

# Réseau bayésien et génétique

Simon de Givry

INRA-UBIA Toulouse

21 juin 2011

# PLAN

## Réseau Bayésien

- définition
- réseaux de génotypes, d'allèles, de ségrégations
- séparation dirigée et indépendance conditionnelle
- principales requêtes
- méthode exacte pour l'inférence probabiliste
- méthode exacte pour l'optimisation combinatoire
- quelques résultats en reconstruction d'haplotypes à partir d'un pedigree et de génotypes

# Supports de cours

Transparents des cours, articles et tutoriaux disponibles à

<http://mulcyber.toulouse.inra.fr/scm/viewvc.php/MAB/?root=mposc>

- Patrick Naim, Pierre-Henri Wuillemin, Philippe Leray, Olivier Pourret et Anna Becker, *Réseaux bayésiens*, 2008  
BNT Matlab (Kevin Murphy) <http://bnt.insa-rouen.fr/>
- Adnan Darwiche, *Modeling and Reasoning with Bayesian Networks*, 2009  
Samlam <http://reasoning.cs.ucla.edu/samiam>
- Christopher M. Bishop, *Pattern Recognition and Machine Learning* (chapitre 8), 2006  
<http://research.microsoft.com/~cmbishop/PRML/index.htm>
- Finn V. Jensen, *Bayesian Networks and Decision Graphs*, 2001
- Association for Uncertainty in Artificial Intelligence (AUAI)  
<http://auai.org/>
- Logiciels payants : (usa) Hugin, Netica, (fr) Bayesia, ProBayes,...

# Autres liens sur internet

- Rina Dechter, *courses on Belief Networks (slides and homework)*, 2009 (with focus on Genetic Linkage Analysis, 2005) <http://www.ics.uci.edu/~dechter>
- Dan Geiger, *Advanced Topics in Bioinformatics (genetics)*, 2006 [http://www.cs.technion.ac.il/~anna\\_bi/cs236633/](http://www.cs.technion.ac.il/~anna_bi/cs236633/)
- Kevin Murphy, *A Brief Introduction to Graphical Models and Bayesian Networks*, 1998  
<http://people.cs.ubc.ca/~murphyk/Bayes/bayes.html>
- Bruno Garcia, *notes de cours sur Recherche Opérationnelle*, 2000  
<http://www.bruno-garcia.net/www/polyro/polyro.html>
- Daphne Koller, Nir Friedman, *Probabilistic Graphical Models*, 2009

# Définition d'un réseau bayésien

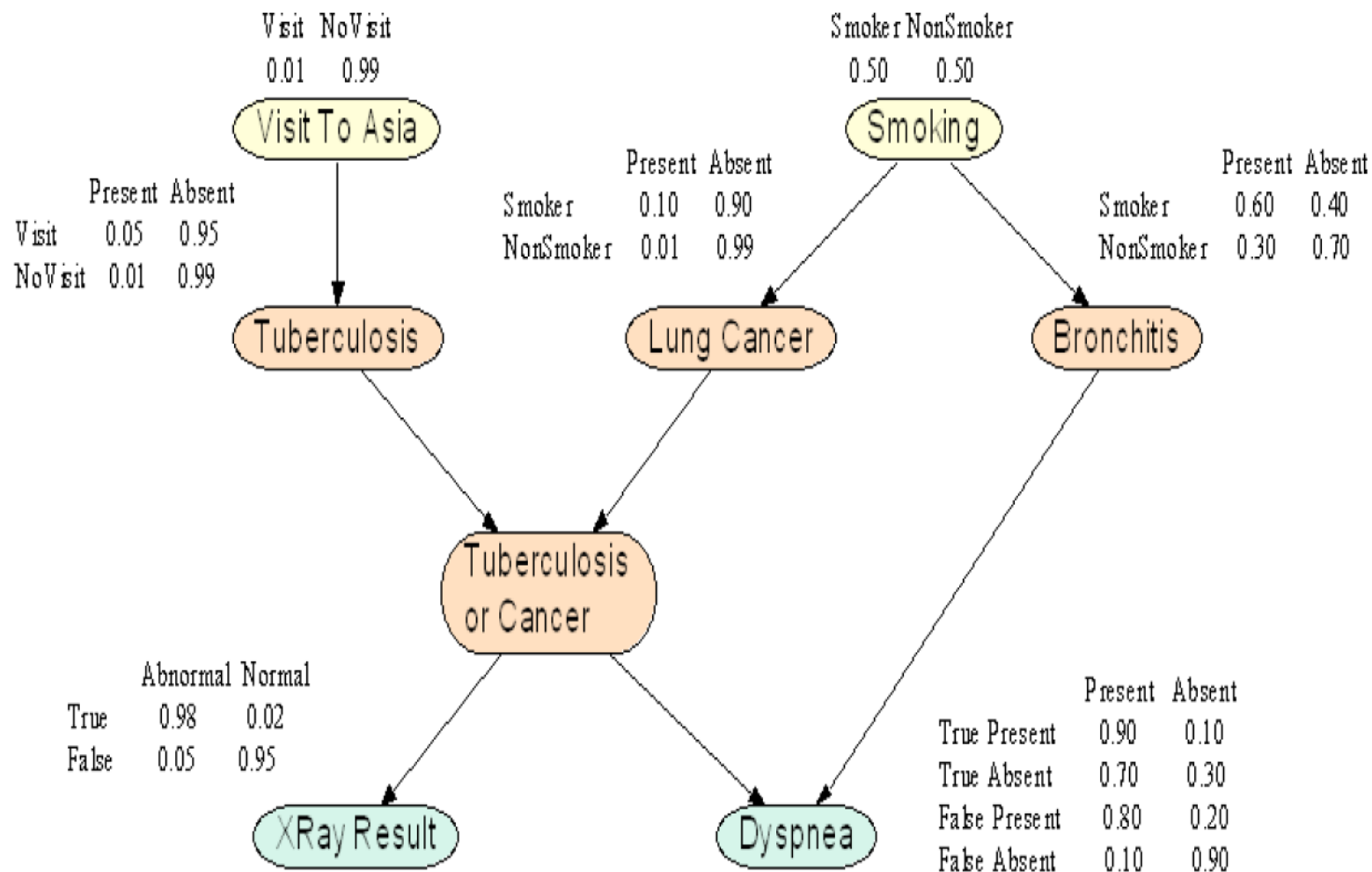
## Définition

- Un réseau bayésien est défini par
  - la description qualitative des dépendances (ou des indépendances conditionnelles) entre des variables  $S_i$   
*graphe orienté sans circuit (DAG)*
  - la description quantitative de ces dépendances  
*probabilités conditionnelles (CPD)*

## Conséquence

- $P(S) = \prod_{i=1}^n P(S_i | \text{parents}(S_i))$
- La loi jointe (globale) se décompose en un produit de lois conditionnelles locales
- RB = représentation compacte de la loi jointe  $P(S)$

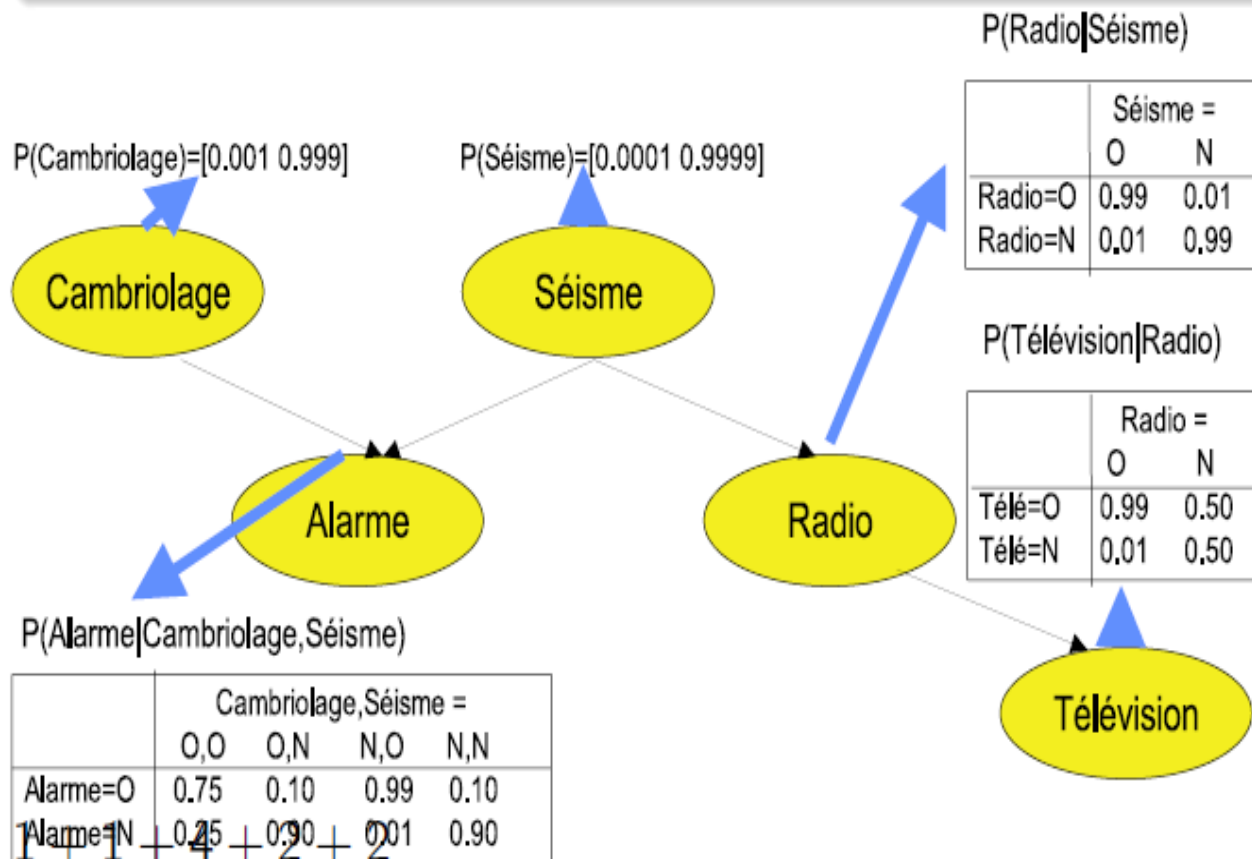
# Exemple



# Dimension d'un réseau bayésien

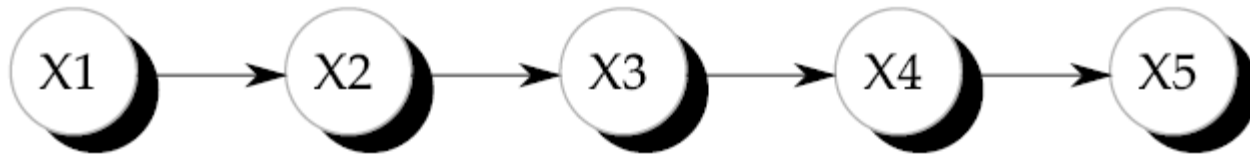
## Définition

Nombre de paramètres (indépendants) nécessaires pour décrire l'ensemble des CPD associées au RB

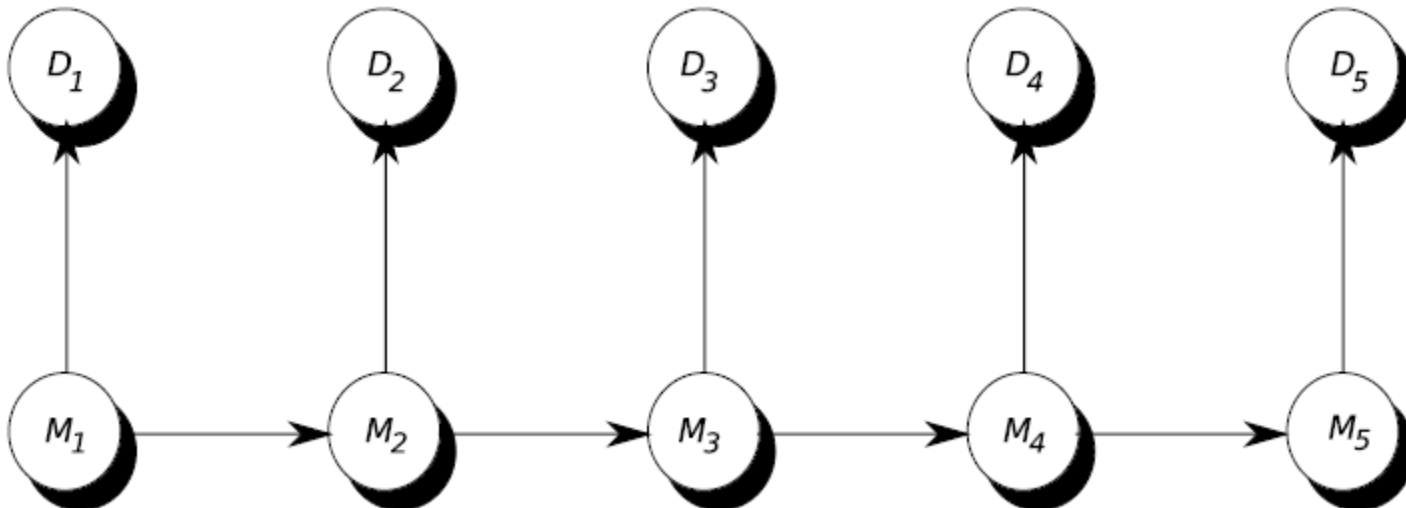


# Cas particulier : HMM à horizon fixé

- Chaîne de Markov homogène d'ordre 1



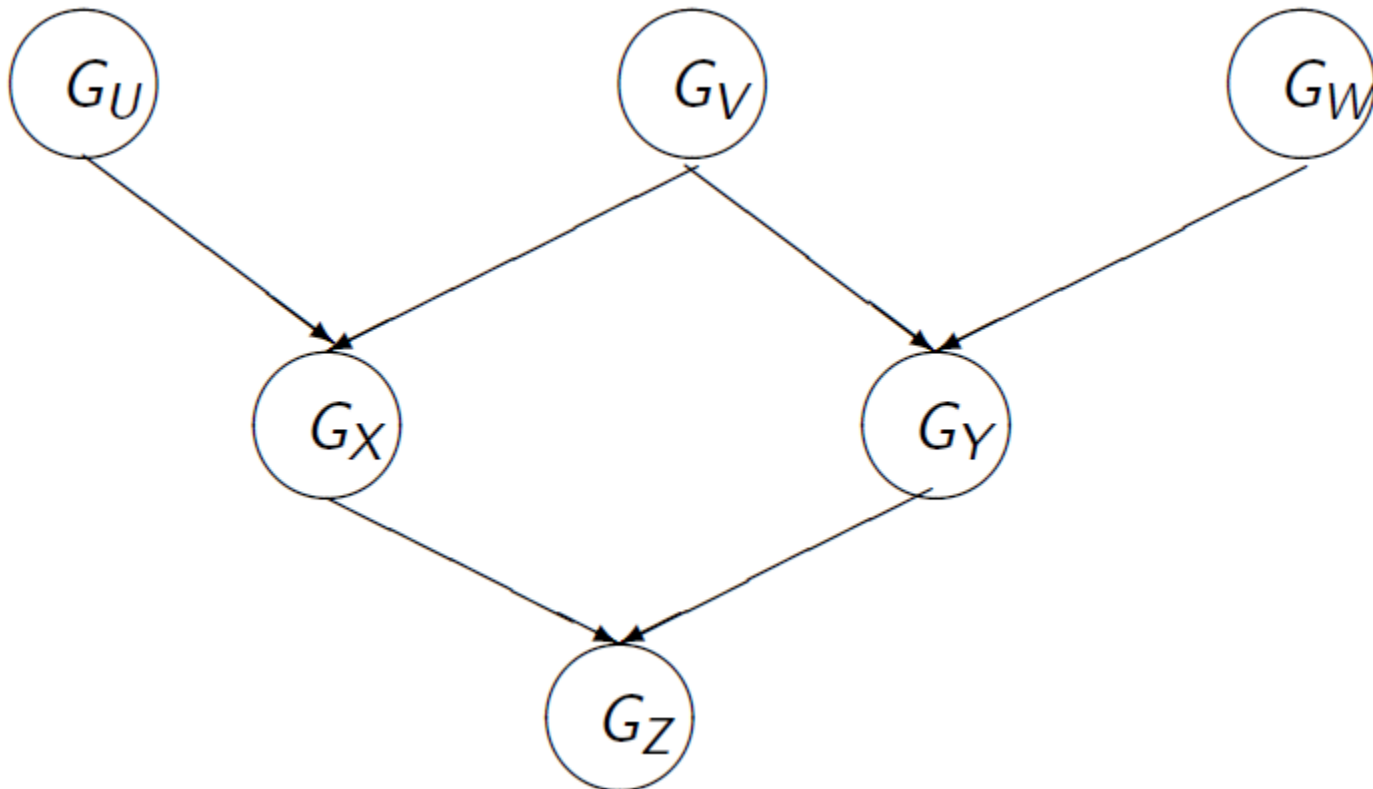
- Chaîne de Markov cachée homogène d'ordre 1





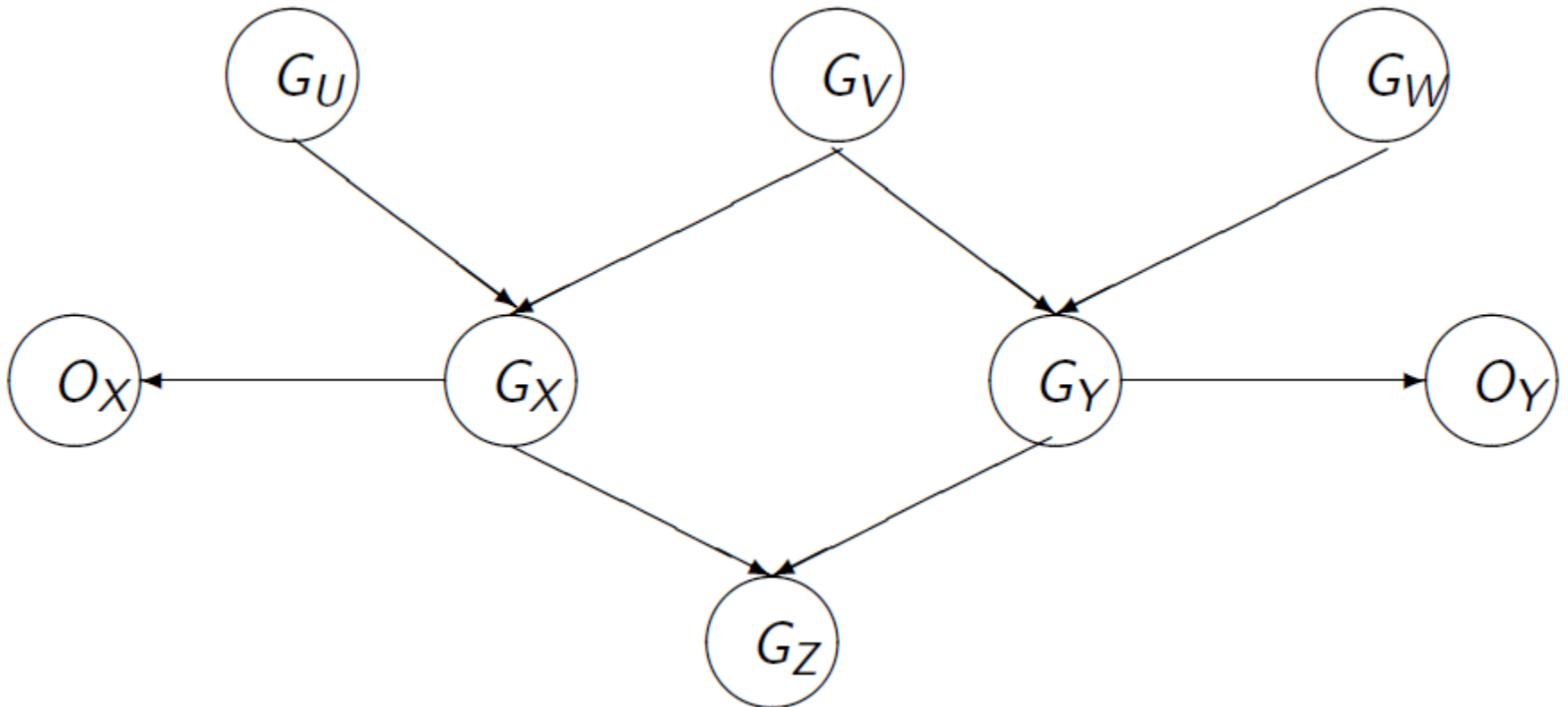
# Réseau de génotypes

Soit un pedigree avec 6 individus, notés  $U, V, W, X, Y, Z$ , tels que les parents de l'individu  $X$  soient  $\{U, V\}$ , de  $Y$  soient  $\{V, W\}$  et de  $Z$  soient  $\{X, Y\}$ .



