# Bioinformatics: Introduction and Methods
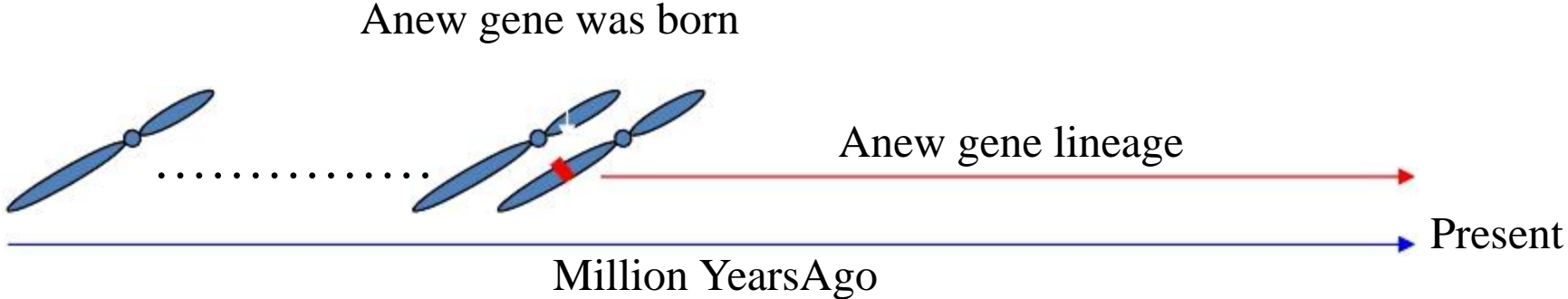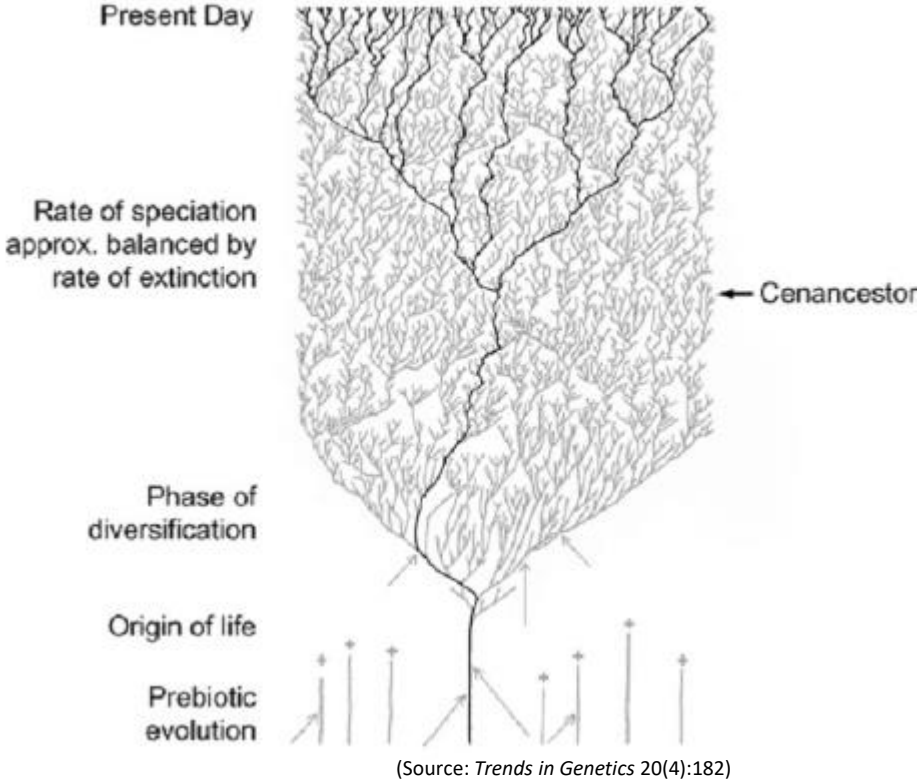## Le Zhang
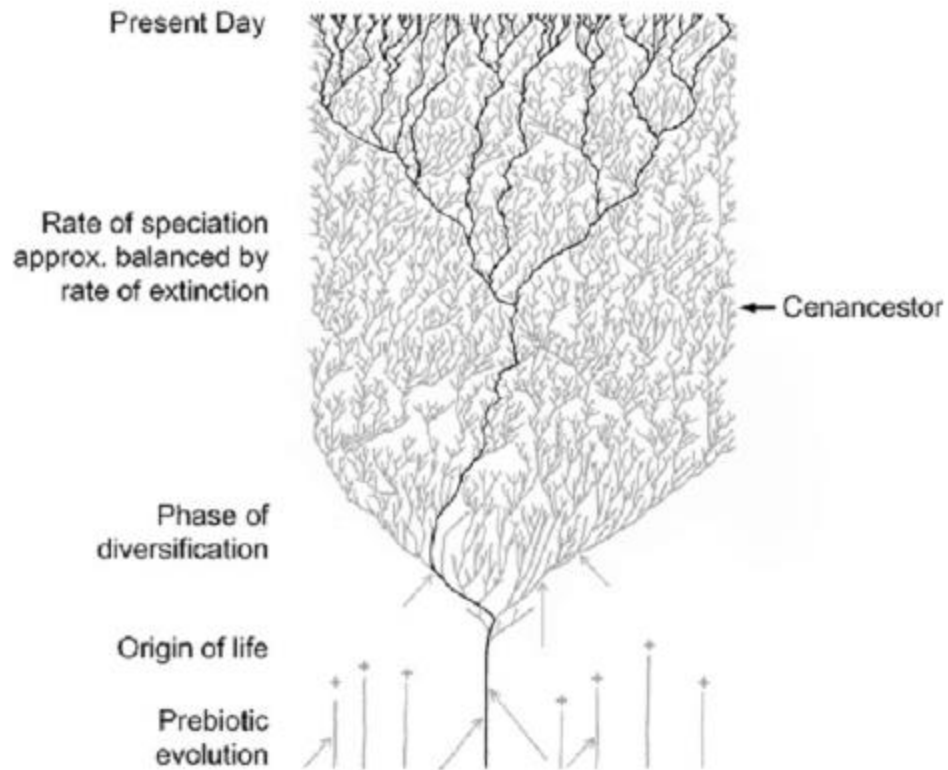
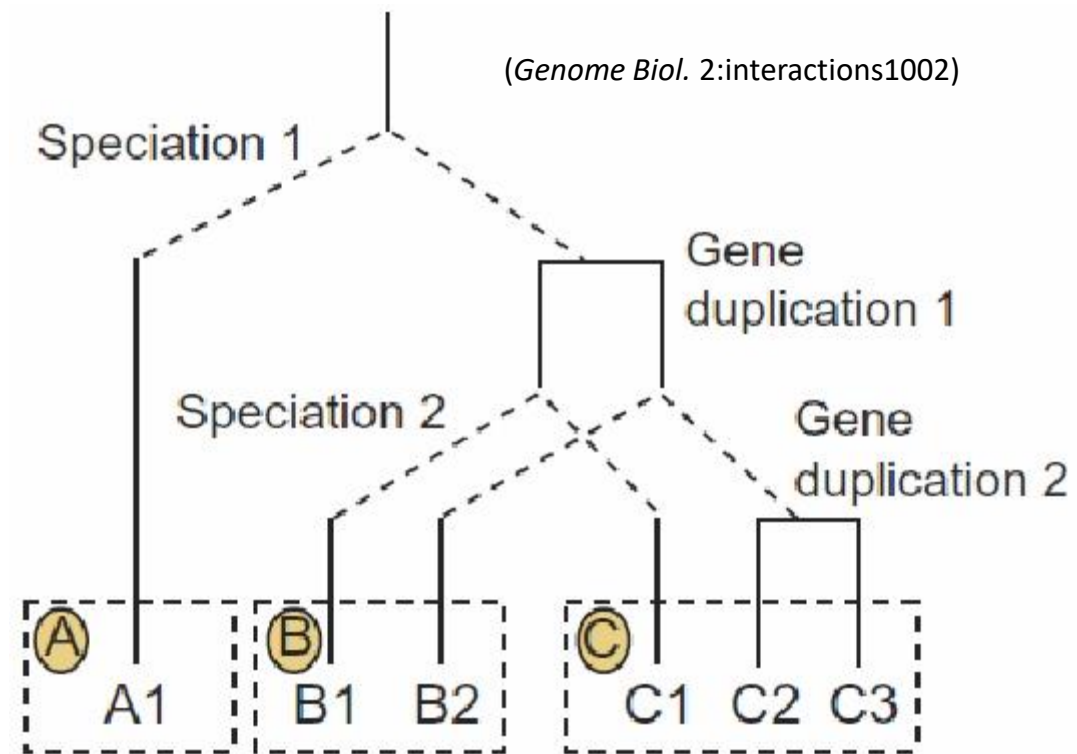## Computer Science Department, Southwest University

# Case Study 2:
# From Dry to Wet, an Evolutionary Story
## Le Zhang, Ph. D.
## Computer Science Department
## Southwest University

# Newly emerging genes are key sources of evolutionary novelty



Present Day

Rate of speciation approx. balanced by rate of extinction

← Cenancestor

Phase of diversification

Origin of life

Prebiotic evolution

(Source: *Trends in Genetics* 20(4):182)



Anew gene was born

Anew gene lineage

Present

Million YearsAgo

Present Day

Rate of speciation approx. balanced by rate of extinction

← Cenancestor

Phase of diversification

Origin of life

Prebiotic evolution

(Source: *Trends in Genetics* 20(4):182)



(*Genome Biol.* 2:interactions1002)

Speciation 1

Speciation 2

Gene duplication 1

Gene duplication 2

A    B    C

A1    B1    B2    C1    C2    C3
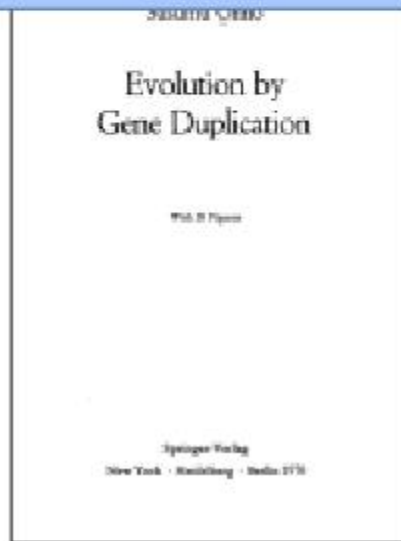
- **Orthologs**: homologuous genes result from **speciation** event

