

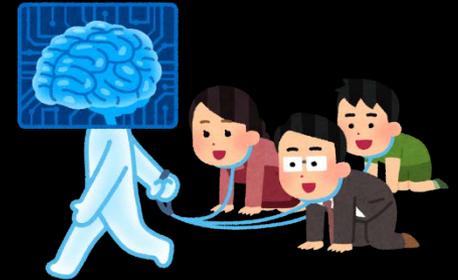
# 言語系統論への 計算的アプローチの可能性

京都大学  
村脇 有吾



# 自己紹介：村脇 有吾

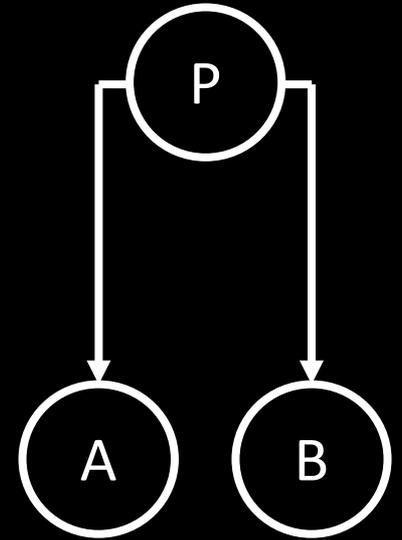
- 京都大学 大学院情報学研究科 知能情報学専攻 助教  
工学部電気電子工学科 兼任
- 専門：計算言語学と自然言語処理
  - 人工知能バブル
  - 言語処理学会：年次大会参加者 >900
  - 2000年代から統計・機械学習を本格導入
  - 表の仕事：ゼロ照応解析、談話構造解析、テキストからの知識獲得ほか
  - 今日前半で紹介する研究は生物学由来で、実は専門外



# かつての言語年代学 (Glottochronology)

$$t = \frac{\log c}{2 \log r}$$

- $t$ : 祖語Pの年代 (単位: 千年)
- $c$ : 言語A, Bの基礎語彙共有率
- $r$ : 基礎語彙の残存率 (200項目で0.81)

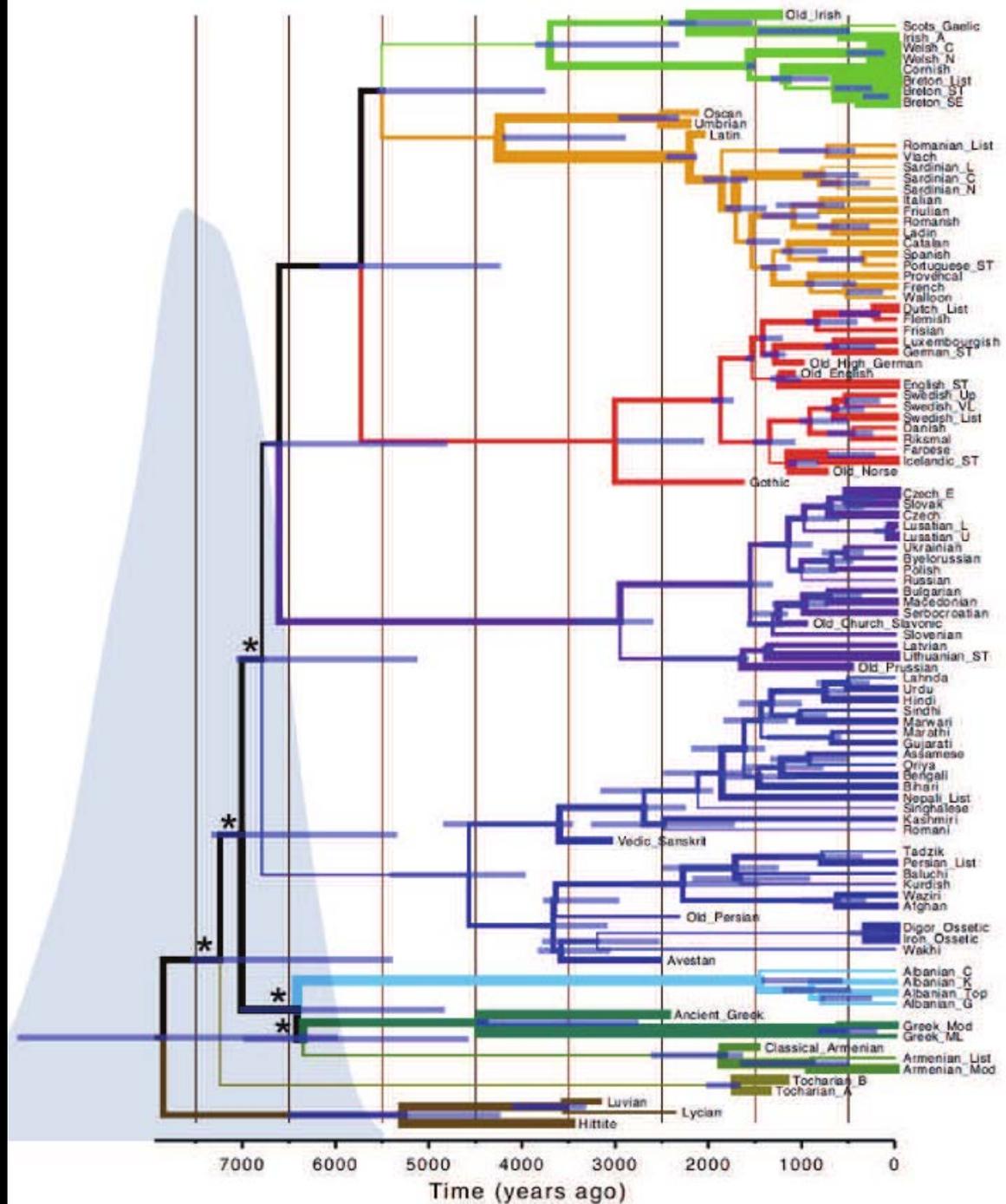


年代	1K	2K	3K	4K	5K	6K	7K	8K	9K
共有率 ( $r=.81$ )	.66	.43	.28	.19	.12	.08	.05	.03	.02

# 今日の目標: →の作り方の 雰囲気をつかむ

- ベイズ統計による確率モデル
- 最大節約、近隣結合は取り上げない

[Bouckaert+, Science 2012]



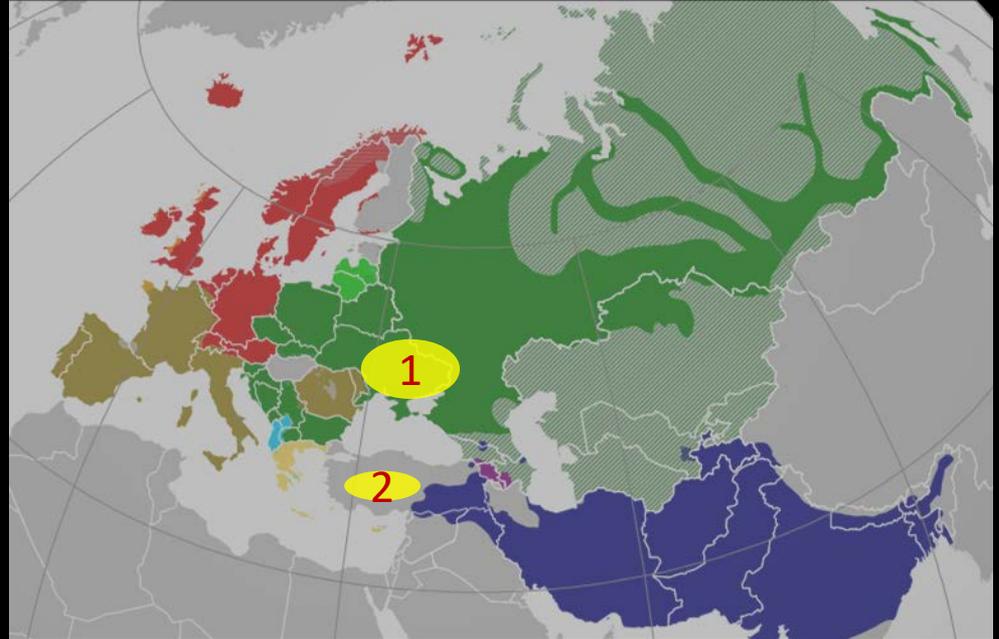
# 印欧祖語の年代と故地

## 1. クルガン仮説

- 5,000-6,000年前
- 黒海周辺のステップ
- 遊牧民の軍事的征服

## 2. アナトリア仮説

- 8,000-9,500年前
- アナトリア
- 農耕とともに拡大
  - Renfrew (考古学者) の農耕・言語同時伝播モデル



Source:  
Map:  
- Author: Alphathon  
- License: CC BY-SA 3.0

# 印欧祖語の年代と故地

## 1. クルガン仮説

- 5,000-6,000年前
- 黒海周辺のステップ
- 遊牧民の軍事的征服

## 2. アナトリア仮説

- 8,000-9,500年前
- アナトリア
- 農耕とともに拡大
  - Renfrew (考古学者) の農耕・言語同時伝播モデル



Source:  
Map:  
- Author: Alphonso  
- License: CC BY-SA 3.0

# 印欧祖語の年代と故地

## 1. クルガン仮説

- 5,000-6,000年前
- 黒海周辺のステップ
- 遊牧民の軍事的征服

## 2. アナトリア仮説

- 8,000-9,500年前
- アナトリア
- 農耕とともに拡大
  - Renfrew (考古学者) の農耕・言語同時伝播モデル



Source:  
Map:  
- Author: Alphathon  
- License: CC BY-SA 3.0

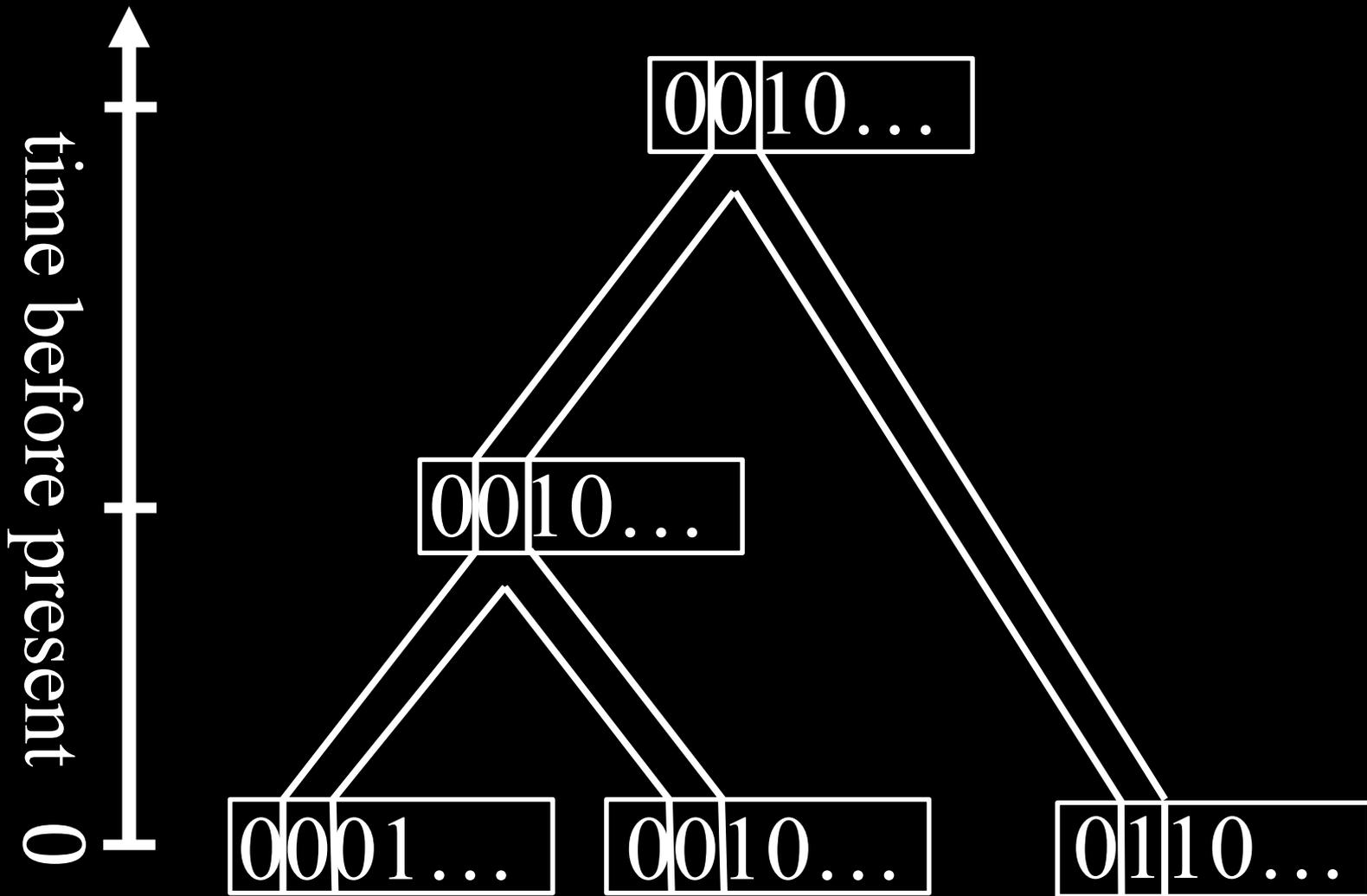
# 年代つき系統樹

0001...