# Réseau de génotypes

Comment étendre ce réseau pour prendre en compte des observations incertaines (possiblement erronées) sur les génotypes  $G_X$  et  $G_Y$  ?



## Probabilités a priori des génotypes fondateurs

$p(G_f = \{a, a\})$	0.25
$p(G_f = \{a, b\})$	0.5
$p(G_f = \{b, b\})$	0.25

## Probabilités de transmission des génotypes

$p(G G_p, G_m)$	$G_m = \{a, a\}$		$G_m = \{a, b\}$		$G_m = \{b, b\}$	
$G_p = \{a, a\}$	$G = \{a, a\}$	1	$G = \{a, a\}$	0.5	$G = \{a, a\}$	0
	$G = \{a, b\}$	0	$G = \{a, b\}$	0.5	$G = \{a, b\}$	1
	$G = \{b, b\}$	0	$G = \{b, b\}$	0	$G = \{b, b\}$	0
$G_p = \{a, b\}$	$G = \{a, a\}$	0.5	$G = \{a, a\}$	0.25	$G = \{a, a\}$	0
	$G = \{a, b\}$	0.5	$G = \{a, b\}$	0.5	$G = \{a, b\}$	0.5
	$G = \{b, b\}$	0	$G = \{b, b\}$	0.25	$G = \{b, b\}$	0.5
$G_p = \{b, b\}$	$G = \{a, a\}$	0	$G = \{a, a\}$	0	$G = \{a, a\}$	0
	$G = \{a, b\}$	1	$G = \{a, b\}$	0.5	$G = \{a, b\}$	0
	$G = \{b, b\}$	0	$G = \{b, b\}$	0.5	$G = \{b, b\}$	1

## Probabilités d'erreur de génotypage

$p(O G)$	$G = \{a, a\}$	$G = \{a, b\}$	$G = \{b, b\}$
$O = \{a, a\}$	$1 - \epsilon$	$\frac{\epsilon}{2}$	$\frac{\epsilon}{2}$
$O = \{a, b\}$	$\frac{\epsilon}{2}$	$1 - \epsilon$	$\frac{\epsilon}{2}$
$O = \{b, b\}$	$\frac{\epsilon}{2}$	$\frac{\epsilon}{2}$	$1 - \epsilon$

# Allele network

Maternal allele at  
locus 1 of person 1



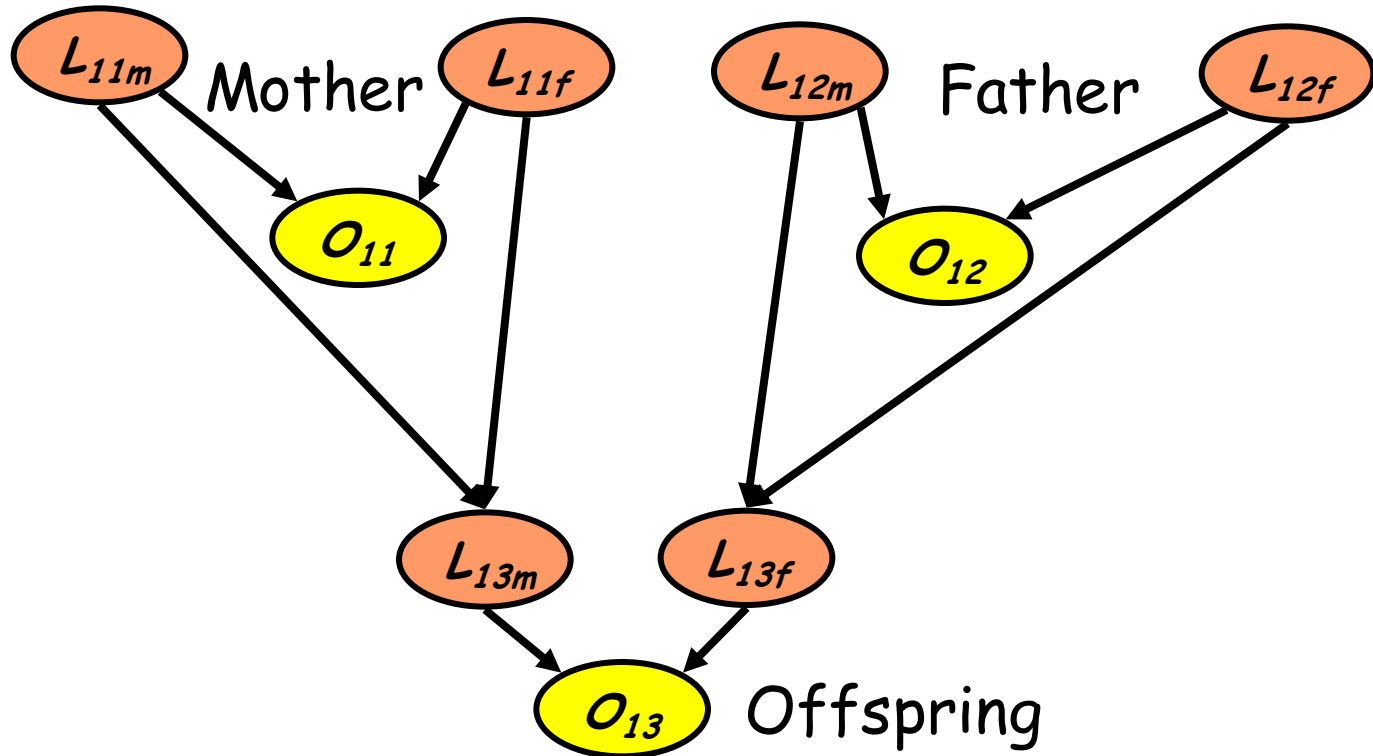
Paternal allele at  
locus 1 of person 1

Unordered allele pair at  
locus 1 of person 1 = data

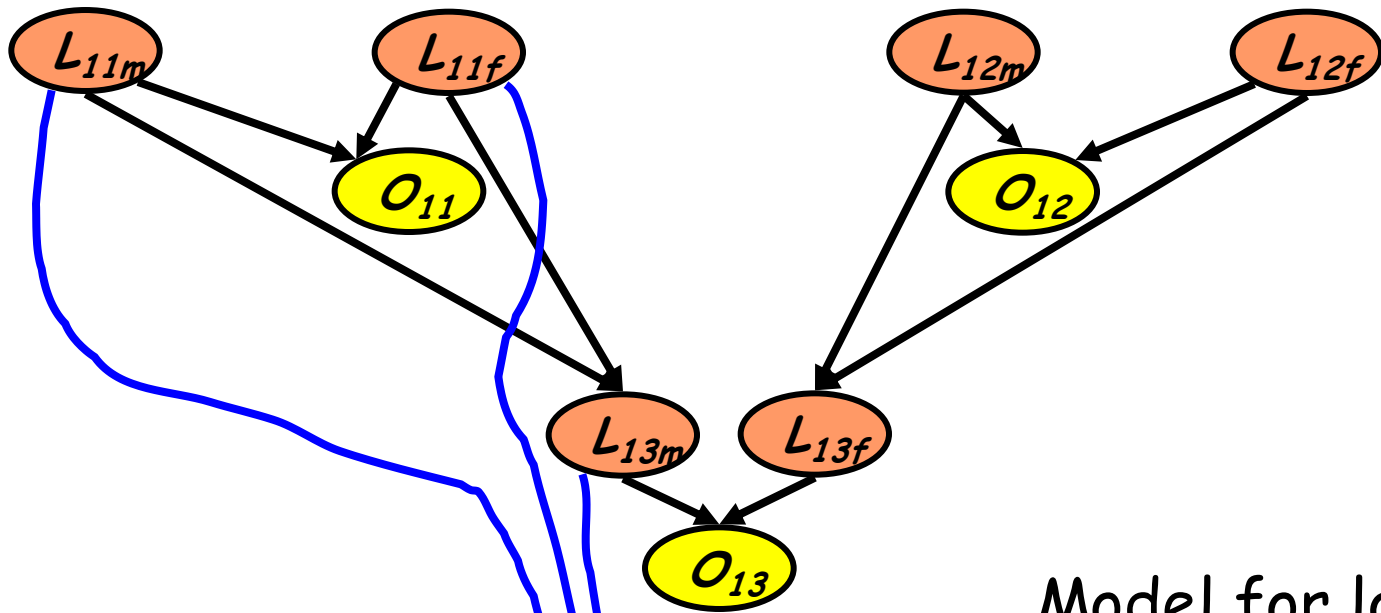
$p(L_{11m} = a)$  is the  
frequency of allele  $a$ .

$p(O_{11} \mid l_{11m}, l_{11f}) = 0$  or  $1$   
depending on consistency

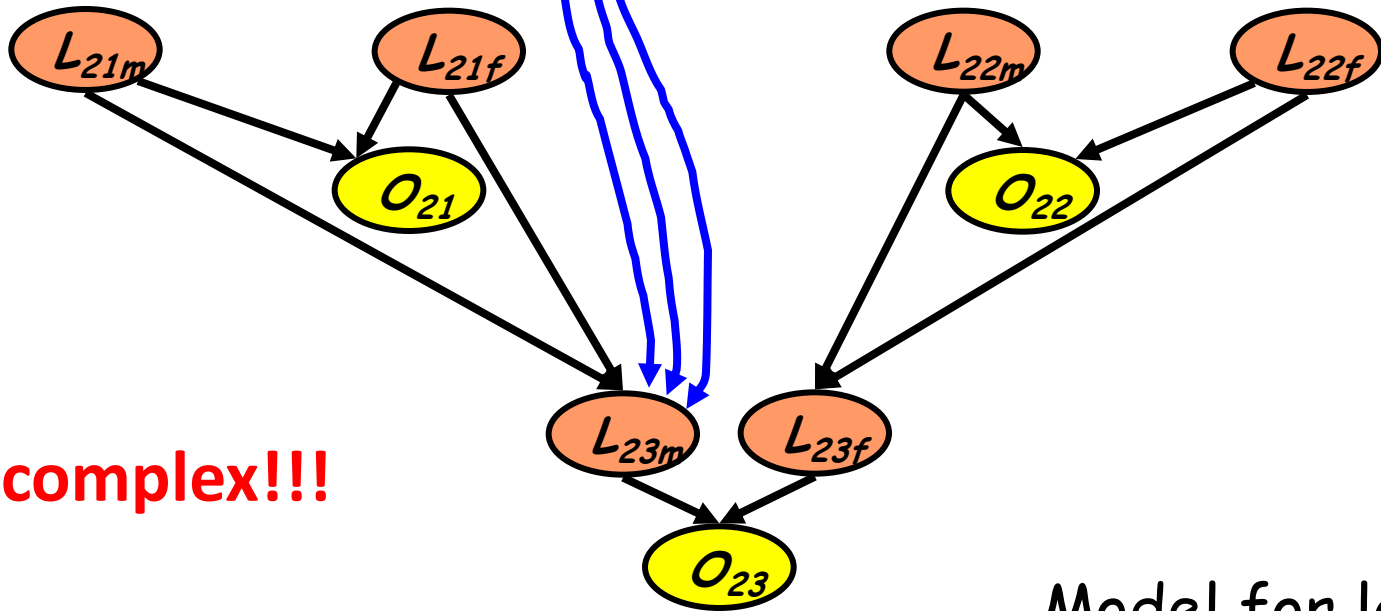
# Allele network



$$p(l_{13m} \mid l_{11m}, l_{11f}) = 1/2 \quad \text{if } l_{13m} = l_{11m} \text{ or } l_{13m} = l_{11f}$$
$$p(l_{13m} \mid l_{11m}, l_{11f}) = 0 \quad \text{otherwise}$$