- **Paralogs**: homologuous genes result from **duplication** event

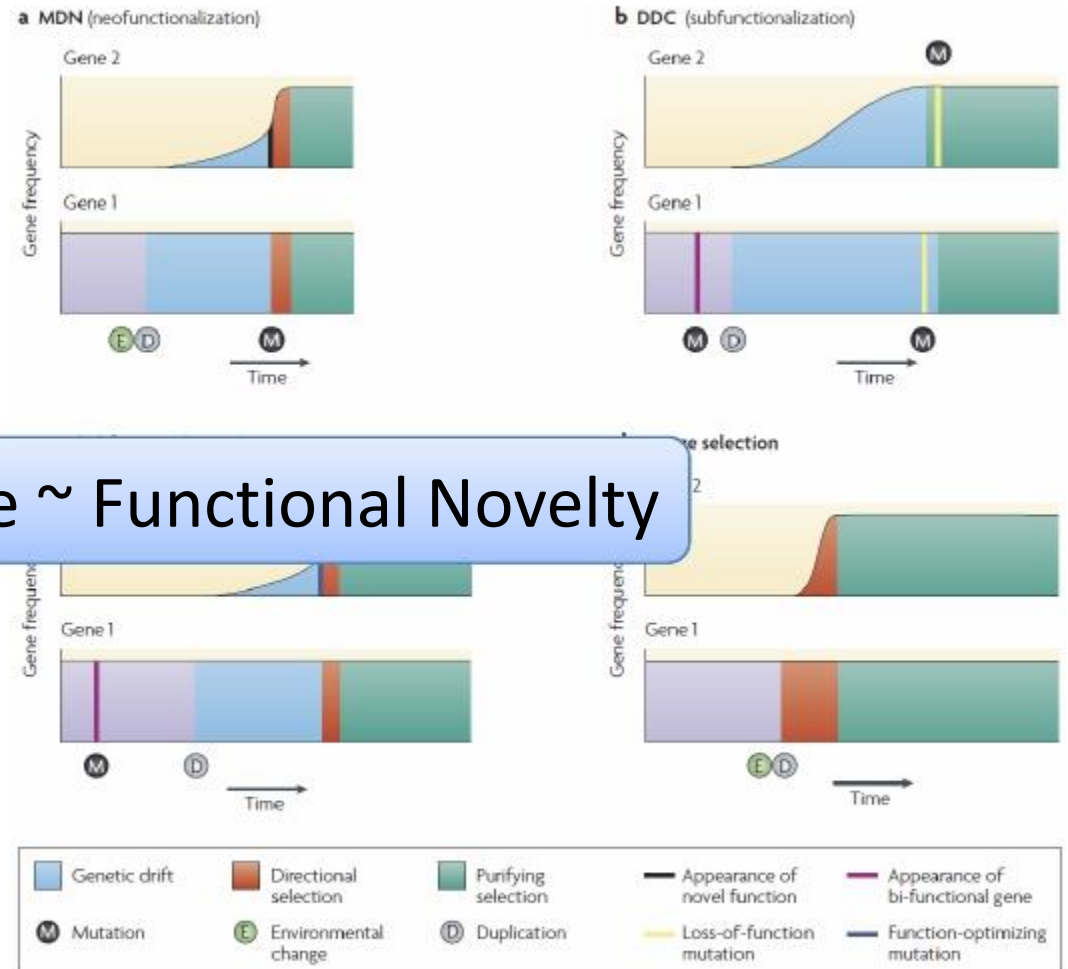Duplicated area

Before duplication

After duplication

(http://www.genome.gov/)

Susumu Ohno
February 1, 1928–January 13, 2000

Evolution by Gene Duplication

**a** MDN (neofunctionalization)

Gene 2

Gene 1

Gene frequency

E D          M

Time

**b** DDC (subfunctionalization)

Gene 2          M

Gene 1

Gene frequency

M D          M

Time

selection

2

Gene frequency

Gene 1

M          D

Time

Gene 1

Gene frequency

E D

Time

Sequence Divergence ~ Functional Novelty

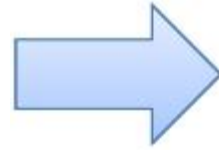| Genetic drift | Directional selection | Purifying selection | — Appearance of novel function | — Appearance of bi-functional gene |
| M Mutation | E Environmental change | D Duplication | Loss-of-function mutation | Function-optimizing mutation |

(Conant *et al*. 2008)

Figure 8. Developmental sequences of various vertebrates shown in phylogenetic context. Note the shared similarities of some closely related taxa, particularly the amniotes (modified from Richardson et al. 1998.) (Figure Source: ncseprojects.org/image/icons-evolution-figure-8)
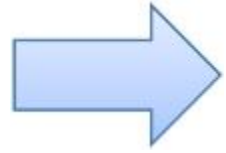
# Computationally screening for function divergent genes involved in early development regulation



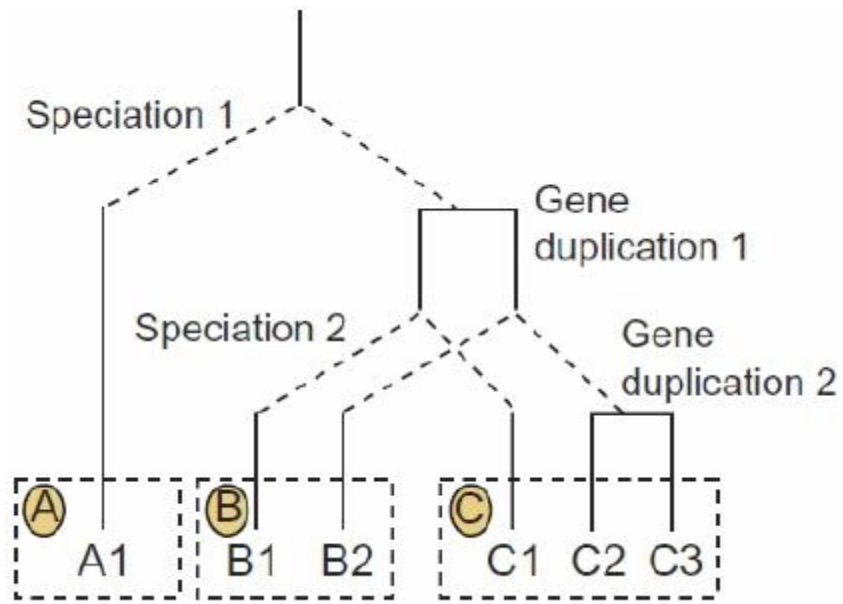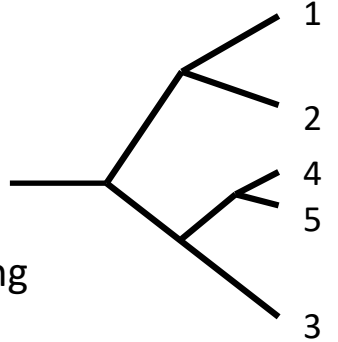Sequences from 14 species

BLAST

Similarity Matrix

Identify similar sequences among inter-species as well as intra-species

Insert an in-video survey here.

Question: Which Scoring Matrix would you like to use here?

A) PAM1
B) BLOSUM80
C) BLOSUM62
D) PAM2

C is the right answer, but B is also okay.
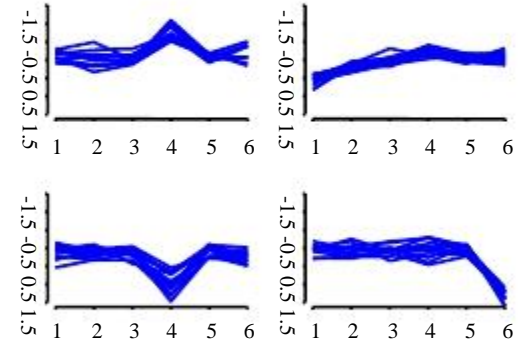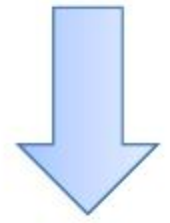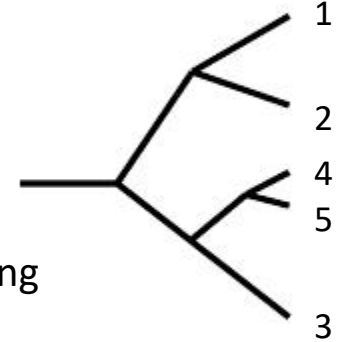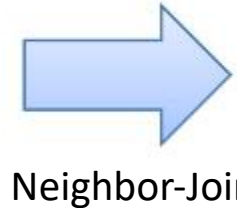
Either the answer is right or wrong, just go on.

# Computationally screening for function divergent genes involved in early development regulation



BLAST

Similarity Matrix

Identify similar sequences among inter-species as well as intra-species

Neighbor-Joining

1
2
4
5
3

Sequences from 14 species

Speciation 1

Gene duplication 1

Speciation 2

Gene duplication 2

A1   B1   B2   C1   C2   C3

(*Genome Biol.* 2:interactions1002)

# Computationally screening for function divergent genes involved in early development regulation

Insert an in-video survey here.

Question: Which database(s) you would NOT use here?

A) PDB
B) NCBI GEO
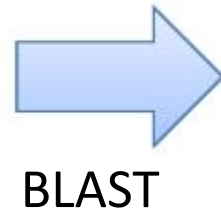C) NCBI SRA
D) EBI ArrayExpress

A is the right answer.

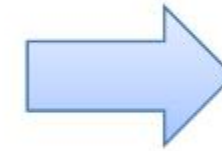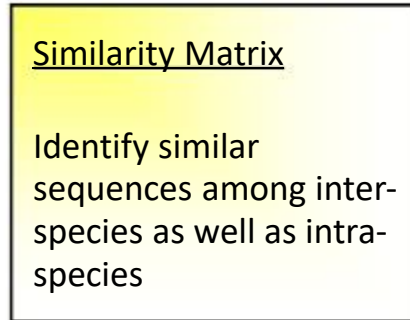Either the answer is right or wrong, just go on.

