# The use of genomic and gene expression large-scale data for the analyses of sexual evolution

## Maria D. Vibranovski

**UNIVERSIDADE DE SÃO PAULO**
**INSTITUTO DE BIOCIÊNCIAS**
**Departamento de Genética e Biologia Evolutiva**

*Instituto de biociências*

colmea

**15 de abril de 2015 - UFRJ**

# Acknowledgments



Manyuan Long
Department of Ecology and Evolution
The University of Chicago



Timothy L. Karr
Drosophila Genetic Resource Center
Kyoto Institute of Technology



Hedibert F. Lopes
INSPER – São Paulo

# Outline

1. Background in Sexual Evolution

2. Biological Problem

3. Large-scale Data

4. Statistical Approaches

5. Ongoing Biological Problems

6. Perspectives
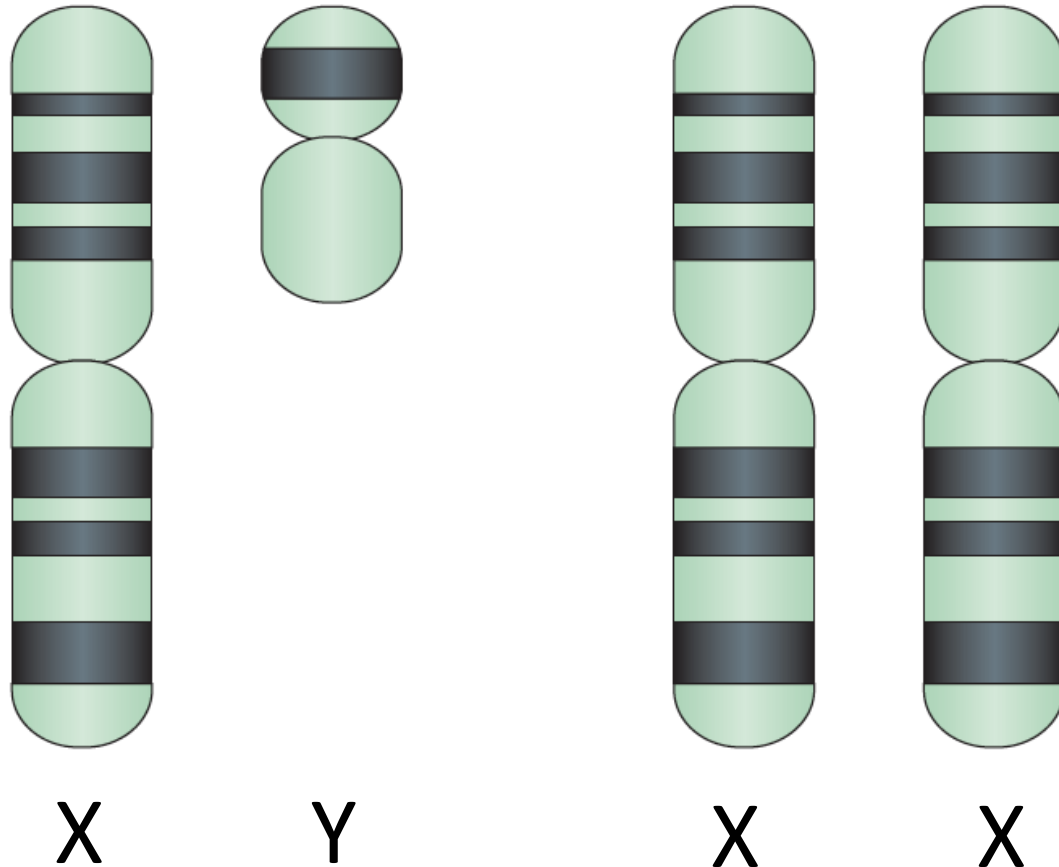
# Why are there phenotypic differences between sexes?

# Sex Chromosomes

Male

