

# Predicting “pain genes”: multi-modal data integration using probabilistic classifiers and interaction networks

Na Zhao<sup>1</sup>, David L Bennett<sup>1</sup>, Georgios Baskozos<sup>1, \*</sup>, Allison M Barry<sup>1, \*, +</sup>

<sup>1</sup>Nuffield Department of Clinical Neurosciences, University of Oxford, OX3 9DU, UK

\*These authors contributed equally to this work

<sup>+</sup>for correspondence: AMB, [allison.barry@ndcn.ox.ac.uk](mailto:allison.barry@ndcn.ox.ac.uk)

---

**Abstract** Accurate identification of pain-related genes remains challenging due to the complex nature of pain pathophysiology and the subjective nature of pain reporting in humans, or inferring pain states in animals on the basis of behaviour. Here, we use a machine learning approach to identify possible “pain genes”. Labelling was based on a gold-standard list of genes with validated involvement across pain conditions, and was trained on a selection of -omics (eg. transcriptomics, proteomics, etc.), protein-protein interaction (PPI) network features, and biological function readouts for each gene. Multiple classifiers were trained, and the top-performing model was selected to predict a “pain score” per gene. The top ranked genes were then validated against pain-related human SNPs to validate against human genetics studies. Functional analysis revealed JAK2/STAT3 signal, ErbB, and Rap1 signalling pathways as promising targets for further exploration, while network topological features contribute significantly to the identification of “pain” genes. As such, a PPI network based on top-ranked genes was constructed to reveal previously uncharacterised pain-related genes including CHRFAM7A and UNC79. These analyses can be further explored using the linked open-source database at <https://livedataoxford.shinyapps.io/drg-directory/>, which is accompanied by a freely accessible code template and user guide for wider adoption across disciplines. Together, the novel insights into pain pathogenesis can indicate promising directions for future experimental research.

---

## Introduction

The identification of pain-related genes remains challenging due to the heterogeneity and multifactorial nature of the disease. “Pain” encompasses an unpleasant sensory and emotional experience, acute and chronic states, and actual or potential tissue damage spanning a spectrum of conditions from UV radiation (sunburn) to diabetic neuropathy (*Raja et al. (2020)*).

High-throughput sequencing techniques (“-omics”), have revolutionised the identification of molecular markers and pathways, with technologies like transcriptomics and translomics providing large gene expression data sets which enable identification of differentially expressed genes and data-driven biomarker discovery. With such diversity in pain characterizations, we expect equally diverse underlying mechanisms. Even so, informative parallels are also seen across conditions. For example, decades of work in the migraine field have culminated to the development of clinical CGRP antibodies – a hallmark treatment for acute and chronic migraine (*Edvinsson et al. (2018)*; *Goadsby et al. (2017)*). This pathway is also being explored in the context of neuropathic and inflammatory pain, as well as non-migraine headache pain due to the underlying mechanism of sensory neuron sensitization (*Schou et al. (2017)*; *Paige et al. (2022)*).

As we continue to accumulate -omics data, we need effective strategies to present and integrate these datasets together, while also considering multi-modal data from external sources.

Machine learning (ML) approaches are well suited to this challenge, and the availability of probabilistic models (that is – models which give the probability something belongs to a specific class) allows us to assign a probability score to each instance. Recently, ML integration with gene expression data has gained popularity in biomarker discovery (reviewed in *Zhang et al. (2021)*) and clinical diagnosis (*Kumar et al. (2023)*).

In the context of pain, this has proven to be powerful: for example, predicting patient classification for painful vs painless diabetic neuropathy has highlighted factors relevant to the painful class (*Baskozos et al. (2022)*). Pre-clinically, there is also a need to integrate multi-modal data, both to give insight to features underlying genes involved in pain, as well as to lay foundation for future studies.

Here, we trained both off-the-shelf probabilistic classifiers and ensembles of classifiers to produce a predictive pain score based on an expansive feature space to address this gap. Features include cross-species transcriptomic and translomics datasets as well as proteomic data, network topology, genetic structure (eg. GC content), and functional pathway assignments.

The top-performing model was employed to predict a class probability score (“pain score”) for each gene, and the gene candidates with highly ranked predicted pain scores are subjected to downstream functional analysis, including validation against human genetics datasets. High ranking features, such as DRG translomics data and protein-protein network interactions were highlighted and further examined in the context of pain, while JAK2/STAT3 signal, ErbB, and Rap1 signalling pathways were identified as promising targets for future exploration.

These scores were curated into an open-access database (<https://livedataoxford.shinyapps.io/drg-directory/>) alongside experimental datasets. Here, we have integrated the STRING DB to facilitate the visualisation of pain-related genes in the context of their network associations (a high ranked feature), building on previous work by Perkins et al., 2013 (*Perkins et al. (2013)*).

In addition to data integration, a fundamental component to effective data use is access. As more big data continues to be produced, better data management practices are needed (*Boeckhout et al. (2018)*): Raw data repositories limit use to those with bioinformatic skills while extensive supplemental data tables quickly become cumbersome.

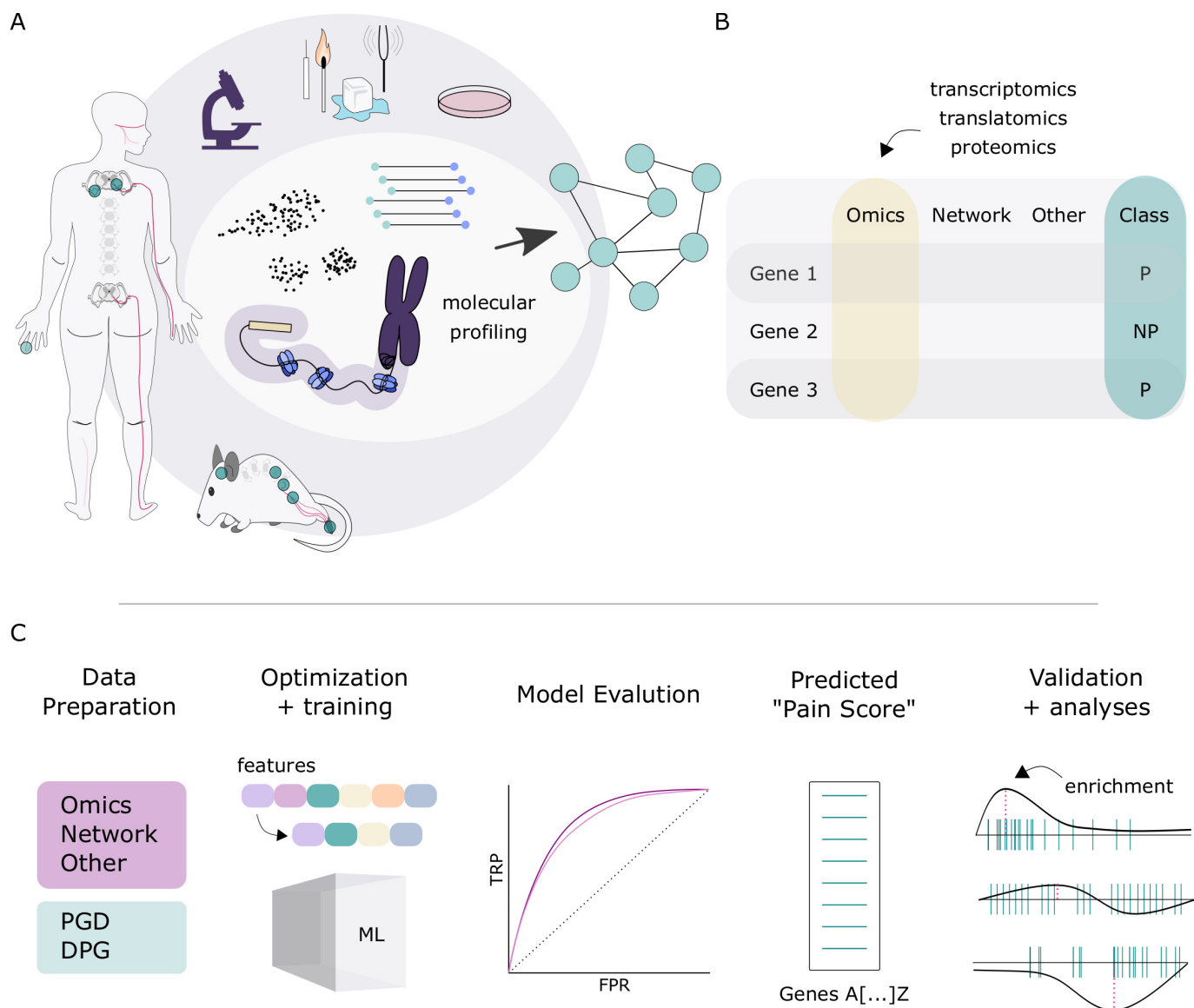
We have thus paired our database to a user-guide for researchers to reimplement these visualization for other omics studies and disease states: This reproducible Shiny-based framework can simply integrate multiple omics datasets and generate composite visualisations and, where relevant, add PPI networks of condition-related genes to help inform candidate selection for downstream experimentation. Additionally, it addresses a legacy gap commonly seen with in-house databases.

Together, this study presents a novel disease-gene identification framework by integrating diverse datasets through machine learning to gain mechanistic insights of pain. Paired to an open access database with an emphasis on PPI networks, this will allow researchers to more effectively select targets, and, ultimately - lead to better data utilisation and increased impact of each study.