# Computationally screening for function divergent genes involved in early development regulation
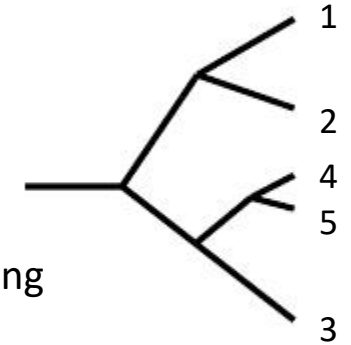


Sequences from 14 species

BLAST

**Similarity Matrix**

Identify similar sequences among inter-species as well as intra-species

Neighbor-Joining

1
2
4
5
3

ARRAYEXPRESS

GEO

Expression Profiles

# Computationally screening for function divergent genes involved in early development regulation

Insert an in-video survey here.

Question: Which database(s) would you like use here?

A) KEGG
B) Gene Ontology Annotation

Both A and B are correct

Either the answer is right or wrong, just go on.

# Computationally screening for function divergent genes involved in early development regulation
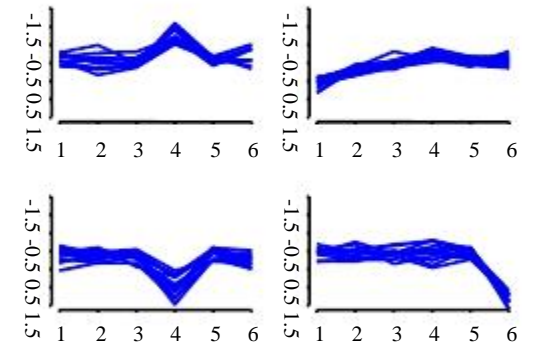


Sequences from 14 species

BLAST

Similarity Matrix
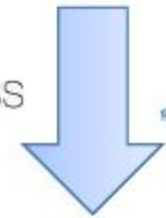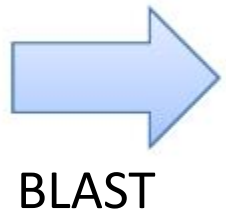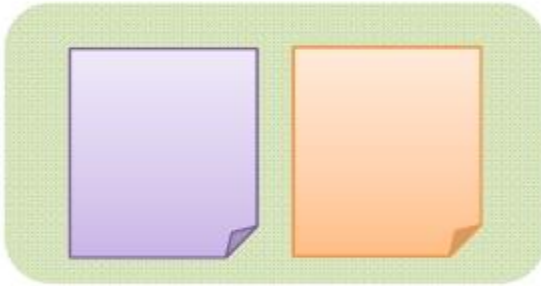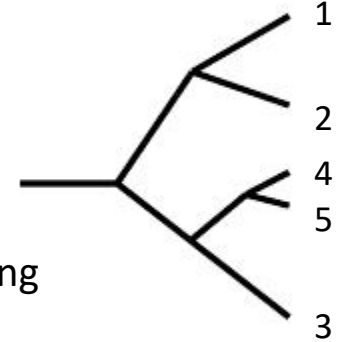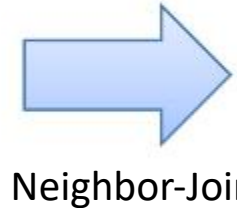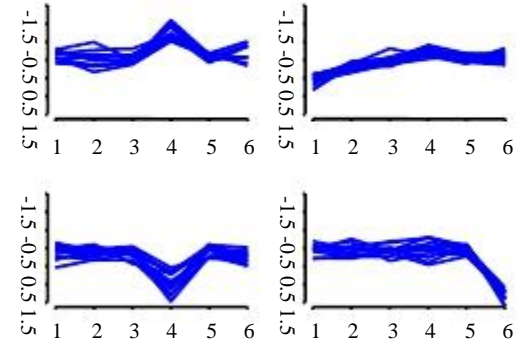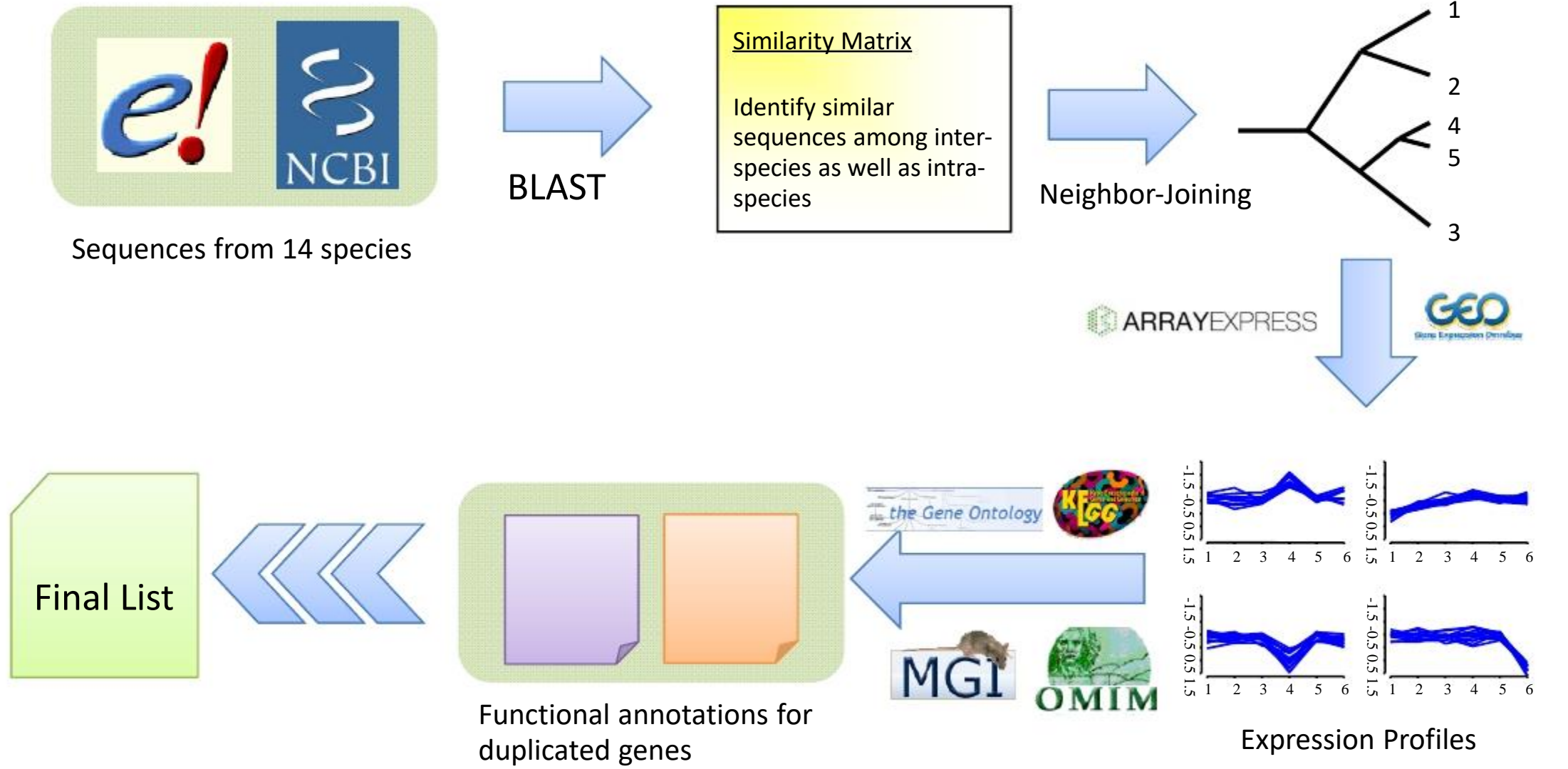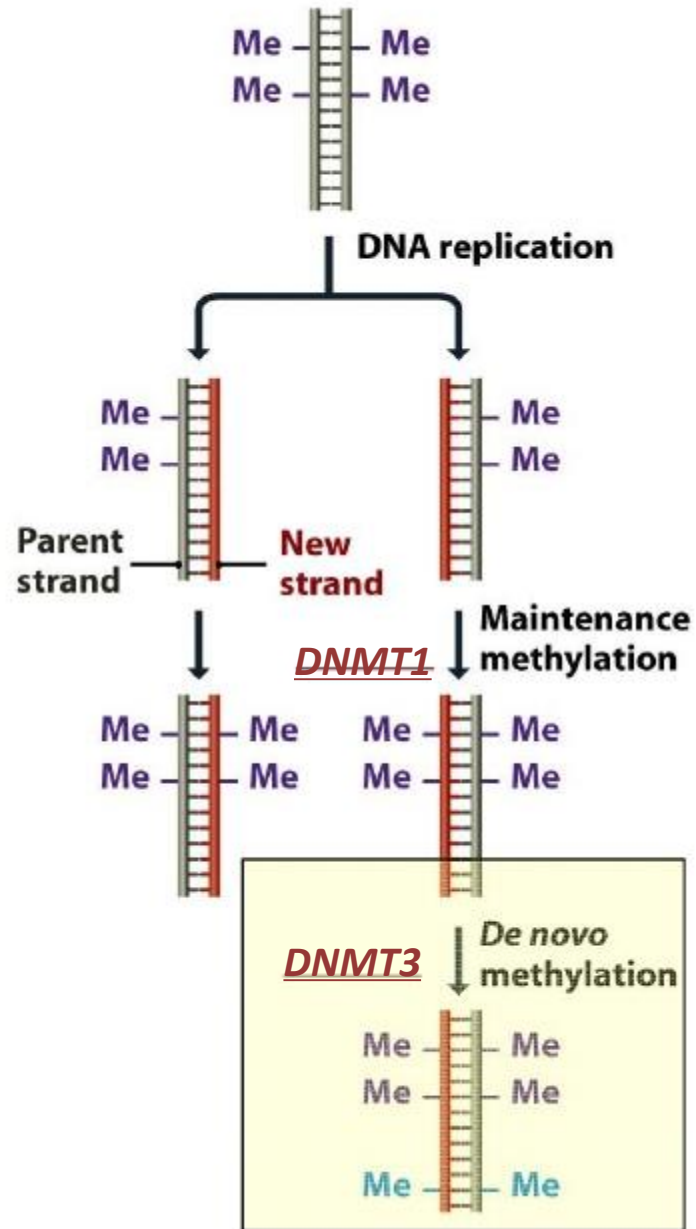
Identify similar sequences among inter-species as well as intra-species