Female



X    Y        X    X

# 1. Background in Sexual Evolution

DNA

Nuclear membrane

mRNA Transcription

Mature mRNA

<>

Transport to cytoplasm

Amino acid

Amino acid chain (protein)
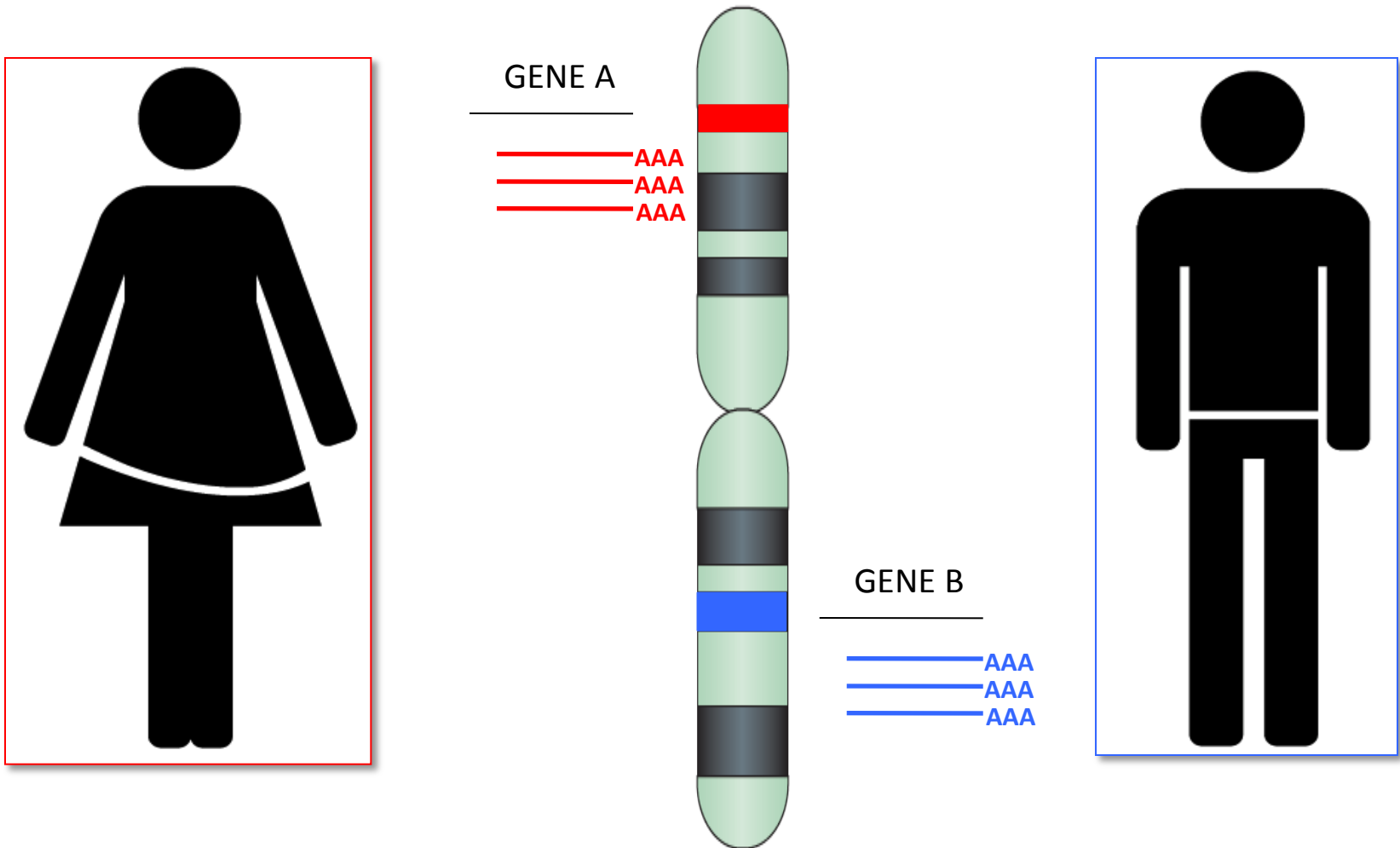
tRNA

Translation

Anti-codon

Codon

mRNA

Ribosome

**Measure the quantity of messenger RNA (mRNA) in the cell**

From Talking Glossary of Genetics

**Gene activity = gene expression = number of mRNA molecules**

# Sex-biased Gene Expression



GENE A

AAA
AAA
AAA

GENE B

AAA
AAA
AAA
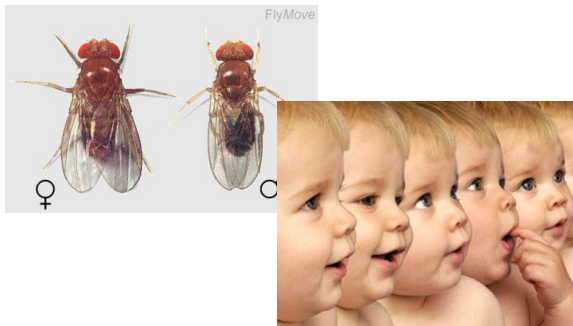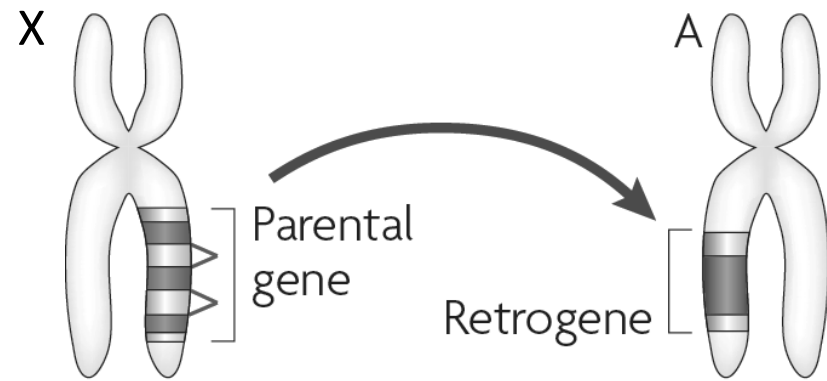
## 2. Biological Problem

What molecular mechanisms and genetics processes are involved?
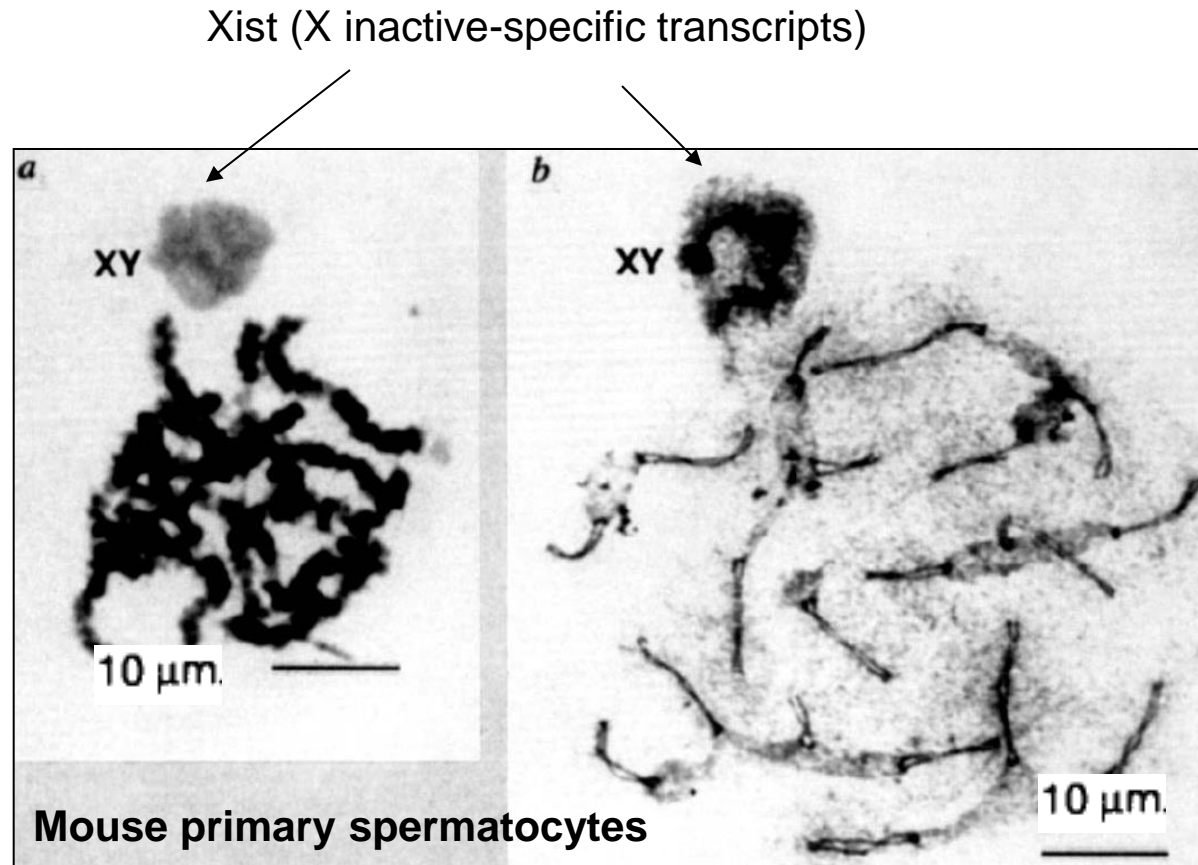
XY systems

