Raazesh Sainudiin $^{\circ}$

^o Biomathematics Research Centre

Department of Mathematics & Statistics, University of Canterbury, NZ,

[‡] Department of Statistics, University of Oxford, UK, and

*Biological Sciences, University of California, Irvine, USA.

Joint with: Peter Donnelly[‡], Bob Griffiths[‡], Gil McVean[‡], and Kevin Thornton^{\circledast}

Outline – Talk Outline

- The Coalescent Models
- Computationally Intensive Likelihoods
- A Paritllay-ordered Coalescent Experiments Graph
- Unlabeled n-Coalescent
- Results
- Summary
- Acknowledgments

Data and Model 1: $\phi \equiv \theta \in \Phi, \theta = 4N_e\mu$ (scaled mutation rate)

The Wright-Fisher Model – Random Mating, Constant Size, No Recombination/Selection

A Population of N = 10 homologous DNA seqns. of length m and the Population History of site i



Data and Model 1: $\phi \equiv \theta \in \Phi, \theta = 4N_e\mu$ (scaled mutation rate)

The Wright-Fisher Model – Random Mating, Constant Size, No Recombination/Selection Ex: Data of 3 homologous DNA sequences at site *i*, its Population History and the Sample History of sampled individuals 1,2, and 3.

: 1 2 3 i : T T A



Model 1: $\phi \equiv \theta \in \Phi, \, \theta = 4 N_e \mu \,$ (scaled mutation rate)

The Coalescent Approximation of the Wright-Fisher (W-F) Model (Kingman, 1982)

A Sample Coalescent Sequence or c-sequence ($\{\{1\}, \{2\}, \{3\}\}, \{\{1, 2\}, \{3\}\}, \{\{1, 2, 3\}\}$)

and coalescent times or epoch times $t_i, i \in \{3, 2\}$.

Offspring "choose" parents uniformly and independently in W-F model







Coalescent Space $\mathcal{A}_n \equiv \mathcal{C}_n \otimes (0, \infty)^{n-1}$ when n = 3 (Model 1)



Realizations from $A_n \equiv C_n \otimes (0, \infty)^{n-1}$ **under Model 1,** n = 6, 32



Raazesh Sainudiin, Department of Mathematics & Statistics, University of Canterbury www.math.canterbury.ac.nz/~r.sainudiin - p. 8/27

Model 2 :
$$\phi \equiv (heta,
u) \in \Phi$$
, $heta = 4N_e\mu$ (scaled mutn. rate) , u (exp. growth rate)



Figures 1-6 of M. Nordburg, Coalescent Theory, 2000

Coalescent Sample Spaces – Partially Ordered Experiments Graph



(1) Every directed acyclic subgraph of the POEG indexes a Martingale

(2) Each node of the POEG is a tri-sequential asymptotic family of Experiments

Likelihood

Likelihood, $P(D|\phi)$, is computed by Integrating Missing-Data:

$$\sum_{c \in \mathcal{C}_n} \int_{t \in (0,\infty)^{n-1}} P(D|c,t,\phi) P(c,t|\phi) dt \, dc$$

Cardinalities of the state spaces of the standard *n*-coalescent on \mathbb{C}_n and the unlabeled *n*-coalescent on \mathbb{F}_n (to be seen in the sequel).

n	4	10	30	60	90
$ \mathbb{C}_n $	15	1.2×10^5	8.5×10^{23}	$9.8 imes 10^{59}$	1.4×10^{101}
$ \mathbb{F}_n $	5	42	5.6×10^3	$9.7 imes 10^5$	5.7×10^7
$ \mathbb{F}_n / \mathbb{C}_n $	0.33	3.6×10^{-4}	6.6×10^{-21}	9.9×10^{-55}	4.0×10^{-94}

Likelihood is computationally prohibitive at MSA/BIM Resolns.



