# A very *practical* MBTA subway map
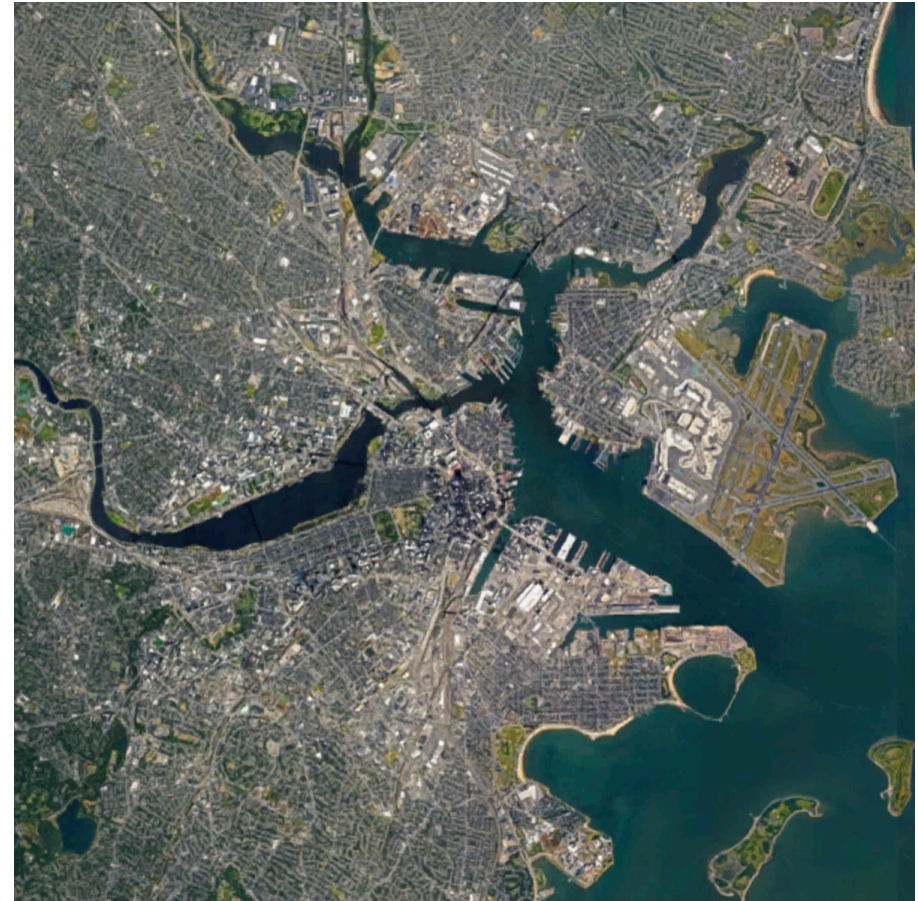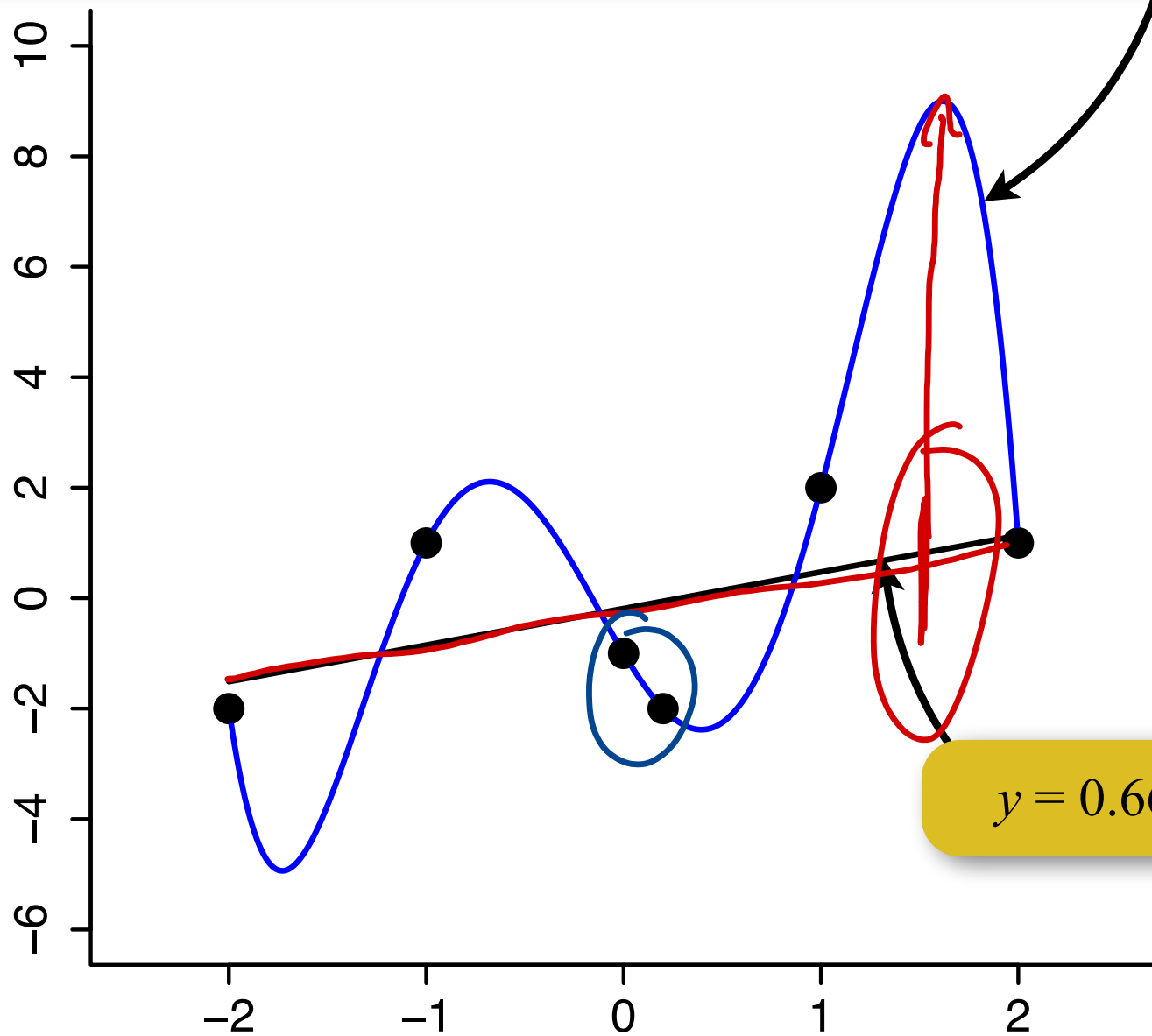
# A very *realistic* MBTA subway map