X — Parental gene

A — Retrogene
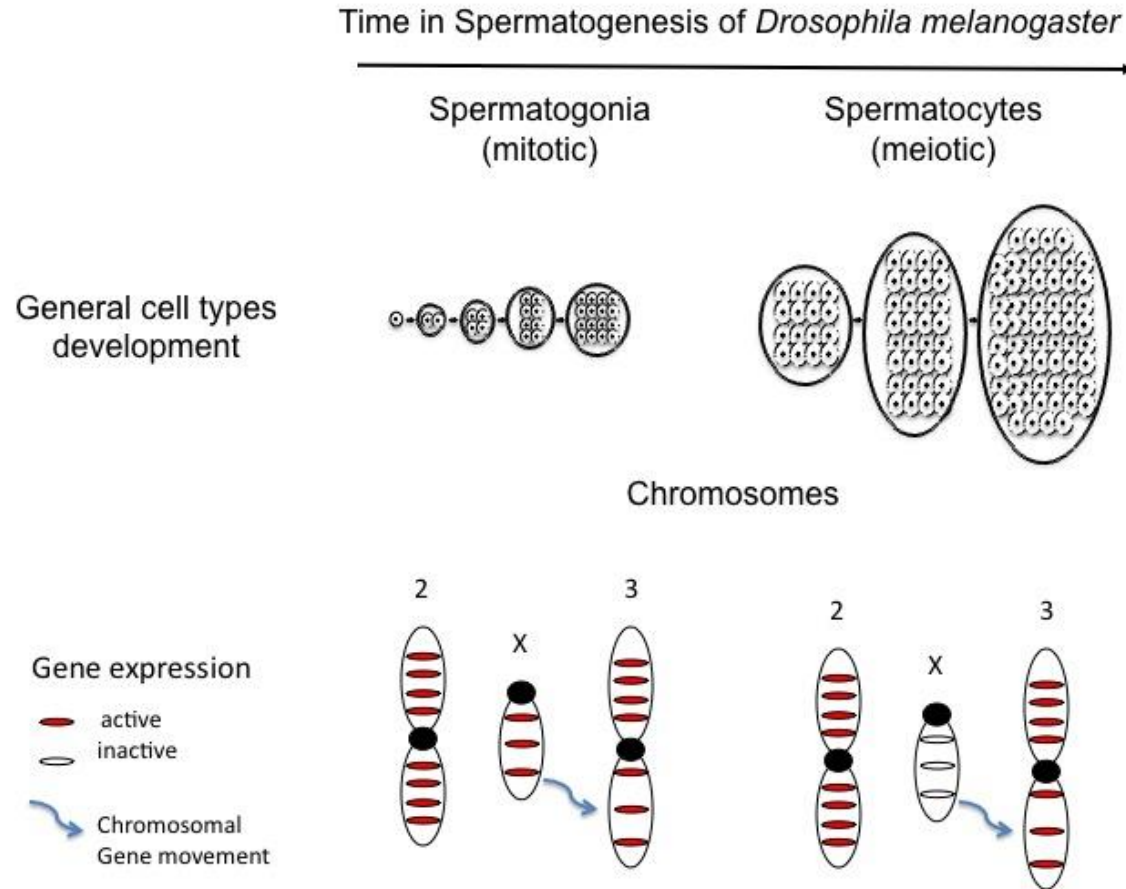
Testis bias

# 2. Biological Problem

Meiotic Sex Chromosome Inactivation (MSCI)

Xist (X inactive-specific transcripts)



**Mouse primary spermatocytes**

Richler C, Soreq H, Wahrman J, 1992. Nat Gene;
Ayoub N, Richler C, Wahrman J, 1997, Chromosoma