MSA 10,000 Auto-validating i.i.d. Posterior Samples in MRS SY2006 - novel (3/4 leaved phylogenetic tree spaces) $pprox 200 \ {
m CPU}$ sec for n <= 3, :-(\rightarrow impractical for n > 4Complete Recursion in PTREE G1980 (1 Locus, $\theta = 10$, C-Model 1) :-(\rightarrow out of stack for n > 4Approximate Methods : MSA MCMC in COALESCE KYF1998 : n < 200 & **BIM** SIS in GENETREE GT1994 : $L(\theta|v) \cong 4$ CPU

- The Bottom Line: Exact Genome Scanning at fine DNA resolution is currently impractical for n > 4
- A Solution: Inference at coarser empirical resolutions, eg. SFS and its sub-experiments – novel

∞ -many-sites M-Model: BIM $v \in \mathcal{V}_n^m \to SFS \ x \in \mathcal{X}_n^m$

Let $v \in \mathcal{V}_n^m \equiv \{0,1\}^{n \times m}$ be a BIM, then the SFS $x \equiv (x_1, \dots, x_{n-1}) \in \mathcal{X}_n^m \equiv \{x \in \mathbb{Z}_+^{n-1} : \sum_{i=1}^{n-1} x_i \le m\}$ $x_i = N_i(v^T \cdot (1, 1, \dots, 1)), \qquad N_i(y_1, y_2, \dots, y_s) = \sum_{i=1}^s \mathbf{1}_{\{i\}}(y_j), \qquad i = 1, \dots, n-1.$ BIM $v \in \mathcal{V}_4^9$ t_2 t_3

 SITES:
 1
 2
 3
 4
 5
 6
 7
 8
 9

 IND 1:
 0
 0
 1
 0
 0
 1
 0
 0
 1
 0
 0
 1
 0
 0
 1
 0
 0
 1
 0
 0
 1
 0
 0
 1
 0
 0
 1
 0
 0
 1
 0
 0
 1
 0
 0
 1
 1
 0
 0
 0
 0
 1
 1
 0
 0
 0
 0
 1
 1
 0
 0
 0
 0
 0
 1
 1
 0
 0
 0
 0
 0
 1
 1
 0
 0
 0
 0
 0
 1
 1
 0
 1
 0
 0
 0
 0
 1
 0
 1
 1
 0
 1
 0
 1
 0
 1
 0
 1
 1
 0
 1
 1
 0
 1
 1
 0
 1
 1
 0
 1
 1
 0
 1
 1
 1
 0
 < 3 9 7 SFS $x = (x_1, x_2, x_3) = (2, 1, 2) \in \mathcal{X}_4^9$

Coalescent Tree Shape, *f***-Sequence and Site Frequency Spectrum**



Examples of *c***-sequence** \rightarrow *f***-sequence, when** n = 4



Ex 1:

$$\begin{split} & [\{1\},\{2\},\{3\},\{4\}], [\{1,2\},\{3\},\{4\}], [\{1,2,3\},\{4\}], [\{1,2,3,4\}] \rightarrow \\ & [(4,0,0,0), (2,1,0,0), (1,0,1,0), (0,0,0,1)] \\ & \text{Ex 2:} \\ & [\{1\},\{2\},\{3\},\{4\}], [\{1,2\},\{3\},\{4\}], [\{1,2\},\{3,4\}], [\{1,2,3,4\}] \rightarrow \\ & [(4,0,0,0), (2,1,0,0), (0,2,0,0), (0,0,0,1)] \end{split}$$

Consider, the integer partitions of n with i blocks:

$$\mathbb{F}_{n}^{i} \equiv \{f_{i} \equiv (f_{i,1}, f_{i,2}, \dots, f_{i,n}) \in \mathbb{Z}_{+}^{n} : \sum_{j=1}^{n} jf_{i,j} = n, \sum_{j=1}^{n} f_{i,j} = i\}.$$

where $f_{i,j}$ denotes the number of lineages subtending j leaves at the *i*-th epoch. **Proposition** (Kingman's Unlabeled *n*-coalescent). It is the continuous time Markov chain on $\mathbb{F}_n \equiv \bigcup_{i=1}^n \mathbb{F}_n^i$, the set of integer partitions of n, whose infinitesimal generator $\mathbf{q}(f_h|f_g)$ for any two states $f_g, f_h \in \mathbb{F}_n$ is:

$$\mathbf{q}(f_{h}|f_{g}) = \begin{cases} -i(i-1)/2 & : \text{if} \quad f_{g} = f_{h}, \, f_{g} \in \mathbb{F}_{n}^{i} \\ f_{g,j}f_{g,k} & : \text{if} \quad f_{h} = f_{g} - e_{j} - e_{k} + e_{j+k}, \, j \neq k, \, f_{g} \in \mathbb{F}_{n}^{i}, \, f_{h} \in \mathbb{F}_{n}^{i-1} \\ (f_{g,j})(f_{g,j} - 1)/2 & : \text{if} \quad f_{h} = f_{g} - e_{j} - e_{k} + e_{j+k}, \, j = k, \, f_{g} \in \mathbb{F}_{n}^{i}, \, f_{h} \in \mathbb{F}_{n}^{i-1} \\ 0 & : \text{ otherwise} \end{cases}$$

Initial state: $f_n = (n, 0, 0, \dots, 0)$ and absorbing state: $f_1 = (0, 0, \dots, 1)$.