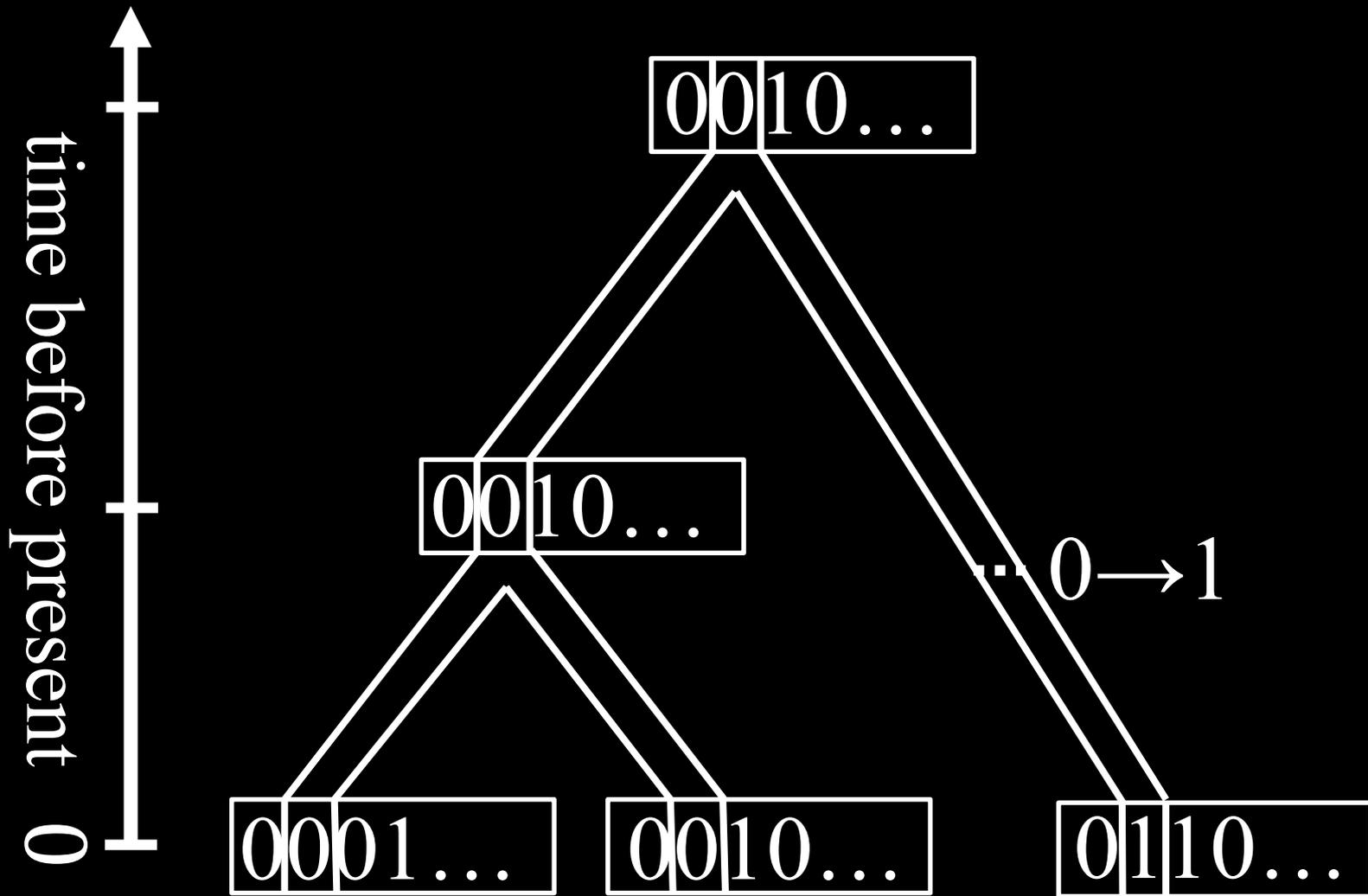
0010...

0110...

# 年代つき系統樹



# 年代つき系統樹



# 言語の状態：符号化された語彙

	WATER	BIG
英語	water	big
ドイツ語	Wasser	gross
ロシア語	вода	большой    великий
フランス語	eau	grand
イタリア語	acqua	grande

# 言語の状態：符号化された語彙

同源語群 WATER

BIG

英語

water

big

ドイツ語

Wasser

gross

ロシア語

вода

большой

великий

フランス語

eau

grand

イタリア語

acqua

grande

# 言語の状態: 符号化された語彙

同源語群 WATER

BIG

英語	1	water	3	big	{1, 3}		
ドイツ語		Wasser	4	gross	{1, 4}		
ロシア語		вода	5	большой	6	великий	{1, 5, 6}
フランス語	2	eau	7	grand	{2, 7}		
イタリア語		acqua		grande	{2, 7}		

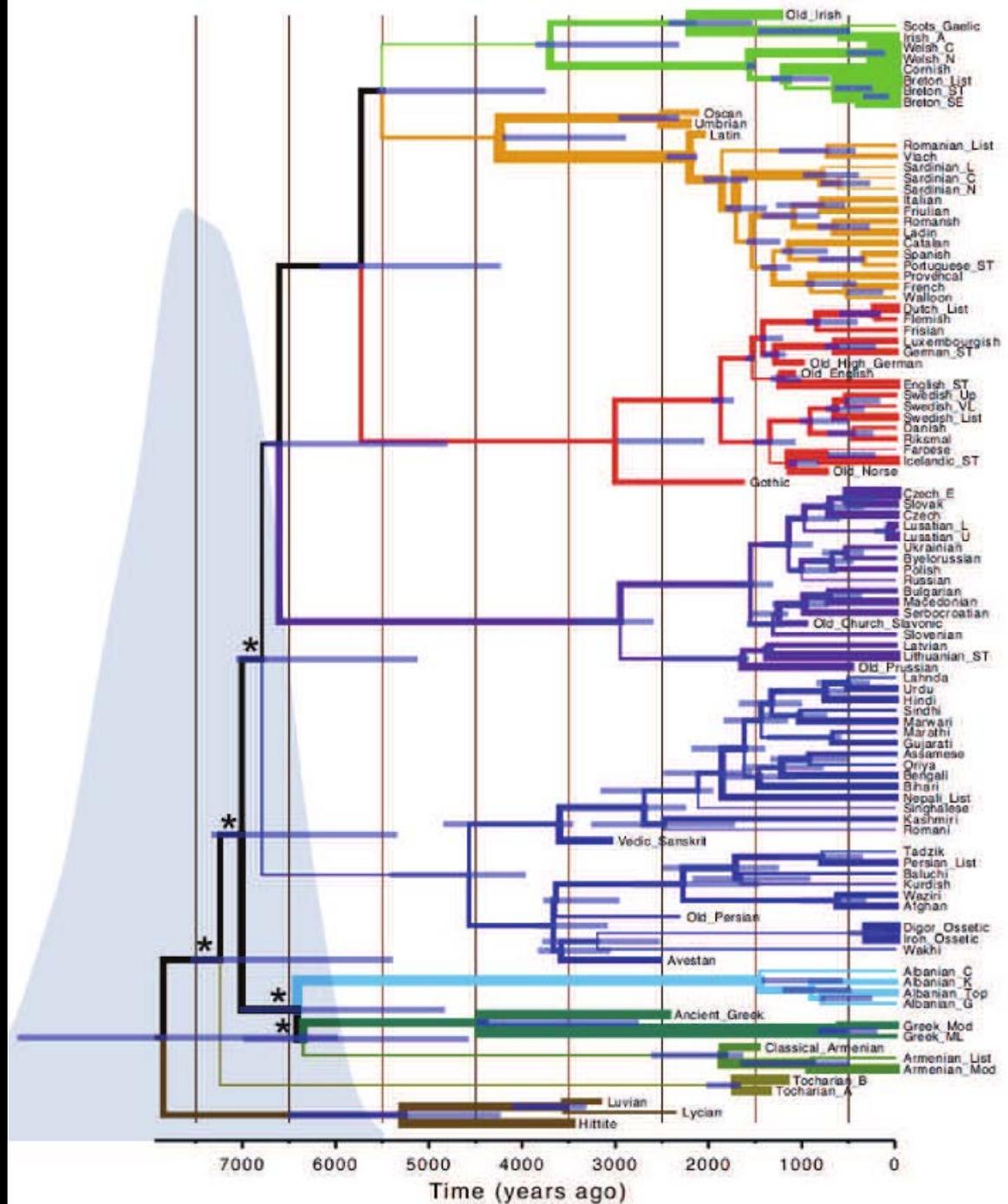
# 言語の状態: 符号化された語彙

同源語群 WATER

BIG

英語	1	water	3	big		1010000	
ドイツ語		Wasser	4	gross		1001000	
ロシア語		вода	5	большой	6	великий	1000110
フランス語	2	eau	7	grand		0100001	
イタリア語		acqua		grande		0100001	

# 問題設定



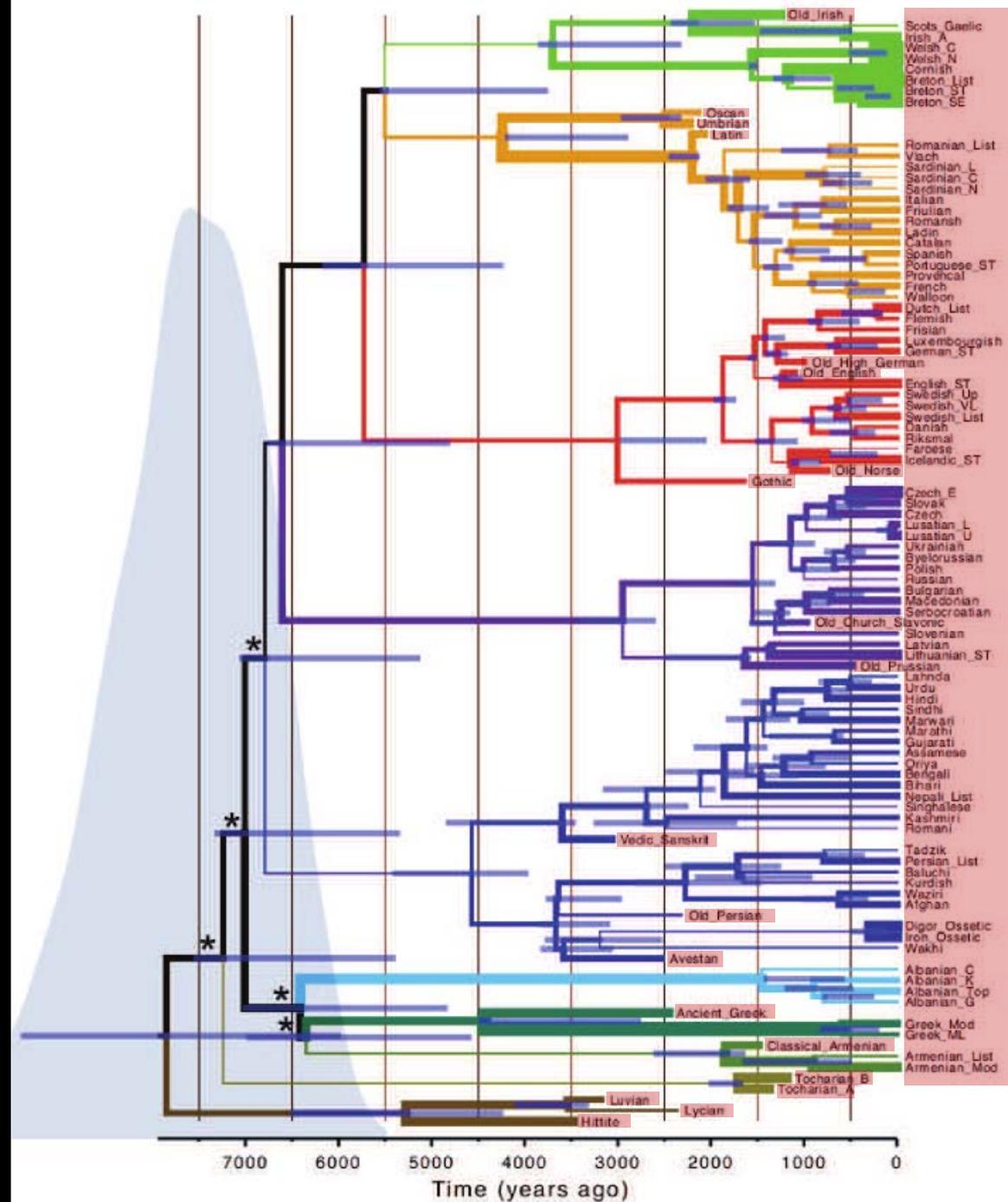
[Bouckaert+, Science 2012]

From Bouckaert et al. 2012. Mapping the Origins and Expansion of the Indo-European Language Family. Science 337(6097). Reprinted with permission from AAAS.

# 問題設定

## 観測データ

- 葉ノードの状態



[Bouckaert+, Science 2012]

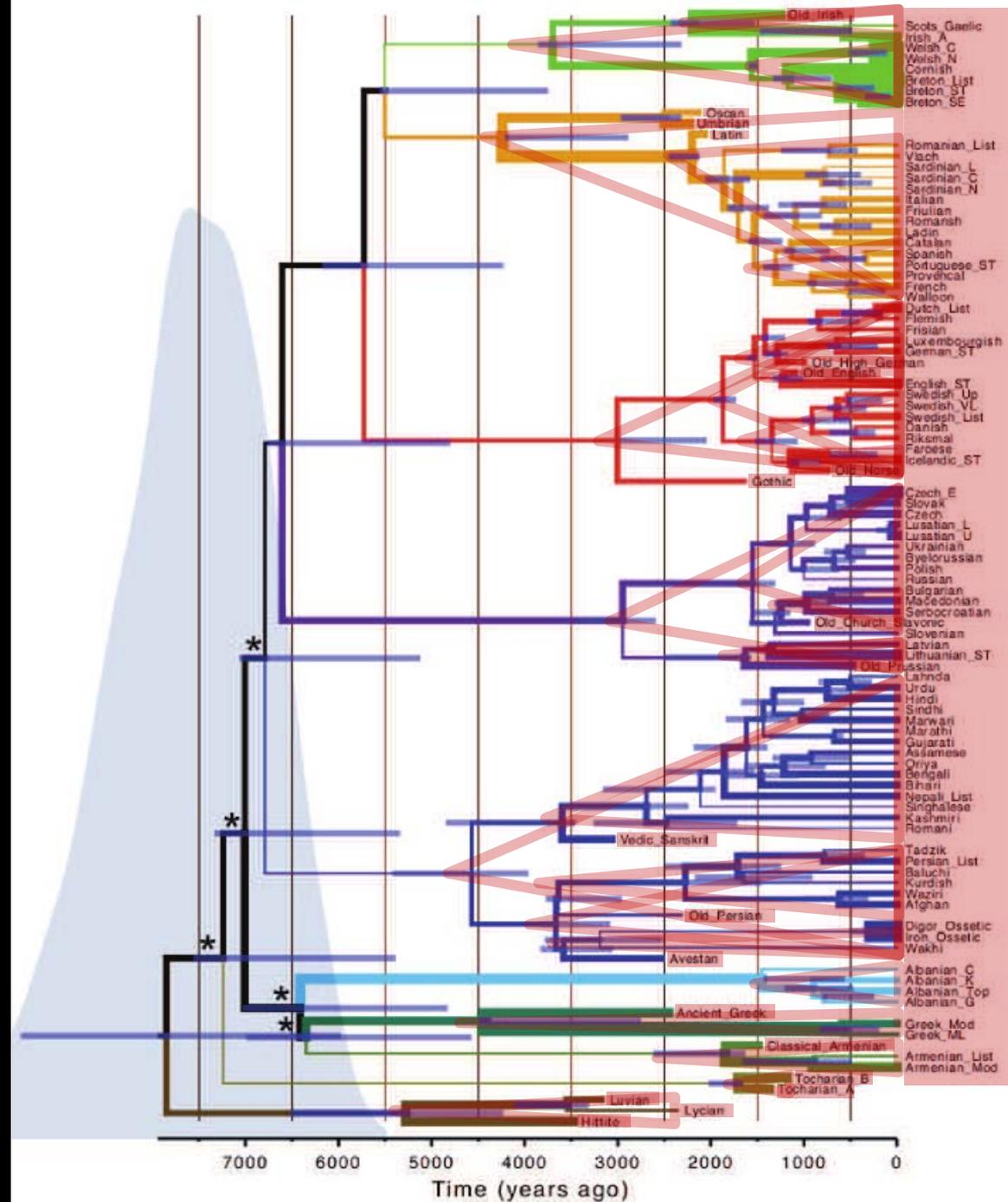
# 問題設定

## 観測データ

- 葉ノードの状態

## 制約

- いくつかの単系統群



[Bouckaert+, Science 2012]