# 2. Biological Problem



Time in Spermatogenesis of *Drosophila melanogaster*

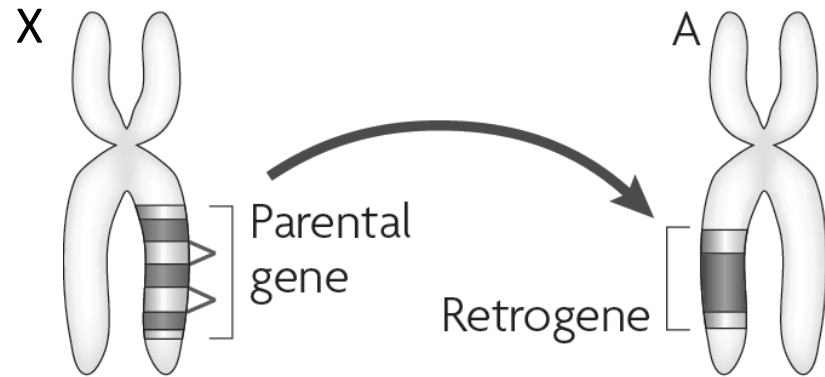Lifschytz and Lindsley, 1972

# 2. Biological Problem



X Y    X X



XY systems
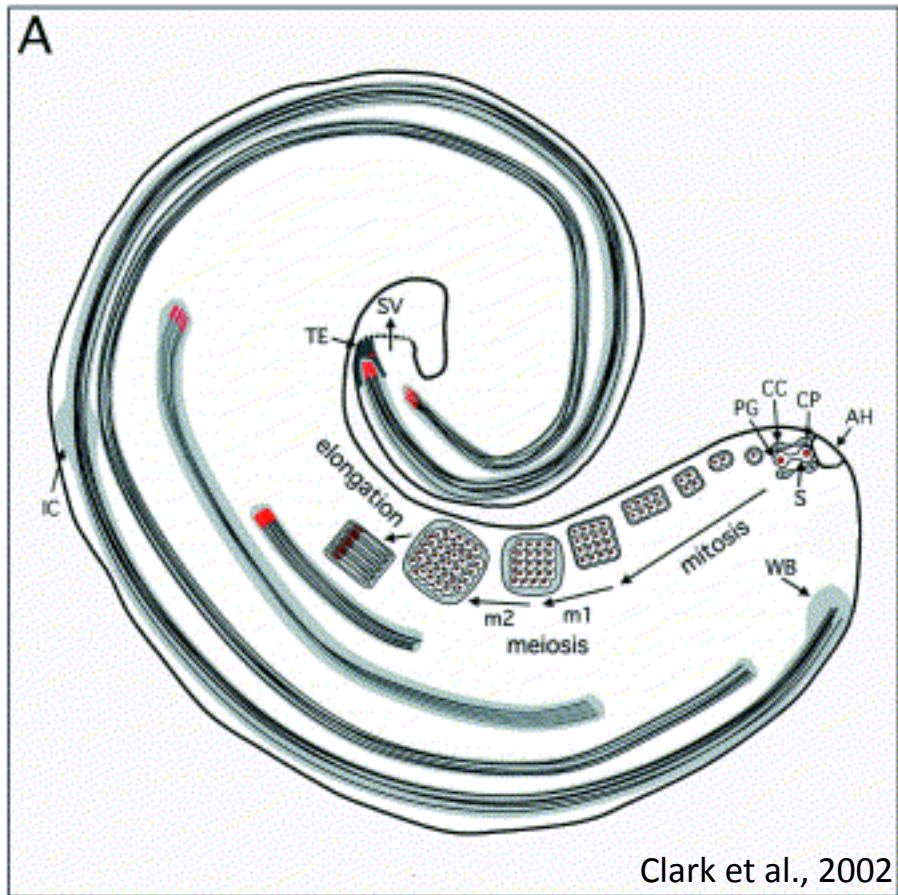
- Does MSCI exist in *Drosophila*?

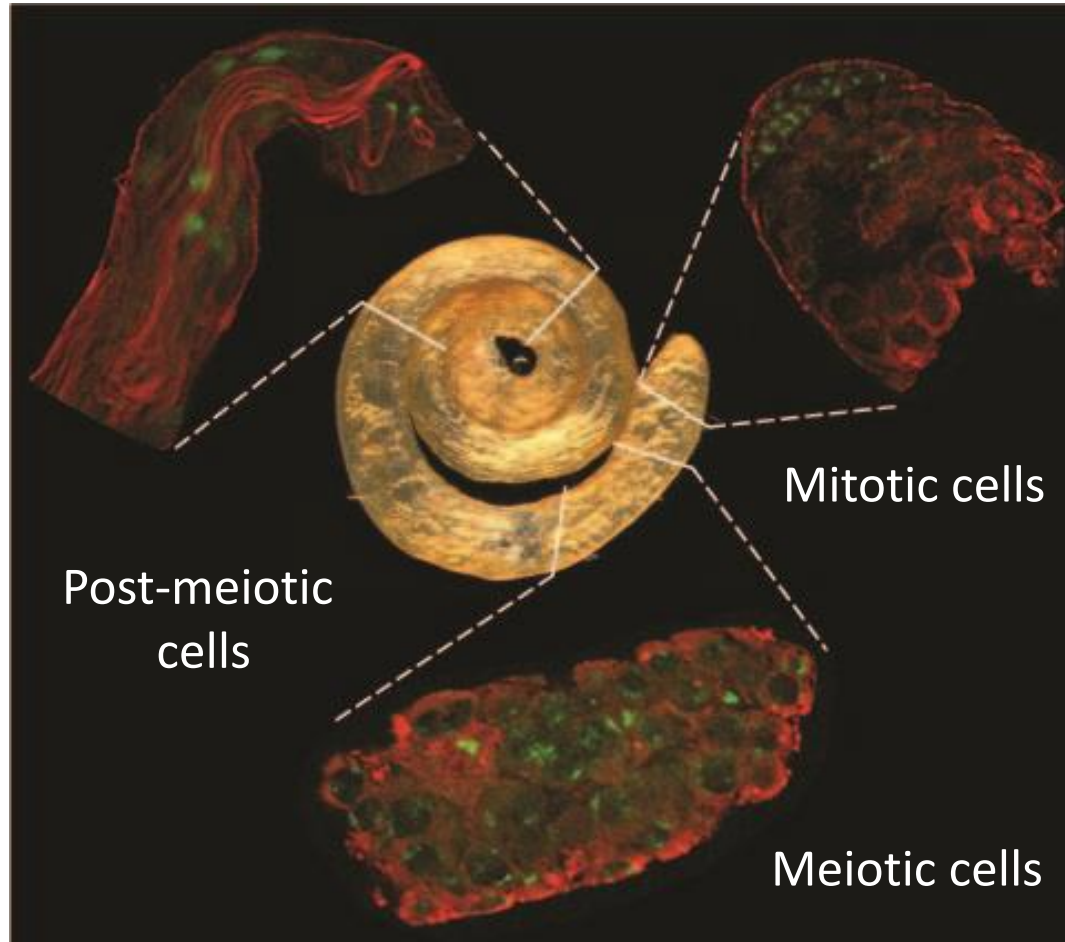- Is MSCI involved in the distribution of sex-biased genes in the genome?



X

A

Parental gene

Retrogene

Testis bias

# 3. Large-scale Data

Drosophila melanogaster Spermatogenesis



Clark et al., 2002

# 3. Large-scale Data

Isolation of Spermatogenic Cells



Post-meiotic cells

Mitotic cells

Meiotic cells

**DNA**
**Cytoplasm**

# 3. Large-scale Data
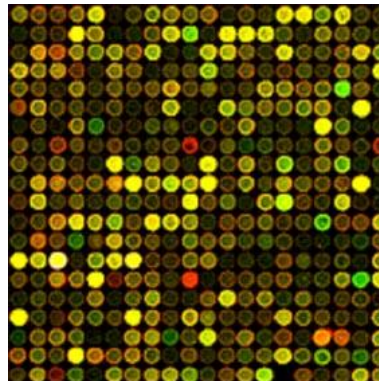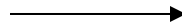
Tissue Isolation (n=3)