# Which is more useful?

$$y = -1.5972\ x^5 + -0.7917\ x^4 + 8.0694\ x^3 + 3.2917\ x^2 + -5.9722\ x + -1.0$$

linear
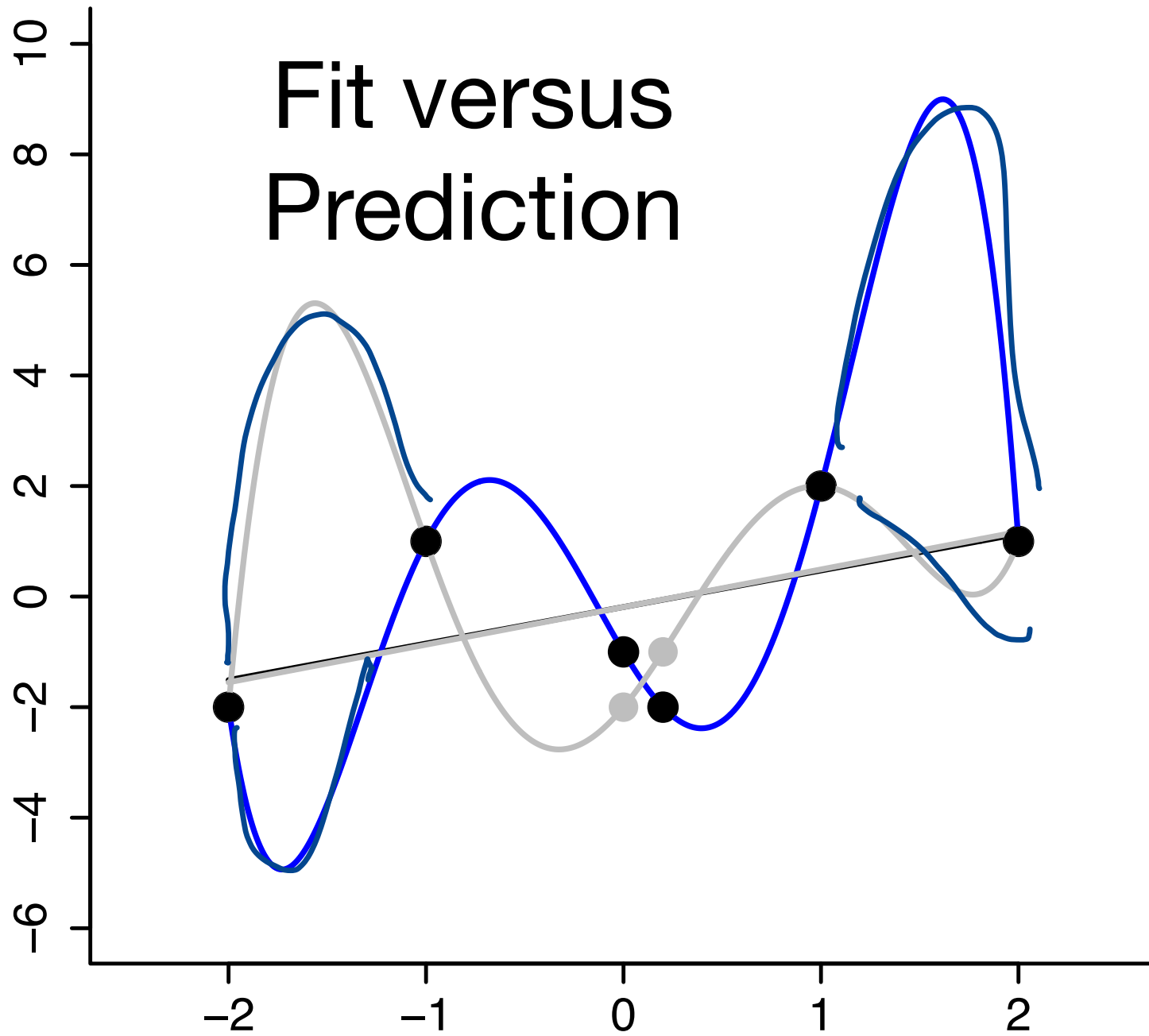
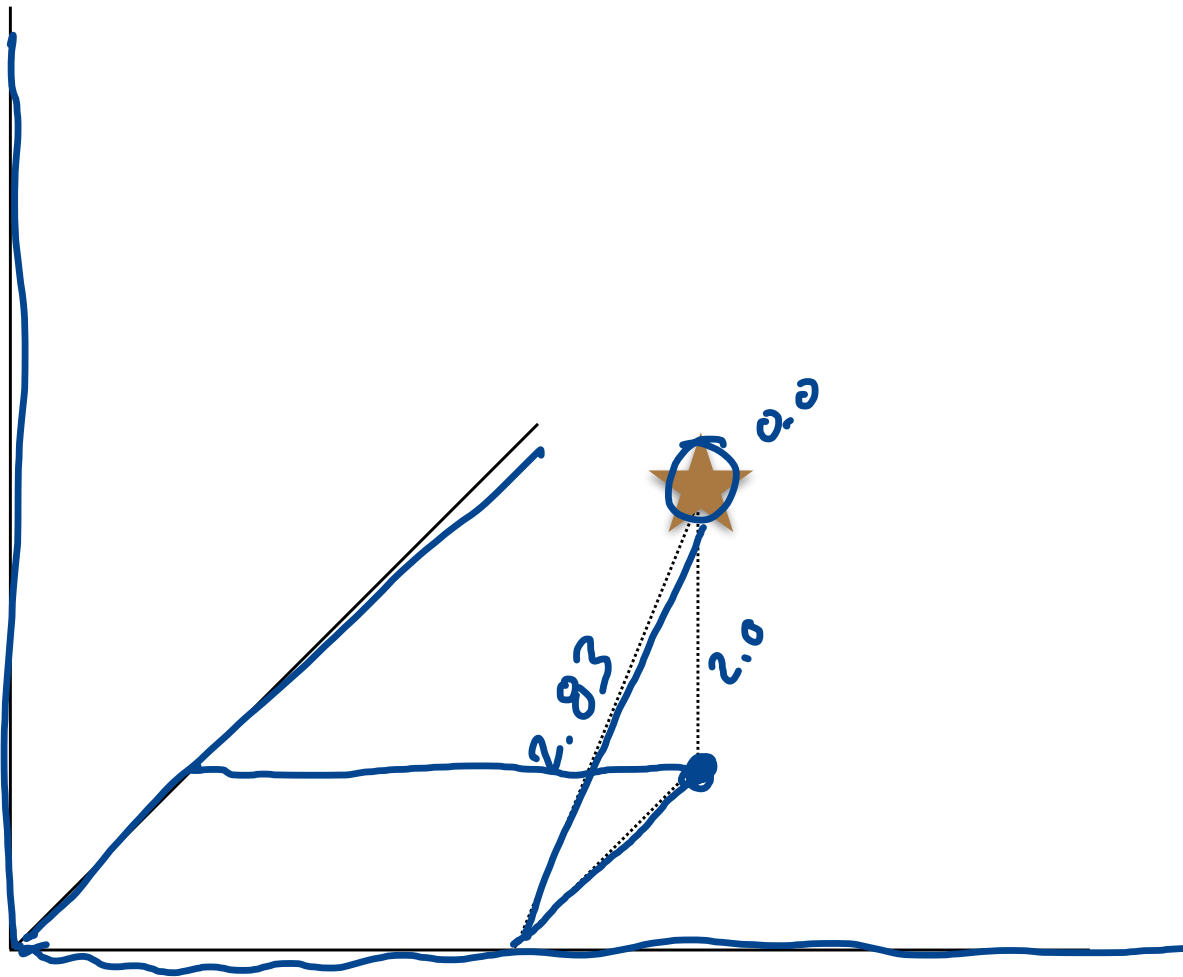$$y = 0.6611x + -0.1887$$