Model for locus 1

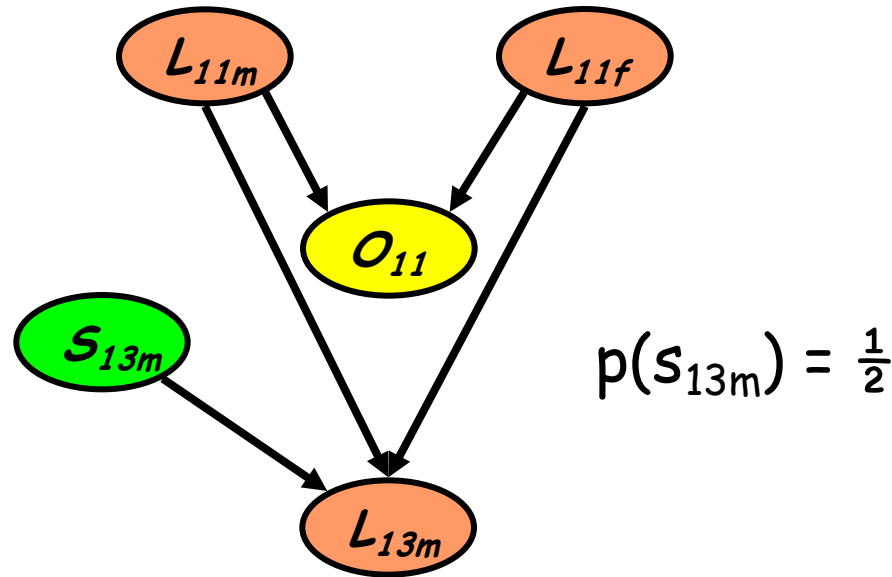


**Too complex!!!**

Model for locus 2

# Adding a selector variable

Selector of maternal allele at locus 1 of person 3



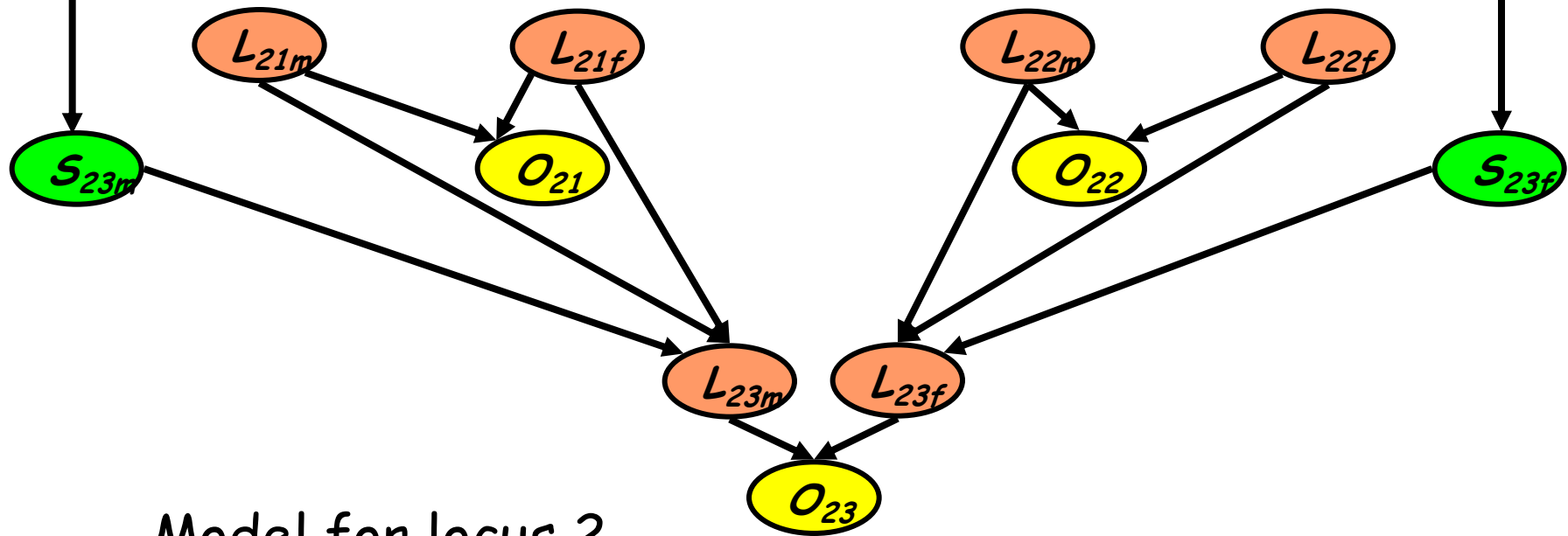
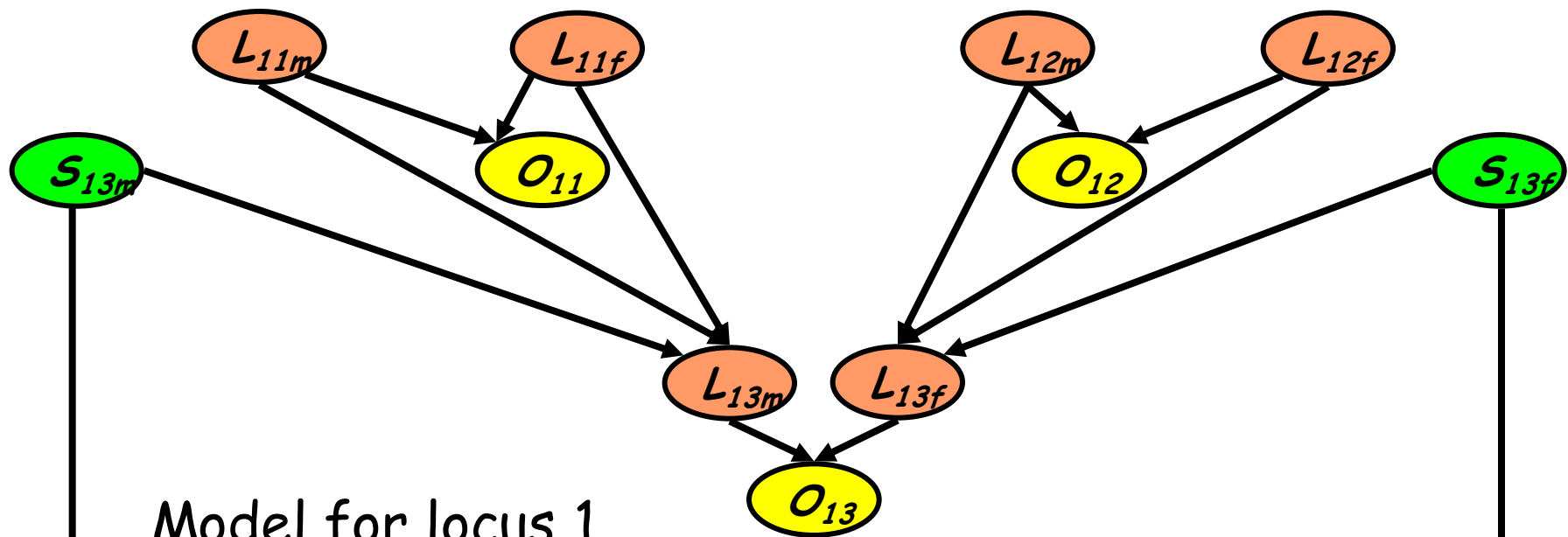
Maternal allele at locus 1 of person 3 (offspring)

Selector variables  $S_{ijm}$  are 0 or 1 depending on whose allele is transmitted to offspring  $i$  at maternal locus  $j$ .

$$p(l_{13m} \mid l_{11m}, l_{11f}, S_{13m}=0) = 1 \text{ if } l_{13m} = l_{11m}$$

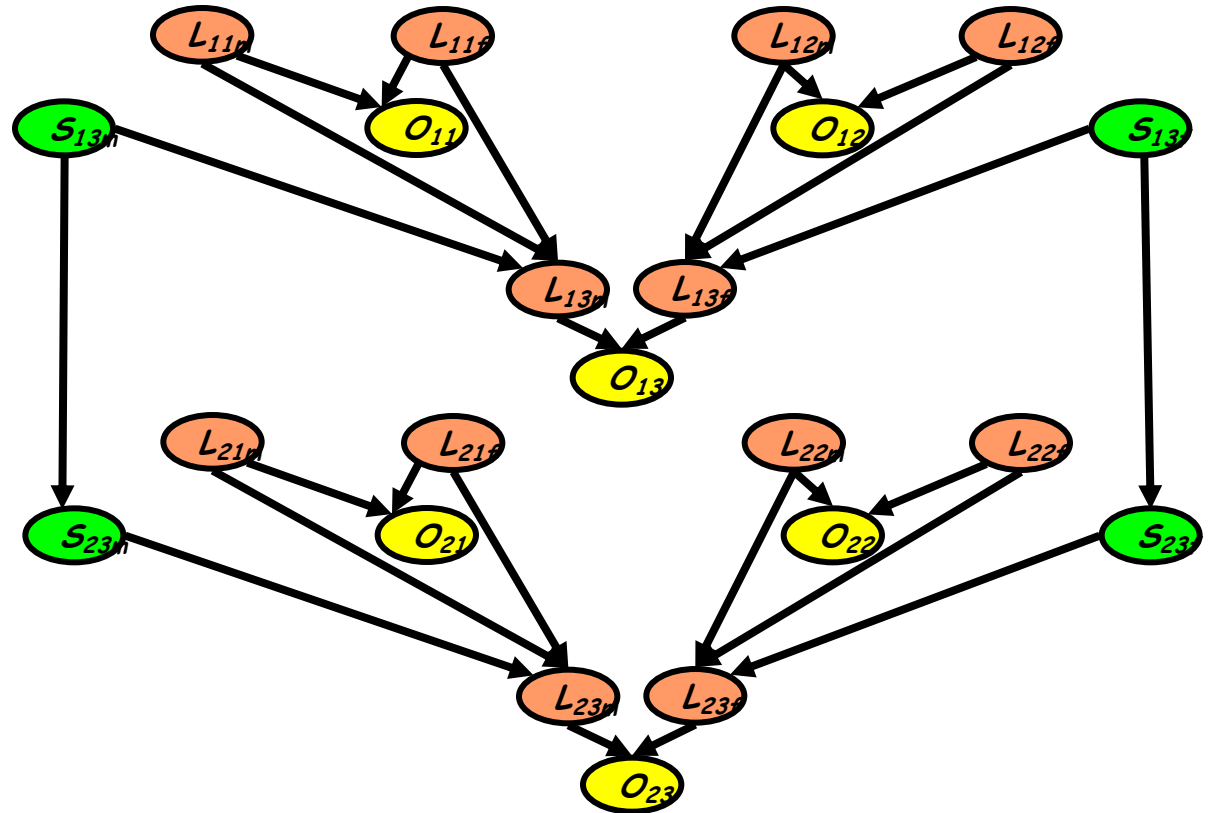
$$p(l_{13m} \mid l_{11m}, l_{11f}, S_{13m}=1) = 1 \text{ if } l_{13m} = l_{11f}$$

$$p(l_{13m} \mid l_{11m}, l_{11f}, s_{13m}) = 0 \text{ otherwise}$$





# Segregation network

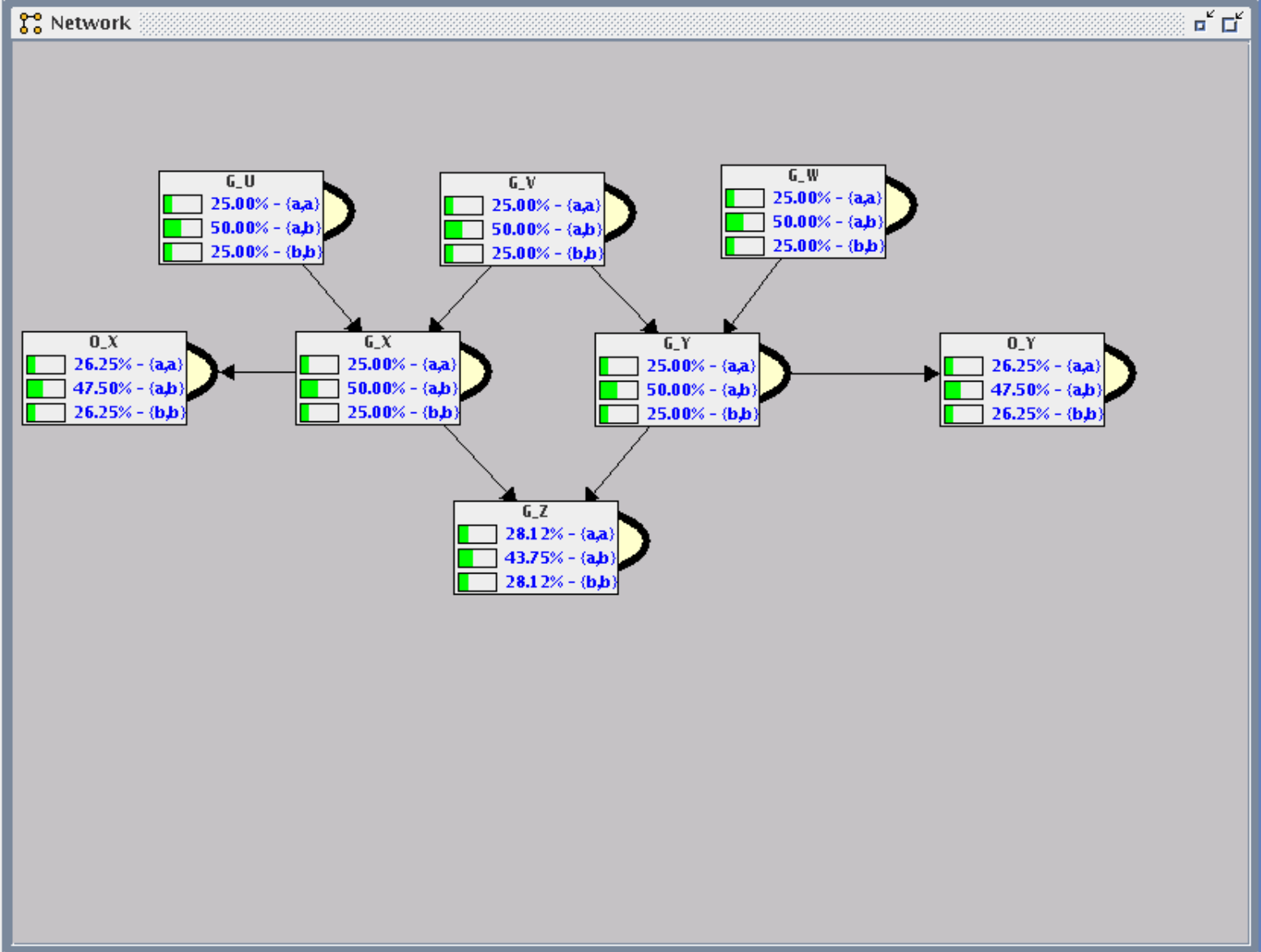


Probabilistic Model for Recombination:

$$p(s_{23t} | s_{13t}, \theta) = \begin{bmatrix} 1 - \theta & \theta \\ \theta & 1 - \theta \end{bmatrix} \quad \text{where } t \in \{m, f\}$$

$\theta$  is the **recombination fraction** between loci 2 & 1.

- in-out degree
- root
    - G\_U
    - G\_V
    - G\_W
  - internal
    - G\_X
    - G\_Y
  - leaf
    - G\_Z
    - O\_X
    - O\_Y



in-out degree

root

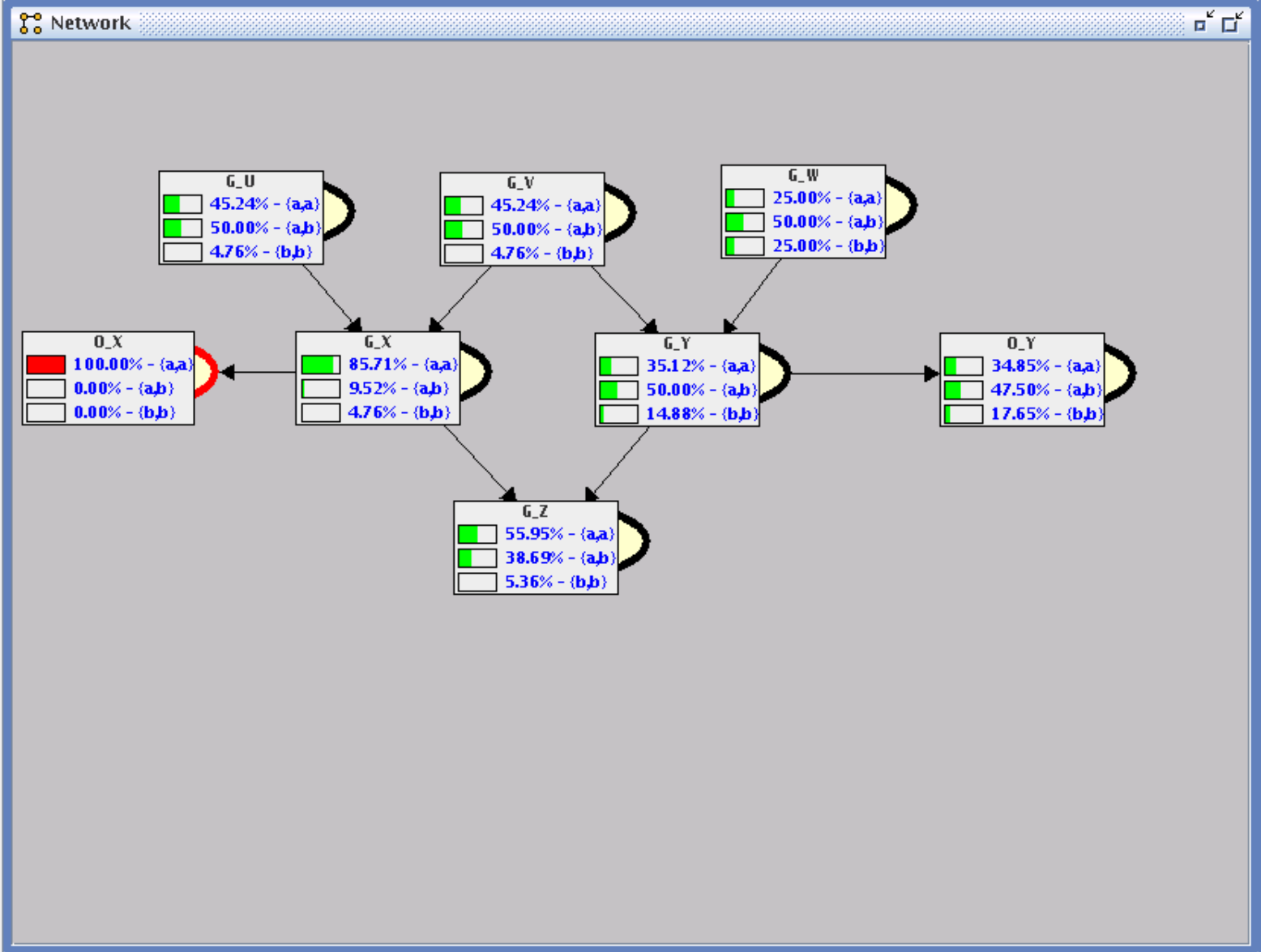
- G\_U
- G\_V
- G\_W

internal

- G\_X
- G\_Y

leaf

- G\_Z
- O\_X = ...
- O\_Y



in-out degree

root

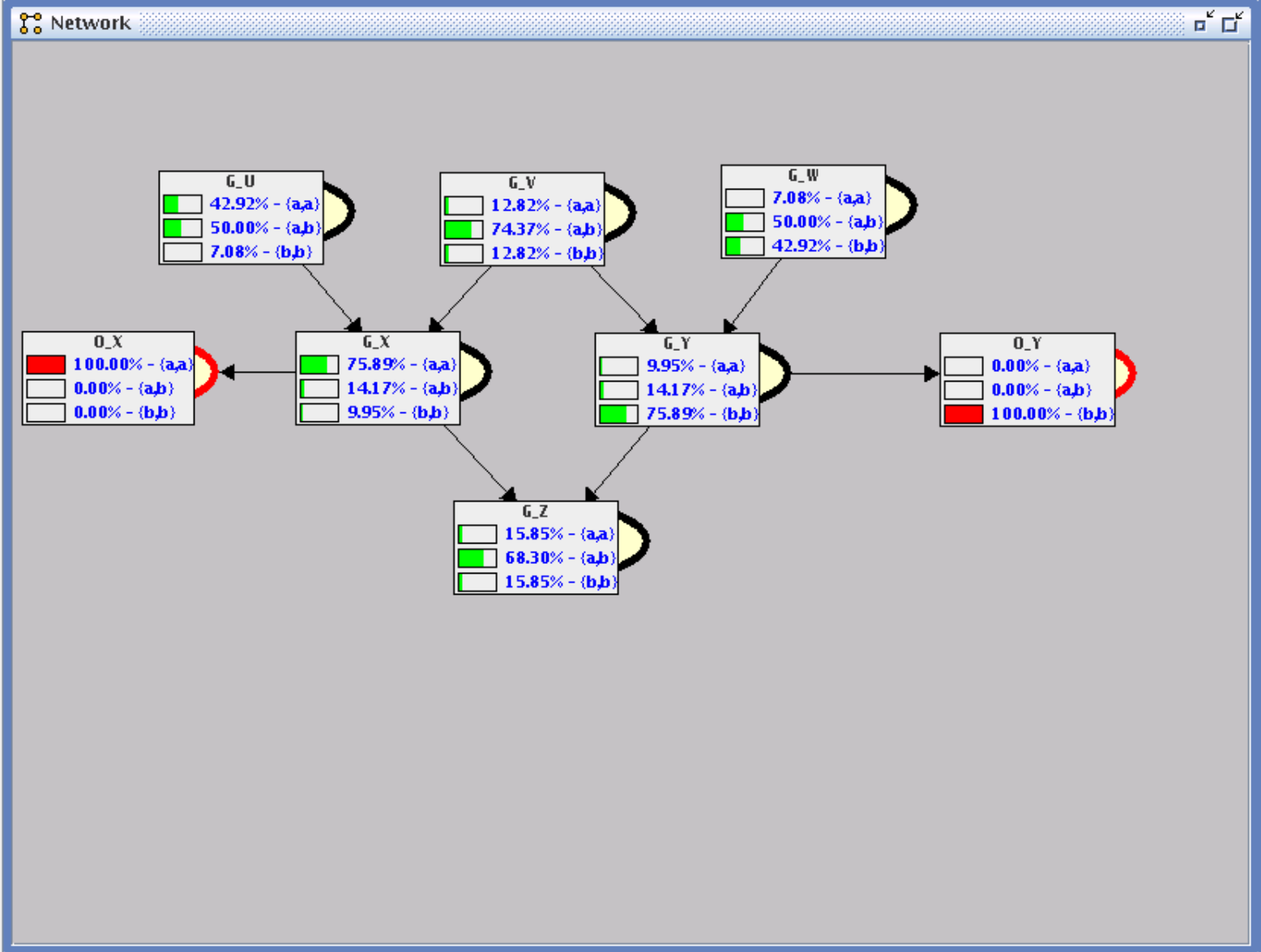
- G\_U
- G\_V
- G\_W

internal

- G\_X
- G\_Y

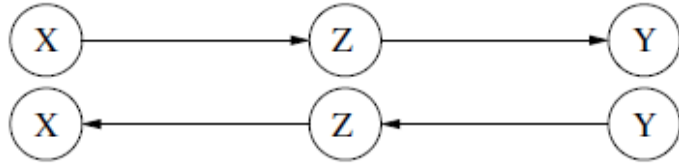
leaf

- G\_Z
- O\_X = ...
- O\_Y = ...