Neighbor-Joining

Final List

Functional annotations for duplicated genes

Expression Profiles

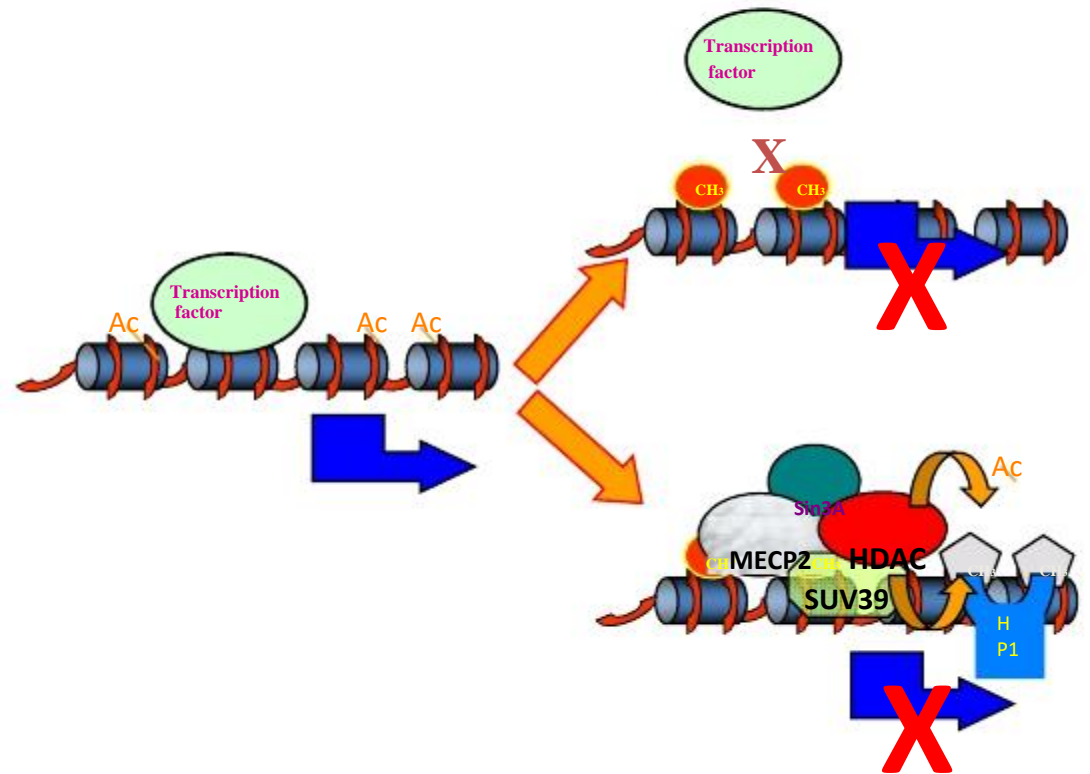# Computationally screening for function divergenced genes involved in early development regulation

**Computational Genomic Analysis and Bioinformatics through *MAPK*:**

1. Sequence databases were constructed directly from the Ensembl website;

2. Each peptide sequence in the database was used to search the database using BLAST package.

3. Phylogenetic trees were constructed and paralogous pairs are identified from the resulting alignments based on a minimal amino acid identity (e.g. 50% and 70%) and an overlap of $\geq$ 35 amino acids in the region of local alignment.

4. Coding regions of pairs that meet these criteria will be aligned with the corresponding region and inspected for putative function divergence hallmarks.

5. Local warehouse were searched for further indicators derived from high-throughput data (esp. genetic, genomic, transcriptomic, proteomic and pathway data).

**7 out of 50000+ new paralogous pairs showed clear functional divergence features involved in early development regulation.**
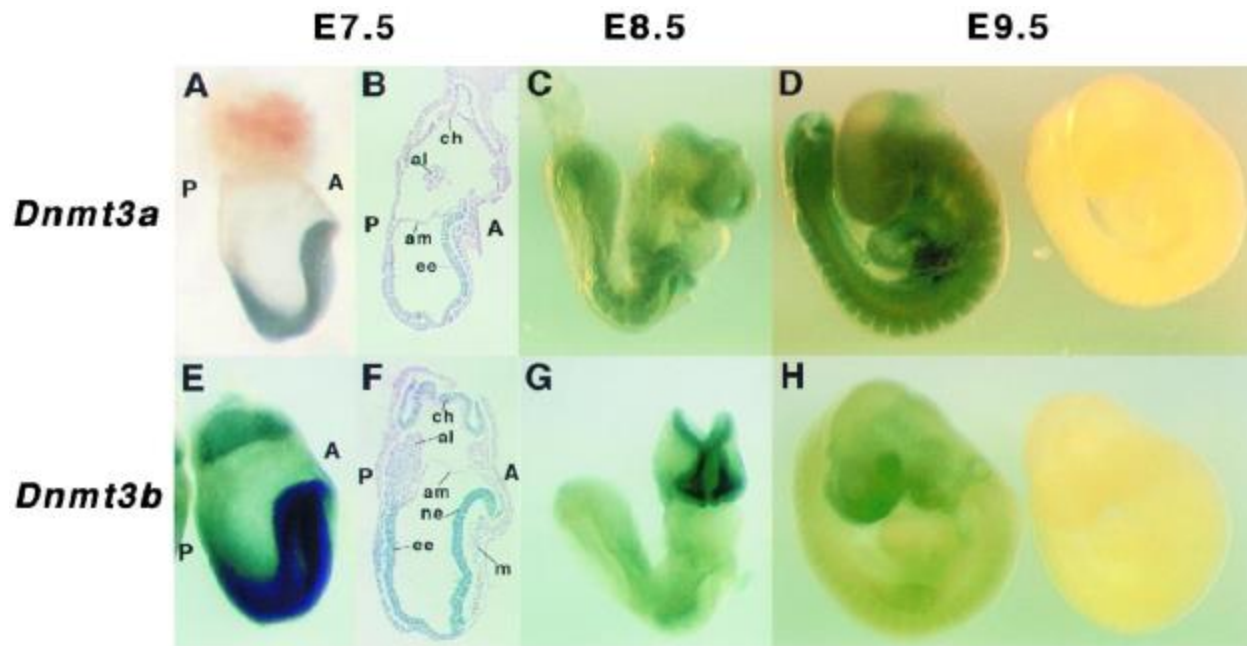
DNA methylation silences gene expression by two mechanisms

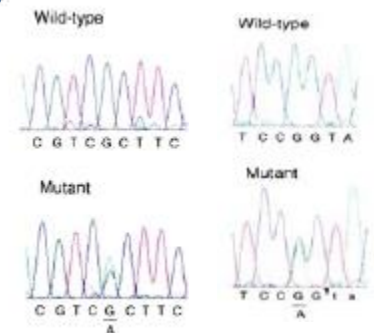(Modified from Moshe Szyf's slide)

DNMT3-induced methylation is critical for early mouse embryo development

...and also for human ICF Syndrome

(http://www.newsoftheworld.co.uk/news/402211/Kiss-of-death-Grace-and-Luke-Hicklin-have-hereditary-ICF-syndrome.html)
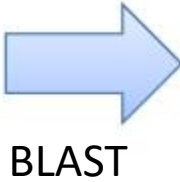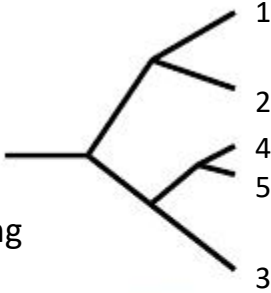
# Computationally screening for function divergent genes involved in early development regulation



Sequences from 14 species

BLAST

**Similarity Matrix**
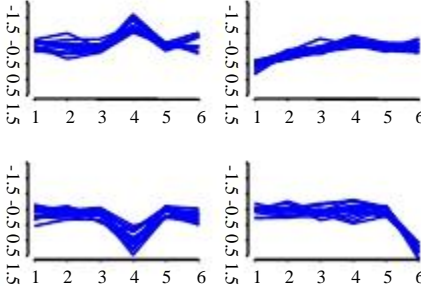
Identify similar sequences among inter-species as well as intra-species

Neighbor-Joining

1
2
4
5
3

ARRAYEXPRESS    GEO

the Gene Ontology    KEGG

MGI    OMIM

Expression Profiles

Functional annotations for duplicated genes
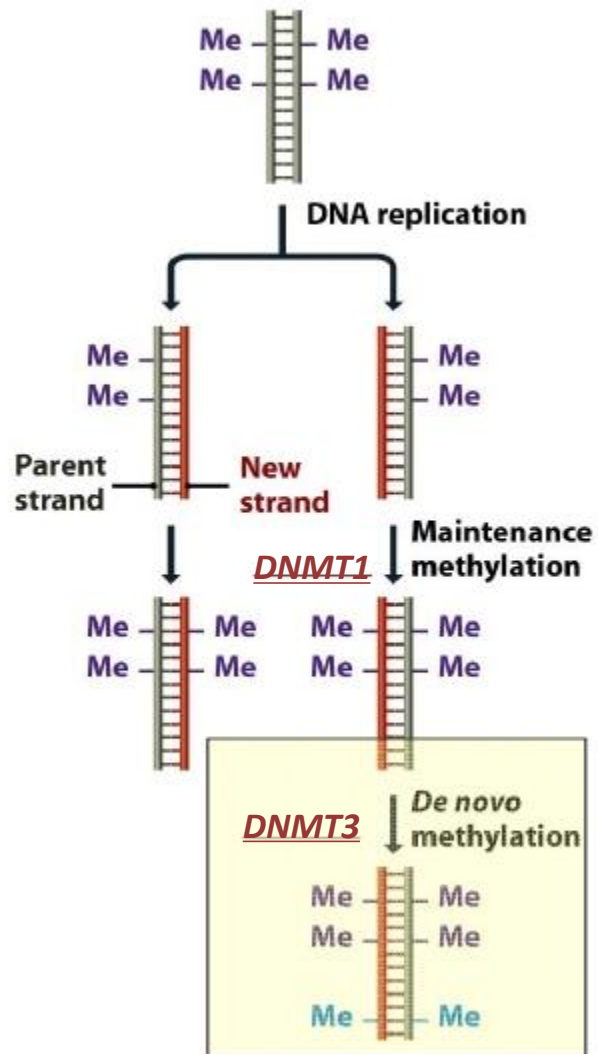
Final List