RNA
extraction



Microarray
Hybridization

# How to count the number of mRNA molecules for each gene in a cell?

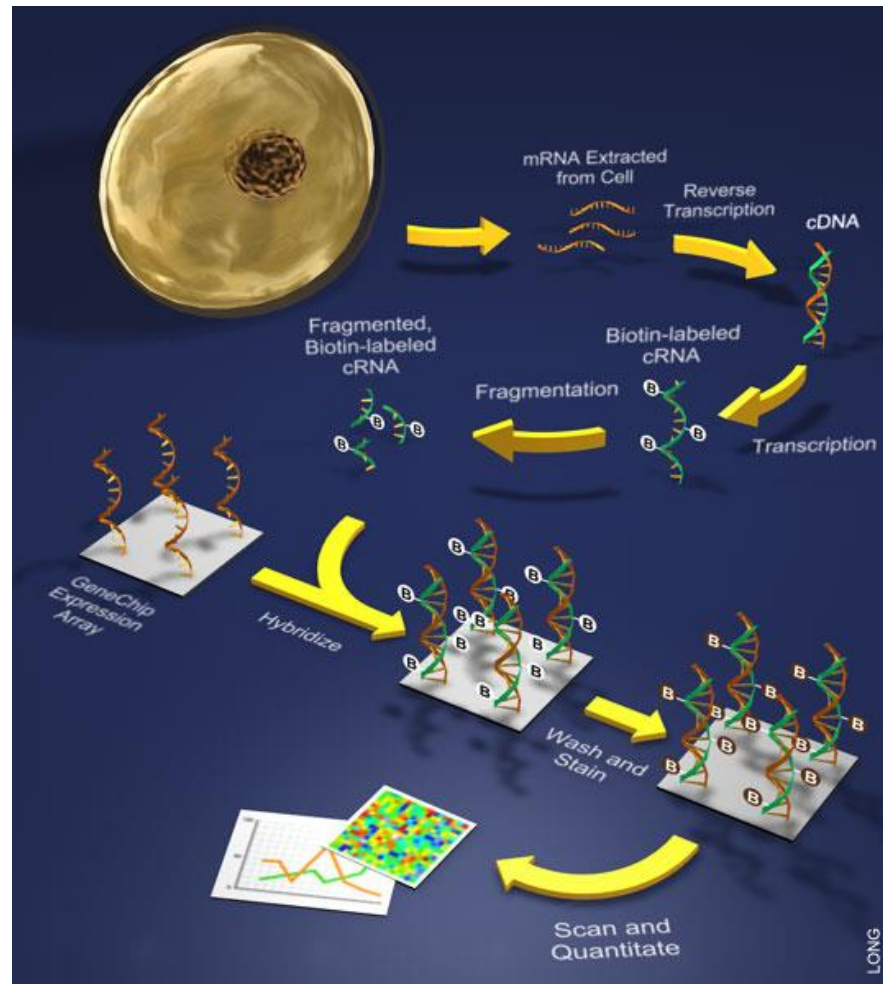**Microarray: A chip containing small pieces of DNA corresponding to all *Drosophila* genes**
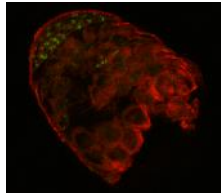
# Microarray

**Specific case for gene expression measure: the GeneChip array**

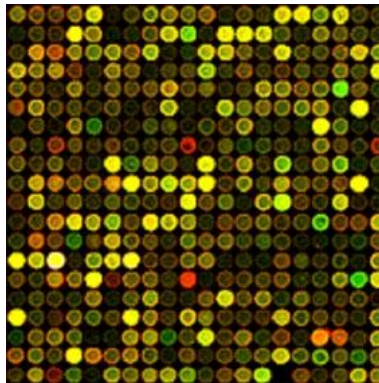# 3. Large-scale Data

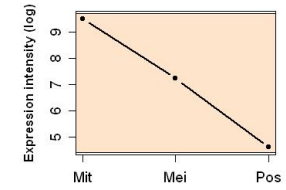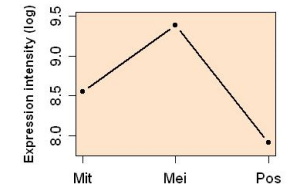Tissue Isolation (n=3)

Gene profile



RNA extraction

Microarray Hybridization

Bam

ms(3)K81

Mst35Bb

# 4. Statistical Approaches

Meiotic Sex Chromosome Inactivation (MSCI)



Post-meiotic cells

Mitotic cells

Meiotic cells

# 4. Statistical Approaches

**X inactivation** occurs when the difference in activity (meiosis-mitosis) in X is lower than the difference in activity in autosome.

$\theta$ = gene expression /activity

$$\theta^X_{meiosis} - \theta^X_{mitosis} < \theta^A_{meiosis} - \theta^A_{mitosis}$$

Gene expression

Mitosis    Meiosis

X inactivation

# 4. Statistical Approaches

**Gene intensity**

  2982 genes (X chromosomes)

  15099 genes (autosomes)

**Spermatogenic replicates**

  3 mitotic

  3 meiotic

  3 post-meiotic



| Gene | Mit1 | Mit2 | Mit3 | Mei1 | Mei2 | Mei3 |
|------|------|------|------|------|------|------|
| 292 | 10.182 | 10.199 | 10.395 | 9.798 | 9.880 | 9.862 |
| 792 | 5.273 | 5.357 | 5.509 | 5.548 | 5.577 | 5.587 |
| 1966 | 5.924 | 5.914 | 5.945 | 6.297 | 6.779 | 6.643 |
| 2811 | 5.177 | 5.197 | 5.168 | 5.092 | 5.166 | 5.169 |

| Gene | Mit1 | Mit2 | Mit3 | Mei1 | Mei2 | Mei3 |
|------|------|------|------|------|------|------|
| 2869 | 9.884 | 9.872 | 9.758 | 11.468 | 11.174 | 11.213 |
| 3939 | 4.833 | 4.755 | 4.775 | 4.827 | 4.839 | 4.804 |
| 10541 | 4.932 | 4.932 | 4.969 | 5.073 | 4.876 | 4.986 |
| 12928 | 10.331 | 10.485 | 10.524 | 11.480 | 11.628 | 11.514 |

# 4. Statistical Approaches

For chromosome $i$ (X or A) and gene $l$, 3 replicates are measured

Mitosis: $mit_{i1l}, mit_{i2l}, mit_{i3l}$ $\qquad$ $E(mit_{ikl}) = \theta_{il}^{mit}$

Meiosis: $mei_{i1l}, mei_{i2l}, mei_{i3l}$ $\qquad$ $E(mei_{ikl}) = \theta_{il}^{mei}$

