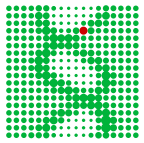


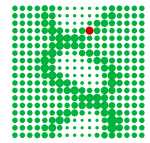
Quantifying the Equilibrium and Irreversibility Properties of the Nucleotide Substitution Process

Federico Squartini and Peter. F. Arndt

Max Planck Institute for Molecular Genetics



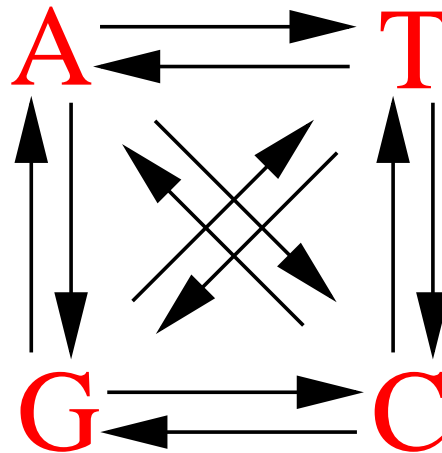
We will talk about **dis**equilibrium and
irreversibility. . .

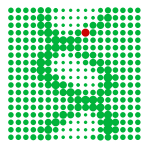


Markovian Sequence Evolution

Nucleotide substitution models: i.i.d Markov models of evolution, i.e. a master equation:

$$\frac{\partial}{\partial t} \rho_{\beta}(t) = \sum_{\alpha} Q_{\beta\alpha} \rho_{\alpha}(t) \quad \alpha, \beta \in \{A, G, C, T\}$$



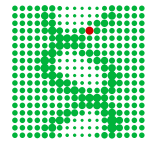


Markovian Sequence Evolution

Nucleotide substitution models: i.i.d Markov models of evolution, i.e. a master equation:

$$\frac{\partial}{\partial t} \rho_{\beta}(t) = \sum_{\alpha} Q_{\beta\alpha} \rho_{\alpha}(t) \quad \alpha, \beta \in \{A, G, C, T\}$$

$$Q = \begin{array}{c} \begin{array}{cccc} & A & C & G & T \\ A & \cdot & Q_{AC} & Q_{AG} & Q_{AT} \\ C & Q_{CA} & \cdot & Q_{CG} & Q_{CT} \\ G & Q_{GA} & Q_{GC} & \cdot & Q_{GT} \\ T & Q_{TA} & Q_{TC} & Q_{TG} & \cdot \end{array} \end{array}.$$



Markovian Sequence Evolution

Nucleotide substitution models: i.i.d Markov models of evolution, i.e. a master equation:

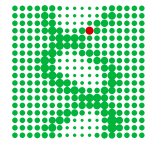
$$\frac{\partial}{\partial t} \rho_{\beta}(t) = \sum_{\alpha} Q_{\beta\alpha} \rho_{\alpha}(t) \quad \alpha, \beta \in \{A, G, C, T\}$$

The solution to this equation, with initial condition ρ_0 , is:

$$\rho_{\beta}(t) = \left[e^{Qt} \rho_0 \right]_{\beta}$$

$$P(t) = e^{Qt}$$

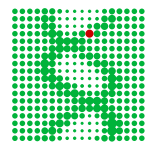
Such a model is not complete...



Choosing Parameters

Specifying an evolutionary mode \Rightarrow postulating a form for the rate matrix:

$$Q = \begin{array}{c} \text{A} \\ \text{C} \\ \text{G} \\ \text{T} \end{array} \begin{array}{c} \text{A} \quad \text{C} \quad \text{G} \quad \text{T} \\ \left(\begin{array}{cccc} \cdot & Q_{AC} & Q_{AG} & Q_{AT} \\ Q_{CA} & \cdot & Q_{CG} & Q_{CT} \\ Q_{GA} & Q_{GC} & \cdot & Q_{GT} \\ Q_{TA} & Q_{TC} & Q_{TG} & \cdot \end{array} \right) \end{array}.$$



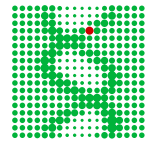
Choosing Parameters

Specifying an evolutionary mode \Rightarrow postulating a form for the rate matrix:

$$Q = \begin{array}{c} \text{A} \\ \text{C} \\ \text{G} \\ \text{T} \end{array} \begin{array}{cccc} \text{A} & \text{C} & \text{G} & \text{T} \\ \left(\begin{array}{cccc} \cdot & Q_{AC} & Q_{AG} & Q_{AT} \\ Q_{CA} & \cdot & Q_{CG} & Q_{CT} \\ Q_{GA} & Q_{GC} & \cdot & Q_{GT} \\ Q_{TA} & Q_{TC} & Q_{TG} & \cdot \end{array} \right) \end{array} \cdot$$

Widely used models:

$$Q = \begin{array}{c} \text{T} \\ \text{C} \\ \text{A} \\ \text{G} \end{array} \begin{array}{cccc} \text{T} & \text{C} & \text{A} & \text{G} \\ \left(\begin{array}{cccc} \cdot & \mu & \mu & \mu \\ \mu & \cdot & \mu & \mu \\ \mu & \mu & \cdot & \mu \\ \mu & \mu & \mu & \cdot \end{array} \right) \end{array} \cdot$$



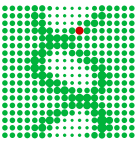
Choosing Parameters

Specifying an evolutionary mode \Rightarrow postulating a form for the rate matrix:

$$Q = \begin{array}{c} \text{A} \\ \text{C} \\ \text{G} \\ \text{T} \end{array} \begin{array}{cccc} \text{A} & \text{C} & \text{G} & \text{T} \\ \left(\begin{array}{cccc} \cdot & Q_{AC} & Q_{AG} & Q_{AT} \\ Q_{CA} & \cdot & Q_{CG} & Q_{CT} \\ Q_{GA} & Q_{GC} & \cdot & Q_{GT} \\ Q_{TA} & Q_{TC} & Q_{TG} & \cdot \end{array} \right) \end{array}.$$

Widely used models:

$$Q = \begin{array}{c} \text{T} \\ \text{C} \\ \text{A} \\ \text{G} \end{array} \begin{array}{cccc} \text{T} & \text{C} & \text{A} & \text{G} \\ \left(\begin{array}{cccc} \cdot & \alpha & \beta & \beta \\ \alpha & \cdot & \beta & \beta \\ \beta & \beta & \cdot & \alpha \\ \beta & \beta & \alpha & \cdot \end{array} \right) \end{array}.$$



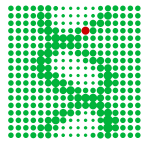
Choosing Parameters

Specifying an evolutionary mode \Rightarrow postulating a form for the rate matrix:

$$Q = \begin{array}{c} \text{A} \\ \text{C} \\ \text{G} \\ \text{T} \end{array} \begin{array}{cccc} \text{A} & \text{C} & \text{G} & \text{T} \\ \left(\begin{array}{cccc} \cdot & Q_{AC} & Q_{AG} & Q_{AT} \\ Q_{CA} & \cdot & Q_{CG} & Q_{CT} \\ Q_{GA} & Q_{GC} & \cdot & Q_{GT} \\ Q_{TA} & Q_{TC} & Q_{TG} & \cdot \end{array} \right) \end{array} \cdot$$

