

Imputación de variantes genéticas no genotipadas

Carla Lluís-Ganella, MSc PhD

Grupo de Epidemiología y Genética Cardiovascular
Parque de investigación biomédica de Barcelona (PRBB)
clluis@imim.es | carla.lluis@gmail.com

Qué es y para qué la necesitamos?

En qué se basa?

El proceso de la meiosis

Recombinación

“Recombination hot spots”

Desequilibrio de ligamiento

Imputación de variantes genéticas

Software para la imputación

Métodos estadísticos para la imputación

Una visión histórica

Qué es?

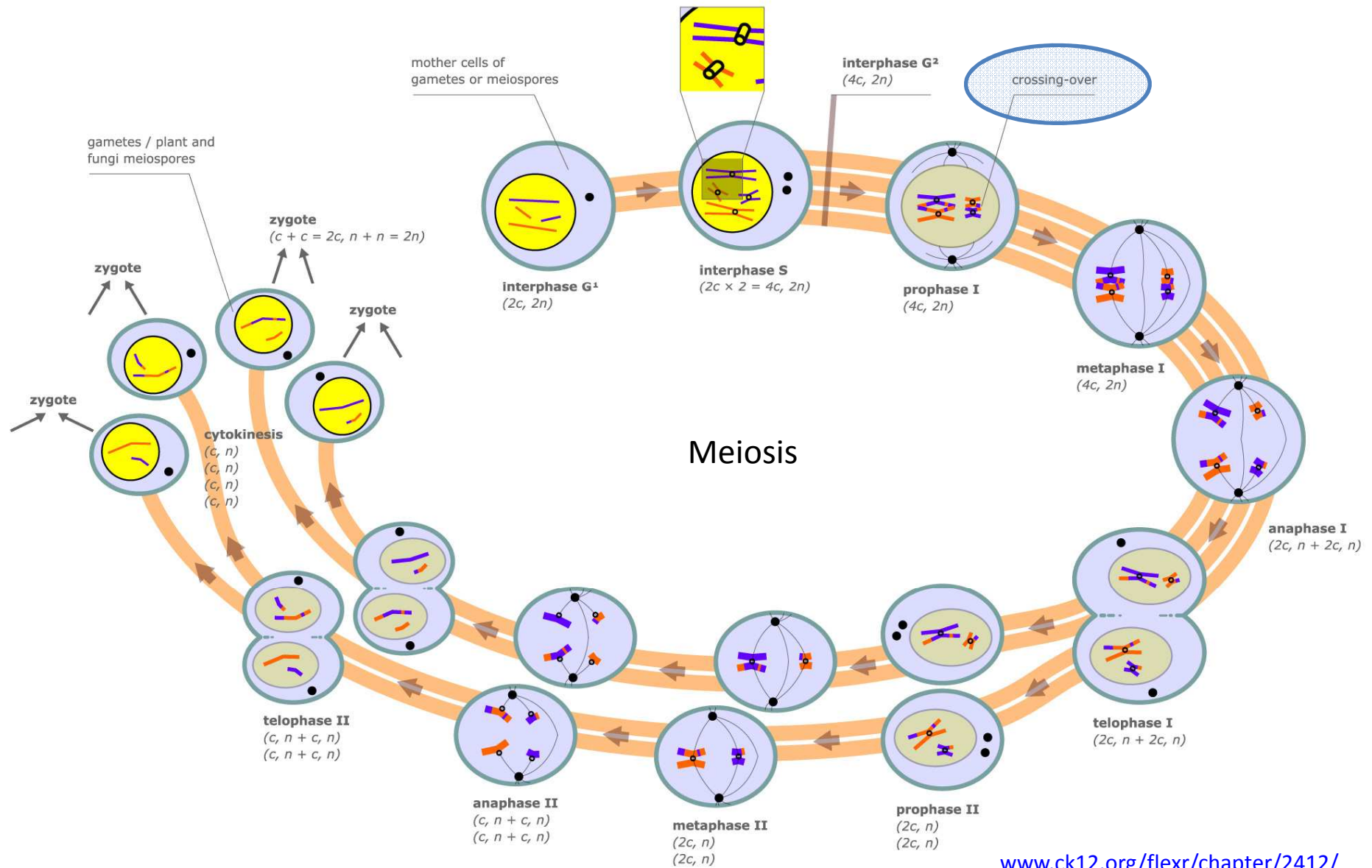
La imputación es un método estadístico que sirve para lidiar con datos no disponibles a partir de la información de otros datos si disponibles.

Por qué la necesitamos en genética?

