

# Cartographie génétique

INRA, Thomas Schiex, Simon deGivry, Brigitte Mangin

Septembre 2016

# Plan

1

## Cartographier

- Quoi, pourquoi
- Comment ?

2

## Construction de cartes

- Estimer le taux de recombinaison
- Premières étapes pour la construction
- Grouper les marqueurs
- Ordonner les marqueurs

# Les cartes: s'orienter dans le génome

## Types

- **Cartes physiques** : distance réelle (Kb, Mb), à partir de fragments d'ADN. Résolution habituellement élevée.
- **Cartes d'hybrides irradiés** : Distance "statistique" liée à la cassure par irradiation, résolution intermédiaire.
- **Cartes génétiques** : s'appuie sur la recombinaison durant la méiose. Distance "statistique".

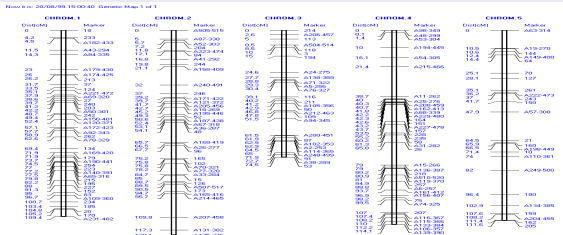
## Carte génétique/hybrides irradiés

Représentation d'un génome positionnant un ensemble de repères (marqueurs) dont on connaît les positions sur des groupes de liaison (chromosomes idéalement).

Quoi, pourquoi

# Exemple

## Carte génétique



Groupes de liaison génétique

## Génome



Chromosomes

# Pourquoi

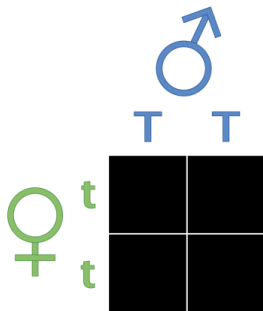
- Identifier les régions du génome influençant un caractère d'intérêt (maladie ou caractère quantitatif plus complexe)
- Positionner et identifier un gène (clonage positionnel)
- Comparer les génomes (étude de la synténie, évolution, transfert d'information)
- Faciliter la construction de cartes physiques, assemblage
- Étudier la méiose

# Les bases: lois de Mendel (modernes, diploïdes)

Cross: **TT** x **tt**

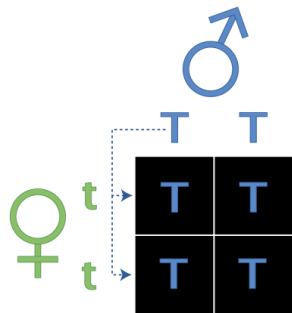
## Loi de ségrégation

Un génome contient un ensemble de paires de gènes. Les paires ségrègent (se séparent) dans les gamètes, la moitié des gamètes portant un gène, l'autre moitié portant l'autre gène. Taille de plante (allèles *Tt*).



# Les bases: lois de Mendel (modernes, diploïdes)

Cross: **TT** x **tt**



**Loi de ségrégation**

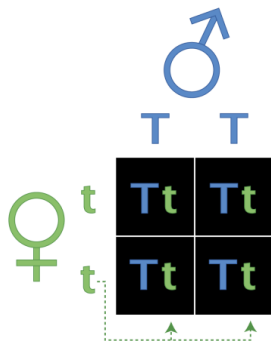
Un génome contient un ensemble de paires de gènes. Les paires ségrègent (se séparent) dans les gamètes, la moitié des gamètes portant un gène, l'autre moitié portant l'autre gène. Taille de plante (allèles *Tt*).

# Les bases: lois de Mendel (modernes, diploïdes)

## Loi de ségrégation

Un génome contient un ensemble de paires de gènes. Les paires ségrègent (se séparent) dans les gamètes, la moitié des gamètes portant un gène, l'autre moitié portant l'autre gène. Taille de plante (allèles  $Tt$ ).