## Results

Here, we use a machine learning approach to identify possible “pain genes” (Figure 1). Because of the diversity in the studies used to generate initial labels, the term “pain” here is used as a broad characterization across acute and chronic states. Here, the predicted pain genes thus represent more generalisable genes across conditions, opposed to, for example, predictions tailored to neuropathic pain specifically.

Selecting gene labels in the experimental design is not a trivial task, due to the variable amount of research surrounding each gene in the context for pain. Here, we opt for a highly stringent approach, requiring functional, in vivo validation in mice or detailed characterization in humans (see methods for full details). With this, there is an underlying expectation that a number of genes being studied in the context of pain do not yet reach our threshold for inclusion as “pain”, even though



**Figure 1.** Experimental overview. A. Data integration schematic, where data can be integrated across species and modalities using machine learning. B. Example data input, when information for each gene was gathered across modalities and labelled either pain ("P") or no-pain ("NP") based on gold-standard lists from prior literature. C. Pipeline schematic for predicting "pain" genes, from data preparation through to validation.

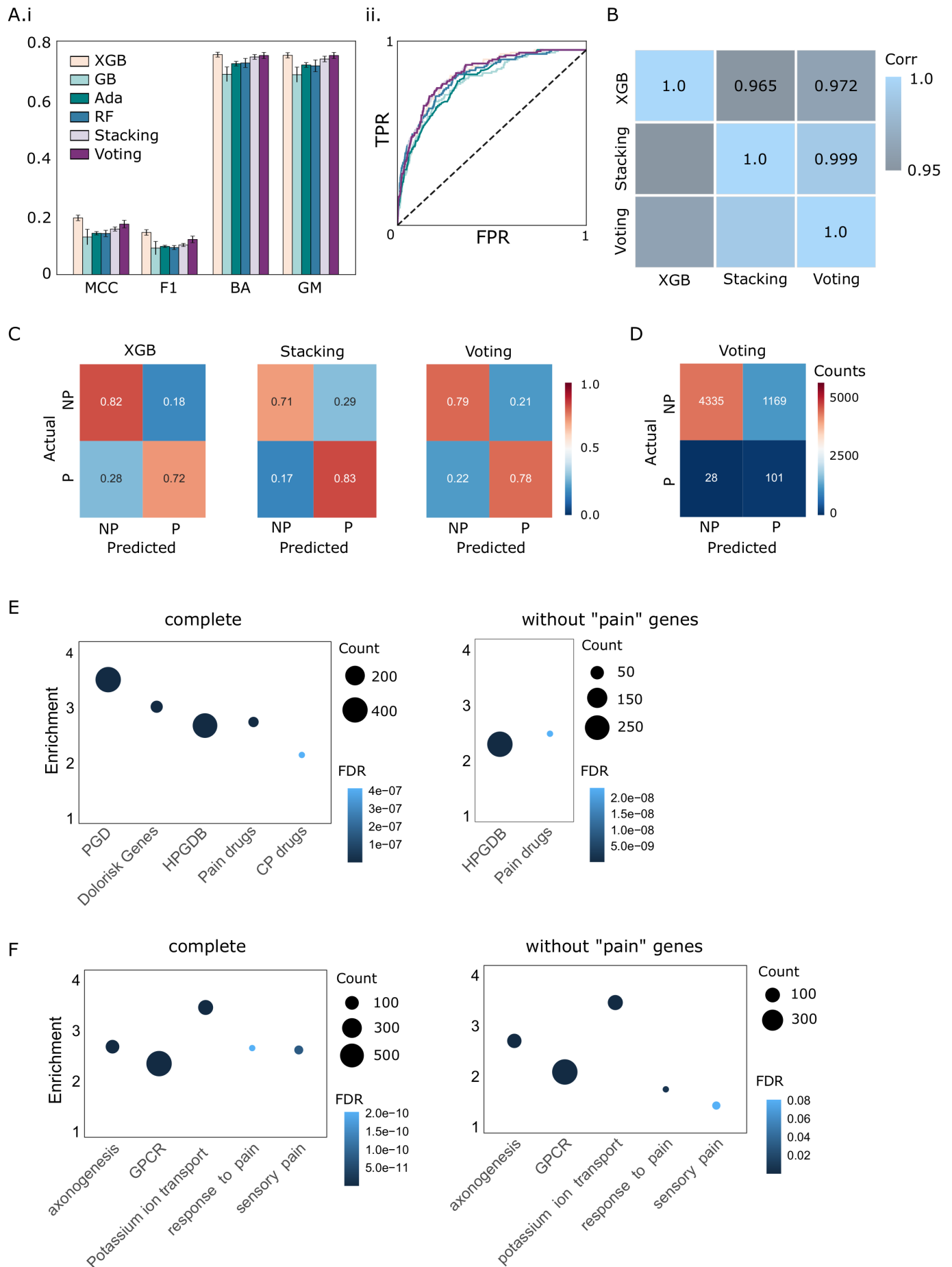
they are likely to be relevant: ie. gene classified as pain when labelled non-pain is not necessarily incorrect biologically, but a comment on our current knowledge in the field.

## Feature Selection and Exploration

Scikit-learn (Sklearn) was used to classify "pain"/"no pain" genes in Python using a variety of algorithms (*Pedregosa et al. (2011)*). We performed initial feature selection using the Gradient Boosted Trees (XGBoost) Classifier after hyperparameter tuning by the Optuna framework (*Akiba et al. (2019)*) and the shap library (*Lundberg and Lee (2017)*). Features were ranked for importance using their SHAP (SHapley Additive exPlanations) values. Using a backward elimination method, the top 23 features were selected and used to train models.

## Model Training and Performance

We used a labelled dataset of known genes, out of which we have labelled 429 genes found in the Pain Genes Database (PGD) (*LaCroix-Fralish et al. (2007)*) and DOLORisk Priority Group (*Themistocleous et al. (2023)*) as Pain (P) and the remaining as



**Figure 2.** Classifier evaluation. Caption continued on the next page.

**Figure 2.** A. Six classifiers were trained, with performance scores for (i) Mathew's Correlation Coefficient (MCC), F1 score (F1), Balanced Accuracy (BA), and G mean (GM) evaluated. (ii) ROC curves show similar true positive rate (TPR) and false positive rate (FPR) across classifiers. B. Ranked prediction score correlations for the top three performing classifiers. C. Confusion matrix for the top three classifier in the test dataset. D. Confusion matrix by count for the top classifier in the test dataset, highlighting the imbalance in the dataset. E-F. Prediction scores for the "pain" class were ranked for enrichment analyses. Left; with all genes. Right; with all "no-pain" genes to prevent leakage. E. Predicted "pain gene" enrichment against curated known pain lists and approved drug targets. HPGBD: Human Pain Genetics Database, containing SNPs against relevant pain states/disorders. CP: chronic pain. F. Enrichment analyses against relevant biological pathways. Related to Supplemental Figure 1.

Non-Pain (NP) genes (Figure 1B). These genes represent the gold-standard of highly confident targets in pain. Six classifiers, including random forest (RF), AdaBoost (Ada), gradient boost (GB), Gradient Boosted Trees (XGBoost); and Stacking and Voting ensembles were trained to classify pain/non-pain genes based on multi-omics, genomic and network topological data (Figure 2) (*Ho (1995); Schapire (2013); Friedman (2001); Chen and Guestrin (2016)*).

### Model Evaluation and Selection

We assessed model performance based on four metrics: F1 score, Matthews Correlation Coefficient (MCC), balanced accuracy (ACC), and geometric mean (GM) (Figure 2A-B) (*Chicco and Jurman (2020); Brodersen et al. (2010); Espíndola and Ebecken (2005)*). The voting classifier achieved the best performance (Voting; MCC =  $0.1787 \pm 0.0108$ , GM =  $0.7581 \pm 0.0132$ ), followed by the XGBoost classifier (XGBoost; MCC =  $0.1995 \pm 0.0149$ , GM =  $0.7535 \pm 0.0192$ ). The next best performance is the Stacking Classifier (Stacking; MCC =  $0.1582 \pm 0.0094$ , GM =  $0.7452 \pm 0.0132$ ) and the Random Forest Classifier (RF; MCC =  $0.1413 \pm 0.0144$ , GM =  $0.7183 \pm 0.0266$ ). Next, AdaBoost achieved a moderate performance (Ada; MCC =  $0.1487 \pm 0.0070$ , GM =  $0.7300 \pm 0.0061$ ). Finally, GradientBoost Classifier achieved the worst performance (GB; MCC =  $0.1325 \pm 0.0171$ , GM =  $0.6951 \pm 0.0201$ ). Results are summarized in Table 1.

Here, we prioritized GM over MCC, while also looking at balanced accuracy and F1 scores. Together, these four metrics provide insight into the performance of imbalanced datasets, while GM controls accuracy in both classes and ranks higher classifiers that are equally good in both classes, regardless of their size (as it uses the ratios TPR, FPR). In our case, this is important as of course not many genes are validated as "pain genes". This allows us to choose a classifier with a high number of true positives while still prioritizing the true negatives. (see Fig 2C).