# Circulation de l'information

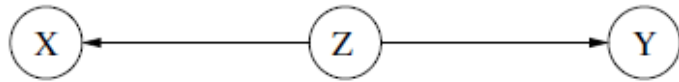
## Connexion en série



*ne circule de X vers Y  
que si Z n'est pas connu*

$$G_X = \{a, a\} \implies p(G_Y = \{a, a\}) > p(G_Y = \{b, b\})$$

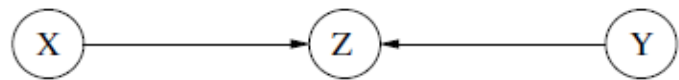
## Connexion divergente



*ne circule de X vers Y  
que si Z n'est pas connu*

$$G_X = \{a, a\} \implies p(G_Y = \{a, a\}) > p(G_Y = \{b, b\})$$

## Connexion convergente



*ne circule de X vers Y  
que si Z est connu*

$$G_X = \{a, a\}, G_Z = \{a, b\} \implies b \in G_Y$$

# Séparation dirigée (d-separation)

$X$  et  $Y$  sont d-séparés par  $Z$  (noté  $\langle X|Z|Y \rangle$ ) si pour **toutes** les chaînes entre  $X$  et  $Y$  :

- La chaîne converge en un sommet  $W$ , tel que  $W \neq Z$  et  $W$  n'est pas une cause directe ou *indirecte* de  $Z$  ;
- *OU* la chaîne passe par  $Z$ , et est soit divergente, soit en série au sommet  $Z$ .

Il y a alors blocage de l'information par la connaissance de  $Z$ .

Extension au cas où  $\mathbf{X}$ ,  $\mathbf{Y}$  et  $\mathbf{Z}$  sont des ensembles disjoints.

L'absence de chaîne entre  $X$  et  $Y$  implique la d-séparation entre  $X$  et  $Y$  quelque soit  $Z$ .

Complexité de tester  $\langle X|Z|Y \rangle$  linéaire en la taille du réseau.

# Lien entre d-separation et indépendance conditionnelle

## Indépendance conditionnelle

Soit des ensembles disjoints  $\mathbf{X}, \mathbf{Y}, \mathbf{Z} \subset \mathbf{V}$ .  $\mathbf{X}$  est indépendant de  $\mathbf{Y}$  conditionnellement à  $\mathbf{Z}$  (noté  $\mathbf{X} \perp\!\!\!\perp \mathbf{Y} | \mathbf{Z}$ ) ssi :

$$\begin{aligned} \mathbf{X} \perp\!\!\!\perp \mathbf{Y} | \mathbf{Z} &\iff \begin{cases} p(\mathbf{X} | \mathbf{Y}, \mathbf{Z}) = p(\mathbf{X} | \mathbf{Z}) \\ \text{et } p(\mathbf{Y} | \mathbf{X}, \mathbf{Z}) = p(\mathbf{Y} | \mathbf{Z}) \end{cases} \\ &\iff p(\mathbf{X}, \mathbf{Y} | \mathbf{Z}) = p(\mathbf{X} | \mathbf{Z}) \cdot p(\mathbf{Y} | \mathbf{Z}) \\ &\iff p(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = p(\mathbf{X} | \mathbf{Z}) \cdot p(\mathbf{Y} | \mathbf{Z}) \cdot p(\mathbf{Z}) \end{aligned}$$

## Théorème

Soit un réseau Bayésien  $\vec{G}$  définissant la loi de probabilité jointe  $p(\mathbf{V})$ . Il vérifie la *propriété orientée de Markov globale (OG)* :

$$\forall \mathbf{X}, \mathbf{Y}, \mathbf{Z} \subset \mathbf{V} \text{ disjoints, } \langle \mathbf{X} | \mathbf{Z} | \mathbf{Y} \rangle \implies \mathbf{X} \perp\!\!\!\perp \mathbf{Y} | \mathbf{Z}$$

# Limites d'expressivité des réseaux bayésiens

- Il existe des lois dont le *modèle d'indépendance* n'est pas représentable par un réseau bayésien  
Exemple:  $M = \{ X \perp\!\!\!\perp W \mid \{Y,Z\}, Y \perp\!\!\!\perp Z \mid \{X,W\} \}$
- *Réseau de Markov* : graphe non-orienté avec séparation directe
- Quelque soit le formalisme (réseau bayésien ou réseau de Markov), il existe des lois non représentables



# Requêtes

## Vraisemblance des observations (PR)

Calculer la probabilité  $p(\mathbf{E} = \mathbf{e})$  ayant observé les variables  $\mathbf{E} \subset \mathbf{V}$ .

## Probabilités Marginales (MAR)

Calculer la probabilité *a posteriori*  $p(V_i | \mathbf{E} = \mathbf{e})$  de toutes les variables  $V_i \in \mathbf{V} \setminus \mathbf{E}$ .

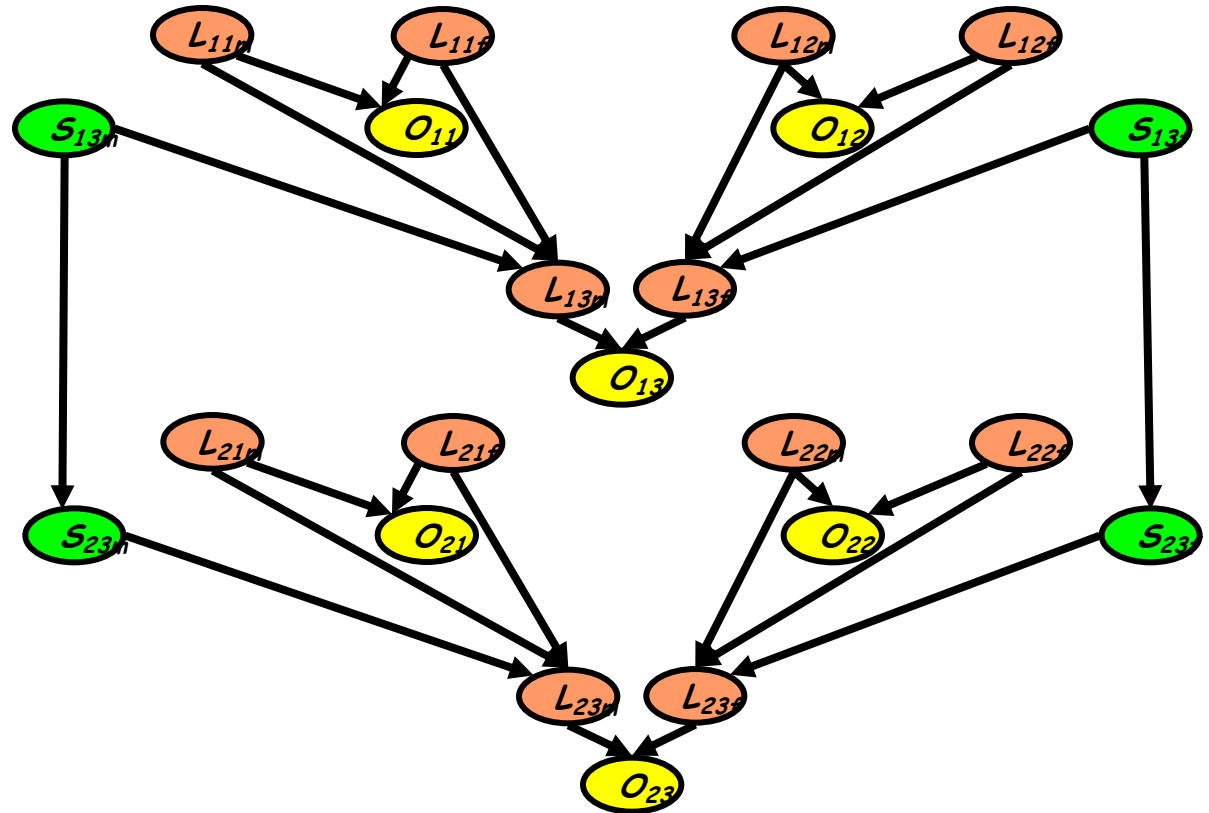
## Maximum a Posteriori hypothesis (MAP)

Rechercher une affectation *partielle*  $\mathbf{u}$  d'un sous-ensemble de variables  $\mathbf{U} \subseteq \mathbf{V} \setminus \mathbf{E}$  de probabilité  $p(\mathbf{U} = \mathbf{u} | \mathbf{E} = \mathbf{e})$  maximum.

## Most Probable Explanation (MPE)

Rechercher une affectation *complète*  $\mathbf{u}$  des variables  $\mathbf{U} = \mathbf{V} \setminus \mathbf{E}$  de probabilité  $p(\mathbf{U} = \mathbf{u} | \mathbf{E} = \mathbf{e})$  maximum.

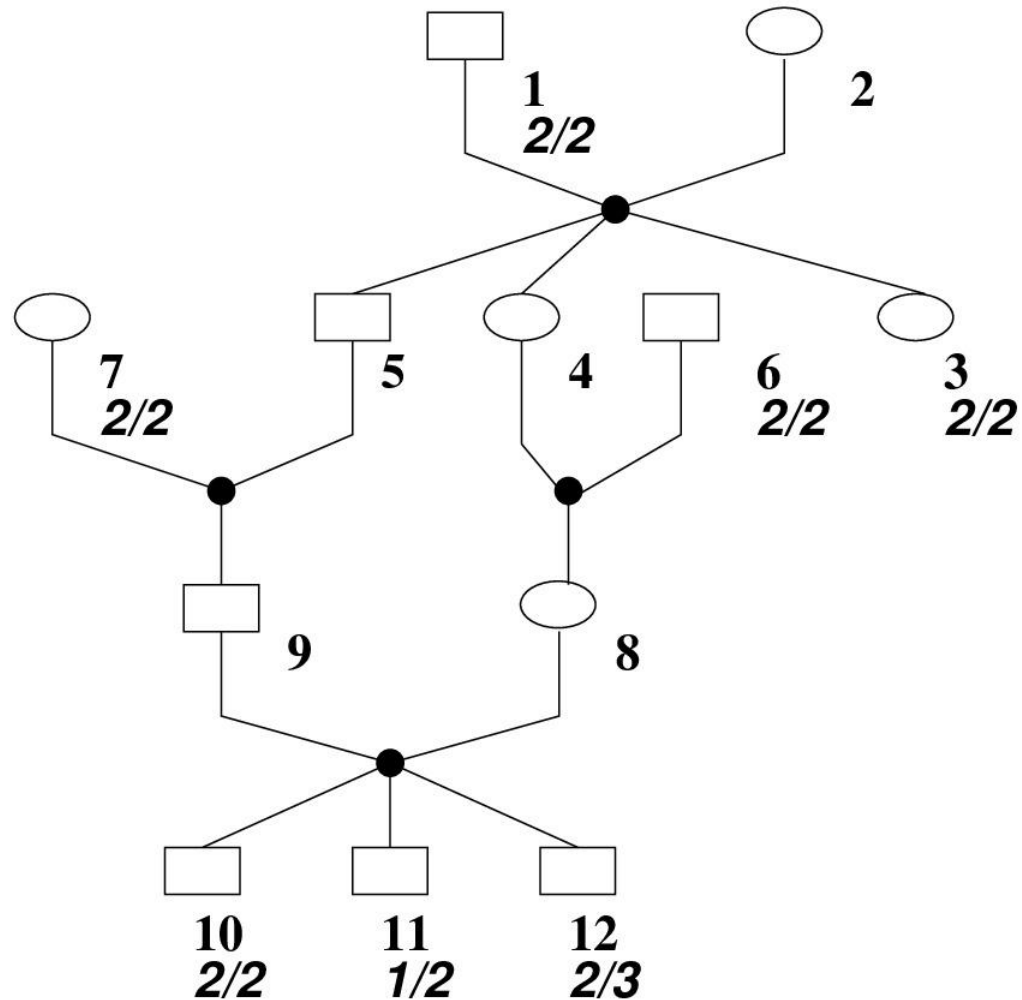
# Segregation network



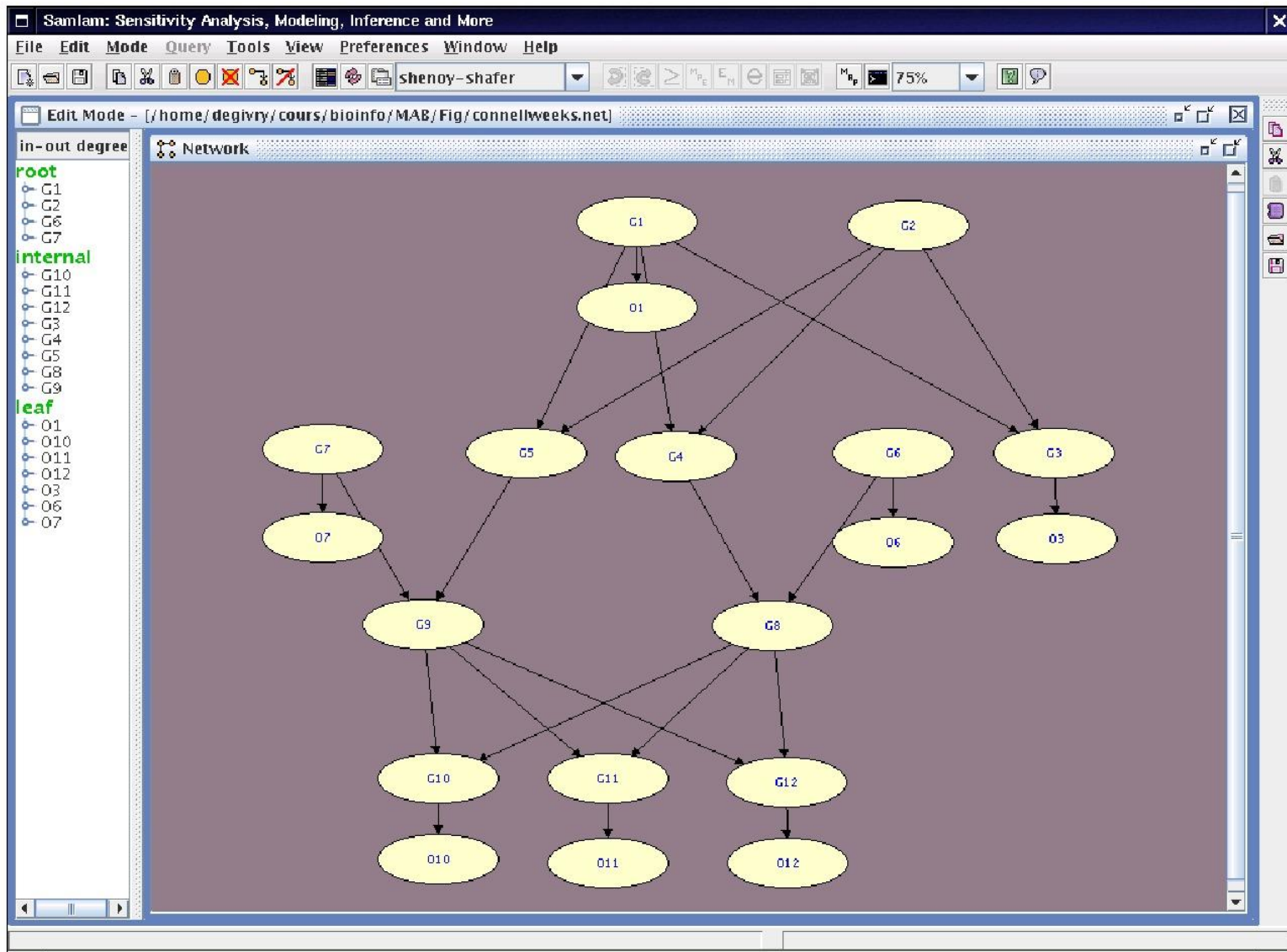
Genetic linkage analysis (**PR**):  $\operatorname{argmax}_{\theta} \sum_{L,S} p(L, S, \mathbf{O} \mid \theta)$