goodness-of-fit

prediction

Fit versus Prediction

# Model dimensions



0.0

2.83

2.0

1-parameter model: 2.83

2-parameter model: 2.00

3-parameter model: 0.00

# Model dimensions

gratuitous complexity
colinearity

# JC69 model

$$\pi_A = \tfrac{1}{4}$$
$$\pi_G = \tfrac{1}{4}$$
$$\pi_C = \tfrac{1}{4}$$
$$\pi_T = \tfrac{1}{4}$$

equilibrium relative frequencies

total rate $= 3\beta$

C $\xleftarrow{\beta}$ A $\xrightarrow{\beta}$ G

A $\xrightarrow{\beta}$ T

A $\xuparrow{\beta}$

C $\xrightarrow{\beta}$ G

C $\xrightarrow{\beta}$ T

Jukes & Cantor (1969)

number of substitutions

substitution rate

total substitution rate

# Edge lengths

$$\text{number} = (\text{rate})(\text{time})$$

$$100 \text{ miles} = \left(50 \frac{\text{miles}}{\text{hour}}\right)(2 \text{ hours})$$

$$\begin{array}{c}\text{expected} \\ \text{no. subst.}\end{array} = \left(\begin{array}{c}\text{rate of} \\ \text{subst.}\end{array}\right)(\text{time})$$

⇑

sequence
data

$$v = (3\beta)t \qquad \leftarrow JC69$$
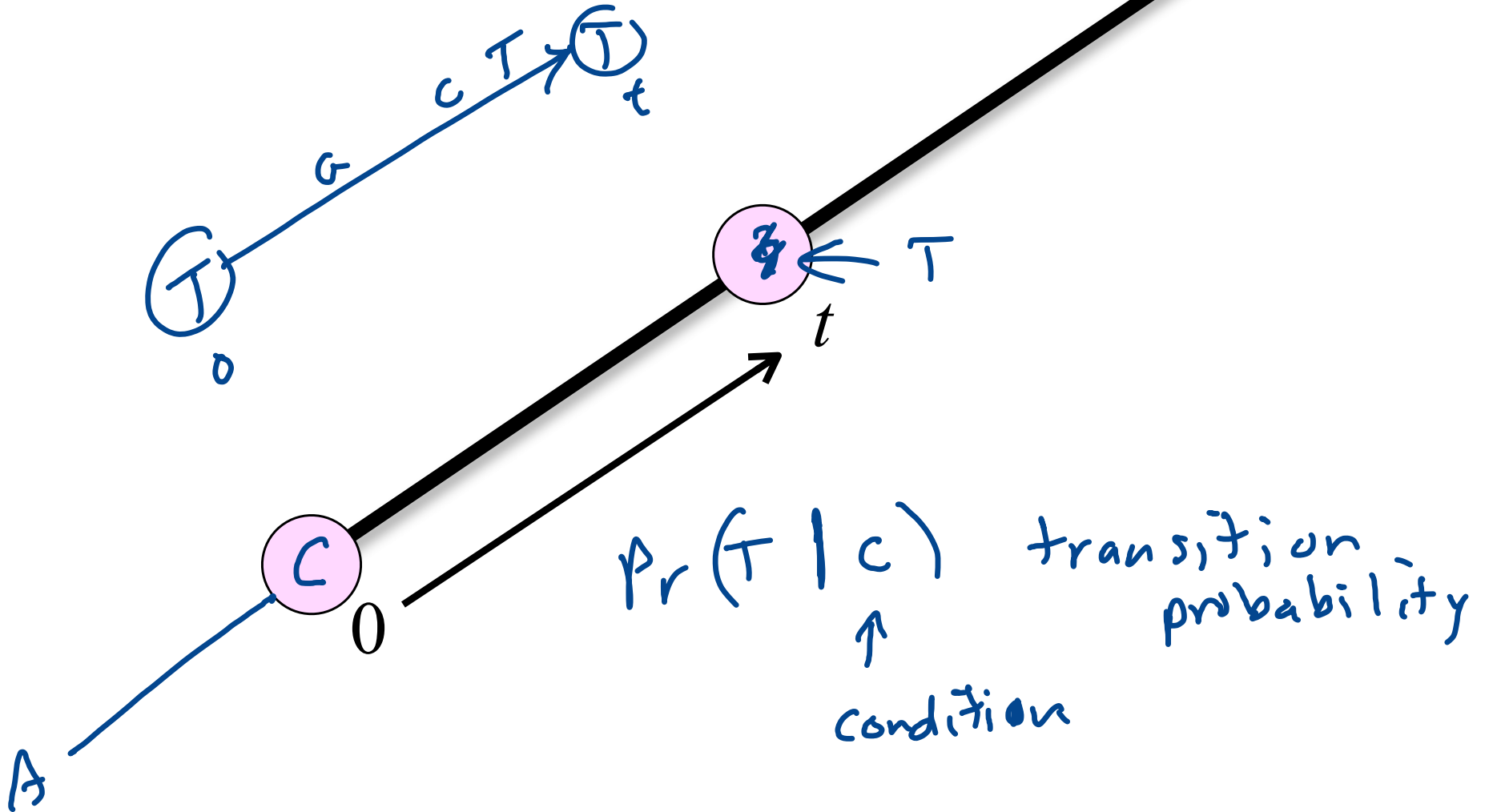
edge length parameters
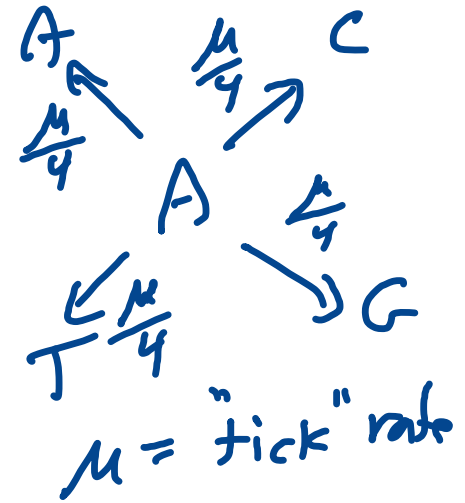
long edge lengths means...