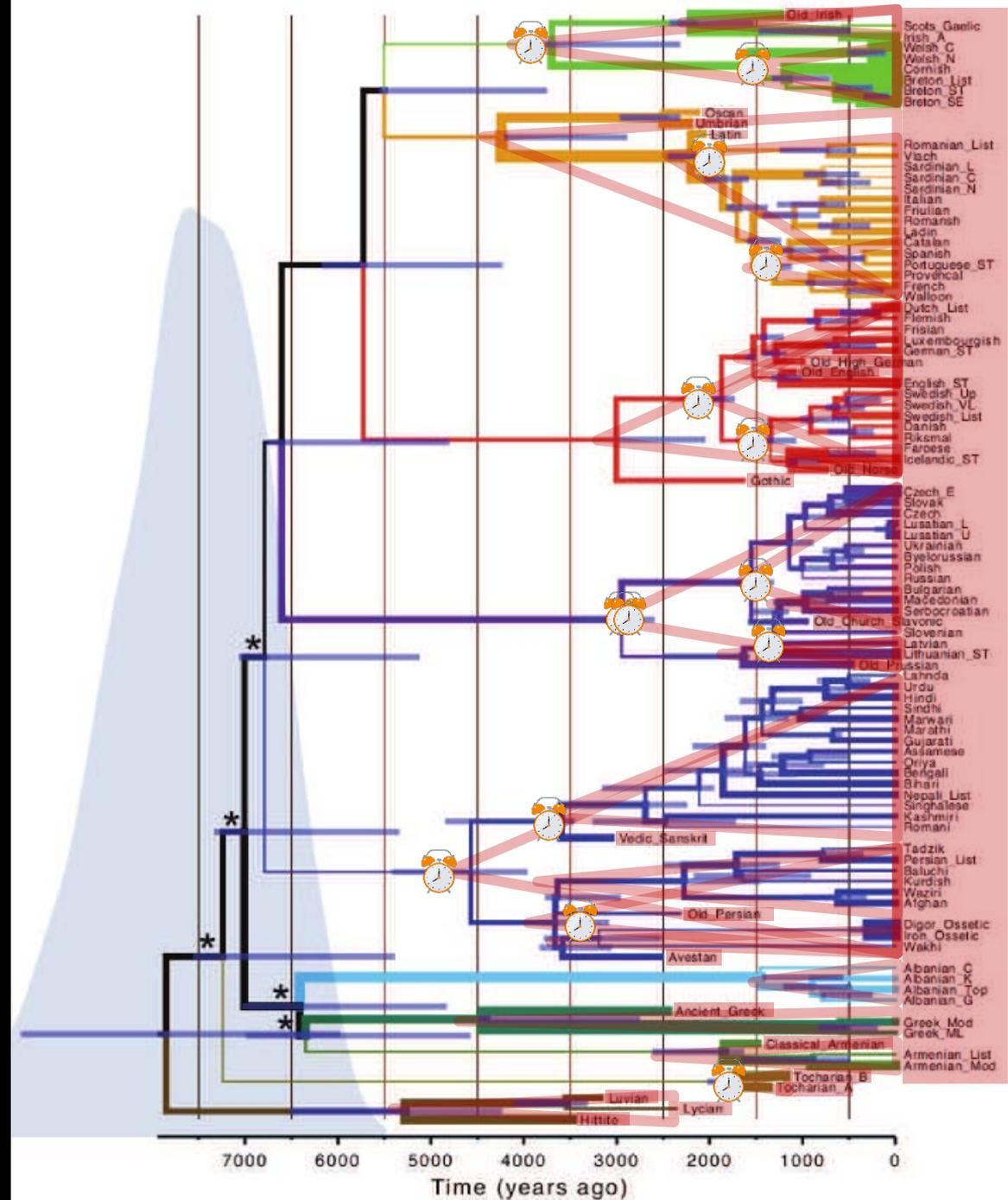
# 問題設定

## 観測データ

- 葉ノードの状態

## 制約

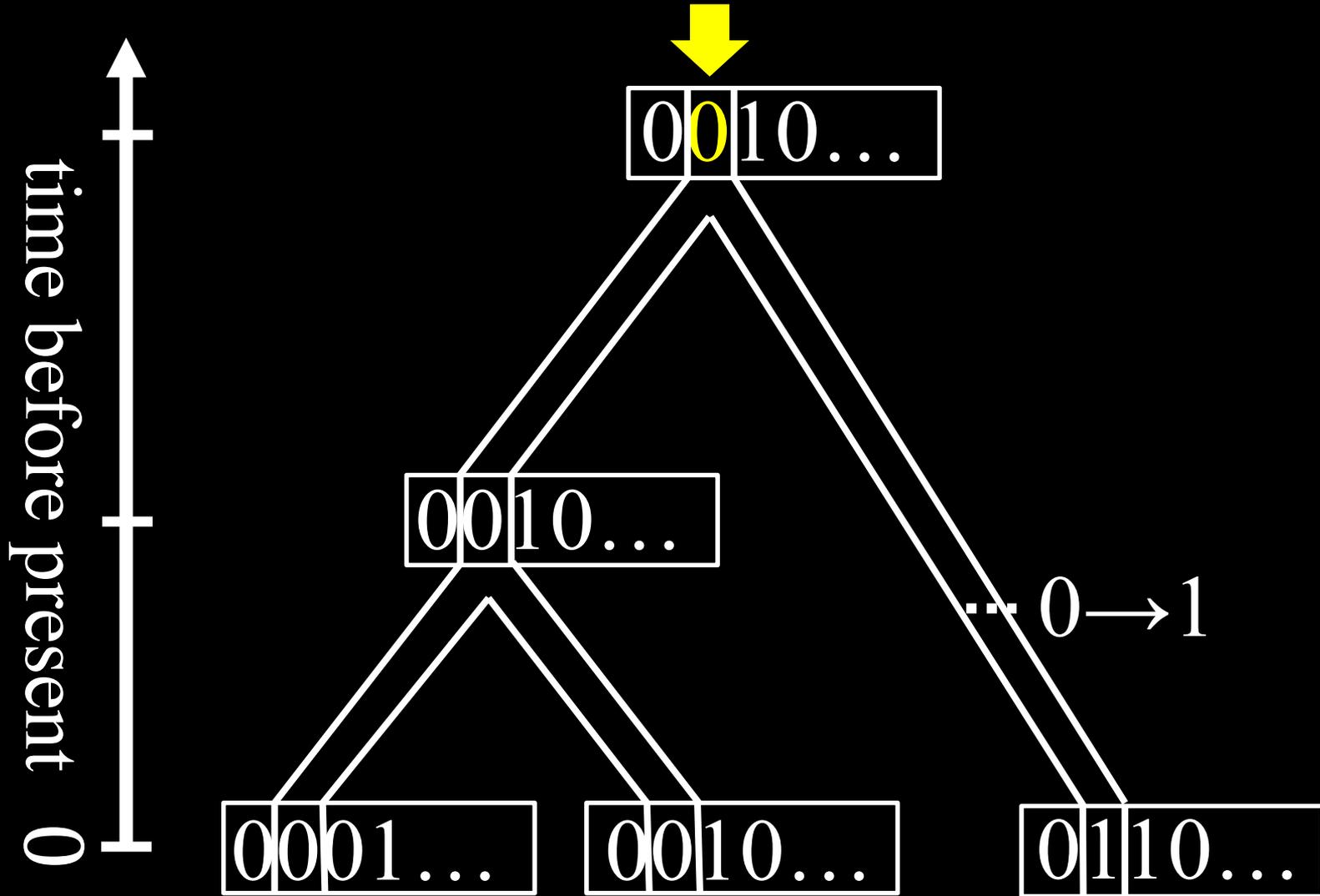
- いくつかの単系統群
- いくつかの系統群のおおよその年代



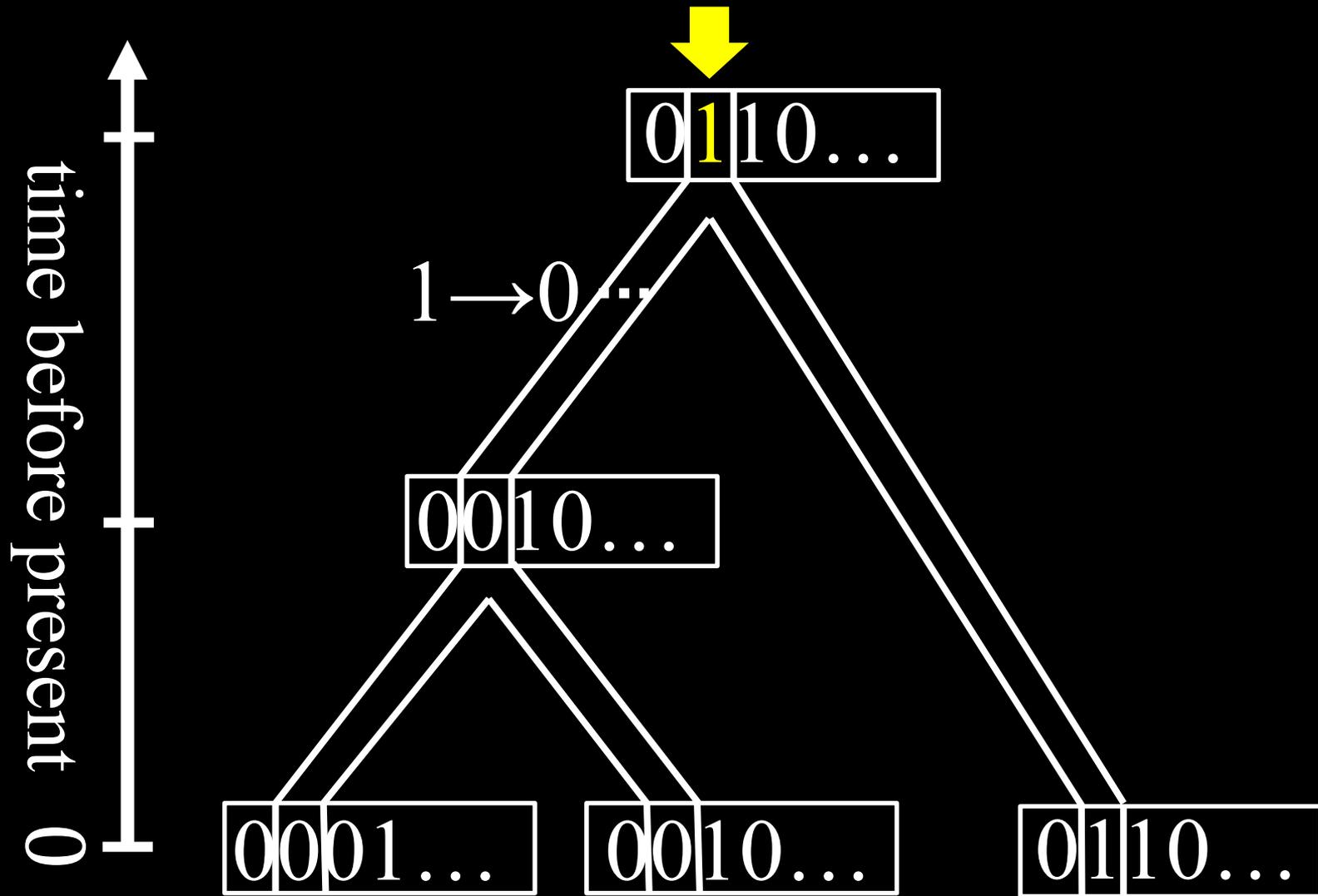
[Bouckaert+, Science 2012]



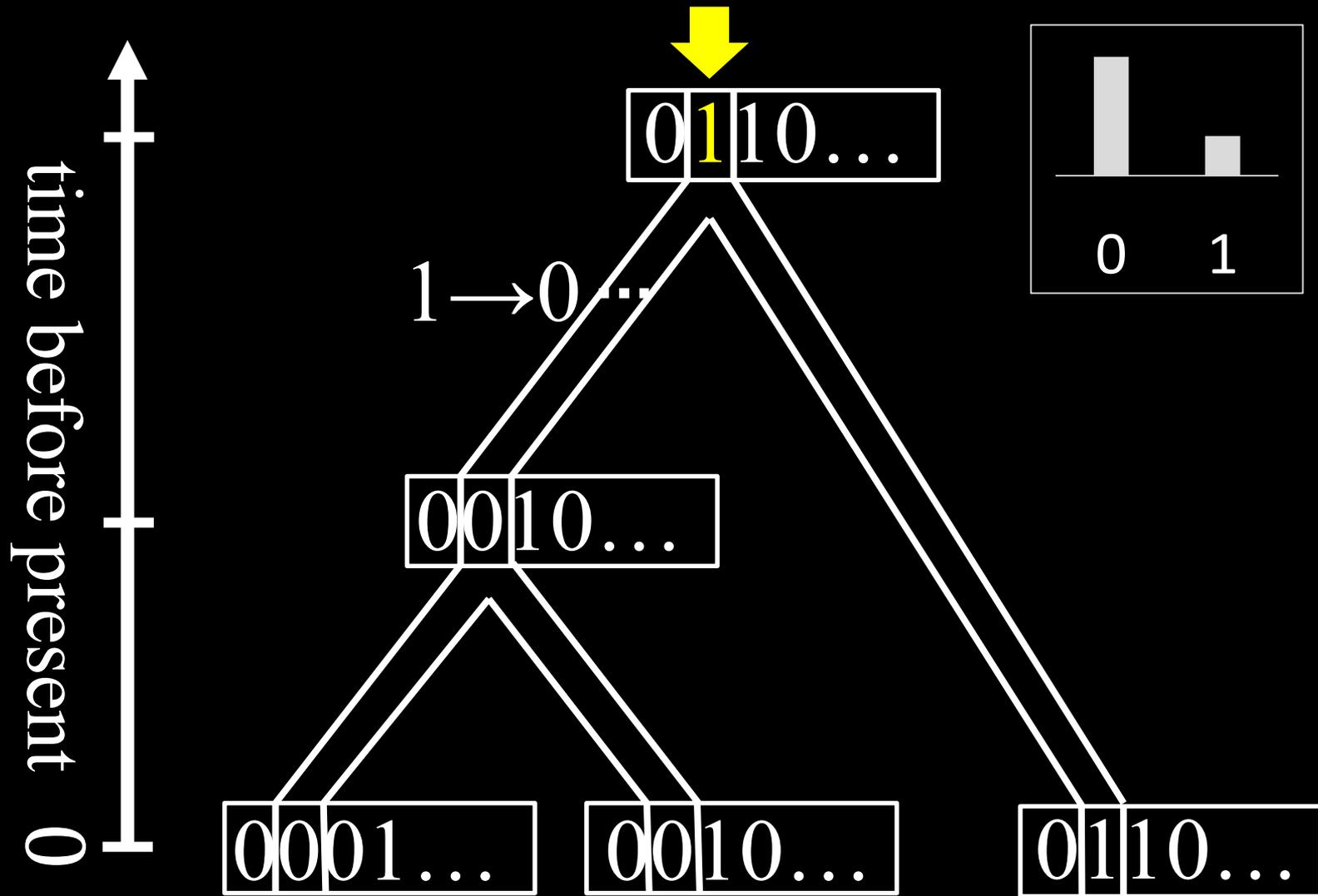
Q: 手持ちの手がかりだけでは  
トポロジー/状態/年代を決めようがないのでは?