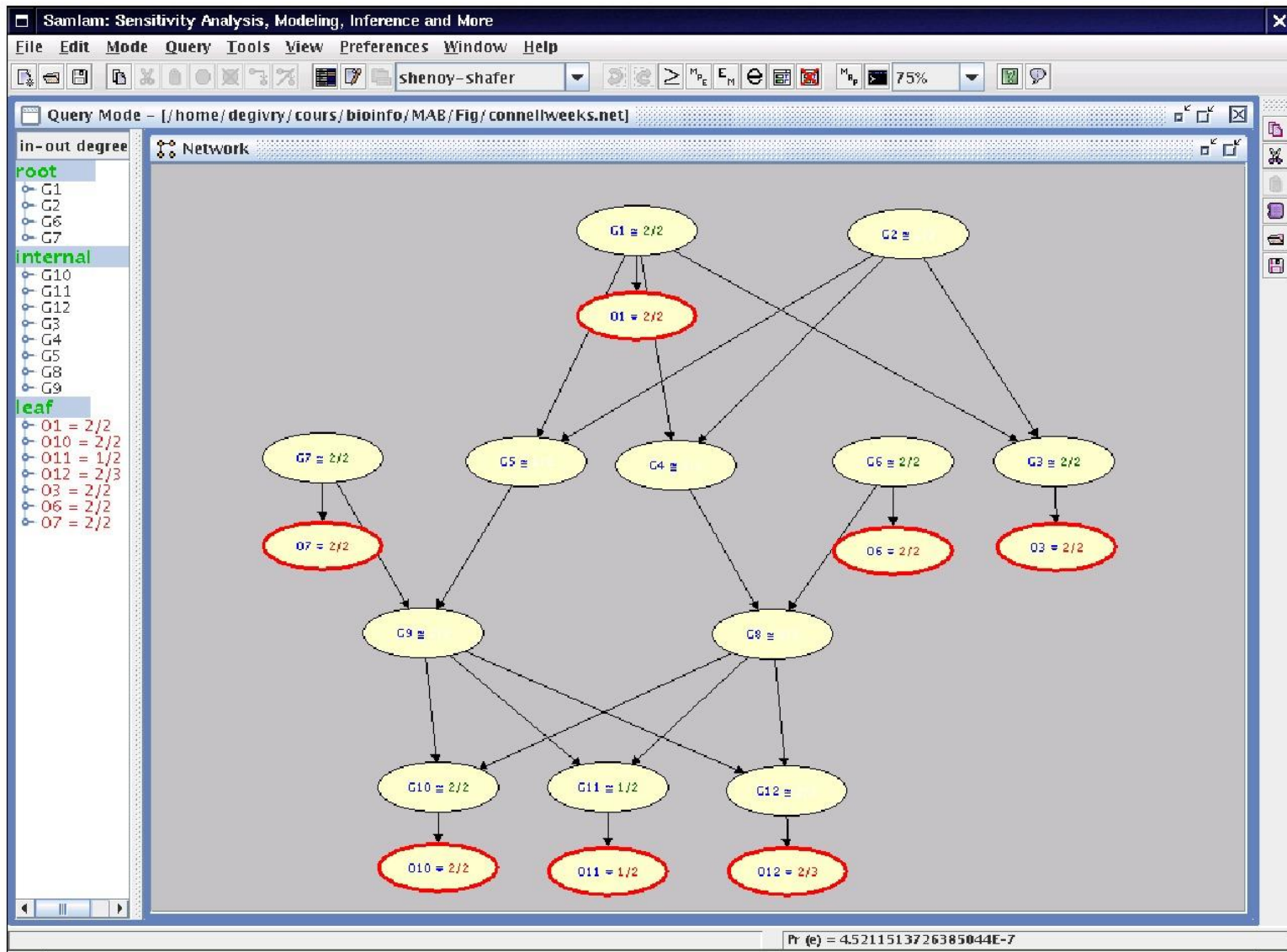
Haplotype reconstruction (**MPE**):  $\operatorname{argmax}_{L,S} p(L, S \mid \mathbf{O}, \theta)$

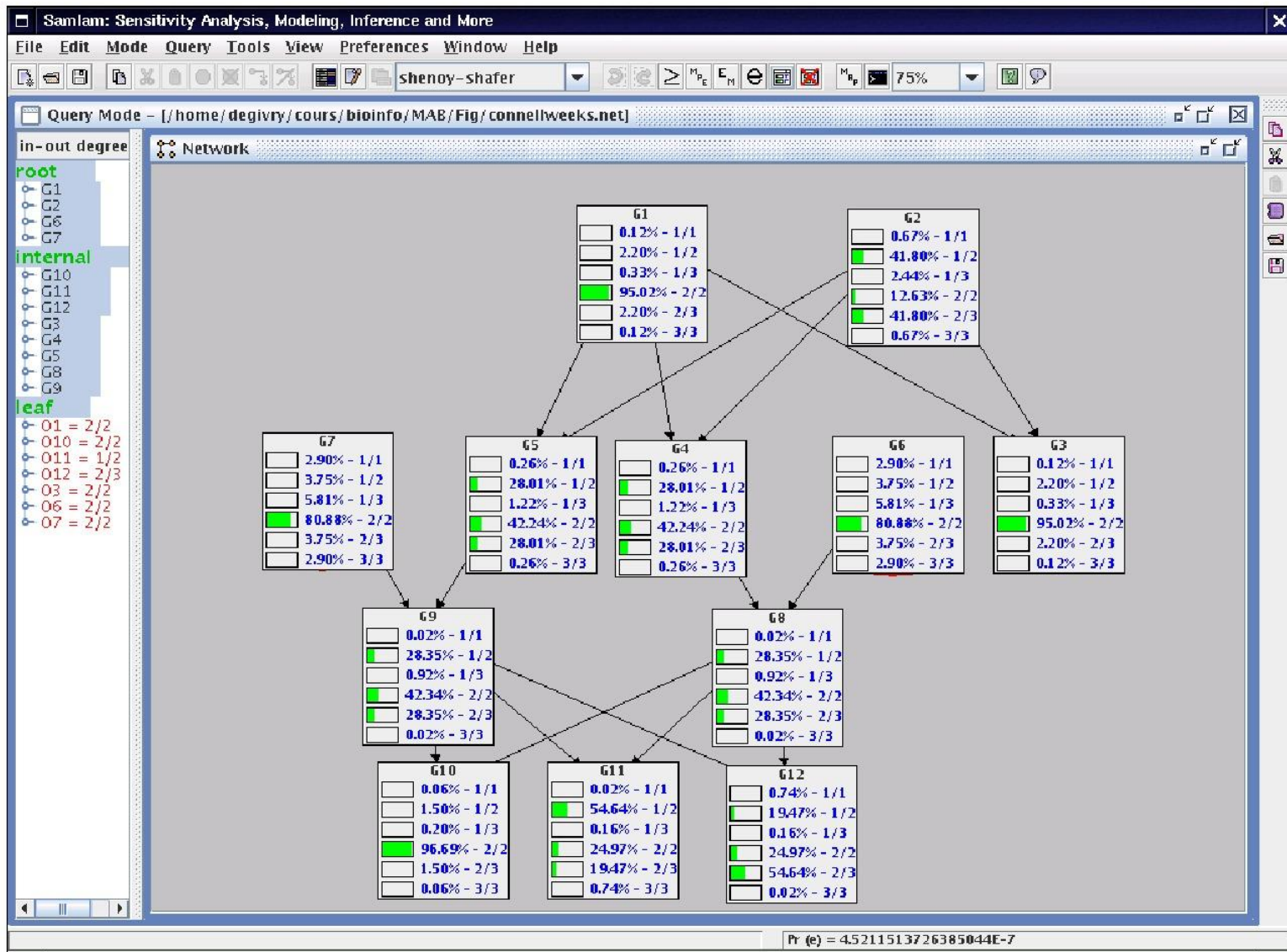
## Exemple de pedigree avec génotypes observés



## Réseau de génotypes







## MAP

Samlam: Sensitivity Analysis, Modeling, Inference and More

File Edit Mode Query Tools View Preferences Window Help

shenoy-shafer

Query Mode - [//home/degivry/cours/bioinfo/MAB/Fig/connellweeks.net]

in-out degree

root

- G1
- G2
- G6
- G7

internal

- G10
- G11
- G12
- G3
- G4
- G5
- G8
- G9

leaf

- O1 = 2/2
- O10 = 2/2
- O11 = 1/2
- O12 = 2/3
- O3 = 2/2
- O6 = 2/2
- O7 = 2/2

MAP Computation

$P(\text{MAP}, e) = 0.0000000568951151818332$   
 $P(\text{MAP} | e) = 0.1258420930698228$   
 Result is exact.

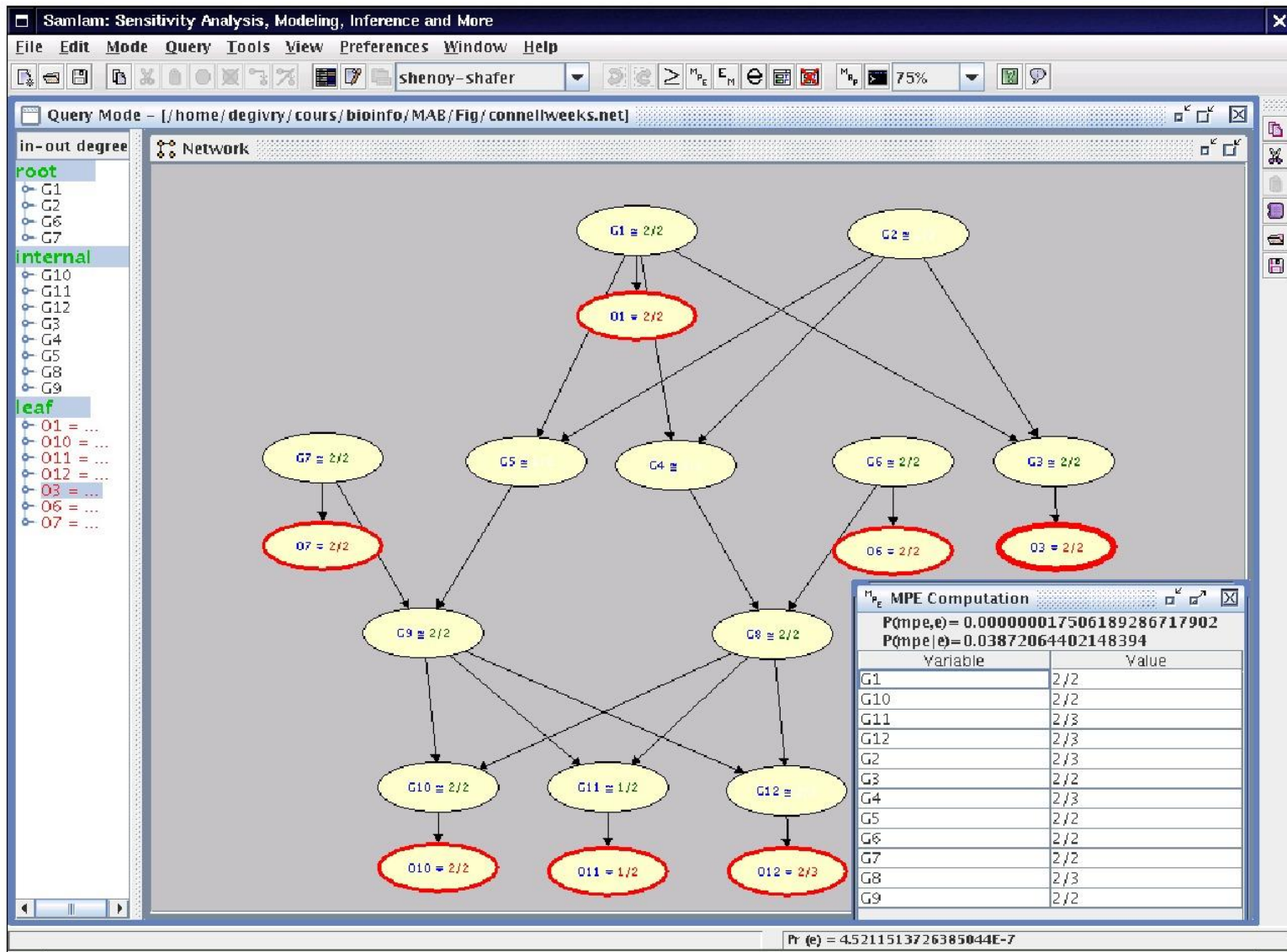
Variable	Value
G1	2/2
G10	2/2
G11	1/2
G12	1/2
G3	2/2
G6	2/2
G7	2/2

Copy Copy (+evidence)

Code Bandit Close

Pr (e) = 4.5211513726385044E-7

## MPE





# Requêtes et algorithmes

## – liens avec HMM

- Recherche de l'état le plus probable (MPE)
  - algorithme de Viterbi (opérateurs max-prod)
- Vraisemblance des observations (PR)
  - algorithme Forward (opérateurs sum-prod)
- Probabilités marginales (MAR)
  - algorithme Forward-Backward (opérateurs sum-prod)


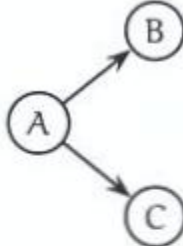
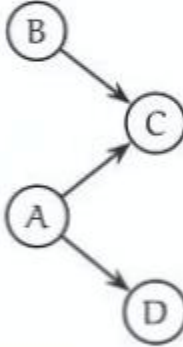
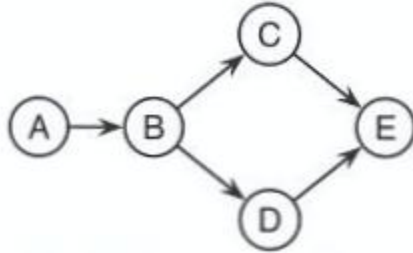
➔ Cadre général :

programmation dynamique non sérielle

# Inférence probabiliste

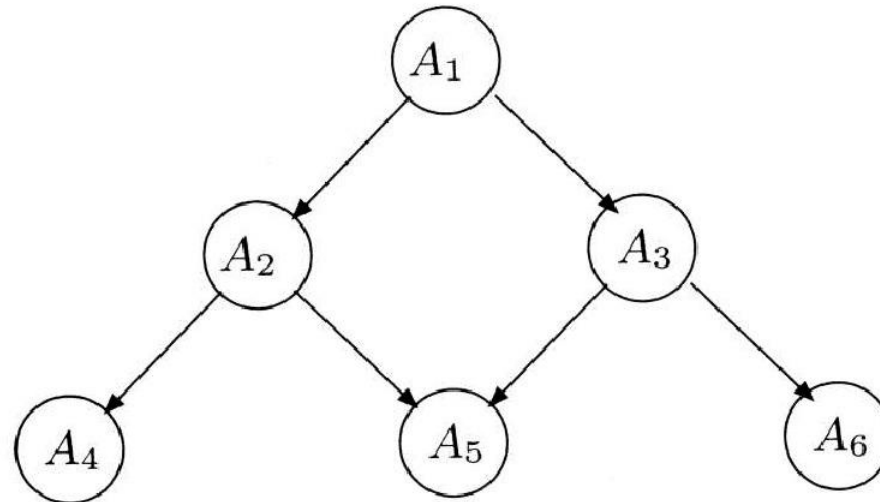
- Méthodes exactes
  - programmation dynamique non sérielle
  - recherche arborescente
  - méthodes hybrides
- Méthodes approchées
  - Méthodes par passage de messages (*loopy belief propagation,...*)
  - Méthodes par simulation (*MCMC,...*)

# Inférence probabiliste

	Chaîne	$p(C   A)?$
	Arbre	$p(C   B)?$
	Polyarbre	$p(D   B)?$
	Réseau avec boucles	$p(E   A)?$

# Algorithme d'inférence exacte fondé sur l'élimination de variable

Soit les variables  $\mathbf{V} = A_1, A_2, A_3, A_4, A_5, A_6$   
et les *potentiels*  $\phi_1(A_1) = p(A_1)$ ,  $\phi_2(A_2, A_1) = p(A_2|A_1)$ ,  
 $\phi_3(A_3, A_1) = p(A_3|A_1)$ ,  $\phi_4(A_4, A_2) = p(A_4|A_2)$ ,  
 $\phi_5(A_5, A_2, A_3) = p(A_5|A_2, A_3)$ ,  $\phi_6(A_6, A_3) = p(A_6|A_3)$ .



Calcul de  $p(A_4)$  ?

# Les deux grandes règles

Règle de la somme - élimination de variable - marginale

$$p(X) = \sum_Y p(X, Y)$$

Règle du produit

$$p(X, Y) = p(Y|X).p(X) = p(X|Y).p(Y)$$

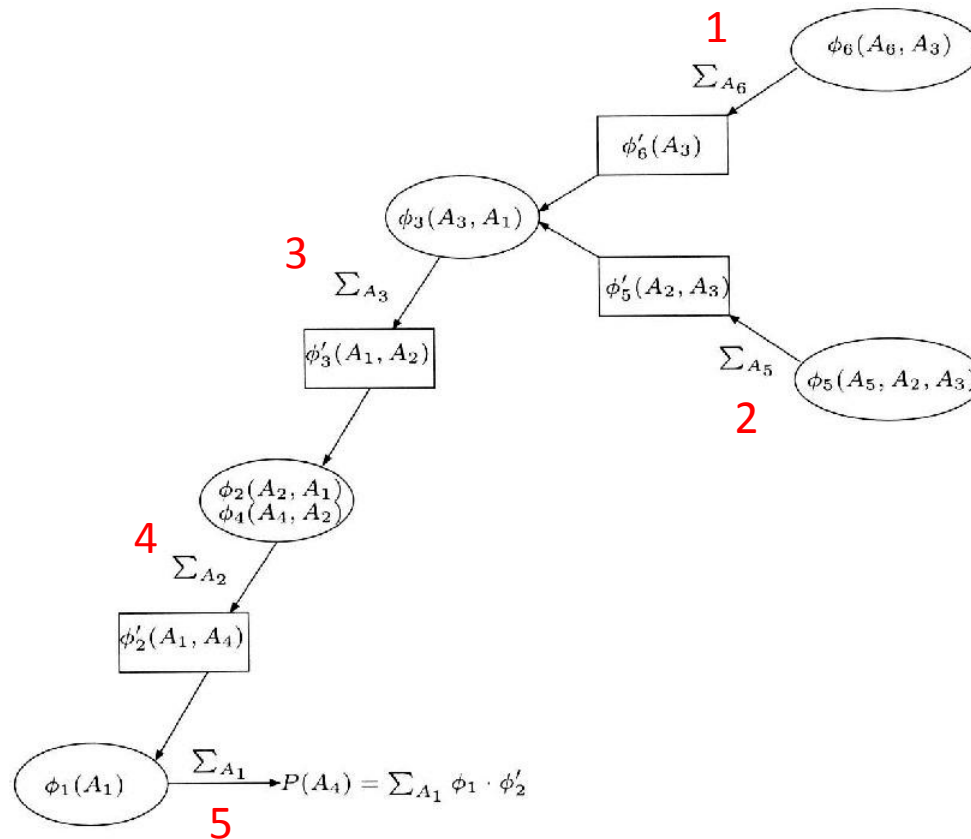
# Algorithme d'inférence par élimination de variable

Utilisation de la règle de distribution entre somme et produit.

$$\begin{aligned}
 p(A_4) &= \sum_{A_1, A_2, A_3, A_5, A_6} p(\mathbf{V}) \\
 &= \sum_{A_1, A_2, A_3, A_5, A_6} \phi_1 \phi_2 \phi_3 \phi_4 \phi_5 \phi_6 \\
 &= \sum_{A_1} \phi_1(A_1) \sum_{A_2} \phi_2(A_2, A_1) \phi_4(A_4, A_2) \sum_{A_3} \phi_3(A_3, A_1) \\
 &\quad \sum_{A_5} \phi_5(A_5, A_2, A_3) \sum_{A_6} \phi_6(A_6, A_3)
 \end{aligned}$$

