

# Introduction à la bioinformatique (UE SSV3U15)

## TP2. Du gène à la protéine

### Diaporama d'accompagnement du TP

Jacques van Helden (Aix-Marseille Université)  
ORCID [0000-0002-8799-8584](https://orcid.org/0000-0002-8799-8584)

### Objectifs

- Utiliser des ressources bioinformatiques pour explorer les génomes d'organismes modèles, afin de comprendre la structuration et la composition de ces génomes.

### Notions mises en pratique

- Structuration des gènes (transcrits, introns, exons, régions codantes, régions non traduites)
- Organisation des génomes (éléments structurels des chromosomes, régions géniques et intergéniques, opérons bactériens).
- Quelques éléments de génomique comparative (conservation / divergence, réarrangements chromosomiques, synténie)
- Les principaux types d'homologie : orthologie, paralogie
- Annotation fonctionnelle des gènes.
- Transcriptomique : expression différentielle des gènes dans différents tissus
- Protéomique : analyse de l'ensemble des protéines codées par un génome

N'oubliez pas que vous pouvez à tout moment consulter le [glossaire du cours](#) pour obtenir une définition sommaire des principaux termes utilisés.

## Etapes

- Annotations génomiques dans la région du gène humain PAX6
- Génomique comparative : homologues de PAX6 chez les métazoaires
- Organisation des gènes en intron chez les bactéries

## Complétion

- Tous les exercices doivent être réalisés par chaque étudiant.
- Les QCM de TP ne sont pas notés.

# Éléments de contexte

# Annotations génomiques

Le tableau indique le nombre d'annotations pour différents types d'éléments du génome humain (gènes codants, non-codants de différents types, transcrits) dans différentes bases de données de référence.

Les nombres précis d'annotations varient d'une base de données à une autre, mais les ordres de grandeur sont indicatifs.

## Constats

- Au début du projet de séquençage du génome humain, on s'attendait à trouver ~100 000 gènes codants. Une vingtaine d'années plus tard, on en répertorie ~20 000
- Le séquençage de l'ARN révèle un nombre à peu près équivalent de "gènes" non-codants (plus précisément, régions transcrites dont on ignore généralement la fonction).
- Au total, on dénombre ~200 000 transcrits, qui incluent les transcrits alternatifs (variants d'épissage) pour ces gènes codants et non-codants.

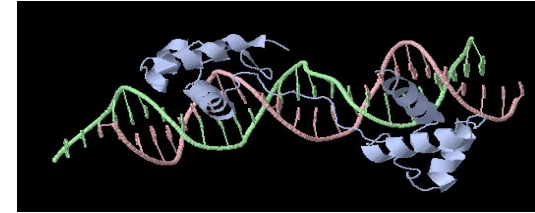
Feature type	Gencode	Ensembl	RefSeq	CHES
<b>Protein-coding genes</b>	19 901	20 376	20 345	21 306
<b>lncRNA genes</b>	15 779	14 720	17 712	18 484
<b>Antisense RNA</b>	5 501		28	2 694
<b>Miscellaneous RNA</b>	2 213	2 222	13 899	4 347
<b>Pseudogenes</b>	14 723	1 740	15 952	
<b>Total transcripts</b>	203 835	203 903	154 484	323 827

## Exemples traités

# Facteur transcriptionnel – PAX6

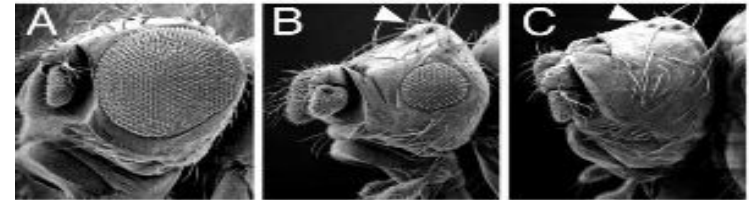
- La protéine PAX6 (violet sur l'image du haut) est un facteur transcriptionnel, qui se lie à des sites spécifiques sur l'ADN génomique (vert et rose).
- PAX6 contrôle l'expression de gènes impliqués dans la formation de l'oeil. Les gènes-cibles de PAX6 sont encore pour la plupart inconnus à ce jour.
- **Perte de fonction:** chez la drosophile, l'inactivation de *eyeless* (= PAX6) provoque une malformation ou une absence d'oeil.
- **Gain de fonction:** des mouches chez lesquelles on force le gène *eyeless* à s'exprimer dans les antennes développent des structures optiques à la place des antennes.

Liaison PAX6 - ADN

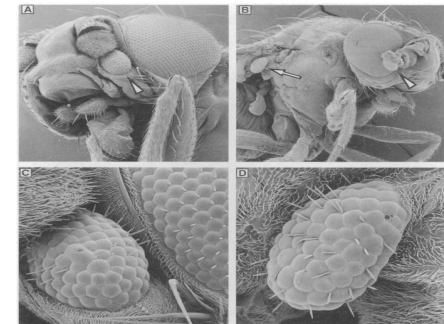


<http://www.rcsb.org/pdb/explore.do?structureId=6PAX>

Phénotype de perte de fonction



Phénotype de gain de fonction



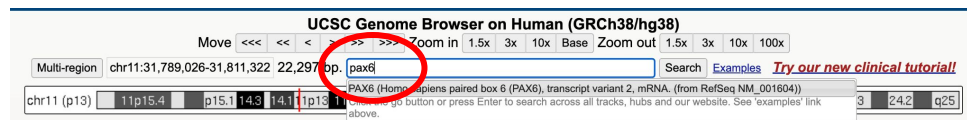
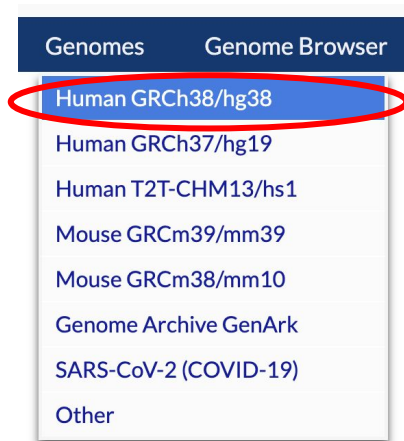
# Tutoriel illustré

## Annotations génomiques dans la région du gène humain PAX6



# PAX6 sur UCSC genome browser

- Connectez-vous au [UCSC Genome Browser](#)
- Dans le menu **Genomes**, sélectionnez la version **hg38** du **génomme humain**.
- Entrez le nom du gène d'intérêt (**PAX6**) dans la boîte de recherche et cliquez sur **Search**.