The objective is to learn whether genes are
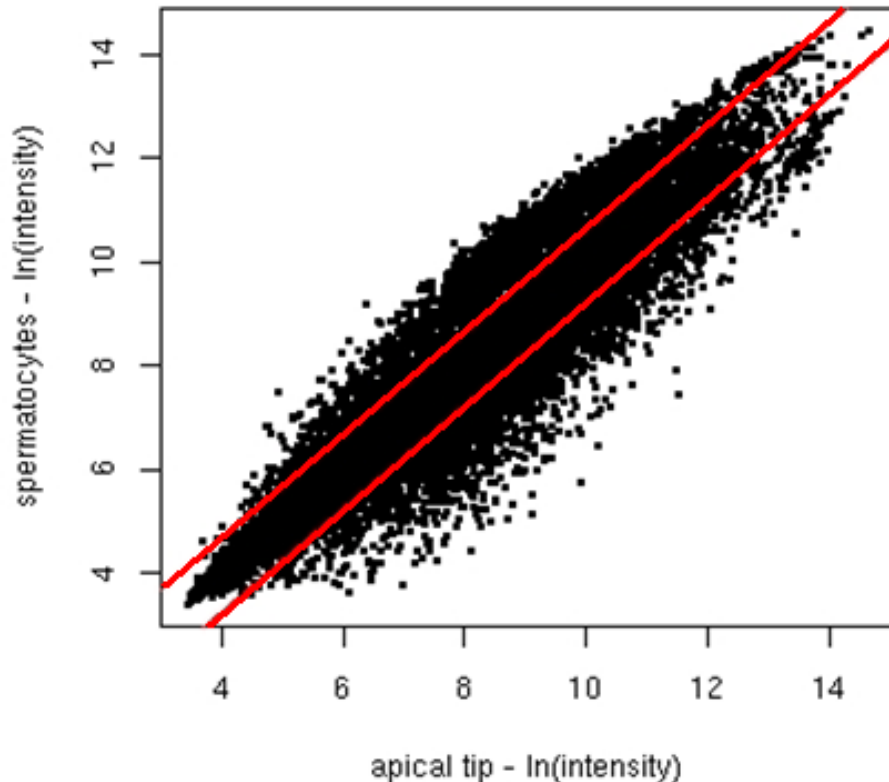
Differently expressed: $\begin{cases} H_{1,il} : \theta_{il}^{mei} > \theta_{il}^{mit} & \text{(OVER)} \\ H_{2,il} : \theta_{il}^{mei} < \theta_{il}^{mit} & \text{(UNDER)} \end{cases}$

or

Equally expressed: $\qquad$ $H_{3,il} : \theta_{il}^{mei} = \theta_{il}^{mit}$ $\quad$ (EQUAL).

# 4. Statistical Approaches

**How do biologists usually approach this problem?**



**APPROACH 1:**
**Identify differently expressed genes based on 2 fold intensity differences (based on homotypic experiments);**

**APPROACH 2:**
**Identify differently expressed genes by controlling type I error;**

**APPROACH 3:**
**Identify differently expressed genes by controlling false discovery rates;**

**APPROACH 4:**
**Hierarchical Bayes.**

*Storey and Tibshirani (2003) Statistical significance for genomewide studies. PNAS, 100(16), 9440-9445.*

# 4. Statistical Approaches

**How do biologists usually approach this problem?**

<span style="color:red">**APPROACH 1, 2 or 3**</span>
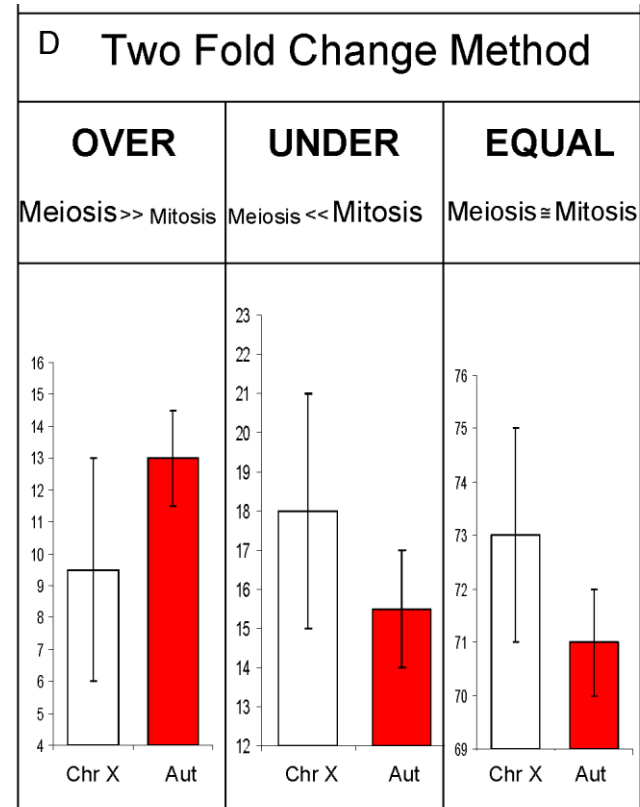**Identify differently expressed genes**

↓

## Counts

| Expression | ChrX | ChrA |
|---|---|---|
| Under expressed | 537(18%) | 2416(16%) |
| Equally expressed | 2147(72%) | 10720(71%) |
| Over expressed | 298(10%) | 1963(13%) |
| Total | 2982 | 15099 |

# 4. Statistical Approaches



Gene expression

Mitosis    Meiosis

↑

X inactivation

**Proportion of genes (%)**

D  Two Fold Change Method

| OVER | UNDER | EQUAL |
|---|---|---|
| Meiosis >> Mitosis | Meiosis << Mitosis | Meiosis ≅ Mitosis |

# 4. Statistical Approaches

## Approach 2: *False positive rate (FPR)*

|  | (E) Evidence against H | (NE)No evidence against H | Total |
|---|---|---|---|
| (A) Over/under expressed | $\tau$ | $\eta_1 - \tau$ | $\eta_1$ |
| (H) Equally expressed | $\eta$ | $\eta_0 - \eta$ | $\eta_0$ |
| Total | x | n-x | n |

Measures error rate (evidence against H) when H is true: $E(\eta / \eta_0)$

FPR=0.05 => $E(\eta) < 0.05 E(\eta_0)$

Bonferroni correction
When evidence is measure by {p-value<0.05/n}, then

$$Pr(\eta > 0) < 0.05$$

Controlling $E(\eta) < 0.05n$   is too liberal!
Controlling $Pr(\eta > 0) < 0.05$ is too conservative!