Cross:  $TT$  x  $tt$



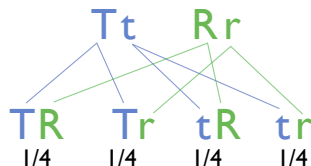


Quoi, pourquoi

# Les bases: lois de Mendel (modernisées, diploïdes)

## Loi de ségrégation indépendante

L'assortiment de plusieurs gènes dans une cellule sexuelle se fait de façon indépendante entre les différents gènes. Taille  $Tt$  et forme  $Rr$  (ridé).

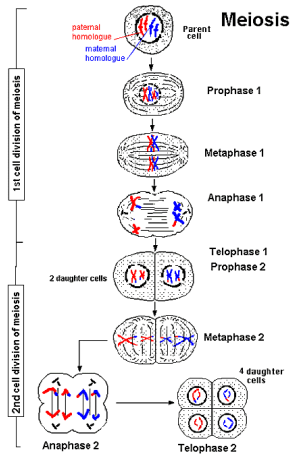


# Le principe historique de la cartographie

## Liaison génétique (Bateson 1905)

Pour certaines paires de gènes, la fréquence des combinaisons parentales dans les gamètes est supérieure à ce que l'on attend. On parle de **liaison génétique**.

Expliqué par Morgan (1911) par l'appartenance à un même chromosome et un éventuel chiasma durant la méiose (crossing-over).





Comment ?

# Bases

- **Loci, gènes, marqueurs** :  $\mathcal{A}, \mathcal{B}$ . Emplacement sur un chromosome

Comment ?

# Bases

- **Loci, gènes, marqueurs** :  $\mathcal{A}, \mathcal{B}$ . Emplacement sur un chromosome
- **Polymorphisme** : présente au moins deux formes différentes (allèles  $Aa$ ).

Comment ?

# Bases

- **Loci, gènes, marqueurs** :  $\mathcal{A}, \mathcal{B}$ . Emplacement sur un chromosome
- **Polymorphisme** : présente au moins deux formes différentes (allèles  $Aa$ ).
- **Homozygote** : paire d'allèles identiques ( $\frac{A}{A}$  ou  $\frac{a}{a}$ ), sinon hétérozygote ( $\frac{A}{a}$ ).

Comment ?

# Bases

- **Loci, gènes, marqueurs** :  $\mathcal{A}, \mathcal{B}$ . Emplacement sur un chromosome
- **Polymorphisme** : présente au moins deux formes différentes (allèles  $Aa$ ).
- **Homozygote** : paire d'allèles identiques ( $\frac{A}{A}$  ou  $\frac{a}{a}$ ), sinon hétérozygote ( $\frac{A}{a}$ ).
- **Haplotype** : séquence des allèles portés par chacun des chromosomes ( $\frac{Ab}{aB}$  par exemple).

Comment ?

# Bases

- **Loci, gènes, marqueurs** :  $\mathcal{A}, \mathcal{B}$ . Emplacement sur un chromosome
- **Polymorphisme** : présente au moins deux formes différentes (allèles  $Aa$ ).
- **Homozygote** : paire d'allèles identiques ( $\frac{A}{A}$  ou  $\frac{a}{a}$ ), sinon hétérozygote ( $\frac{A}{a}$ ).
- **Haplotype** : séquence des allèles portés par chacun des chromosomes ( $\frac{Ab}{aB}$  par exemple).
- **Génotype** : séquence des paires d'allèles (non ordonnées) portés par les chromosomes homologues ( $\frac{A}{a} \frac{B}{b}$  par exemple).



Comment ?

# Bases

- **Loci, gènes, marqueurs** :  $\mathcal{A}, \mathcal{B}$ . Emplacement sur un chromosome
- **Polymorphisme** : présente au moins deux formes différentes (allèles  $Aa$ ).
- **Homozygote** : paire d'allèles identiques ( $\frac{A}{A}$  ou  $\frac{a}{a}$ ), sinon hétérozygote ( $\frac{A}{a}$ ).
- **Haplotype** : séquence des allèles portés par chacun des chromosomes ( $\frac{Ab}{aB}$  par exemple).
- **Génotype** : séquence des paires d'allèles (non ordonnées) portés par les chromosomes homologues ( $\frac{A}{a} \frac{B}{b}$  par exemple).
- **Phase**: information suffisante pour déterminer les deux haplotypes à partir du génotype.