Widely used models:

$$Q = \begin{array}{c} \text{T} \\ \text{C} \\ \text{A} \\ \text{G} \end{array} \begin{array}{cccc} \text{T} & \text{C} & \text{A} & \text{G} \\ \left(\begin{array}{cccc} \cdot & \pi_T & \pi_T & \pi_T \\ \pi_C & \cdot & \pi_C & \pi_C \\ \pi_A & \pi_A & \cdot & \pi_A \\ \pi_G & \pi_G & \pi_G & \cdot \end{array} \right) \end{array} \cdot$$



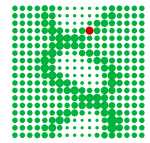
Choosing Parameters

Specifying an evolutionary mode \Rightarrow postulating a form for the rate matrix:

$$Q = \begin{array}{c} \text{A} \\ \text{C} \\ \text{G} \\ \text{T} \end{array} \begin{array}{cccc} \text{A} & \text{C} & \text{G} & \text{T} \\ \left(\begin{array}{cccc} \cdot & Q_{AC} & Q_{AG} & Q_{AT} \\ Q_{CA} & \cdot & Q_{CG} & Q_{CT} \\ Q_{GA} & Q_{GC} & \cdot & Q_{GT} \\ Q_{TA} & Q_{TC} & Q_{TG} & \cdot \end{array} \right) \cdot \end{array}$$

Widely used models:

$$Q = \begin{array}{c} \text{T} \\ \text{C} \\ \text{A} \\ \text{G} \end{array} \begin{array}{cccc} \text{T} & \text{C} & \text{A} & \text{G} \\ \left(\begin{array}{cccc} \cdot & k\pi_{\text{T}} & \pi_{\text{T}} & \pi_{\text{T}} \\ k\pi_{\text{C}} & \cdot & \pi_{\text{C}} & \pi_{\text{C}} \\ \pi_{\text{A}} & \pi_{\text{A}} & \cdot & k\pi_{\text{A}} \\ \pi_{\text{G}} & \pi_{\text{G}} & k\pi_{\text{G}} & \cdot \end{array} \right) \cdot \end{array}$$



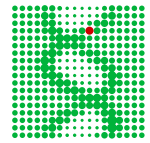
Choosing Parameters

Specifying an evolutionary mode \Rightarrow postulating a form for the rate matrix:

$$Q = \begin{array}{c} \text{A} \\ \text{C} \\ \text{G} \\ \text{T} \end{array} \begin{array}{cccc} \text{A} & \text{C} & \text{G} & \text{T} \\ \left(\begin{array}{cccc} \cdot & Q_{AC} & Q_{AG} & Q_{AT} \\ Q_{CA} & \cdot & Q_{CG} & Q_{CT} \\ Q_{GA} & Q_{GC} & \cdot & Q_{GT} \\ Q_{TA} & Q_{TC} & Q_{TG} & \cdot \end{array} \right) \end{array} \cdot$$

Widely used models:

$$Q = \begin{array}{c} \text{T} \\ \text{C} \\ \text{A} \\ \text{G} \end{array} \begin{array}{cccc} \text{T} & \text{C} & \text{A} & \text{G} \\ \left(\begin{array}{cccc} \cdot & k_1 \pi_T & \pi_T & \pi_T \\ k_1 \pi_C & \cdot & \pi_C & \pi_C \\ \pi_A & \pi_A & \cdot & k_2 \pi_A \\ \pi_G & \pi_G & k_2 \pi_G & \cdot \end{array} \right) \end{array} \cdot$$



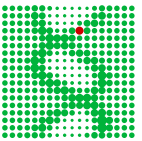
Two Evolutionary Models

All preceding models are nested into the following:

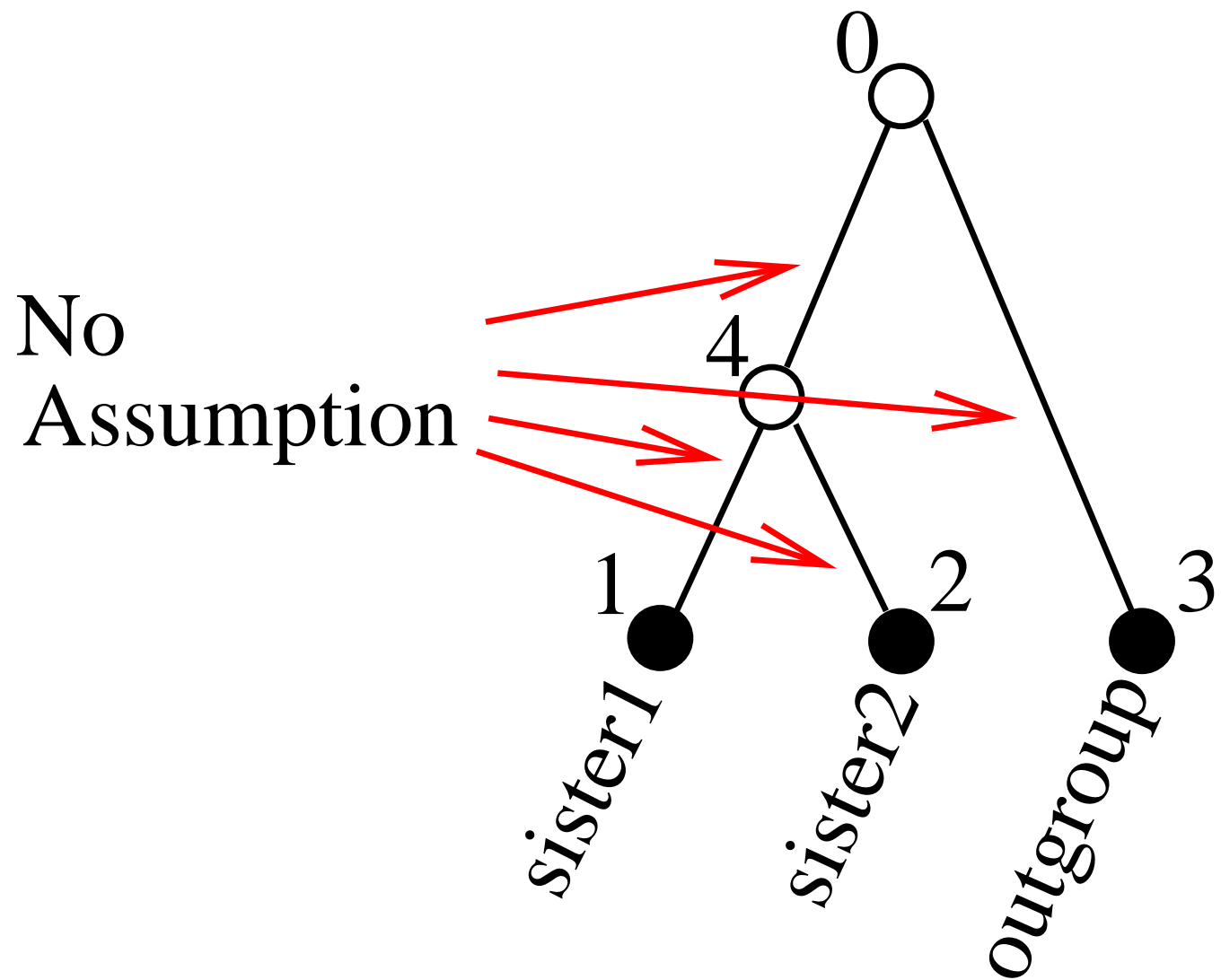
$$Q_{\text{GTR}} = \begin{matrix} & \text{A} & \text{G} & \text{T} & \text{C} \\ \text{A} & \left(\begin{array}{cccc} \cdot & a\pi_{\text{A}} & b\pi_{\text{A}} & c\pi_{\text{A}} \\ a\pi_{\text{G}} & \cdot & d\pi_{\text{G}} & e\pi_{\text{G}} \\ b\pi_{\text{T}} & d\pi_{\text{T}} & \cdot & f\pi_{\text{T}} \\ c\pi_{\text{C}} & e\pi_{\text{C}} & f\pi_{\text{C}} & \cdot \end{array} \right) \end{matrix} \cdot$$