- Connectez-vous au [UCSC Genome Browser](https://genome-browser.ensembl.org/)
- Sélectionnez la version **hg38** du **génom**e **humain**
- Entrez le nom du gène d'intérêt (**PAX6**) dans la boîte de recherche et cliquez sur **Search**.

La page de résultat affiche une série d'annotations de PAX6 dans différentes bases de données de référence pour le génome humain. Comment choisir ? En première instance, le mieux est de se fier aux annotations du consortium international HUGO, responsable de la nomenclature des gènes humains.

- Sous le titre "HUGO Gene Nomenclature", cliquez sur le lien **PAX6** - chr11:31789026-31817960

Search Results on hg38 (Human Dec. 2013 (GRCh38/hg38))

☐ MANE Select Plus Clinical: Representative transcript from RefSeq & GENCODE:

☒ HUGO Gene Nomenclature:

- **PAX6** - chr11:31789026-31817960
- **PAX6-AS1** - chr11:31816266-32002405

☐ Gencode Genes:

- **PAX6** (ENST00000640368.2) - chr11:31789026-31811322 - **PAX6** ENST00000640368.2 Homo sapiens paired box 6 PAX6 transcript variant
- **PAUPAR** (ENST00000644607.1) - chr11:31816266-32002405 - PAUPAR ENST00000644607.1 **PAX6** upstream antisense RNA from HGNC PAUPAR BX648962 uc285izg.1 uc285izg.1
- **BCL2L15** (ENST00000393316.8) - chr11:113876816-113887581 - ... Q5TBC7 P50222 MEOX2 NbExp 3 IntAct EBI-10247136 EBI-748397 Q5TBC7 P26367 **PAX6** NbExp 3 IntAct EBI-10247136 EBI-747278 Q5TBC7 P62487 POLR2G NbExp
- **LYSMD1** (ENST00000368908.10) - chr1:151159748-151165902 - ... Q96S90 Q5JR59 MTUS2 NbExp 3 IntAct EBI-10293291 EBI-742948 Q96S90 P26367 **PAX6** NbExp 3 IntAct EBI-10293291 EBI-747278 Q96S90 Q6NUQ1 RINT1 NbExp
- **CCDC103** (ENST00000417826.3) - chr17:74899729-44905390 - ... Q8IW40 Q6FHY5 MEOX2 NbExp 3 IntAct EBI-10261970 EBI-16439278 Q8IW40 P26367 **PAX6** NbExp 3 IntAct EBI-10261970 EBI-747278 Q8IW40 Q9NRD5 PICK1 NbExp
- **SLC12A8** (ENST00000469902.6) - chr3:125082644-125212748 - ... a role in the control of keratinocyte proliferation A0A02V P26367 **PAX6** NbExp 3 IntAct EBI-11737524 EBI-747278 Membrane Multi-pass membrane protein
- **TCP11L1** (ENST00000334274.9) - chr11:33039572-33073550 - ... Q9NUJ3 P50221 MEOX1 NbExp 3 IntAct EBI-2555179 EBI-2864512 Q9NUJ3 P26367 **PAX6** NbExp 3 IntAct EBI-2555179 EBI-747278 Q9NUJ3 Q5SXH7-1 PLEKHS1 NbExp
- **ZNFX13** (ENST00000327033.11) - chr2:27377235-27380734 - ... Binds DNA Can associate with the proximal promoter regions of **PAX6** and SP4 and their known targets including ARR3 RHO
- **SPDYC** (ENST00000377185.3) - chr11:65170233-65173374 - ... Q5MJ68 Q9Y250 LZTS1 NbExp 3 IntAct EBI-12162209 EBI-1216080 Q5MJ68 P26367 **PAX6** NbExp 3 IntAct EBI-12162209 EBI-747278 Cytoplasm Note Colocalizes with
- **CXorf38** (ENST00000327877.10) - chrX:40626921-40647561 - ... Q8TB03 Q02548 PAX5 NbExp 3 IntAct EBI-12024320 EBI-296331 Q8TB03 P26367 **PAX6** NbExp 3 IntAct EBI-12024320 EBI-747278 Q8TB03 Q9Y3C5 RNFI1 NbExp

[Show 98 more matches for Gencode Genes](#)

☐ NCBI RefSeq genes, curated subset (NM\_\*, NR\_\*, NP\_\* or YP\_\*):

- **NR\_033971.1** - chr11:31816566-31887041
- **NM\_001258463.2** - chr11:31789026-31812203
- **NM\_001258464.2** - chr11:31789026-31811322
- **NM\_001368926.2** - chr11:31789026-31811322
- **NM\_001368914.2** - chr11:31789026-31811322
- **NM\_001604.6** - chr11:31789026-31811322
- **NM\_001368893.2** - chr11:31789026-31811322
- **NM\_001368917.2** - chr11:31789026-31811322
- **NM\_000280.6** - chr11:31789026-31811322
- **NM\_001368927.2** - chr11:31789026-31811322

[Show 44 more matches for NCBI RefSeq genes, curated subset \(NM\\_\\*, NR\\_\\*, NP\\_\\* or YP\\_\\*\)](#)

☐ RefSeq Genes:

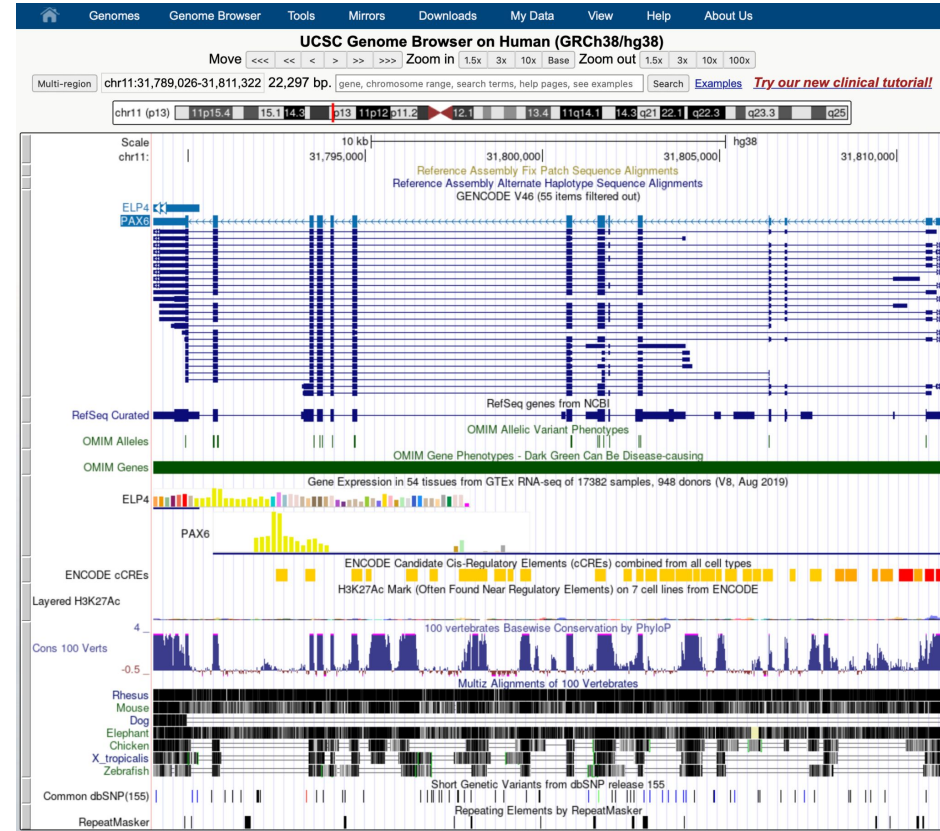
- **PAX6** - chr11:31789026-31811121 - (NM\_000280) paired box protein Pax-6 isoform a
- **PAX6** - chr11:31789026-31811322 - (NM\_001258464) paired box protein Pax-6 isoform a
- **PAX6** - chr11:31789026-31812203 - (NM\_001310158) paired box protein Pax-6 isoform b
- **PAX6** - chr11:31789026-31817961 - (NM\_001127612) paired box protein Pax-6 isoform a
- **PAX6** - chr11:31789026-31811322 - (NM\_001604) paired box protein Pax-6 isoform b
- **PAX6** - chr11:31789026-31817961 - (NM\_001258462) paired box protein Pax-6 isoform b
- **PAX6** - chr11:31789026-31804059 - (NM\_001310161) paired box protein Pax-6 isoform d
- **PAX6** - chr11:31789026-31811121 - (NM\_001258465) paired box protein Pax-6 isoform a
- **PAX6** - chr11:31789026-31812203 - (NM\_001258463) paired box protein Pax-6 isoform b
- **PAX6** - chr11:31793205-31806925 - (NM\_001310159) paired box protein Pax-6 isoform c

[Show 44 more matches for RefSeq Genes](#)

# Choix de pistes d'annotations du UCSC Genome Browser

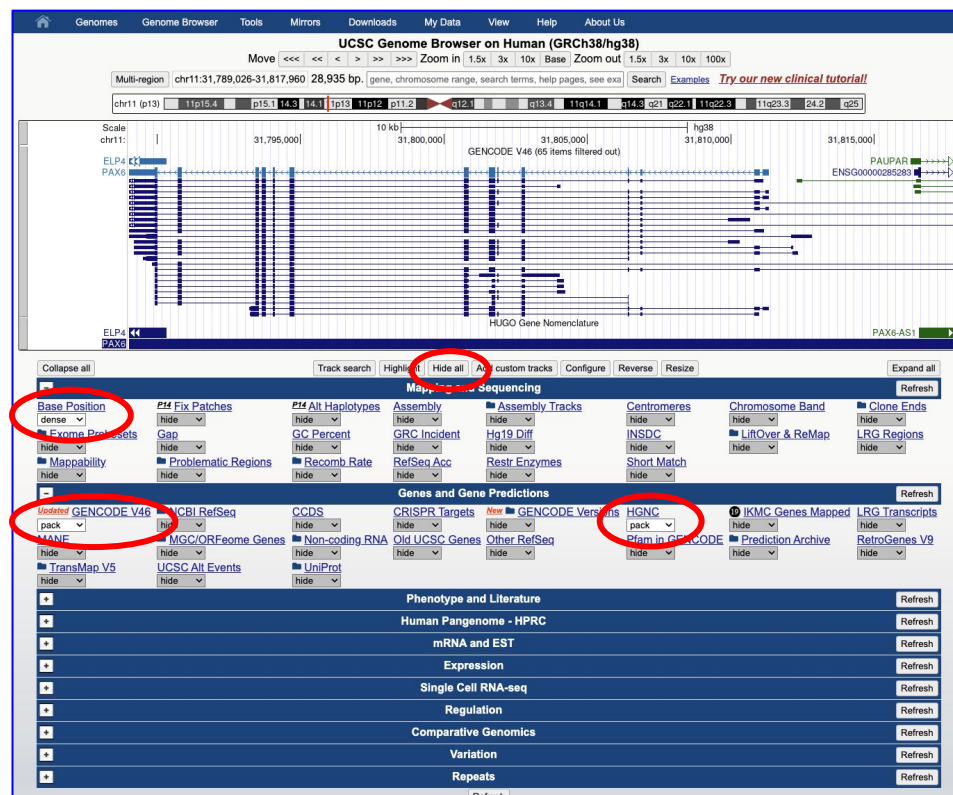
Le navigateur de génomes [UCSC Genome Browser](#) affiche un vaste choix de pistes d'annotation. La carte génomique en affiche un sous-ensemble, qui s'adaptent en fonction de vos consultations précédentes.

Nous allons restreindre la visualisation aux pistes d'annotations utilisées pour ce TP.