**Table 1.** Model performance (mean  $\pm$  SD) of six classifiers

Algorithm	MCC	Balanced Accuracy	GM	F1
XGBoost	<b><math>0.1995 \pm 0.0149</math></b>	$0.7574 \pm 0.0169$	$0.7535 \pm 0.0192$	<b><math>0.1530 \pm 0.0137</math></b>
GradientBoost	$0.1325 \pm 0.0171$	$0.6996 \pm 0.0155$	$0.6951 \pm 0.0201$	$0.0998 \pm 0.0165$
AdaBoost	$0.1487 \pm 0.0070$	$0.7331 \pm 0.0063$	$0.7300 \pm 0.0061$	$0.1022 \pm 0.0074$
RandomForest	$0.1413 \pm 0.0144$	$0.7290 \pm 0.0195$	$0.7183 \pm 0.0266$	$0.0934 \pm 0.0098$
Stacking	$0.1582 \pm 0.0094$	$0.7519 \pm 0.0115$	$0.7452 \pm 0.0132$	$0.1034 \pm 0.0072$
Voting	$0.1787 \pm 0.0108$	<b><math>0.7586 \pm 0.0129</math></b>	<b><math>0.7581 \pm 0.0132</math></b>	$0.1267 \pm 0.0091$

The Voting Classifier ensemble was selected as the top performing model: Even so, we see a high correlation between prediction scores from top three highly ranked models (XGBoost, Voting Classifier, Stacking Classifier), suggesting that all three algorithms predict pain genes to a similar degree and that our predictions are sufficiently robust to not depend on a single classifier or set of parameters (Fig 2B). Although the voting classifier has a lower MCC and F1 score than XGBoost (2A), it predicts more true positive genes (Fig 2C-D), and has the highest GM and balanced accuracy, which are two important metrics in evaluating imbalanced datasets.

When building classifiers, it is important to test the external validity of the model. In this design, we do not have a

separate cohort to probe, as that would require a separate genome. Instead, we rely on relevant human genetics data, taking advantage of a curated list of single nucleotide polymorphisms (SNPs) relevant in pain from the Human Pain Genetics Database (HPGDB) *Meloto et al. (2018)*. These represent relevant genetic polymorphisms in the context of pain across a broad range of conditions, in line with our original dataset. While some have been functionally validated and overlap with our gold-standard list of "pain genes", many have not been followed up, representing likely, but unvalidated targets of pain. As such, these were not labelled as true pain genes in the original experimental design.

Functional GSEA was used to internally validate prediction scores from the top three highly ranked classifiers (XGBoost, Voting, and Stacking) using the known pain-related genes sets. In addition to internal validation against the sets used for labelling (Pain Genes Database and DOLORisk pain genes), prediction scores were externally validated by comparing them against the independent Human Pain Genetics Database as well as a curated list of pain-related drug-targets (Fig 2E). GSEAs were repeated for gene lists with true labels removed (that is, removing the 429 "pain" labelled genes where overlap occurs).

Crucially, external validation using genes in Human Pain Genetic Database after removing overlaps with our training data supports the use of the Voting Classifier, as it showed the highest enrichment for pain-related SNPs in the predicted pain ranking across genes in an unbiased dataset, as well as to drug targets for approved drugs relevant to pain and chronic pain (CP) (Figure 2E). GSEA plots for each classifier against the HPGDB without "pain" labels is shown in Supplemental Figure 1.

As expected, prediction results from the voting classifier also shows a high enrichment score for pain-related pathways including response to pain (GO:0048265) and sensory pain (GO:0051930) with and without the inclusion of pain-labelled genes (Fig. 2F).

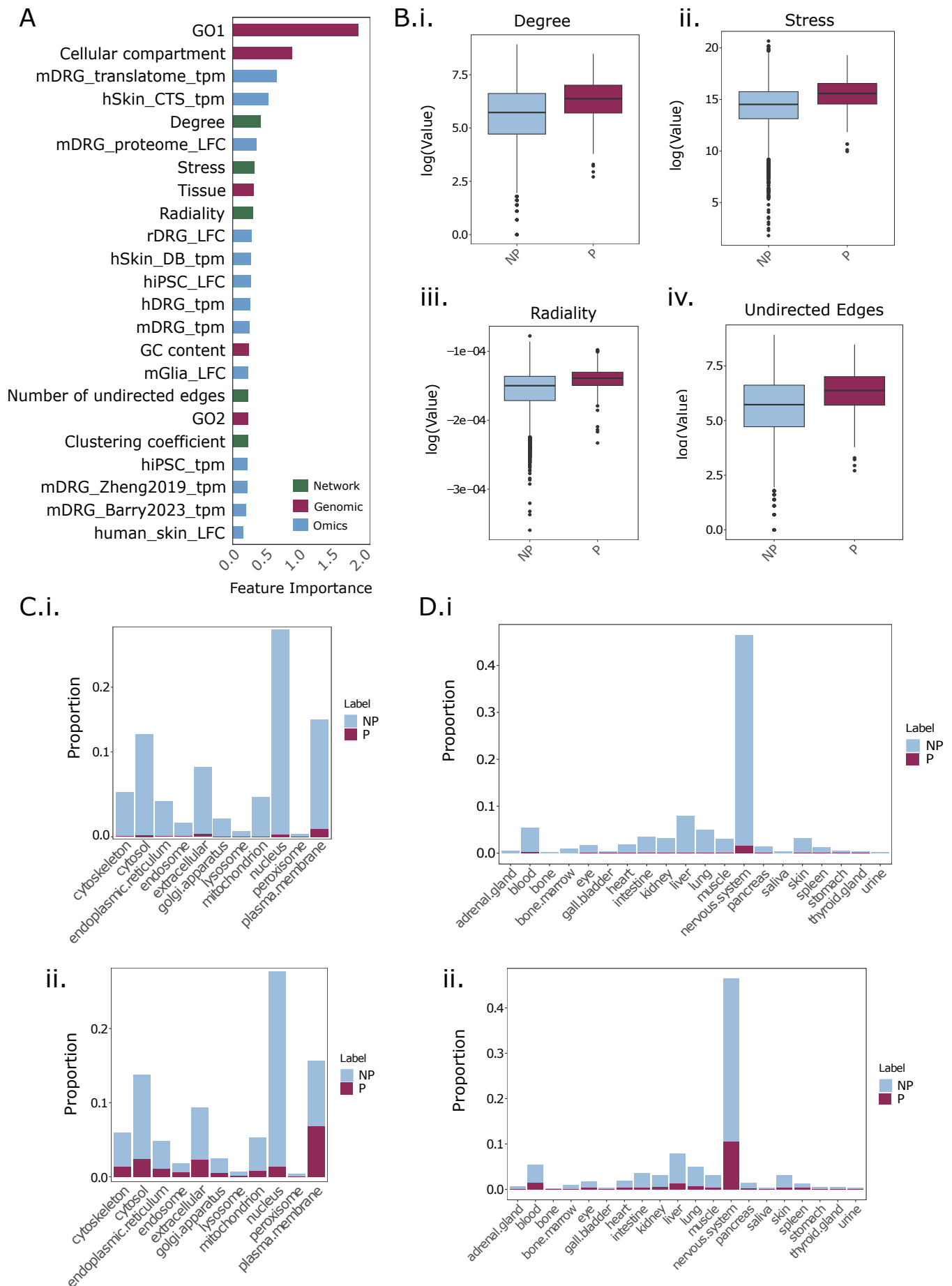
## Feature analysis

In addition to studying which genes have the highest predictive pain scores, which features make up these predictions are also of interest. Because the voting classifier is an ensemble, SHAP values across classifiers were weighted and combined to reflect the most relevant features underlying the classification (Figure 3).

The most highly informative feature was GO pathways. This is intuitive, given that well known pain genes used as true positive labels are associated with relevant GO pathways, thus classification by pathway is likely to be effective. Even so, we see here that even high order GO pathways from the GO slim collection are important, such that classification is not dependant on the smaller, pain-specific terms (e.g. "response to pain" and/or "sensory pain", which are not present in the GO slim collection used to build the feature space)".

The cellular component of the corresponding protein is the next highly ranked, with tissue expression following shortly there after. Notably, three network-based features extracted from the STRING database were also seen in the top 10 features (radiality, stress, and the number of undirected edges), suggesting that highly connected proteins are more likely to be involved in pain responses. These finding support a guilt-by-association approach showing that "pain" genes are likely to share similar functions as their interacting partners and aggregate in local interactome neighbours. Further exploration of network features shows that pain genes demonstrate higher radiality, stress, degree, and edge count (Fig. 3B). This suggest that they are likely to be hub proteins, which are the most highly connected central proteins in PPI networks (*Higurashi et al. (2008)*).

In terms of predictive "-omics" datasets, log fold changes in mDRG proteomics, as well as translome murine expression studies also come up in the top 10, even above human DRG expression. Here, this may reflect that as the "functional building blocks" of a cell, changes in protein expression carries significant weight in predicting pain relevance. Alternatively, this may be a stronger comment towards a bias in the labels used, as some of our best documented "pain genes" were initially highlighted through rodent -omics studies, and/or that studies focus commonly focus on candidates with available antibodies due to technical limitations. As such, this also hints at the limitation of the true positive labels, which is discussed in below.



**Figure 3.** Voting classifier feature analysis. Caption continued on next page



**Figure 3.** A. Ranked feature importance (weighted SHAP values) of the voting classifier. B. distribution of four top ranked network features in pain (P) and non-pain (NP) genes, extracted from the STRING DB through Cytoscape. 'Degree' refers to the number of edges linked to a node whereas 'stress' counts the number of shortest paths passing through a node. C. Proportion of genes in each cellular compartment (i) before and (ii) after prediction by the voting classifier. D. Proportion of genes in each tissue type (i) before and (ii) after prediction by the voting classifier. Species denoted as h (human), m (mouse), r (rat) (eg. hDRG\_tpm = human DRG). TPM, transcripts per million; LFC, log fold change. Healthy human skin RNA-seq from two cohorts, CTS (carpal tunnel syndrome) and DB (diabetic) cohorts are highlighted as "hSkin\_CTS\_tpm" and "hSkin\_DB\_tpm" respectively.