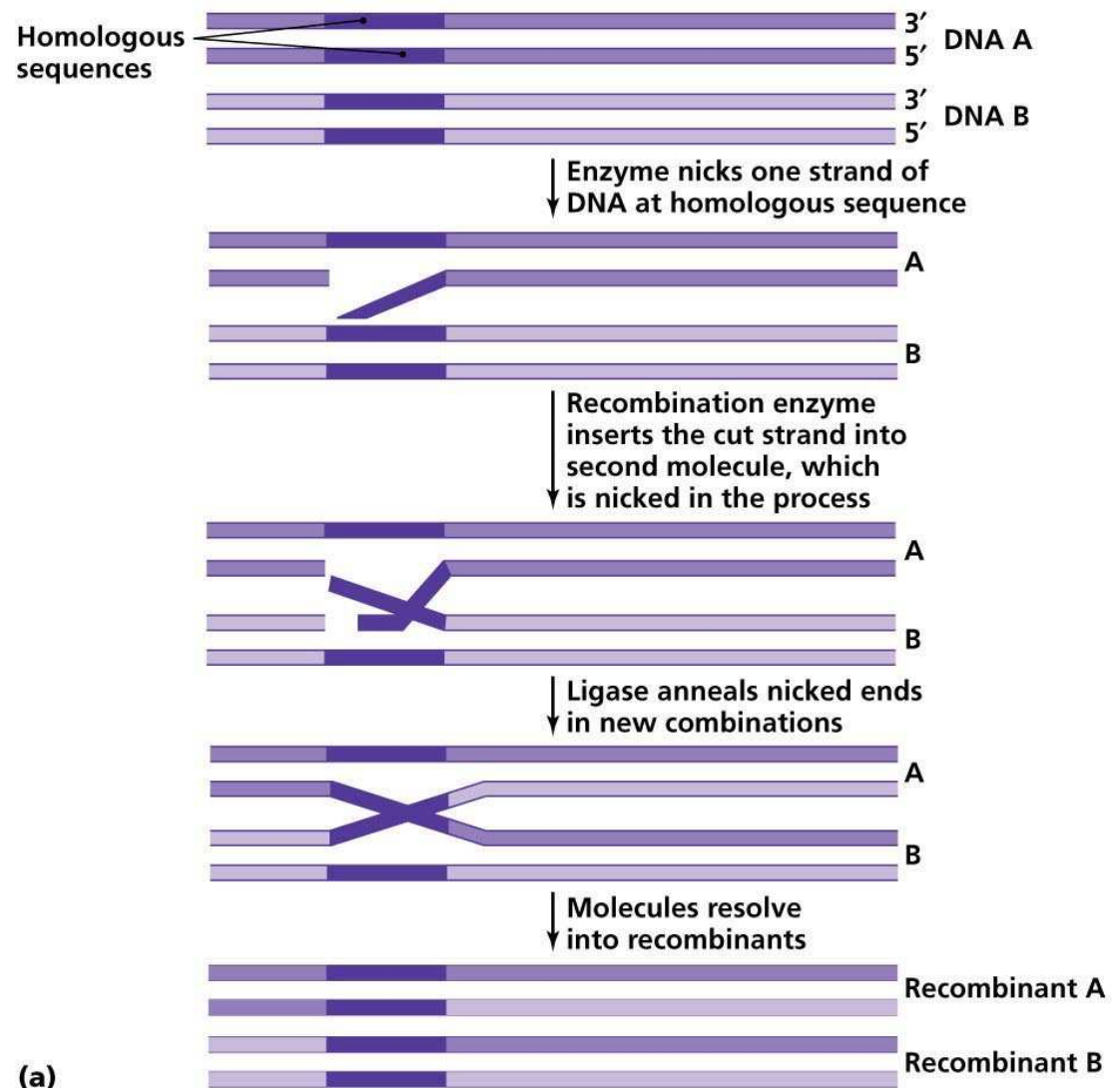
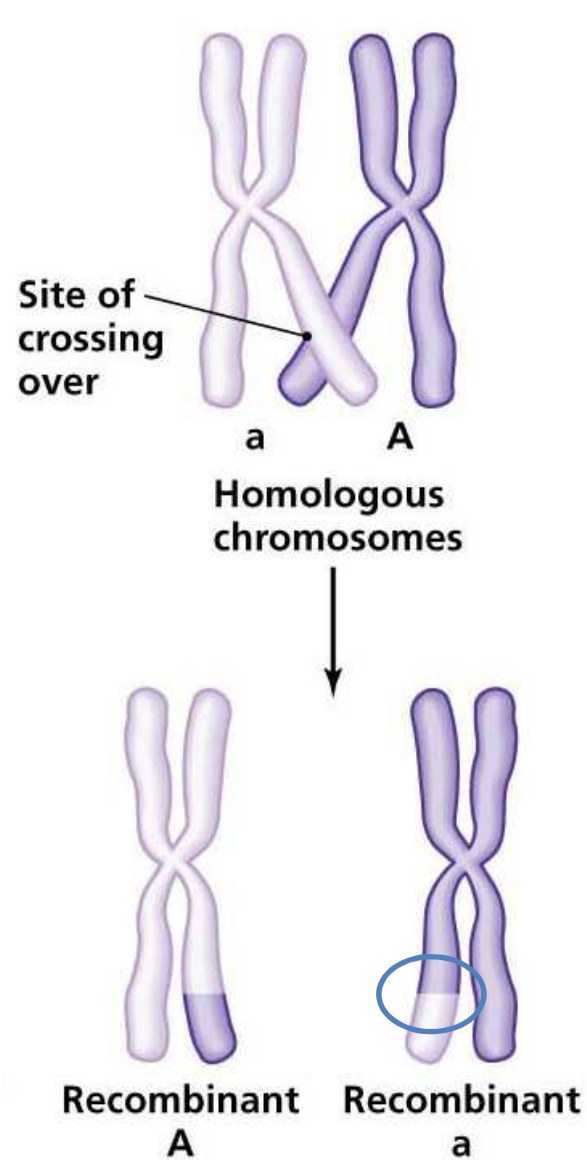
Genotipar un gran número de variantes no es habitualmente posible debido a cuestiones económicas.

