# Methylation of Transposable Elements and Gene Expression in *Arabidopsis Thaliana*

Katia Antonenko[1, 2, 3], Marie Dogo[1, 2, 3], Jérémy Cohen[1, 2, 3], Sylvain Caillaud[1, 2, 3], Louna De Oliveira[4], Aurélien Petit[4], Vincent Colot[4], Chloé-Agathe Azencott[1, 2, 3], Pierre Baduel[4]

1. CBIO-Centre for Computational Biology, Mines Paris, PSL Research University 2. Institut Curie, PSL Research University 3. U1331 Institut National de la Santé et de la Recherche Médicale - INSERM 4. Institut de Biologie de l'École Normale Supérieure (IBENS) ENS, CNRS UMR8197, Inserm U1024, Paris
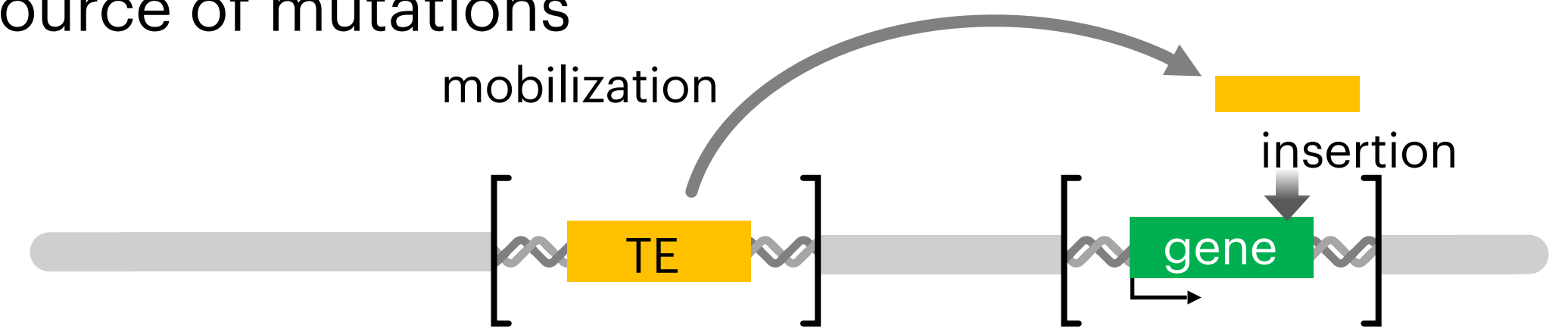
27 November 2025

# Contents

- Background: transposable elements and methylation

- Motivation: to explain GWAS findings

- Model: to understand methylation spreading
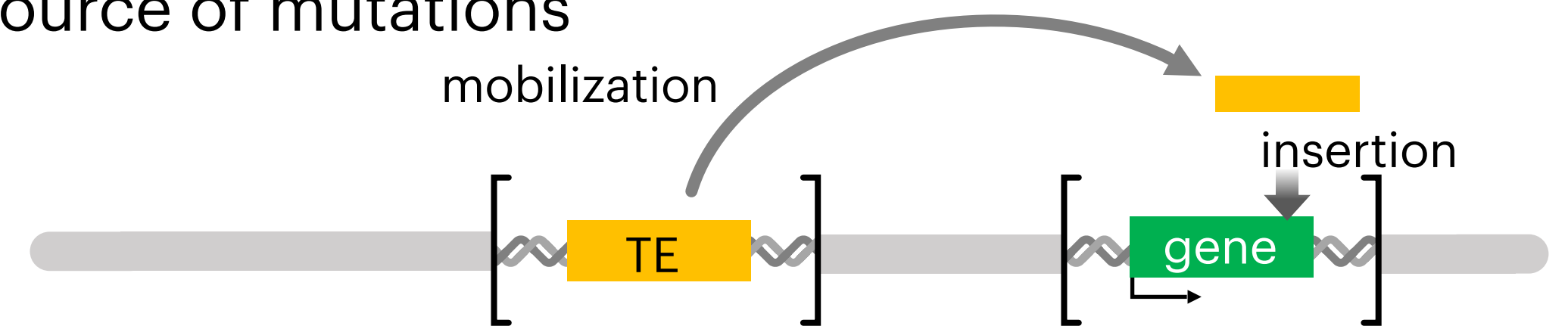
- Conclusions

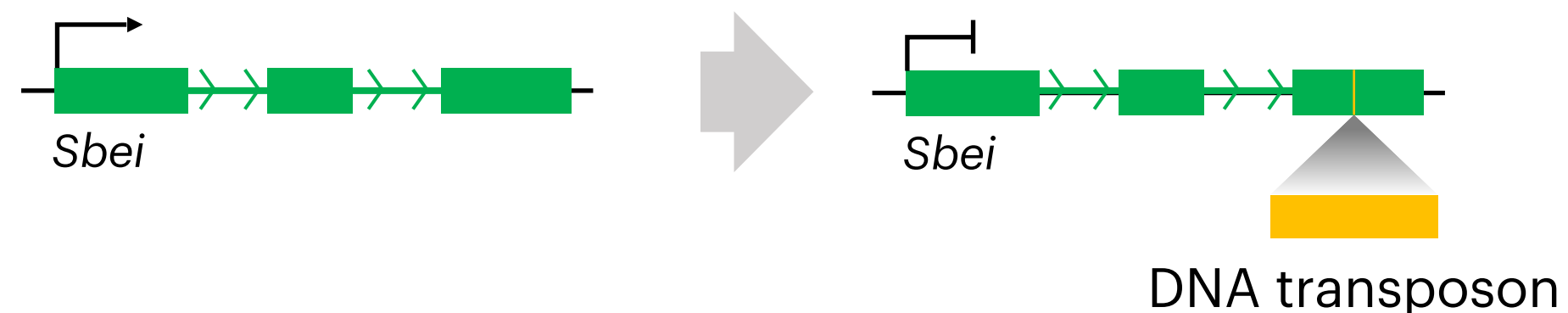# Transposable Elements

# Transposable Elements

- Transposable Elements (TEs, "jumping genes") are an important source of mutations

- TEs transpose by cut-and-paste or copy-and-paste mechanisms

- BUT: most TEs are degraded and do not transpose

# Transposable Elements

- Transposable Elements (TEs, "jumping genes") are an important source of mutations

- TEs transpose by cut-and-paste or copy-and-paste mechanisms

- BUT: most TEs are degraded and do not transpose
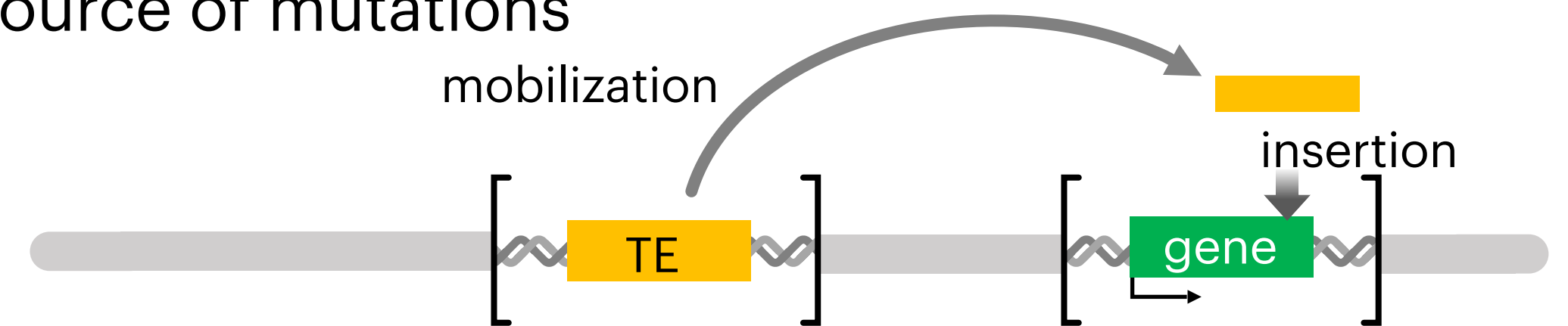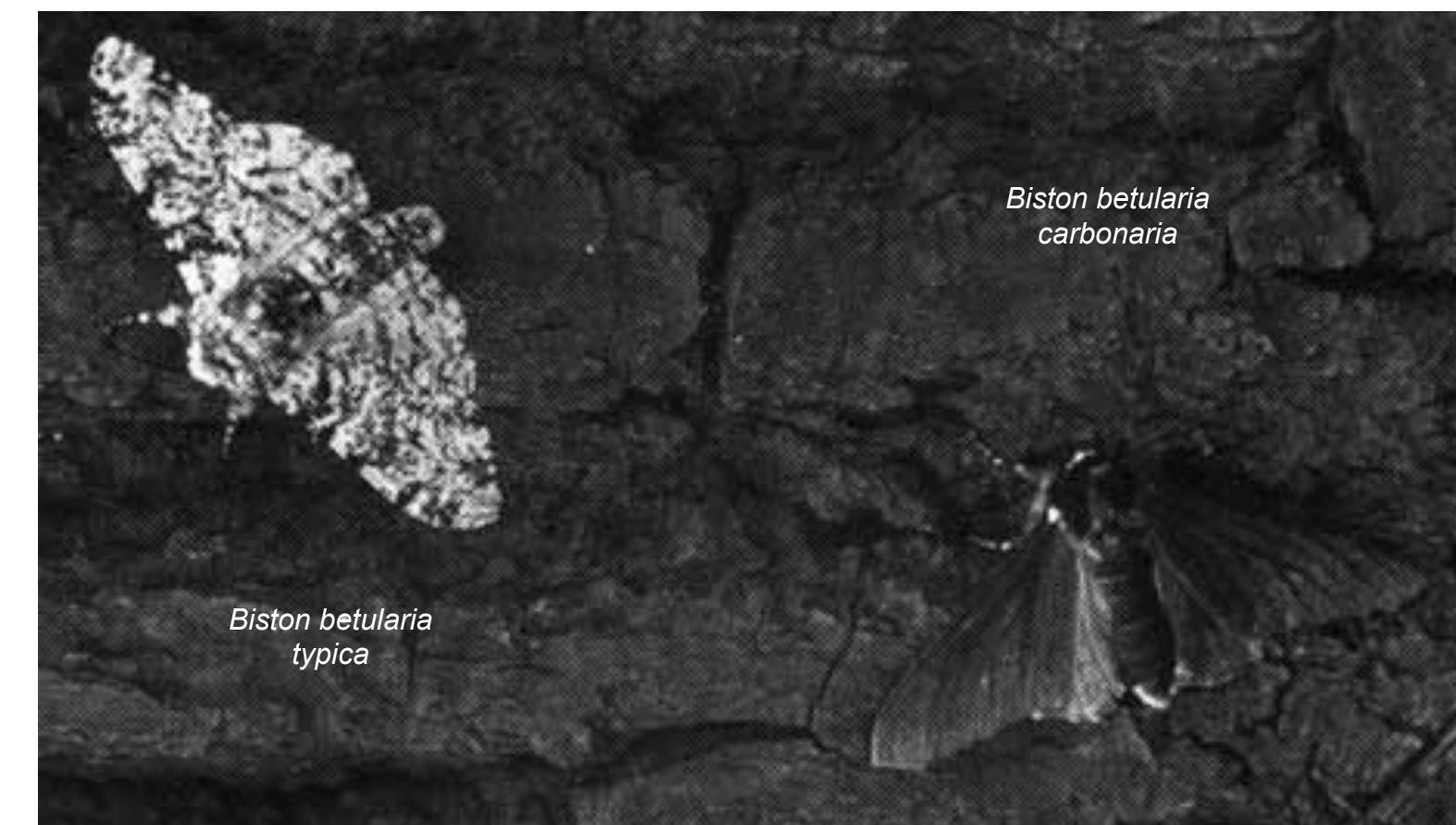


Mutations may be **deleterious**...



Bhattacharyya *et al. Cell* 1990



DNA transposon

# Transposable Elements

- Transposable Elements (TEs, "jumping genes") are an important source of mutations

- TEs transpose by cut-and-paste or copy-and-paste mechanisms

- BUT: most TEs are degraded and do not transpose
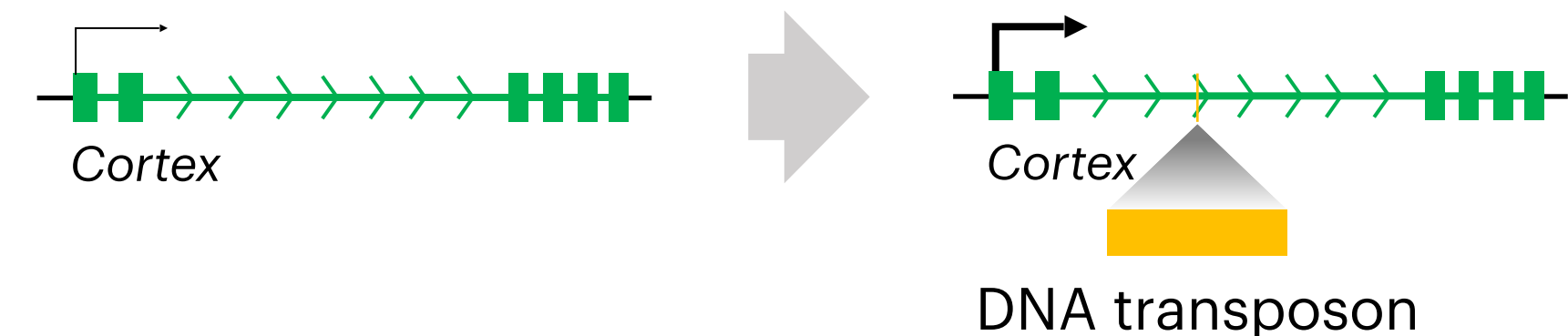


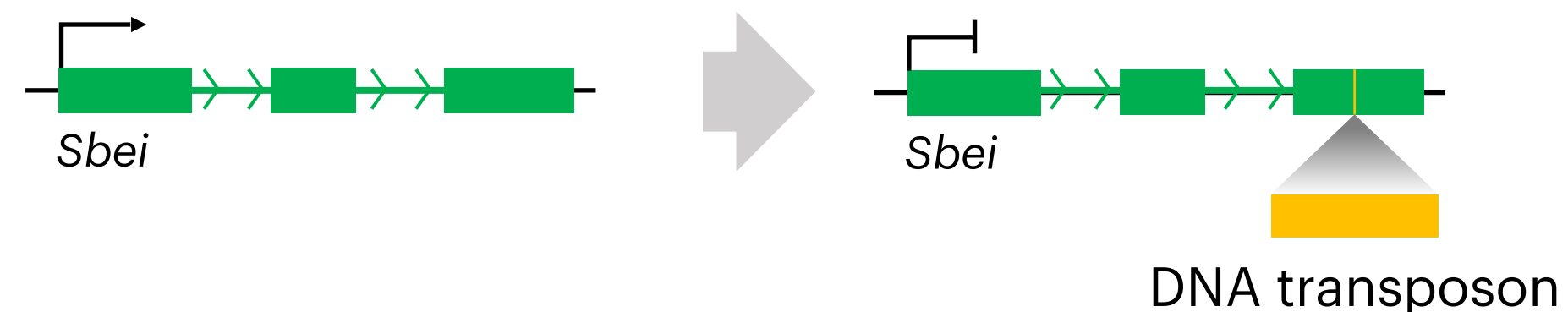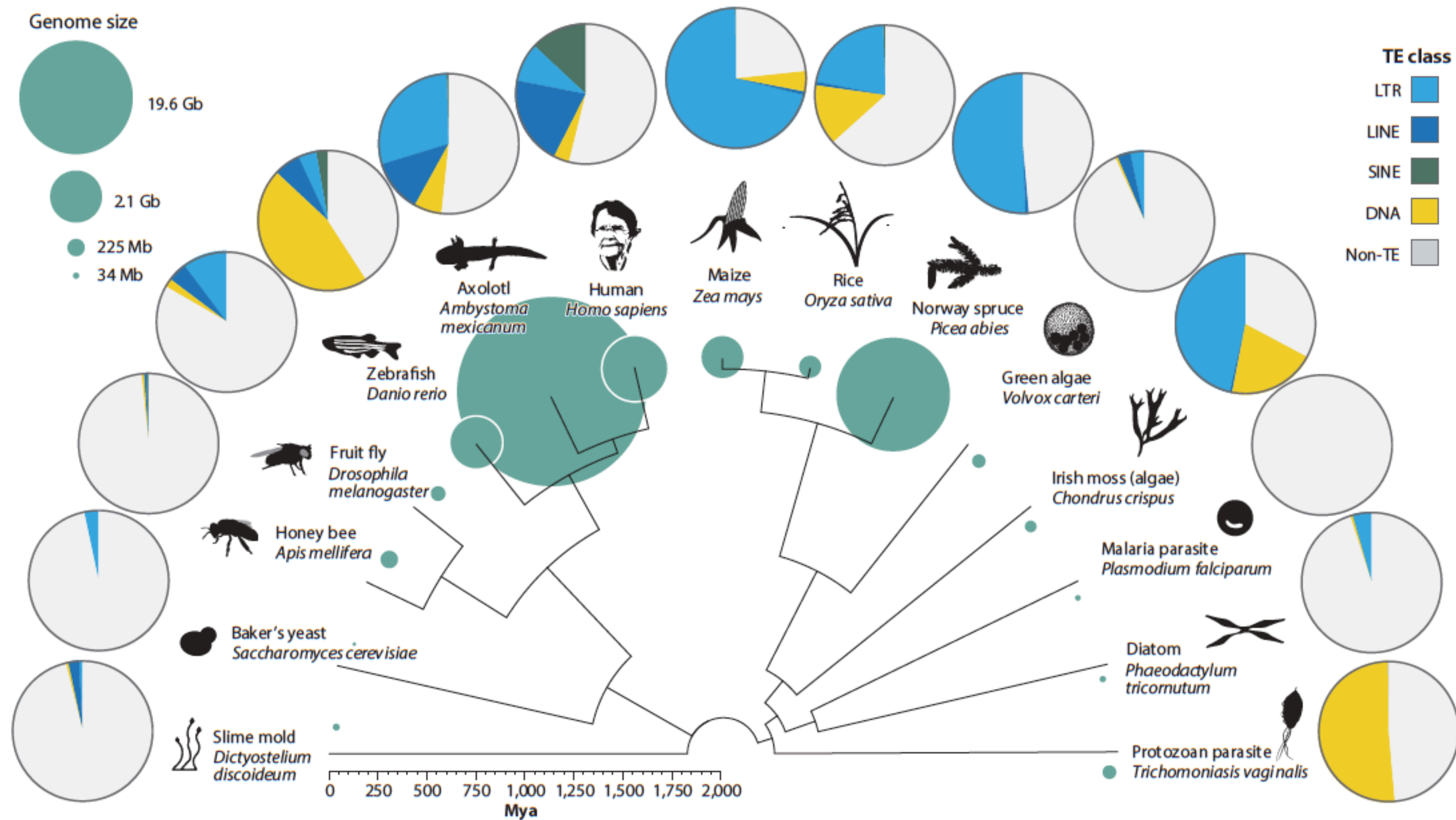## Mutations may be **deleterious**...