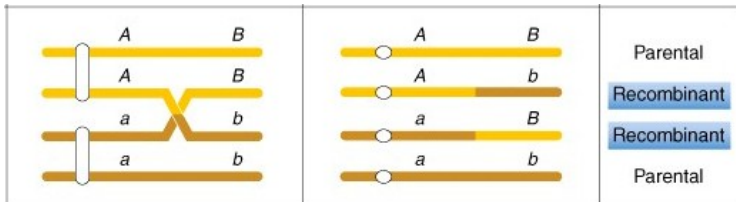
Comment ?

# Recombinants et Non recombinants

## Marqueurs $A, B$

- une cellule diploïde portant les haplotypes  $AB/ab$ ,
- on peut avoir les gamètes porteuses des haplotypes  $AB, ab, Ab, aB$

Les deux premiers sont *parentaux* ou *non recombinants*. Les deux autres *recombinants* (nombre impair de cross-overs).



Comment ?

# Taux de recombinaison

Taux de recombinaison  $\leq \frac{1}{2}$

Le taux de recombinaison  $\rho_{AB}$  entre les deux marqueurs  $A$  et  $B$  est la proportion de *recombinants*.

## Exemple

Entre 3 gènes  $\mathcal{Y}$  (yellow),  $\mathcal{W}$  (white),  $\mathcal{M}$  (miniature) de la drosophile, on observe  $\rho_{\mathcal{Y},\mathcal{W}} = 1.3\%$ ,  $\rho_{\mathcal{W},\mathcal{M}} = 32.6\%$  et  $\rho_{\mathcal{Y},\mathcal{M}} = 33.8\%$ . On peut penser que les marqueurs sont dans l'ordre  $\mathcal{Y} - \mathcal{W} - \mathcal{M}$

Du fait des doubles crossing-overs, pour un ordre  $\mathcal{Y} - \mathcal{W} - \mathcal{M}$  :

$$\rho_{\mathcal{Y},\mathcal{M}} < \rho_{\mathcal{Y},\mathcal{W}} + \rho_{\mathcal{W},\mathcal{M}} \quad (\text{non additif})$$

# Distance génétique

## Définition

La distance génétique  $d_{AB}$  entre deux marqueurs  $\mathcal{A}$  et  $\mathcal{B}$  est le nombre moyen de *crossing-overs* entre les deux marqueurs par méiose.

## Propriétés

- Additif
- 1cM (centiMorgan) correspond à un crossover sur un haplotype pour 100 méioses.
- Les cross-overs ne sont pas facilement observables.



## Fonction de distance - map functions

Entre deux marqueurs. La première fonction de distance s'appuie sur un modèle de recombinaison simplifié (sans interférence, deux chromatides).

### Fonction de Haldane - sans interférence (1919)

$$\rho = \frac{1}{2}(1 - e^{-2d}) \quad d = -\frac{1}{2} \log(1 - 2\rho)$$

Beaucoup d'autres fonctions pour l'interférence:

### Fonction de Kosambi - interférence (1944)

$$\rho = \frac{1}{2} \tanh(2d) \quad d = \frac{1}{2} \tanh^{-1}(2\rho)$$

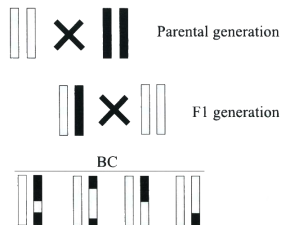
Pour de faibles distances/taux de recombinaison,  $d \approx \rho$ .