Permite recuperar algunas variantes genéticas con genotipos inciertos en algunos individuos, permitiendo un aumento del número de variantes y muestras disponibles.

La imputación aprovecha un proceso biológico llamado “recombinación cromosómica” para calcular la probabilidad de que un individuo tenga un genotipo determinado.

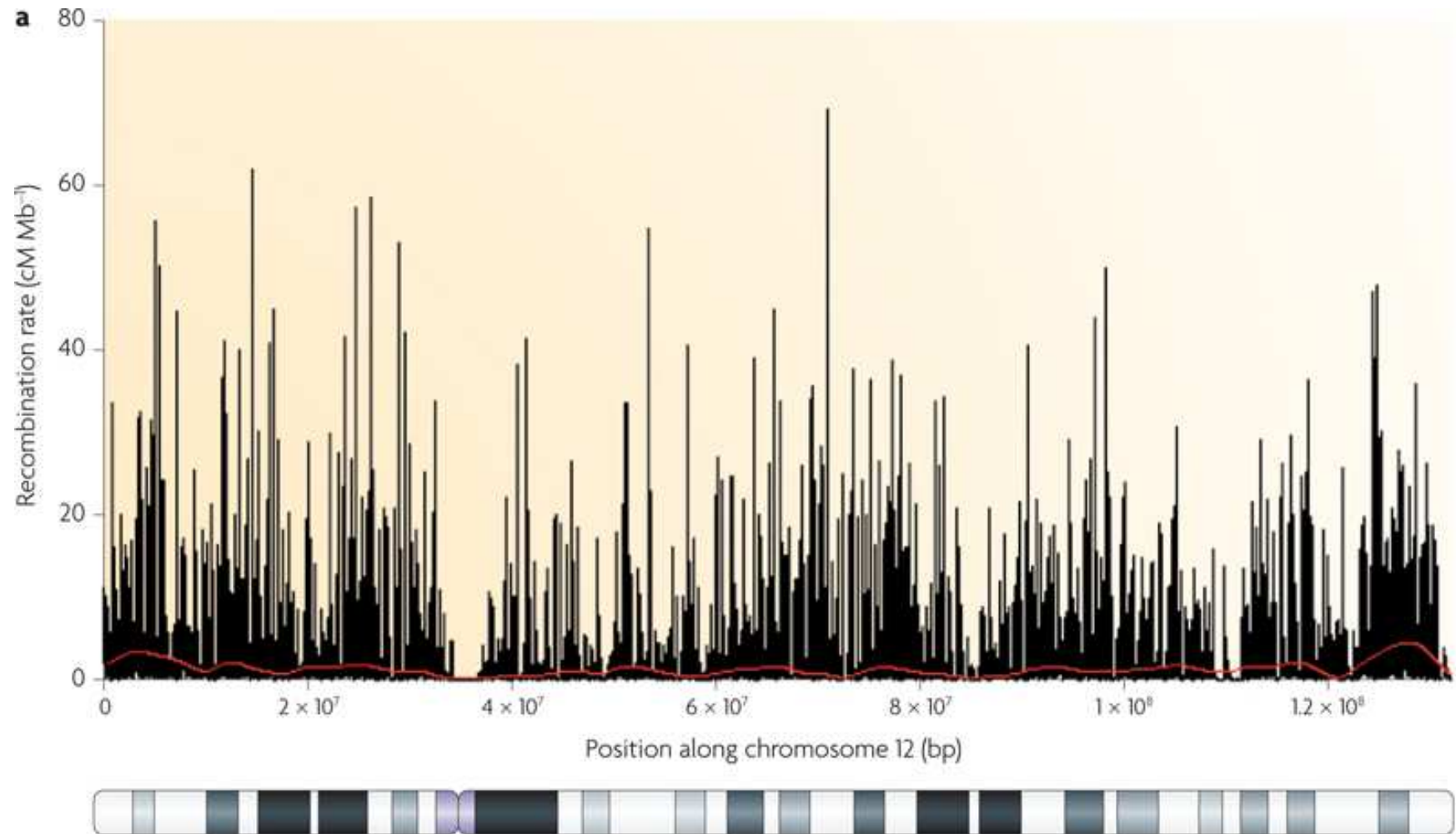