Bhattacharyya *et al.* *Cell* 1990

## ...yet sometimes **adaptive**



*Biston betularia carbonaria*

*Biston betularia typica*

Kettelwell. *Heredity* 1956; van't Hof *et al.* *Nature* 2016



*Sbei* → *Sbei*

DNA transposon

*Cortex* → *Cortex*

DNA transposon

# Transposable Elements



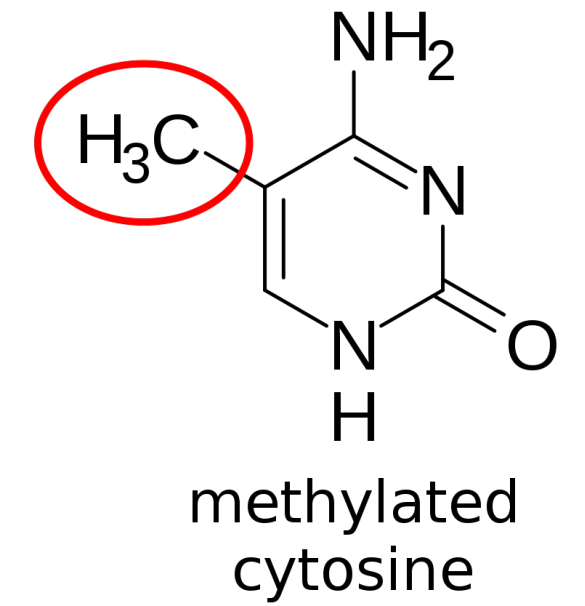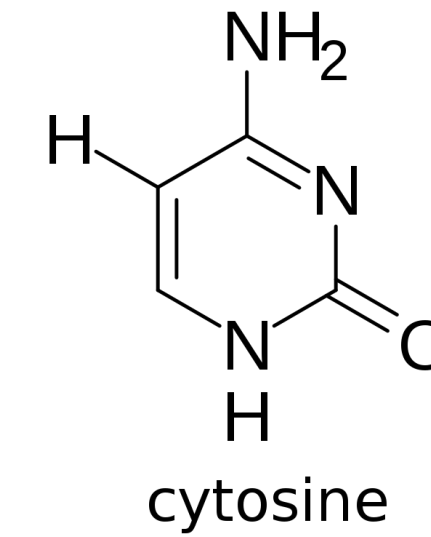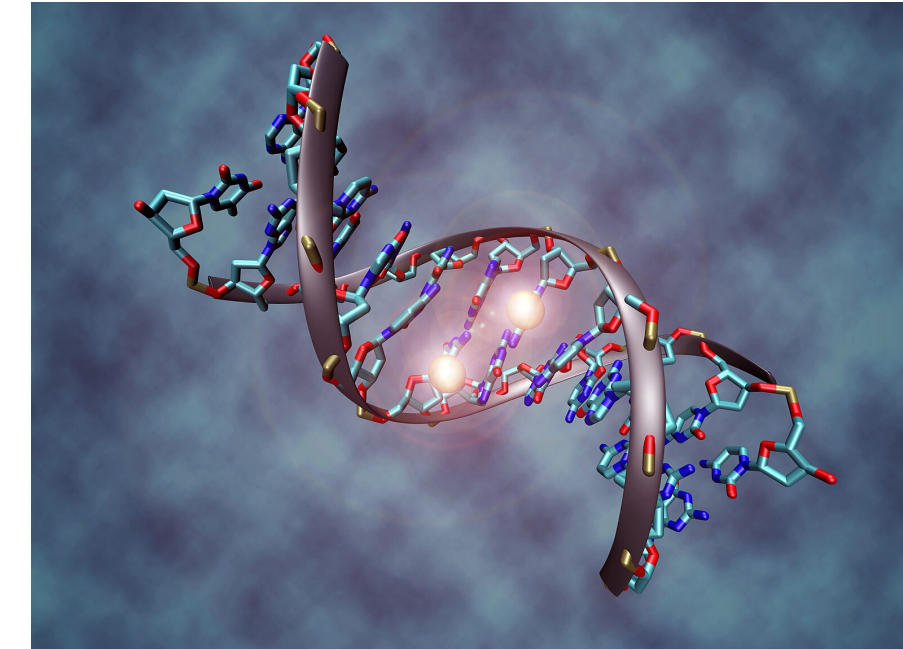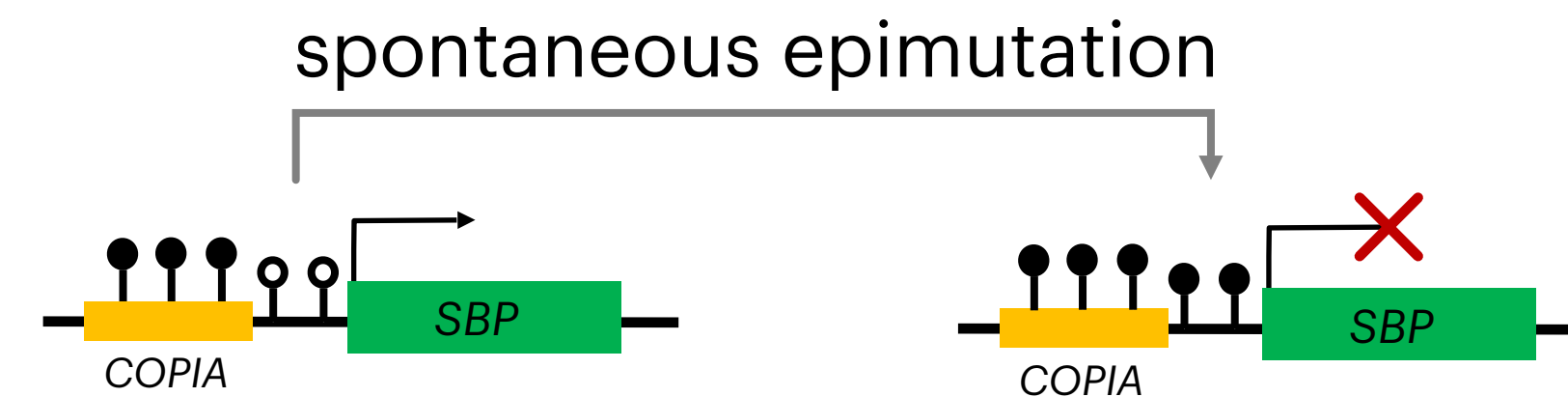Wells & Feschotte, *Annual Review of Genetics 2020*

# Epigenetic Regulation of Transposable Elements

**DNA methylation:**

- is an essential regulatory mechanism of TEs activity

- targets CG / CHG / CHH in plants
  [H = anything besides G]

- affects TE / gene expression (silencing)

- may spread to flanking regions

- example:

  methylated promoter $\implies$ no RNA $\implies$

  $\implies$ no protein $\implies$ no function

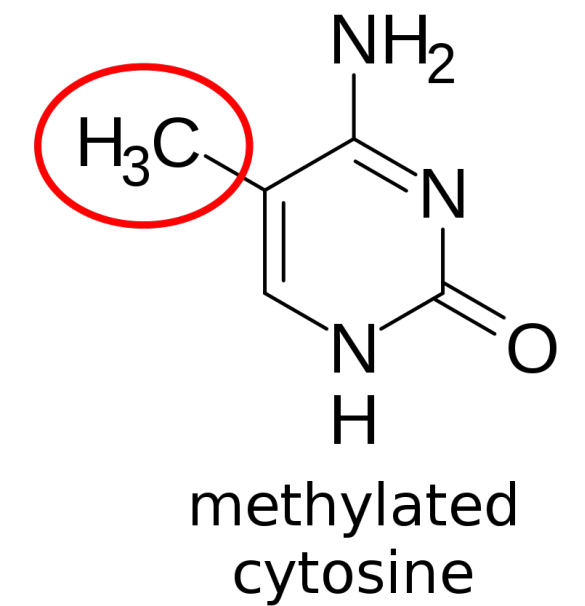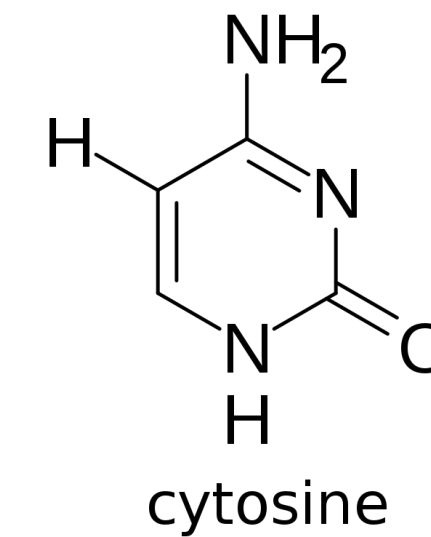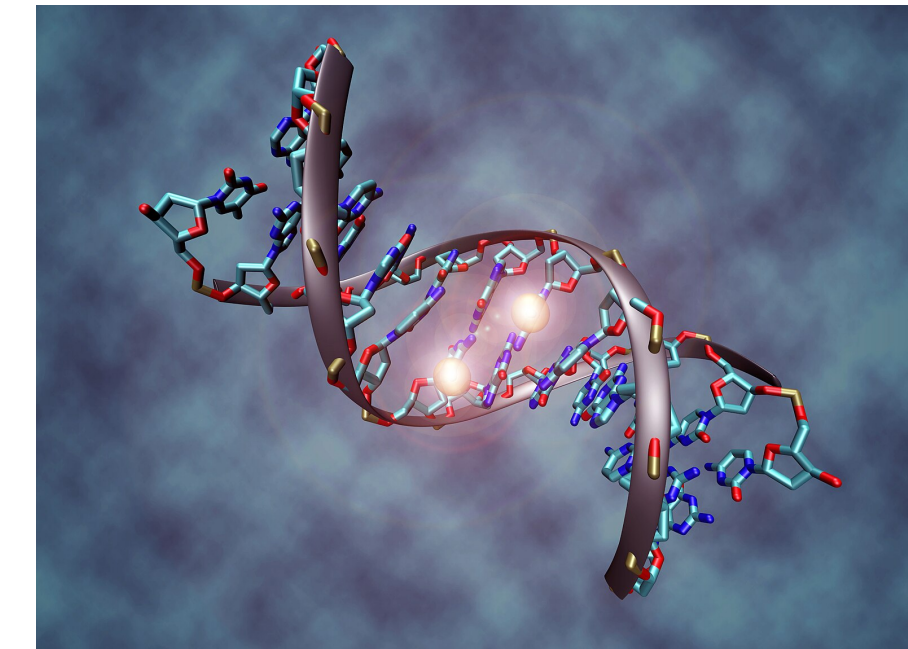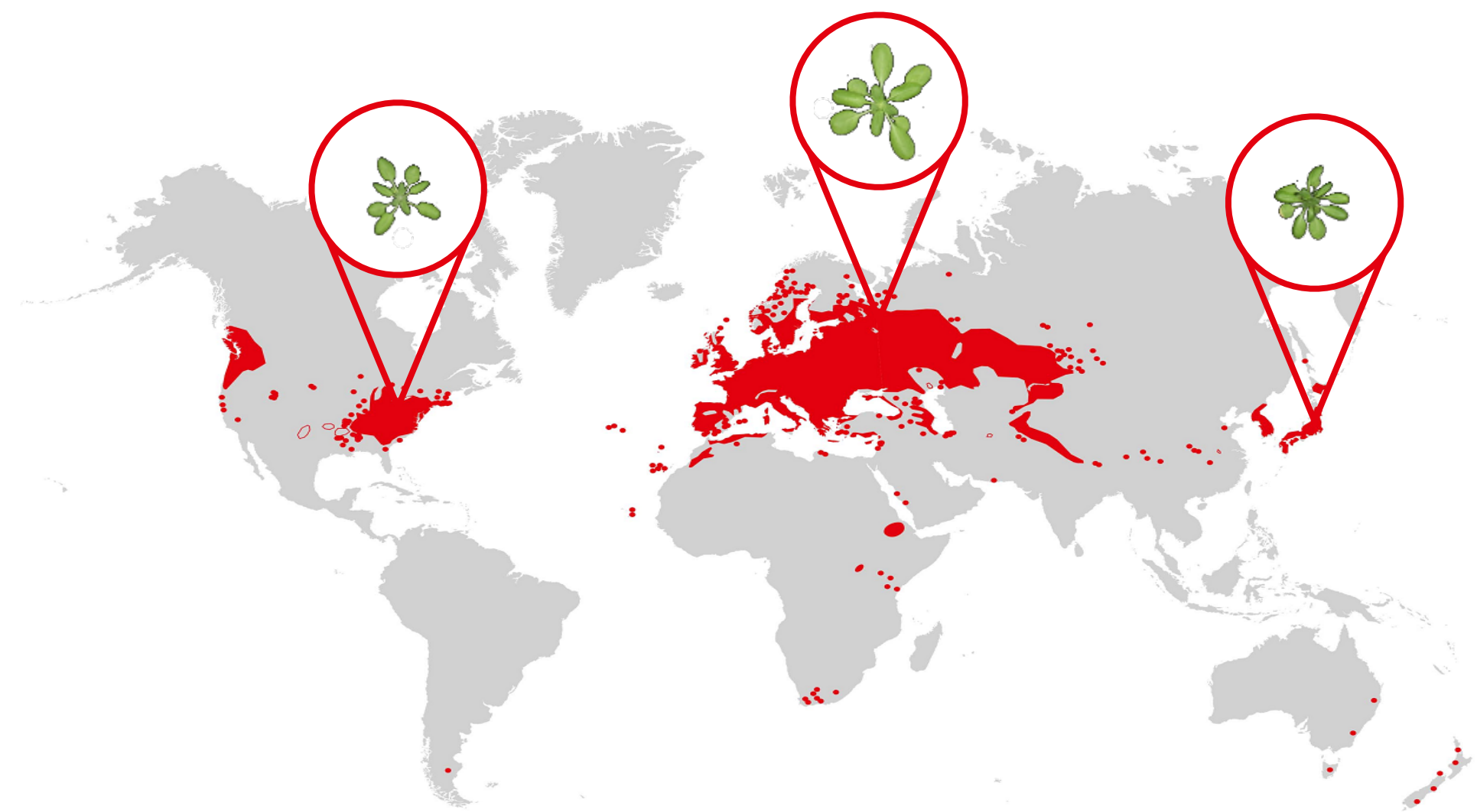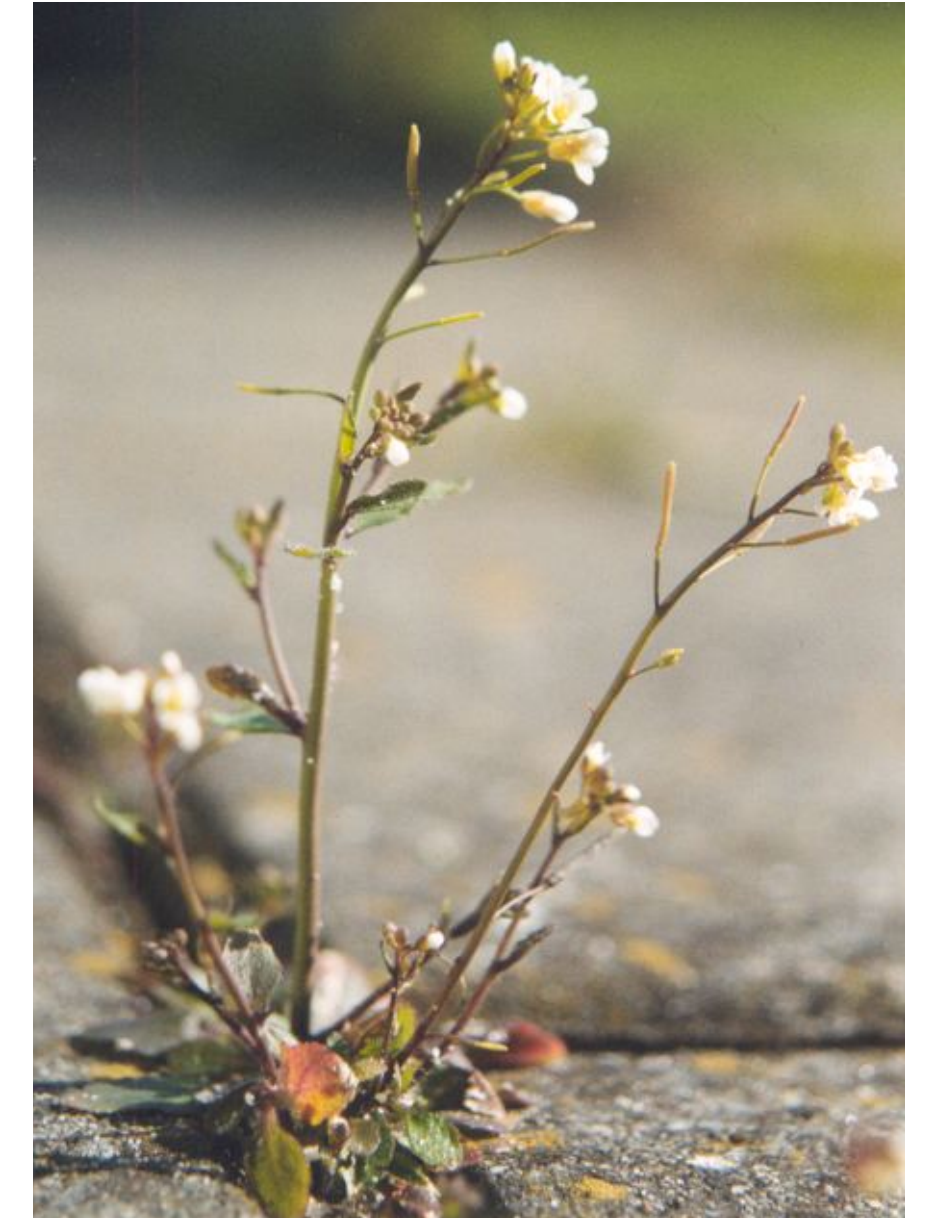cytosine

methylated
cytosine

# Epigenetic Regulation of Transposable Elements

**DNA methylation:**

- is an essential regulatory mechanism of TEs activity

- targets CG / CHG / CHH in plants
  [H = anything besides G]

- affects TE / gene expression (silencing)

- may spread to flanking regions

- example:

  methylated promoter $\Longrightarrow$ no RNA $\Longrightarrow$

  $\Longrightarrow$ no protein $\Longrightarrow$ no function

cytosine

methylated cytosine

spontaneous epimutation

COPIA    SBP

COPIA    SBP

cnr

Manning et al., *Nat Genet 2006*

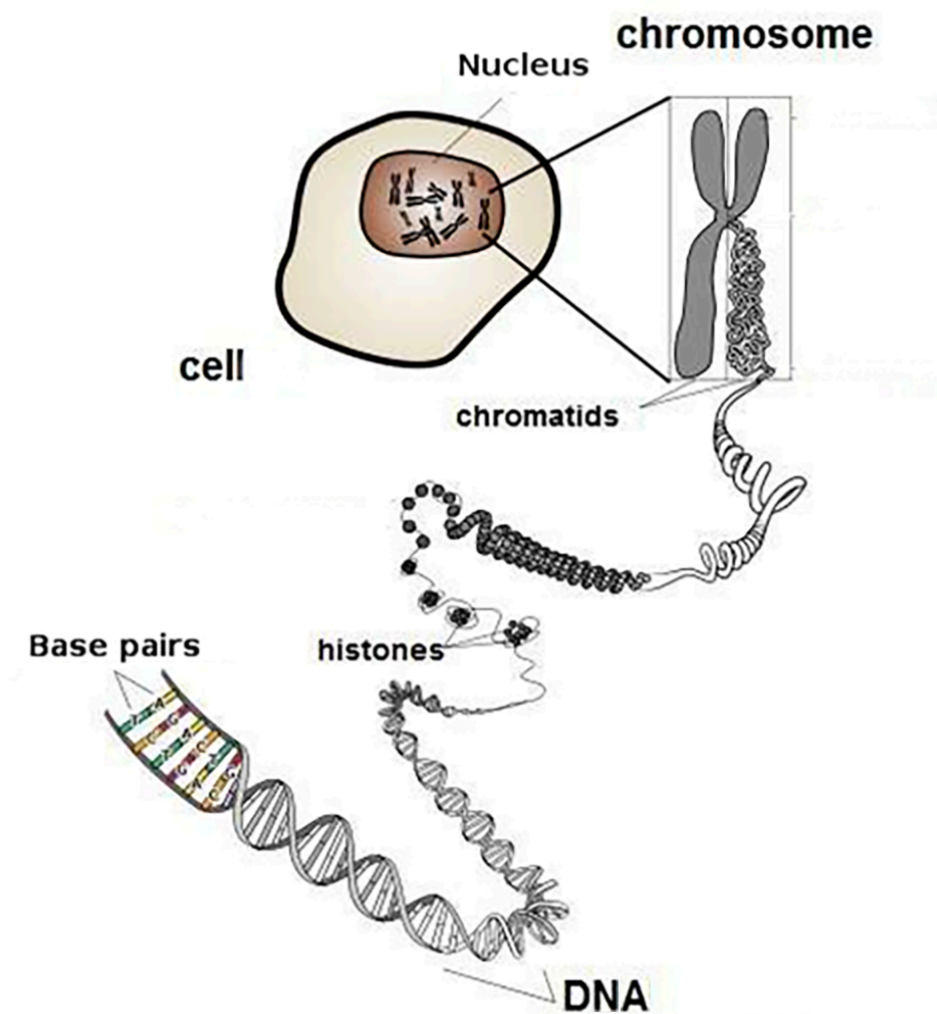$\Longrightarrow$ perfect Mendelian segregation though no DNA changes observed