# Factorisation des calculs selon un arbre

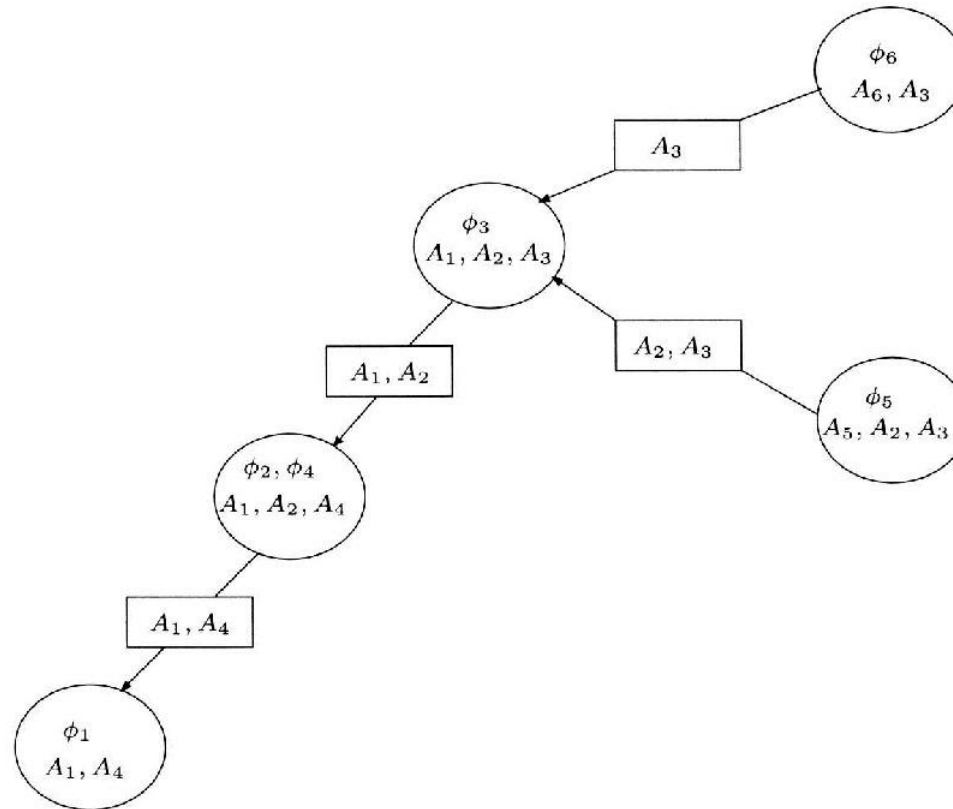
Introduction de **nouveaux** potentiels obtenus par marginalisation  
 (par ex.,  $\phi'_6(A_3) = \sum_{A_6} \phi_6(A_6, A_3)$ )



*Multiplication des potentiels dans les ellipses avec les potentiels entrants.*

# Impact de l'ordre d'élimination de variable

Ordre d'élimination :  $(A_6, A_5, A_3, A_2, A_1)$



Complexité exponentielle en le nombre de variables. Ici, temps  $O(d^3)$  et espace  $O(d^2)$ .



# Algorithme d'inférence exacte fondé sur un regroupement des variables

## Principe général

- Phase 1 : **Compilation**

Idée : regrouper les variables en noeuds de calcul (*clusters*) de manière à définir un arbre de calcul.

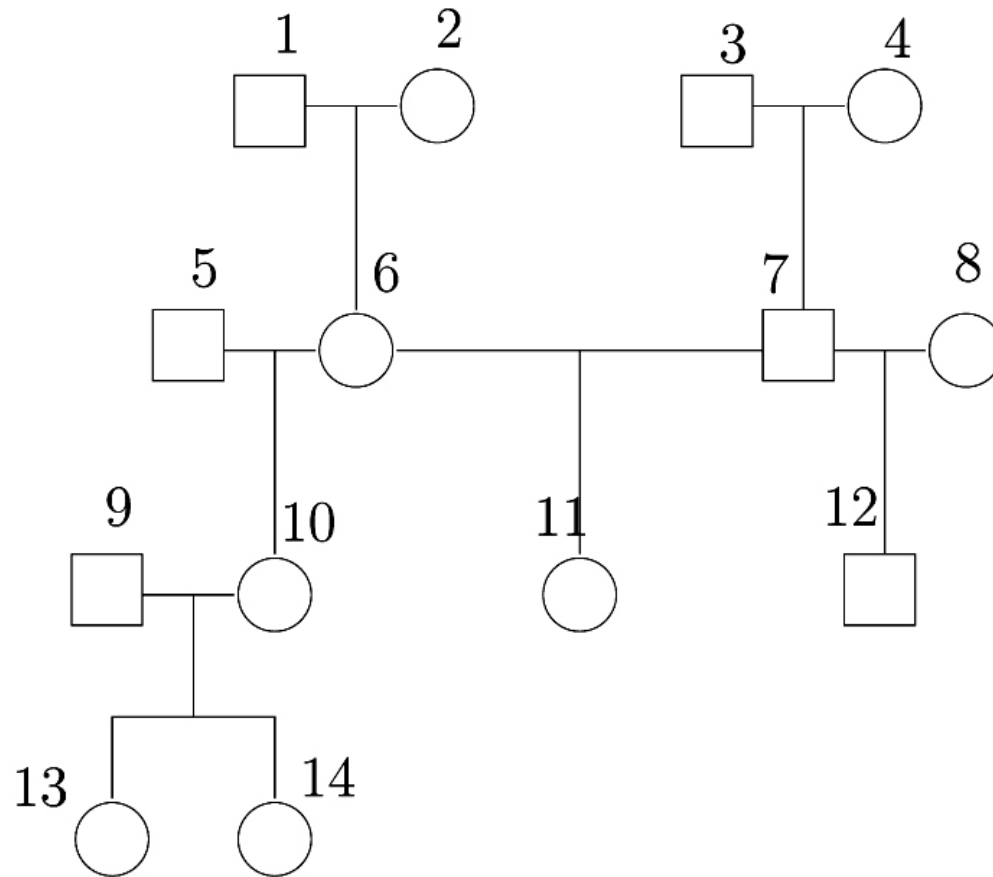
- Transformation du réseau Bayésien en graphe non-orienté
- Trouver un jeu d'élimination des variables
- Construire l'arbre de clusters
- Remplir les clusters avec les probabilités conditionnelles
- Ajouter les observations

*Remarque : manipulations de la phase 1 uniquement sur la structure du réseau*

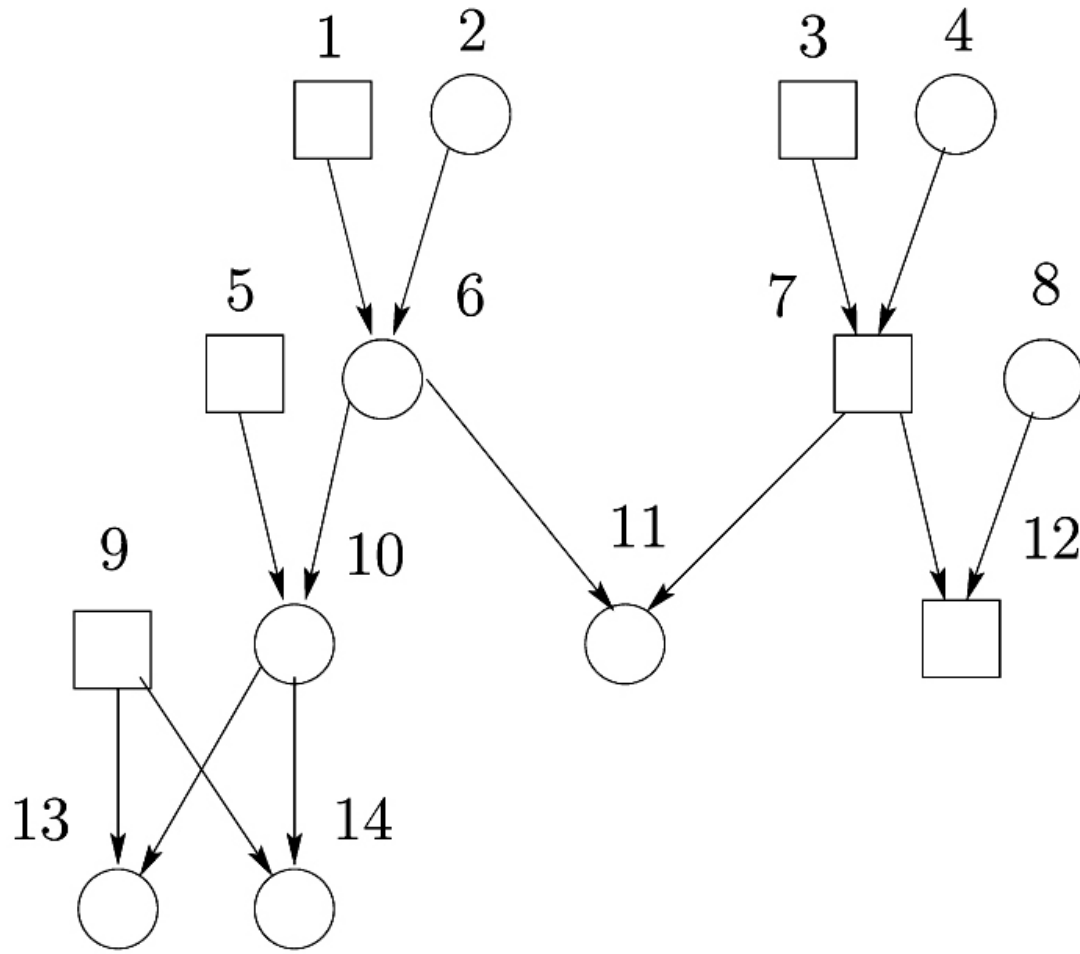
- Phase 2 : **Propagation**

Idée : appliquer l'algorithme de propagation des polyarbres (version simplifiée au cas d'un arbre) sur l'arbre de clusters