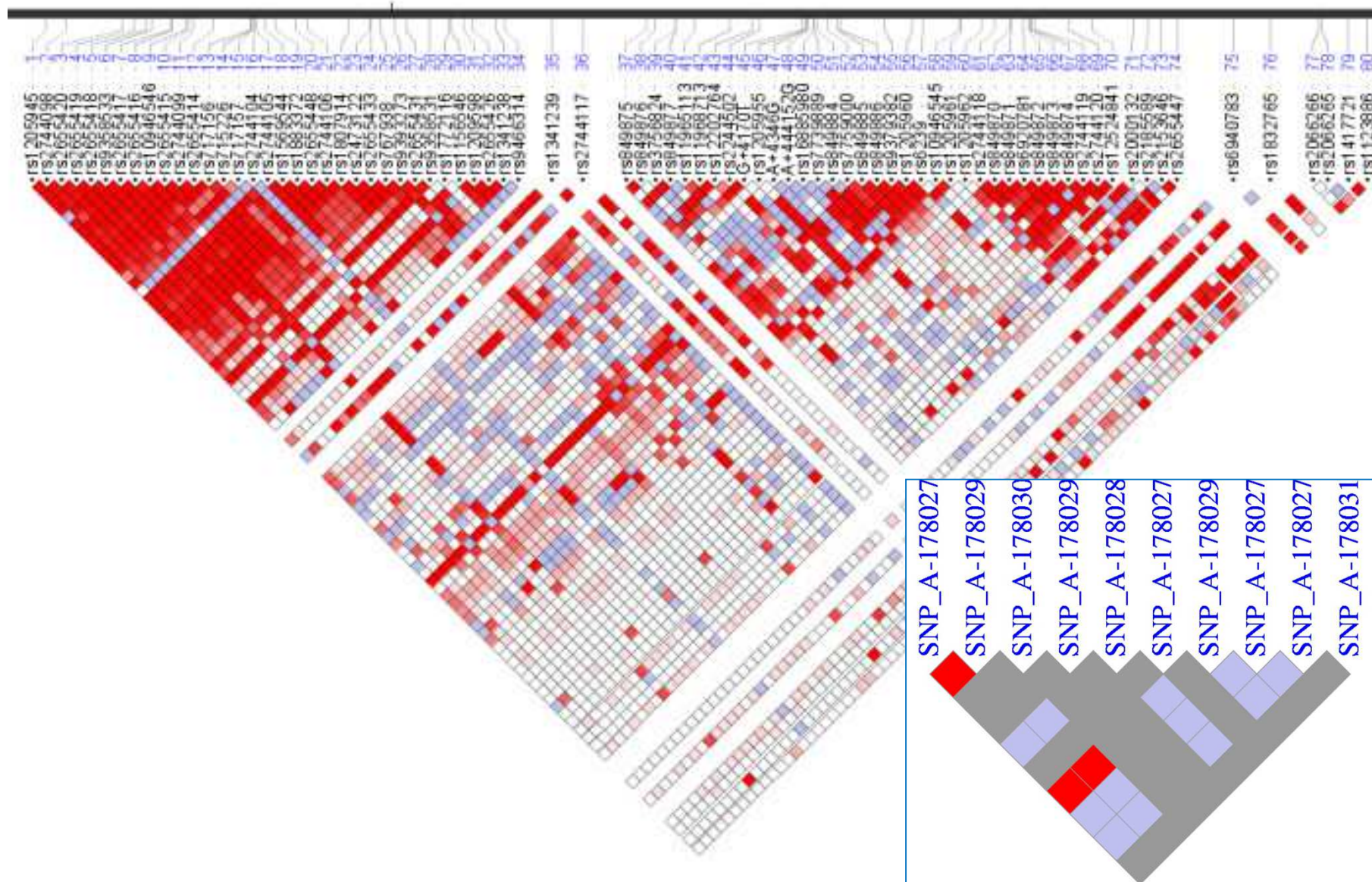
www.ck12.org/flexr/chapter/2412/

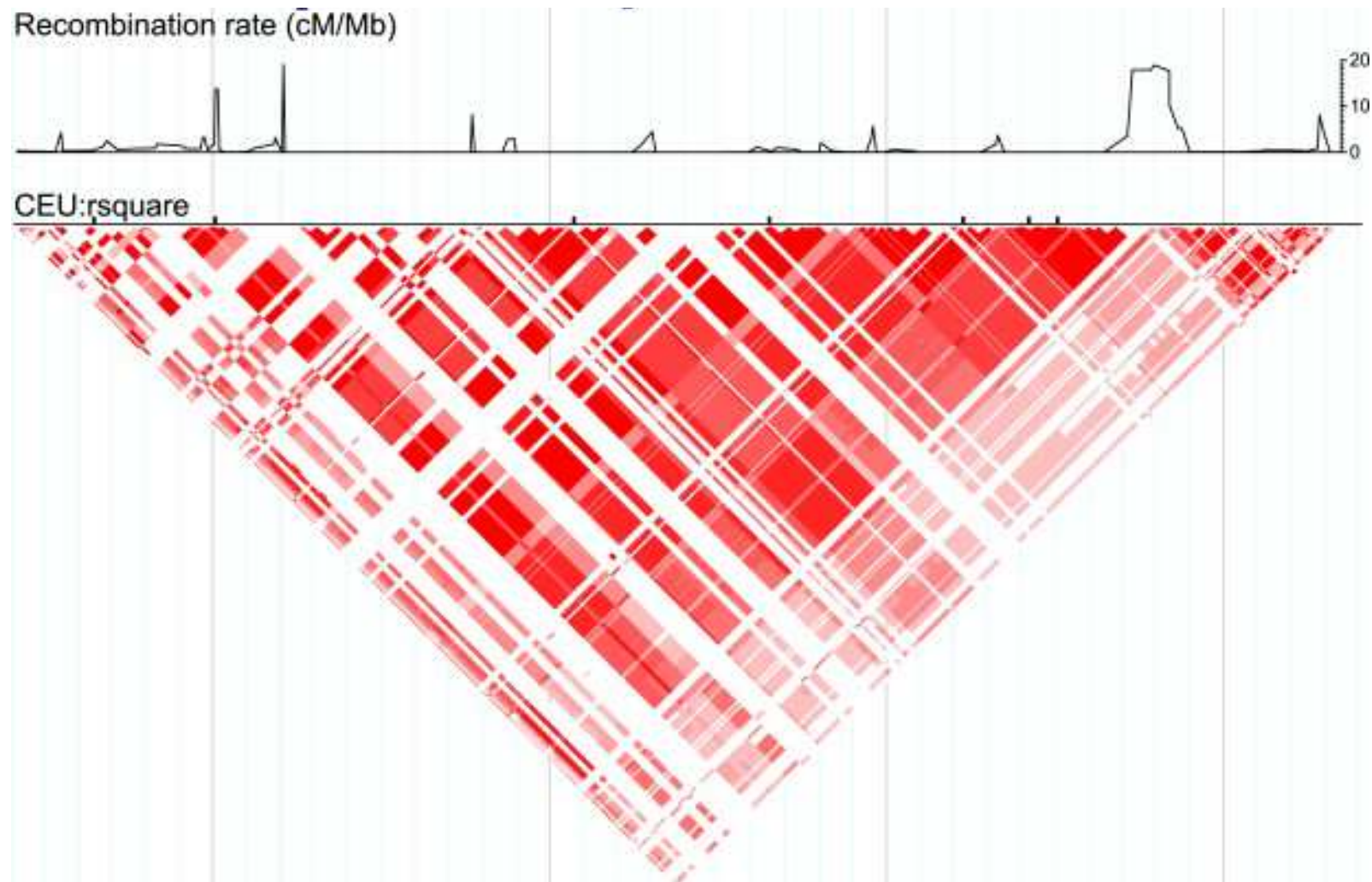


Copyright © 2006 Pearson Education, Inc., publishing as Benjamin Cummings

<http://academic.pgcc.edu>







IMPUTE

Marchini *et al.* Nat Genetics 2007

MACH

Li *et al.* Genetic Epidemiology 2010

BEAGLE

Browning and Browning. Hum Genet 2008