# Our data: *Arabidopsis Thaliana*

- 89 strains from throughout the world, **sequenced with ultra-long reads (Nanopore)**

- **TE annotation + Full methylation profiles** (for all contexts CG, CHG, CHH)
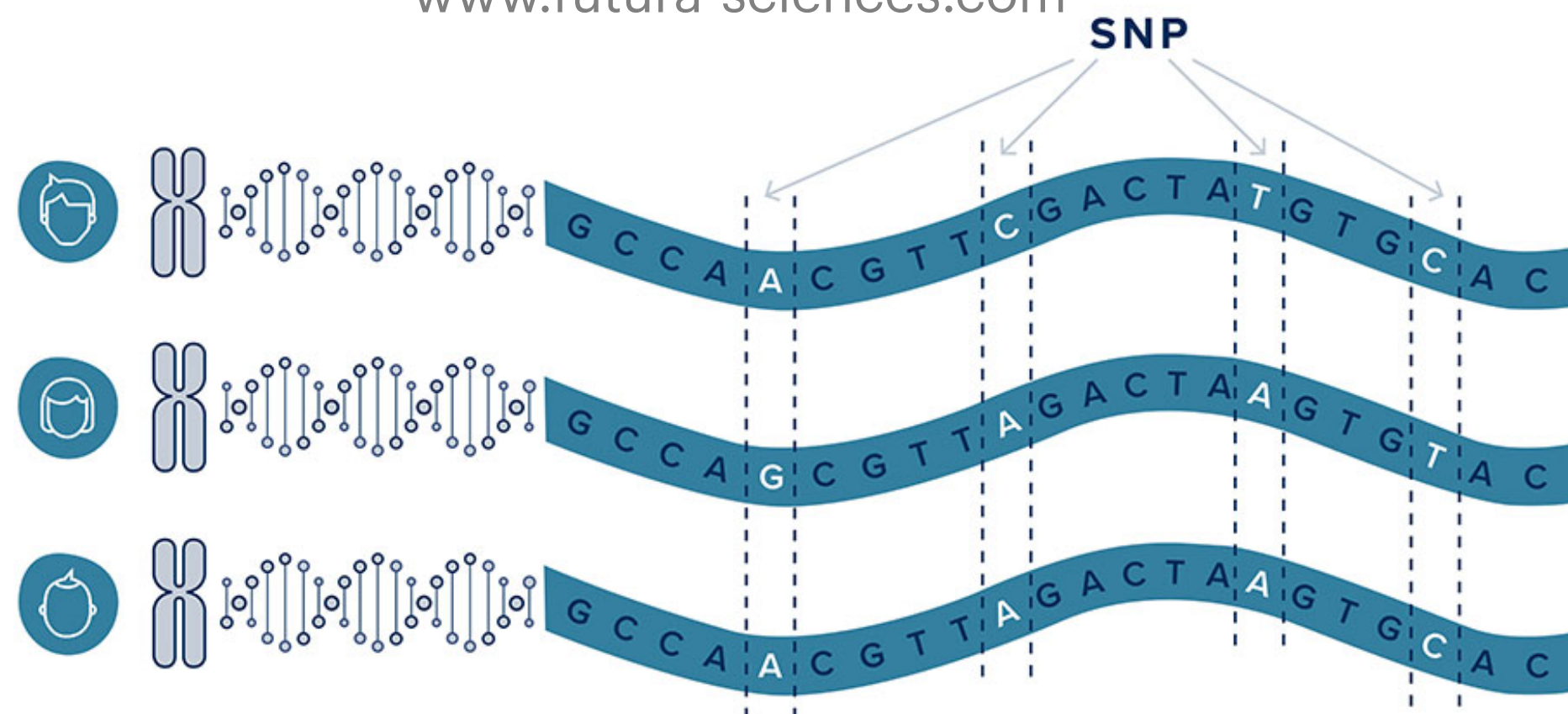
- Gene expression data



Kawakatsu *et al. Cell 2016*, Alonso-Blanco *et al. Cell 2016*, Quadrana *et al. eLife 2016*
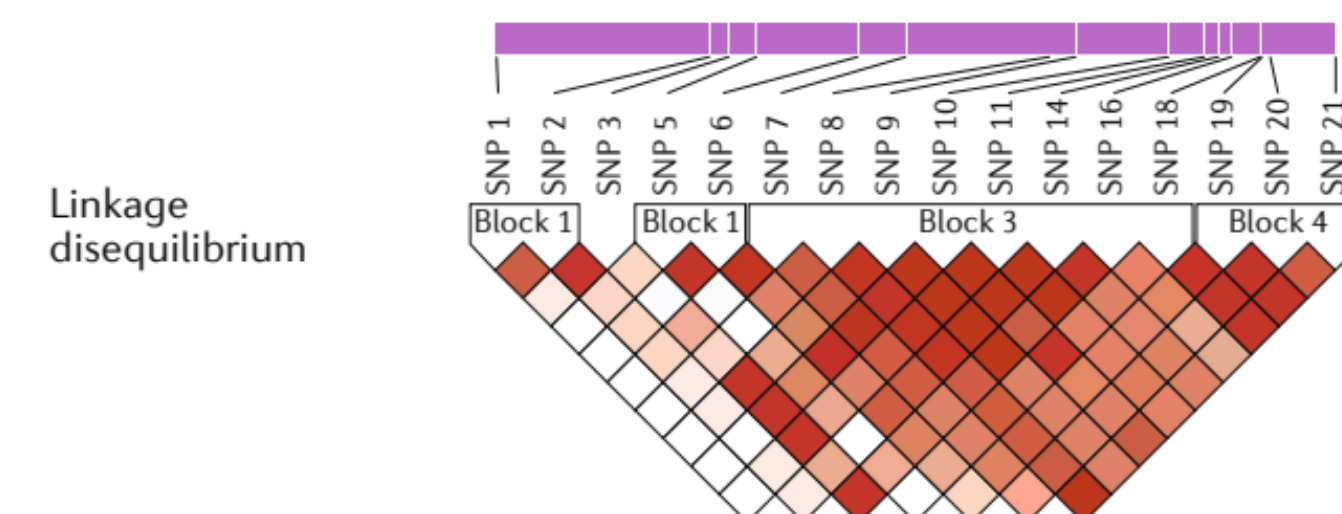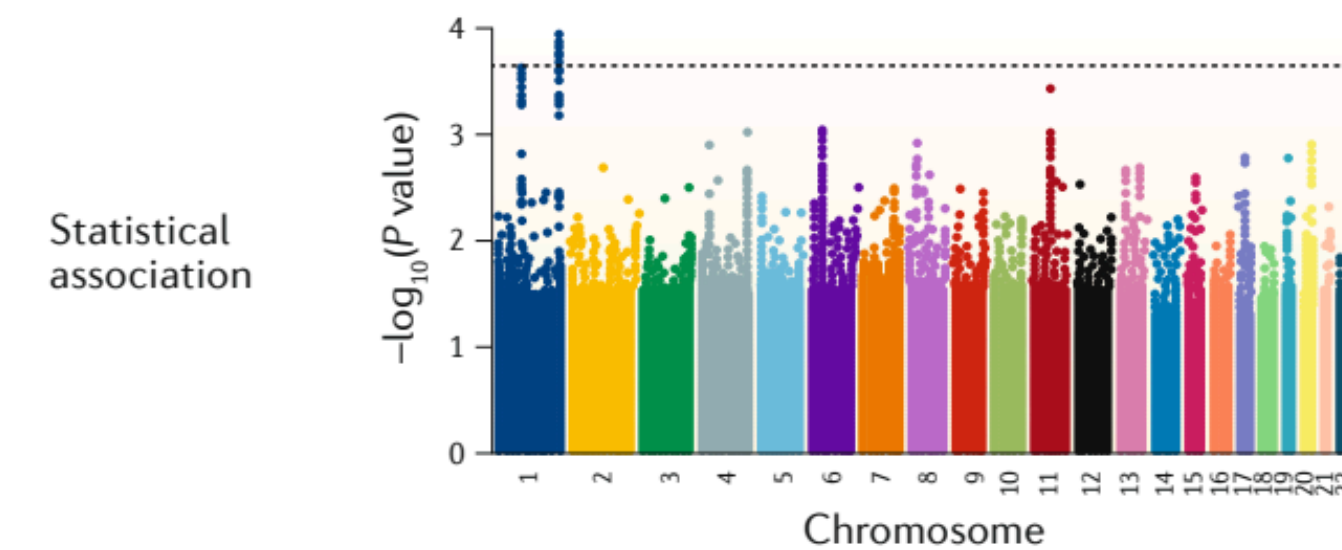
# From genotype to phenotype



Genome-Wide Association Study

www.futura-sciences.com

Scientific DX GmbH, 2020

Tam et al., *Nature Reviews Genetics 2019*

# From epi-genotype to phenotype

# From epi-genotype to phenotype

# From epi-genotype to phenotype

Genome 1

Genome 2

Genome 3

**Genome-Wide Association Study**

**3 groups:**
0 = absent
1 = present and not methylated (< 5%)
2 = present and methylated (> 5%)

|  | Gene A | Gene B | Gene C |
|---|---|---|---|
| **Genome 1** | | | |
| **Genome 2** | | | |
| **Genome 3** | | | |

## For the moment, cis- effects only:

|  | P_tip | P_meth | TIP | Chr | start | end | Distance from gene |
|---|---|---|---|---|---|---|---|
| **2780** | 0.516462 | 0.000002 | fixed.DEL6462 | Chr3 | 9783357 | NaN | 0.0 |



fixed.DEL6462 (0.0 bp away from AT3G26612)

# From epi-genotype to phenotype



Genome 1

Genome 2

Genome 3

**Genome-Wide Association Study**

**3 groups:**
0 = absent
1 = present and not methylated (< 5%)
2 = present and methylated (> 5%)

|  | Gene A | Gene B | Gene C |
| --- | --- | --- | --- |
| **Genome 1** |  |  |  |
| **Genome 2** |  |  |  |
| **Genome 3** |  |  |  |

## For the moment, cis- effects only:

|  | P_tip | P_meth | TIP | Chr | start | end | Distance from gene |
| --- | --- | --- | --- | --- | --- | --- | --- |
| **2780** | 0.516462 | 0.000002 | fixed.DEL6462 | Chr3 | 9783357 | NaN | 0.0 |

### fixed.DEL6462 (0.0 bp away from AT3G26612)

# From **epi**-genotype to phenotype



Genome 1

Genome 2

Genome 3

**Genome-Wide Association Study**

**3 groups:**
0 = __absent__
1 = __present__ and __not methylated__ (< 5%)
2 = __present__ and __methylated__ (> 5%)

|  | Gene A | Gene B | Gene C |
|---|---|---|---|
| Genome 1 |  |  |  |
| Genome 2 |  |  |  |
| Genome 3 |  |  |  |

## For the moment, cis- effects only:

|  | P_tip | P_meth |  | TIP | Chr | start | end | Distance from gene |
|---|---|---|---|---|---|---|---|---|
| 2780 | 0.516462 | 0.000002 |  | fixed.DEL6462 | Chr3 | 9783357 | NaN | 0.0 |



fixed.DEL6462 (0.0 bp away from AT3G26612)

|  | P_tip | P_meth |  | TIP | Chr | start | end | Distance from gene |
|---|---|---|---|---|---|---|---|---|
| 535 | 0.000068 | 0.002327 |  | IP_Her12.svim_asm.DEL.774 | Chr2 | 15110051 | 15110322.0 | 170.0 |



IP_Her12.svim_asm.DEL.774 (170.0 bp away from AT2G35980)

# From **epi**-genotype to phenotype



**Genome-Wide Association Study**

**3 groups:**
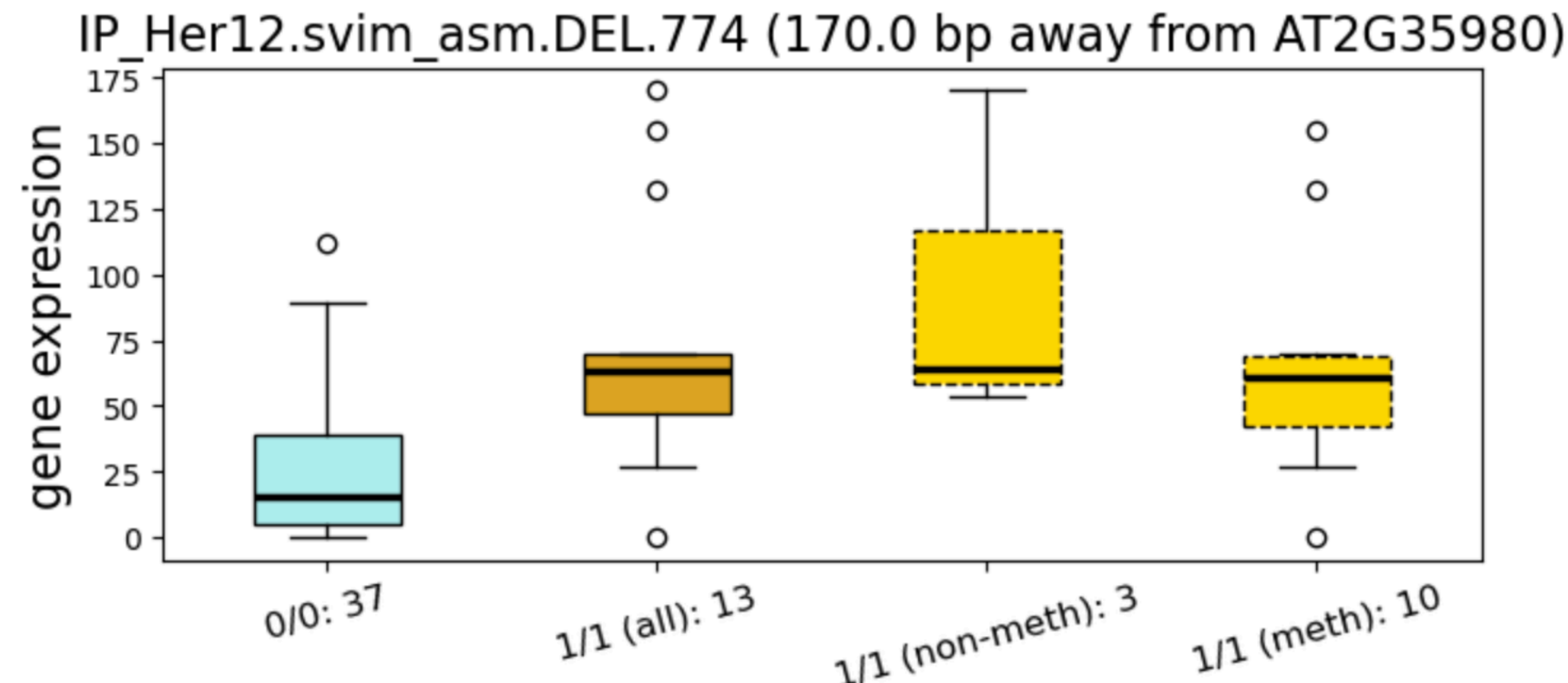0 = **absent**
1 = **present** and **not methylated** (< 5%)
2 = **present** and **methylated** (> 5%)