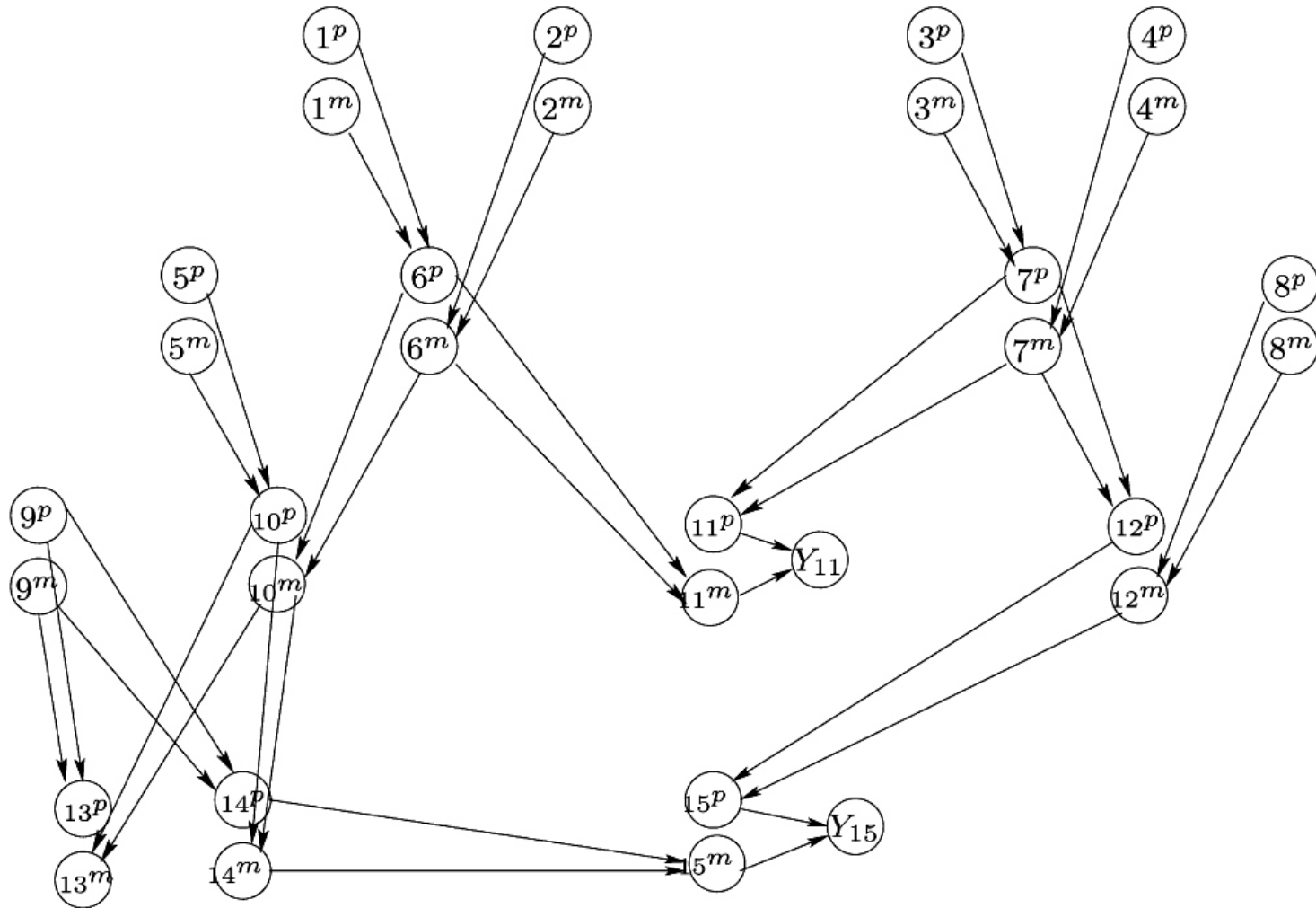
# Simple pedigree example



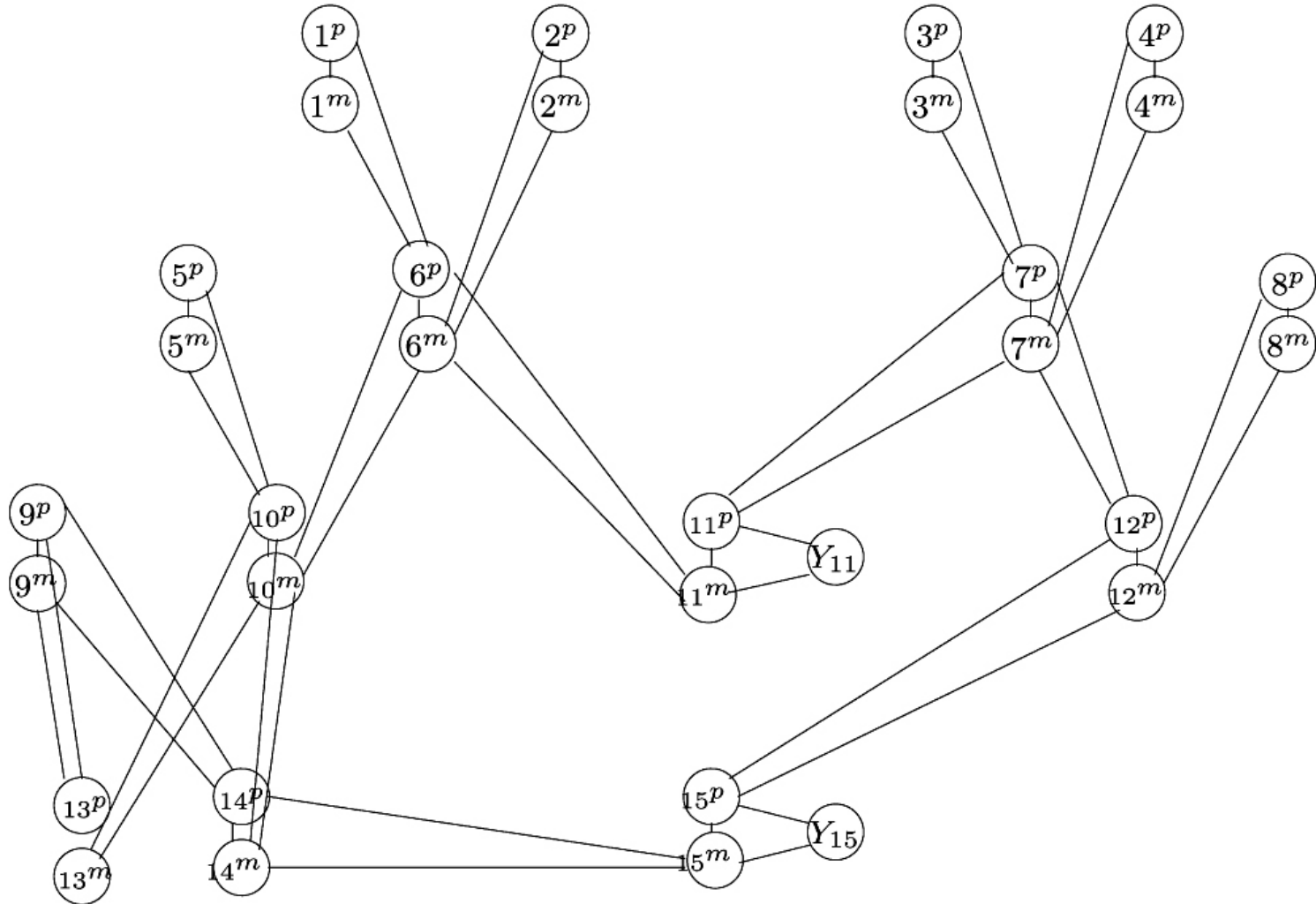
# Genotype network



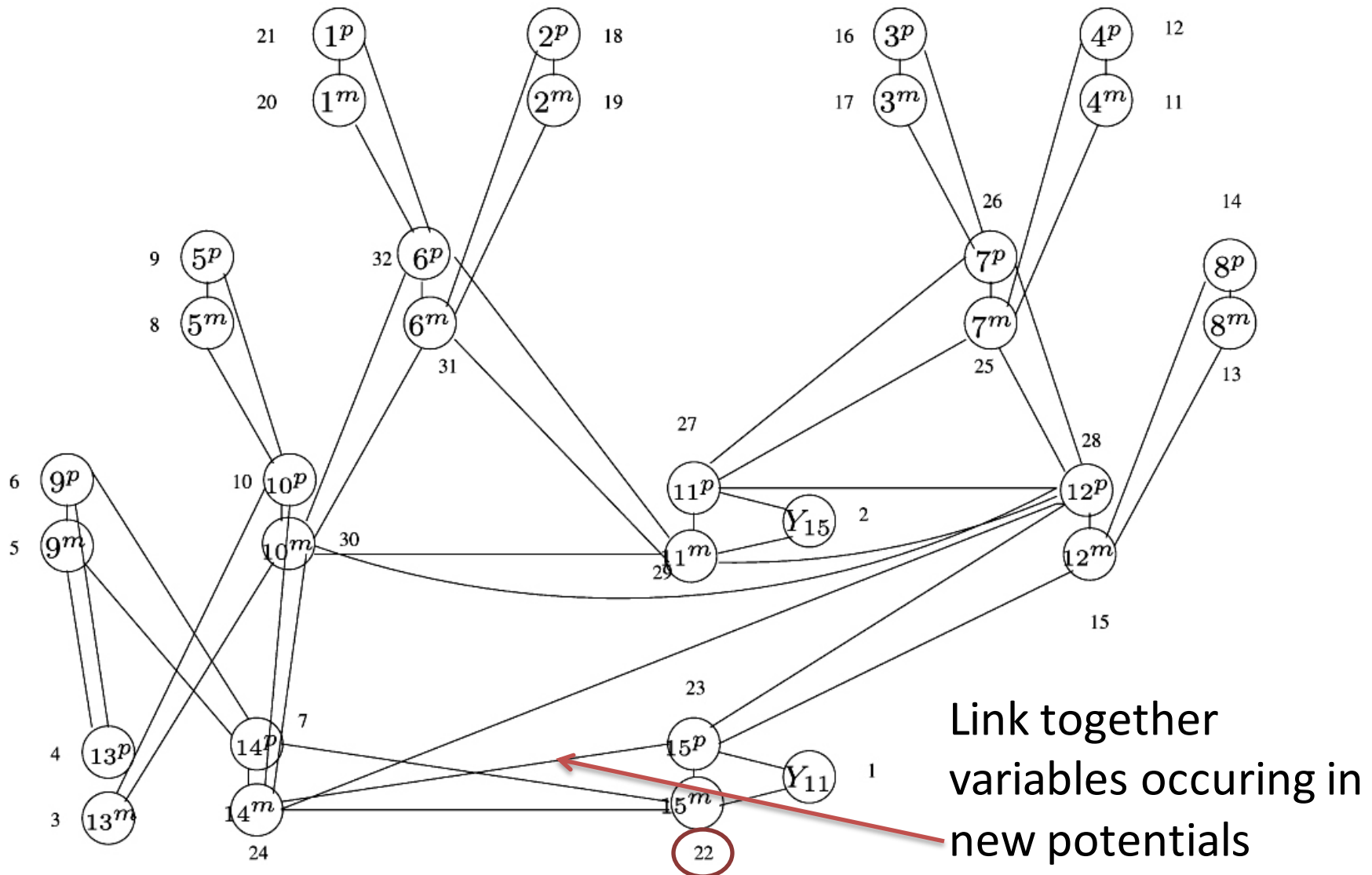
# Allele network with phenotypic information on individual 11 and new individual 15



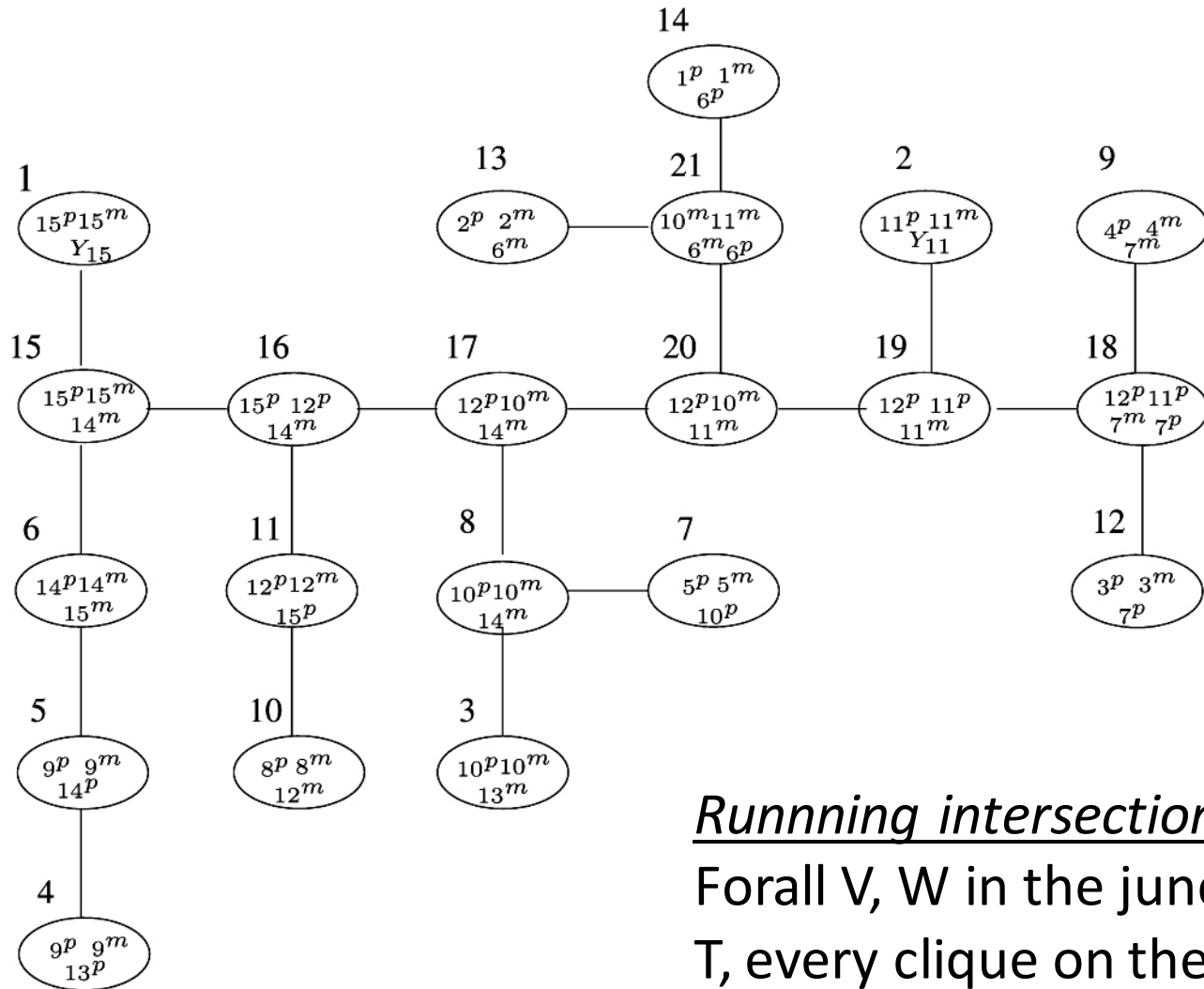
# Primal graph (*moral graph*)



# Triangulated graph (by following an elimination order)



Junction tree (*find all maximal cliques, build a tree from clique intersections*)



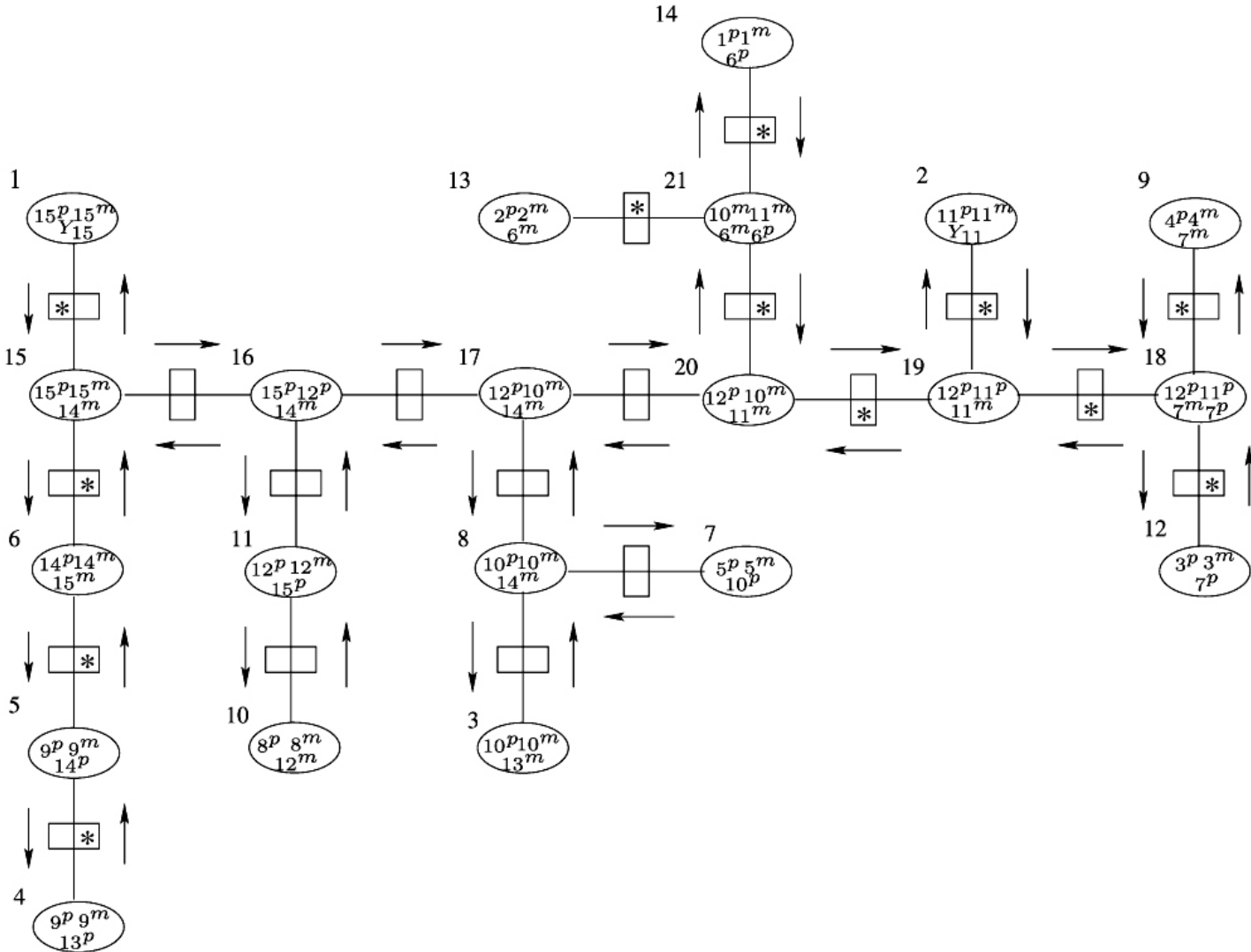
Running intersection property  
 For all  $V, W$  in the junction tree  $T$ , every clique on the path from  $V$  to  $W$  in  $T$  must contain  $V \cap W$

# Assignments of potentials

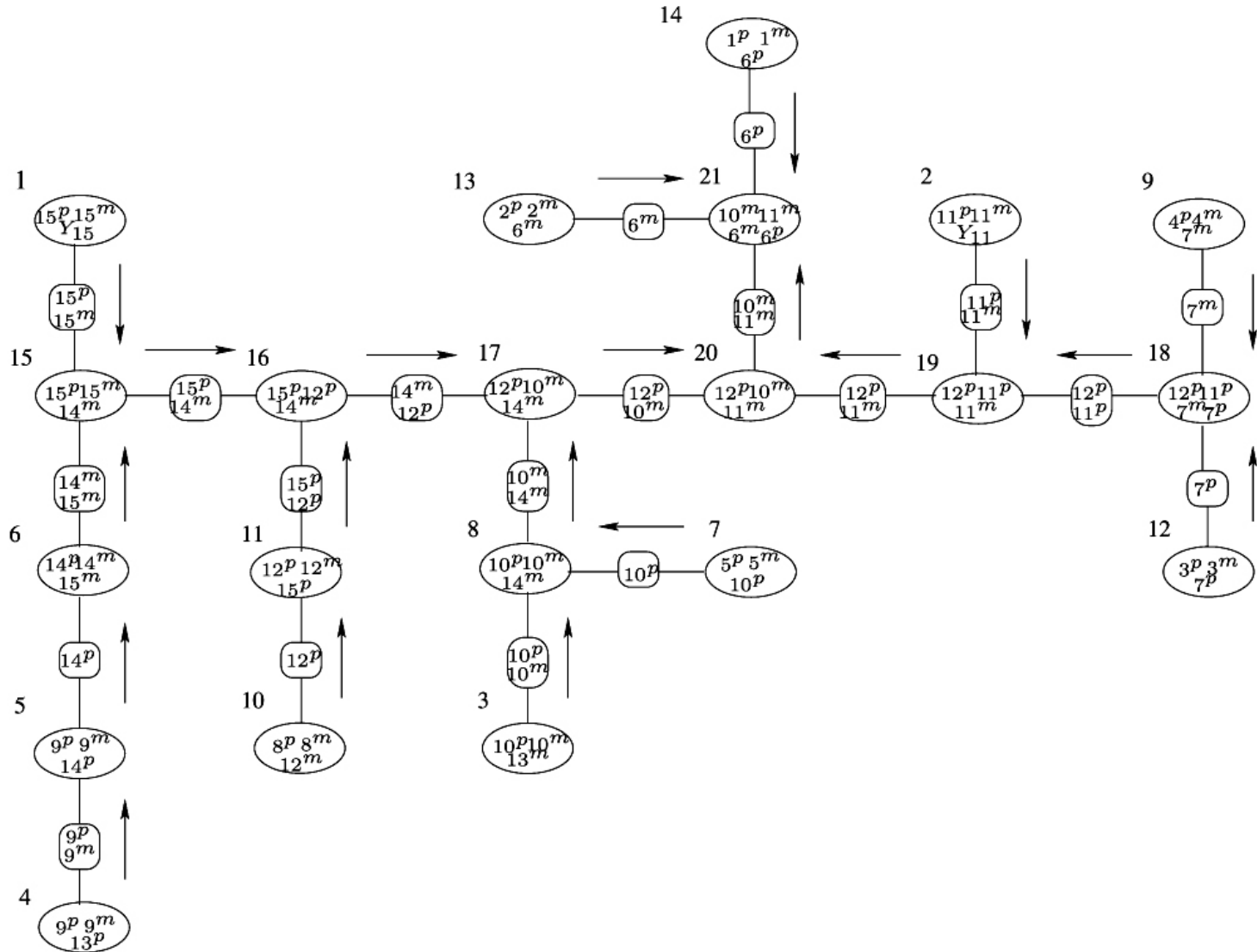
Number	Elements	Assignments	Potential
1	$15^P, 15^m, Y_{15}$	$Y_{15}$	$f(Y_{15} 15^P, 15^m)$
2	$11^P, 11^m, Y_{11}$	$Y_{11}$	$f(Y_{11} 11^P, 11^m)$
3	$10^P, 10^m, 13^m$	$13^m$	$f(13^m 10^P, 10^m)$
4	$9^P, 9^m, 13^P$	$13^P$	$f(13^P 9^P, 9^m)$
5	$9^P, 9^m, 14^P$	$14^P, 9^m, 9^P$	$f(14^P 9^P, 9^m) f(9^P) f(9^m)$
6	$14^P, 14^m, 15^m$	$15^m$	$f(15^m 14^P, 14^m)$
7	$5^P, 5^m, 10^P$	$10^P, 5^m, 5^P$	$f(10^m 5^P, 5^m) f(5^P) f(5^m)$
8	$10^P, 10^m, 14^m$	$14^m$	$f(14^m 10^P, 10^m)$
9	$4^P, 4^m, 7^m$	$7^m, 4^m, 4^P$	$f(7^m 4^P, 4^m) f(4^P) f(4^m)$
10	$8^P, 8^m, 12^m$	$12^m, 8^m, 8^P$	$f(12^m 8^P, 8^m) f(8^P) f(8^m)$
11	$12^P, 12^m, 15^P$	$15^P$	$f(15^P 12^P, 12^m)$
12	$3^P, 3^m, 7^P$	$7^P, 3^m, 3^P$	$f(7^P 3^P, 3^m) f(3^P) f(3^m)$
13	$2^P, 2^m, 6^m$	$6^m, 2^m, 2^P$	$f(6^m 2^P, 2^m) f(2^P) f(2^m)$
14	$1^P, 1^m, 6^P$	$6^P, 1^m, 1^P$	$f(6^P 1^P, 1^m) f(1^P) f(1^m)$
15	$15^P, 15^m, 14^m$		1
16	$15^P, 12^P, 14^m$		1
17	$12^P, 10^m, 14^m$		1
18	$12^P, 11^P, 7^m, 7^P$	$11^P, 12^P$	$f(11^P 7^P, 7^m) f(12^P 7^P, 7^m)$
19	$12^P, 11^P, 11^m$		1
20	$12^P, 10^m, 11^m$		1
21	$10^m, 11^m, 6^m, 6^P$	$11^m, 10^m$	$f(11^m 6^P, 6^m) f(10^m 6^P, 6^m)$