PLINK

Purcell *et al.* Am J Hum Genet 2007

at may show
³, microsatel-
 different study
 l mapping by
 of statistical

requires care.
 ensity as well
 used and the
 ompare head-

at will not be
 ie population
 y features of
 or context is
 differences in
 sample used
 ond, whether
 ure within a
 the analysis
 to impute

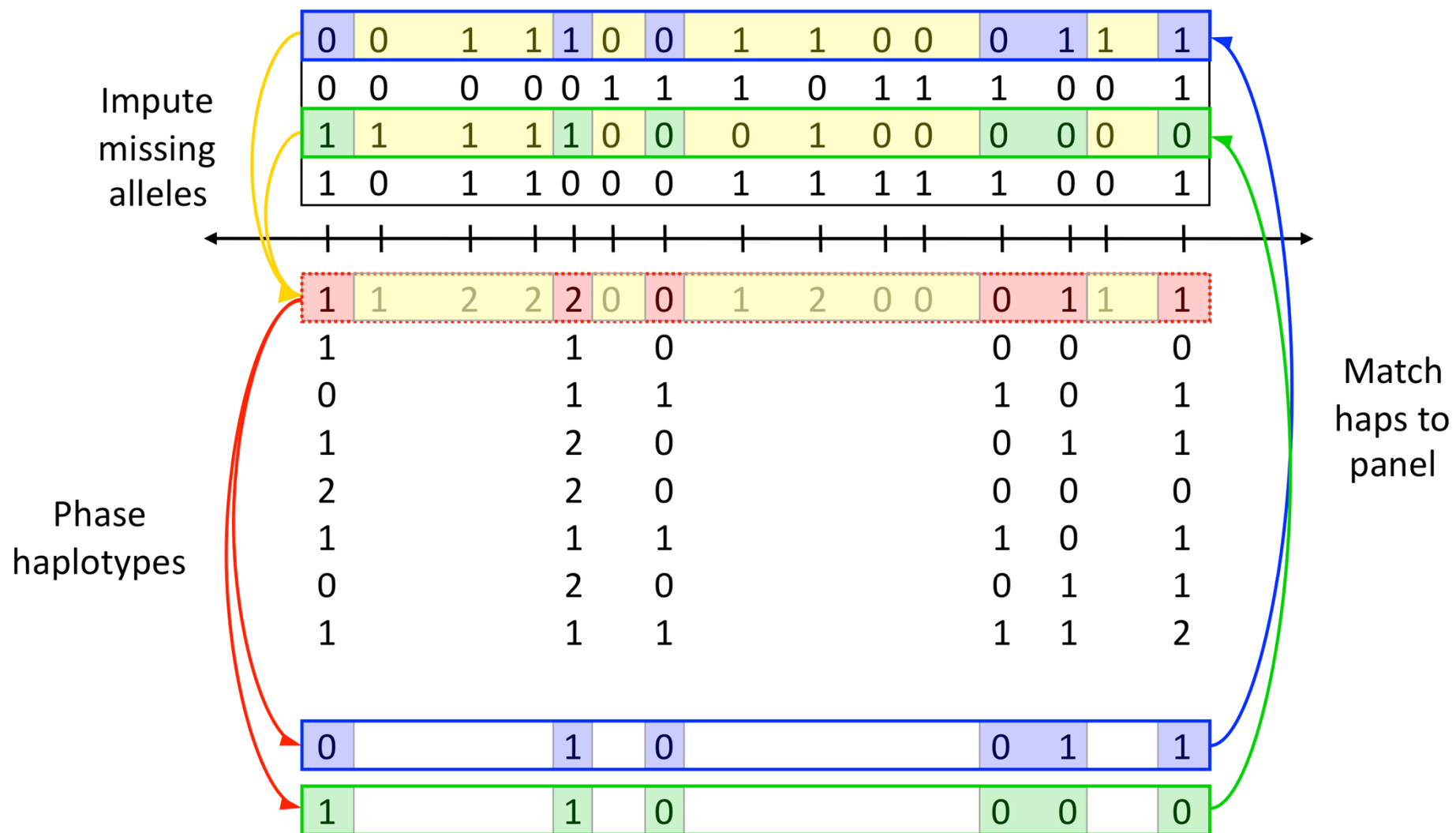
gent $G = \{G_O, G_M\}$. To impute the missing genotypes, we require the joint distribution of observed and missing genotype data, and we make the modeling assumption that each individual's genotype vector can be considered independently of the others. That is,

$$\Pr(G_M|G_O, H) \propto \Pr(G_M, G_O|H) = \Pr(G|H) = \prod_{i=1}^K \Pr(G_i|H)$$