## Functional analyses

The predicted pain scores from the voting classifier were subjected to functional enrichment analysis.

### GO and KEGG analysis

Given the high importance of GO in predicting pain-genes (Fig. 3A), we conducted GO functional and KEGG pathway enrichment analyses to explore which GO terms are enriched in pain genes (Fig 4A).

GO term analysis of top 10% genes highlights GO term enrichment for potassium ion transport, regulation of membrane potential, positive regulation of cell proliferation, positive regulation of calcium ion concentration, positive regulation of ERK1 & ERK2 cascade, chemical synaptic transmission (Fig 4A.i). GO analysis of bottom 10% ranked genes show non-pain related GO terms, as expected (Fig 4A.ii).

KEGG analysis of top 10% genes shows up-regulation of Ras-, MAPK- and Erb-signalling pathways as promising targets for future study. It also involves the Natural Killer-mediated cytotoxicity pathway, suggesting this analysis can capture the critical roles of the immune system in pain, which was recently reviewed at *Kim et al. (2023)*.

### Network contextualization

Network features are highly ranked for predictive importance (Figure 3A-B). As such, we have integrated predictive scores of the voting classifier with the STRING DB through <https://livedataoxford.shinyapps.io/drg-directory/> ("network analyses" tab, Figure 4C).

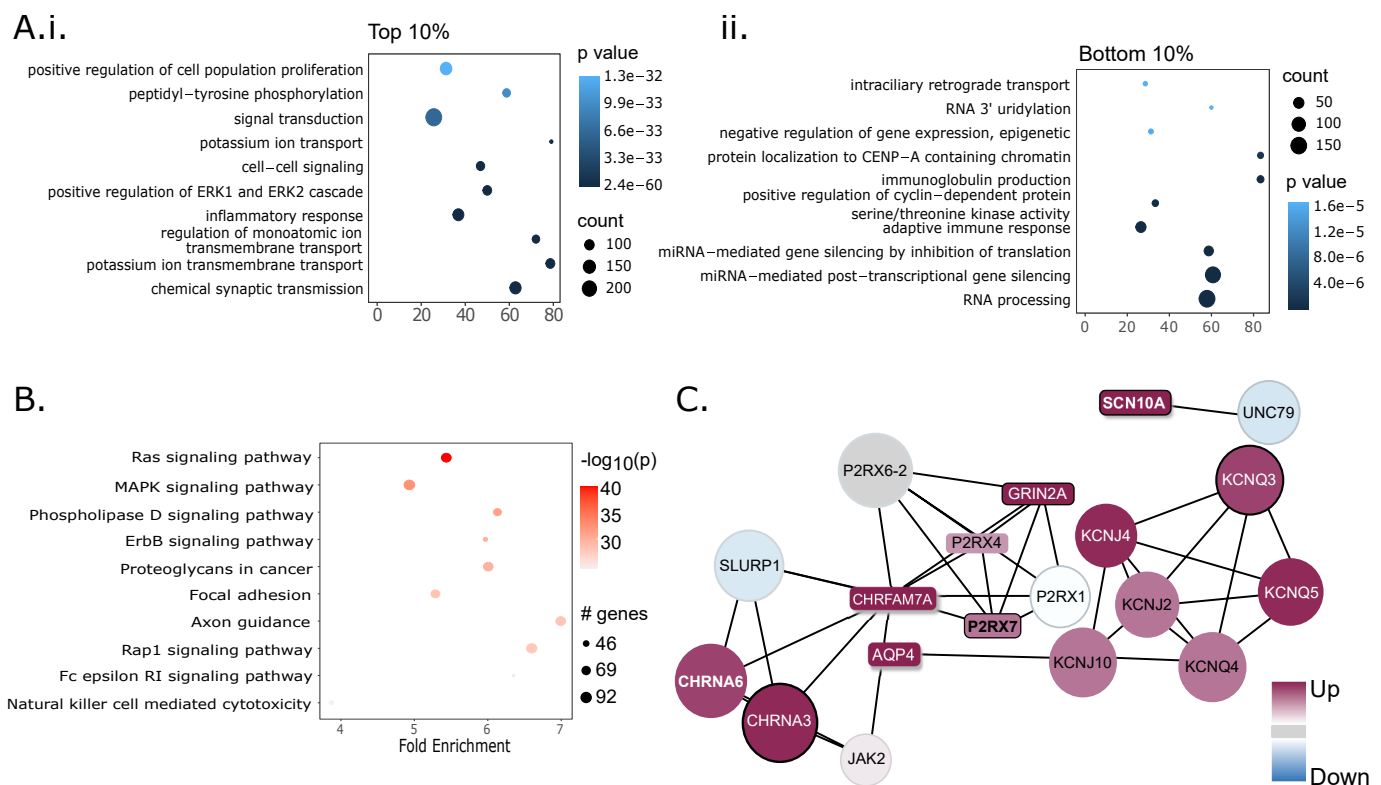
The network can be annotated using information on known pain-associated genes from several sources: DOLORisk Priority Group, Human Pain Genetics Database, and the Pain Genes Database (*Themistocleous et al. (2023)*; *Meloto et al. (2018)*; *LaCroix-Fralish et al. (2007)*). In addition, users can enrich these networks by using data from pain-focused gene expression studies to highlight genes that change expression in each condition or pairs of genes showing correlated expression patterns across different experiments.

The top 5 highly ranked genes are input to construct a PPI network (Fig 4C). Most of the extracted proteins demonstrate a high score, emphasising the significance of PPI networks in pinpointing disease-related proteins by elucidating their functional correlations with known pain genes.

## Discussion

The current study proposed a gene-centred machine learning approach to identify pain-related genes using multi-omics, PPI network and genomic data. ML has been used to identify disease biomarkers using multi-omics data (*Reel et al. (2021)*). However, due to the small sample sizes with high-dimensional features, training a large-scale generalizable ML model with multi-omics data alone can be challenging. Moreover, ML has also been successfully used to predict disease biomarkers using PPI network topological data (*Wu and Wang (2023)*). A study by Yu. et al used network topological features to predict proteins that may cause neurodegenerative disease, and used multi-omics data for validation and target selection (*Yu et al. (2020)*). Combining these approaches, we adopted a gene-centred approach to derive pain-related proteins from multi-omics, PPI networks, and genomic features. While previous studies focus more on the performance of the prediction algorithms, we brought attention to the biological explanation of predictive features and provided a list of genes as





**Figure 4.** Functional analysis of high-ranked prediction scores from the Voting classifier. A. Top 10 GO term enrichment terms for top 10% and bottom 10% of ranked genes. B. Top 10 KEGG pathway terms for top 10% of ranked genes. C. PPI network of top 5 highly ranked genes by the Voting Classifier. Annotation from <https://livedataoxford.shinyapps.io/drg-directory/>: Colour gradient = Enrichment (red for PES > 0.02, blue for PES < 0.02); grey background = non-significant; Grey colour = absent from current dataset; Box shape = genes contained in the Pain Genes Database; Shadow = genes contained in the DOLORisk Priority Group; Black border = genes contained in the Human Pain Genetics Database.

promising targets for future pain studies.

Out of six classifiers used, the voting classifier achieved the best performance (MCC = 0.1787 ± 0.0108, GM = 0.7581 ± 0.0132). It predicts the class label based on the argmax of the sums of the predicted probabilities by four individual classifiers: XGBoost, AdaBoost, Stacking Classifier, and GradientBoost Classifier.

Because many pain genes are known from rodent studies in the Pain Gene Database, external validation against the HPGDB and drug targets relevant to pain was crucial to establishing how to extrapolate probabilistic scores. The strong enrichment of the HPGDB, even after removing the gold-standard pain genes suggest these results can a) be extrapolated to humans, and b) are relevant in the context of human genetic studies. The enrichment of “pain drug” targets in high-ranking predictive scores further suggests that we are capturing a population of potentially-druggable targets (Figure 2E).

This study included data from multiple species, so we are likely capturing targets relevant cross-species, opposed to human-specific targets. Even so, gene “conservation scores” are themselves not predictive (Figure 3).

## Feature exploration

GO and KEGG functional analysis were conducted to identify enriched GO terms associated with pain genes. Analysis of the top 10% ranked genes revealed molecular functions closely linked to neuropathic pain pathogenesis. Among the enriched terms, five are linked to membrane transport processes (potassium ion transport, chemical synaptic transmission, regulation of membrane potential, regulation of ion transmembrane transport, and positive regulation of cytosolic calcium ion concentration), likely contributing to neuronal hyperexcitability associated with neuropathic pain (Choi et al. (2023)). The molecular functions of inflammatory response and cell-cell signalling are associated with production of pro-inflammatory molecules that sensitise nociceptive neurons underlying pain sensation (Ji et al. (2016)). The ERK1/2 is a characterised impor-

tant signalling pathway in pain, and its activation is engaged in modulating the pain sensitivity (*Kondo and Shibuta (2020)*). GO analysis of bottom 10% ranked genes show non-pain related GO terms, as expected (Fig. 4A.ii).

KEGG analysis of the top 10% genes also revealed pathways associated with pain pathogenesis and progression. Ras signalling, implicated in NGF signalling through TrkA via the PI3K/Ras pathway leading to TRPV1 activation, is crucial in pain sensation (*Bonnington and McNaughton (2003)*). Ras is also frequently discussed in relation to the Ras/Raf/MAPK pathway: both pathways are enriched here. MAPK has a long-standing role in pathophysiology *Ji et al. (2009)*, thus dissecting the interaction with Ras could lead to new avenues of druggable targets. Additionally, Rap1 signalling has also been described in the context of inflammatory pain through Epac1 (*Singhmar et al. (2016)*).