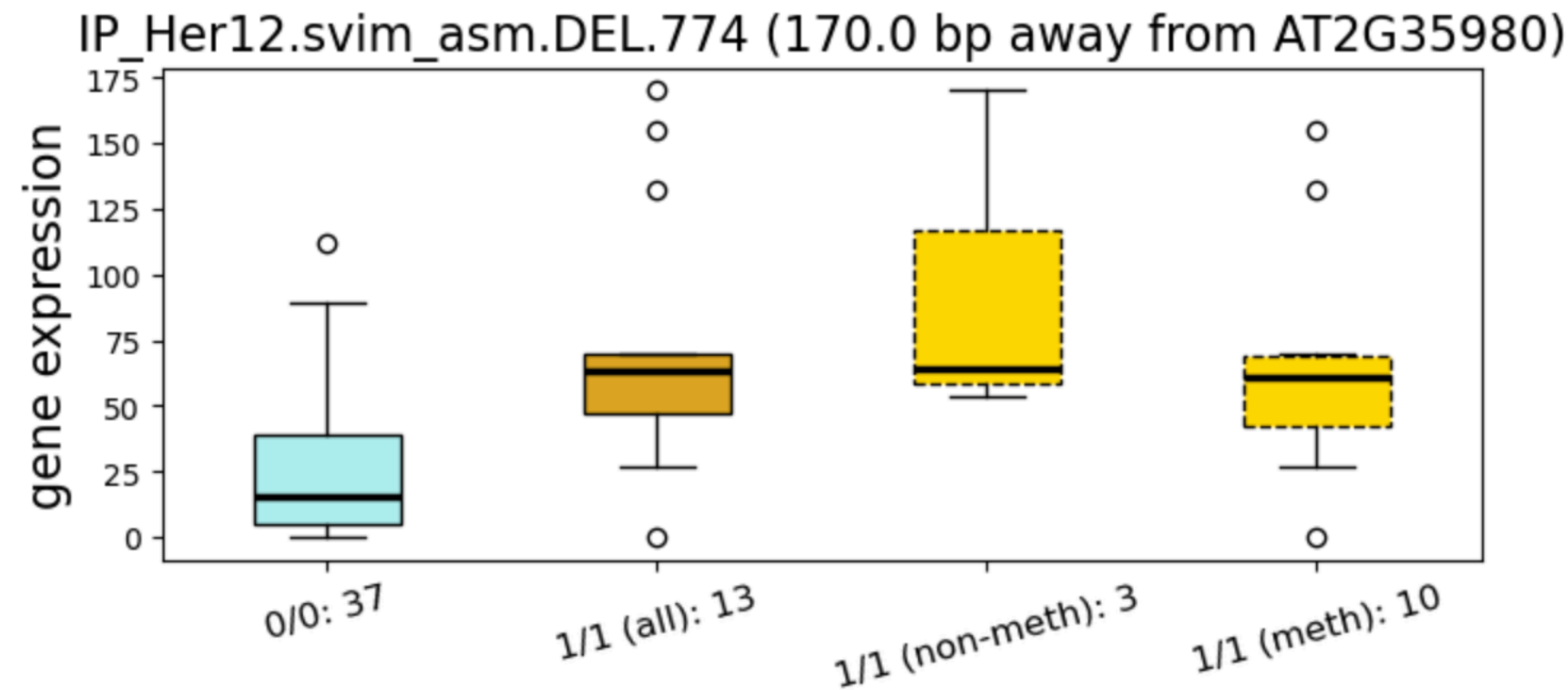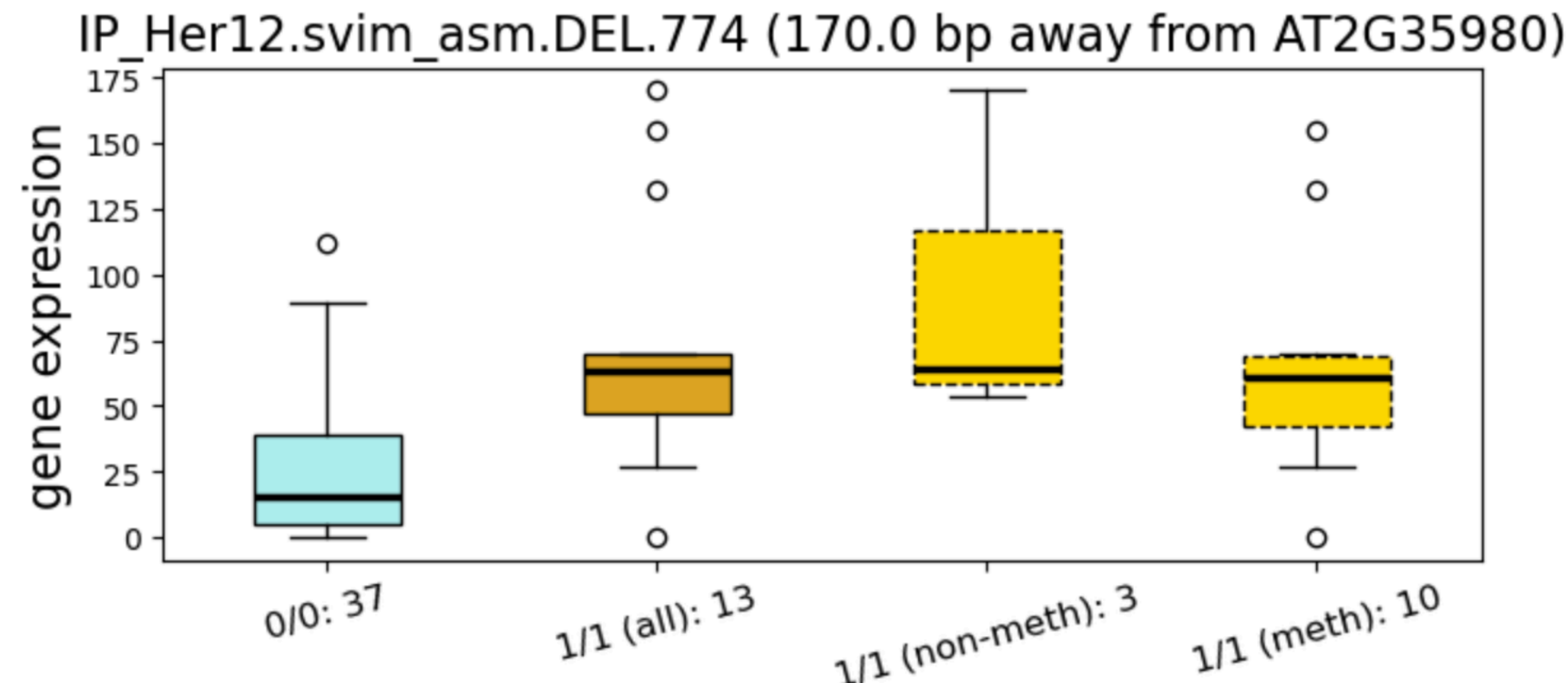|  | Gene A | Gene B | Gene C |
|---|---|---|---|
| Genome 1 | | | |
| Genome 2 | | | |
| Genome 3 | | | |

## For the moment, cis- effects only:

| | P_tip | P_meth | | TIP | Chr | start | end | Distance from gene |
|---|---|---|---|---|---|---|---|---|
| 2780 | 0.516462 | 0.000002 | | fixed.DEL6462 | Chr3 | 9783357 | NaN | 0.0 |

fixed.DEL6462 (0.0 bp away from AT3G26612)

0/0: 7    1/1 (all): 43    1/1 (non-meth): 18    1/1 (meth): 25

| | P_tip | P_meth | | TIP | Chr | start | end | Distance from gene |
|---|---|---|---|---|---|---|---|---|
| 535 | 0.000068 | 0.002327 | | IP_Her12.svim_asm.DEL.774 | Chr2 | 15110051 | 15110322.0 | 170.0 |

IP_Her12.svim_asm.DEL.774 (170.0 bp away from AT2G35980)

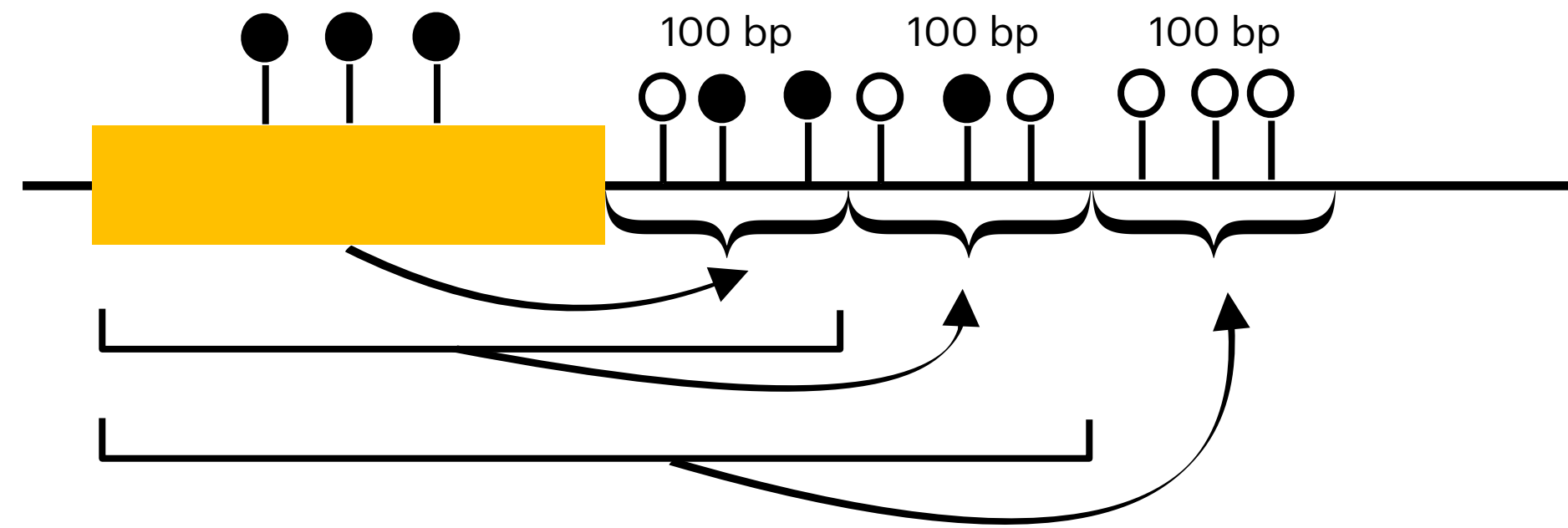0/0: 37    1/1 (all): 13    1/1 (non-meth): 3    1/1 (meth): 10

# Prediction of methylation spreading



**Model:**
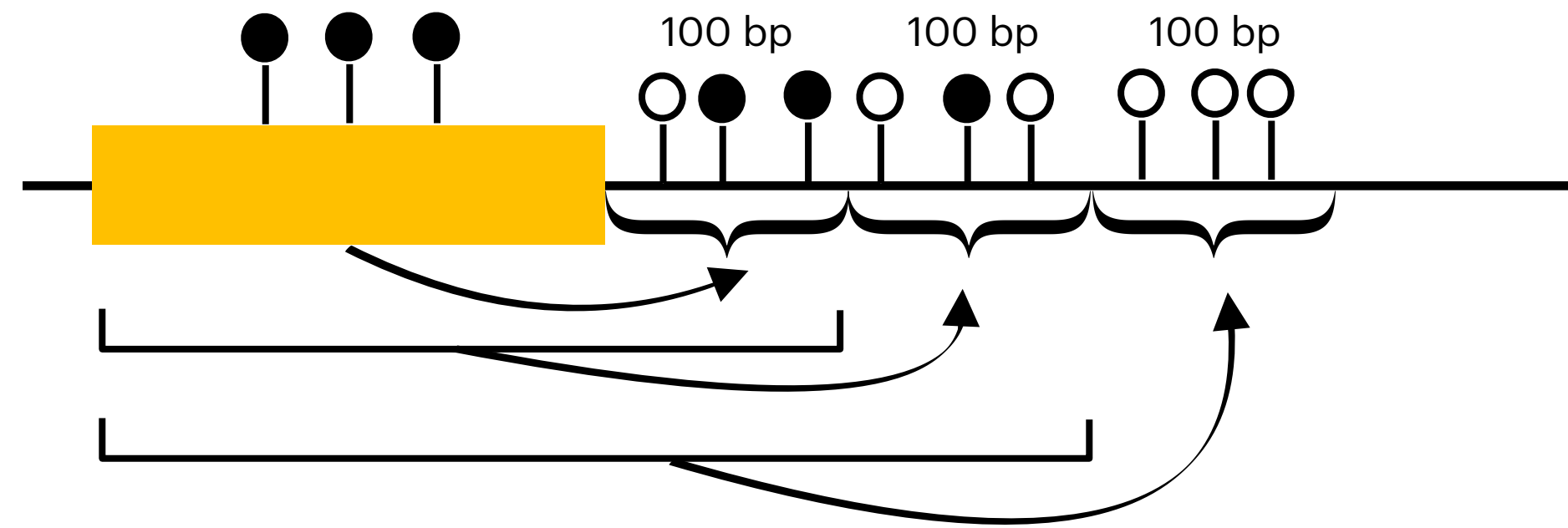Random Forest

# Prediction of methylation spreading



**Model:**
Random Forest

**Features:**
- **TE** (length, distance to pericentromere, superfamily, insertion frequency, divergence, if inside a gene)
- **Nearest 2 genes** (length, distance, relative direction)
- **Methylation** in CG, CHG, CHH contexts (average genome-wide, TE, on the ends of TE, previous windows)
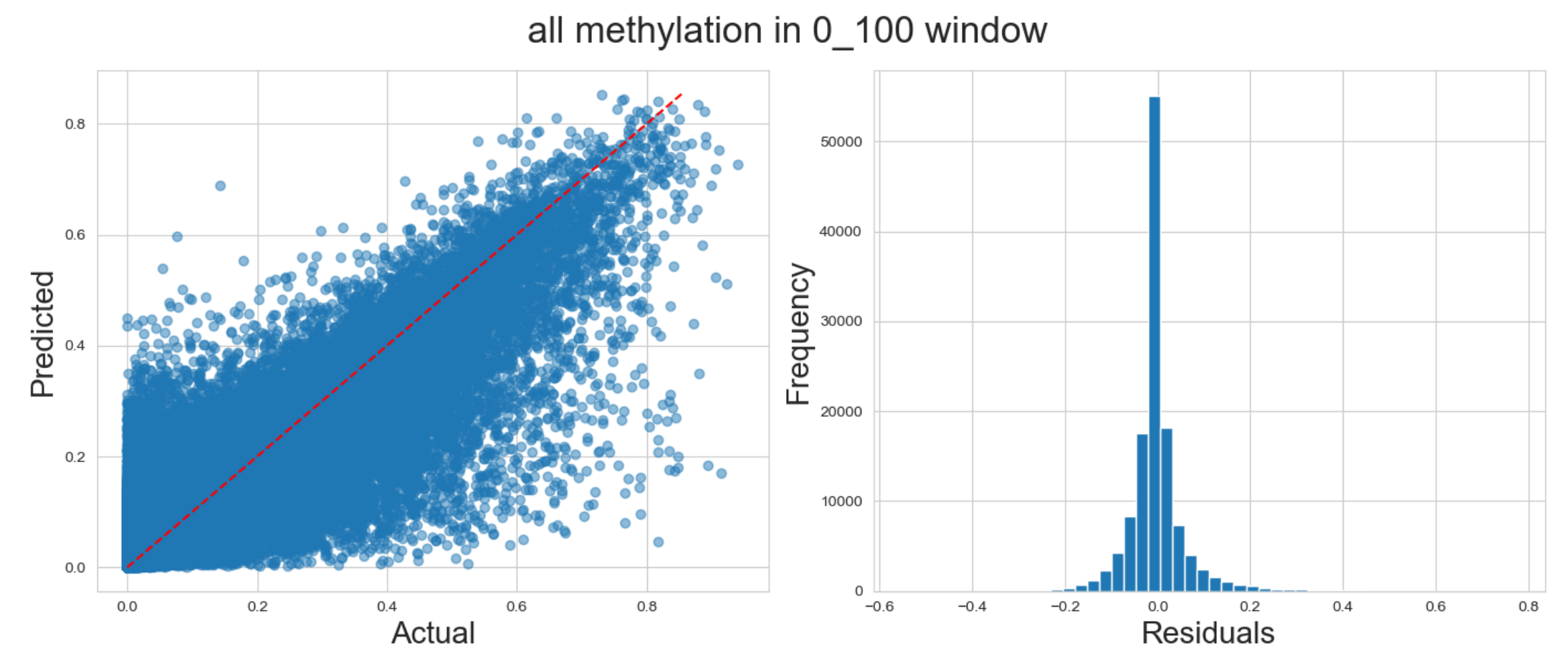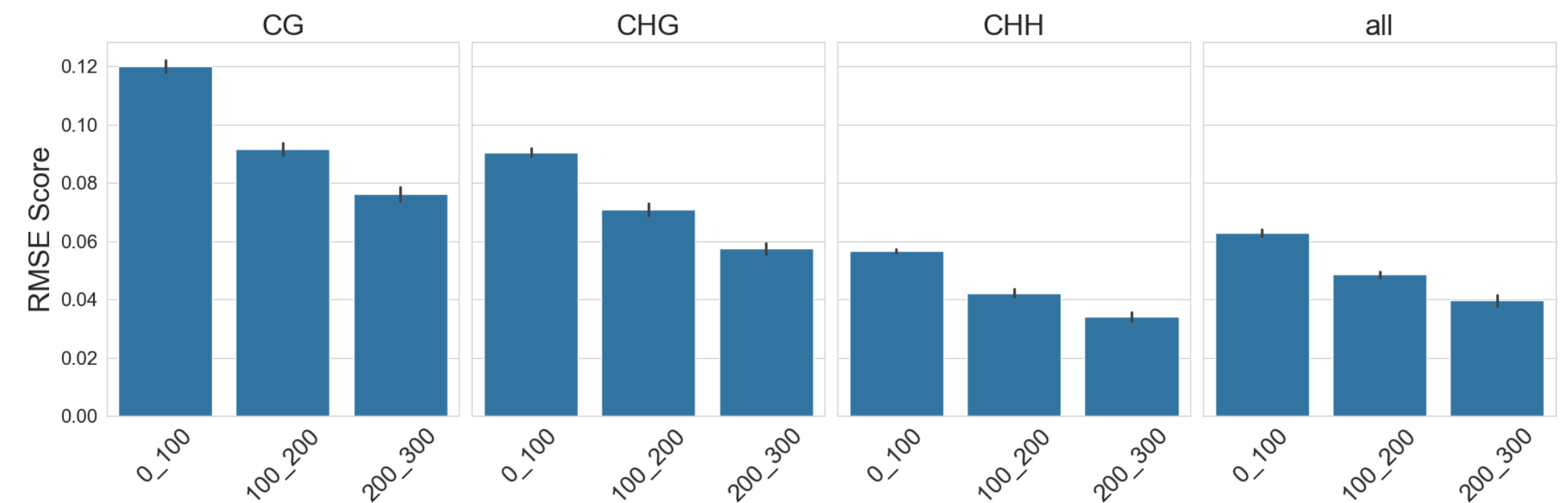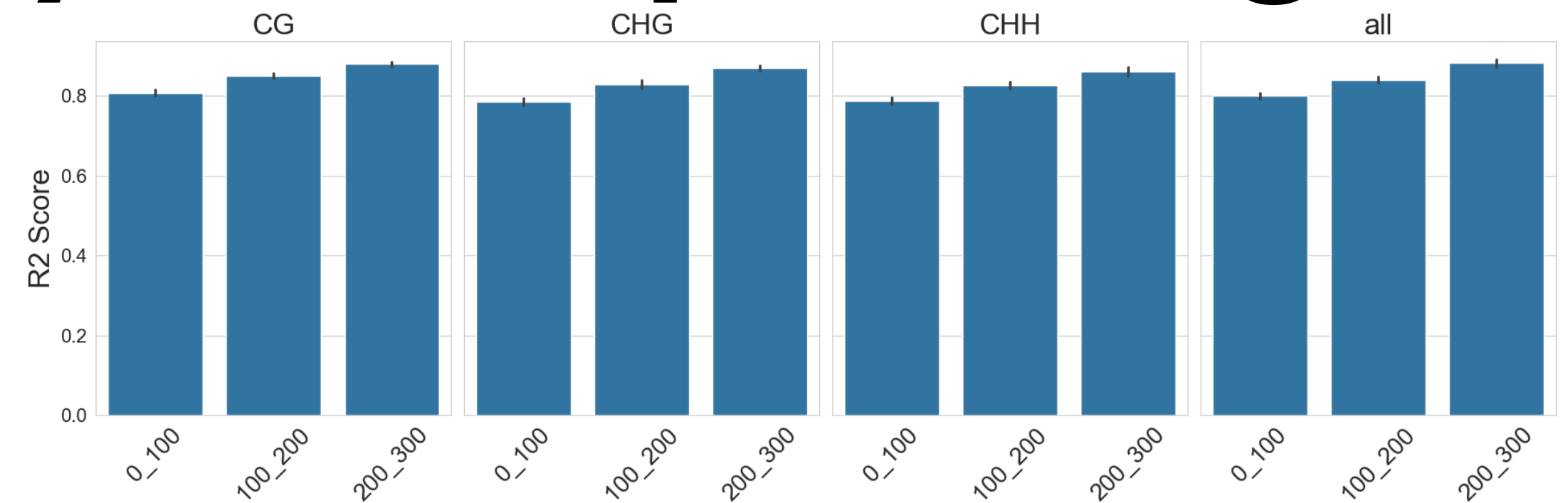- **Densities** of CG, CHG, CHH contexts