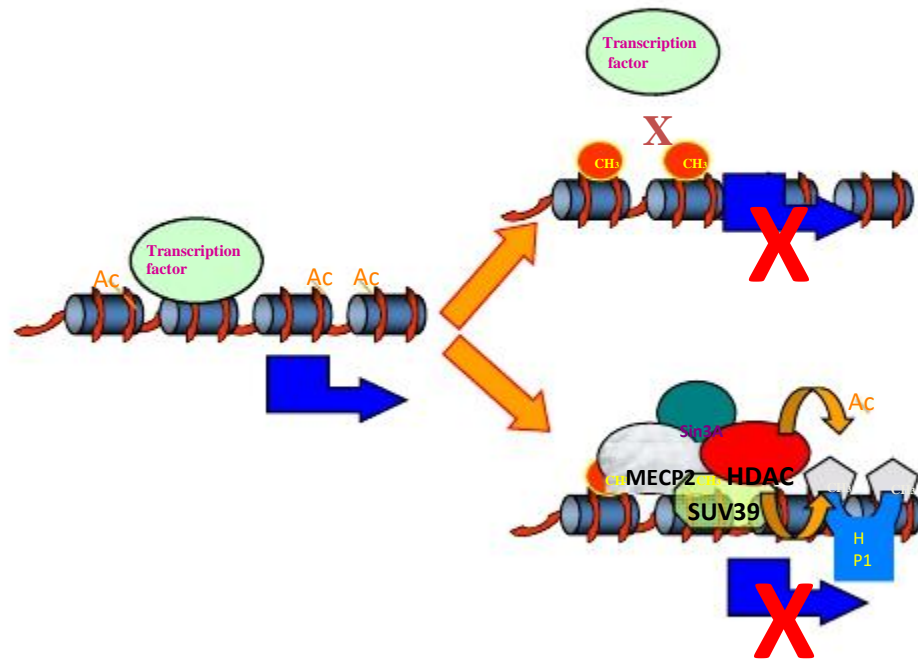# Computationally screening for function divergenced genes involved in early development regulation

**Computational Genomic Analysis and Bioinformatics through *MAPK*:**

1. Sequence databases constructed directly from the Ensembl website;

2. Each peptide sequence in the database used to search the database using BLAST package.

3. Phylogenetic trees were constructed and paralogous pairs are identified from the resulting alignments based on a minimal amino acid identity (e.g. 50% and 70%) and an overlap of $\geq$ 35 amino acids in the region of local alignment.

4. Coding regions of pairs that meet these criteria will be aligned with the corresponding region and inspected for putative function divergence  hallmarks.

5. Local warehouse were searched for further indicators derived from high-throughput data (esp. genetic, genomic, transcriptomic, proteomic and pathway data).

**7 out of 50000+ new paralogous pairs showed clear functional divergence features involved in early development regulation.**
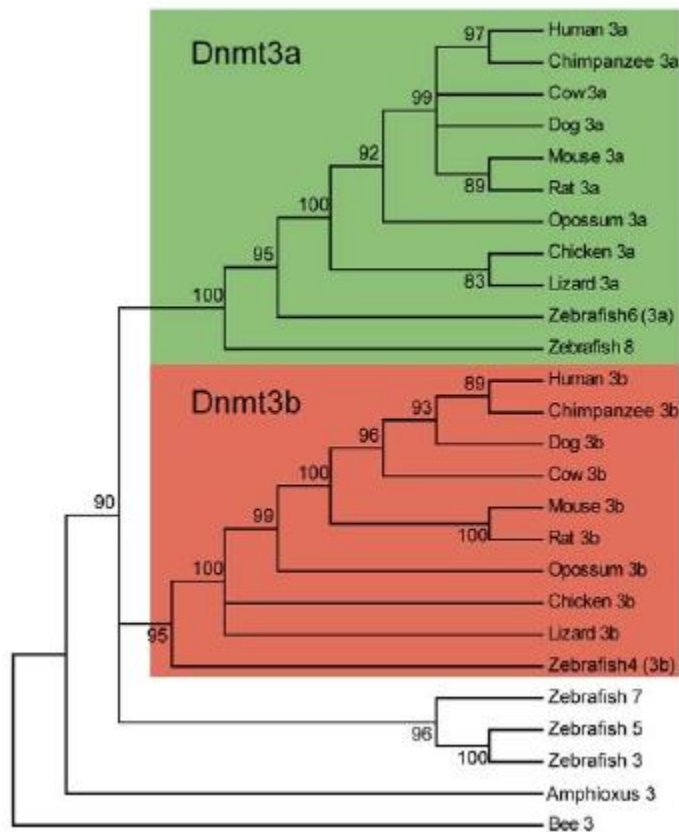
DNA methylation silences gene expression by two mechanisms
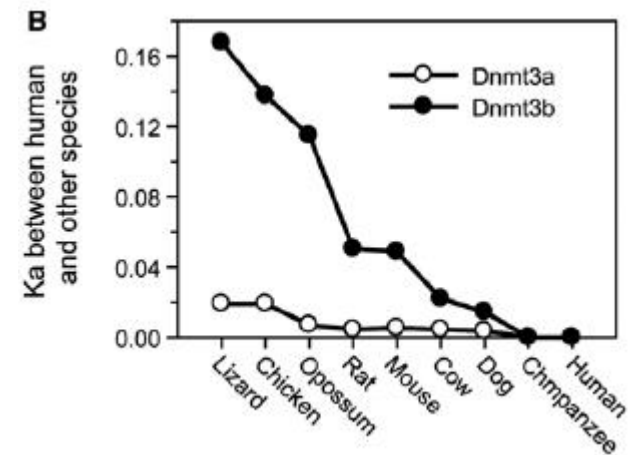
(Modified from Moshe Szyf's slide)

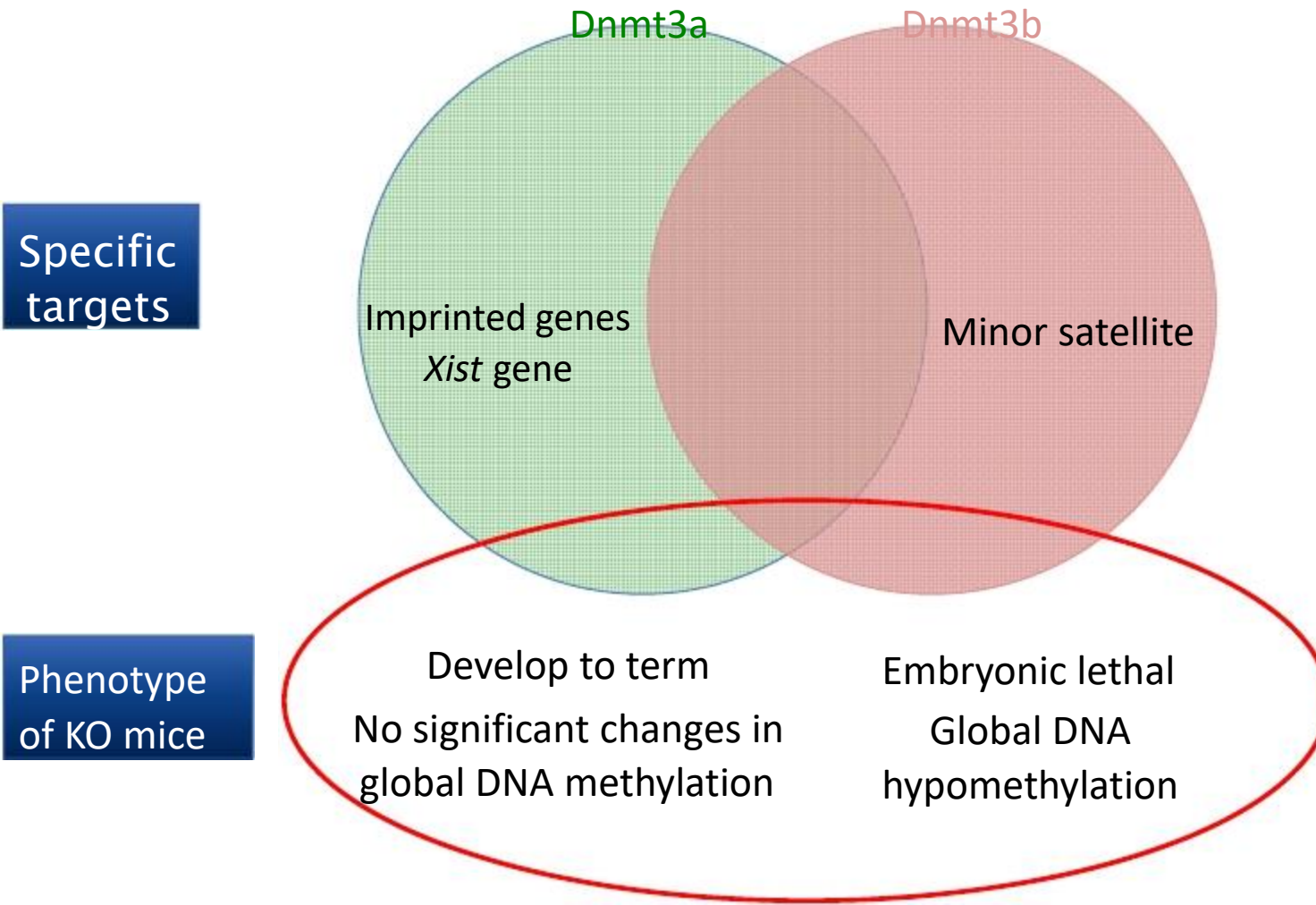# DNMT3 raised around the separation of vertebrates



(Shen *et al. Nucle Acids Res.* 2010)