# PAX6 sur UCSC genome browser

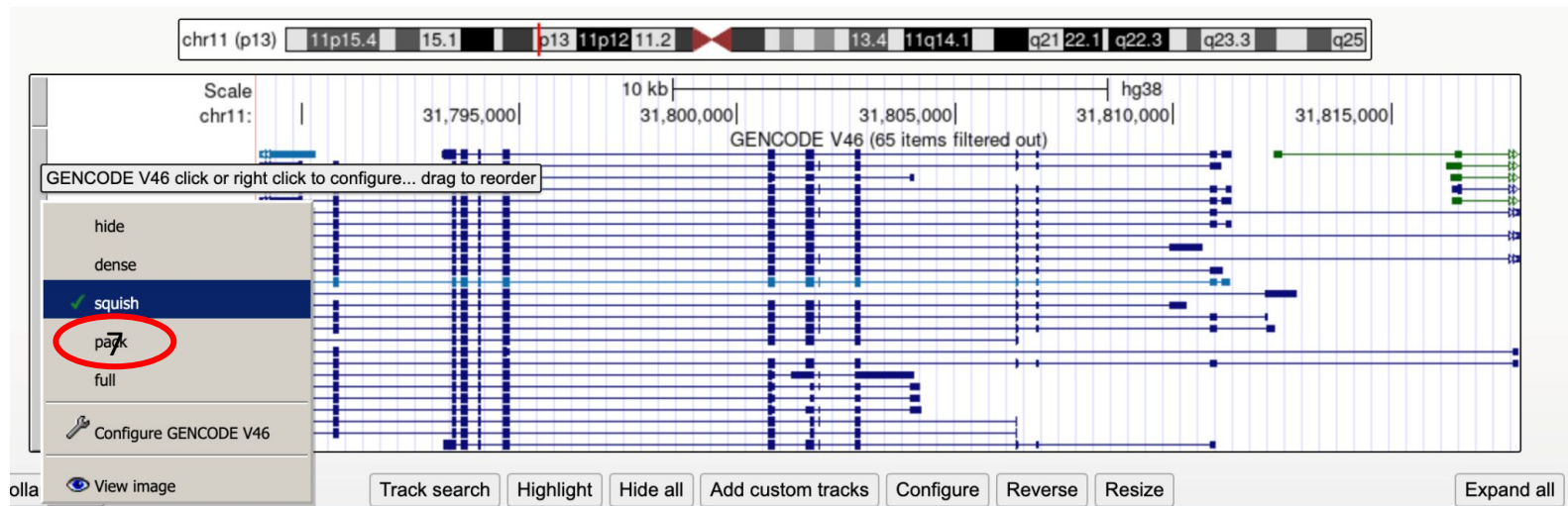
- Connectez-vous au [UCSC Genome Browser](https://genome.ucsc.edu/)
- Dans le menu **Genomes**, sélectionnez la version **hg38** du **génomme humain**.
- Entrez le nom du gène d'intérêt (**PAX6**) dans la boîte de recherche et cliquez sur **Search**.
- Sous le titre "HUGO Gene Nomenclature", cliquez sur le lien **PAX6** - chr11:31789026-31817960.
- Descendez sous la carte génomique pour afficher les choix de pistes d'annotations.
- Entre la carte et les options, cliquez **Hide all** pour masquer les pistes génomiques par défaut.
- Dans la catégorie "**Mapping and Sequencing**", sélectionnez le mode d'affichage "**dense**" pour la piste d'annotation "**Base position**".
- Dans la catégorie "**Genes and Gene Prediction**", sélectionnez le mode "**pack**" pour les pistes "**GENCODE\_V46**" et **HGNC**.
  - HGNC indique les limites des gènes, tandis que GENCODE\_V46 fournit des informations plus détaillées sur la structure des gènes (introns, exons, transcrits alternatifs, ...).
- Cliquez "**Refresh**" à droite d'une des catégories.
- Cliquez "**Resize**" sous la carte pour ajuster la largeur à celle de votre écran.



## Reconfigurer le mode d'affichage

Vous pouvez à tout moment reconfigurer le mode d'affichage d'une piste d'annotation, en cliquant droit (**contrôle-clic**) sur la figure. Ceci vous affichera un menu avec des modes d'affichages de plus en plus détaillés : hide, dense, squish, pack, full.

- Testez les différents niveaux de détail avec la piste **GENCODE\_V46**, puis sélectionnez le mode **pack**, qui vous permet généralement de visualiser les transcrits alternatifs en occupant une place raisonnable.