# Prediction of methylation spreading



**Model:**
Random Forest

**Features:**
- **TE** (length, distance to pericentromere, superfamily, insertion frequency, divergence, if inside a gene)
- **Nearest 2 genes** (length, distance, relative direction)
- **Methylation** in CG, CHG, CHH contexts (average genome-wide, TE, on the ends of TE, previous windows)
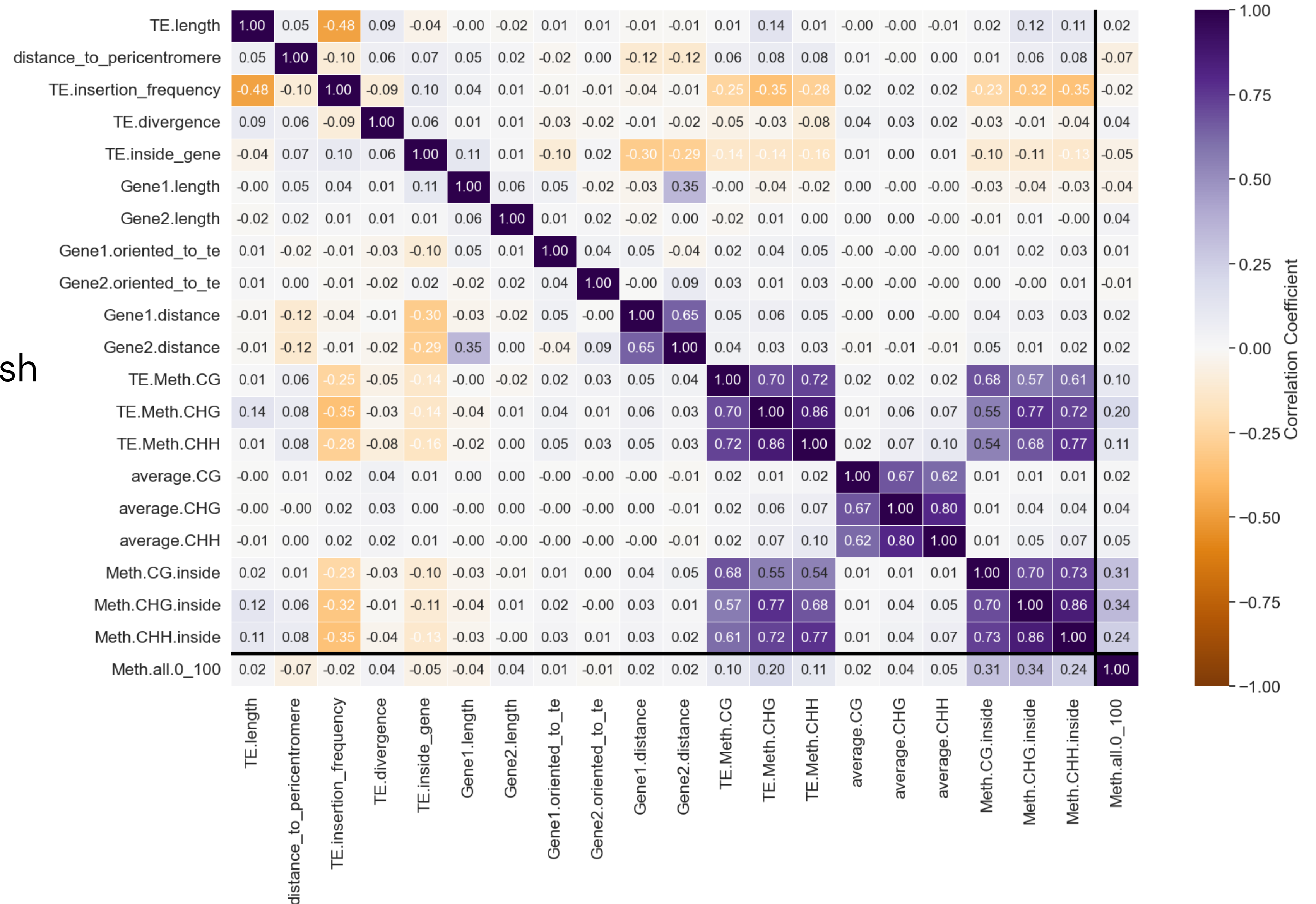- **Densities** of CG, CHG, CHH contexts

# Prediction of methylation spreading

- The model predicts well, but which features define methylation level?
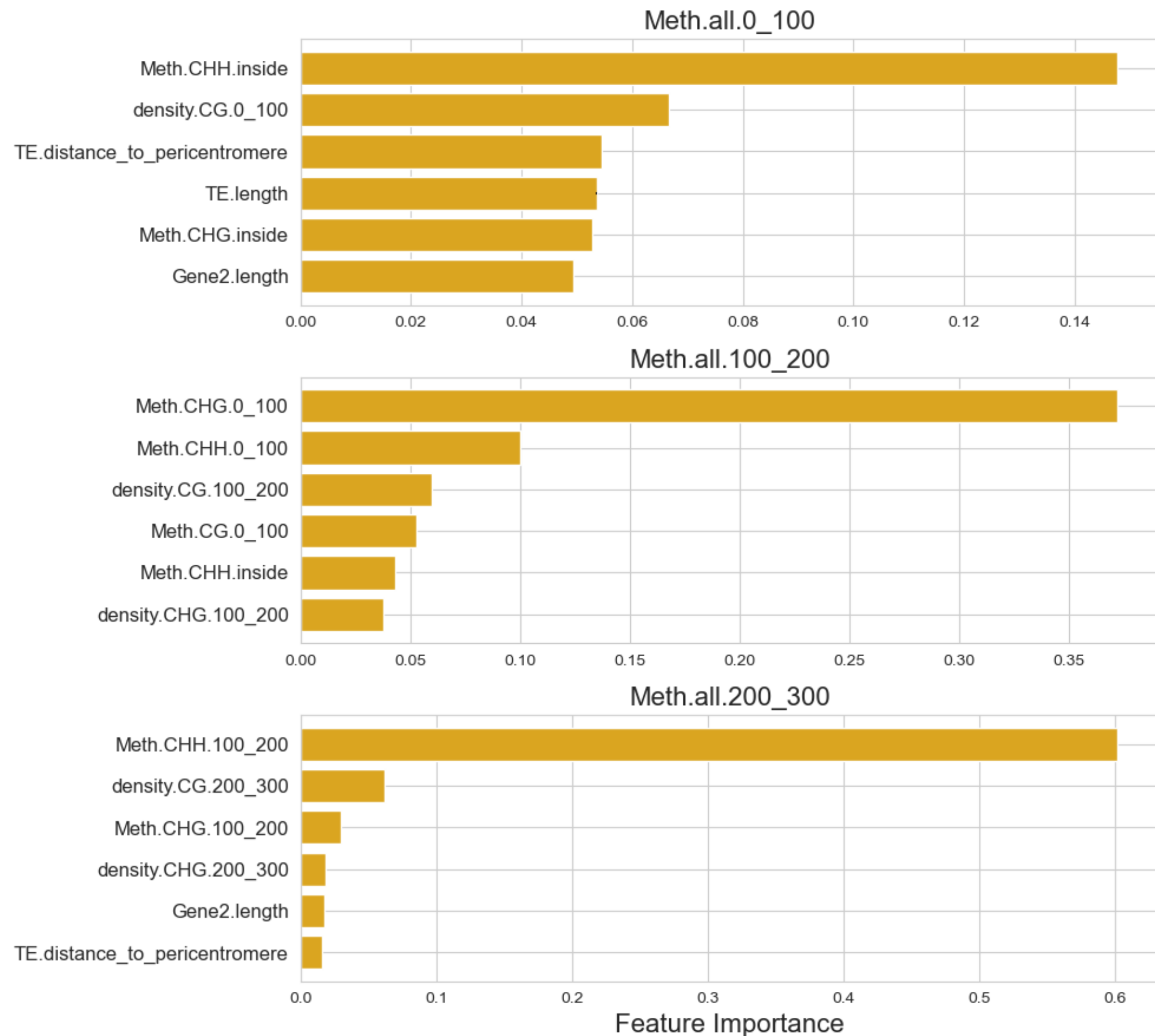
# Prediction of methylation spreading

- The model predicts well, but which features define methylation level?

- Some features are highly correlated $\implies$ hard to distinguish between them

# Prediction of methylation spreading
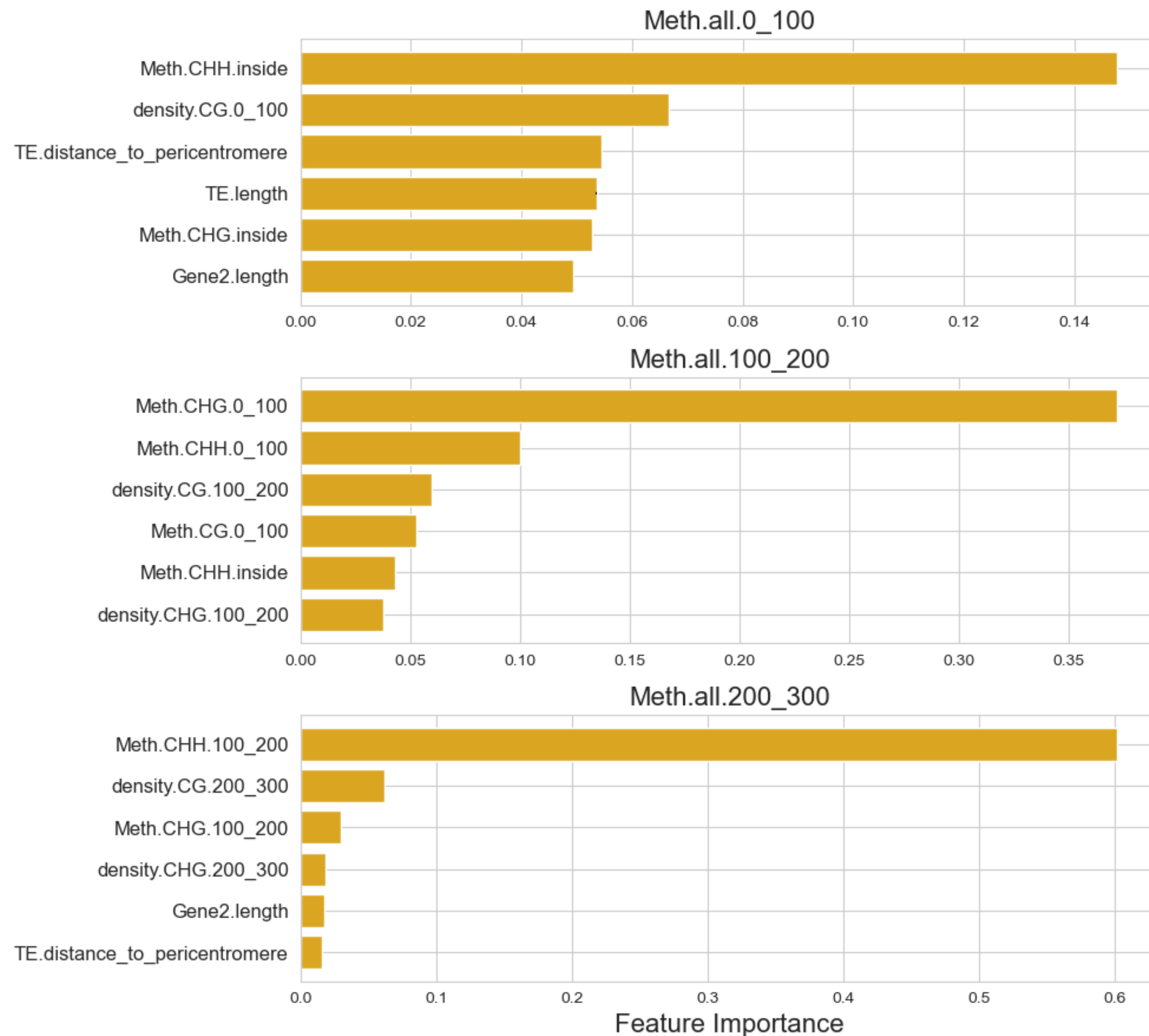
**Impurity-based feature importances**

* 10 independent runs with different random seeds
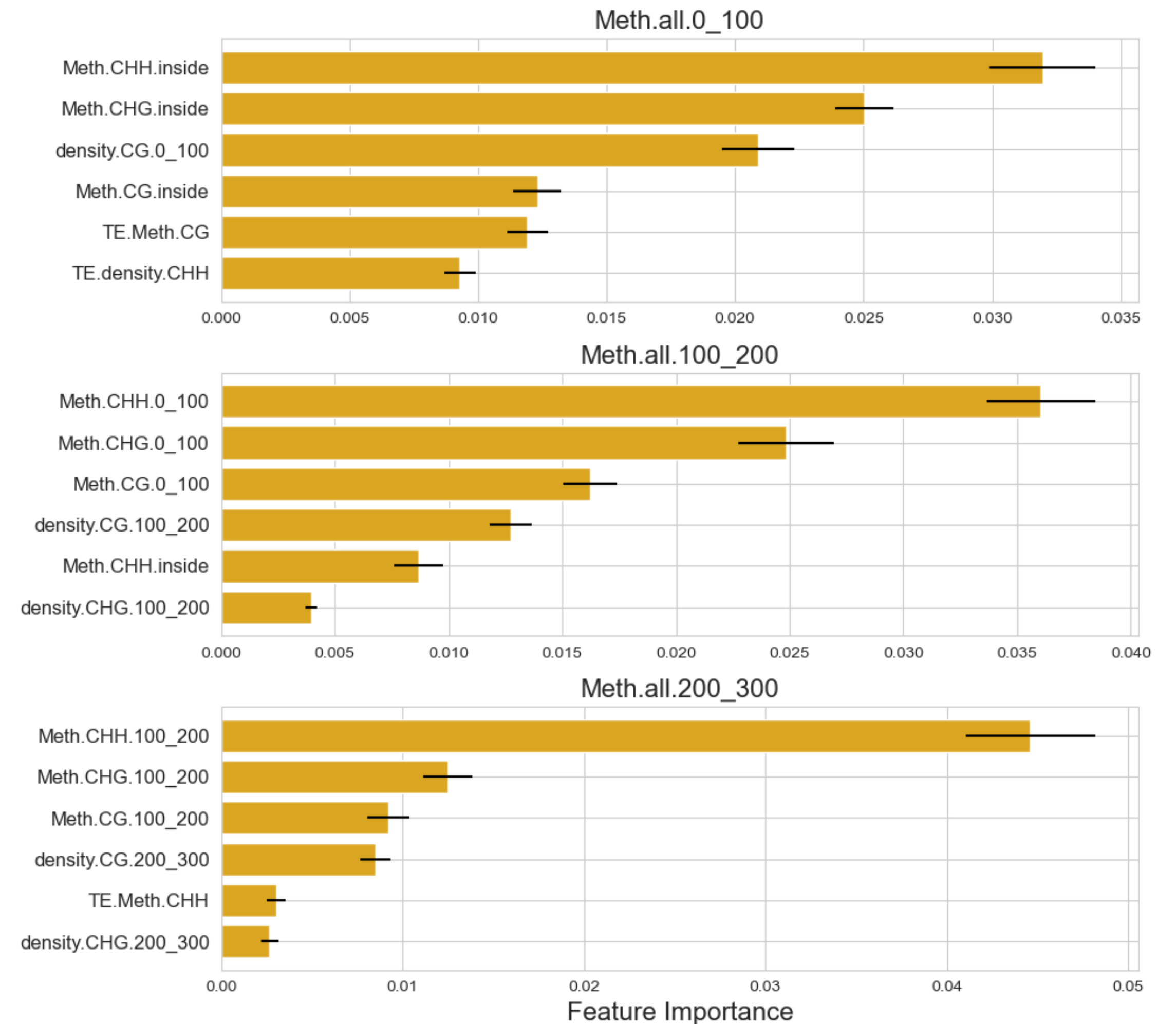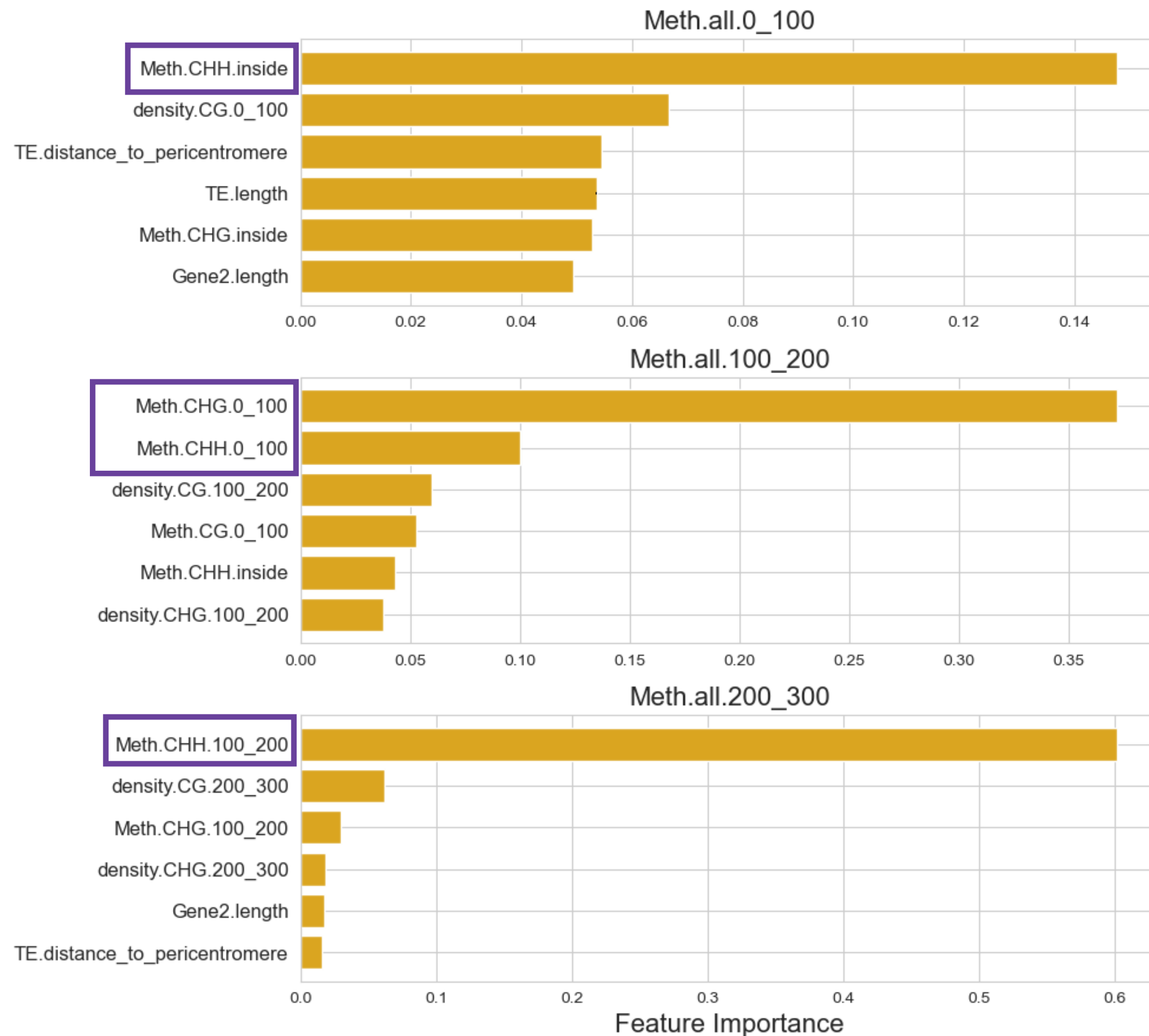
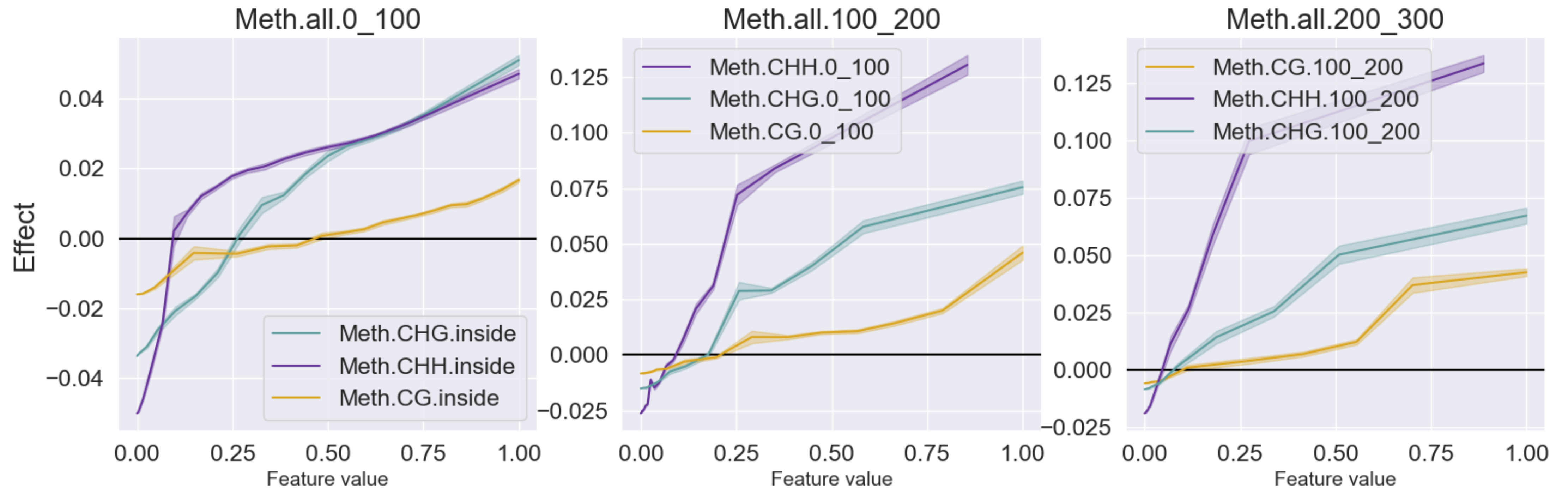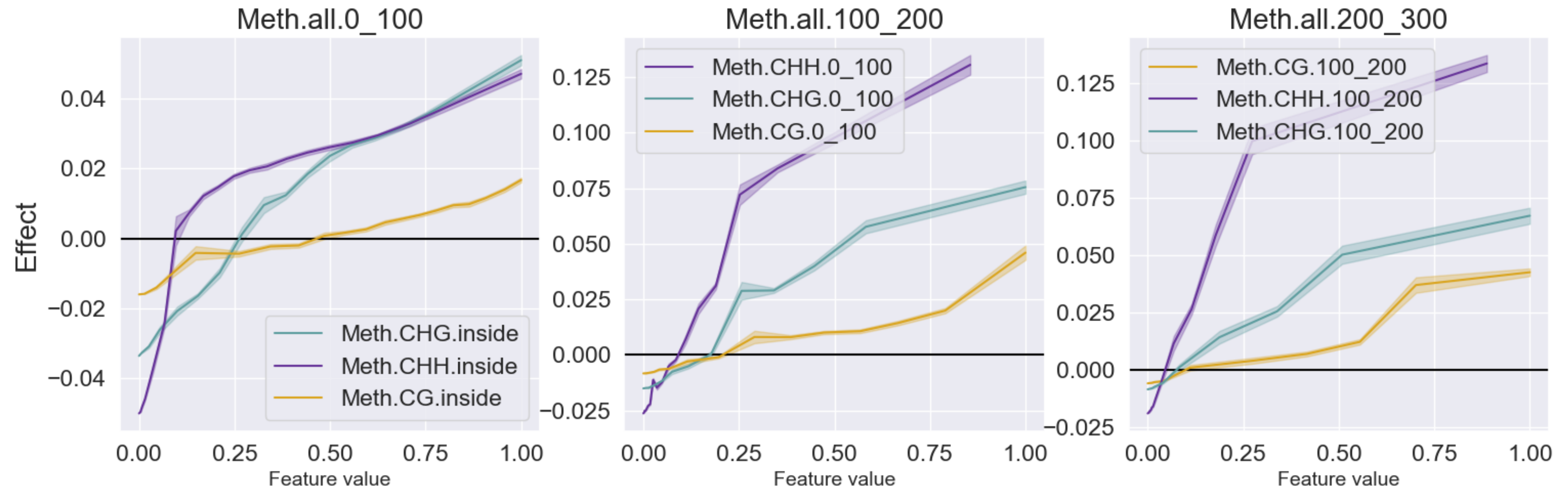# Prediction of methylation spreading