**A**

| | Synonymous | | | Non-Synonymous | | |
|---|---|---|---|---|---|---|
| | $S_{3a}$ | $S_{3b}$ | P-value | $N_{3a}$ | $N_{3b}$ | P-value |
| Human | 221 | 228 | 0.377 | 294 | 343 | 0.025[*] |
| Chimpanzee | 222 | 227 | 0.413 | 294 | 343 | 0.025[*] |
| Dog | 227 | 236 | 0.340 | 298 | 342 | 0.041[*] |
| Cow | 230 | 222 | 0.366 | 299 | 347 | 0.029[*] |
| Mouse | 220 | 228 | 0.342 | 300 | 345 | 0.040[*] |
| Rat | 227 | 229 | 0.460 | 298 | 340 | 0.049[*] |
| Opossum | 247 | 257 | 0.319 | 298 | 357 | 0.011[*] |
| Chicken | 226 | 208 | 0.196 | 307 | 350 | 0.047[*] |
| Lizard | 231 | 229 | 0.463 | 305 | 353 | 0.031[*] |

**B**

Specific targets

Dnmt3a

Dnmt3b

Imprinted genes
*Xist* gene

Minor satellite

Phenotype of KO mice

Develop to term

No significant changes in global DNA methylation
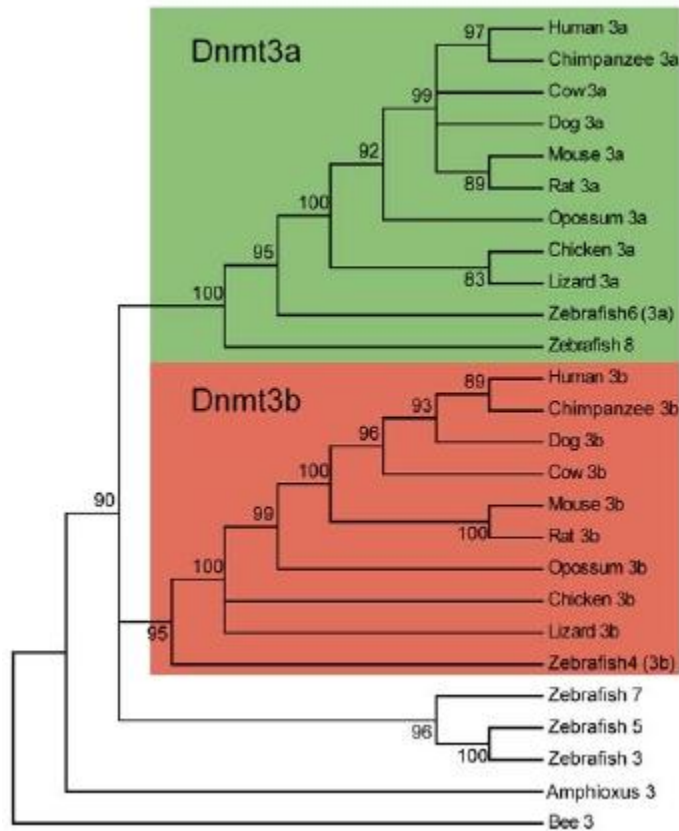
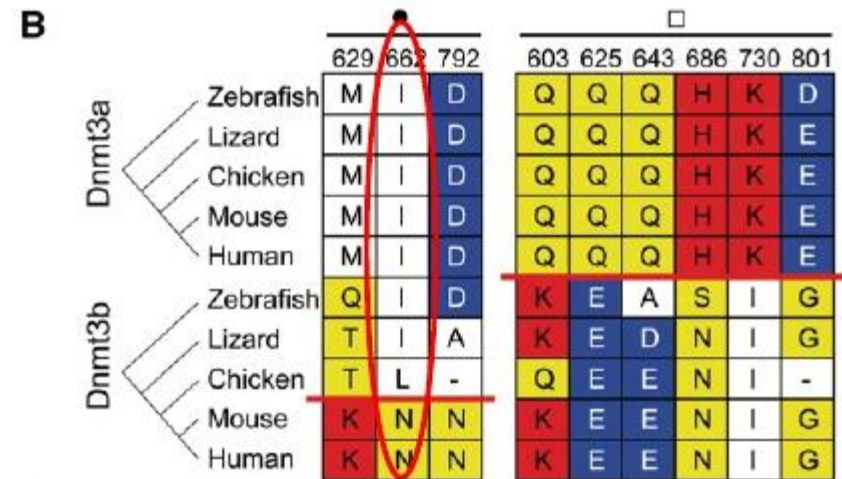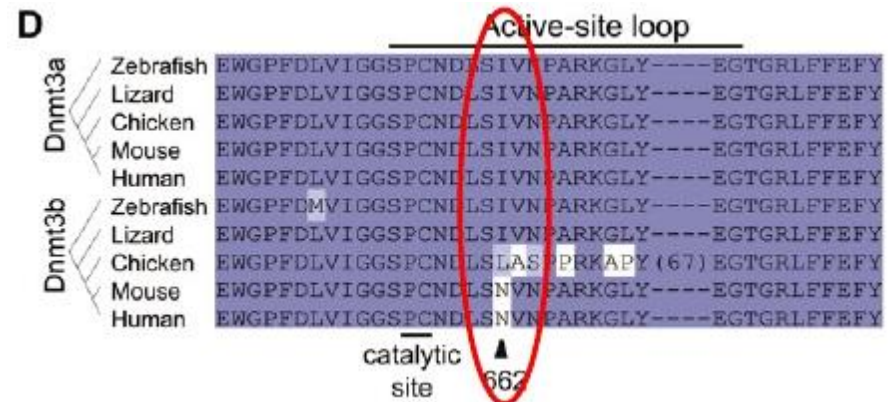Embryonic lethal

Global DNA hypomethylation
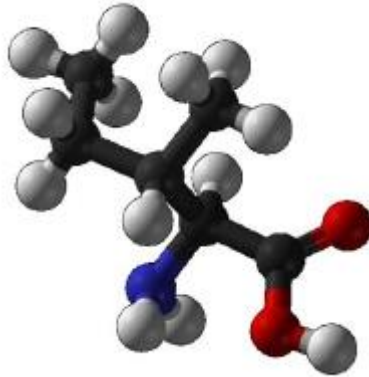
(Source: Li Shen, 2010)

# A mammalian DNMT3b-specific amino acid change appeared near catalytic site
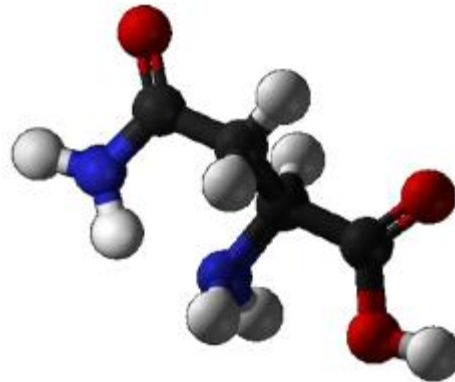


(Shen *et al. Nucle Acids Res.* 2010)
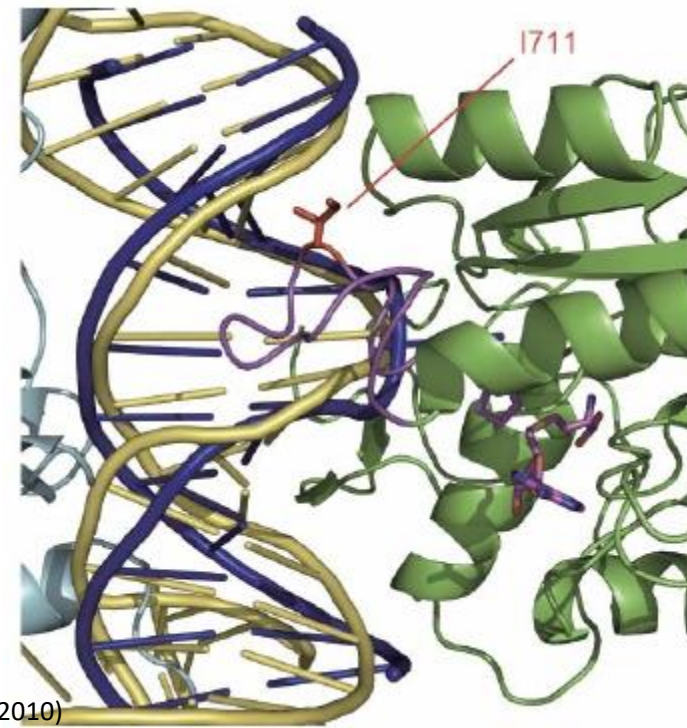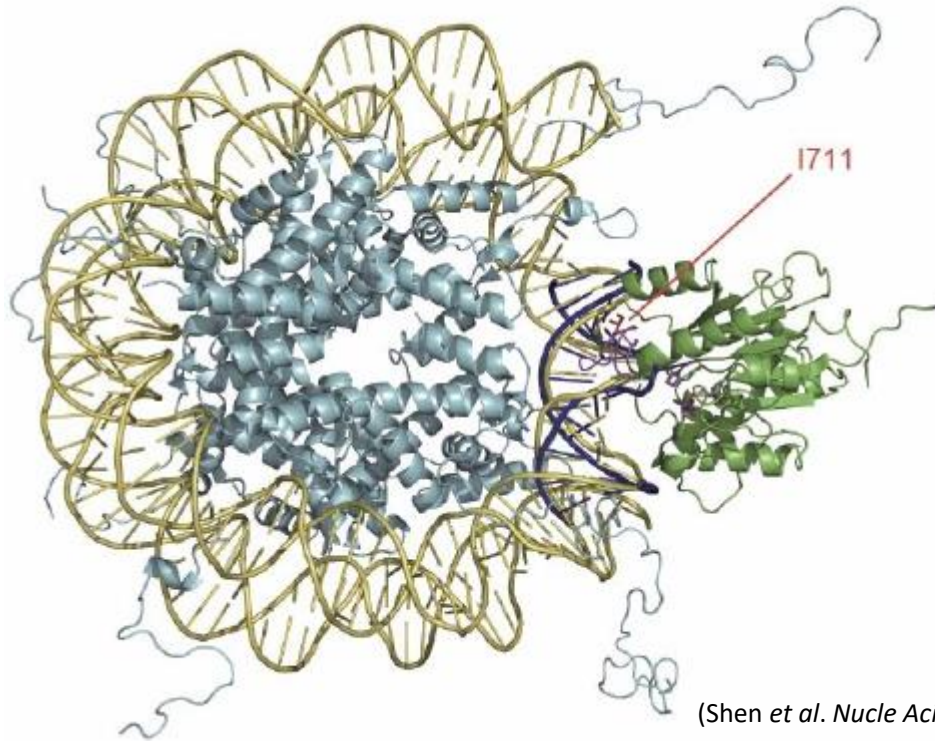
**Isoleucine (I)**
异亮氨酸