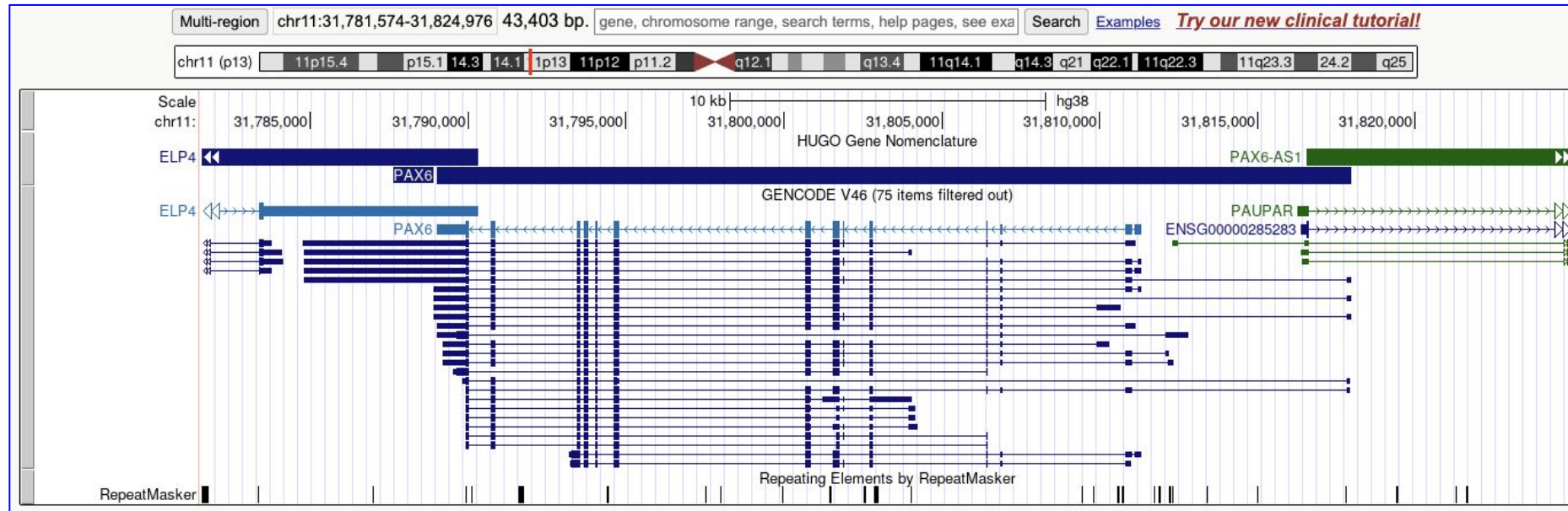




## Régions répétitives

- Dans la catégorie "**Repeats**", activez l'affichage de "**Repeatmasker**" en format "**dense**" et cliquez "**Refresh**".
- Dézoomez (**Zoom out**) d'un facteur **x1.5** pour voir les environs du gène
- **Déplacez** la piste HGCN au-dessus de la piste GENCODE (en positionnant la souris vers le coin supérieur gauche d'une piste, une flèche apparaît qui permet de la déplacer)

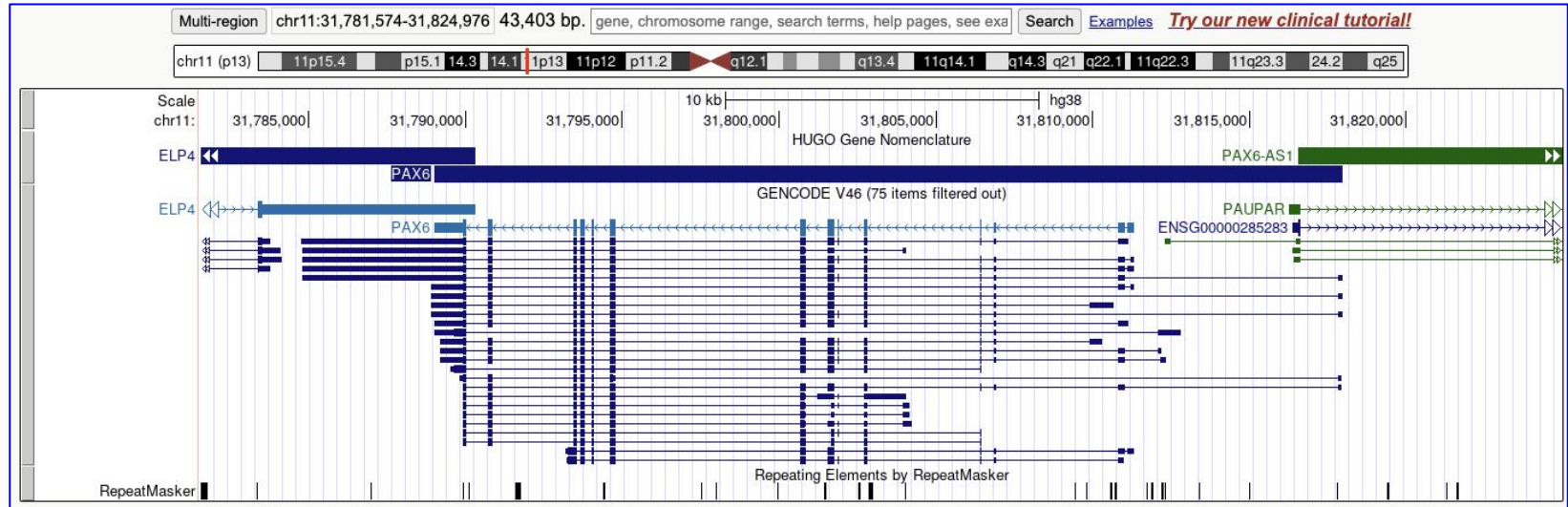
Observez la disposition du gène PAX6. Notez qu'il chevauche ses voisins de gauche (ELP4) et de droite (PAX6-AS1, où AS indique qu'il s'agit d'un gène antisens).



# Exercice 1. Annotations génomiques dans la région du gène humain PAX6

Sur Ametice, ouvrez le questionnaire du TP3 et répondez aux questions de l'**Exercice 1**.

- Sur quel brin chromosomique se situe le gène PAX ? +, -
- Sur quel chromosome est situé le gène PAX6 ? 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, X, Y
- Sur quel bras chromosomique (gauche / droite) ? Gauche, droite
- Sur quelle région chromosomique ? p15.4, 15.1, p13, p11.2, q2.1, q23.3
- Quelle est sa longueur en kilobases ? 12, 22, 29, 31 789, 31 817
- Combien de régions répétitives distinguez-vous sur la région du gène PAX6 ? aucune, 15, 21, 22, 23, 25, 30
- Quelle est la densité approximative (nombre de régions répétitives / kilobase) ? 0.1, 0.8, 1, 1.2, 1.3, 2, 5, 10
- Sur la longueur du gène PAX6, les régions répétitives coïncident généralement avec les introns, régions exoniques codantes, régions exoniques non-codantes