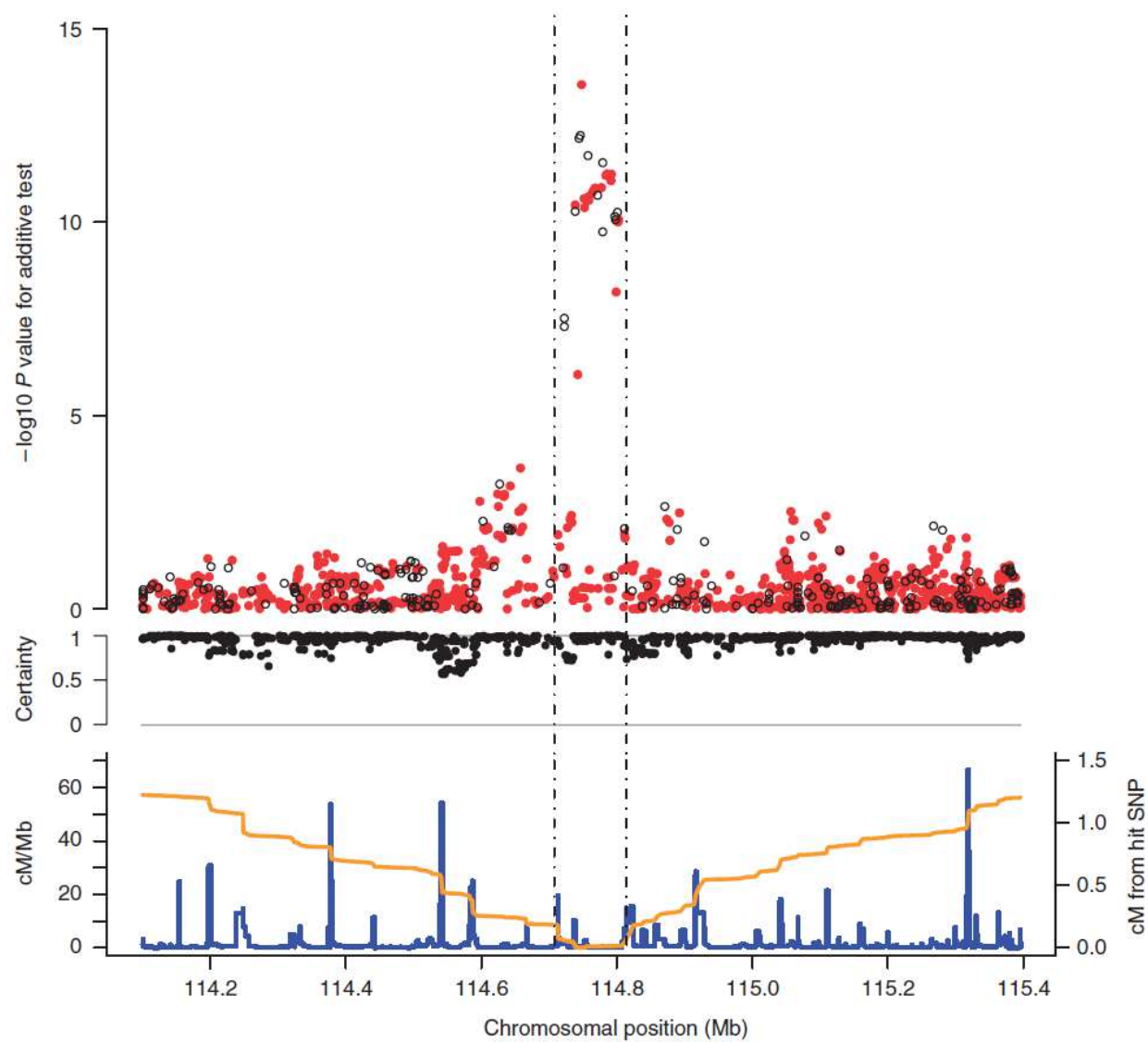
Our model for each individual's genotype vector, $\Pr(G_i|H)$, is a Hidden Markov Model in which the hidden states are a sequence of pairs of the N known haplotypes in the set H . That is,