$C_4H_8N_2O_3$

- non-Polar,
- Hydropathy index = 4.5
- pI at 25° C = 6.04

**Asparagine (N)**
天冬酰胺



- Polar,
- Hydropathy index = −3.5
- pI at 25° C = 10.76

(Shen *et al. Nucle Acids Res.* 2010)

Structural analysis suggested that the I→N results in a tighter enzyme-DNA interaction

# Test the hypothesis with wet experiences



Mammalian

Dnmt3a (*wild-type*)  →  I → N  →  Dnmt3a (*mutation*)

Dnmt3b (*wild-type*)  →  N → I  →  Dnmt3b (*mutation*)

"Site-directed mutagenesis is a molecular biology method that is used to make specific and intentional changes to the DNA sequence of a gene and any gene products. Also called site-specific mutagenes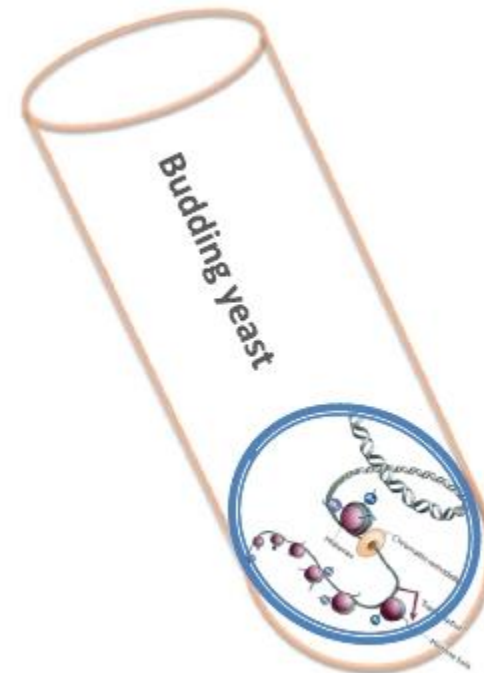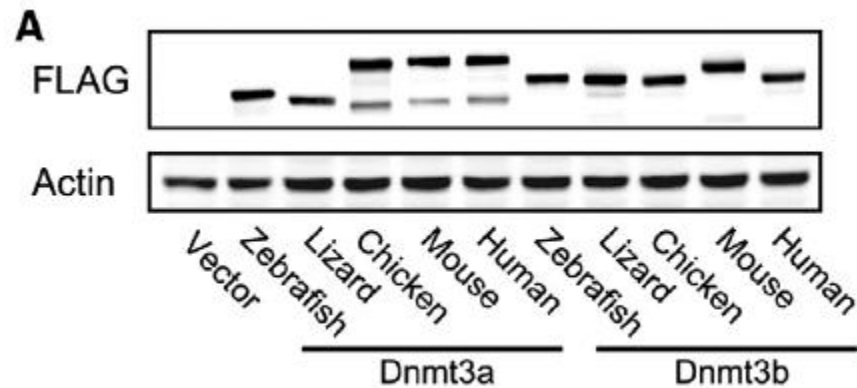is or oligonucleotide-directed mutagenesis, it is used for investigating the structure and biological activity of DNA, RNA, and protein molecules, and for protein engineering. With decreasing costs of oligonucleotide synthesis, artificial gene synthesis is now occasionally used as an alternative to site-directed mutagenesis." (Source: wikipedia.org)
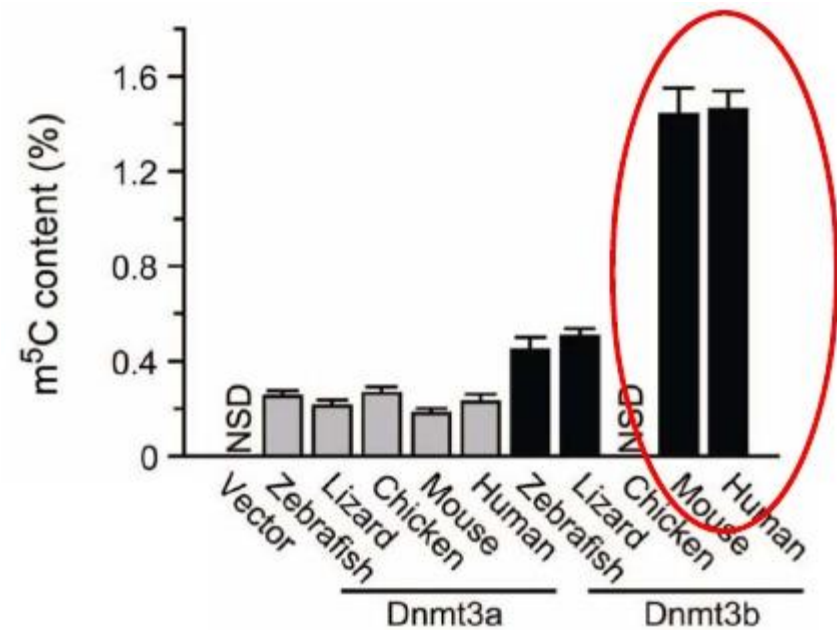
endogenous DNMITs ??

# Using budding yeast as the "*in vivo* test tube"

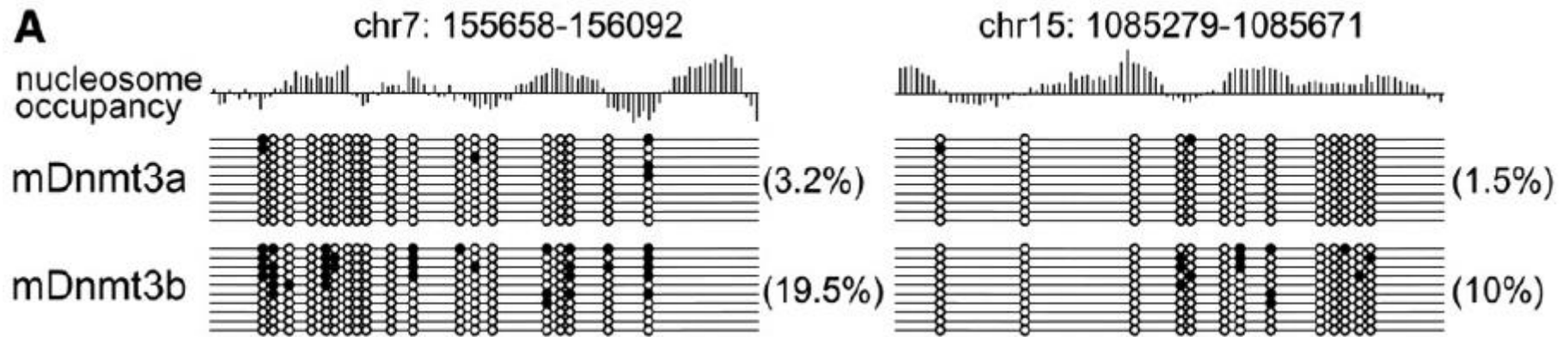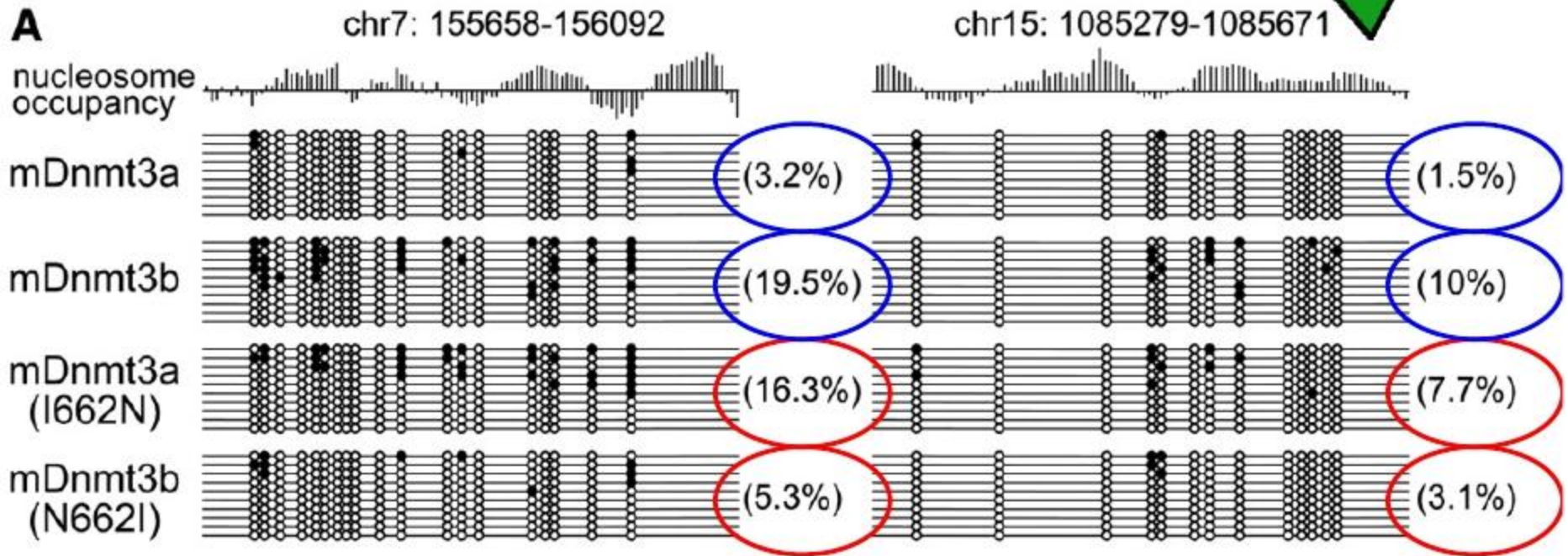| Epigenetic features | mammals | budding yeast |
|---|---|---|
| Chromatin | Yes | Yes |
| Histone acetylation | Yes | Yes |
| H3K4 | Yes | Yes |
| H3K36 | Yes | Yes |
| H3K79 | Yes | Yes |
| SWI/SNF complexes | Yes | Yes |
| CHD1 ATPase | Yes | Yes |
| SWR1 ATPase | Yes | Yes |
| ISWI ATPase | Yes | Yes |
| Endogenous methylation | Yes | No |