NRG1/ErbB signalling is significant in spinal cord injury (SCI)-induced chronic neuropathic pain (*Tao et al. (2013)*), and NRG1 has been implicated in axonal development as well as regeneration after nerve injury in the periphery *Fricker et al. (2009, 2011)*. Given the relevance to the neuropathic pain model of SCI, as well as recent work highlighting neuronal death after traumatic nerve injury in mice *Cooper et al. (2024)*, one can speculate how this pathway may be highly relevant across neuropathic pain conditions where axonal damage occurs. In line with this, the Natural Killer-mediated cytotoxicity pathway underscores cytotoxic immunity's role in response to nerve injury (*Davies et al. (2019, 2020)*). Together, this suggests these are highly relevant candidates for further study.

Cellular compartment is the feature with second highest importance. Membrane proteins are heavily involved in pain reception. During nociception, high threshold stimuli that could lead to injury result in activation of ligand gated ion channels such as Transient Receptor Potential (TRP) channels in nociceptors in the periphery; subsequent opening of cation channels (potassium and sodium channels) results in depolarization and action potential propagation along afferent sensory fibres to the dorsal horn synapse (*Middleton et al. (2021)*). Consequently, the cellular compartment containing the highest number of pain-related genes is the plasma membrane (Fig. 3C). Interestingly, there are also a large number of pain-related genes in the nucleus, cytosol, and extracellular space, which serve as interesting insights for future target discovery.

While some features are highly predictive, others, such as species conservation rank less important in the current study. Here, it is unclear if this is a true trend, or if our bias in rodent studies covers a trend.

## Network Analysis

We developed and employed the Pain RNAseq Hub (<https://livedataoxford.shinyapps.io/drg-directory/>) to visualize the top 5 pain genes within their protein-protein interaction (PPI) network context, incorporating multi-omics data and pain-related annotations.

The extracted proteins include various characterised pain-related genes including CHRNA6, SCN10A, P2RX7, KCNQ5, KCNQ4, KCNQ3, CHRNA3, AQP4, JAK2, P2RX4. More importantly, we found several previously uncharacterised pain-related genes, which will be discussed below.

One such example is ADORA2A, a gene that encodes the adenosine A2A receptor. Binding of adenosine to the adenosine A2A receptor during stress initiates potentially destructive inflammatory cascades that lead to the activation of immune cells and release of proinflammatory mediators *Flögel et al. (2012)*. It was found that prolonged accumulated circulating adenosine contributes to chronic pain by promoting immune-neuronal interaction and revealed multiple therapeutic targets (*Hu et al. (2016)*). A2A receptor agonists have been shown to block adenosine and thus inhibit the release of proinflammatory mediators (*Cekic and Linden (2016)*) while related genes, ADORA2B and ADORA3, can cause nociceptor hyperexcitability and promote chronic pain (*Wahlman et al. (2018); Middleton et al. (2021)*). Together, this makes ADORA2A an attractive target to follow up.

Another interesting gene for future exploration is CHRFA7A, a uniquely human fusion gene that functions as a dominant negative regulator of alpha 7 acetylcholine nicotinic receptors ( $\alpha 7nAChR$ ). Recently, CHRFA7A was found to contribute to exacerbating inflammation and tissue damage associated with osteoarthritis, and thus being a novel genetic risk factor and therapeutic target for pain (*Courties et al. (2023)*).

Lastly, UNC79 is an auxiliary subunit of the NALCN channel, which carries depolarizing sodium (Na<sup>+</sup>) leak currents to regulate the resting membrane potential of many neurons to modulate pain sensitivity (*Ren (2011)*). UNC79 and UNC80 are HEAT-repeat proteins that docks intracellularly onto the NALCN-FAM155A pore-forming subcomplex and are important for regulating the gating of NALCN. A recent study shows that the NALCN channel contributes to neuronal sensitization in neuropathic pain (*Zhang et al. (2021)*), and this result may lead experimental research to examine the regulatory mechanisms of NALCN by UNC79 and their associations with neuropathic pain in detail.

To encourage researchers from other fields to integrate multiple datasets, we provided a flexible, reproducible, and easy-to-understand code template (<https://github.com/aliibarry/omics-database>), as well as a tutorial, paving the way for better data utilisation and increased impact of individual -omics studies in biomarker discovery.

## Limitations

This research faces certain limitations, with one significant constraint being the limited number of pain-related genes that have been labelled, resulting in a class imbalance (P/NP = 1/40). Additionally, the relative difficulty in validating candidates poses a challenge in obtaining additional labels. To address these limitations, we employed multiple strategies: 1) utilising MCC and GM as the training metrics for model optimization, as they are more resilient to imbalance, and 2) incorporating ensemble classifiers into our approach.

## Future directions

As additional pain-related genes are discovered, it becomes possible to categorize different pain types, such as neuropathic, inflammatory, or cancer pain, which will ultimately refine gene prediction models specific to each subtype, enhancing accuracy. Building on this, the ability to probe tissue-specific mechanisms, opposed to a broad peripheral nervous system focus, will further enhance our knowledge. As new -omics datatypes continue to evolve and improve, such as epigenomics and human genetic studies, our ability to predict relevant candidates will also advance.

Furthermore, future research can focus on employing the existing machine learning framework to analyze gene expression patterns across a range of pathological conditions, such as neurodegenerative diseases, psychiatric disorders, and cancer, thereby broadening clinical impact and therapeutic understanding.

## Conclusions

This study uses large-scale multi-omics, PPI network, and genomic data to predict potential pain-related genes. Based on predicted pain scores, a number of hub proteins were selected as promising studies for future studies. A shiny app, accompanied by code template, is developed for further exploration of pain genes in the context of their PPI networks. Together, the findings and methodology presented in this study not only shed light on future directions in pain research, but also offer a valuable framework that can be adapted and applied to other fields for biomarker discovery.

## Methods and Materials

### Classification labels

The purpose of this study was two-fold: 1) explore which factors help predict if a gene is involved in pain, and 2) identify novel candidates for follow-up studies. To these effect, we used a rigorous labelling system to denote Pain (P) and Non-Pain (NP) genes, encompassing a variety of pain conditions. Here, we required high confidence levels in the true "pain" designation and assume that a subset of pain genes to be discovered within the list of "non-pain" genes.

Together, 429 "pain genes" from gold-standard, experimentally validated lists were used, including those from the Pain Genes Database (PGD) (*LaCroix-Fralish et al. (2007)*) and from the DOLORisk Priority Group (*Themistocleous et al. (2023)*) to label genes as Pain (P), the remaining known genes were labelled as Non-Pain (NP) genes. The PGD contains pain-related phenotypes (both acute and injury-induced) of transgenic mice, whereas the DOLORisk Priority Group includes genes shown

to have a causal (Tier 1) role in human neuropathic pain based on casual variants in multiple, independent families. Tier 2 genes have been implicated in human neuropathic pain, but do not meet the criteria set out in Tier 1 and were also included as "pain genes". Tier 3 genes, which as described as "of interest, as determined by expert consensus" but have little or no published data in the context of human neuropathic pain were omitted from the "pain" label (ie. labelled "non-pain") but were used as a comparator for downstream analyses.

Genes from the Human Pain Genetics Database (HPGDB) (*Meloto et al. (2018)*), which contains pain-associated SNPs from human GWAS studies were labelled as NP: While SNPs against painful conditions are considered relevant to pain, there is a lack of experimental validation specifically highlighting a functional role. Future work may implicate these more strongly, but in this study they are labelled "NP" due to the lack of experimental validation unless there was data from the Pain Genes Database or DOLORisk priority genes suggesting that they were pain genes. These were instead used for external validation of the classification "pain score" in the form of an GSEA enrichment analysis (described below).

## Data input and pre-processing

Fifty-eight input features were selected from three categories: genomic features, experimental -omics datasets, and network topological coefficients, which were mapped from the rodent to human genome using biomaRt where necessary (Table 2 and 3).

### Genomic features

The genomic features used in the model include the cellular compartment with the highest gene expression, the GC content (percentage of GC nucleotides) in the gene sequence, the chromosome name, and the conservation score of the gene indicating its evolutionary conservation. Conservation is calculated as the total conservation score of each base divided by the number of DNA bases. These features are retrieved using the biomaRt package (*Durinck et al. (2005)*) as well as the STRING DB plugin through Cytoscape, as discussed below. Both cellular compartment and tissue expression were extracted through the Cytoscape plug-in, representing the compartment and tissue in which the protein is the most highly expressed. These, along with the chromosome names were then vectorized using 'OrdinalEncoder()'.

GO terms of each gene are also retrieved from Ensembl using 'goslim\_goa\_accession', with terms appearing in less than 20% of the genes retained. The GO "slim" dataset used provides an overview of biological function, retaining high order classifications without the specific terms such as "response to pain" which may affect classification through data leakage. Here, we verified that no terms containing the word "pain" were included during feature space generation. These terms were vectorised via term frequency - inverse document frequency (TF-IDF), which reduces the weight of frequent terms and increases the weight of rare ones. Dimensionality was reduced by PCA, where the data is projected in the subspace of a few principal components that explains most of the observed variance. With this, they can be visualised to assess the clustering of P vs NP genes.