## Back-cross: estimer $\rho$ entre 2 marqueurs

### Un individu peut avoir

- deux marqueurs homozygotes ( $AA, BB$ ) ou hétérozygotes ( $Aa, Bb$ ) : **non recombinant** ( $NR$ ).
- un hétérozygote, un homozygote ( $Aa, BB$  ou  $AA, Bb$ ) : **recombinant** ( $R$ ).



### Vraisemblance - probabilité des données

Si  $\rho$  est le taux de recombinaison (à estimer), la probabilité d'observer les données de typage  $Data$  (la vraisemblance) est :

$$\text{Prob}(Data|\rho) = \rho^R(1 - \rho)^{NR}$$



# Back-cross : estimer $\rho$ entre 2 marqueurs

## Vraisemblance - probabilité des données

$$\text{Prob}(\text{Data}|\rho) = \rho^R(1 - \rho)^{NR}$$

## Maximum de vraisemblance

La valeur estimée  $\hat{\rho}$  de  $\rho$  choisie est celle qui maximise la probabilité d'observer les données (estimateur convergent).

Par un passage au logarithme et une étude de la dérivée on obtient :

$$\hat{\rho} = \frac{R}{R + NR}$$

## En pratique

- Individus non typés sur un marqueur. Données manquantes.
- On n'observe pas toujours les génotypes. Si un allèle  $A$  est "dominant",  $A$  est compatible avec  $AA, Aa$  en back-cross.
- Erreurs de typages

La vraisemblance de données incomplètes est compliquée. De même que celle lorsque les marqueurs ne sont pas codominants, ou lorsque le pedigree est plus complexe que le back-cross.

Pour la maximiser on utilise des algorithmes d'optimisation dédiés (par exemple EM - Expectation Maximisation).

Dempster et al., JRSS, 1977

# Nettoyage des données : distorsion

## Marqueur distordu

Allèle sur-représenté dans la descendance / à la fréquence attendue (gène lié à la reproduction/croissance, réarrangements ou problème d'échantillonnage)

## Test de $\chi^2$ de Pearson $T_{\chi^2}$

Sous l'hypothèse nulle: { les données observées sont tirées de la distribution théorique attendue }.

Pour un risque  $\alpha = 0.05$ ,  $\chi^2_{1ddl} = 3.84$

si  $T_{\chi^2} > 3.84$  on rejette l'hypothèse de non distorsion.

# Nettoyage des données : “marqueurs confondus”

## Jeux de données modernes

- typage de plusieurs dizaines de milliers de marqueurs
- distance minimale inter-marqueurs très faible
- pas de recombinaison/cassure observée : même génotypes (ou génotypes compatibles avec les données manquantes).

⇒ Supprimer ou fusionner des marqueurs



Ordonner les marqueurs

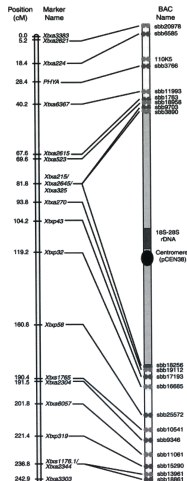
# Construction de la carte

## Cartographeur

Pour chaque groupe de liaison (partie de chromosome), déterminer l'ordre des marqueurs et les distances (taux de recombinaisons) qui séparent deux marqueurs adjacents

## Carte saturée

autant de groupes que de chromosomes, tous les marqueurs de la carte sont liés à un groupe.



# Trouver une bonne carte

## Un problème combinatoire

Pour  $n$  marqueurs, il y a  $n!/2$  ordres de marqueurs définissant des cartes différentes.  $\frac{10!}{2} = 1.810^6$ .

- Impossible d'énumérer les ordres.
- Problème d'optimisation difficile (même dans ses versions les plus simples).

# Logiciels de cartographie

- végétaux: MapMaker, CarthaGene, JoinMap
- animaux: CRIMAP,
- homme: MapMaker
- hybrides irradiés: RHMAP, RHO, CarthaGene

Voir <http://linkage.rockefeller.edu/soft/>