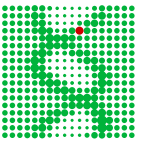
A possible alternative:

$$Q_{\text{RCS}} = \begin{matrix} & \text{A} & \text{C} & \text{G} & \text{T} \\ \text{A} & \left(\begin{array}{cccc} \cdot & r_{\text{AC}} & r_{\text{AG}} & r_{\text{AT}} \\ r_{\text{GT}} & \cdot & r_{\text{CG}} & r_{\text{CT}} \\ r_{\text{CT}} & r_{\text{CG}} & \cdot & r_{\text{GT}} \\ r_{\text{AT}} & r_{\text{AG}} & r_{\text{AC}} & \cdot \end{array} \right) \end{matrix} \cdot$$

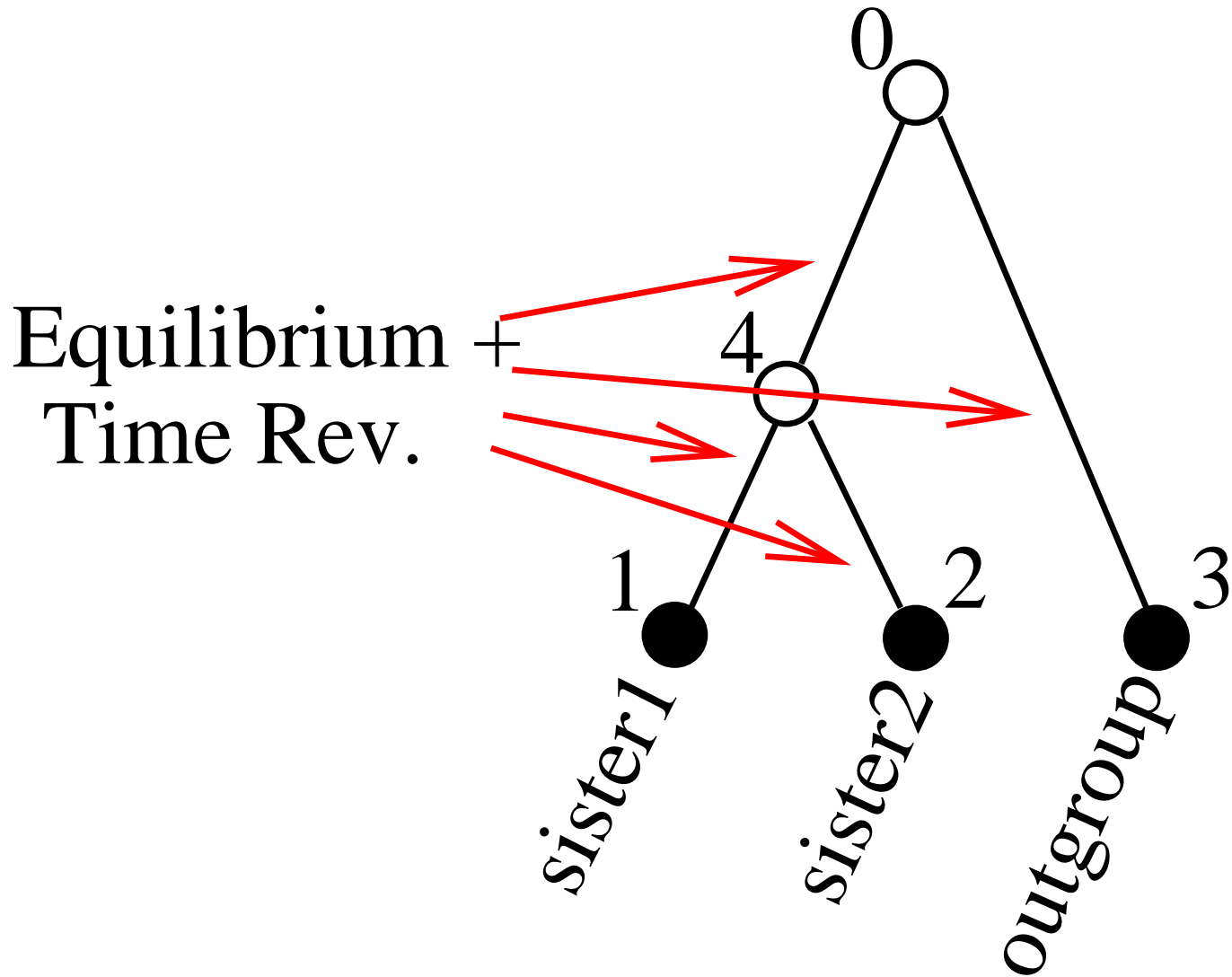


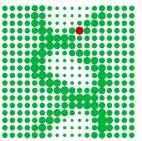
Two Evolutionary Models - 2



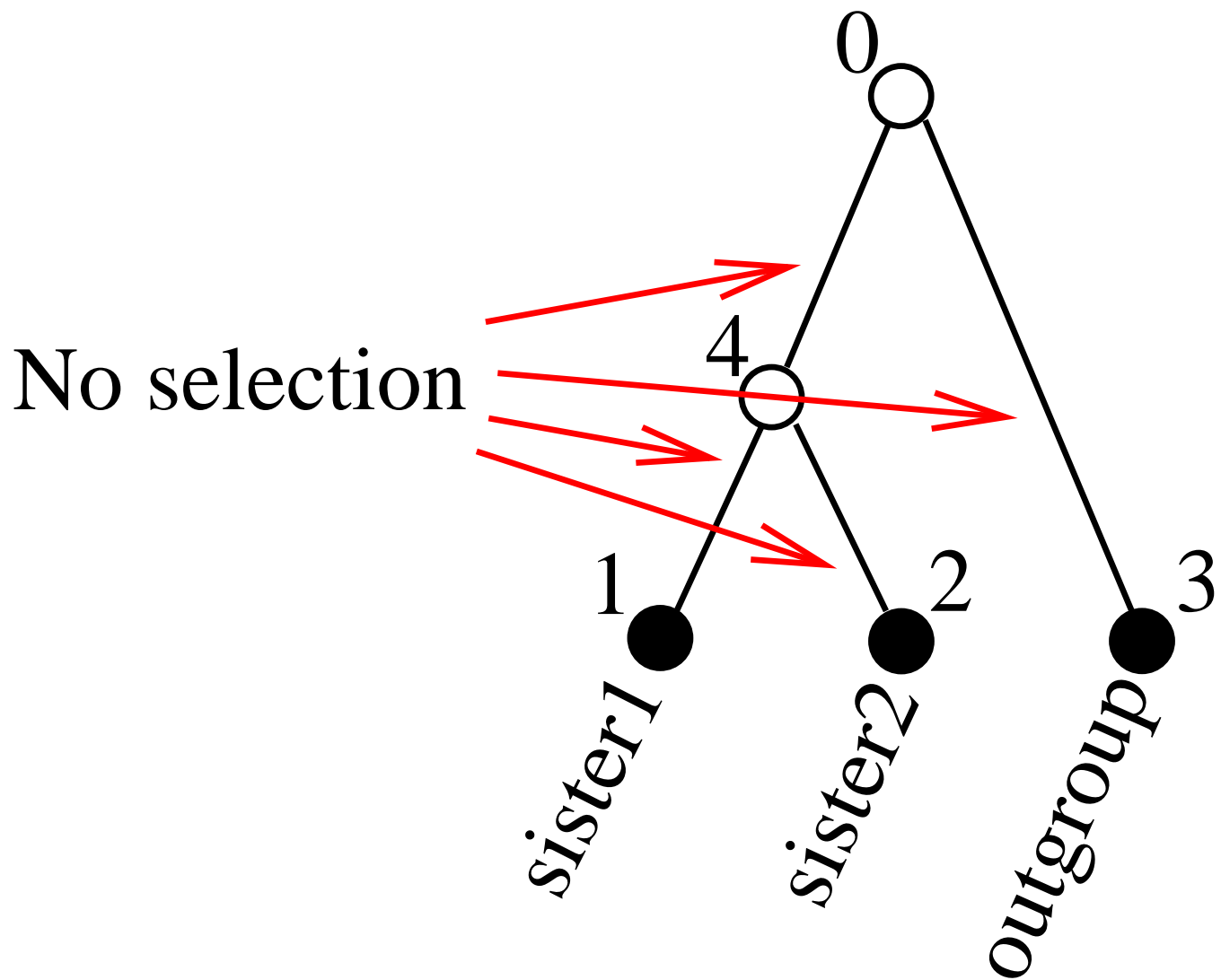


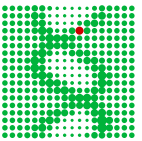
Two Evolutionary Models - 2





Two Evolutionary Models - 2



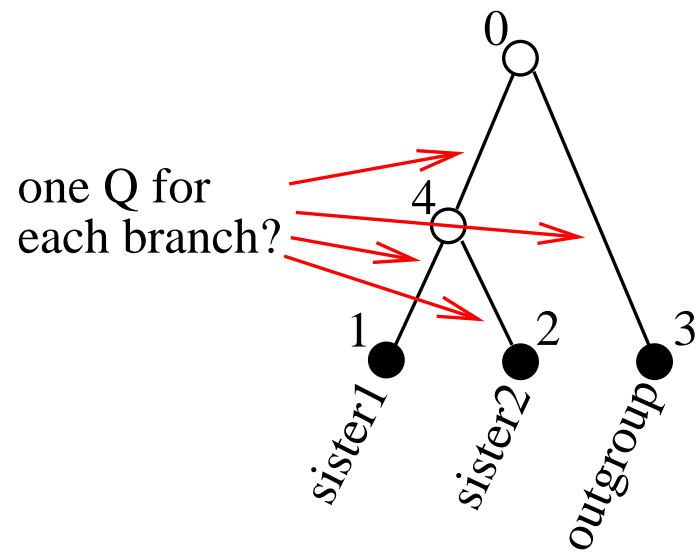


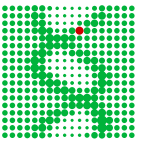
Estimating Parameters

For a given triple alignment $\vec{\alpha}^i$ of nucleotide sequences from 3 species, the likelihood of the alignment is:

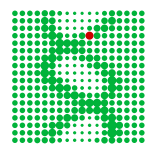