### Omics datasets

A curated selection of high quality -omics datasets were used to reflect diversity in species, relevant tissue, and high-throughput method. Together, this includes a mix of transcriptomic, translatomic, and proteomic datasets from mouse, rat, and human. Various tissues were included across the somatosensory pathway, both in the context of naïve expression as well as fold changes in injured states. This includes skin, nerve, dorsal root ganglia (DRG), and spinal cord, with a slight bias towards mouse DRG due to the sheer number of studies across pain conditions by the field. Where possible, human expression data and differential gene expression was included. A full list of datasets are available in Table 2 and 3.

### Network Topological Features

The STRING DB was used to calculate network protein-protein interaction scores for human protein-coding genes (*Szklarczyk et al. (2021)*). Here, 11 general topological features were calculated using the software Cytoscape plugin NetworkAnalyzer

using default parameters (Table 2 and 3). Features with zero variance (eg. PartnerOfMultiEdgedNodePairs and IsSingleNode) were not included. More detailed explanations and mathematical formulae can be found in the online help document of NetworkAnalyzer (<https://med.bioinf.mpi-inf.mpg.de/netanalyzer/help/2.7/>).

## Data pre-processing and Feature Engineering

Data processing was performed in Python, using Sklearn (*Pedregosa et al. (2011)*). Numerical features, including network topology, conservation scores, and -omics data were scaled using *MinMaxScaler(feature\_range = (-1, 1))*. Categorical features were converted as described above prior to scaling with using *MinMaxScaler(feature\_range = (-1, 1))*.

For pairs of highly correlative features (correlation coefficients > 0.75), the feature with smaller variance was removed. Including highly correlated features can lead to overfitting and decrease the model's interpretability.

The data was split into training (70%) and validation (30%) sets, stratified by label. Numerical features were centred and normalised separately for the training and validation set using min-max normalisation from the Sklearn library. This step ensures unbiased feature comparisons and improves the stability and convergence of ML algorithms.

To improve the model's effectiveness by reducing complexity, we calculated composite LFC values for bulk transcriptomics for every tissue and species. This involved categorizing datasets according to tissue and species. Within each group, we determined the mean LFC value for every gene, after setting non-significant entries (FDR > 0.05) to zero.

## Evaluation Metrics

This dataset is featured by a severe class imbalance due to our stringent labelling of true "pain" genes. To tackle the class imbalance, four metrics (Matthews Correlation Coefficient (MCC, 1), geometric mean (GM, 2), F1 score (F1, 3), and balanced accuracy (BA, 4)) were chosen to be maximised during model selection and for benchmarking the best performing models during validation. These metrics have been shown to be robust in class imbalances. Equations for each metric are provided, with TP = true positive, TN = true negative, FP = false positive, FN = false negative.

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP) \times (TP + FN) \times (TN + FP) \times (TN + FN)}} \quad (1)$$

$$GM = \sqrt{\text{Sensitivity} \times \text{Specificity}} = \sqrt{\frac{TP}{TP + FN} \times \frac{TN}{TN + FP}} \quad (2)$$

$$F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = 2 \times \frac{\frac{TP}{TP + FP} \times \frac{TP}{TP + FN}}{\frac{TP}{TP + FP} + \frac{TP}{TP + FN}} \quad (3)$$

$$BA = \frac{1}{2} \left( \frac{TP}{TP + FN} + \frac{TN}{TN + FP} \right) \quad (4)$$

## Hyperparameter Tuning & Feature Selection

Initial feature selection was performed using the XGBoost Classifier after hyperparameter tuning by Optuna framework (*Akiba et al. (2019)*) and the shap library (*Lundberg and Lee (2017)*). The optimal number of features was determined through backward elimination, removing one feature at a time. The combinations of features with the greatest GM and the greatest MCC score were then used to train models. The final 23 features are highlighted in Figure 3A.

Features were ranked for importance using their SHAP (SHapley Additive exPlanations) values. These values are derived from cooperative game theory and represent the importance of each feature to a given model. They can be approximated using 'shap.TreeExplainer()' in Python by extracting the 'shap\_value' from the test data.

## Model Training

Six machine learning models, including RandomForest Classifier, AdaBoostClassifier, GradientBoostingClassifier, XGBoost Classifier, Stacking Classifier, and Voting Classifier were used for classifying pain/non-pain genes (*Ho (1995); Schapire (2013); Friedman (2001); Chen and Guestrin (2016)*). After using the backward elimination method, optimizing for GM and MCC, 23 features were used (Figure 3A). The Stacking Classifier (final\_estimator = 'logistic regression') and Voting Classifier(voting = 'soft') consists in stacking the output of RandomForest Classifier, AdaBoostClassifier, GradientBoostingClassifier, and XGBoost Classifier.

Hyperparameter tuning was conducted to optimise the summed GM, MCC, BA, and F1 scores of models using the Optuna framework *Akiba et al. (2019)*. The weights of the four individual classifiers in the voting classifier were tuned using gridsearch. During the training step, a 10-fold cross-validation was utilised to assess the models' performance. The full list of parameters are available in a Jupyter Notebook at <https://github.com/alibarry/omics-classifier/>.

The voting classifier is a weighted ensemble of the base classifiers. A permutation function generates all possible orderings of the list, and the weight is selected by calculating and selecting the highest GM score of the classifier with all permutations of weights. Here, the weights [4,2,1,3] correspond to the estimators [xgbm, gb, ada, rf]. We set the voting to 'soft', by which the classifier predicts the class label based on the argmax of the sums of the predicted probabilities by each individual classifier. The class probabilities predicted by each classifier are multiplied by the weight before averaging (soft voting).

## Model Validation

The gene ranking is derived based on the class probability of the best performing classifier. Internal validation of gene ranking is achieved by GSEA enrichment against the Pain Genes Database, which are derived from results of pain-relevant knockout studies, and the DOLORisk pain Genes. External validation is achieved by GSEA enrichment against Human Pain Genetics Database (HPGDB), which contains pain-associated genes from human GWAS studies (*Meloto et al. (2018)*). This was run with and without overlapping "pain" labelled genes to prevent any leakage from the original data, as some SNPs from the HPGDB have been validated functionally and are thus contained in our "pain gene" list. Data were also compared to drugs relevant to "pain", "neuropathic pain", and "chronic pain", extracted from opentargets.org (*Ochoa et al. (2023)*). Here, only approved drugs were used for GSEA analyses.

## Functional Analysis

Gene Ontology (GO) analysis was conducted on the top and bottom 10% ranked genes respectively using the goseq package in R to identify enriched GO terms associated with top ranking pain genes, using the complete set of GO terms (*Young et al. (2010)*). A significance threshold of adjusted p-value (Benjamini-Hochberg corrected) < 0.05 was applied to determine significantly enriched GO terms. KEGG analysis was done on the top 10% ranked genes using the pathfindR package (*Ulgen et al. (2019)*). Furthermore, the top 10% genes were subjected to GSEA against the HSPGB, pain-related GO terms (GO: 0051930, GO: 0071805, GO:0007409, GO:0048265, GO:0007186) and pain-unrelated GO terms (GO:0006357 and GO:0006355) using the clusterProfiler R package (*Wu et al. (2021)*).



**Table 2.** Genomic and Network Features

Engineered Features	Description
<i>Genomic Features</i>	<i>Durinck et al. (2009, 2005)</i>
Cellular compartment	The cellular compartment by which the gene has the highest expression
GC	% GC content
Chromosome name	Name of chromosome
Conservation score	The conservation score of the gene; calculated as the total conservation score of each base divided by the number of DNA bases.
GO1 & GO2	PCA components after vectorisation and PCA transformation of Gene Ontology terms of each gene; The GO terms are filtered such that only terms that appear in less than 20% of the genes are retained.
Tissue	The tissue with the highest expression
<i>Network Topological features</i>	<i>Gustavsen et al. (2019); Shannon et al. (2003)</i>
Average Shortest Path Length	Expected distance between two connected nodes
Betweenness Centrality	Control that this node exerts over the interactions of other nodes in the network
Closeness Centrality	How fast information spreads from a given node to other reachable nodes in the network
Clustering Coefficient	The number of triangles (3-loops) that pass through this node, relative to the maximum number of 3-loops that could pass through the node
Degree	The number of edges linked to a node
Eccentricity	The maximum non-infinite length of a shortest path between a node and another in the network
Neighborhood Connectivity	The connectivity of a node is the number of its neighbours. The neighbourhood connectivity of a node n is defined as the average connectivity of all neighbours of n
Undirected Edges	The number of undirected edges that are connected to a node
Radiality	A node centrality index
Stress	Counts the number of shortest paths passing through a node
Topological Coefficient	Relative measure for the extent to which a protein in the network shares interaction partners with other proteins