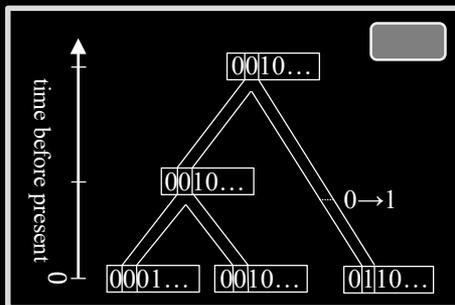
Q: 手持ちの手がかりだけでは  
トポロジー/状態/年代を決めようがないのでは?



Q: 手持ちの手がかりだけでは  
トポロジー/状態/年代を決めようがないのでは?

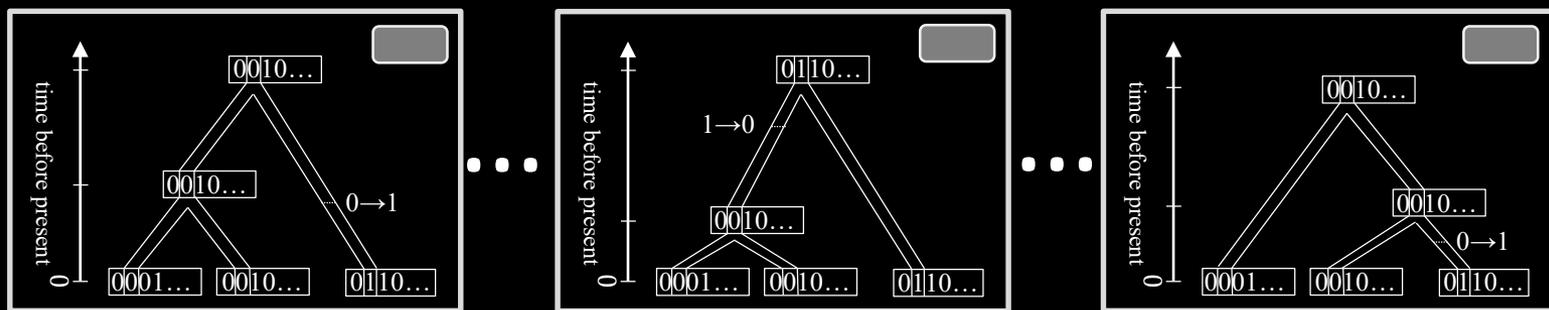


# マルコフ連鎖モンテカルロ法 (MCMC)



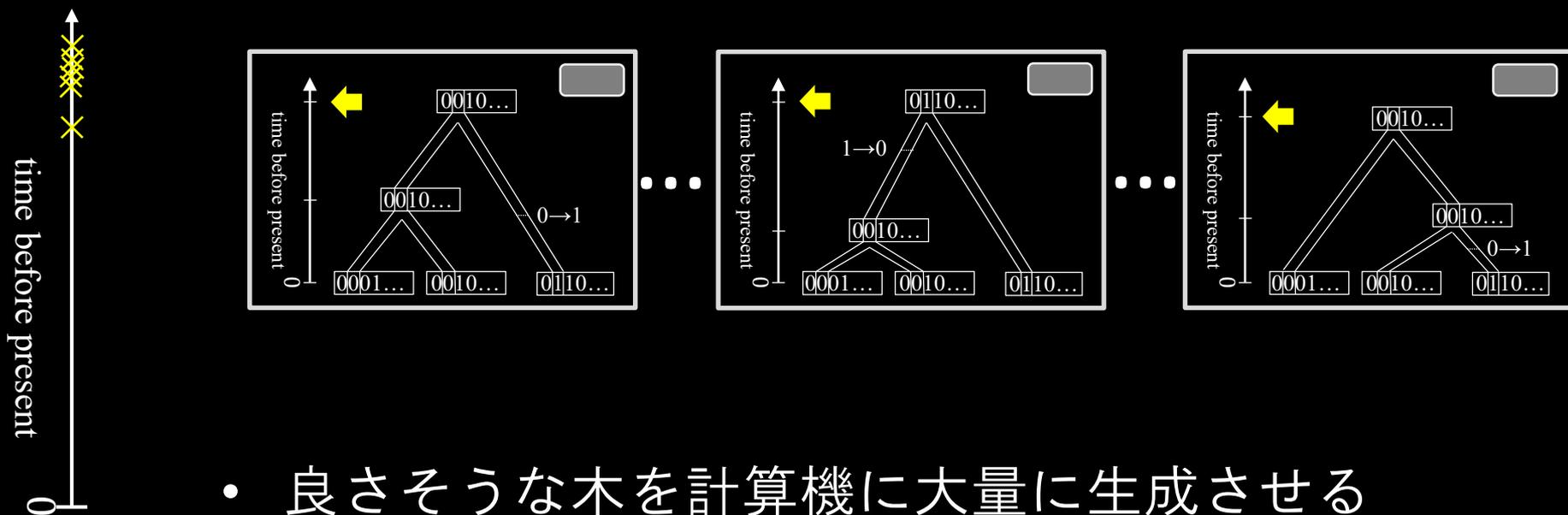
- 良さそうな木を計算機に大量に生成させる
  - 観測データは固定
  - 推論対象データを一定の条件にしたがってランダムに動かす

# マルコフ連鎖モンテカルロ法 (MCMC)



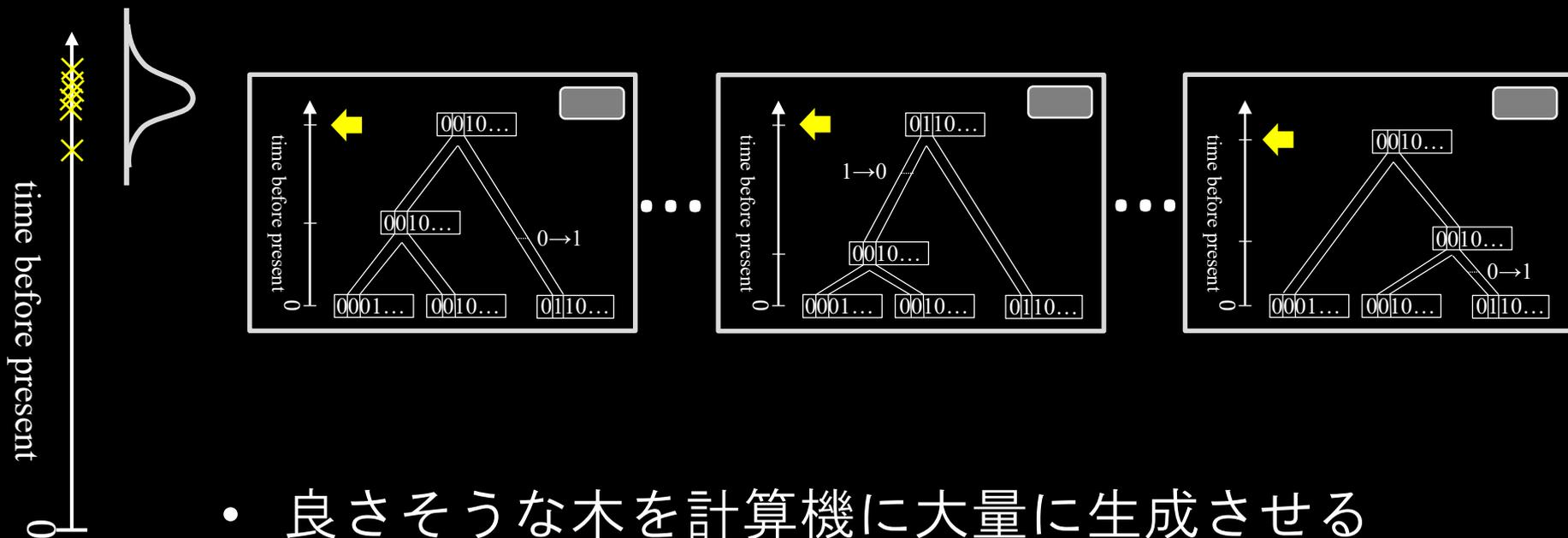
- 良さそうな木を計算機に大量に生成させる
  - 観測データは固定
  - 推論対象データを一定の条件にしたがってランダムに動かす

# マルコフ連鎖モンテカルロ法 (MCMC)



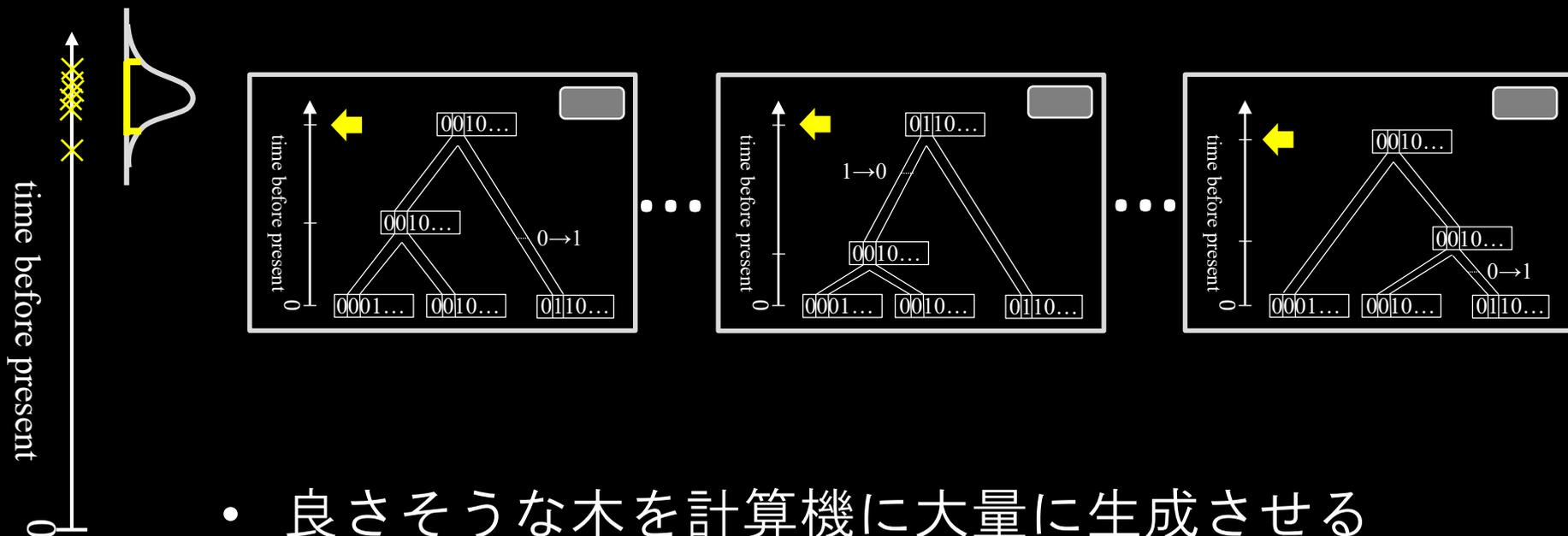
- 良さそうな木を計算機に大量に生成させる
  - 観測データは固定
  - 推論対象データを一定の条件にしたがってランダムに動かす
- 生成された大量の木を人間に見せるために要約

# マルコフ連鎖モンテカルロ法 (MCMC)



- 良さそうな木を計算機に大量に生成させる
  - 観測データは固定
  - 推論対象データを一定の条件にしたがってランダムに動かす
- 生成された大量の木を人間に見せるために要約