# Markov Models



$$Pr(T \mid c)$$

transition probability

↑
condition

Markov property
transition probability ←
conditional probability ←
multiple hits ←

# JC69 Transition Probability

$$\beta = \frac{\mu}{4}$$

$$P(G|T) = \boxed{P_{TG}(t)} = \left(\frac{1}{4}\right)\left(1 - e^{-4\beta t}\right)$$

$\rightarrow$ prob. at least one tick

prob. last tick resulted in $G$ $\rightarrow \boxed{\frac{1}{4}}$

$\mu$ = "tick" rate

Poisson $p(k) = \dfrac{\lambda^k e^{-\lambda}}{k!}$

$\lambda = \mu t$

$0! = 1$

$5! = 5 \cdot 4 \cdot 3 \cdot 2 \cdot 1$

$= \dfrac{(\mu t)^k e^{-\mu t}}{k!}$

$p(0) = \dfrac{(\mu t)^0 e^{-\mu t}}{0!}$

e = 2.718281828459045...

"perturbation" rate vs. substitution rate

$$e^{-\mu t} = \text{prob. no tick marks from 0 to } t$$

$$1 - e^{-\mu t} = \text{prob. at least one tick}$$

$$\frac{\mu}{4} = \beta \quad \mu = 4\beta$$

$$1 - e^{-4\beta t}$$

$$P(G|T) = \frac{1}{4}\left(1 - e^{-4\beta t}\right)$$

# JC69 Transition Probability

$$P_{TA}(t) = \frac{1}{4} - \frac{1}{4}e^{-4\beta t}$$

$$P_{TC}(t) = \frac{1}{4} - \frac{1}{4}e^{-4\beta t}$$

$$P_{TG}(t) = \frac{1}{4} - \frac{1}{4}e^{-4\beta t}$$

$$P_{TT}(t) = \frac{1}{4} - \frac{1}{4}e^{-4\beta t}$$

$$\frac{1}{4} + \frac{3}{4}e^{-4\beta t}$$

$$+ \frac{4}{4}e^{-4\beta t}$$

$$1 - e^{-4\beta t} + e^{-4\beta t}$$

# JC69 model assumptions

equal frequencies

$$\pi_A = \pi_C = \pi_G = \pi_T = \frac{1}{4}$$

equal rates

C

G

$\beta$

$\boxed{\beta}$

A

$\beta$

$\beta$

T

$\nu_1$ $\beta_1$

$\beta_4$

$\nu_4$

$\beta_3$

$\nu_3$

$\nu_5$ $\beta_5$

$\beta_2$

$\nu_2$

1 parameter for each
edge length in tree

Jukes & Cantor
(1969)

# Equilibrium Frequencies



(architect: Joe Bielawski)