**Table 3.** -omics features

Engineered Features	Tissue	Species	Original Features	Description	Reference
Transcriptome DEG, mDRG	DRG	Mouse	LFC_CGRT_3D	LFC, subpopulations in mice after nerve injury	<i>Barry et al. (2023)</i>
Transcriptome DEG, mDRG	DRG	Mouse	LFC_CGRT_4W	LFC, subpopulations in mice after nerve injury	<i>Barry et al. (2023)</i>
Transcriptome DEG, mDRG	DRG	Mouse	LFC_CRTH_3D	LFC, subpopulations in mice after nerve injury	<i>Barry et al. (2023)</i>
Transcriptome DEG, mDRG	DRG	Mouse	LFC_CRTH_4W	LFC, subpopulations in mice after nerve injury	<i>Barry et al. (2023)</i>
Transcriptome DEG, mDRG	DRG	Mouse	LFC_MRTD_3D	LFC, subpopulations in mice after nerve injury	<i>Barry et al. (2023)</i>
Transcriptome DEG, mDRG	DRG	Mouse	LFC_MRTD_4W	LFC, subpopulations in mice after nerve injury	<i>Barry et al. (2023)</i>
Transcriptome DEG, mDRG	DRG	Mouse	LFC_TBAC_3D	LFC, subpopulations in mice after nerve injury	<i>Barry et al. (2023)</i>
Transcriptome DEG, mDRG	DRG	Mouse	LFC_TBAC_4W	LFC, subpopulations in mice after nerve injury	<i>Barry et al. (2023)</i>
Transcriptome DEG, mDRG	DRG	Mouse	LFC_TDNV_3D	LFC, subpopulations in mice after nerve injury	<i>Barry et al. (2023)</i>
Transcriptome DEG, mDRG	DRG	Mouse	LFC_TDNV_4W	LFC, subpopulations in mice after nerve injury	<i>Barry et al. (2023)</i>
Transcriptome DEG, mDRG	DRG	Mouse	LFC_balb	LFC values of mouse DRG gene expression after SNI	<i>Baskozos et al. (2019)</i>
Transcriptome DEG, mDRG	DRG	Mouse	LFC_b10d2	LFC values of mouse DRG gene expression after SNI	<i>Baskozos et al. (2019)</i>
Transcriptome expression, mDRG subpopulations	DRG	Mouse	subtype_tpm	TPM counts of genes from transgenically-labelled subpopulations of neurons in male and female mice	<i>Barry et al. (2023)</i>
Transcriptome expression, hSkin (healthy)	Skin	Human	hSkin_DB_tpm	Naive gene expression counts in human skin	<i>Baskozos et al. (2022)</i>
Transcriptome DEG, hSkin from painful vs painless DB patients	Skin	Human	LFC_Diabetes	LFC values of human DPN patients vs control	<i>Baskozos et al. (2022)</i>
Transcriptome DEG, hSkin from painful vs painless DB patients	Skin	Human	LFC_Diabetes_male	LFC values of human DPN patients vs control, males	<i>Baskozos et al. (2022)</i>
Transcriptome DEG, hSkin from painful vs painless DB patients	Skin	Human	LFC_Diabetes_female	LFC values of human DPN patients vs control, females	<i>Baskozos et al. (2022)</i>
Transcriptome expression, hSkin (healthy)	Skin	Human	hSkin_CTS_tpm	Naive gene expression counts in human skin	<i>Baskozos et al. (2020)</i>
Transcriptome DEG, hSkin from carpal tunnel patients, pre- and post- surgery	Skin	Human	LFC_skin	LFC of the gene expression in skin	<i>Baskozos et al. (2020)</i>
Transcriptome expression, iPSC	iPSC	Human	iPSC_tpm	Gene expression counts of iPSC cells	<i>Clark et al. (2021)</i>

Transcriptome DEG, iPSC	iPSC	Human	LFC_iPSC_young	LFC values of iPSC cells following nerve injury, young	<i>Clark et al. (2021)</i>
Transcriptome DEG, iPSC	iPSC	Human	LFC_iPSC_old	LFC values of iPSC cells following nerve injury, old	<i>Clark et al. (2021)</i>
Transcriptome DEG, iPSC	iPSC	Human	LFC_iPSC	LFC values of iPSC cells following nerve injury	<i>Clark et al. (2021)</i>
Transcriptome expression, mDRG	DRG	Mouse	mDRG_sham_tpm	TPM, sham nerve injury DRG	<i>Baskozos et al. (2019)</i>
Transcriptome DEG, rat	DRG	Rat	SNI_rat_LFC	LFC, rat DRG model of SNI	<i>Maratou et al. (2008)</i>
Transcriptome DEG, rat	DRG	Rat	SNT_rat_LFC	LFC, rat DRG model of SNT	<i>Vega-Avelaira et al. (2009)</i>
Transcriptome DEG, rat	DRG	Rat	LFC_hiv	LFC, rat DRG model of HIV-associated neuropathic pain	<i>Maratou et al. (2008)</i>
Transcriptome DEG, rat	DRG	Rat	LFC_bone_cancer	LFC, rat DRG model of bone cancer	<i>Perkins et al. (2013)</i>
Transcriptome DEG, glial injury	DRG SG	Mouse	LFC_d3_glial	LFC, transcriptional fingerprint of satellite glial cells following peripheral nerve injury, day 3	<i>Jager et al. (2020)</i>
Transcriptome DEG, glial injury	DRG SG	Mouse	LFC_d8_glial	LFC, transcriptional fingerprint of satellite glial cells following peripheral nerve injury, day 8	<i>Jager et al. (2020)</i>
Transcriptome DEG, glial injury	DRG satellite glia	Mouse	LFC_d14_glial	LFC, transcriptional fingerprint of satellite glial cells following peripheral nerve injury, day 14	<i>Jager et al. (2020)</i>
Translatome expression	DRG nociceptors	Mouse	trans_tpm	TPM, nociceptor translatomes data of in male and female mice following nerve injury	<i>Tavares-Ferreira et al. (2022)</i>
Translatome expression, normalized	DRG nociceptors	Mouse	transnorm_tpm	TPM, nociceptor translatomes data of in male and female mice following nerve injury, normalised (as published)	<i>Tavares-Ferreira et al. (2022)</i>
Proteome DEG (FDR), mSC	SC	Mouse	sc_lfc	LFC, protein expression in spinal cord after SNI in male mice	<i>Barry et al. (2018)</i>
Proteome DEG (LFC), mDRG	DRG	Mouse	drg_lfc	LFC, protein expression in DRG after SNI in male mice	<i>Barry et al. (2018)</i>
Proteome DEG (LFC), mSN	SN	Mouse	sn_lfc	LFC, protein expression in the sciatic nerve (SN) after SNI in male mice	<i>Barry et al. (2018)</i>
Proteome DEG (FDR), mSC	SC	Mouse	sc_padj	FDR, protein expression changes in the spinal cord (SC) after SNI in male mice	<i>Barry et al. (2018)</i>
Proteome DEG (FDR), mDRG	DRG	Mouse	drg_padj	FDR, protein expression changes in the dorsal root ganglia (DRG) after SNI in male mice	<i>Barry et al. (2018)</i>

Proteome DEG (FDR), mSN	SN	Mouse	sn_padj	FDR, protein expression changes in the sciatic nerve (SN) after SNI in male mice	<b><i>Barry et al. (2018)</i></b>
Transcriptome DEG, hiPSC-SN vs hiPSC	hiPSC-SN	Human	LFC_iPSC	LFC, gene expression in iPSC cells	<b><i>McDermott et al. (2019)</i></b>
Transcriptome expression, hDRG	DRG	Human	hdrg_tpm	TPM, naïve human DRG	<b><i>Ray et al. (2022)</i></b>
Transcriptome expression, mDRG subpopulation	DRG	Mouse	TPM_zheng	TPM, naïve mouse nociceptor	<b><i>Zheng et al. (2019)</i></b>

---

## Data and code availability

No new datasets were generated from this study. Predicted pain scores and searchable datasets are searchable at <https://livedataoxford.shinyapps.io/drg-directory/>. A Jupyter Notebook for building the classifiers is available at <https://github.com/aliibarry/omics-classifier>. Database code, with a simplified example are available at <https://github.com/aliibarry/omics-database>. A manual for database development is available at <https://aliibarry.github.io/database-book/>.

## Author contributions

AMB, NZ, and GB designed the study with input from DB. NZ trained classifiers with input from AMB and GB. AMB and NZ developed the database and associated code/manual to improve data access with input from GB and DB. AMB and NZ drafted the manuscript and figures that were reviewed and approved by all authors.

## Acknowledgments

This work was funded in part by a Wellcome Investigator Grant to DB (223149/Z/21/Z), as well as the MRC (MR/T020113/1), and with funding from the MRC and Versus Arthritis to the PAINSTORM consortium as part of the Advanced Pain Discovery Platform (MR/W002388/1). GB is funded by Diabetes UK, grant number 19/0005984, MRC and Versus Arthritis through the PAINSTORM consortium as part of the Advanced Pain Discovery Platform (MR/W002388/1) and by the Wellcome Trust (223149/Z/21/Z). AMB is funded by the MRC and Versus Arthritis through the PAINSTORM consortium as part of the Advanced Pain Discovery Platform (MR/W002388/1). ChatGTP was used to help format equations in LaTeX for submission.

This research was funded in part by the Wellcome Trust [223149/Z/21/Z]. For the purpose of open access, the author has applied a CC BY public copyright license to any Author Accepted Manuscript version arising from this submission.