Budding Yeast

(Shen *et al. Nucle Acids Res.* 2010)

Mammalian Dnmt3b posses higher chromatin DNA methylation activity than Dnmt3a and non-mammalian Dnmt3b

Would I662N substitution accounts for the increased nucleosome DNA methylation activity in mammalian DNMT3b?
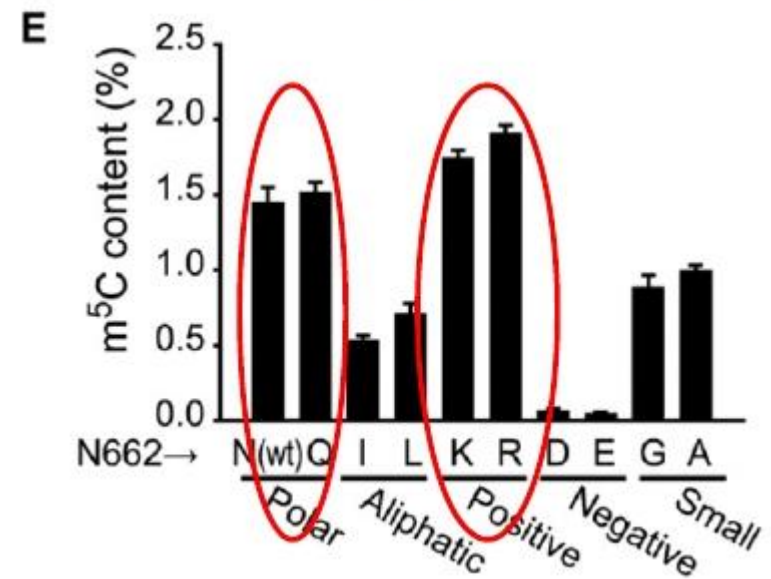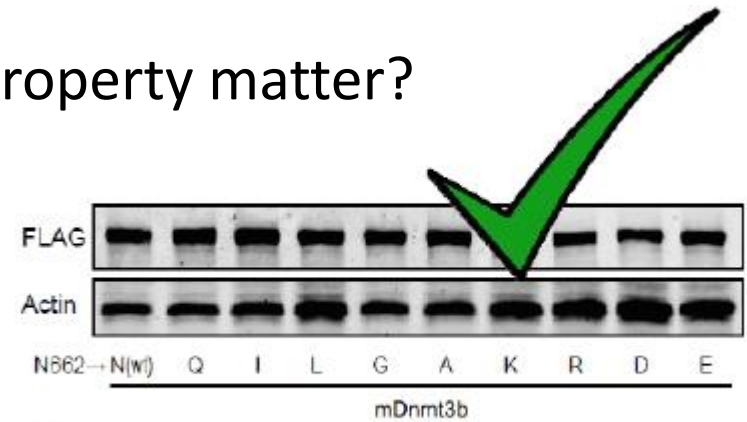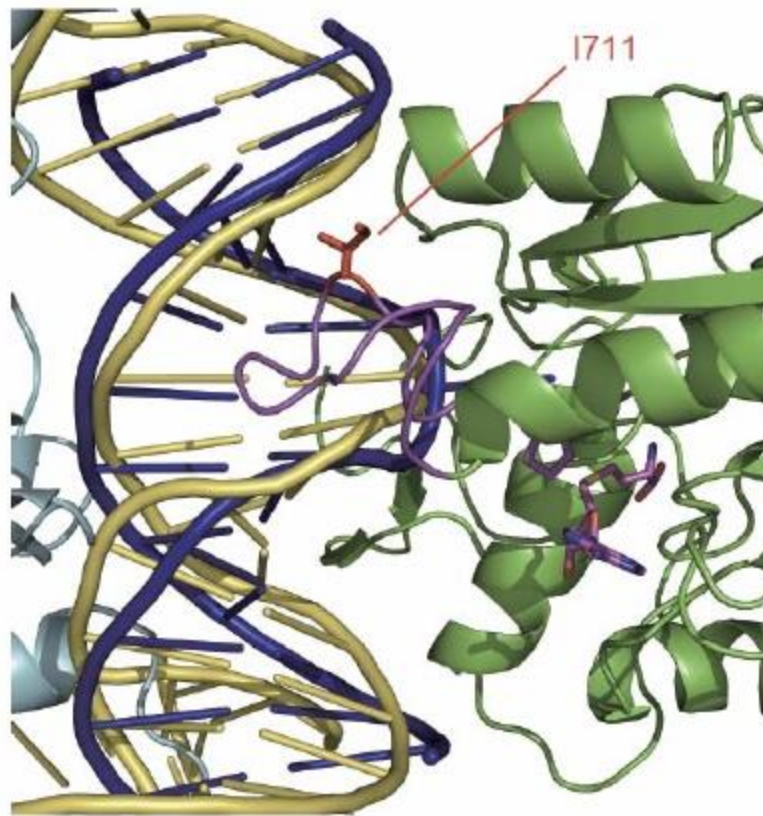


A

chr7: 155658-156092   chr15: 1085279-1085671

nucleosome occupancy

mDnmt3a   (3.2%)   (1.5%)

mDnmt3b   (19.5%)   (10%)

Would I662N substitution accounts for the increased nucleosome DNA methylation activity in mammalian DNMT3b?



**A**

chr7: 155658-156092        chr15: 1085279-1085671

nucleosome occupancy

| | chr7 | chr15 |
|---|---|---|
| mDnmt3a | (3.2%) | (1.5%) |
| mDnmt3b | (19.5%) | (10%) |
| mDnmt3a (I662N) | (16.3%) | (7.7%) |
| mDnmt3b (N662I) | (5.3%) | (3.1%) |

Does the changed physicochemical property matter?

A (hypothesis) connection of the chromatin DNA methylation activity of Dnmt3b with the density of repetitive sequences in the genome?

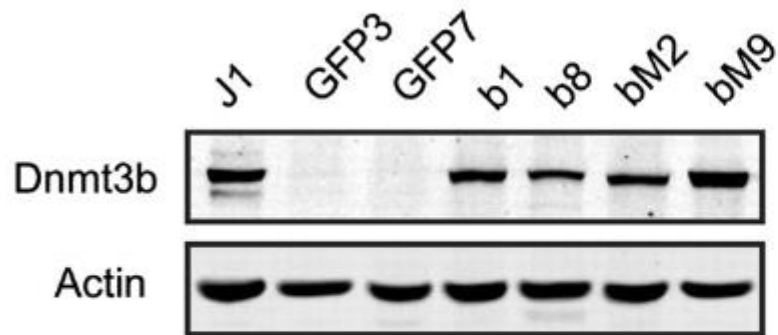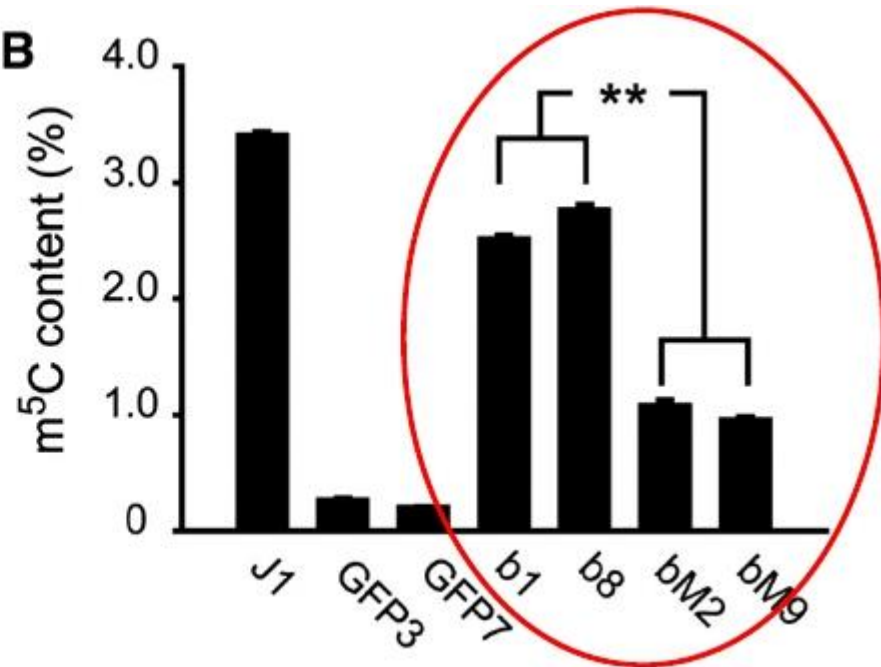| | Percentage of repeats in the genome | Chromatin DNA methylation activity of Dnmt3b |
|---|---|---|
|  | ~ 40-50% | highest |
|  ~ | ~10% | inactive |

GFP: GFP3, GFP7
mDnmt3b: b1, b8
mDnmt3b(N662I): bM2, bM9

The substitution (I662N) is crucial for mammalian Dnmt3b to efficiently methylate repetitive sequences in mammalian cells

# Summary

- Evolution-guided bioinformatics analysis successfully identified interesting genes involved in early development regulation showed clear functional novelty during evolution, and also provided strong hints for the key substitution, its biochemical effort, and the eventually functional significance.

- Key single substitution could result in significant functional novelty and help novel gene (re-)wired itself into existing circuits.

- An integrated, genome-scale bioinformatic analysis combined with targeted experimental assay is effective in studying complex biological system.

**Bioinformatics: an interdisciplinary field that develop and apply computer and computational technologies to study biomedical questions**

- As a technology, bioinformatics is a powerful technology to manage, search, and analyze big data in life sciences.

- As a methodology, bioinformatics is a top-down, holistic, data-driven, genome-wide, and systems approach that generates new hypotheses, find new patterns, and discover new functional elements.

知其道 用其妙 THIS IS HOW:

# Bioinformatics: Introduction and Methods

## Computer Science Department, Southwest University

# Thank you