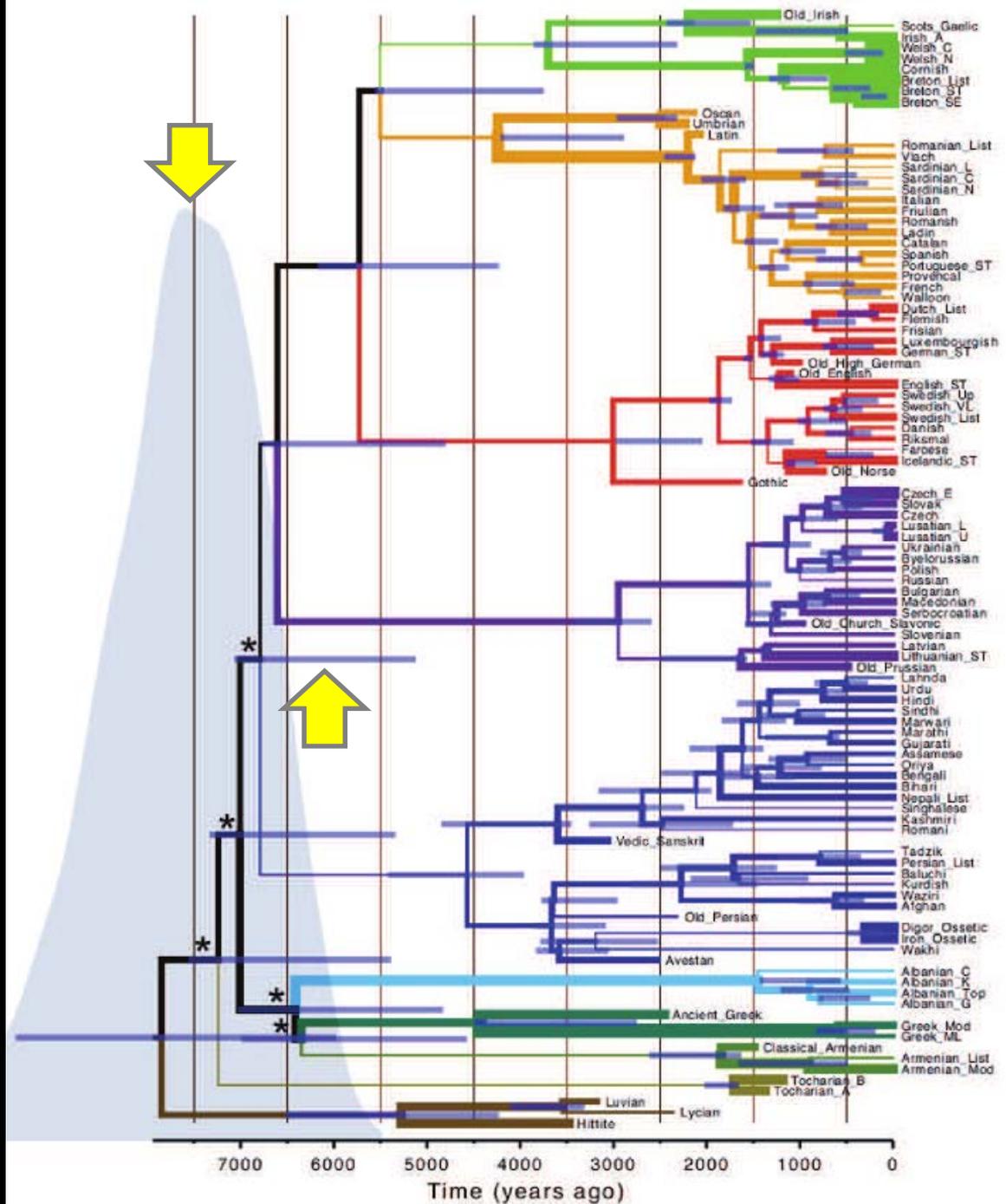
# マルコフ連鎖モンテカルロ法 (MCMC)



- 良さそうな木を計算機に大量に生成させる
  - 観測データは固定
  - 推論対象データを一定の条件にしたがってランダムに動かす
- 生成された大量の木を人間に見せるために要約

# 印欧語族年代つき 系統樹 (再掲)

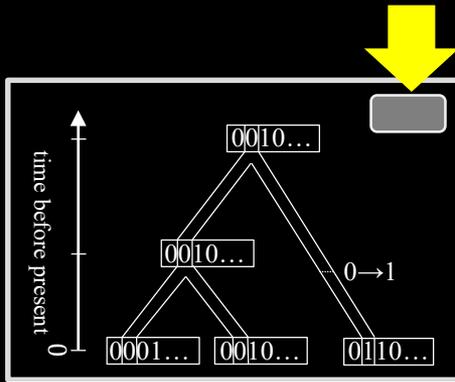
- 木のトポロジーも複数の木を要約したもの (最大系統群信頼度木)



[Bouckaert+, Science 2012]

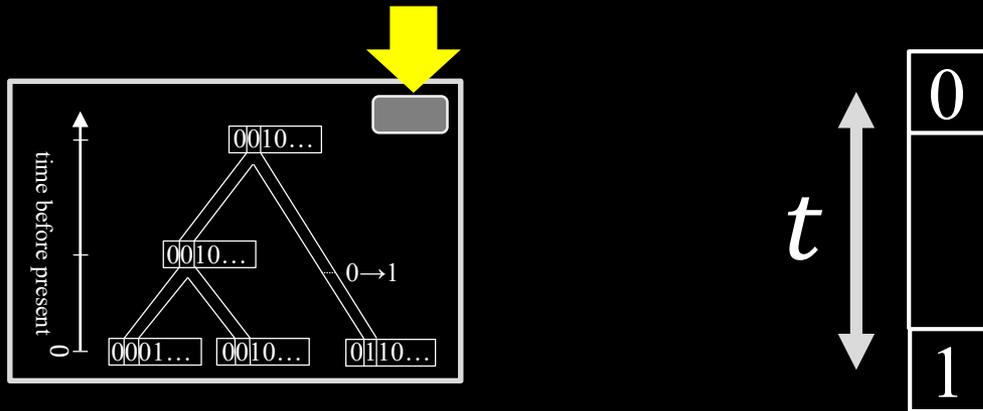
From Bouckaert et al. 2012. Mapping the Origins and Expansion of the Indo-European Language Family. Science 337(6097). Reprinted with permission from AAAS.

# 年代つき系統樹の良さ



- 確率的生成モデルは年代つき系統樹の良さを採点

# 年代つき系統樹の良さ



- 確率的生成モデルは年代つき系統樹の良さを採点

# 年代つき系統樹の良さ



- 確率的生成モデルは年代つき系統樹の良さを採点

# 年代つき系統樹の良さ



- 確率的生成モデルは年代つき系統樹の良さを採点

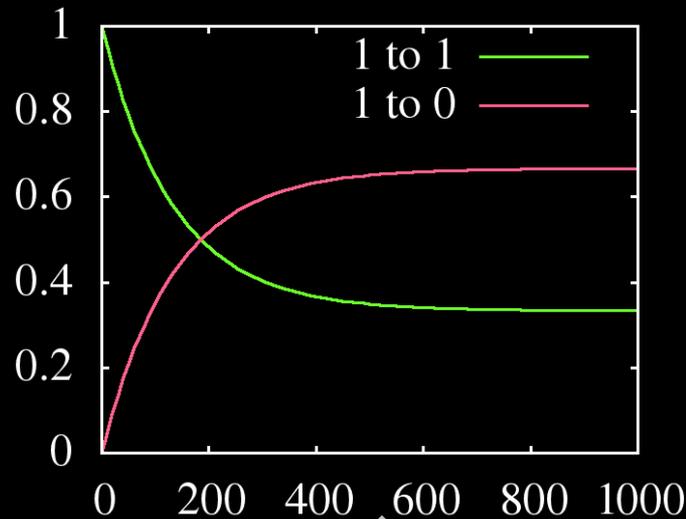
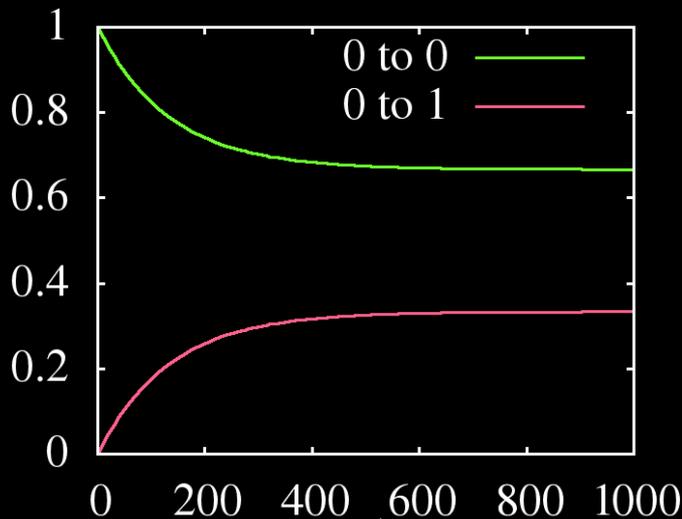
# 年代つき系統樹の良さ



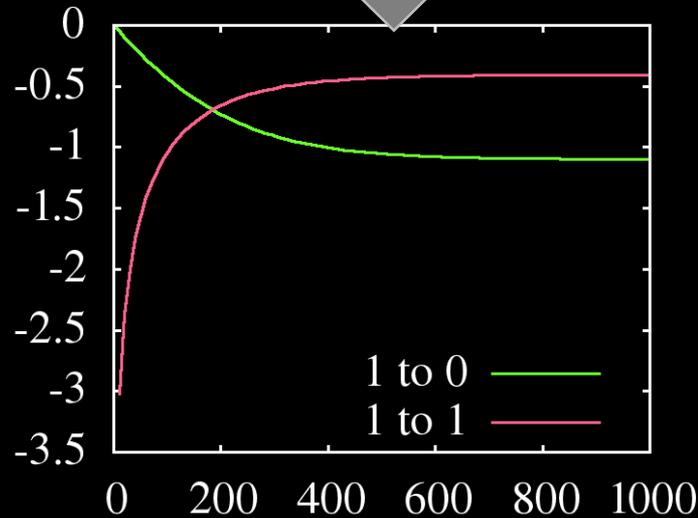
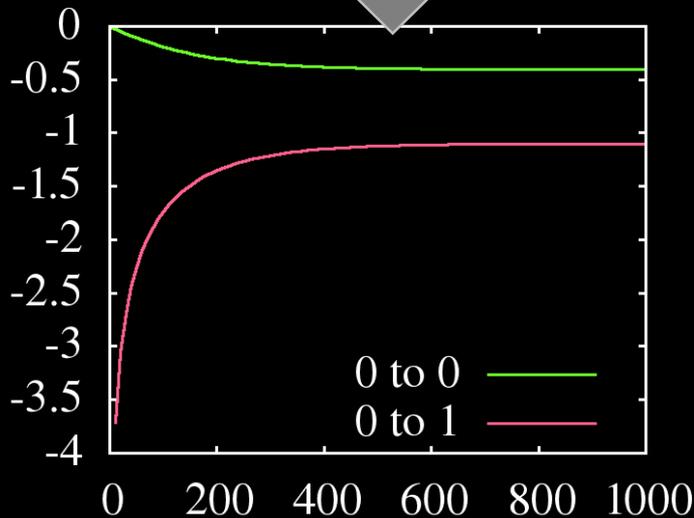
- 確率的生成モデルは年代つき系統樹の良さを採点
- モデルはデータの振る舞いに対する仮定  
– 年代つき系統樹モデルは数多くの仮定の産物

# 置換モデル (連続時間マルコフモデル)

確率

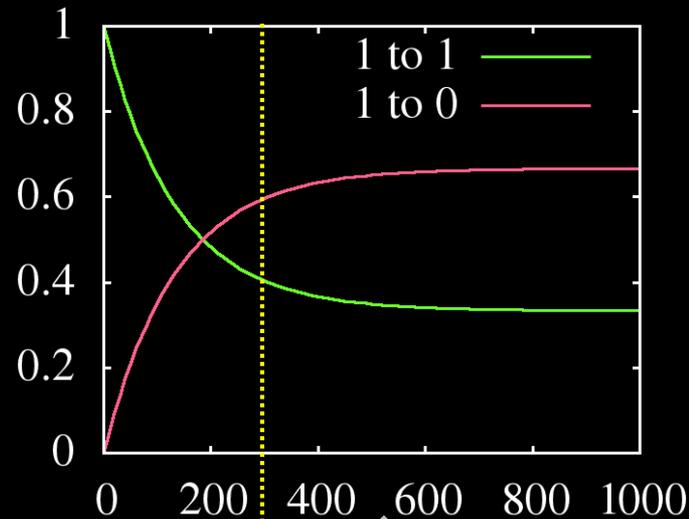
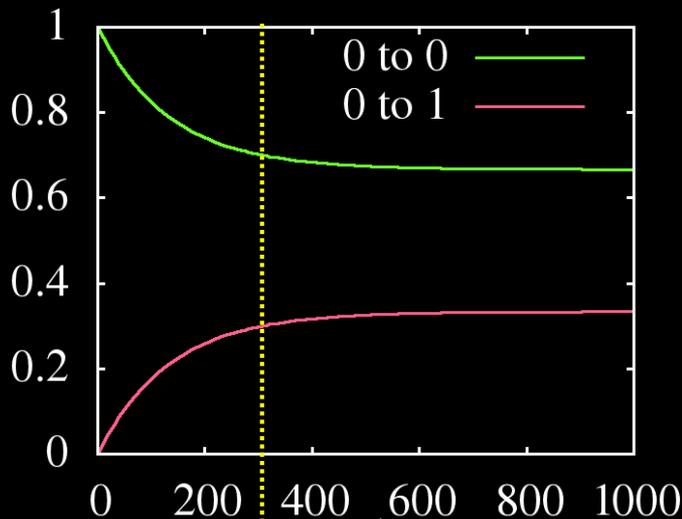


点数 (対数確率)

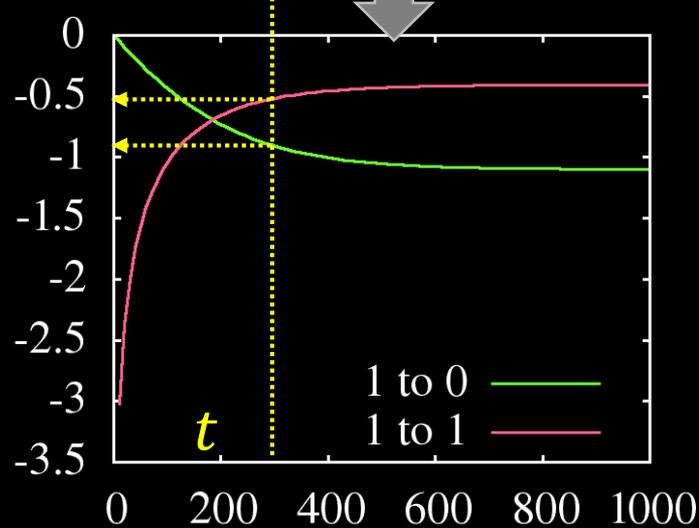
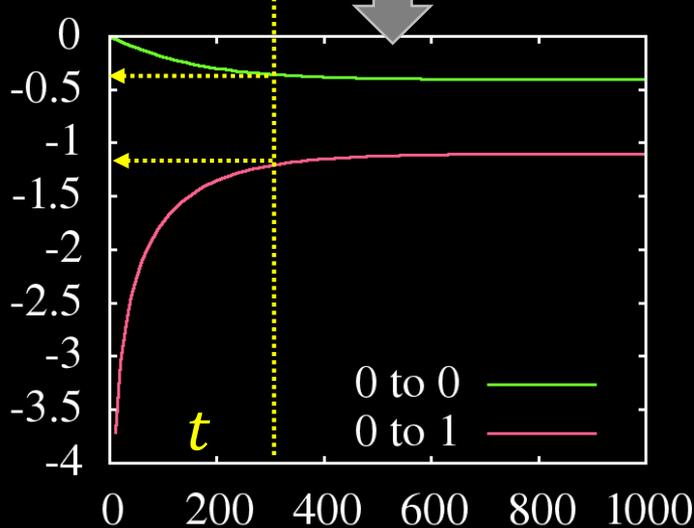