Any realization of the chain is an f-sequence: $f = (f_n, f_{n-1}, \ldots, f_1) \in \mathcal{F}_n$.

Proposition (Probability of an f_i). The probability of an $f_i \in \mathbb{F}_n^i$ is:

$$P(f_i) = \frac{i!}{\prod_{j=1}^{i} f_{i,j}!} {\binom{n-1}{i-1}}^{-1}$$

Proposition (Probability of an f_i). The probability of an $f_i \in \mathbb{F}_n^i$ is:

$$P(f_i) = \frac{i!}{\prod_{j=1}^{i} f_{i,j}!} {\binom{n-1}{i-1}}^{-1}$$

Proposition (Probability of an f-sequence).

$$P(f) = \prod_{i=2}^{n} P(f_i | f_{i-i}) = \frac{2^{\neg(f)}}{(n-1)!} \prod_{i=2}^{n} \ddot{f}_i$$

where,

- $\neg(f)$ is the number of distinctly-sized lineage splits
- f_i is the number of lineages at the beginning of the *i*-th epoch that subtend the same number of leaves as the lineage that was split then.

c-sequence,
$$c \in \mathcal{C}_n \to c$$
-shape, $\tilde{c} \in \mathcal{C}_n \to f$ -sequence, $f \in \mathcal{F}_n$

 \sim

$$\tilde{c}^{\wedge} = ((\cdot, _{3} \cdot), _{1} (\cdot, _{2} \cdot))$$

$$\tilde{c}^{\wedge} = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 2 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 4 & 0 & 0 & 0 \end{pmatrix}$$

$$\mathbf{f}^{\wedge} = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 0 & 0 \\ 4 & 0 & 0 & 0 \end{pmatrix}$$

$$\mathbf{f}^{\wedge} = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 0 & 0 \\ 4 & 0 & 0 & 0 \end{pmatrix}$$

Raazesh Sainudiin, Department of Mathematics & Statistics, University of Canterbury www.math.canterbury.ac.nz/~r.sainudiin - p. 17/27

c-sequence,
$$c \in \mathcal{C}_n \to c$$
-shape, $\tilde{c} \in \mathcal{C}_n \to f$ -sequence, $f \in \mathcal{F}_n$

 \sim

$$n = 5$$

$$\tilde{c}^{(a)} = ((((\cdot, 4 \cdot), 3 \cdot), 2 \cdot), 1 \cdot)$$

$$f^{a} = \begin{pmatrix} 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \\ 2 & 0 & 1 & 0 & 0 \\ 3 & 1 & 0 & 0 & 0 \\ 5 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$f^{b} = \begin{pmatrix} 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \\ 1 & 2 & 0 & 0 & 0 \\ 3 & 1 & 0 & 0 & 0 \\ 5 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$\tilde{c}^{(c)} = ((\cdot, 4 \cdot), 1 ((\cdot, 3 \cdot), 2 \cdot)) \quad \tilde{c}^{(d)} = (((\cdot, 4 \cdot), 2 \cdot), 1 (\cdot, 3 \cdot)) \quad \tilde{c}^{(e)} = (((\cdot, 4 \cdot), 3 \cdot), 1 (\cdot, 2 \cdot))$$

$$\tilde{c}^{(c)} = \begin{pmatrix} 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 5 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$\tilde{c}^{(e)} = (((\cdot, 4 \cdot), 3 \cdot), 1 (\cdot, 2 \cdot)) \quad \tilde{c}^{(d)} = (((\cdot, 4 \cdot), 2 \cdot), 1 (\cdot, 3 \cdot)) \quad \tilde{c}^{(e)} = (((\cdot, 4 \cdot), 3 \cdot), 1 (\cdot, 2 \cdot))$$

$$\tilde{c}^{(e)} = ((0 \cdot 4 \cdot), 1 (\cdot, 3 \cdot), 1 (\cdot, 2 \cdot)) \quad \tilde{c}^{(e)} = ((0 \cdot 4 \cdot), 3 \cdot), 1 (\cdot, 2 \cdot))$$