# Equilibrium Frequencies

# Equilibrium Frequencies



STOPPED HERE 2024-01-30

$$\pi_A = \pi_C = \pi_G = \pi_T = \frac{1}{4}$$

# JC69 Distance Formula

$$p = \frac{3}{4}\left(1 - e^{-4\beta t}\right)$$

$$\log e^{a} = a$$

$$e^{\log a} = a$$

$$\frac{4}{3} p = \frac{4}{3} \cdot \frac{3}{4}\left(1 - e^{-4\beta t}\right)$$

$$0 = \frac{4}{3}p - \frac{4}{3}p = 1 - e^{-4\beta t} - \frac{4}{3}p$$

$$e^{-4\beta t} = 1 - e^{-4\beta t} + e^{-4\beta t} - \frac{4}{3}p$$

$$\log e^{-4\beta t} = \log\left(1 - \frac{4}{3}p\right)$$

$$\frac{-4\beta t}{-4} = \frac{\log\left(1 - \frac{4}{3}p\right)}{-4}$$

$$v = -\frac{3}{4}\log\left(1 - \frac{4}{3}p\right)$$

$$.2326$$

$$\uparrow$$

$$.2$$

$$\boxed{v} = 3\beta t = -\frac{3}{4}\log\left(1 - \frac{4}{3}p\right)$$

# JC69 rate matrix

1-parameter model
$\beta$

"To" state

| | | A | C | G | T |
|---|---|---|---|---|---|
| "From" state | A | $-3\beta$ | $\beta$ | $\beta$ | $\beta$ |
| | C | $\beta$ | $-3\beta$ | $\beta$ | $\beta$ |
| | G | $\beta$ | $\beta$ | $-3\beta$ | $\beta$ |
| | T | $\beta$ | $\beta$ | $\beta$ | $-3\beta$ |

# K80 (K2P) rate matrix

2 parameters: $\alpha, \beta$

$$
\begin{array}{cccc}
& A & C & G & T \\
A & -2\beta-\alpha & \beta & \boxed{\alpha}\ k\beta & \beta \\
C & \beta & -2\beta-\alpha & \beta & k\beta\ \boxed{\alpha} \\
G & \boxed{\alpha}\ k\beta & \beta & -2\beta-\alpha & \beta \\
T & \beta & \boxed{\alpha}\ k\beta & \beta & -2\beta-\alpha
\end{array}
$$

$\alpha = \beta$
$= JC69$

$k = \alpha/\beta$
$\uparrow$
Kappa

Kimura (1980)

rate of transitions = $\alpha$

rate of transversions = $\beta$

transition/transversion rate ratio $\alpha/\beta$

no. parameters

equivalence to JC69

# "Transition/transversion ratio" vs. "transition/transversion *rate* ratio"

Possible transitions:                    Possible transversions:

$$\frac{E[\text{No. transitions}]}{E[\text{No. transversions}]} = \frac{\phantom{xxxxxxxxxxxxxxx}}{\phantom{xxxxxxxxxxxxxxx}} =$$

# F81 rate matrix

4 parameters
$\mu, \pi_A, \pi_C, \pi_G$

$\pi_T = 1 - \pi_A - \pi_C - \pi_G$

$JC69: \pi_A = \pi_C = \pi_G = \pi_T = \frac{1}{4}$

|   | A | C | G | T |
|---|---|---|---|---|
| A | — | $\pi_C \mu$ | $\pi_G \mu$ | $\pi_T \mu$ |
| C | $\pi_A \mu$ | — | $\pi_G \mu$ | $\pi_T \mu$ |
| G | $\pi_A \mu$ | $\pi_C \mu$ | — | $\pi_T \mu$ |
| T | $\pi_A \mu$ | $\pi_C \mu$ | $\pi_G \mu$ | — |

$\frac{1}{4}\mu = \beta$

no. parameters

equivalence to JC69

# HKY85 rate matrix

5 parameters
$\beta, \alpha, \pi_A, \pi_C, \pi_G$

$$
\begin{array}{cccc}
& A & C & G & T \\
A & \overline{\phantom{---}} & \pi_C \beta & \pi_G \alpha & \pi_T \beta \\
C & \pi_A \beta & \overline{\phantom{---}} & \pi_G \beta & \pi_T \alpha \\
G & \pi_A \alpha & \pi_C \beta & \overline{\phantom{---}} & \pi_T \beta \\
T & \pi_A \beta & \pi_C \alpha & \pi_G \beta & \overline{\phantom{---}}
\end{array}
$$