# 置換モデル (連続時間マルコフモデル)

確率

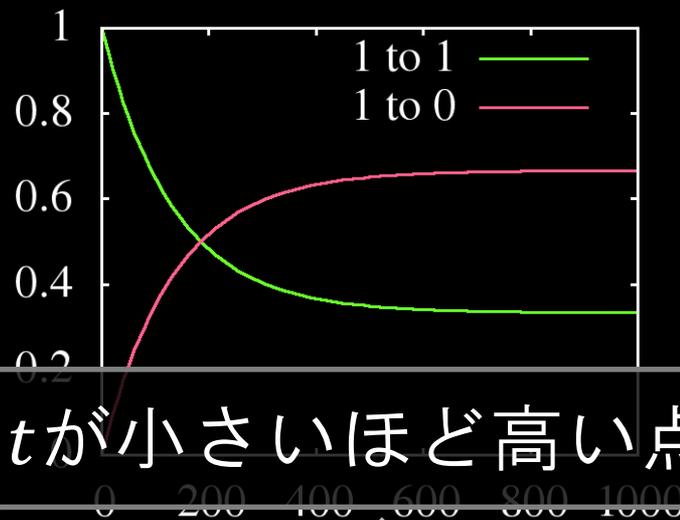
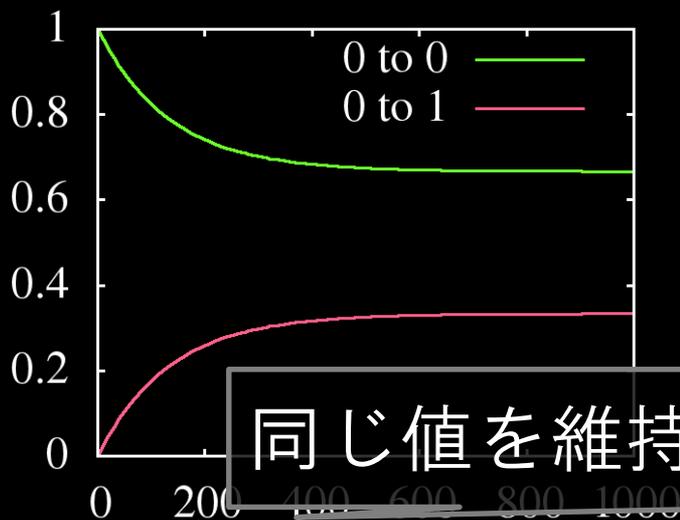


点数 (対数確率)



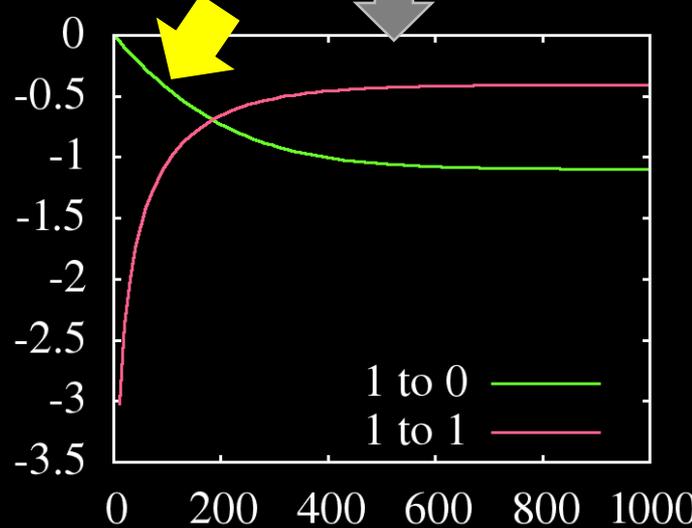
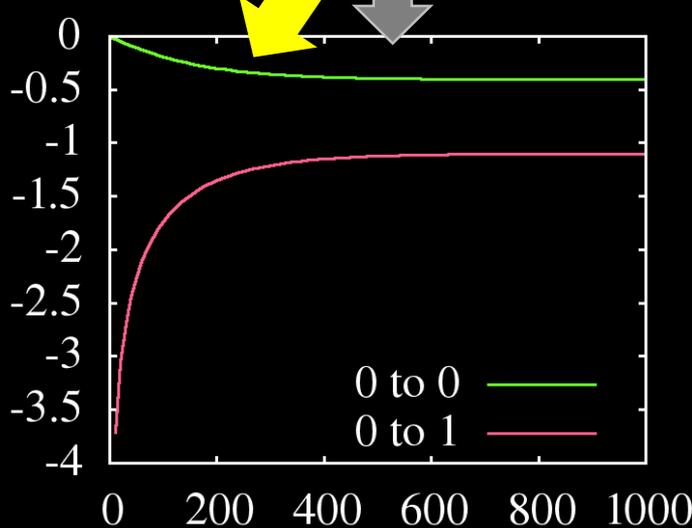
# 置換モデル (連続時間マルコフモデル)

確率



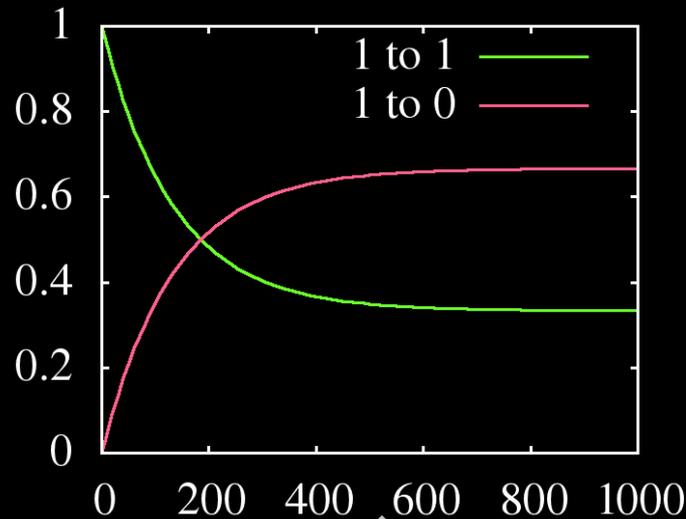
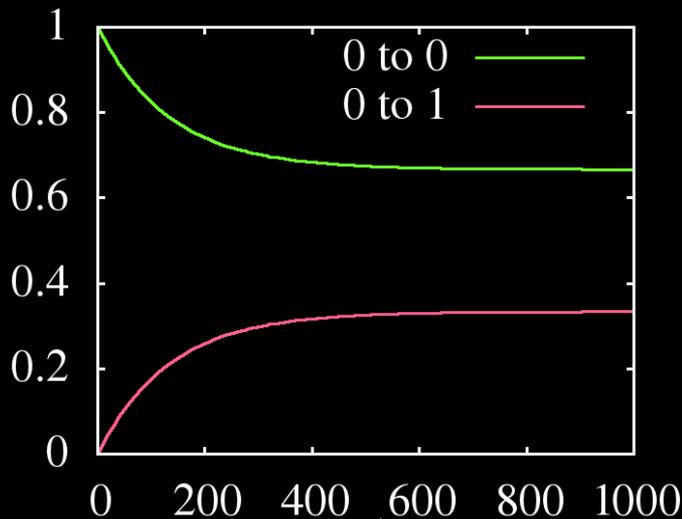
同じ値を維持  $\Rightarrow$   $t$  が小さいほど高い点数

点数 (対数確率)

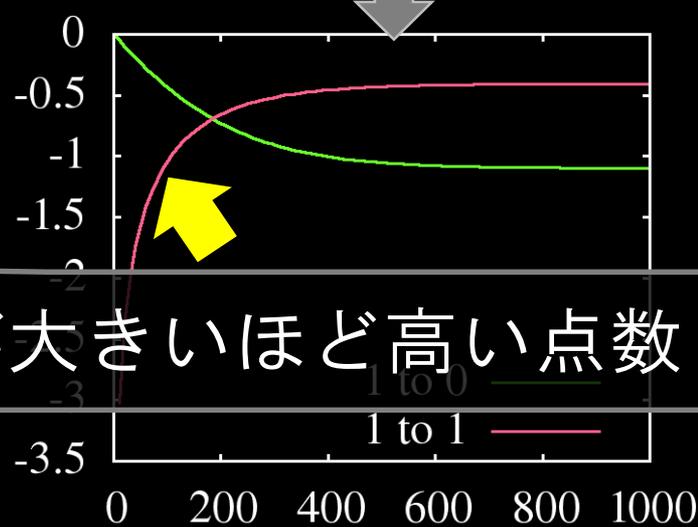
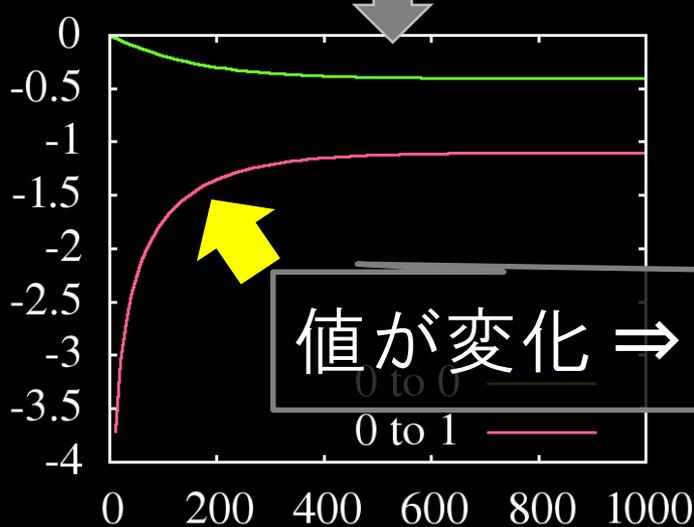


# 置換モデル (連続時間マルコフモデル)

確率

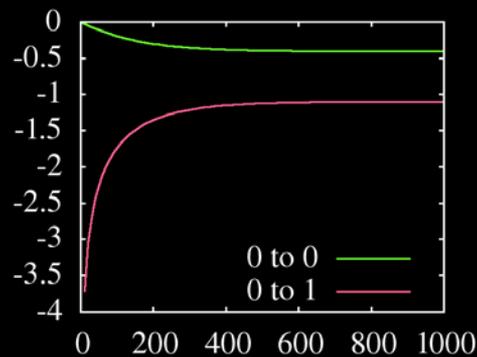


点数 (対数確率)

