



Edition 0.3

⚠ Obsolete

Author date: 2023-03-06

Citation:

E. Castedo Ellerman (2023) "Gametic lineage space" perm.pub
https://perm.pub/dsi:Fs4CjTor07Ssbb69Qcggp5yiBA0/0.3

Copyright:

creativecommons.org/licenses/by/4.0/
2023 © The Authors. This document is distributed under a Creative Commons Attribution 4.0 International license.

Gametic lineage space

E. Castedo Ellerman  (castedo@castedo.com)

Abstract

STAGE: WORKING DRAFT

DOCUMENT TYPE: Definition Document

OBJECTIVES

- Formal mathematical definition to be used in stochastic process model under development.
- Solicit feedback, in particular, on choice of wording for definitions.
- Precise mathematical definitions for technical discussions relating to ancestral recombination graphs

Gametic genealogy

A *gametic genealogy* is a convenient mathematical formalism of the genealogy of a population from the perspective of gametes. Mathematically, it is a quadruple

$$(\text{Gam}, \text{Mate}, \text{Par}, \text{Fert})$$

with components

- **Gam**, the set of underlying gametes,
- **Mate**, the set of zygotes formed by the fusion of egg gametes and sperm gametes,
- **Par**, a mapping from child gametes to parent zygotes, and
- **Fert**, a mapping from zygotes to fertilization time.

For convenience, given a *gametic genealogy*,

- Gam_0 denotes the set of egg gametes,
- Gam_1 denotes the set of sperm gametes, and
- Mate_* denotes the mapping from gametes to the zygotes they formed during fertilization.

Formally, a *gametic genealogy* must satisfy the following conditions.

1. $\text{Gam}_0 \cup \text{Gam}_1 = \text{Gam}$ and $\text{Gam}_0 \cap \text{Gam}_1 = \emptyset$.
2. $\text{Mate} \subset \text{Gam}_0 \times \text{Gam}_1$ and forms a one-to-one mapping between Gam_0 and Gam_1 .

3. Par is a function $C \mapsto \text{Mate}$, where C is a subset of Gam representing child gametes.

4. Fert is a function $\text{Mate} \mapsto \mathbb{R}$ such that for all child gametes $g \in \text{dom Par}$,

$$\text{Fert}(\text{Mate}_*(g)) > \text{Fert}(\text{Par}(g)) .$$

dom Par denotes the domain of Par , that is, the set of child gametes.

Gametic lineage space

A *gametic lineage space* is a mathematical formalism representing the lines of transmission of genetic information via gametes of a population over time. It is a triplet

$$(\text{Loc}, G, \text{Lin})$$

where

- Loc is the set of all genomic locations,
- G is a gametic genealogy $(\text{Gam}, \text{Mate}, \text{Par}, \text{Fert})$, and
- Lin is a function $\text{Loc} \times \text{Gam} \mapsto 2^{\text{Gam}}$ mapping a genomic position in a gamete to the set of gametes that transmitted genetic information to that position in that gamete.

For every location $\ell \in \text{Loc}$ and gamete $g \in \text{Gam}$, $\text{Lin}(\ell, g)$ is the lineage ending at gamete g via locus ℓ and it must satisfy the condition

$$\text{Lin}(\ell, g) = \{g\} \cup \text{Lin}(\ell, \text{Par}(g)_i) \text{ for either } i = 0 \text{ or } i = 1$$

when $g \in \text{dom Par}$, otherwise $\text{Lin}(\ell, g) = \{g\}$.

$\text{Par}(g)_0$ and $\text{Par}(g)_1$ are the maternal and paternal gametes, respectively, that fertilized the parent of g .

Stochastic gametic lineage space

A *stochastic gametic lineage space* is a gametic lineage space extended to model a random gametic lineage.

From a stochastic gametic lineage space, a time-indexed family of probability distributions $\{P_t\}_{t \in I}$ is induced.

TO DO: Need to rework space formalism to clarify over what is the σ -algebra: [issues #42](#).