**Accumulated Local Effects (ALE)**

# Prediction of methylation spreading

**Accumulated Local Effects (ALE)**



**Conclusion:**

- ⊙ **Methylation of the TE inside ends** consistently comes as the most important feature with monotonous effect increase

- ⊙ TE is **methylated on the ends** $\Longrightarrow$ more **likely to spread**

# Prediction of methylation spreading

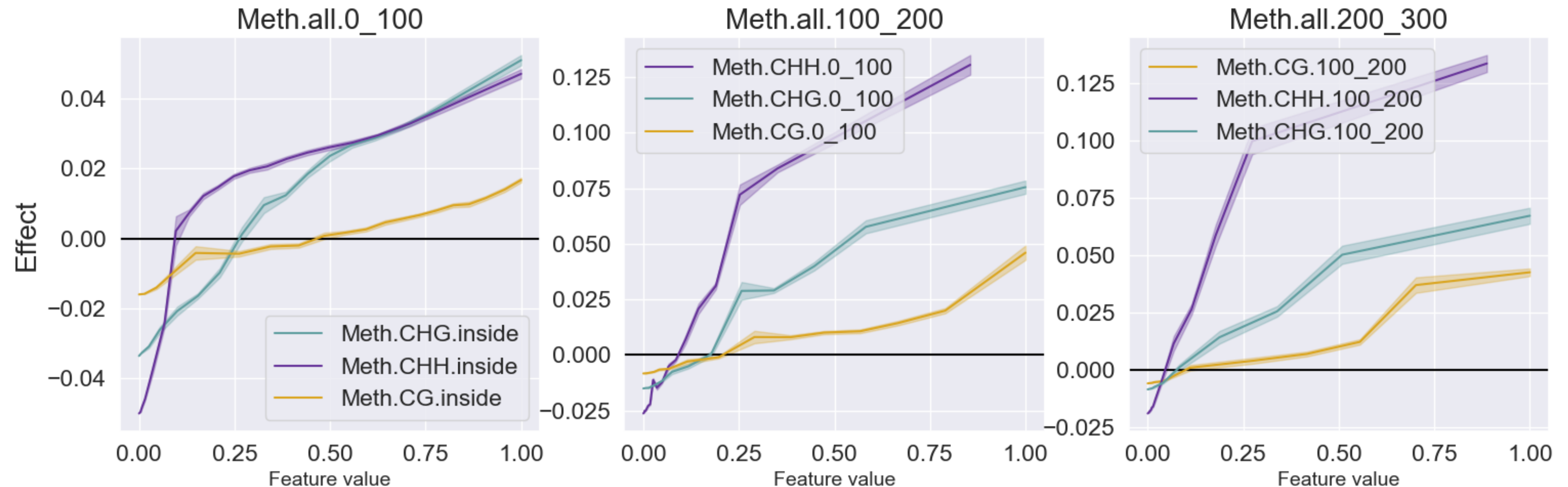**Accumulated Local Effects (ALE)**



**Conclusion:**

⬤ **Methylation of the TE inside ends** consistently comes as the most important feature with monotonous effect increase

⬤ TE is **methylated on the ends** $\implies$ more **likely to spread**

**Question:**

⬤ What defines the **methylation of the TE inside ends**?

# Prediction of inside methylation

**Features:**

- **TE** (length, distance to pericentromere, superfamily, insertion frequency, divergence, if inside a gene)
- **Nearest 2 genes** (length, distance, relative direction)
- **Methylation** in CG, CHG, CHH contexts (average genome-wide)
- **Densities** of CG, CHG, CHH contexts
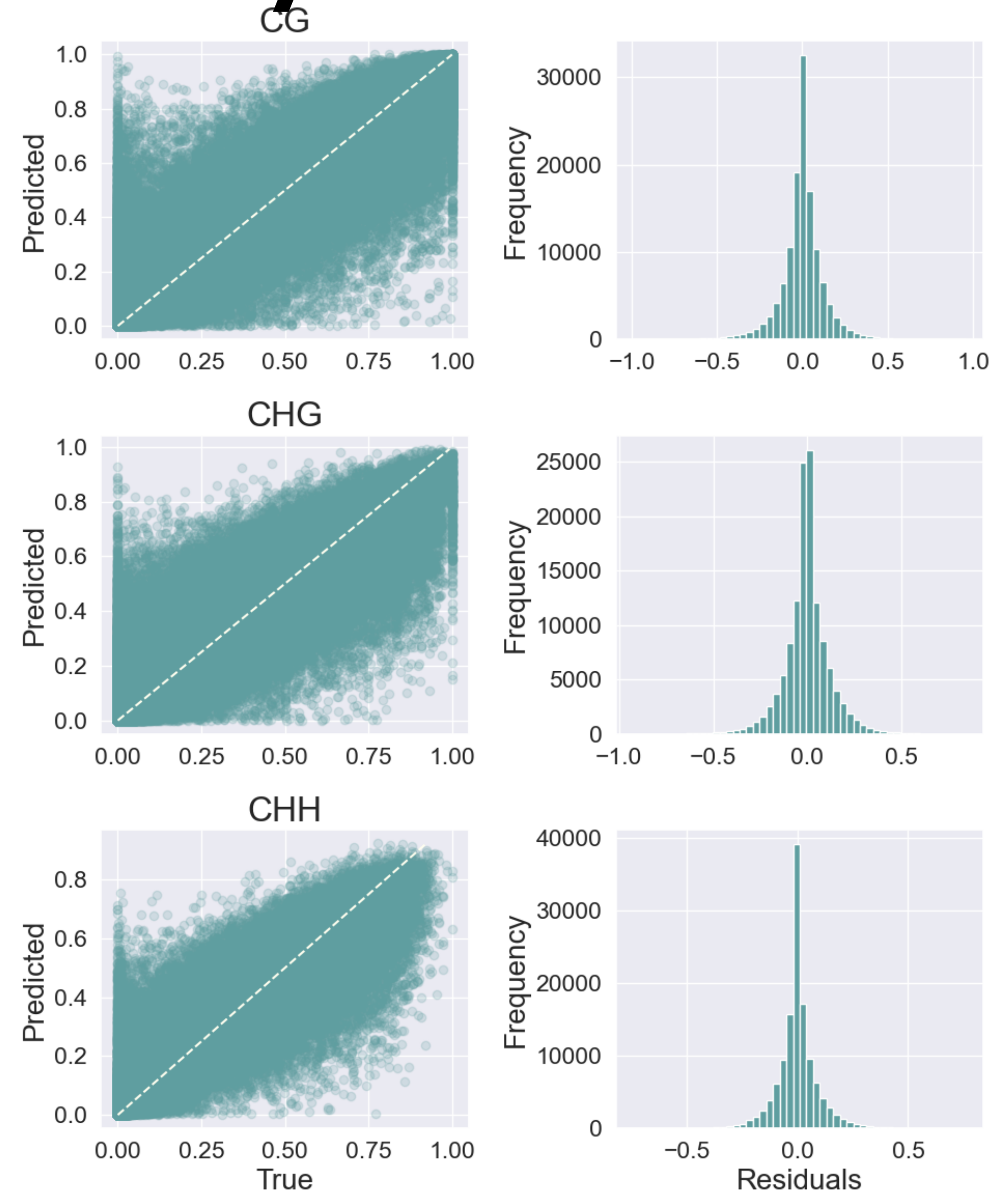
# Prediction of inside methylation

**Features:**

- **TE** (length, distance to pericentromere, superfamily, insertion frequency, divergence, if inside a gene)
- **Nearest 2 genes** (length, distance, relative direction)
- **Methylation** in CG, CHG, CHH contexts (average genome-wide)
- **Densities** of CG, CHG, CHH contexts
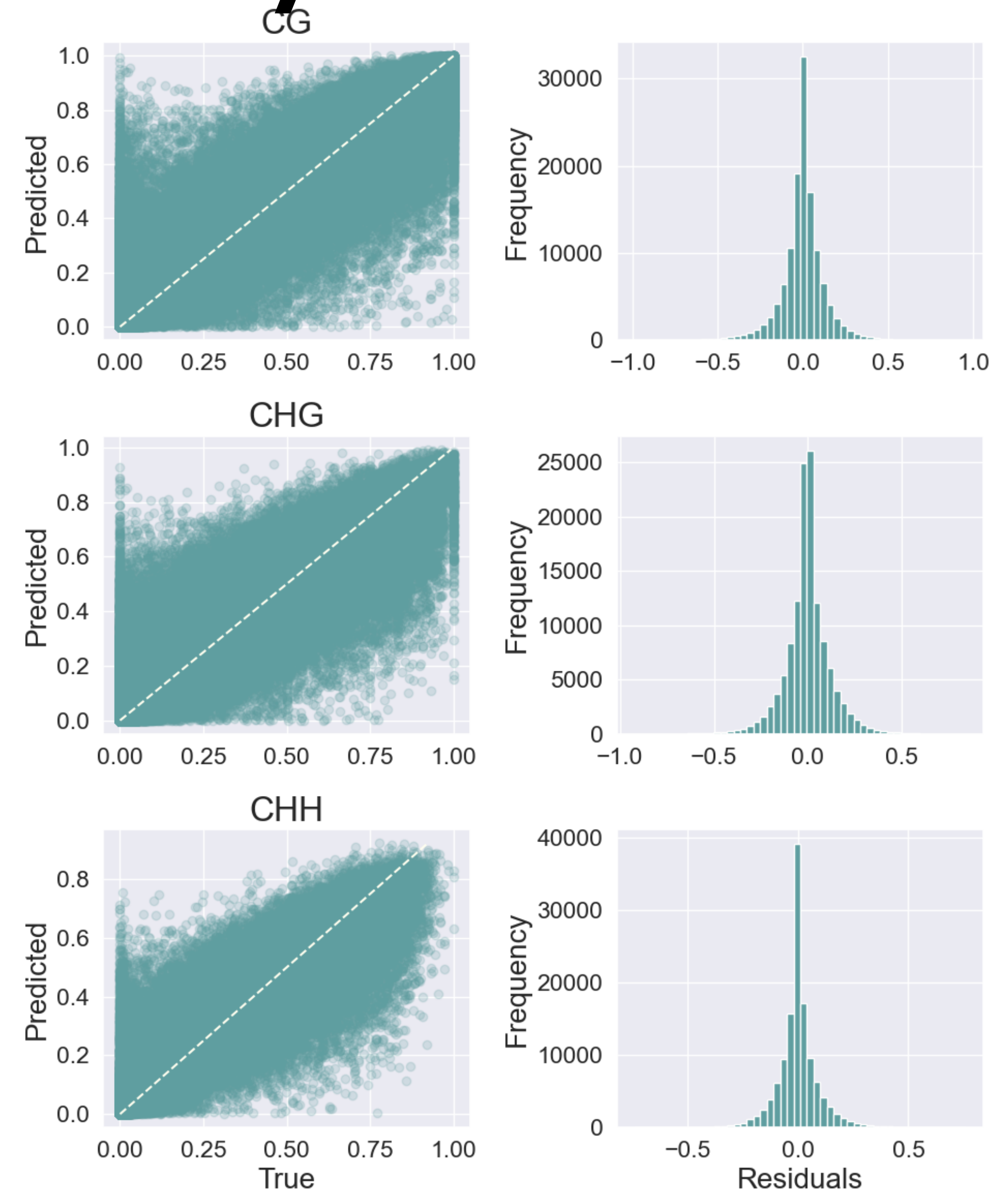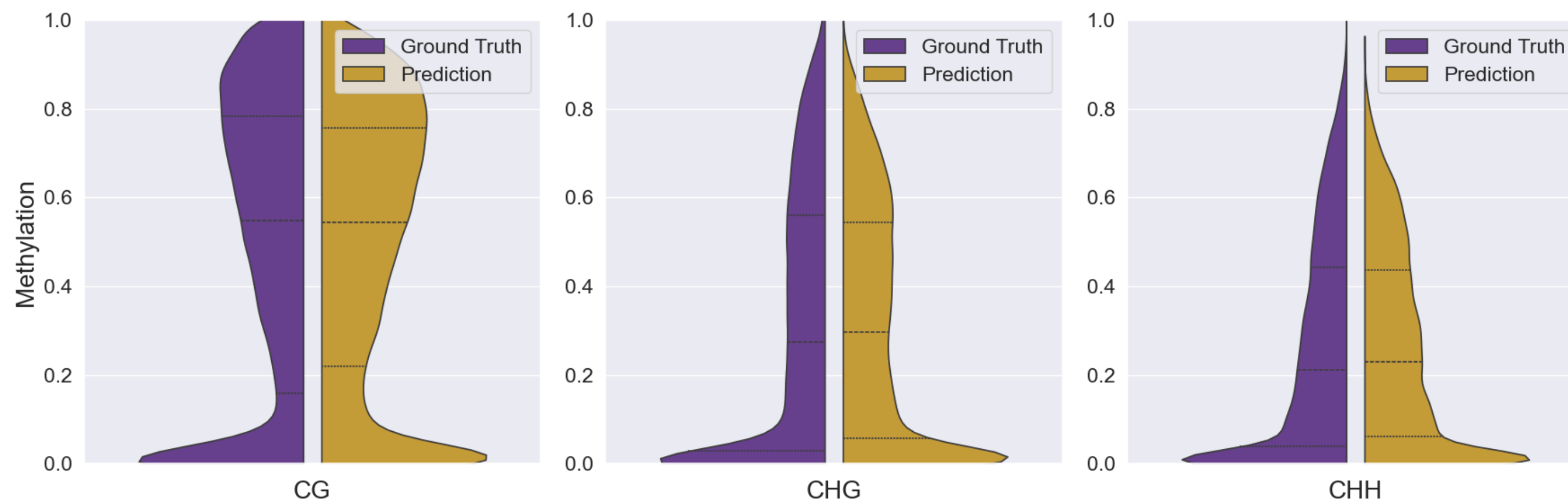
# Prediction of inside methylation
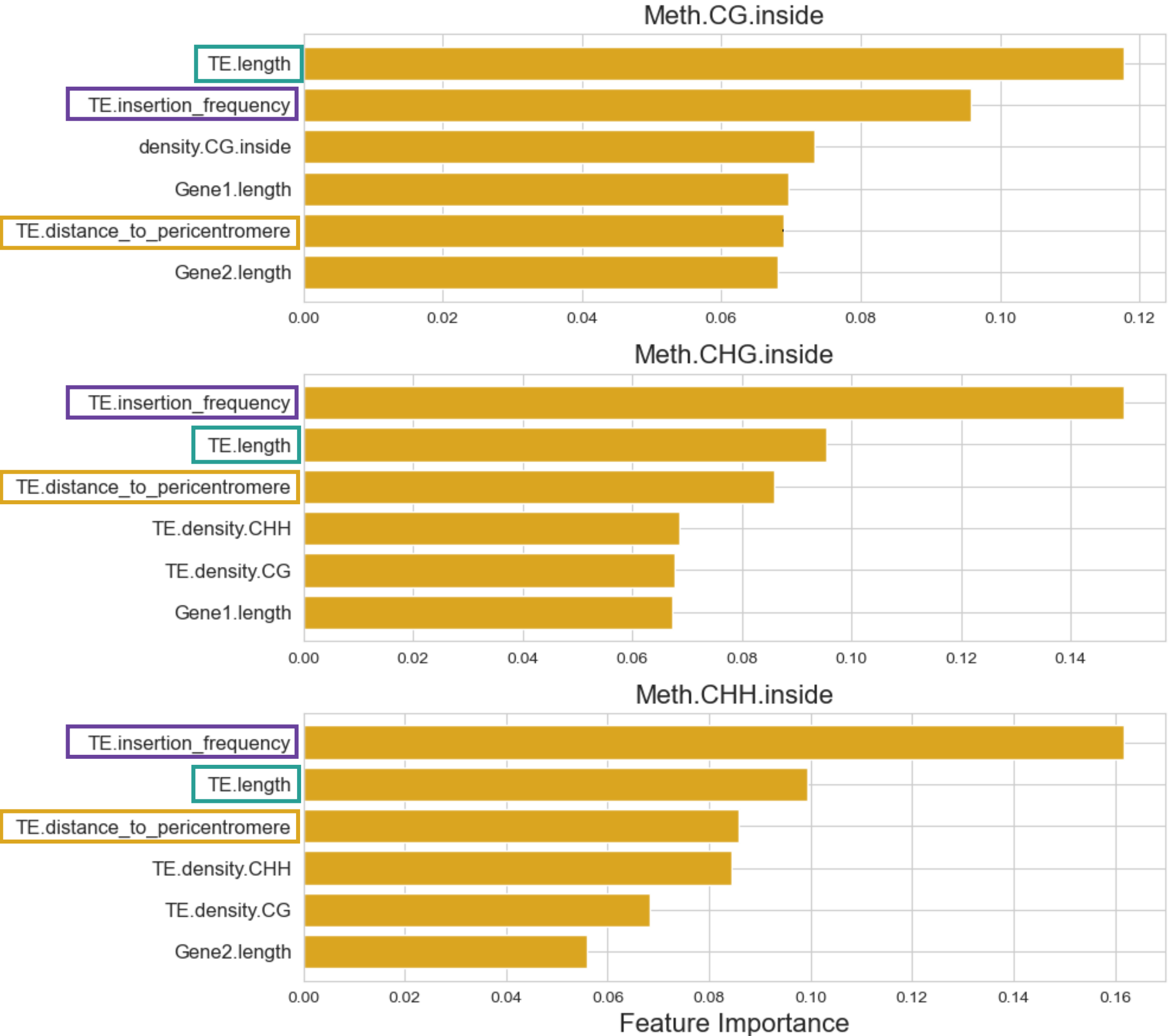
**Features:**