Hasegawa, Kishino, & Yano (1985)   no. parameters

equivalence to JC69, F81

# F84 vs. HKY85

## F84 model:

$\mu$       rate of process generating *all types of substitutions*

$k\mu$      rate of process generating *only transitions*

Becomes F81 model if $k = 0$

## HKY85 model:

$\beta$       rate of process generating *only transversions*

$\kappa\beta$      rate of process generating *only transitions*

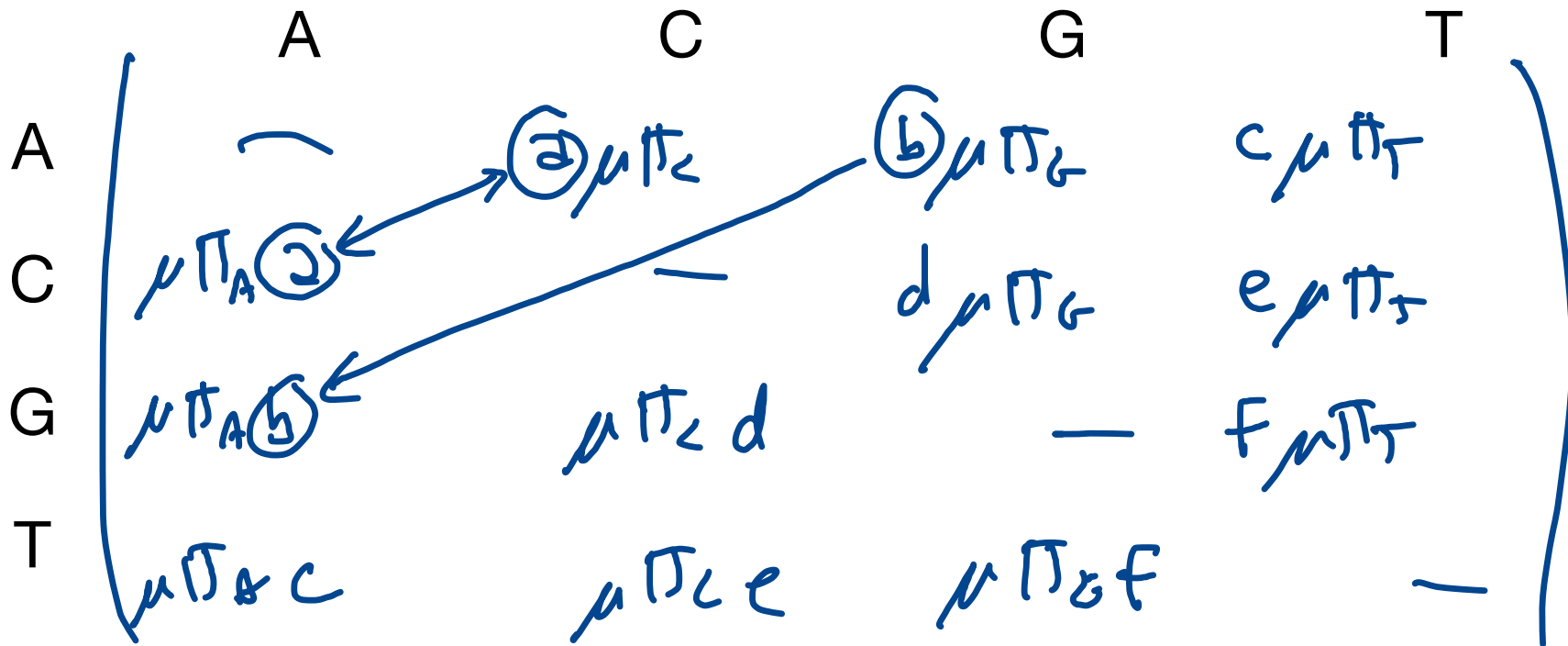Becomes F81 model if $\kappa = 1$

F84 first used in Joe Felsenstein's PHYLIP in 1984

F84 published by Kishino & Hasegawa (1989)

# GTR rate matrix

9 parameters
$\mu, \pi_A, \pi_C, \pi_G$
$a, b, c, d, e$

$$
\begin{array}{cccc}
 & A & C & G & T \\
A & \sim & \textcircled{a}\,\mu\pi_C & \textcircled{b}\,\mu\pi_G & c\,\mu\pi_T \\
C & \mu\pi_A\,\textcircled{a} & \sim & d\,\mu\pi_G & e\,\mu\pi_T \\
G & \mu\pi_A\,\textcircled{b} & \mu\pi_C\,d & \sim & f\,\mu\pi_T \\
T & \mu\pi_A\,c & \mu\pi_C\,e & \mu\pi_G\,f & \sim
\end{array}
$$

Paul O. Lewis ~Phylogenetics, Spring 2024          Lanave et al. (1984)