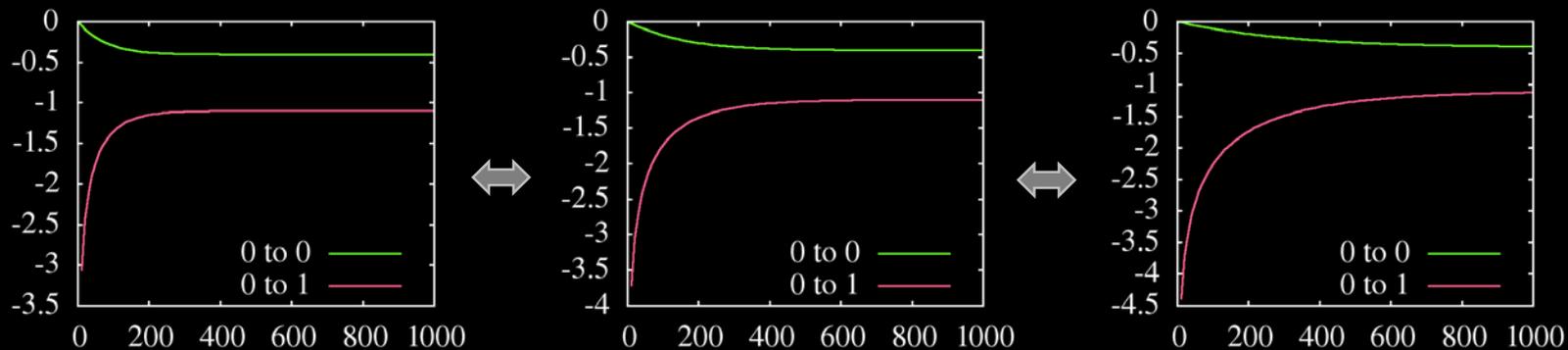


値が変化  $\Rightarrow$   $t$  が大きいほど高い点数

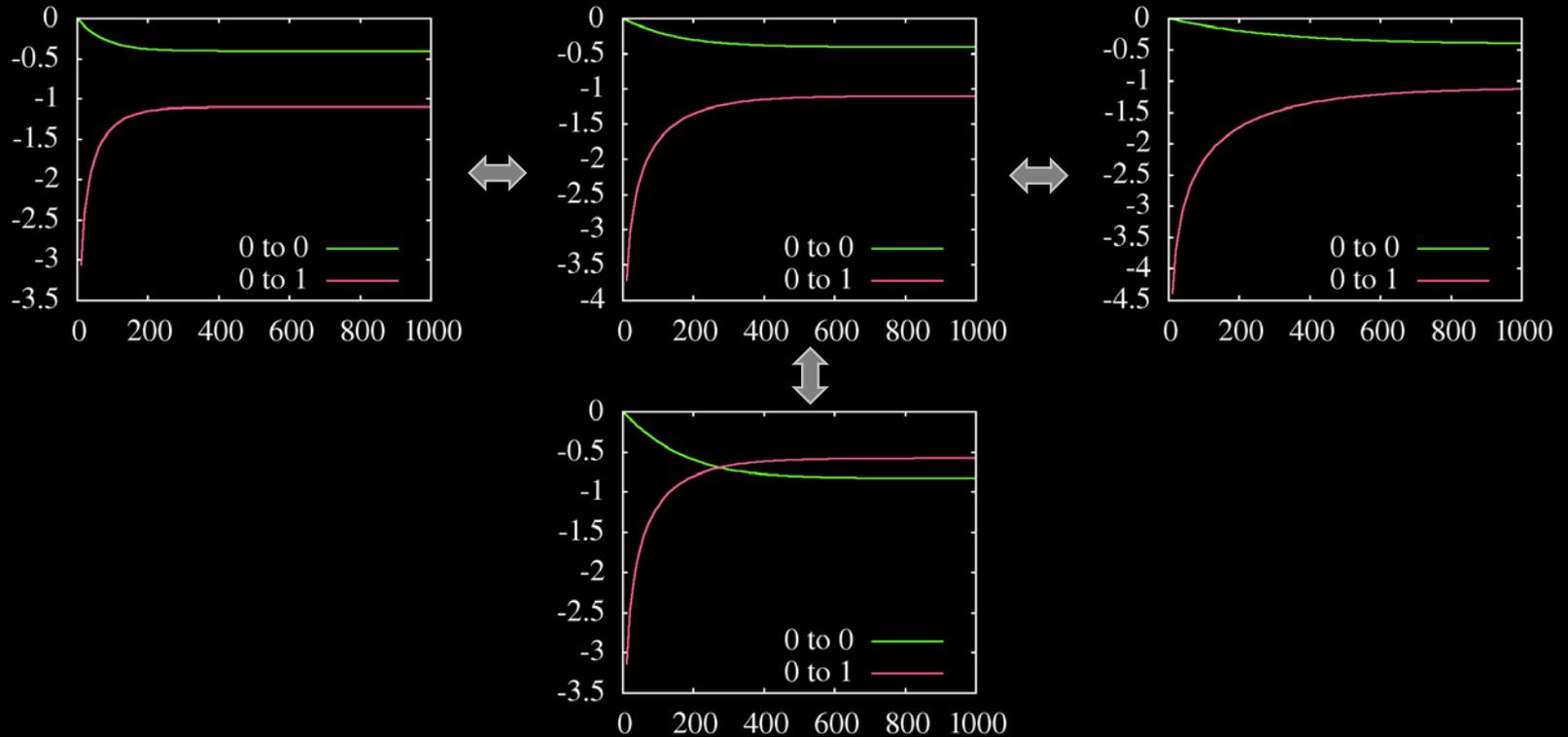
# 置換モデル自体も採点対象



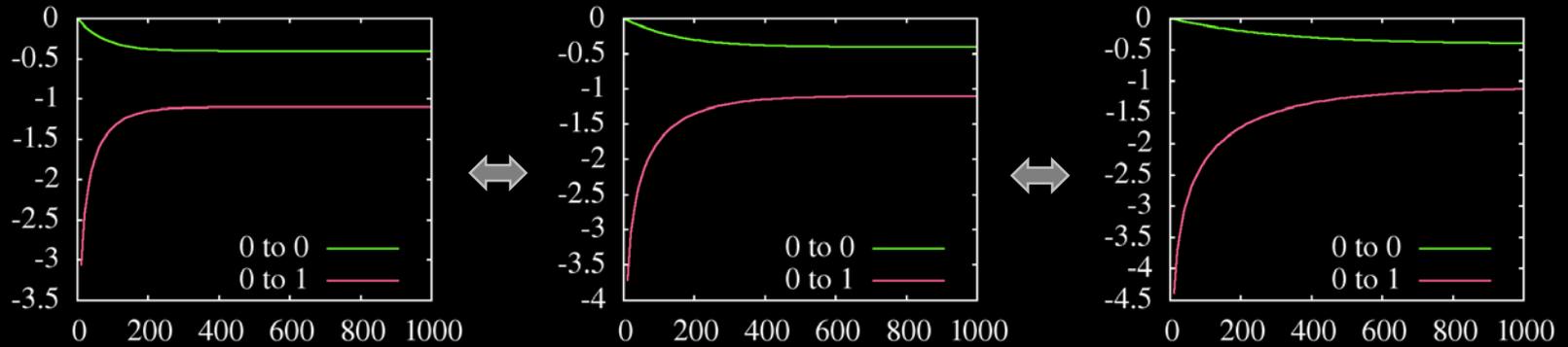
# 置換モデル自体も採点対象



# 置換モデル自体も採点対象

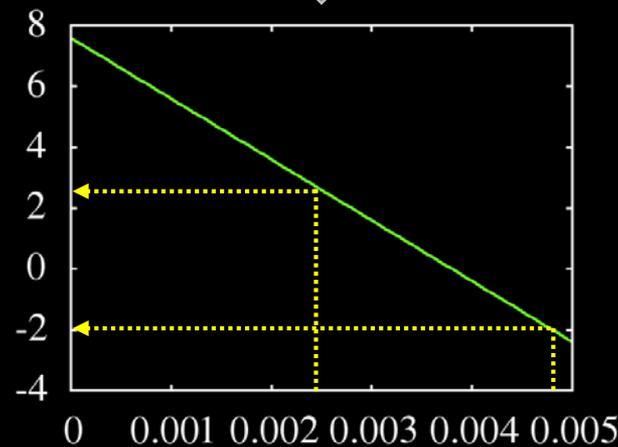
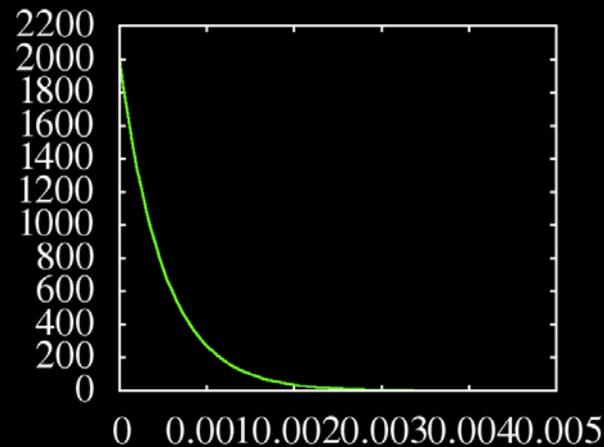
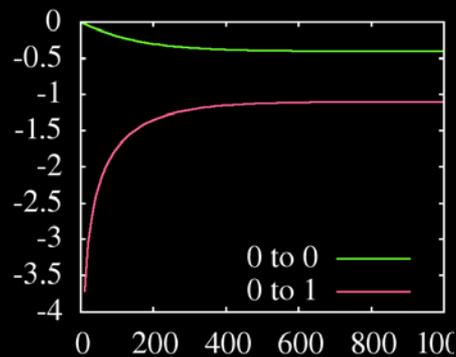
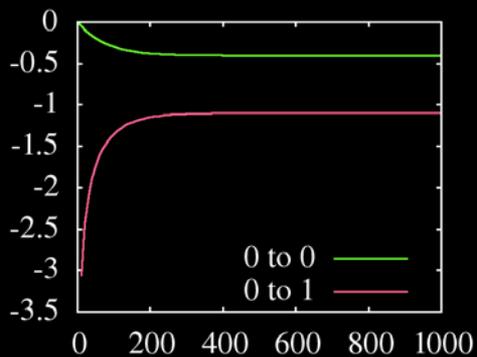


# 置換モデル自体も採点対象



$$\begin{pmatrix} * & \alpha \\ \beta & * \end{pmatrix}$$

# 置換モデル自体も採点対象



$$\begin{pmatrix} * & \alpha \\ \beta & * \end{pmatrix}$$

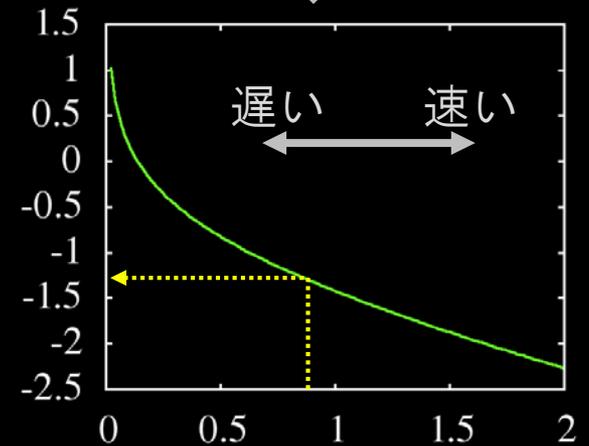
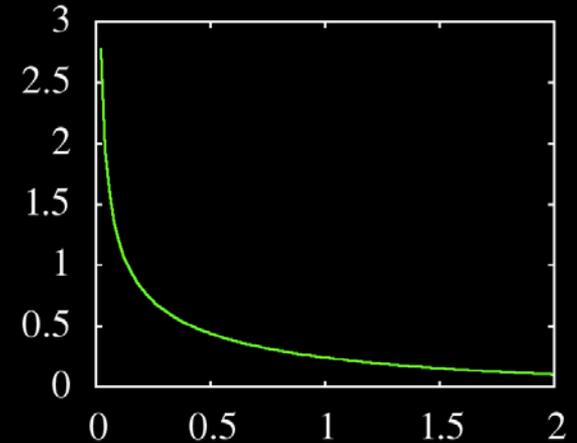
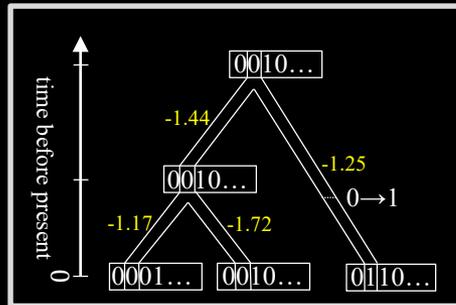
# モデルやデータに対する 様々な疑問

- 変化の速さは共通?
  - 言語によって異なるのでは?
  - 語によって異なるのでは?
- 系統樹モデルは接触の影響を無視するけど、本当は無視できないのでは?
- 借用語は事前に取り除いているというけど、本当に充分?
- 単純な置換モデルでは、誕生→死亡→誕生という変化があり得るけど、一度消えた復活するのは不自然では?
- 何か体系的なバイアスがデータに潜んでいない?

# Q: 変化の速さはすべての言語に共通?

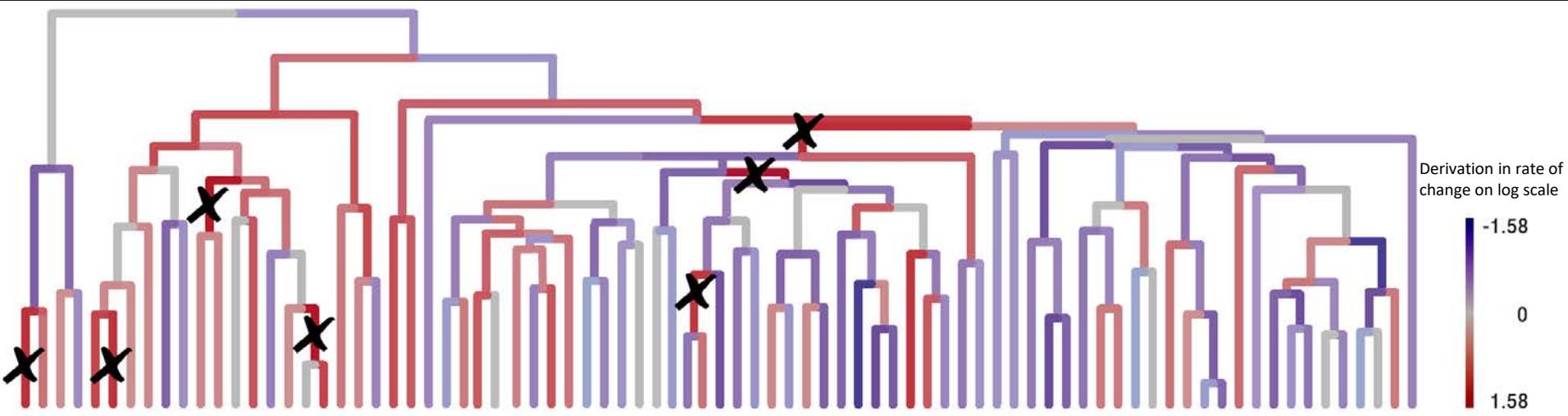
- アイスランド語は超保守的 [Bergsland+, 1962]
- 系統樹の枝 $k$ ごとに変異率 $r_k$ を導入 (緩和時計モデル)

$$r_k \begin{pmatrix} * & \alpha \\ \beta & * \end{pmatrix}$$



# Q: 変化の速さはすべての言語に共通?

- オーストロネシア語族における変化の緩急 [Greenhill+, PNAS 2017]



- 古い時代は変化の速さを決める手がかりにとぼしいので、緩和時計モデルを用いると、共通祖語の年代の分布はすそが広がるはず

# Q: 意味変化による体系的バイアス?

[Chang+, Language 2015]

- 成因的相同 (homoplasy): 同じ特徴が独立に発生
- 実は無視できないほど頻出
  - IELEXのロマンス諸語の基礎語彙の8.1%

\**d<sup>h</sup>ǵ<sup>h</sup>om-*, ADULT MALE

現代アイル  
ランド語  
*duine*

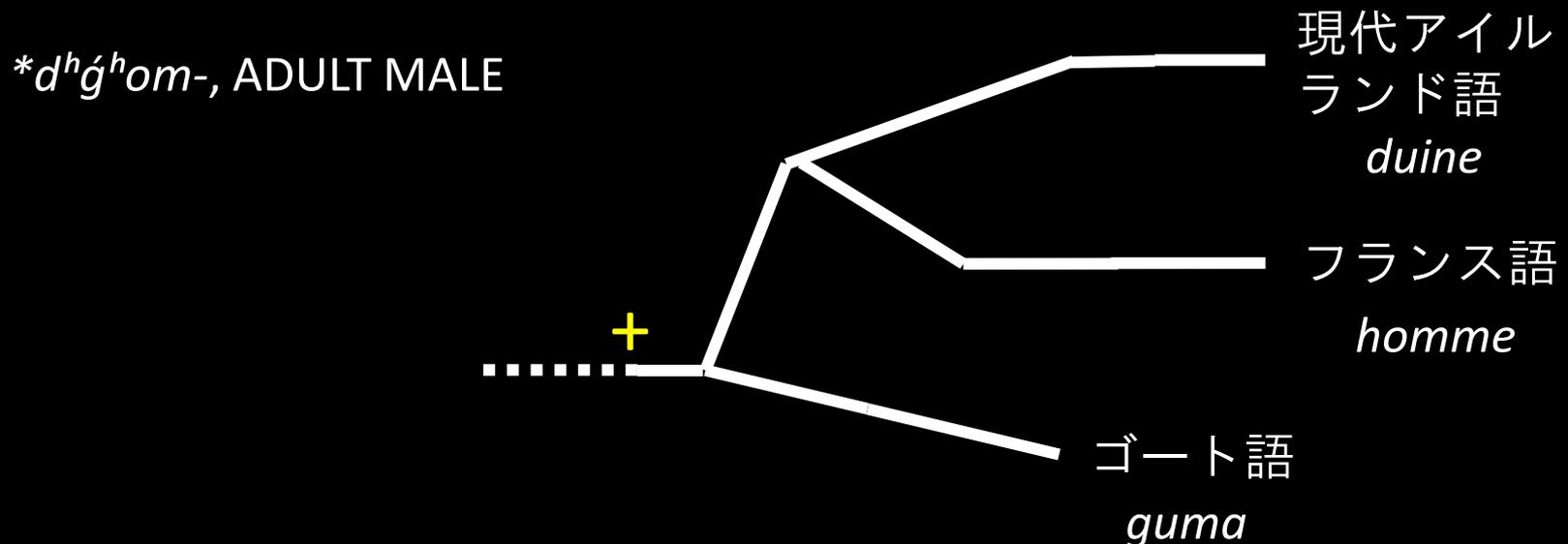
フランス語  
*homme*

ゴート語  
*guma*

# Q: 意味変化による体系的バイアス?

[Chang+, Language 2015]