$$L = \prod_{k=1}^N \sum_{\alpha^0, \alpha^4 \in \{A, C, G, T\}} \rho_{\alpha^0}^0 [P^{30}]_{\alpha_k^3 \alpha^0} [P^{40}]_{\alpha^4 \alpha^0} [P^{24}]_{\alpha_k^2 \alpha^4} [P^{14}]_{\alpha_k^1 \alpha^4}$$

The vector ρ^0 represents the ancestral nucleotide distribution at the root node.





Equilibrium



The stationarity index

The equilibrium distribution of a Markov process is defined by:

$$Q\pi = 0$$

Just taking the difference between present and stationary distribution:

$$\Delta_{\alpha} = \rho_{\alpha} - \pi_{\alpha}$$

And rearrange the terms:

$$\text{STI}_1 = \Delta_{\text{C}} + \Delta_{\text{G}} = \rho_{\text{GC}} - \pi_{\text{GC}}$$

$$\text{STI}_2 = \Delta_{\text{A}} - \Delta_{\text{T}}$$

$$\text{STI}_3 = \Delta_{\text{C}} - \Delta_{\text{G}},$$



The STI - Reverse complement symmetry

Substituting the equilibrium distribution:

$$(1 - \pi_{CG}, \pi_{CG}, \pi_{CG}, 1 - \pi_{CG})$$

Where:

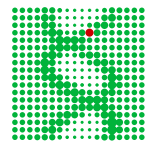
$$\pi_{CG} = \frac{r_{GT} + r_{CT}}{r_{AC} + r_{AG} + r_{GT} + r_{CT}}$$

For the reverse complement symmetric model the STI has a simple form:

$$STI_1 = \rho_{GC} - \pi_{GC}$$

$$STI_2 = (\rho_A - \rho_T)$$

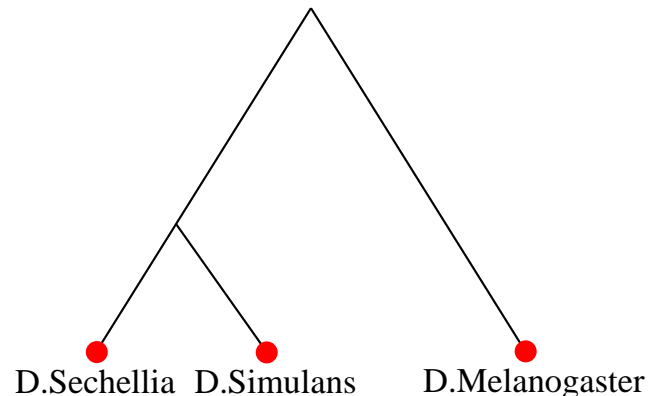
$$STI_3 = (\rho_C - \rho_G).$$



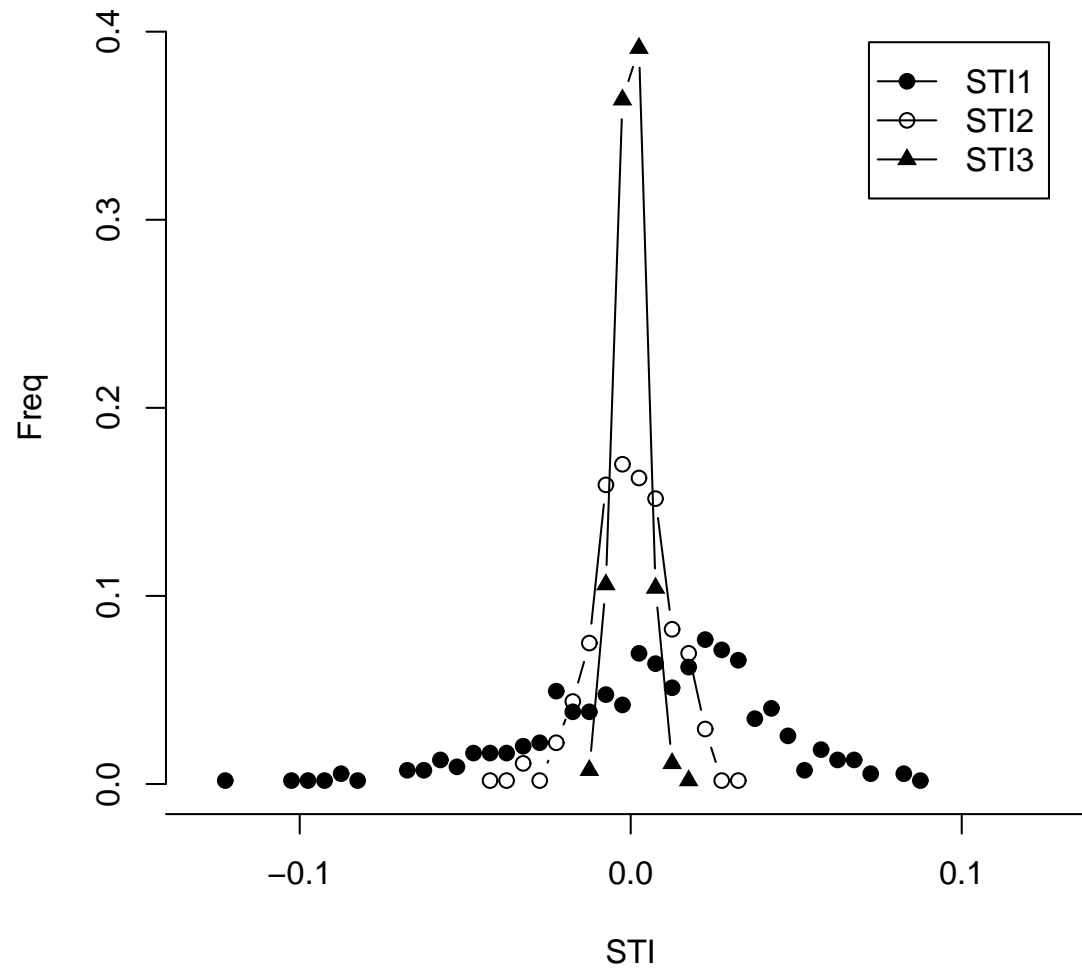
Analysis of the Fly Genome

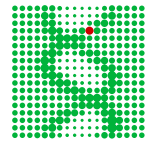
Results about the time reversal properties for the evolution of the fly genome:

- ▶ Alignment of 3 Drosophilas: sechellia, simulans and melanogaster
- ▶ Removed annotated coding regions
- ▶ Rates have been estimated using a maximum likelihood algorithm
- ▶ Sliding window analysis, 50kbp length
- ▶ For each window we have calculated the stationarity index in the simulans lineage



Analysis of the Fly Genome - Stationarity





Reversibility



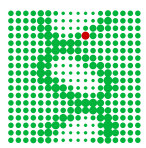
Time Reversibility: the Detailed Balance

Time reversibility is usually defined in terms of the **detailed balance conditions**:

$$Q_{ji}\pi_i = Q_{ij}\pi_j$$

From which one can derive the General Time Reversible (GTR) Parameterization:

$$Q_{\text{GTR}} = \begin{array}{c} \text{A} \\ \text{G} \\ \text{T} \\ \text{C} \end{array} \begin{array}{c} \text{A} \quad \text{G} \quad \text{T} \quad \text{C} \\ \left(\begin{array}{cccc} \cdot & a\pi_{\text{A}} & b\pi_{\text{A}} & c\pi_{\text{A}} \\ a\pi_{\text{G}} & \cdot & d\pi_{\text{G}} & e\pi_{\text{G}} \\ b\pi_{\text{T}} & d\pi_{\text{T}} & \cdot & f\pi_{\text{T}} \\ c\pi_{\text{C}} & e\pi_{\text{C}} & f\pi_{\text{C}} & \cdot \end{array} \right)$$

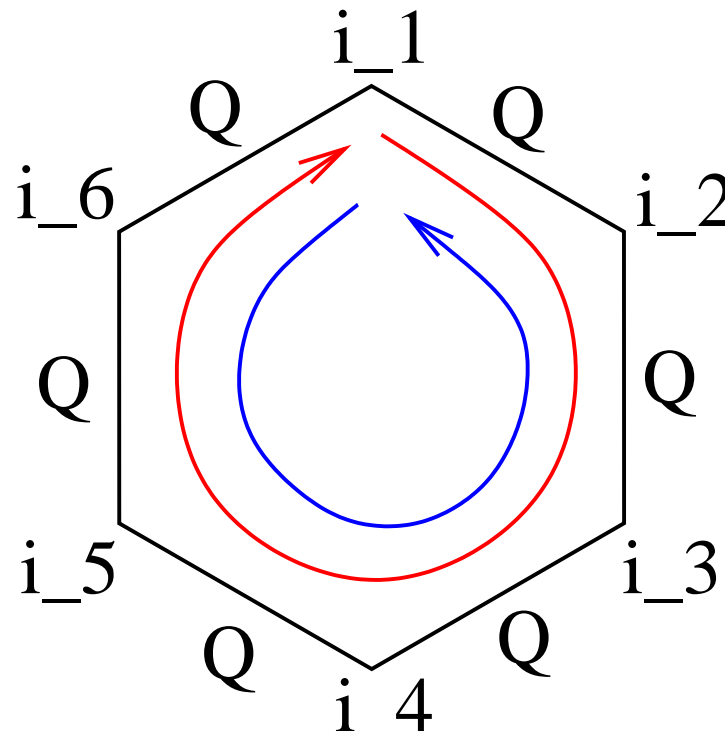


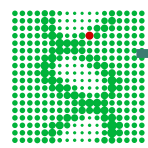
Time reversibility: Kolmogorov Cycle Conditions

A lesser known formulation of time reversibility:

Definition. A Markov process is said to satisfy the Kolmogorov cycle conditions if the following equality on generators holds:

$$Q_{i_1 i_n} Q_{i_n i_{n-1}} \cdots Q_{i_2 i_1} = Q_{i_1 i_2} \cdots Q_{i_{n-1} i_n} Q_{i_n i_1} \quad \forall i_1, \dots, i_n \in \mathcal{C} \quad (-2)$$





Time reversibility: Kolmogorov Cycle Conditions - 2

Moreover the following proposition (relevant when analyzing biological sequences) holds:

Proposition. *If the coefficients of the rate matrix are strictly positive and if Kolmogorov conditions hold for three cycles then they hold for cycles of arbitrary length.*

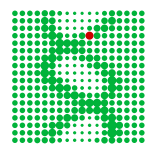
Proposition. *Given a four states Markov process with strictly positive rate matrix coefficients, if the conditions:*

$$Q_{\alpha\delta}Q_{\delta\gamma}Q_{\gamma\beta}Q_{\beta\alpha} = Q_{\alpha\beta}Q_{\beta\gamma}Q_{\gamma\delta}Q_{\delta\alpha}, \quad (-2)$$

hold for $(\alpha, \beta, \gamma, \delta)$ equal to (A, G, C, T) , (A, G, T, C) and (A, C, G, T) then Kolmogorov conditions hold for 3-cycles.