# 4. Statistical Approaches

## Approach 3: *False discovery rate (FDR)*

|  | (E) Evidence against H | (NE)No evidence against H | Total |
|---|---|---|---|
| (A) Over/under expressed | $\tau$ | $\eta_1 - \tau$ | $\eta_1$ |
| (H) Equally expressed | $\eta$ | $\eta_0 - \eta$ | $\eta_0$ |
| Total | x | n-x | n |

$$FDR(\alpha) = E\left(\frac{\eta(\alpha)}{x(\alpha)}\right) \approx \frac{E(\eta(\alpha))}{E(x(\alpha))} = \frac{\eta_0 \alpha}{E(x(\alpha))} \approx \frac{n\hat{\pi}_0 \alpha}{\sum_{i=1}^{n} 1(p_i < \alpha)}$$

$$\hat{\pi}_0 = \lim_{\lambda \to 1} \frac{\sum_{i=1}^{n} 1(p_i > \lambda)}{n(1-\lambda)} \qquad \text{(estimate of true nulls)}$$

$$FDR(\alpha) \approx \frac{\pi_0 \Pr(E \mid H)}{\Pr(E)} = \Pr(H \mid E) \qquad \text{(Looks Bayesian!!)}$$

## Approach 4: mixture model for differences

Because $\sigma_X^2$ and $\sigma_A^2$ (within gene variability) in the hierarchical model are negligible when compared $\tau^2$ (between genes variability), the next model we implemented is a mixture of normals model that accounts for the excess of intensities around zero:

Model:

$$d_l = (\bar{x}_l^{mei} - \bar{x}_l^{mit}) \sim \pi N(\theta_1, \tau_1^2) + (1-\pi)N(\theta_2, \tau_2^2) \qquad l=1,...,n$$

Prior:

$$p(\pi, \theta_1, \theta_2, \tau_1^2, \tau_2^2) \propto \tau_1^{-2}\tau_2^{-2}$$

# Approach 4: <span style="color:blue">mixture</span> model for differences

• Bayesian Mixture Model classifies genes as over or under expressed

- Detect differential gene expression while comparing chromosomes distributions

- Avoids the arbitrariness of folding

- Avoids multiple testing distortions

EQUAL

UNDER                    OVER

Lopes, Mueller and Rosner (2003)

## Approach 4: mixture model for differences

$$p(d \mid data) = \int \left( \pi f_N(d; \theta_1, \tau_1^2) + (1-\pi) f_N(d; \theta_2, \tau_2^2) \right) p(\pi, \theta_1, \theta_2, \tau_1^2, \tau_2^2 \mid data) d\pi d\theta_1 d\theta_2 d\tau_1^2 d\tau_2^2$$



LOG(meiosis/Mitosis)

## Approach 4: <span style="color:blue">mixture</span> model for differences

Classification (Prior):         $Z_l \sim Ber(\pi)$

Classification (Posterior):     $(Z_l = 1) \, \& \, (d_l > 0)$ — over expressed

$(Z_l = 1) \, \& \, (d_l < 0)$ — under expressed

$(Z_l = 0)$ — equally expressed

## Approach 4: Bayesian mixture model for differences

# 4. Statistical Approaches

Meiotic Sex Chromosome Inactivation (MSCI)



Vibranovski et al, 2009, PLoS Genetics

# 4. Statistical Approaches

Meiotic Sex Chromosome Inactivation (MSCI)

Vibranovski et al, 2009, PLoS Genetics

## 2. Biological Problem

Vibranovski et al, 2009, PLoS Genetics

# 2. Biological Problem



X Y    X X



XY systems

- Does MSCI exist in *Drosophila*? YES

- Is MSCI involved in the distribution of sex-biased genes in the genome?



X

Parental gene

A

Retrogene

Testis bias

36

## 2. Biological Problem

Out-of-the-X retroposition

X

A

Parental gene

Retrogene

Parental gene and retrogene transcription levels during distinct stages of spermatogenesis

~ Parental gene mRNA
~ Retrogene mRNA

Spermatogonia (Mitotic)

Spermatocytes (Meiotic)

Intensity of sex-chromosome inactivation during and after meiosis

37

Kaessmann et al, 2009, Nat Rev Genetics

# 2. Biological Problem

Complementary expression

Tom40: Translocase of outer membrane 40 / Chr X
Tomboy40: protein transmembrane transporter activity / Chr 2R

## 2. Biological Problem

Complementary expression

No Complementary expression



Expression Intensity

Mitosis          Meiosis

Parental gene

Retrogene gene

Mitosis          Meiosis

Out of X

Out of A

# 4. Statistical Approaches

## Approach 1: Counting YESES and NOS

| Counts | NO | YES |     | Proportions | NO | YES |
|--------|----|-----|-----|-------------|-----|-----|
| Out of X | 9 | 20 |     | Out of X | 0.31 | 0.69 |
| Out of A | 28 | 34 |     | Out of A | 0.45 | 0.55 |

**Fisher's Exact Test for Count Data**
Null hypothesis: true odds ratio = 1
Alternative hypothesis: true odds ratio is not equal to 1
p-value = 0.2546
Estimated odds ratio=0.55  (roughly (34/28)/(20/9)=0.546)
95% confidence interval: (0.1888,1.5086)

95% confidence interval for YES (normal approximation)
Out of X: (0.518,0.862)
Out of A: (0.424,0.676)

# 4. Statistical Approaches

## Approach 2: Controling for FDR

Considering only significant change of expression 2 fold difference p<0.05, q < 0.05

| Counts | NO | YES | | Proportions | NO | YES |
|---|---|---|---|---|---|---|
| Out of X | 9 | 18 | | Out of X | 0.33 | 0.66 |
| Out of A | 37 | 25 | | Out of A | 0.60 | 0.40 |

**Fisher's Exact Test for Count Data**
Null hypothesis: true odds ratio = 1
Alternative hypothesis: true odds ratio is not equal to 1
p-value = 0.03687
Estimated odds ratio=0.342 (roughly (25/37)/(18/9)=0.339)
95% confidence interval: (0.11,0.95)

95% confidence interval for YES (normal approximation)
Out of X: (0.478,0.842)
Out of A: (0.276,0.524)