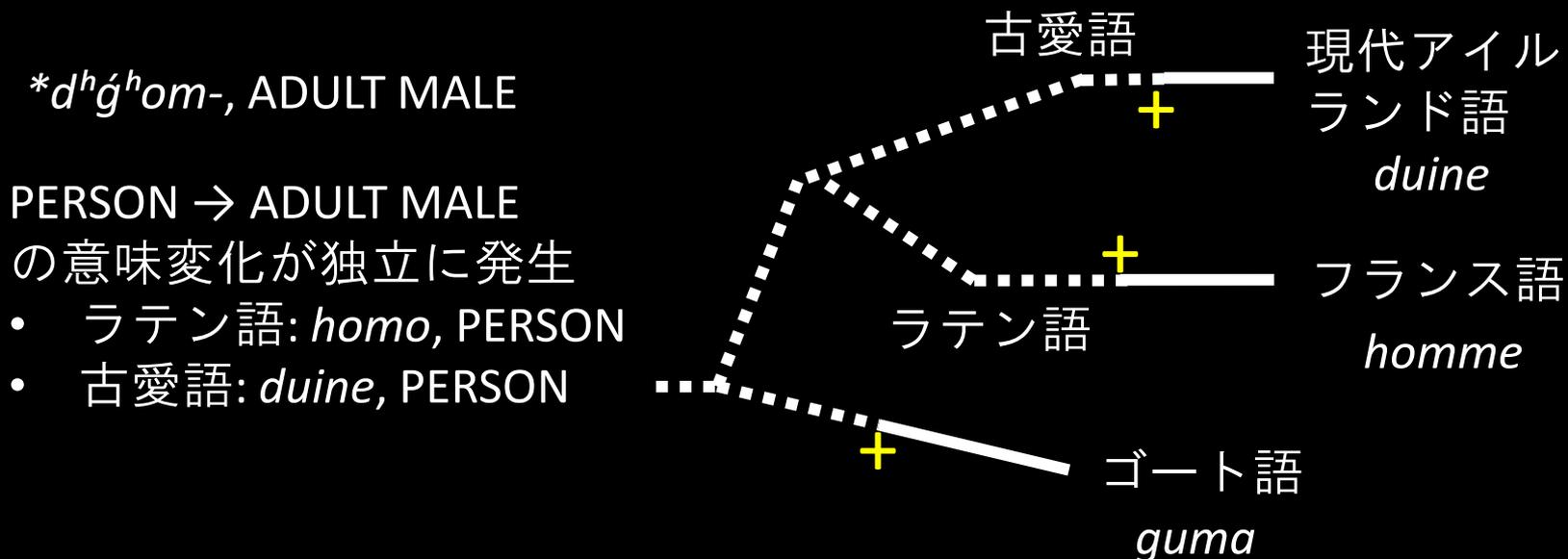
- 成因的相同 (homoplasy): 同じ特徴が独立に発生
- 実は無視できないほど頻出
  - IELEXのロマンス諸語の基礎語彙の8.1%



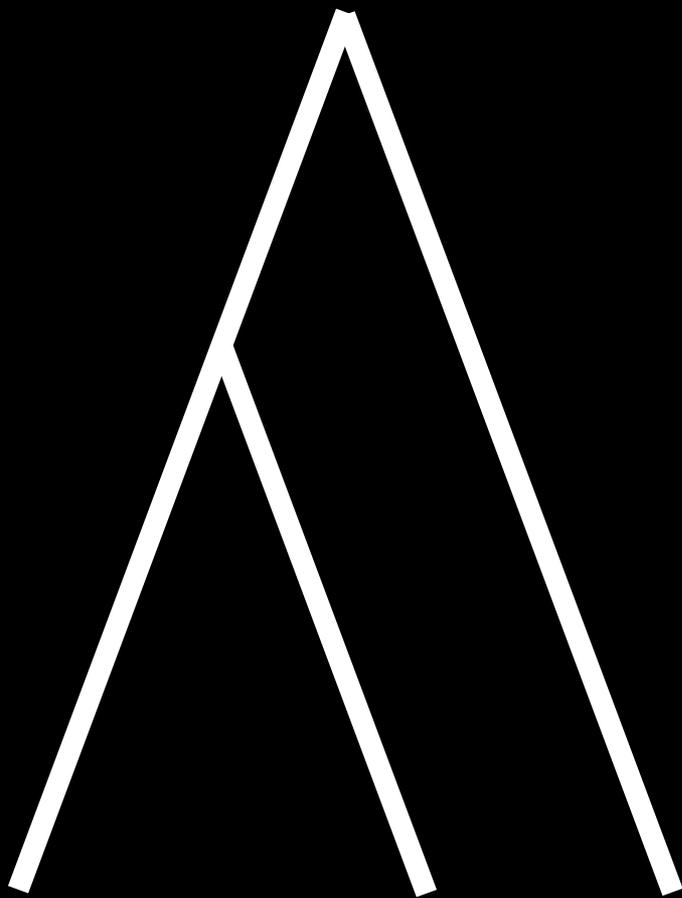
# Q: 意味変化による体系的バイアス?

[Chang+, Language 2015]

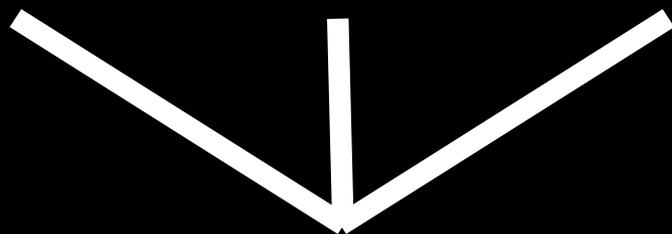
- 提案手法: 古代語を制約として使う
- 結果: 印欧祖語の年代は約6,500年前に繰り上がり、クルガン仮説寄りに



時間

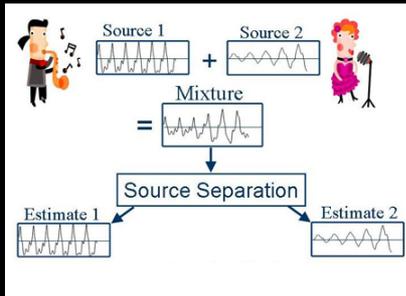


系統樹モデル

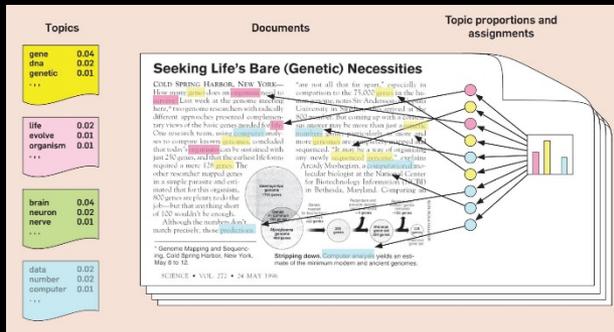
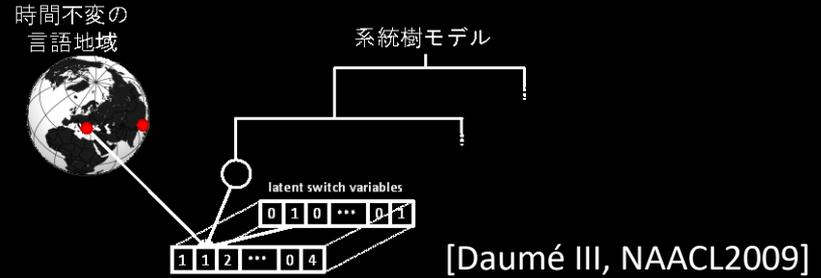


混合モデル

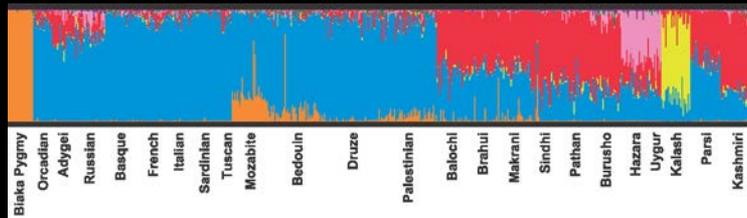
# 混合モデル



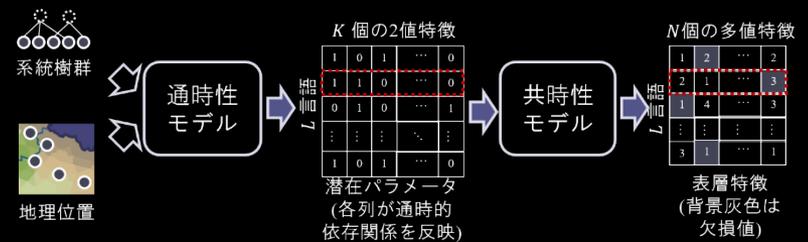
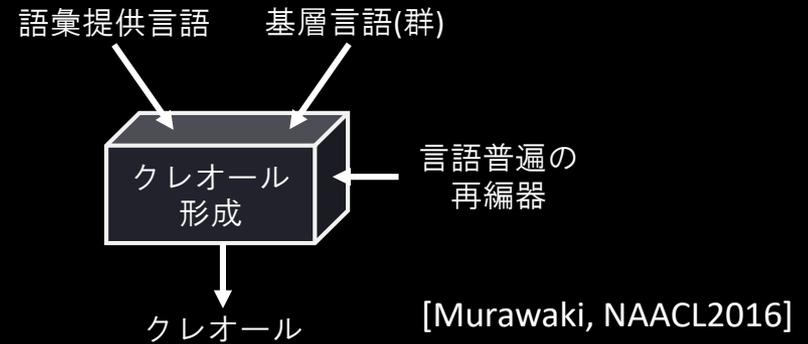
音源分離



文書のトピックモデル [Blei+, 2012]

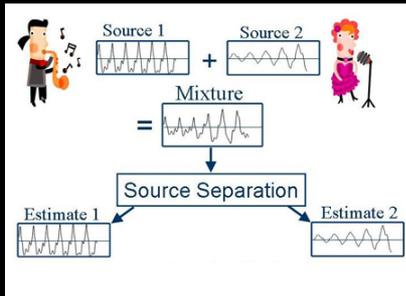


集団遺伝学における混合分析 [Pritchard+, 2000]

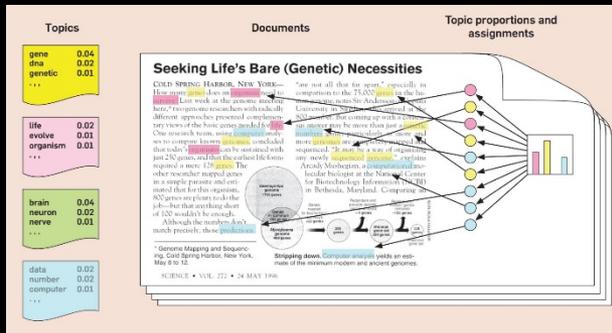
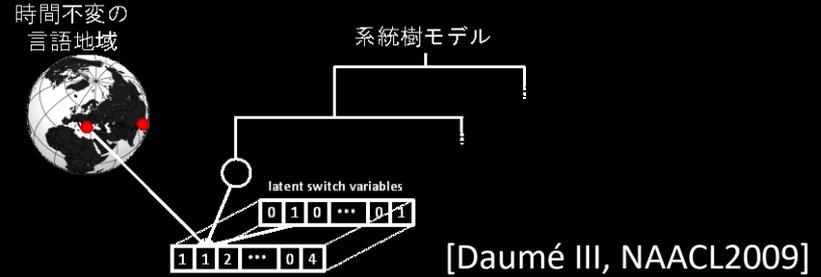


[Yamauchi+, COLING2016] [Murawaki, IJCNLP2017]

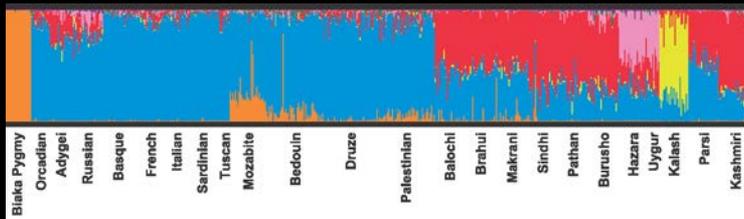
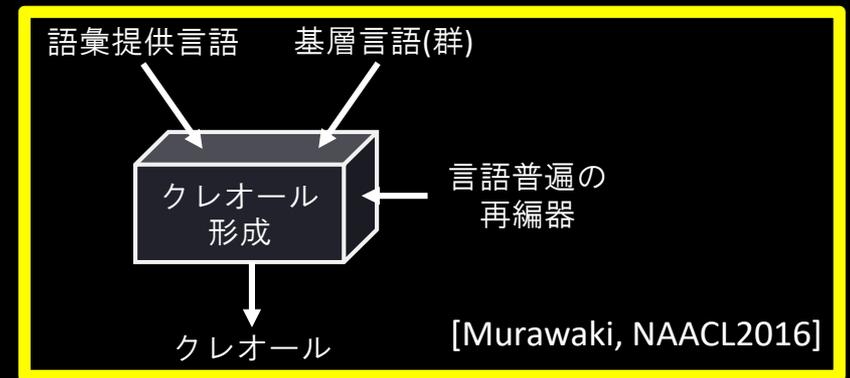
# 混合モデル



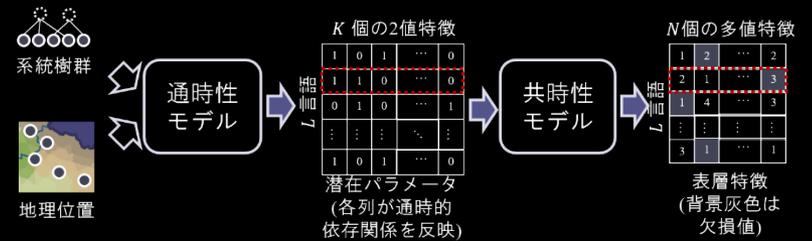
音源分離



文書のトピックモデル [Blei+, 2012]



集団遺伝学における混合分析 [Pritchard+, 2000]





# まとめ

- 人間は苦手だが計算機なら解ける問題は計算機にやらせればよい
  - ベイズ統計は不確実性を自然に扱う枠組み
- モデルとはデータに対する仮定
  - 仮定次第で結果が大きく変わる (場合がある)
  - すべてが明示されていて検証可能
- 個人的な課題は系統樹モデルの克服
- 統一的な基準で作成された言語データを機械可読な形でどんどん公開してほしい