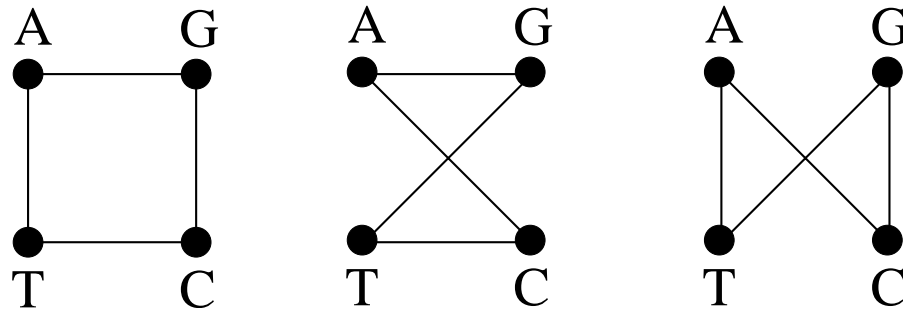
Ans lastly:

Proposition. *If the coefficients of the rate matrix are strictly positive and if Kolmogorov conditions hold for four cycles then they hold for cycles of arbitrary length.*



IRI - The general iid case

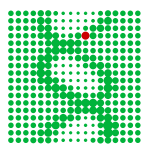
To check reversibility for nucleotide sequences we need to check the following conditions on four cycles:



$$\text{IRI}_1 := \frac{Q_{AG}Q_{GC}Q_{CT}Q_{TA} - Q_{AT}Q_{TC}Q_{CG}Q_{GA}}{Q_{AG}Q_{GC}Q_{CT}Q_{TA} + Q_{AT}Q_{TC}Q_{CG}Q_{GA}}$$

$$\text{IRI}_2 := \frac{Q_{AC}Q_{CT}Q_{TG}Q_{GA} - Q_{AG}Q_{GT}Q_{TC}Q_{CA}}{Q_{AC}Q_{CT}Q_{TG}Q_{GA} + Q_{AG}Q_{GT}Q_{TC}Q_{CA}}$$

$$\text{IRI}_3 := \frac{Q_{AC}Q_{CG}Q_{GT}Q_{TA} - Q_{AT}Q_{TG}Q_{GC}Q_{CA}}{Q_{AC}Q_{CG}Q_{GT}Q_{TA} + Q_{AT}Q_{TG}Q_{GC}Q_{CA}}$$



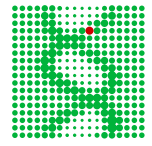
Iri for the Reverse Complement Symmetric Model

Out of the previous indices we get a specialized version of the IRI:

$$\begin{aligned} \text{IRI}_1 &= \frac{r_{\text{AG}}^2 r_{\text{GT}}^2 - r_{\text{AC}}^2 r_{\text{CT}}^2}{r_{\text{AG}}^2 r_{\text{GT}}^2 + r_{\text{AC}}^2 r_{\text{CT}}^2} \\ \text{IRI}_2 &= 0 \\ \text{IRI}_3 &= 0 \end{aligned}$$

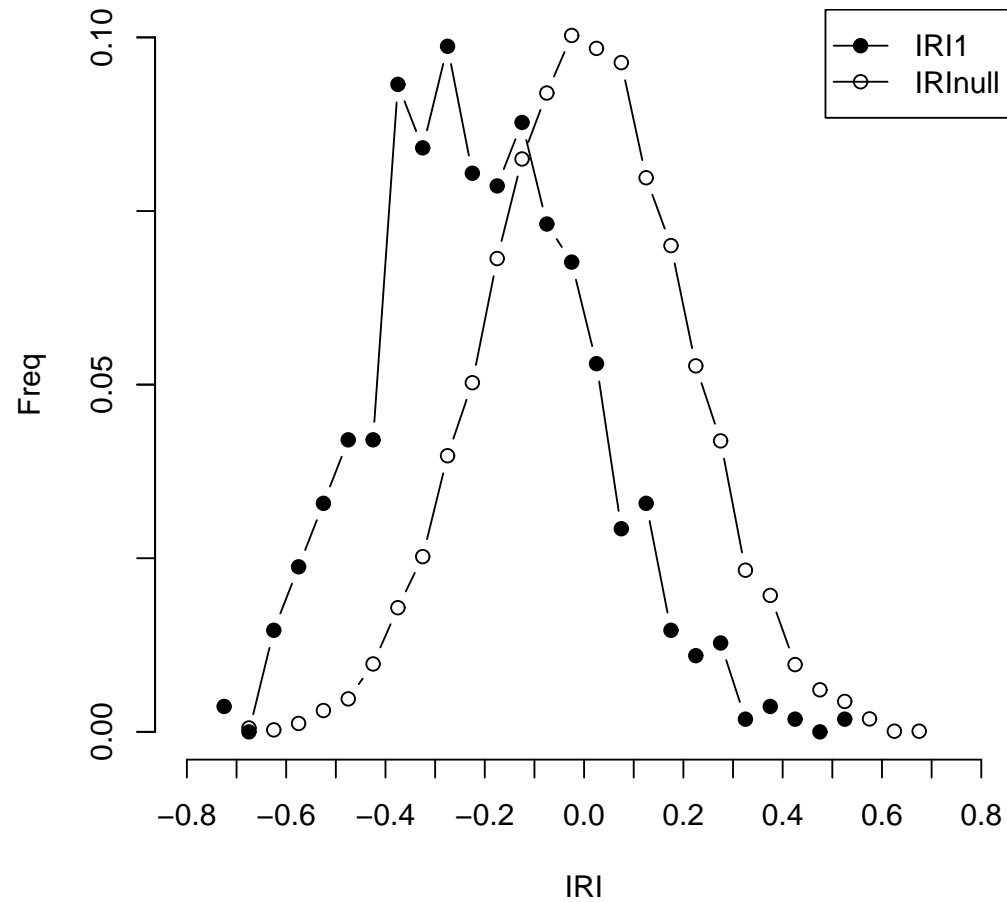
The IRI_1 will thus be comprised in the interval $[-1, 1]$ and if the system under study evolves time symmetrically:

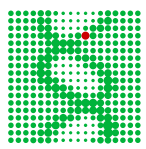
$$\text{IRI}_1 = 0$$



Irreversibility in the Fly Genome

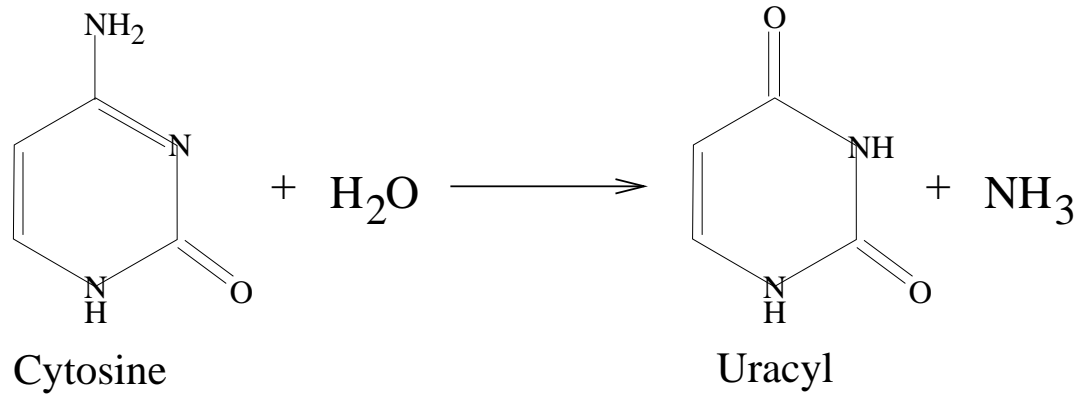
Plots of the IRI for the *Drosophila simulans* genome and for the null model:



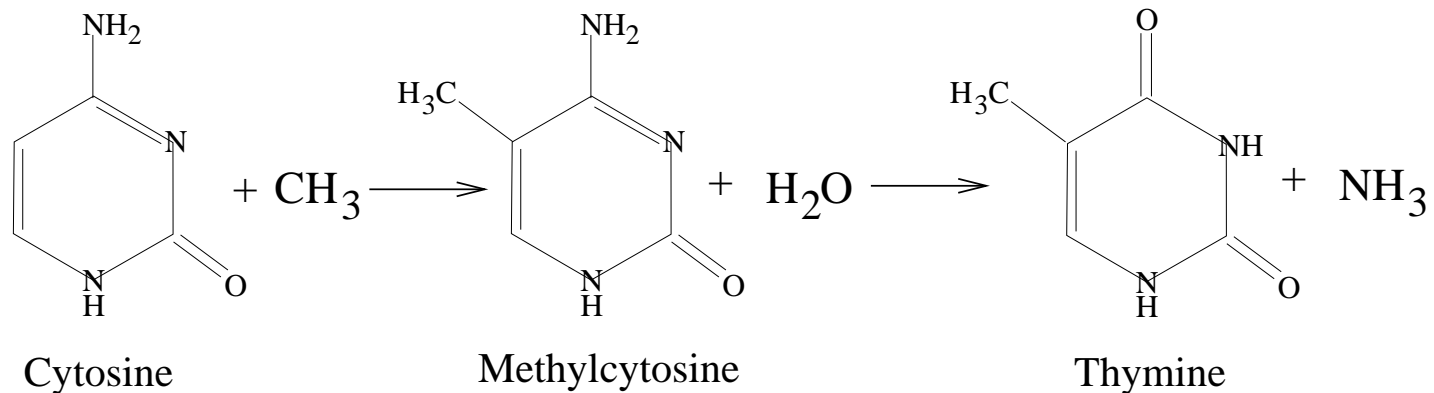


If water is around...

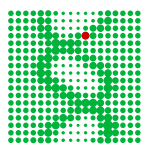
Cytosine can easily decay into Uracyl:



On the other hand GpC pairs often occur in a methylated form:



The net effect is the decay of CpG pairs into TpG and CpA pairs.



A Nucleotide Substitution Model with CpG Decay

We need to extend the configuration space:

$$\mathcal{C} = s_1 \times \dots \times s_N \quad s_i \in \{A, C, G, T\}.$$

We assume the following form for the generator:

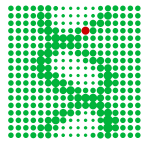
$$Q = \sum_{i=1}^N Q_i + \sum_{i=1}^{N-1} Q_{i,i+1}^{\text{CpG}}.$$

Where:

$$Q_i = \underbrace{\mathbb{I} \otimes \dots \otimes \mathbb{I}}_{i-1} \otimes Q \otimes \underbrace{\mathbb{I} \otimes \dots \otimes \mathbb{I}}_{N-i}.$$

And:

$$Q_{i,i+1}^{\text{CpG}} = \underbrace{\mathbb{I} \otimes \dots \otimes \mathbb{I}}_{i-1} \otimes Q^{\text{CpG}} \otimes \underbrace{\mathbb{I} \otimes \dots \otimes \mathbb{I}}_{N-i-1}.$$

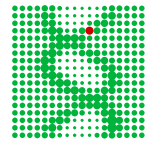


The IRI of a Process with CpG Decay

We get two IRI's in this case:

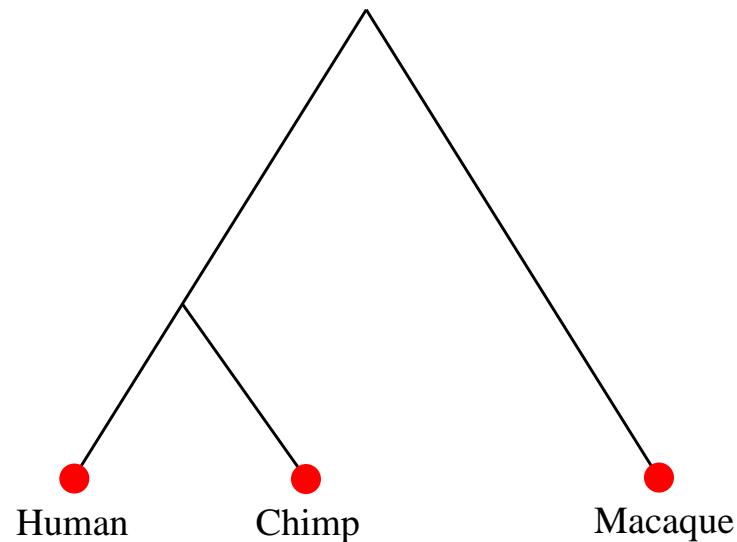
$$\text{IRI}_1 := \frac{r_{\text{AG}}^2 r_{\text{GT}}^2 - r_{\text{AC}}^2 r_{\text{CT}}^2}{r_{\text{AG}}^2 r_{\text{GT}}^2 + r_{\text{AC}}^2 r_{\text{CT}}^2}$$

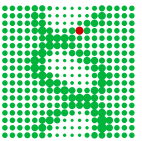
$$\text{IRI}_{\text{CpG}} := \frac{r_{\text{GT}}^2 (r_{\text{AG}} + r_{\text{CpG}})^2 - (r_{\text{CT}} + r_{\text{CpG}}^{\text{rev}})^2 r_{\text{AC}}^2}{r_{\text{GT}}^2 (r_{\text{AG}} + r_{\text{CpG}})^2 + (r_{\text{CT}} + r_{\text{CpG}}^{\text{rev}})^2 r_{\text{AC}}^2}$$