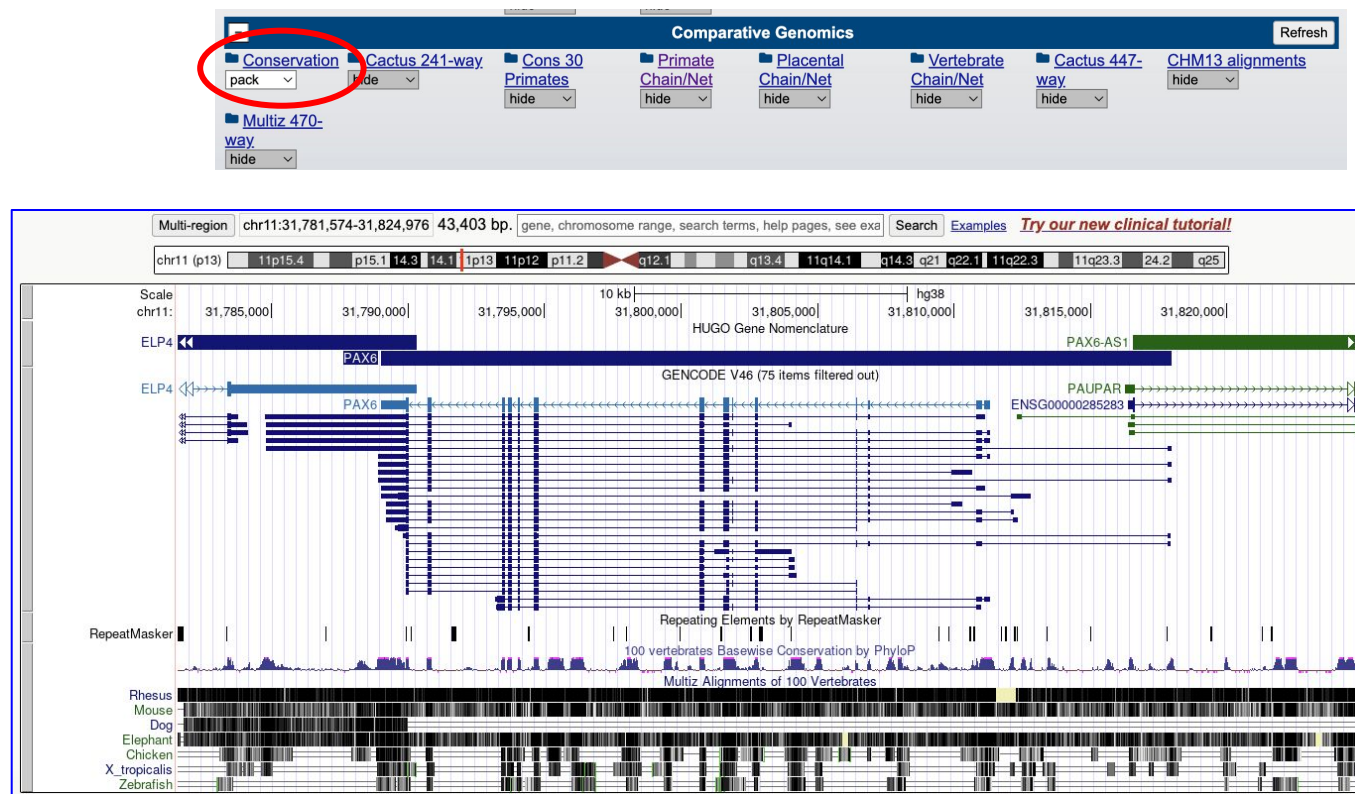


# Conservation du gène PAX6 dans les génomes de vertébrés



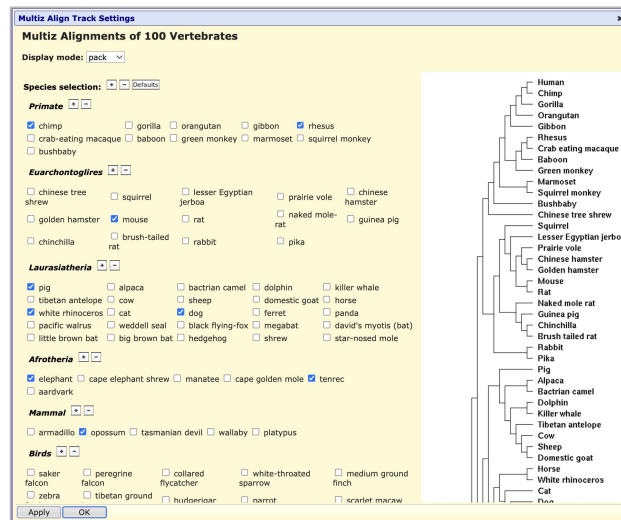
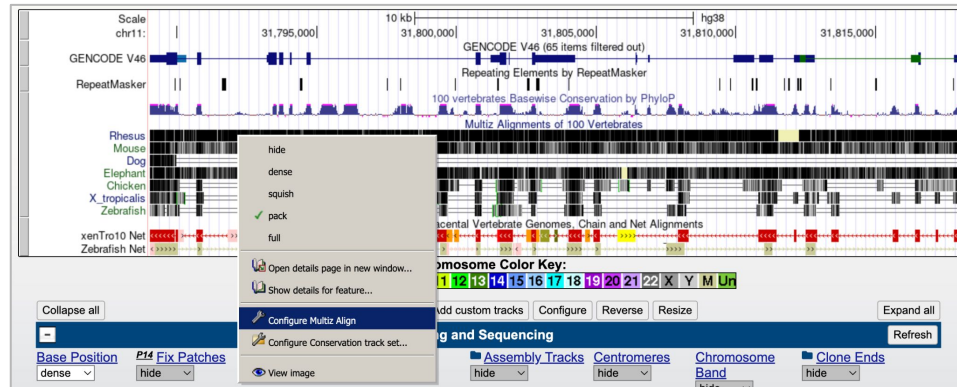
# Génomique comparative : régions conservées chez les vertébrés

- Dans la catégorie **Comparative genomics**, activez l'affichage **pack** de la piste **Conservation**. Cette piste s'affiche entre les annotations GENCODE\_V46 et les régions répétitives
- Faites remonter la piste **RepeatMasker** pour la placer entre les pistes GENCODE\_V46 et Conservation



# Configuration de la piste d'annotation Conservation

- Cliquez droit (contrôle-clic) sur l'image de conservation à la hauteur où s'affichent les espèces et sélectionnez **Configure MultiZ Align**.
- Dans la fenêtre d'options qui apparaît, **cochez quelques espèces de votre choix**
  - **Veillez à panacher** (essayez d'avoir une ou deux espèces de chaque groupe plutôt qu'un tas d'espèces du même groupe).
  - Pour une raison technique, le génome du chien présente des lacunes à cet endroit du génome. **Désactivez l'affichage du chien (dog)** et activez celui d'un ou deux autres mammifères du même groupe.
  - Dans la catégorie **Mammal**, **cochez toutes les espèces**. Notez que les catégories précédentes contiennent également des mammifères (Primates, Euarchontoglires, Laurasiatheria). La catégorie Mammal présente des espèces plus éloignées (marsupiaux, monotrèmes), qui sont utiles pour visualiser les régions les plus conservées entre mammifères.
- Cliquez **Apply**.
- **Dézoomez d'un facteur 1.5** pour observer le contexte aux alentours du gène.



