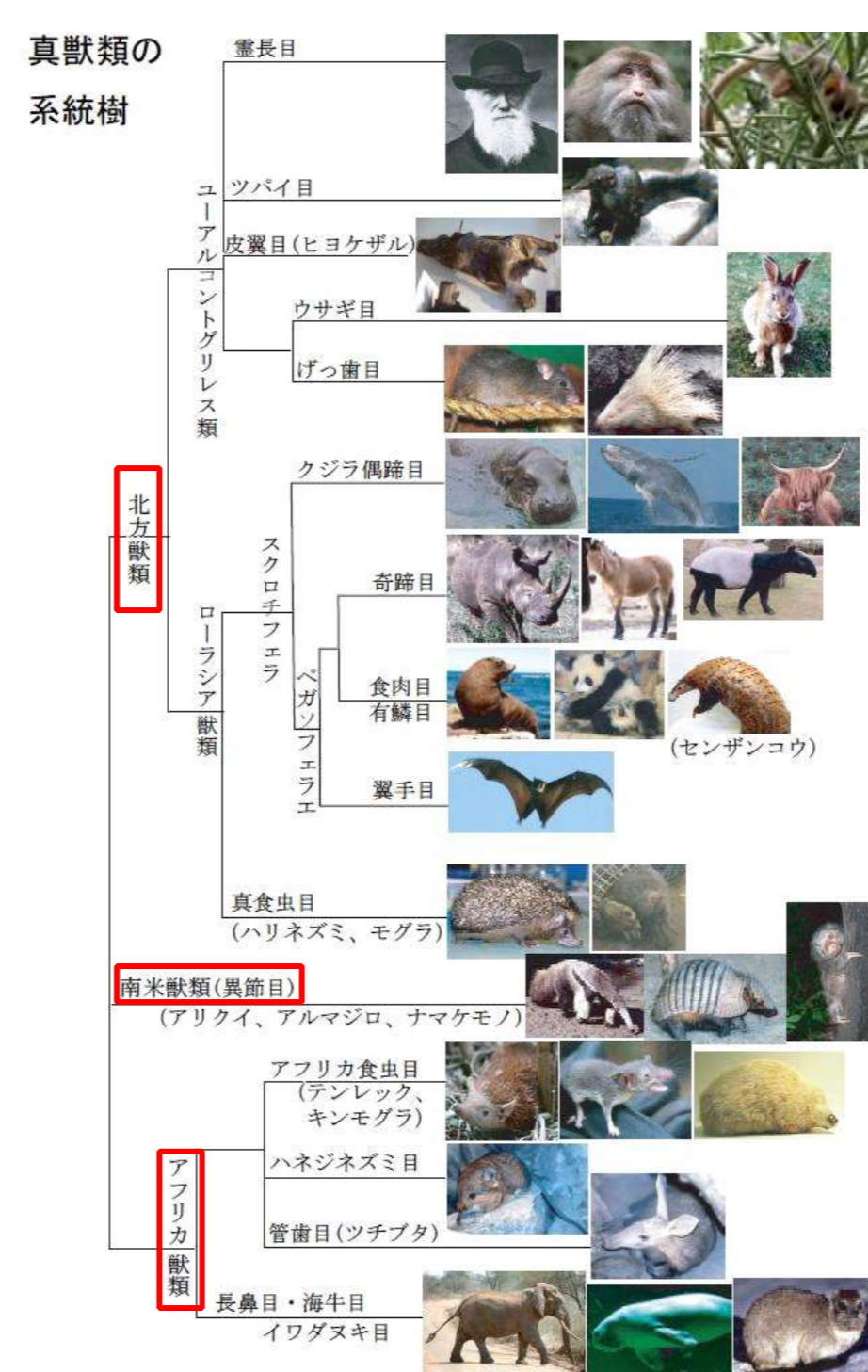
The space of all 15 possible unrooted trees with 5 tips

New star decomposition



Miss-clustering is hard to be appeared in the early stages

1億年前の大陸配置と真獣類の進化



遺伝子系統 × ゲノム構造比較 = 進化ダイナミクス