# 4. Statistical Approaches

## Approach 3: Bayesian hierarchical model

$$\begin{pmatrix} mit_{ijkl} \\ mei_{ijkl} \end{pmatrix} \sim N \left[ \begin{pmatrix} \theta_{ijl}^{mit} \\ \theta_{ijl}^{mei} \end{pmatrix}, \sigma_i^2 I_2 \right]$$

$$\begin{pmatrix} \theta_{ijl}^{mit} \\ \theta_{ijl}^{mei} \end{pmatrix} \sim N \left[ \begin{pmatrix} \theta_{ij}^{mit} \\ \theta_{ij}^{mei} \end{pmatrix}, \begin{pmatrix} \tau_{mit}^2 & 0 \\ 0 & \tau_{mei}^2 \end{pmatrix} \right]$$

Pairs $mit_{ijkl}$ and $mei_{ijkl}$, for each gene $l$, each classification group $i$ (out of X, out of A) and gene type $j$ (parental, offspring), have individual means $\theta_{ijl}^{mit}$ and $\theta_{ijl}^{mei}$, respectively, and common classification group variances $\sigma_i^2$. Then, the objective is to compute,

$$\Pr\left[ \text{YES}_{\text{Out of X},l} \right] = \Pr\left[ \left\{ \theta_{\text{Out of X},par,l}^{mit} > \theta_{\text{Out of X},par,l}^{mei} \right\} \text{ I } \left\{ \theta_{\text{Out of X},off,l}^{mit} < \theta_{\text{Out of X,},off,l}^{mei} \right\} \right]$$

$$\Pr\left[ \text{YES}_{\text{Out of A},l} \right] = \Pr\left[ \left\{ \theta_{\text{Out of A},par,l}^{mit} > \theta_{\text{Out of A},par,l}^{mei} \right\} \text{ I } \left\{ \theta_{\text{Out of A,},off,l}^{mit} < \theta_{\text{Out of A,},off,l}^{mei} \right\} \right]$$

## Approach 3: Bayesian  hierarchical model

Hyperparameters

$$\Theta = (\theta_{11}^{mit}, \theta_{12}^{mit}, \theta_{21}^{mit}, \theta_{22}^{mit}, \theta_{11}^{mei}, \theta_{12}^{mei}, \theta_{21}^{mei}, \theta_{22}^{mei})$$

$$\Lambda = (\sigma_1^2, \sigma_2^2, \tau_{mit}^2, \tau_{mei}^2)$$

Prior distribution

$$p(\Theta, \Lambda) \propto \sigma_1^{-2} \sigma_2^{-2} \tau_{mit}^{-2} \tau_{mei}^{-2}$$

This represents vague/noninformative prior views.

## Approach 3: Bayesian  hierarchical model



46

## Approach 3: Bayesian hierarchical model

# 2. Biological Problem



X Y    X X



XY systems

- Does MSCI exist in *Drosophila*? YES

- Is MSCI involved in the distribution of sex-biased genes in the genome? YES



X

A

Parental gene

Retrogene

Testis bias

# 2. Biological Problem

## Stage-Specific Expression Profiling of *Drosophila* Spermatogenesis Suggests that Meiotic Sex Chromosome Inactivation Drives Genomic Relocation of Testis-Expressed Genes

Maria D. Vibranovski[1], Hedibert F. Lopes[2], Timothy L. Karr[3], Manyuan Long[1]*

1 Department of Ecology and Evolution, The University of Chicago, Chicago, Illinois, United States of America, 2 The University of Chicago Booth School of Business, Chicago, Illinois, United States of America, 3 The Biodesign Institute, Arizona State University, Tempe, Arizona, United States of America



| Título   1–20 | Citado por | Ano |
|---|---|---|
| Stage-specific expression profiling of Drosophila spermatogenesis suggests that meiotic sex chromosome inactivation drives genomic relocation of testis-expressed genes<br>MD Vibranovski, HF Lopes, TL Karr, M Long<br>PLoS genetics 5 (11), e1000731 | 109 | 2009 |

# 2. Biological Problem

## Male infertility



In Humans:

- 40% of infertility
- High cost treatment
- Associated with gene expression/function in spermatogenesis
- 30% of infertility are sperm deficiency

Drosophila Model:

- Spermatozoa is one of the few cell types that has homologous function for all sexual organism including humans
- Post-meiotic transcription

# 2. Biological Problem

http://pondside.uchicago.edu/~longlab/spermpress/

# 5. Ongoing Biological Problems



Post-meiotic cells

Mitotic cells

Meiotic cells

# 5. Ongoing Biological Problems

## Post-meiosis Over-expression



~ 20% of the genes

RNA level

Mitosis | Meiosis | Post-meiosis

Spermatogenesis

# 5. Ongoing Biological Problems

RNA synthesis does <span style="color:red">NOT</span> occur in post-meiotic stages

Testis Apical region

Spermatocytes



RNA synthesis: incorporation of [$^3$H]uridine in the *Drosophila melanogaster* testes has been studied by autoradiography.

# 5. Ongoing Biological Problems



56

# 5. Ongoing Biological Problems

| Transcription |
|---|

| Translation |
|---|

| RNA |
|---|



RNA level

Yes for post-meiotic transcription

No post-meiotic transcription

| Mitosis | Meiosis | Post-meiosis |
|---|---|---|

Spermatogenesis

Bromo-uridine (BrU) incorporation: Direct evidence

Domitille Chalopin

(1) Incorporation of BrU

(2) Labeling nascent RNAs with Br-UTP

(3) Binding anti-BrdU antibody

(4) Imaging

58

Vibranovski et al, 2010, Genetics

# 5. Ongoing Biological Problems

Bromo-uridine (BrU) incorporation:
Post-meiotic transcription - Direct evidence

Domitille Chalopin



Spermatid cyst

**BrU**
**DNA**
n: nuclei

Vibranovski et al, 2010, Genetics

# 5. Ongoing Biological Problems



Post-meiotic cells

Mitotic cells

Meiotic cells

# 5. Ongoing Biological Problems



Zhang, Vibranovski, Krinsky and Long, 2010

Júlia Raices

# 5. Ongoing Biological Problems

Júlia Raices



**Proportion of Genes by Spermatogenesis Phase and by age**

Júlia Raices

# 5. Ongoing Biological Problems

Júlia Raices

# Perspectives

Development

Molecular Biology

Bioinformatics

Cell Biology

Statistics

Genetics

# Acknowledgments

X➜A

A

A\A

B

NO Out of the X movement

Out of the X movement

D. melanogaster    Aedes aegypti    Anopheles gambiae

**Toups and Hahn, 2010, Genetics**