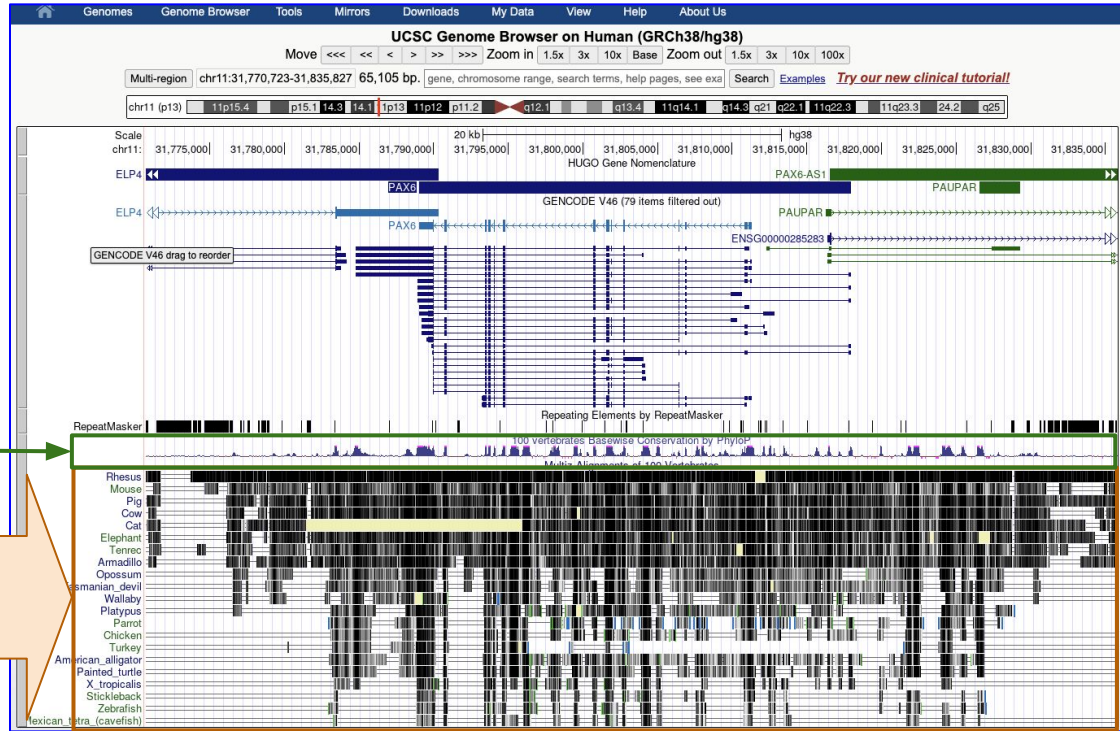
## Exercice 2. Conservation de la région génomique PAX6 chez les vertébrés

La carte de conservation génomique comporte deux parties.

1. La partie supérieure affiche un **profil de conservation dans 100 génomes de vertébrés**. La hauteur du profil indique le pourcentage de positions identiques (PPI) à chaque position du génome. Notez que l'échelle verticale va de 50% à 100%, pour mieux faire ressortir les régions conservées.
2. La partie inférieure indique, sous forme d'une échelle de gris, le **profil de conservation par espèce** (pour celles que vous avez sélectionnées).
3. Les zones marquées en jaune correspondent à des trous de séquençage (on ne dispose pas de la séquence).

Conservation chez  
100 vertébrés

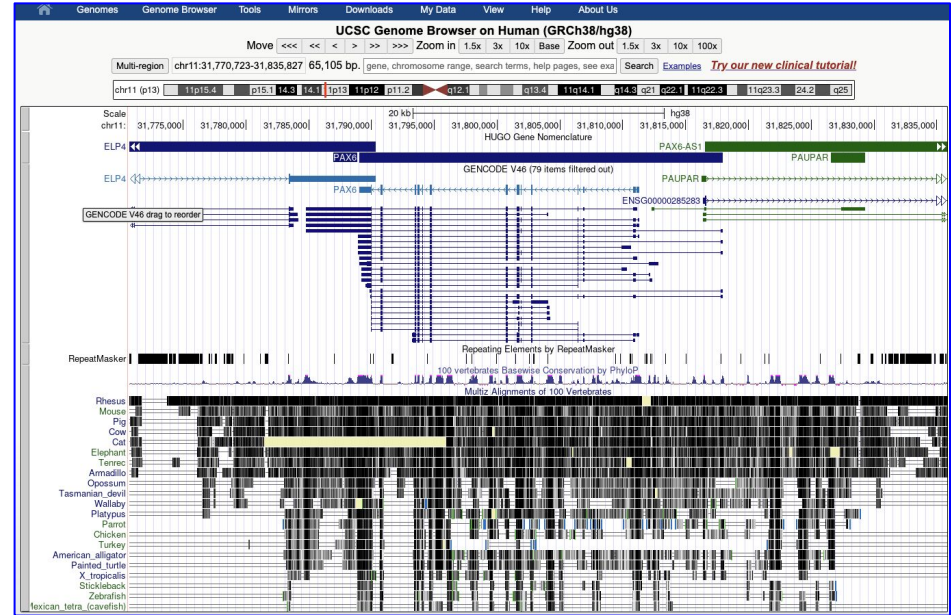
Profils par  
espèce



## Exercice 2. Conservation de la région génomique PAX6 chez les vertébrés

Sur Ametice, ouvrez le questionnaire du TP3 et répondez aux questions de l'**Exercice 2**.