no. parameters
equivalence to JC69, F81, K80, HKY   25

# Other kinds of models
## (we'll get to some of these later)

- Amino acid substitution models

- Codon models

- Secondary structure models

- Insertion/deletion models

- Relaxed molecular clock models

- Correlated evolution models

- Discrete morphological character models

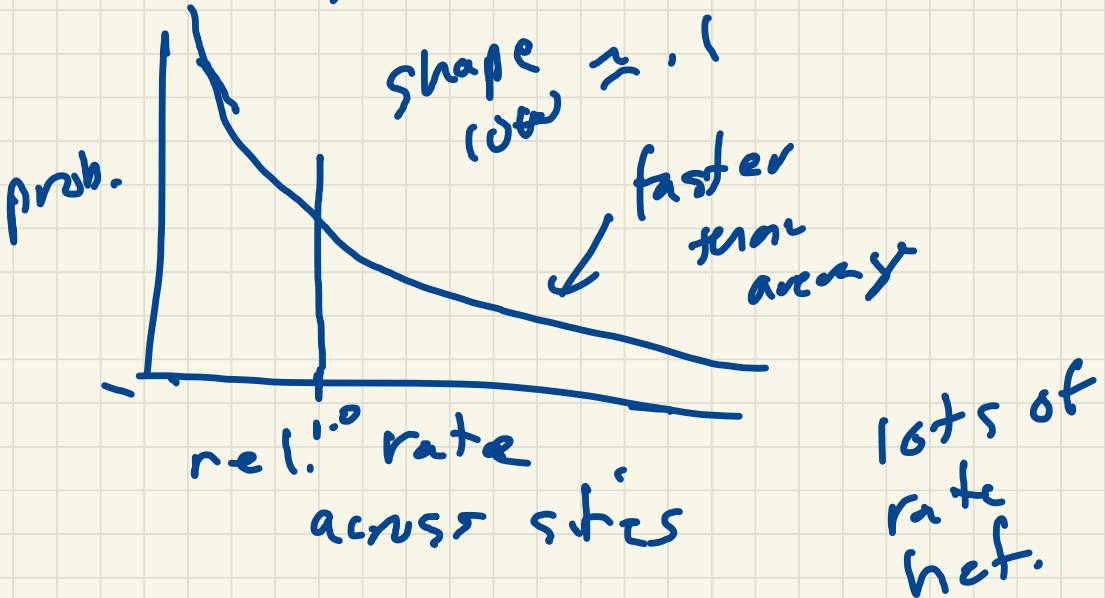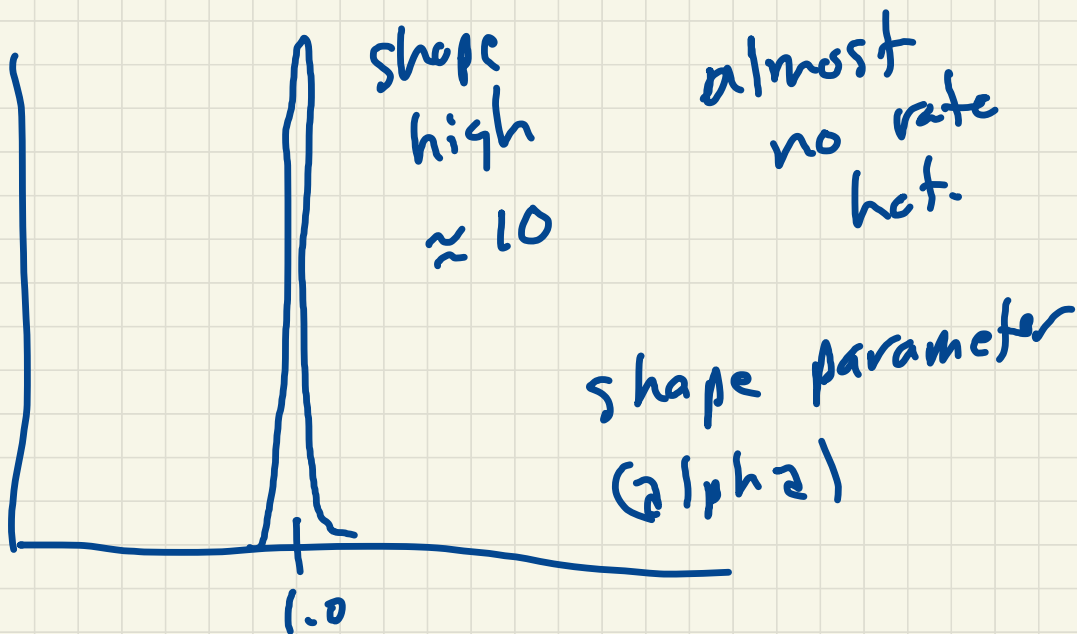- Brownian motion models for continuous traits

# HKY85

I model
G models

I = invariable sites

(Pinvar) = proportion of sites
that are invariable
(rate = 0)

G = gamma model



prob.

shape ≈ .1
(α)

faster
turn
array

rel.⁰ rate
across sites

lots of
rate
het.

shape
high
≈ 10

almost
no rate
het.

shape parameter
(alpha)

1.0

---

in lab

K 80 rate matrix

trs:trv
rate $= K = \dfrac{\alpha}{\beta}$

$$\begin{pmatrix} - & \beta & K\beta & \beta \\ \beta & - & \beta & K\beta \\ K\beta & \beta & - & \beta \\ \beta & K\beta & \beta & \underline{\beta} \end{pmatrix}$$

$$\frac{\text{trs:trv}}{\text{ratio}} = \frac{\text{Expected no. trs.}}{\text{Expected no. trv}}$$

$$= \frac{K\beta\ell}{2\beta\ell} = \frac{K}{2} \quad \left\} \begin{array}{l} \text{true for} \\ \text{each row} \\ \text{of rate} \\ \text{matrix} \end{array}\right.$$

Thus, rate <u>ratio</u> not same as <u>ratio</u>