- **TE** (length, distance to pericentromere, superfamily, insertion frequency, divergence, if inside a gene)
- **Nearest 2 genes** (length, distance, relative direction)
- **Methylation** in CG, CHG, CHH contexts (average genome-wide)
- **Densities** of CG, CHG, CHH contexts
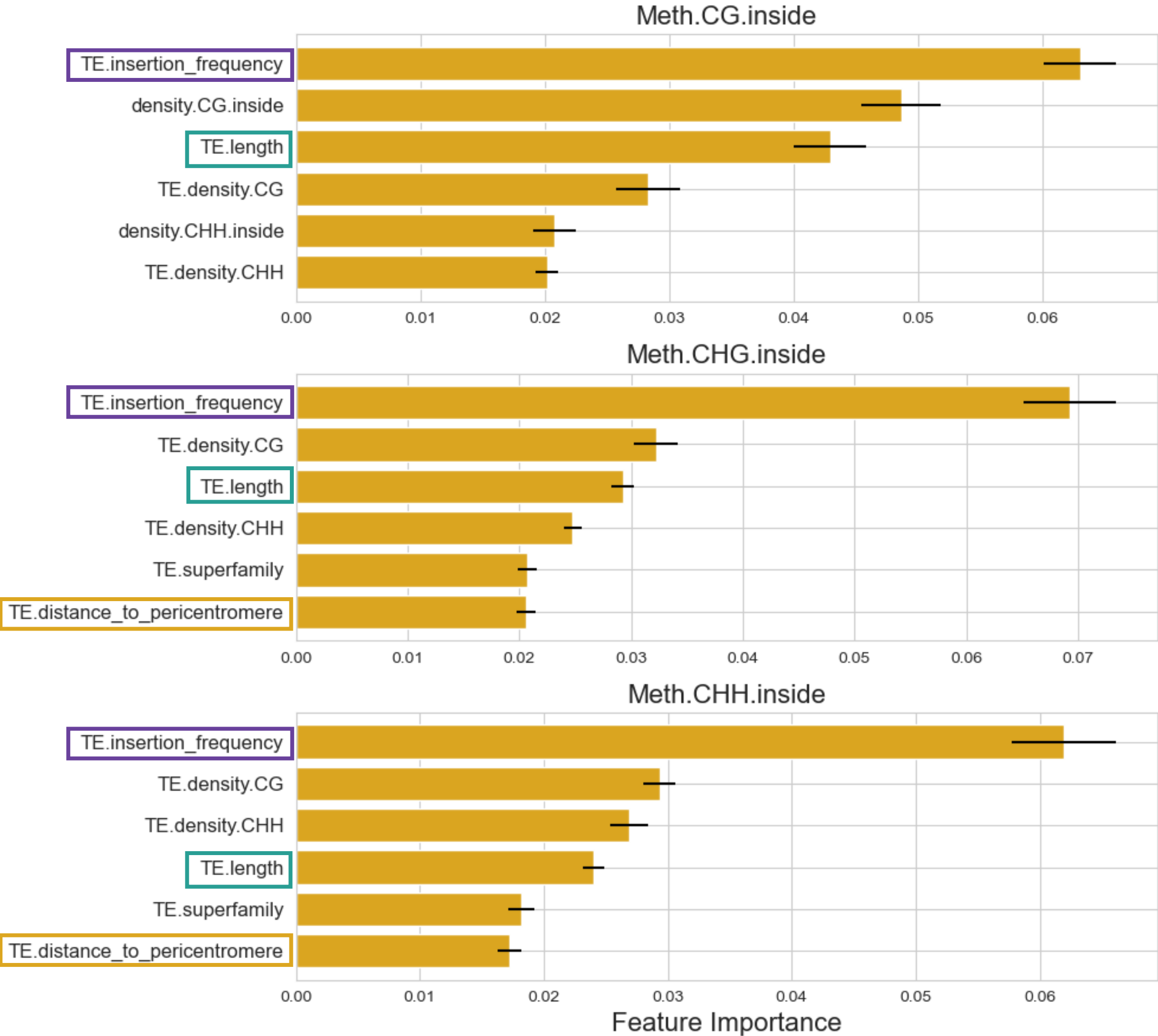
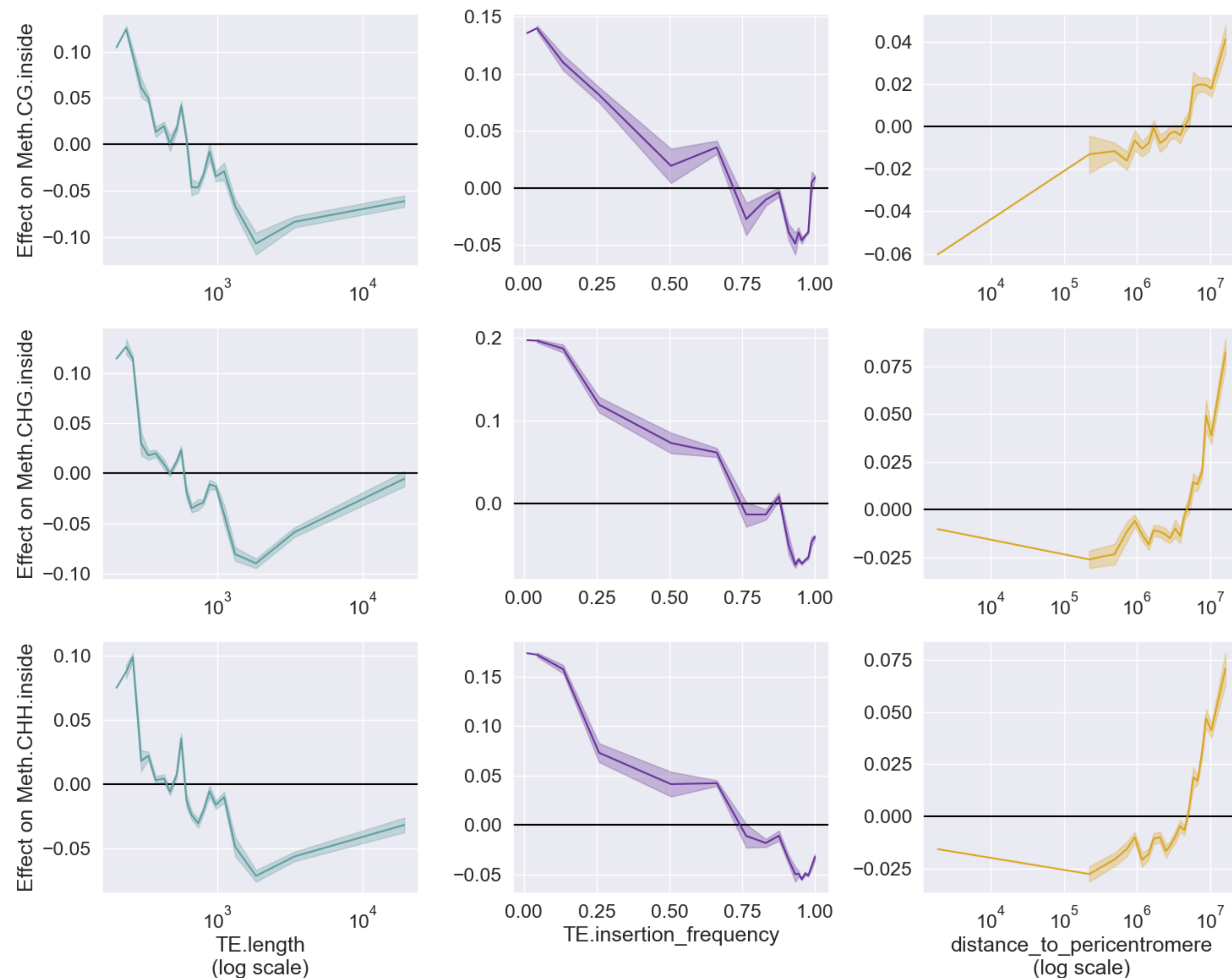# Prediction of inside methylation

**Accumulated Local Effects**

# Prediction of inside methylation

## Accumulated Local Effects

## Conclusions



- Length of TE is unlikely to be a driving factor

  (rather a confounder)

- **Insertion frequency**:

  Rare (= new) TEs are targeted by methylation machinery

- **Distance to pericentromere**:

  More distant TEs are more likely to be targeted
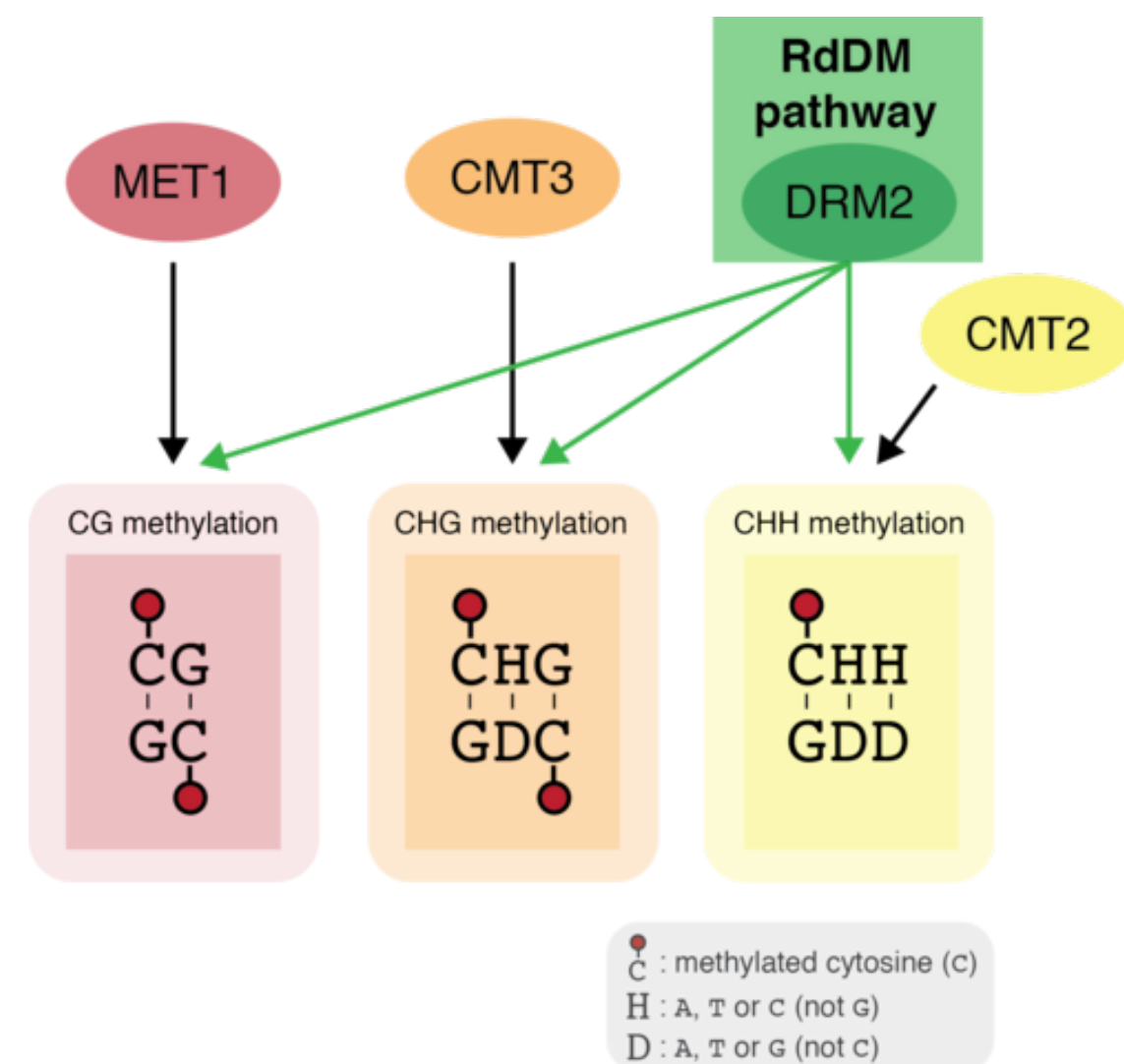
# Biological conclusions

- **The most important factors:**

  - methylation of insides in the CHG and CHH contexts

  - insertion frequency

  - distance to pericentromere

# Biological conclusions

**The most important factors:**

- methylation of insides in the CHG and CHH contexts

- insertion frequency

- distance to pericentromere



Wikipedia

**Hypothesis:** the **RdDM** machinery is responsible for spreading
- targets all contexts
- the <u>only</u> pathway capable of adding DNA methylation *de novo*
- targets rather chromosome arms

# Biological conclusions

○ **The most important factors:**

- methylation of insides in the CHG and CHH contexts

- insertion frequency
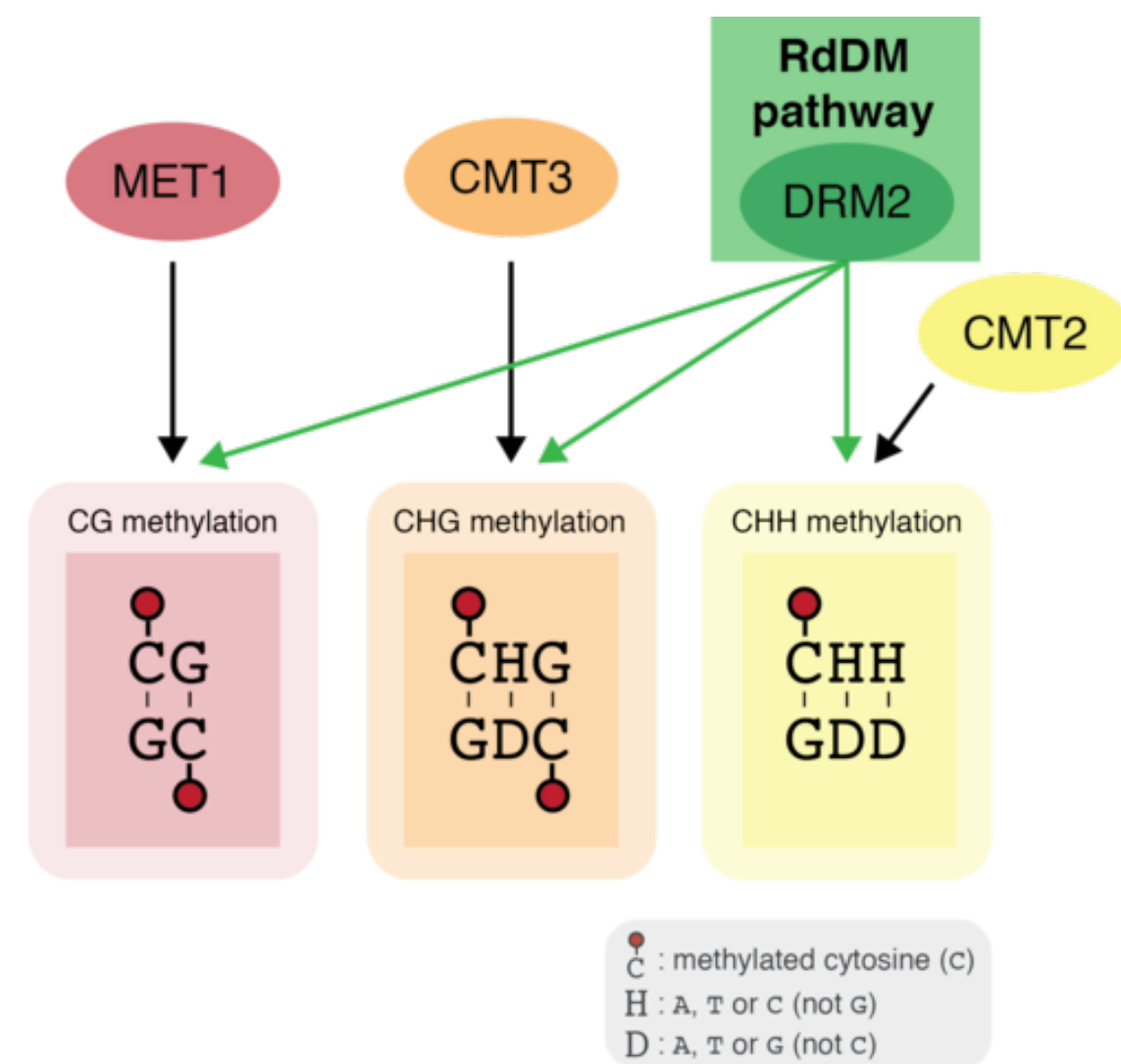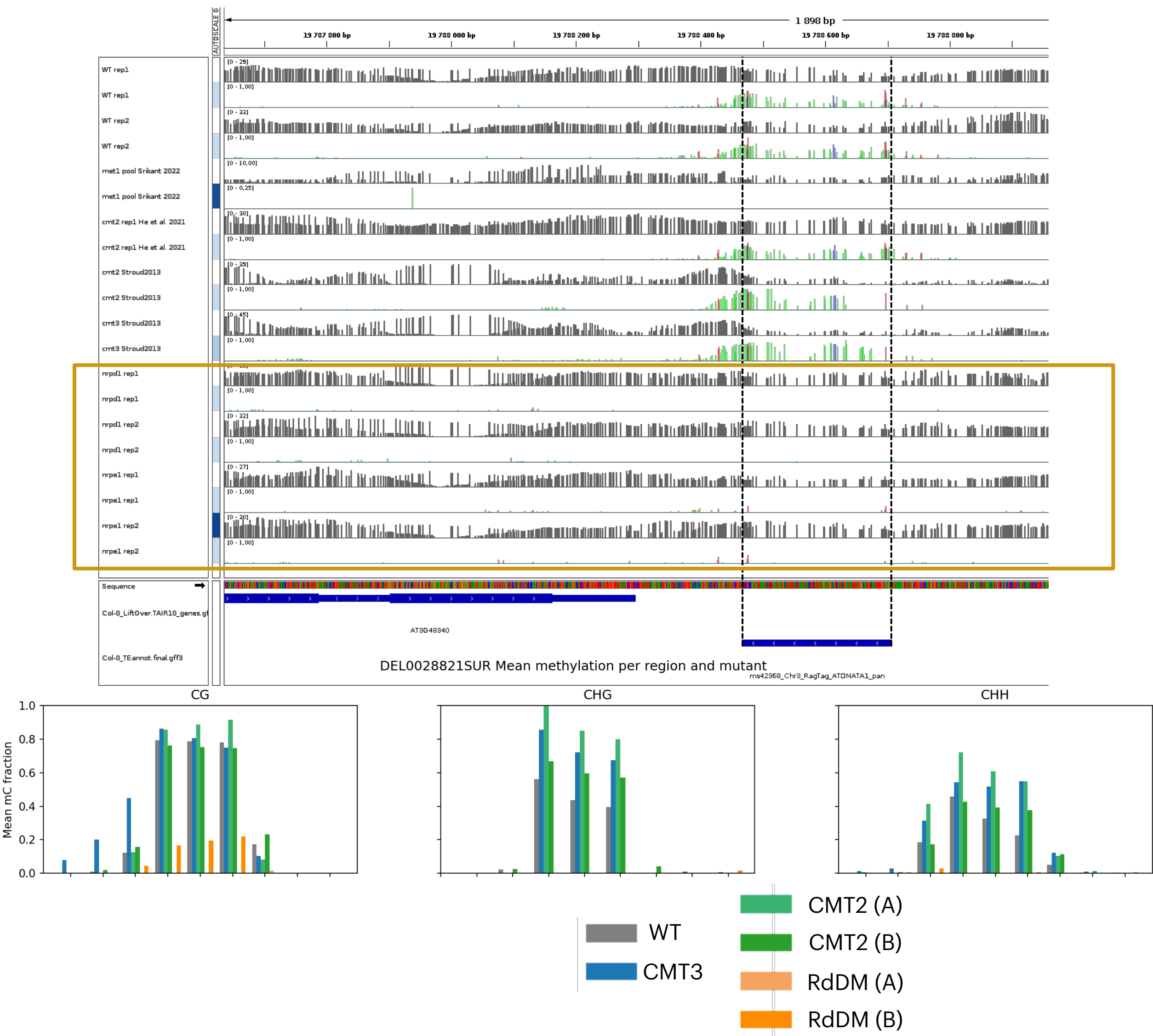
- distance to pericentromere