1. Zoomez sur la carte génomique que vous avez générée et parcourez le gène PAX6 sur toute sa longueur. Sur le profil de conservation génomique 100 vertébrés, dans quelles régions du gène retrouve-t-on généralement les régions conservées (plusieurs réponses possibles):
  - ☐ 5'UTR
  - ☐ partie codante des exons
  - ☐ Introns
  - ☐ 3'UTR
2. Sur la carte des profils par espèces, quel est l'ordre des degrés de conservation (une seule réponse)
  - Primates > Poissons > Mammifères non primates
  - Poissons > Mammifères non primates > Primates
  - Primates > Mammifères non primates > Poissons
  - Poissons > Primates > Mammifères non primates



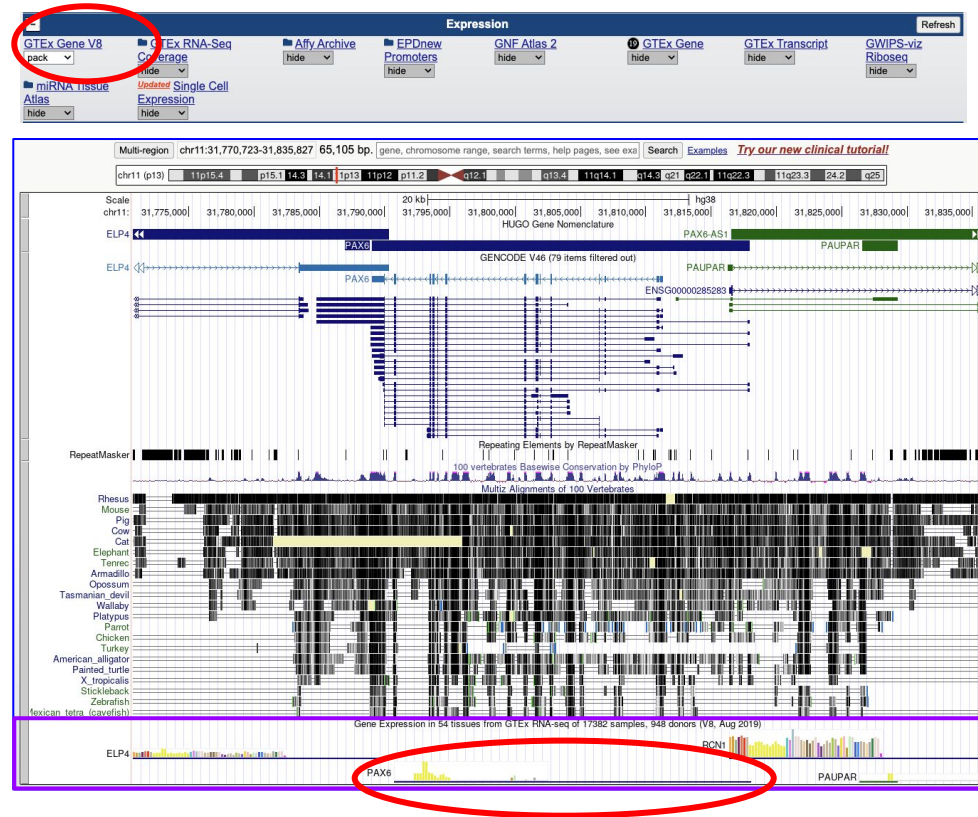
# Expression tissulaire du gène PAX6



# Profils tissulaires de transcription dans GTEx (Genotype-Tissue Expression)

Nous allons maintenant ajouter à notre carte génomique une piste d'annotation de la base de données [GTEx](#) (Genotype-Tissue Expression). GTEx contient des données de transcriptome (mesure quantitative de tous les transcrits produits par un génome) dans des échantillons de 54 tissus prélevés chez 948 personnes adultes.

- Dans la catégorie **Expression**, activez l'affichage **pack** de **GTEx\_Gene\_V8**.
- Cliquez sur l'icône du gène PAX6 sur la piste GTEx\_Gene\_V8, et examinez le profil d'expression tissulaire.

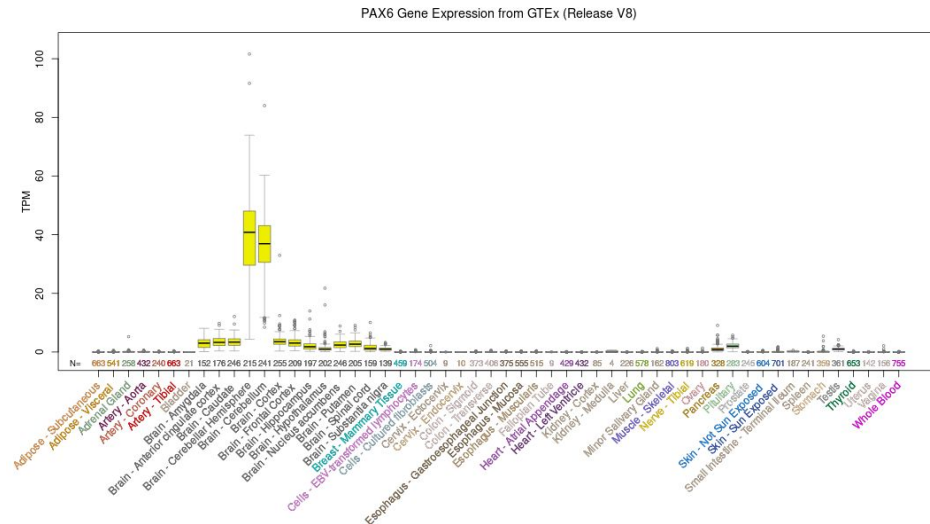




### Exercice 3. Profil d'expression tissulaire de PAX6

Sur Ametice, ouvrez le questionnaire du TP3 et répondez aux questions de l'**Exercice 1**.

1. Indiquez le niveau d'expression de PAX6 dans les tissus suivants (fort/faible/indéetectable)
  - ☐ Cervelet
  - ☐ Autres tissus du cerveau
  - ☐ Muscles
  - ☐ Poumons
  - ☐ Sang
2. GTEx se base sur des échantillons adultes. D'après la fonction de PAX6, on s'attend à observer également un profil d'expression très spécifique durant le développement embryonnaire. Vrai / Faux





## Informations complémentaires

- **Page “training” de UCSC Genome Browser ([genome.ucsc.edu/training](http://genome.ucsc.edu/training))**
  - Guides d'utilisation (en anglais)
  - Courtes vidéos pour apprendre à manipuler les nombreuses fonctionnalités du site Web