$$\Pr(G_i|H) = \sum_{Z_i^{(1)}, Z_i^{(2)}} \Pr(G_i|Z_i^{(1)}, Z_i^{(2)}, H) \Pr(Z_i^{(1)}, Z_i^{(2)}|H)$$

where $Z_i^{(1)} = \{Z_{i1}^{(1)}, \dots, Z_{iL}^{(1)}\}$ and $Z_i^{(2)} = \{Z_{i1}^{(2)}, \dots, Z_{iL}^{(2)}\}$ are the two sequences of hidden states at the L sites and $Z_{il}^{(j)} \in \{1, \dots, N\}$. These hidden states can be thought of as the pair of haplotypes in the set H that are being copied to form the genotype vector G_i . The term $\Pr(Z_i^{(1)}, Z_i^{(2)}|H)$ defines our prior probability on how sequences of hidden states change along the sequence, and $\Pr(G_i|Z_i^{(1)}, Z_i^{(2)}, H)$ models how the observed genotypes will be close to but not exactly the same as the haplotypes being copied. This model extends a

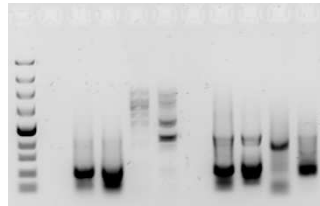


Extracted from Bryan Howie website

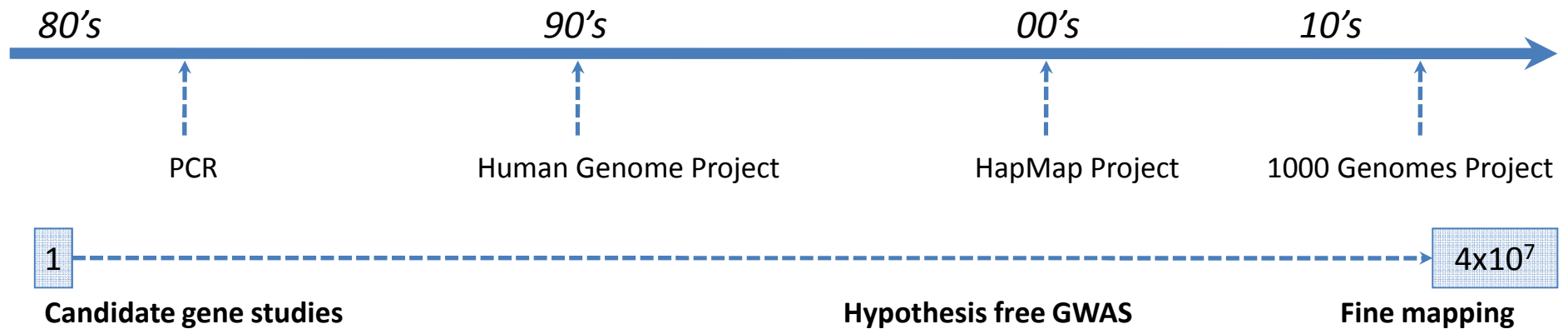
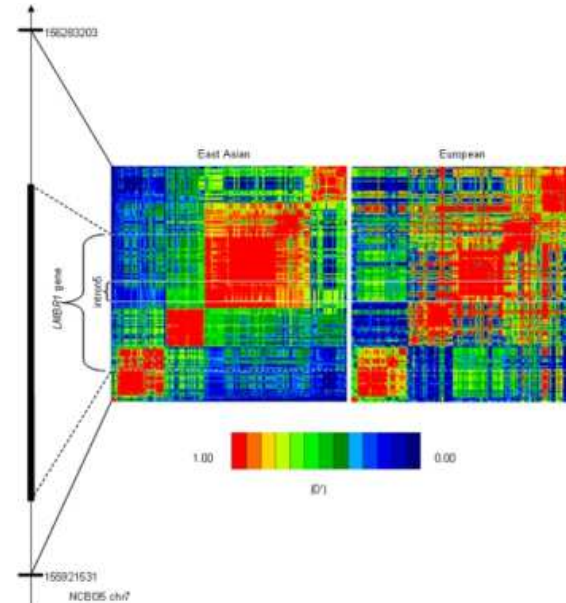


Marchini et al. Nat Genetics 2007

Una visión histórica



EcoRI
XbaI
PvuII
BamII
...



Si queréis más información:



clluis@imim.es | carla.lluis@gmail.com



www.linkedin.com/in/carlalluis



Lluis-Ganella C[Author]

Fi de la presentació