$$\tilde{c}^{(e)} = ((0 \cdot 4 \cdot), 1 (\cdot, 3 \cdot), 1 (\cdot, 2 \cdot)) \quad \tilde{c}^{(e)} = ((0 \cdot 4 \cdot), 3 \cdot), 1 (\cdot, 2 \cdot))$$

$$\tilde{c}^{(e)} = ((0 \cdot 4 \cdot), 1 (\cdot, 3 \cdot), 1 (\cdot, 2 \cdot)) \quad \tilde{c}^{(e)} = ((0 \cdot 4 \cdot), 3 \cdot), 1 (\cdot, 2 \cdot))$$

c-sequence,
$$c \in \mathcal{C}_n \to c$$
-shape, $\tilde{c} \in \mathcal{C}_n \to f$ -sequence, $f \in \mathcal{F}_n$

The number of c-sequences corresponding to the given f is

$$|F^{-1}(f)| = 2^{1-n} n! (n-1)! P(f) = n! 2^{\neg (f)+1-n} \prod_{i=2}^{n} \ddot{f}_i$$

Let $\exists (\tilde{c})$ be the number of cherries of a *c*-shape $\tilde{c} \in \widetilde{C}$.

$$|\tilde{C}^{-1}(\tilde{c})| = 2^{1-n} n! (n-1)! P(\tilde{c}) = n! 2^{-\Im(\tilde{c})}$$
 (*Tajima*, 1983)

The number of c-shapes corresponding to the given f is

$$|\widetilde{C}(F^{-1}(f))| = 2^{-\beth(f)} \prod_{i=2}^{n} \ddot{f}_i$$

 $\beth(f) \equiv n - 1 - \urcorner(f) - \beth(f)$, the number of balanced splits that are not

cherries.

Hasse Diagram of the Poset making \mathcal{F}_n (n = 4, ..., 12)



Simulating *f*-sequences: for SFS, Shape Stats, ...

- 1: input:
 - 1. scaled mutation rate θ
 - 2. sample size n
- 2: **output:** a SFS sample x from the n-coalescent
- 3: generate an *f*-sequence under the unlabeled *n*-coalescent
- 4: draw $t \sim T = (T_2, T_3, ..., T_n)$, where T_i 's are

independently distributed as Exponential $\binom{i}{2}$

5:
$$l \leftarrow t^{\mathrm{T}} \cdot \mathbf{f}$$
 and $l_{\bullet} = \sum_{i=1}^{n-1}$

- 6: draw x from Poisson-Multinomial distribution $e^{-\theta l_{\bullet}}(\theta l_{\bullet})^{\sum_{i=1}^{n-1} x_i} \prod_{i=1}^{n-1} \overline{l}_i^{x_i} / \prod_{i=1}^{n-1} x_i!$
- **7**: return: *x*

Various tree shape statistics are further summaries of the *f***-sequence**

 \tilde{s} -sequence or Aldous shape statistic (Aldous, 2001) $\tilde{S}(f_n, f_{n-1}, \dots, f_1) = \tilde{s} \equiv (\tilde{s}_n, \tilde{s}_{n-1}, \dots, \tilde{s}_2) : \mathcal{F}_n \to \tilde{\mathcal{S}}_n$:

 $\tilde{s}_{i} \equiv (\tilde{s}_{i,1}, \tilde{s}_{i,2}) \equiv \left(\max\left(\|f\|_{i} \right), \min\left(\|f\|_{i} \right) 2^{-\mathbf{1}_{\{0\}}(\max\left(\|f\|_{i} \right) - \min\left(\|f\|_{i} \right))} \right), \\ \|f\|_{i} \equiv \{ j | f_{i,j} - f_{i-1,j} | \in \mathbb{N} : j \in \{1, 2, \dots, n\} \}.$

$$\mathfrak{Q}_n \equiv \{ Q_{\mathfrak{I}}(\tilde{s}) = q_{\mathfrak{I}} \equiv \sum_{i=n}^2 \tilde{s}_{i,1} \mathbf{1}_{\mathfrak{I}}(\tilde{s}_{i,1}) : \widetilde{\mathcal{S}}_n \to \mathcal{Q}_{\mathfrak{I}n}, \ \mathfrak{I} \in \mathbf{2}^{\{2,3,\dots,n\}} \setminus \emptyset \}$$