For convenience we define the set of all lineages that contain a gamete $g \in \text{Gam}$ as

$$B(g) := \{\text{Lin}(\ell, s) : g \in \text{Lin}(\ell, s), \ell \in \text{Loc}, s \in \text{Gam}\} .$$

Formally, a *stochastic gametic lineage space* is a quintuple

$$(G, I, \{S_i\}_{i \in I}, \mathcal{F}, \mu)$$

where

- G is a gametic lineage space $(\text{Loc}, (\text{Gam}, \text{Mate}, \text{Par}, \text{Fert}), \text{Lin})$

- I is an index set of points in time with $\text{rng Fert} \subset I$,
- $\{S_t\}_{t \in I}$ is a time-indexed collection of sets of living zygotes,
- \mathcal{F} is a σ -algebra (*sigma-field*) over rng Lin , and
- μ is a measure on \mathcal{F}

which satisfy the following conditions

- $B(g) \in \mathcal{F}$ for all $g \in \text{Gam}$, and
- $\mu(S_t)$ is defined and finite for all $t \in I$.

Every gametic lineage space induces a time-indexed family $\{P_t\}_{t \in I}$ of probabilities spaces measurable on σ -algebra F . This defines the probability of lineages which end in a zygote alive at time t .

TO DO: Need to clarify relationship between F and Loc and Mate for when they are uncountable.

An embedded Ancestral Recombination Graph

An ancestral recombination graph [1] [2] [3] of a sampled population is embedded in a gametic lineage space. We formally show the exact embedding using the gARG formalism [4].

We start by defining the *genetic legacy* of a gamete $g \in \text{Gam}$ for sample population $S \subseteq \text{Gam}$ to be

$$\text{Leg}(g, S) := \{(\ell, d) \in \text{Loc} \times S : g \in \text{Lin}(\ell, d)\} .$$

This genetic legacy is the genetic material that survives in the sample population S originally copied from ancestral gamete g (with or without mutations).

Genetic legacy for a sample population S induces the following equivalence relationship over pairs of gametes g_1 and g_2 in Gam :

$$g_1 \simeq_S g_2 := \text{Leg}(g_1, S) = \text{Leg}(g_2, S) .$$

We denote the resulting equivalence class containing $g \in \text{Gam}$ as

$$[g]_S := \{g' : \text{Leg}(g', S) = \text{Leg}(g, S)\} .$$

In this equivalence relationship, gametes are considered equivalent if they have the same genetic legacy for the sample population S .

A convenient choice for an embedded gARG [4] is to set the gARG nodes (vertices) to be the equivalence classes:

$$\text{Nodes}(S) := \{[g]_S : g \in \text{Gam}\} .$$

The (unannotated) graph edges of the gARG are chosen as child-parent node pairs $(C, P) \in \text{Nodes}(S)^2$ where

$$\text{Par}(g)_i \in P \text{ for some } g \in C \text{ and some } i \in \{0, 1\} .$$

In the gARG, annotations are added for each graph edge (pair of child and parent nodes). This annotation is the set of locations through which genetic information has been copied from

parent to child. In the following interpretation, the only locations of interest are those for which genetic information has been transmitted into the sample population S . With this interpretation, the annotation for edge (C, P) is

$$\{\ell \in \text{Loc} : C \cup P \subseteq \text{Lin}(\ell, g) \text{ for some } g \in S\}.$$

Acknowledgements

Thanks to Daria Shipilina and Nick Barton for sharing their preprint [5] and discussing the conjecture in edition 0.1 of this document relating to their preprint.

Changes from edition 0.1

- add section about embedded ARG
- removed conjecture relating to [5]

References

1. Griffiths RC, Marjoram P. An Ancestral Recombination Graph. In: Friedman A, Miller W, Donnelly P, Tavaré S, editors. *Progress in Population Genetics and Human Evolution*. New York, NY: Springer New York; 1997. pp. 257–270. doi:[10.1007/978-1-4757-2609-1_16](https://doi.org/10.1007/978-1-4757-2609-1_16)
2. Hein J, Schierup MH, Wiuf C. *Gene genealogies, variation and evolution: A primer in coalescent theory*. Oxford ; New York: Oxford University Press; 2005.
3. Wakeley J. *Coalescent theory: An introduction*. Greenwood Village, Colo: Roberts & Co. Publishers; 2009.
4. Wong Y, Ignatieva A, Koskela J, Gorjanc G, Wohns AW, Kelleher J. A general and efficient representation of ancestral recombination graphs. <https://archive.softwareheritage.org/swh:1:rev:7df4f1995028cc676a6c1b231e8d7a024666b5fc>; 2022.
5. Shipilina D, Stankowski S, Pal A, Chan YF, Barton N. On the origin and structure of haplotype blocks. Preprints; 2022 Feb. doi:[10.22541/au.164425910.09070763/v1](https://doi.org/10.22541/au.164425910.09070763/v1)