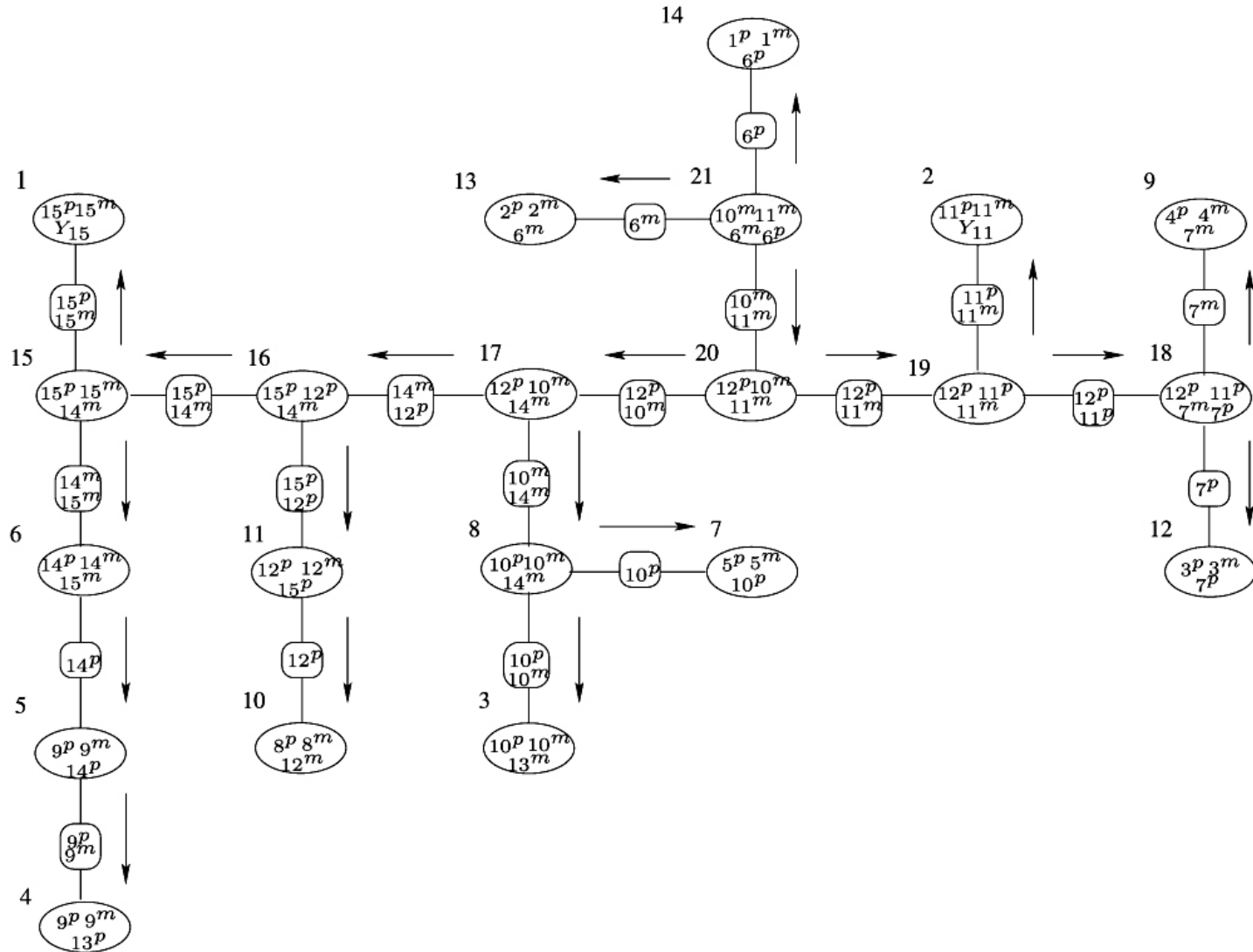
# Propagation in the junction tree



# Collect evidence



# Distribute evidence



# Algorithme de propagation dans l'arbre de jointure

## Théorème

$$p(V_i, \mathbf{e}) = \prod \phi_i \prod \psi^{sep(V_i)} \prod_{V_j \in \Xi(V_i)} \prod \psi_{sep(V_j)}$$

$$p(S_k, \mathbf{e}) = \prod \psi_k \prod \psi^k$$

*Remarque : choix du plus petit ensemble  $V_i$  ou  $S_k$  contenant  $X$  pour calculer  $p(X, \mathbf{e})$ . Choix du plus petit  $S_k$  pour calculer  $p(\mathbf{e})$ .*

## Complexité

Temps  $O(n^2 d^{w+1})$ , espace  $O(nsd^s)$  avec  
 $n = |\mathbf{A}|$ ,  $d = \max_j \text{dom}(A_j)$ ,  $w = \max_i |V_i| - 1$ ,  $s = \max_k |S_k|$

*Remarque :  $w$  est appelé **largeur d'arbre** (tree-width).*

*Dans le cas d'un polyarbre,  $w = r$ ,  $s = 1$ , avec  $r = \max_j |\Pi(A_j)|$ .*

# Graphical model formalisms

## Graphical Model (GM)

A GM is a triplet  $(\mathcal{X}, \mathcal{D}, \mathcal{F})$ .

- $\mathcal{X} = \{X_1, \dots, X_n\}$  a set of variables,
- $\mathcal{D} = \{D_{X_1}, \dots, D_{X_n}\}$  a set of finite domains
- $\mathcal{F} = \{f_1, \dots, f_e\}$ , a set of nonnegative functions, each defined over a subset of variables  $\mathbf{S}_i \subseteq \mathcal{X}$  (i.e. the scope)

Probabilistic  
joint distribution :

$$\mathbb{P}(\mathcal{X}) \propto \prod_{i=1}^e f_i(\mathbf{S}_i)$$

Non probabilistic : WCSP

$$\text{score} = \sum_{i=1}^e f_i(\mathbf{S}_i)$$

# Combinatorial Optimization

## Probabilistic

The *Most Probable Explanation* (MPE) problem is to find the most likely assignment to all variables in  $\mathcal{X}$  maximizing  $\mathbb{P}(\mathcal{X})$ .

## Non probabilistic : WCSP

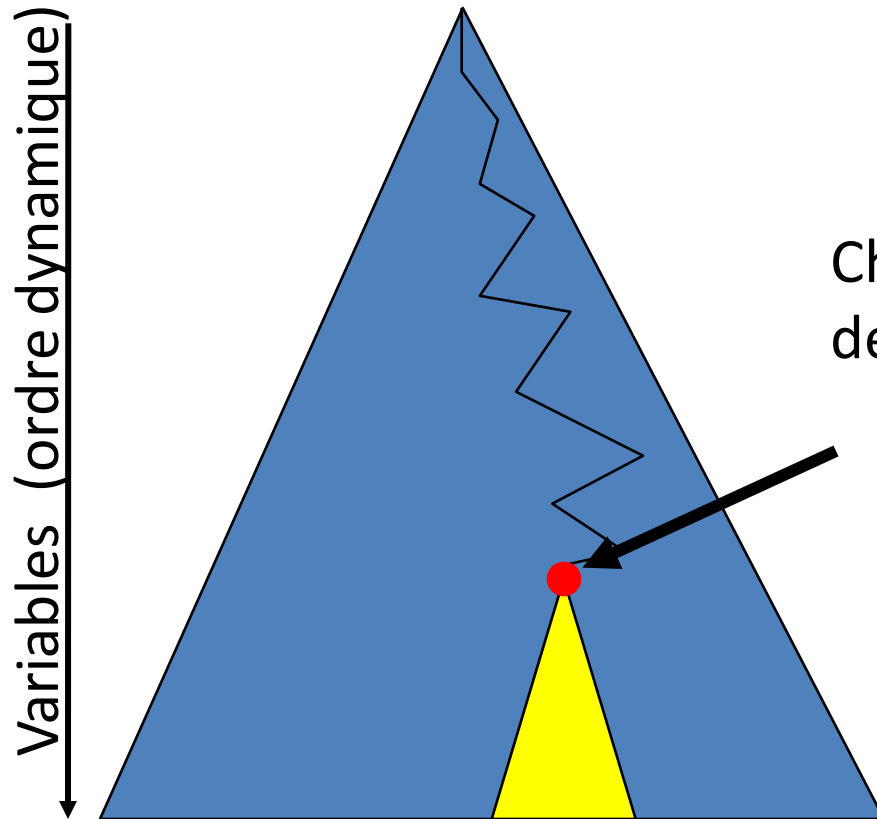
To find a feasible assignment of  $\mathcal{X}$  minimizing score.

## Equivalence

$$\operatorname{argmax}_{\mathcal{X}} \mathbb{P}(\mathcal{X}) = \operatorname{argmin}_{\mathcal{X}} - \sum_{i=1}^e \log(f_i(\mathbf{S}_i))$$

- Inference (Variable Elimination : VE, Join Tree)
- Depth First Branch and Bound
- Hybrids
  - ▶ DFBB + local reasoning : bounded VE + soft local consistencies

# Algorithme de séparation et évaluation



Chaque nœud de l'arbre est un WCSP défini par l'affectation courante

**(LB) Minorant** =  $f_{\emptyset}$

= sous-estimation du coût de la meilleure solution dans le sous-arbre courant

**si  $f_{\emptyset} \geq k$  alors couper**

**(UB) Majorant**

= coût de la meilleure solution trouvée =  $k$

$$\alpha +_k \beta = \min \{k, \alpha + \beta\}$$

# Local operators

Projection :  $f[\mathbf{S}']$

$\mathbf{S}' \subseteq \mathbf{S}$  and  $\forall t' \in D_{\mathbf{S}'}, f[\mathbf{S}'](t') = \min_{t \in D_{\mathbf{S}}} \text{ s.t. } t[\mathbf{S}'] = t' f(t)$

X	Y	f
0	0	5
0	1	7
1	0	6
1	1	4

X	f[X]
0	5
1	4

Join :  $f = f_1 + f_2$

$f(t) = f_1(t[\mathbf{S}_1]) + f_2(t[\mathbf{S}_2]), \forall t \in D_{\mathbf{S}_1 \cup \mathbf{S}_2}$

$i$ -bounded Variable Elimination (VE( $i$ )) :

$X$  such that  $|\bigcap_{f_j: X \in S_j} S_j| < i$   $(\sum_{f_j: X \in S_j} f_j)$   $[\bigcap_{f_j: X \in S_j} S_j \ X]$



# Local operators

Projection :  $f[\mathbf{S}']$

$\mathbf{S}' \subseteq \mathbf{S}$  and  $\forall t' \in D_{\mathbf{S}'}, f[\mathbf{S}'](t') = \min_{t \in D_{\mathbf{S}}} \text{ s.t. } t[\mathbf{S}'] = t' f(t)$

X	Y	f
0	0	5
0	1	7
1	0	6
1	1	4

X	f[X]
0	5
1	4

Join :  $f = f_1 + f_2$

$f(t) = f_1(t[\mathbf{S}_1]) + f_2(t[\mathbf{S}_2]), \forall t \in D_{\mathbf{S}_1 \cup \mathbf{S}_2}$

Subtraction :  $f = f_1 - f_2$

$\mathbf{S}_2 \subseteq \mathbf{S}_1$  and  $f(t) = f_1(t) - f_2(t[\mathbf{S}_2]), \forall t \in D_{\mathbf{S}_1}$

soft Directed Arc Consistency (DAC) [Cooper, *Fuzzy Set Sys.*, 2003] :

$f_1(X) \leftarrow f_1(X) + (f(X, Y) + f_2(Y))[X] \quad X < Y$

$f(X, Y) \leftarrow (f(X, Y) + f_2(Y)) - (f(X, Y) + f_2(Y))[X] \quad f(Y) \leftarrow 0$

- toulbar2 v0.5 avec last conflict
- Temps CPU en secondes pour trouver et prouver l'optimum

# MENDEL

## DFBB-VE(2)

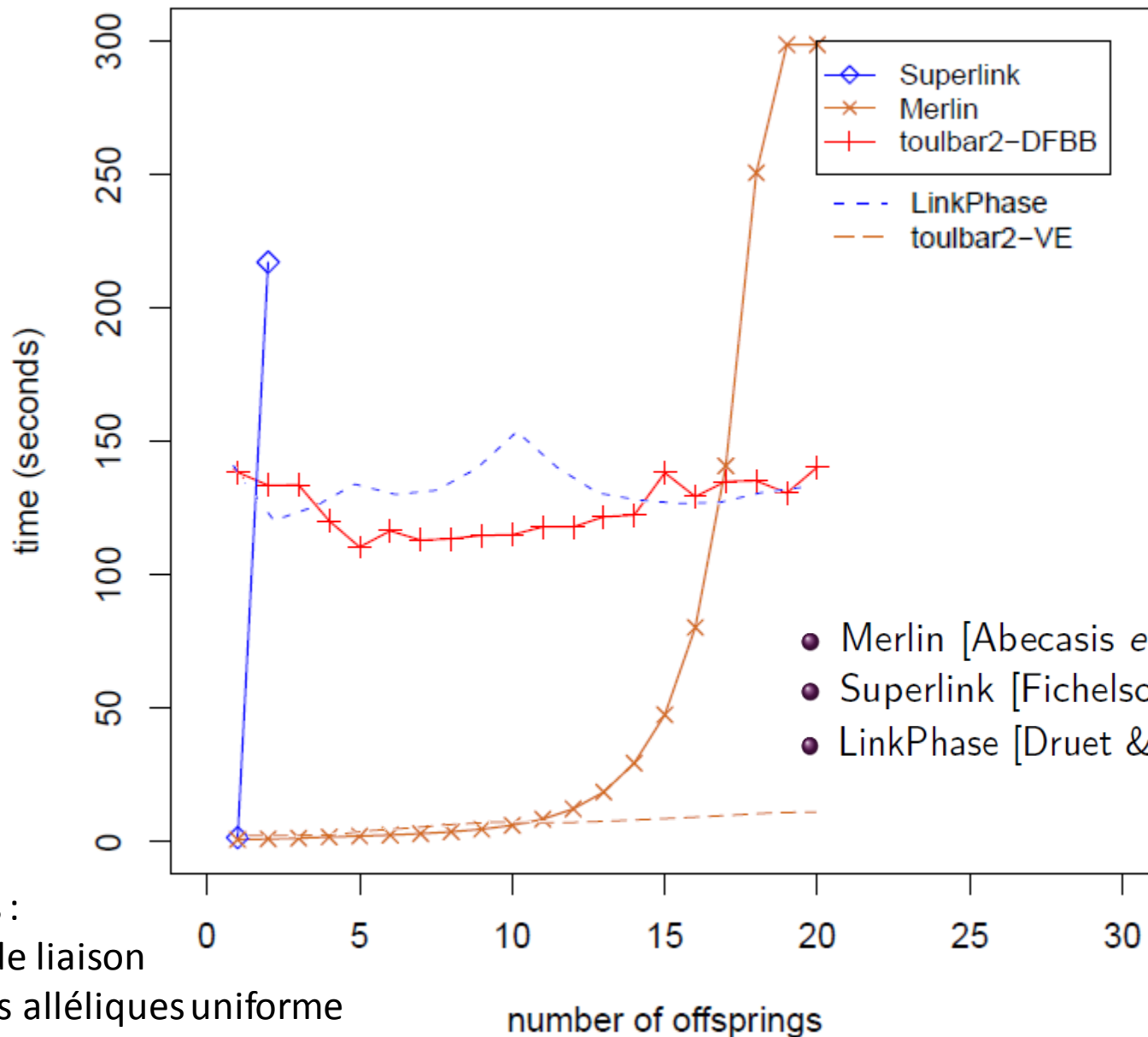
	ind	vars	genotyped	alleles	nf	ngen	treewidth ub	errors	time	nodes
<i>eye</i>	36	36	28	6	11	4	2	1	0.02	0
<i>cancer</i>	49	48	37	8	18	5	2	1	0.21	0
<i>parkinson</i>	37	34	13	4	7	7	5	0	0	6
<i>berrichon<sub>1nc</sub></i>	129516	9947	2448	4	8821	17	262	2	4.73	8805
<i>berrichon<sub>1</sub></i>	129516	10017	2483	4	8786	17	330	23	5.81	8384
<i>berrichon<sub>2nc</sub></i>	27255	19337	10215	4	4719	19	-	41	5.89	6170
<i>berrichon<sub>2</sub></i>	27255	19562	10215	4	2381	19	-	106	17.23	15445
<i>langlade<sub>1</sub></i>	1355	1209	711	9	298	13	84	38	12.28	391
<i>langlade<sub>2</sub></i>	1355	1223	715	7	298	13	82	89	60.56	17857
<i>langlade<sub>3</sub></i>	1355	1258	787	5	298	13	85	39	14.19	6731
<i>langlade<sub>4</sub></i>	1355	1186	672	8	298	13	83	43	59.7	3520
<i>moissac<sub>1</sub></i>	283	260	183	2	81	5	6	0	0	5
<i>moissac<sub>2</sub></i>	283	244	167	7	81	5	6	0	0.51	6
<i>moissac<sub>3</sub></i>	283	225	151	3	81	5	6	0	0	4
<i>moissac<sub>4</sub></i>	283	256	179	2	81	5	6	0	0	5
<i>moissac<sub>5</sub></i>	283	237	161	8	81	5	6	0	1.02	5
<i>moissac<sub>6</sub></i>	283	201	131	11	81	5	5	0	5.64	6

# Reconstruction d'haplotypes dans des pedigrees *en arbre* (Favier et al, IJCAI 2011)

Indiv. x Loci	Problème	DFBB-VE( <i>i</i> )		AOBB-C+SMB( <i>j</i> )+VE( <i>i</i> )					
	DFBB-VE( <i>i</i> )	+dec+ps				+dec+ps			
	time (s)	<i>i</i>	time (s)	<i>i</i>	<i>j</i>	<i>i</i>	time (s)	time (s)	
	<b>Linkage</b>								
25x20	ped7	4.04	2	1.18	4	20	4	1915.56	131.14
	ped9	-		3.36	6	20	6	-	104.62
	ped18	149.71	4	3.19	5	20	5	54.67	18.82
	ped19	-		-	-	20	6	-	-
	ped20	3.46	4	0.39	6	16	6	407.6	66.53
	ped23	0.09	4	0.05	3	12	3	6.05	1.52
	ped25	1207.97	4	0.65	6	20	6	25.91	32.51
	ped30	543.20	2	5.44	5	20	5	28.39	10.10
	ped34	1.13	2	0.36	5	20	5	-	24.56
57x6	ped37	0.21	5	0.11	4	10	4	87.18	9.58
	ped39	16.68	4	0.24	5	18	5	6.87	2.60
	ped41	-		302.05	4	20	4	-	1271.31
	ped42	1.94	4	0.34	6	16	6	240.94	155.39
20x20	ped44	-		505.46	5	20	5	2631.71	333.89
	ped50	0.90	4	0.18	4	12	4	316.92	521.92

# Pedigree de demi-frères (Favier et al, WCB'10)

Chromosome X humain, 36000 marqueurs SNP sur 1,64 Morgan(hapmap)



Hypothèses :

- équilibre de liaison
- fréquences alléliques uniforme

- Merlin [Abecasis *et al.*(2002)]
- Superlink [Fichelson *et al.*(2005)]
- LinkPhase [Druet & Georges(2010)]

# markers on 2 Morgan