 $Q_{\{2,3,\ldots,n\}}(\tilde{s}) = q_{\{2,3,\ldots,n\}} = \sum_{i=n}^{2} \tilde{s}_{i,1}$ is the Sackin's index $Q_{\{2\}}/2 = q_{\{2\}}/2$ is the number of cherries $(n^2 - 3n + 2)^{-1} \sum_{i=n}^{2} (\tilde{s}_{i,1} - 2\tilde{s}_{i,d})$ is the Colless' index Note: There are $2^{n-1} - 3$ others in the family \mathfrak{Q}_n

Likelihood of a Site Frequency Spectrum

Proposition (Likelihood of SFS). Let $a \in A_n$ be a given coalescent tree, c be its c-sequence, f = F(c) be its f-sequence, $t \equiv (t_2, t_3, t_n) \in (0, \infty)^{n-1}$ be its epoch times and let

$$l \equiv (l_1, \dots, l_{n-1}) = t^{\mathrm{T}} f = \left(\sum_{i=2}^n t_i f_{i,1}, \dots, \sum_{i=2}^2 t_i f_{i,n-1} \right), \quad l_{\bullet} \equiv \sum_{i=2}^n l_i, \quad \bar{l}_i \equiv \frac{l_i}{l_{\bullet}}$$

be its lineage lengths subtending 1, 2, ..., n - 1 leaves, the total tree-size, and relative lineage lengths respectively.

Likelihood of a Site Frequency Spectrum

Proposition (Likelihood of SFS). Let $a \in A_n$ be a given coalescent tree, c be its c-sequence, f = F(c) be its f-sequence, $t \equiv (t_2, t_3, t_n) \in (0, \infty)^{n-1}$ be its epoch times and let

$$l \equiv (l_1, \dots, l_{n-1}) = t^{\mathsf{T}} f = \left(\sum_{i=2}^n t_i f_{i,1}, \dots, \sum_{i=2}^2 t_i f_{i,n-1} \right), \quad l_{\bullet} \equiv \sum_{i=2}^n l_i, \quad \bar{l}_i \equiv \frac{l_i}{l_{\bullet}}$$

be its lineage lengths subtending 1, 2, ..., n - 1 leaves, the total tree-size, and relative lineage lengths respectively.

$$P(x|\phi, a) = P(x|\phi, l = t^{\mathsf{T}}f) = e^{-\theta l_{\bullet}}(\theta l_{\bullet})^{S} \prod_{i=1}^{n-1} \bar{l}_{i}^{x_{i}} / \prod_{i=1}^{n-1} x_{i}!$$

Likelihood of a Site Frequency Spectrum

Proposition (Likelihood of SFS). Let $a \in A_n$ be a given coalescent tree, c be its c-sequence, f = F(c) be its f-sequence, $t \equiv (t_2, t_3, t_n) \in (0, \infty)^{n-1}$ be its epoch times and let

$$l \equiv (l_1, \dots, l_{n-1}) = t^{\mathrm{T}} f = \left(\sum_{i=2}^n t_i f_{i,1}, \dots, \sum_{i=2}^2 t_i f_{i,n-1} \right), \quad l_{\bullet} \equiv \sum_{i=2}^n l_i, \quad \bar{l}_i \equiv \frac{l_i}{l_{\bullet}}$$

be its lineage lengths subtending 1, 2, ..., n - 1 leaves, the total tree-size, and relative lineage lengths respectively.

$$\begin{split} P(x|\phi,a) &= P(x|\phi,l=t^{\mathrm{T}}f) = e^{-\theta l_{\bullet}}(\theta l_{\bullet})^{S} \prod_{i=1}^{n-1} \bar{l}_{i}^{x_{i}} / \prod_{i=1}^{n-1} x_{i}! \\ P(x|\phi) &= \frac{1}{\prod_{i=1}^{n-1} x_{i}!} \sum_{f \in F_{n}^{c}(x^{\circledast})} P(f) \left(\int_{t \in (0,\infty)^{n-1}} \left(e^{-\theta l_{\bullet}}(\theta l_{\bullet})^{S} \prod_{i=1}^{n-1} \bar{l}_{i}^{x_{i}} \right) P(t|\phi) \right) \\ & \text{where,} \quad F_{n}(x^{\circledast}) \equiv \bigcup_{\{h:x_{h}^{\circledast}=1\}} \{f \in \mathcal{F}_{n}: \sum_{i=1}^{n} f_{i,h} = 0\} \\ & X^{\circledast}(x) = x^{\circledast} \equiv (x_{1}^{\circledast}, \dots, x_{n-1}^{\circledast}) \equiv (\mathbf{1}_{\mathbb{N}}(x_{1}), \dots, \mathbf{1}_{\mathbb{N}}(x_{n-1})) \in \{0,1\}^{n-1} \end{split}$$