Wikipedia

○ **Hypothesis:** the **RdDM** machinery is responsible for spreading
- targets all contexts
- the <u>only</u> pathway capable of adding DNA methylation *de novo*
- targets rather chromosome arms

○ **Test:** mutants of Col-0 strain of *A. Thaliana* where different methylation pathways are knocked out
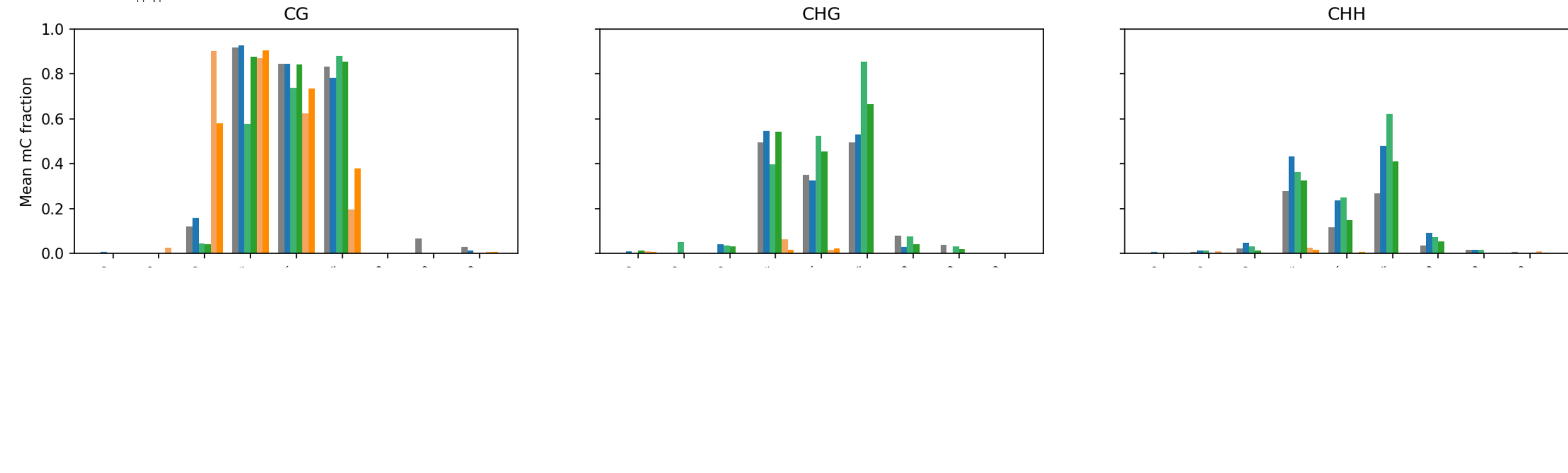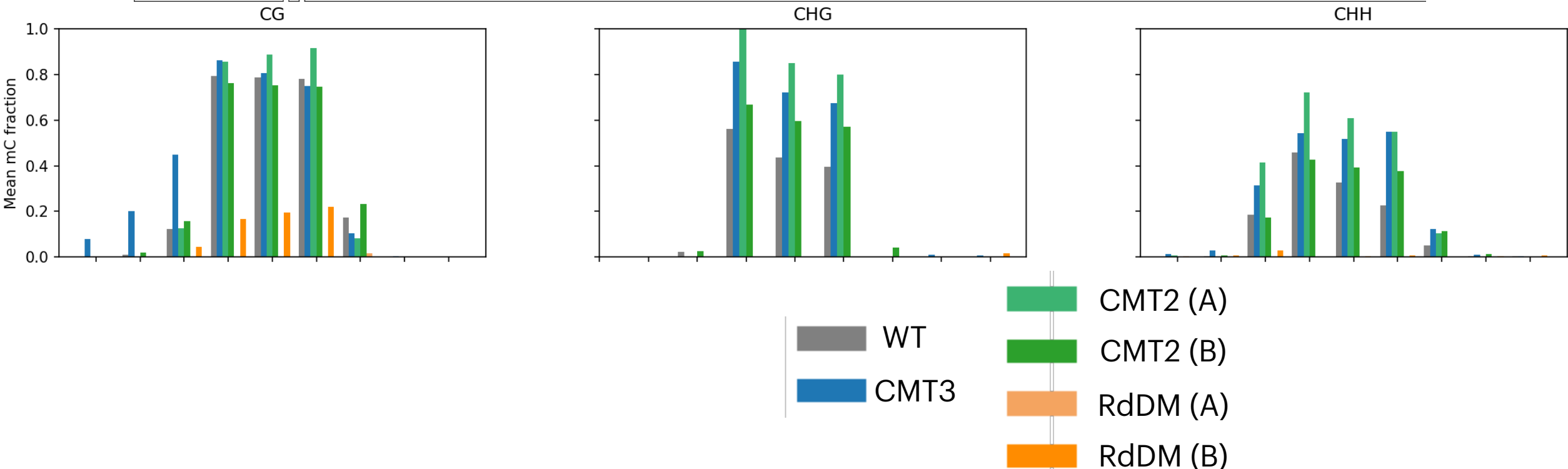
# Biological confirmation
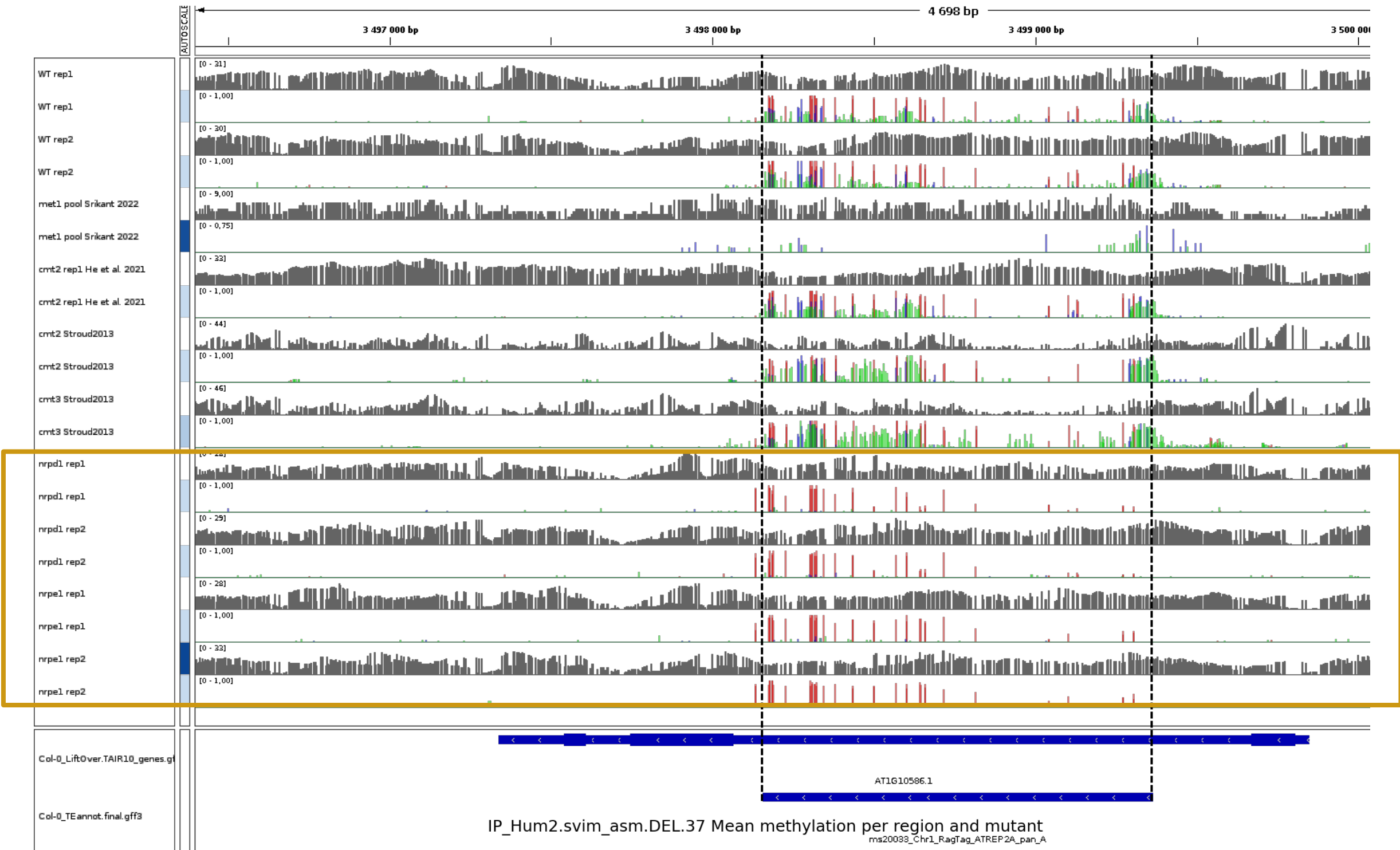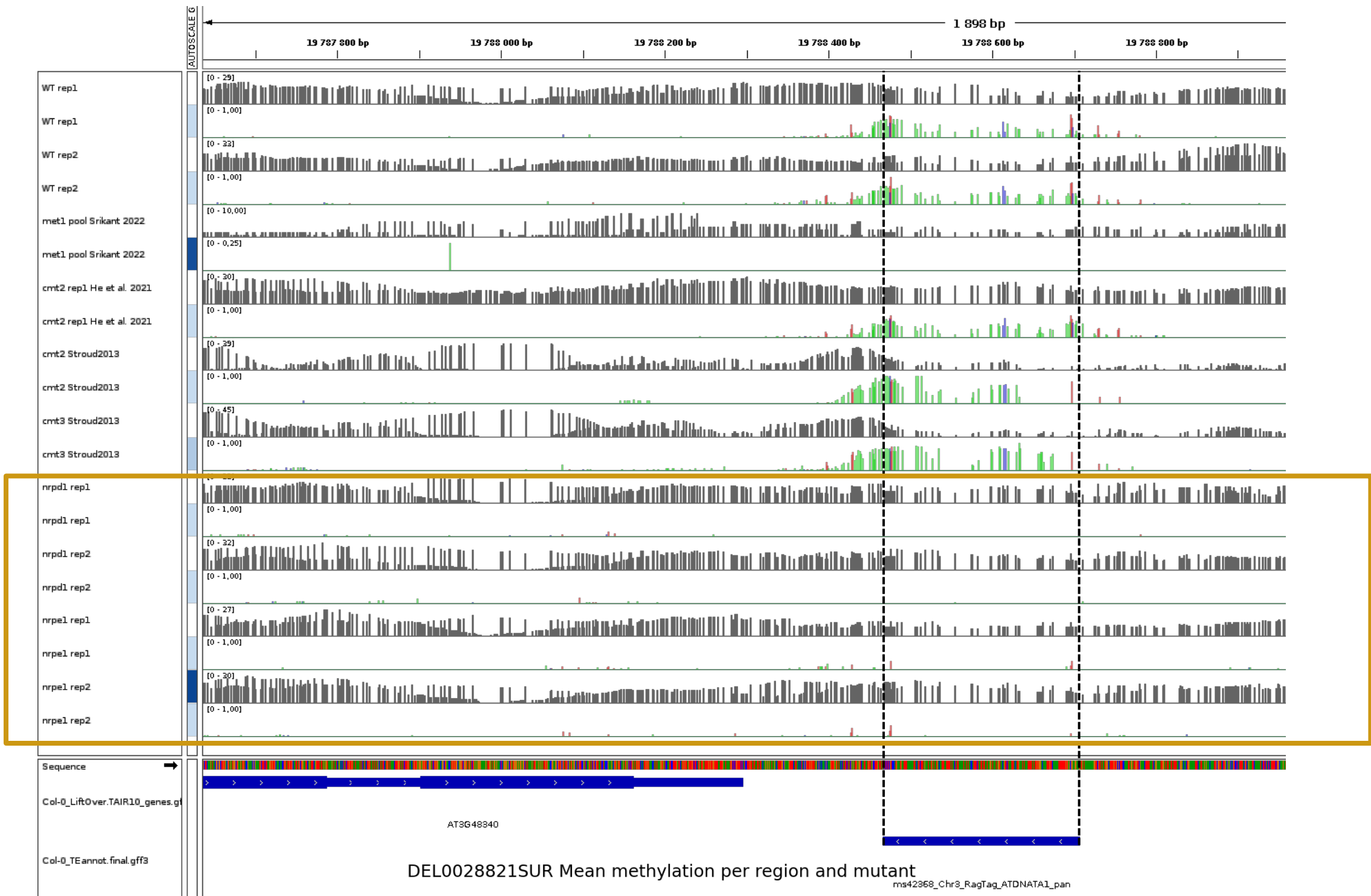
**DEL0028821SUR**



DEL0028821SUR Mean methylation per region and mutant

# Biological confirmation

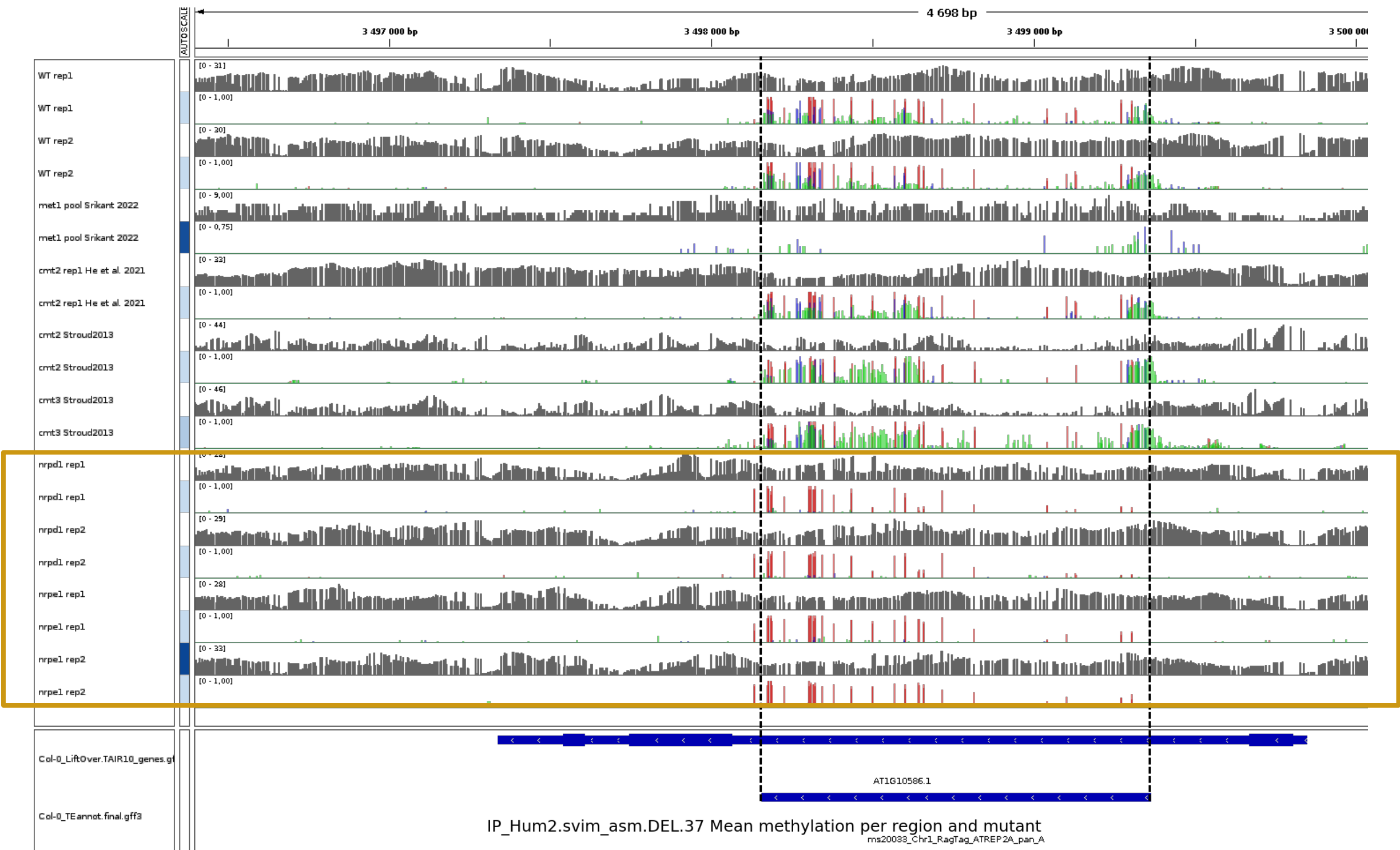# Biological confirmation



DEL0028821SUR

IP_Hum2.svim_asm.DEL.37

CHG and CHH methylation (and spreading!) disappear in RdDM mutants

# Conclusions

- An example of the workflow:

**biological phenomenon $\Longrightarrow$ machine learning model $\Longrightarrow$ explanations $\Longrightarrow$ real biological mechanisms**

# Conclusions

- An example of the workflow:

  **biological phenomenon** $\Longrightarrow$ **machine learning model** $\Longrightarrow$ **explanations** $\Longrightarrow$ **real biological mechanisms**

- Different **explainability tools** have been explored, and they provide **consistent conclusions**

# Conclusions

- An example of the workflow:

  **biological phenomenon $\Longrightarrow$ machine learning model $\Longrightarrow$ explanations $\Longrightarrow$ real biological mechanisms**

- Different **explainability tools** have been explored, and they provide **consistent conclusions**

- Potential actor (**RdDM**) is identified

# Conclusions

- An example of the workflow:

  **biological phenomenon $\implies$ machine learning model $\implies$ explanations $\implies$ real biological mechanisms**

- Different **explainability tools** have been explored, and they provide **consistent conclusions**

- Potential actor (**RdDM**) is identified

- We understand better one of the factors to **explain GWAS signals**

# Acknowledgements



## CBIO

Chloé-Agathe Azencott

Marie Dogo

Jérémy Cohen

Sylvain Cailloud

(and everyone else)

## IBENS

Vincent Colot

Pierre Baduel

Louna De Oliveira

Aurélien Petit

(and everyone else)