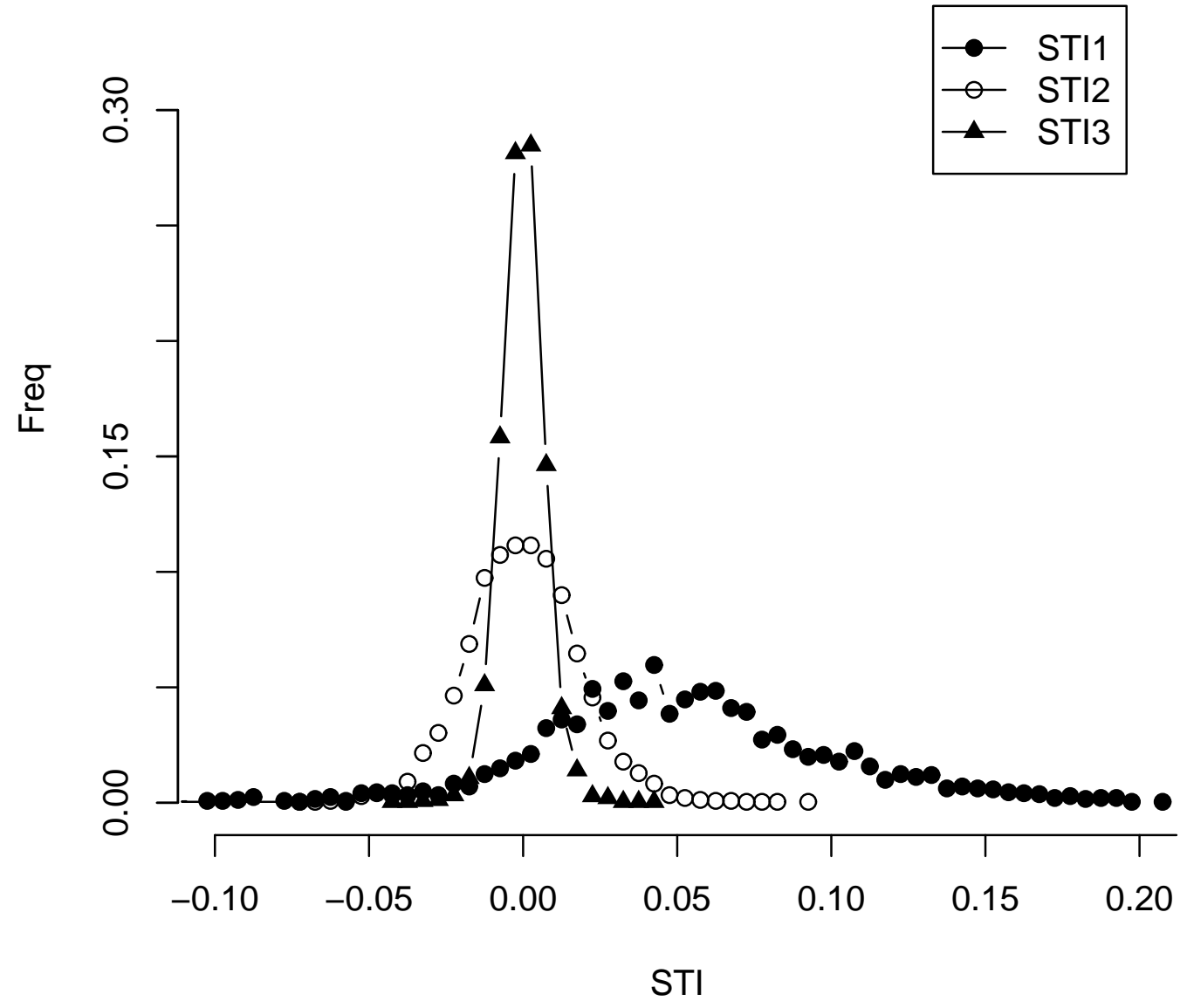
Analysis of the Human Genome

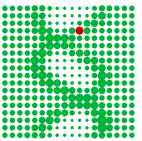
- ▶ Alignment of Human, Chimp and Rhesus Macaque genomes
- ▶ Rates have been estimated using a maximum likelihood algorithm
- ▶ Sliding window analysis, 1 Mbp length
- ▶ For each window we have calculated the STIs, IRI_{RC} and IRI_{CpG} in the human lineage



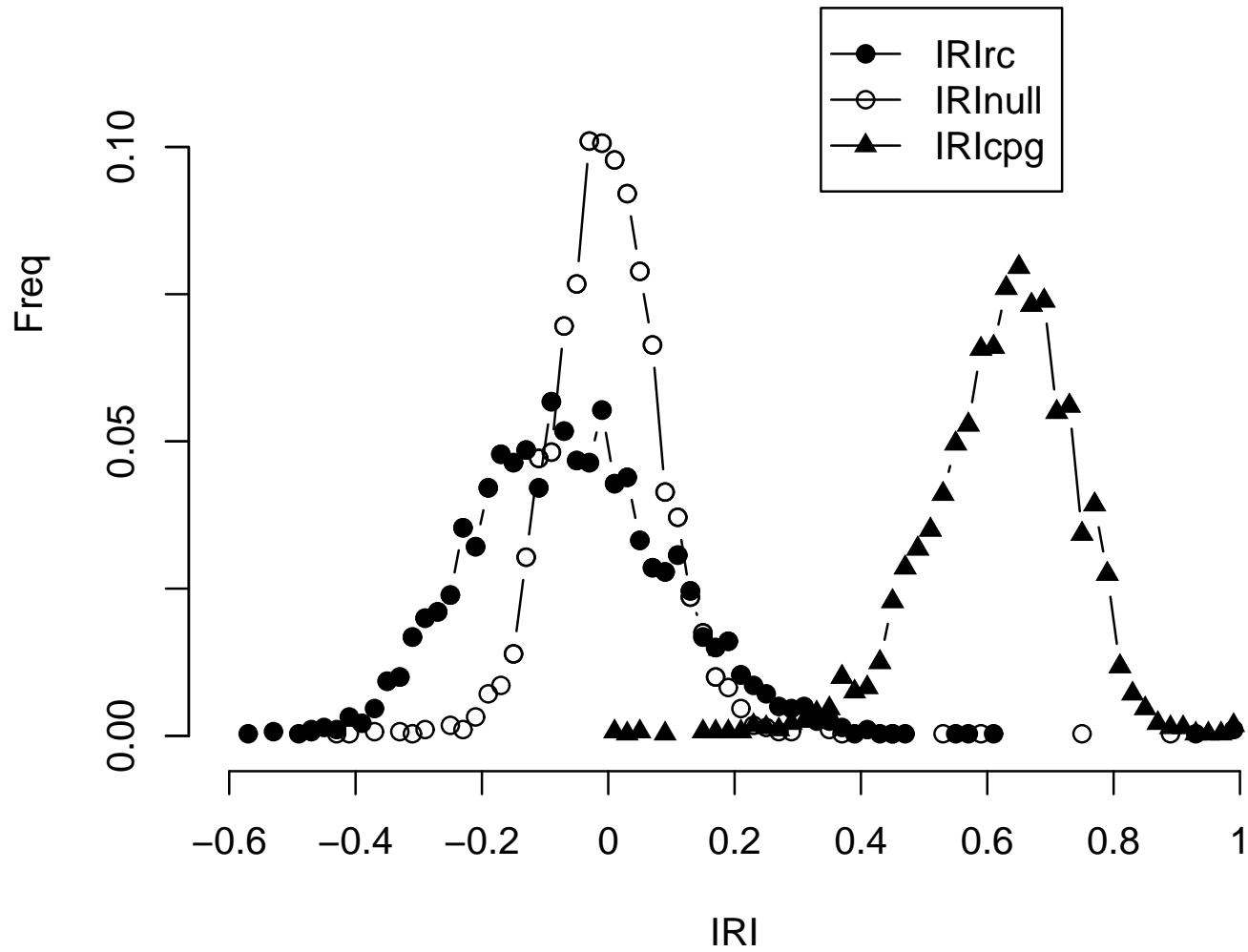


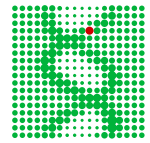
STI Human





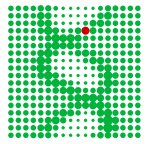
IRI Human





Summary

- ▶ Commonly used evolutionary models assume equilibrium and reversibility
- ▶ We have introduced indices to test for equilibrium (STI) and reversibility (IRI) on each single branch of a given phylogeny
- ▶ Analysis in *Drosophila* and Human show clear violation of the equilibrium/reversibility.
- ▶ Further work has to be done to assess how these violations affect specific bioinformatic algorithms.



It's Evolution Baby...

Thank you!