An Importance Sampler over $F_n^c(x^{\circledast})$

Proposition (A Proposal over $F_n^c(x^{\circledast})$). For a given $x \in \mathcal{X}_n^m$, consider the following discrete time Markov chain on the augmented state space $\mathbb{F}_n \times \{0,1\}^{n-1} \ni (f_h, z_h)$:

$$P^*((f_h, z_h)|(f_g, z_g)) = \begin{cases} P(f_h|f_g) / \Sigma(f_g, z_g) & : \text{ if } (f_h, z_h) \prec_{f, z} (f_g, z_g), \\ 0 & : \text{ otherwise} \end{cases}$$

where,

$$\Sigma(f_g, z_g) = \sum_{(j,k) \in H(f_g, z_g)} P(f_g - e_{j+k} + e_j + e_k | f_g),$$

 $\begin{aligned} H(f_g, z_g) &= \{(j, k) : f_{g, j+k} > 0, \ 1 \le j \le \max\{\min\{\hat{g}, j+k-1\}, \lceil \frac{j+k}{2} \rceil\} \le k \le j+k-1\}, \\ \hat{g} &= \max\{i : z_{g,i} = 1\}, \\ (f_h, z_h) \prec_{f, z} (f_g, z_g) \Leftrightarrow f_h = f_g + e_j + e_k - e_{j+k}, \ z_h = z_g - \mathbf{1}_{\{1\}}(z_{g,j}) e_j - \mathbf{1}_{\{1\}}(z_{g,k}) e_k \end{aligned}$

where, the initial state is $(f_1, X^{\circledast}(x)) = ((0, 0, ..., 1), x^{\circledast})$ and the final absorbing state is $(f_n, (0, 0, ..., 0)) = ((n, 0, ..., 0), (0, 0, ..., 0)).$

Maximum Aposteriori Estimates of θ and ν by \sum over $f \in \mathbb{F}_n^c(x^{\circledast})$

n	$\widehat{\nu}$			$\widehat{ heta}$			$(\widehat{ heta},\widehat{ u})$	
	$\sqrt{\overline{se}}$	bs	$C_{99\%}$	$\sqrt{\overline{se}}$	bs	$C_{99\%}$	$C_{99\%}$	$Qrt(\breve{\mathcal{K}})$
4	46	30	42	43	30	53	98	$\{0.061, 0.079, 0.13\}$
5	32	19	42	31	22	63	96	$\{0.074, 0.098, 0.16\}$
6	31	18	41	35	23	69	93	$\{0.082, 0.11, 0.17\}$
7	34	19	48	32	20	68	87	$\{0.090, 0.12, 0.21\}$
8	26	12	66	21	11	72	92	$\{0.098, 0.14, 0.26\}$
9	27	12	65	18	10	70	93	$\{0.097, 0.14, 0.21\}$
10	23	11	64	17	10	66	95	$\{0.091, 0.14, 0.30\}$

Topological Unfolding of SFS and Tajima's D when n = 4



Simulated Vs. Gen. Fisher's Exact Test with Tajima's D



Left panel: Distribution of p-values from the simulated test (left) and the generalized Fisher's exact test (right) for three values of $\theta = \{1, 10, 50\}$ per 1000 bp with n = 30. Right panel: The almost zero correlation of p-values between the two tests.

Summary



- Limits on Inference from Finest Empirical Resolutions
 - Inference from Coarser Site Frequency Spectrum is Possible via a Collapsed Kingman's *n*-coalescent Markov chain
- Algebraic Geometry is useful to infer from classical summaries of SFS.
- MSEs are smaller the exponential growth model
- Helps speed-up intensive SIS methods (Particle filtering on Experiment Graph)
- Topological unfolding of SFS and $D \Rightarrow$ Tree-less Genome Scans are essentially meaningless
- A Decision-theoretic formalism partially-ordered coalescent experiments graph
- Possible to generalize
 - Saves electricity and slows down global warming!



Research Fellow of the Royal Commission for the Exhibition of 1851.