# References

- Akiba T**, Sano S, Yanase T, Ohta T, Koyama M. Optuna: A Next-generation Hyperparameter Optimization Framework. In: *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*; 2019. doi: 10.1145/3292500.3330701.
- Barry AM**, Sondermann JR, Sondermann JH, Gomez-Varela D, Schmidt M. Region-resolved quantitative proteome profiling reveals molecular dynamics associated with chronic pain in the PNS and spinal cord. *Frontiers in Molecular Neuroscience*. 2018; 11:259. doi: 10.3389/fn-mol.2018.00259.
- Barry AM**, Zhao N, Yang X, Bennett DL, Baskozos G. Deep RNA-seq of male and female murine sensory neuron subtypes after nerve injury. . 2023; <https://doi.org/10.1101/2022.11.21.516781>, doi: 10.1101/2022.11.21.516781.
- Baskozos G**, Dawes JM, Austin JS, Antunes-Martins A, Mcdermott L, Clark AJ, Trendafilova T, Lees JG, McMahon SB, Mogil JS, Orenco C, Bennett DL. Comprehensive analysis of long noncoding RNA expression in dorsal root ganglion reveals cell-type specificity and dysregulation after nerve injury. *Pain*. 2019 2; 160:463–485. <https://pubmed.ncbi.nlm.nih.gov/30335683/>, doi: 10.1097/j.pain.0000000000001416.
- Baskozos G**, Sandy-Hindmarch O, Clark AJ, Windsor K, Karlsson P, Weir GA, McDermott LA, Burchall J, Wiberg A, Furniss D, Bennett DLH, Schmid AB. Molecular and cellular correlates of human nerve regeneration: ADCYAP1/PACAP enhance nerve outgrowth. *Brain*. 2020 7; 143:2009. /pmc/articles/PMC7462094//pmc/articles/PMC7462094/?report=abstract<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7462094/>, doi: 10.1093/brain/awaa163.
- Baskozos G**, Themistocleous AC, Hebert HL, Pascal MMV, John J, Callaghan BC, Laycock H, Granovsky Y, Crombez G, Yarnitsky D, Rice ASC, Smith BH, Bennett DLH. Classification of painful or painless diabetic peripheral neuropathy and identification of the most powerful predictors using machine learning models in large cross-sectional cohorts. *BMC Medical Informatics and Decision Making*. 2022 12; 22. doi: 10.1186/s12911-022-01890-x.
- Boeckhout M**, Zielhuis GA, Bredenoord AL, The FAIR guiding principles for data stewardship: Fair enough?; 2018. doi: 10.1038/s41431-018-0160-0.
- Bonnington JK**, McNaughton PA. Signalling pathways involved in the sensitisation of mouse nociceptive neurones by nerve growth factor. *Journal of Physiology*. 2003; 551. doi: 10.1113/jphysiol.2003.039990.
- Brodersen KH**, Ong CS, Stephan KE, Buhmann JM. The balanced accuracy and its posterior distribution. In: *Proceedings - International Conference on Pattern Recognition*; 2010. doi: 10.1109/ICPR.2010.764.
- Cekic C**, Linden J, Purinergic regulation of the immune system; 2016. doi: 10.1038/nri.2016.4.
- Chen T**, Guestrin C. XGBoost: A scalable tree boosting system. In: *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, vol. 13-17-August-2016; 2016. doi: 10.1145/2939672.2939785.
- Chicco D**, Jurman G. The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC Genomics*. 2020; 21. doi: 10.1186/s12864-019-6413-7.
- Choi D**, Goodwin G, Stevens EB, Soliman N, Namer B, Denk F. Spontaneous activity in peripheral sensory nerves: a systematic review. *Pain*. 2023; doi: 10.1097/j.pain.0000000000003115.
- Clark AJ**, Kugathasan U, Baskozos G, Priestman DA, Fugger N, Lone MA, Othman A, Chu KH, Blesneac I, Wilson ER, Laurà M, Kalmar B, Green-smith L, Hornemann T, Platt FM, Reilly MM, Bennett DL. An iPSC model of hereditary sensory neuropathy-1 reveals L-serine-responsive deficits in neuronal ganglioside composition and axoglial interactions. *Cell Reports Medicine*. 2021 7; 2. /pmc/articles/PMC8324498//pmc/articles/PMC8324498/?report=abstract<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8324498/>, doi: 10.1016/j.xcrm.2021.100345.
- Cooper AH**, Barry AM, Chrysostomidou P, Lollignier R, Titterton HF, Bennett DL, Weir GA, Weir G. Peripheral nerve injury results in a biased loss of sensory neuron sub-populations. *Pain*. 2024; <https://doi.org/10.1101/2023.11.14.566863>, doi: 10.1101/2023.11.14.566863.
- Courties A**, Olmer M, Myers K, Ordoukhanian P, Head SR, Natarajan P, Berenbaum F, Sellam J, Lotz MK. Human-specific duplicate CHRFAM7A gene is associated with more severe osteoarthritis and amplifies pain behaviours. *Annals of the Rheumatic Diseases*. 2023; 82. doi: 10.1136/ard-2022-223470.
- Davies AJ**, Kim HW, Gonzalez-Cano R, Ugolini S, Costigan M, Bae S, Correspondence O, Choi J, Back SK, Roh SE, Johnson E, Gabriac M, Kim MS, Lee J, Lee JE, Kim YS, Bae YC, Kim SJ, Lee KM, Na HS, et al. Natural Killer Cells Degenerate Intact Sensory Afferents following Nerve Injury Article Natural Killer Cells Degenerate Intact Sensory Afferents following Nerve Injury. *Cell*. 2019; 176:716–728. <https://doi.org/10.1016/j.cell.2018.12.022>, doi: 10.1016/j.cell.2018.12.022.



460 **Davies AJ**, Rinaldi S, Costigan M, Oh SB, Cytotoxic Immunity in Peripheral Nerve Injury and Pain; 2020. doi: 10.3389/fnins.2020.00142.

461 **Durinck S**, Moreau Y, Kasprzyk A, Davis S, Moor BD, Brazma A, Huber W. BioMart and Bioconductor: A powerful link between biological  
462 databases and microarray data analysis. *Bioinformatics*. 2005; 21. doi: 10.1093/bioinformatics/bti525.

463 **Durinck S**, Spellman PT, Birney E, Huber W. Mapping identifiers for the integration of genomic datasets with the R/ Bioconductor package  
464 biomaRt. *Nature Protocols*. 2009; 4. doi: 10.1038/nprot.2009.97.

465 **Edvinsson L**, Haanes KA, Warfvinge K, Krause DN, CGRP as the target of new migraine therapies - Successful translation from bench to  
466 clinic; 2018. doi: 10.1038/s41582-018-0003-1.

467 **Espindola RP**, Ebecken NFF. On extending F-measure and G-mean metrics to multi-class problems. In: *Data Mining VI*, vol. 1; 2005. doi:  
468 10.2495/data050031.

469 **Flögel U**, Burghoff S, Lent PLEMV, Temme S, Galbarz L, Ding Z, El-Tayeb A, Huels S, Bönner F, Borg N, Jacoby C, Müller CE, Berg WBVD,  
470 Schrader J. Selective activation of adenosine A2A receptors on immune cells by a CD73-dependent prodrug suppresses joint inflammation  
471 in experimental rheumatoid arthritis. *Science Translational Medicine*. 2012; 4. doi: 10.1126/scitranslmed.3003717.

472 **Fricker FR**, Lago N, Balarajah S, Tsantoulas C, Tanna S, Zhu N, Fageiry SK, Jenkins M, Garratt AN, Birchmeier C, Bennett DLH. Development/-  
473 Plasticity/Repair Axonally Derived Neuregulin-1 Is Required for Remyelination and Regeneration after Nerve Injury in Adulthood. . 2011;  
474 doi: 10.1523/JNEUROSCI.2568-10.2011.

475 **Fricker FR**, Zhu N, Tsantoulas C, Abrahamsen B, Nassar MA, Thakur M, Garratt AN, Birchmeier C, McMahon SB, Wood JN, Bennett DLH. De-  
476 velopment/Plasticity/Repair Sensory Axon-Derived Neuregulin-1 Is Required for Axoglial Signaling and Normal Sensory Function But Not  
477 for Long-Term Axon Maintenance. *J Neurosci*. 2009; 29:7667–7678. www.ncbi.nlm.nih.gov/blast/Blast.cgi, doi: 10.1523/JNEUROSCI.6053-  
478 08.2009.

479 **Friedman JH**. Greedy function approximation: A gradient boosting machine. *Annals of Statistics*. 2001; 29. doi: 10.1214/aos/1013203451.

480 **Goadsby PJ**, Reuter U, Hallström Y, Broessner G, Bonner JH, Zhang F, Sapra S, Picard H, Mikol DD, Lenz RA. A Controlled Trial of Erenumab  
481 for Episodic Migraine. *New England Journal of Medicine*. 2017; 377. doi: 10.1056/nejmoa1705848.

482 **Gustavsen JA**, Pai S, Isserlin R, Demchak B, Pico AR. Rcy3: Network biology using cytoscape from within r [version 1; peer review: 2 approved].  
483 F1000Research. 2019; 8. doi: 10.12688/f1000research.20887.1.

484 **Higurashi M**, Ishida T, Kinoshita K. Identification of transient hub proteins and the possible structural basis for their multiple interactions.  
485 *Protein Science*. 2008; 17. doi: 10.1110/ps.073196308.

486 **Ho TK**. Random decision forests. In: *Proceedings of the International Conference on Document Analysis and Recognition, ICDAR*, vol. 1; 1995.  
487 doi: 10.1109/ICDAR.1995.598994.

488 **Hu X**, Adebiyi MG, Luo J, Sun K, Le TTT, Zhang Y, Wu H, Zhao S, Karmouty-Quintana H, Liu H, Huang A, Wen YE, Zaika OL, Mamenko M,  
489 Pochynyuk OM, Kellems RE, Eltzschig HK, Blackburn MR, Walters ET, Huang D, et al. Sustained Elevated Adenosine via ADORA2B Promotes  
490 Chronic Pain through Neuro-immune Interaction. *Cell Reports*. 2016; 16. doi: 10.1016/j.celrep.2016.05.080.

491 **Jager SE**, Pallesen LT, Richner M, Harley P, Hore Z, McMahon S, Denk F, Vægter CB. Changes in the transcriptional fingerprint of satellite  
492 glial cells following peripheral nerve injury. *Glia*. 2020 7; 68:1375–1395. <https://onlinelibrary.wiley.com/doi/full/10.1002/glia.23785><https://onlinelibrary.wiley.com/doi/abs/10.1002/glia.23785><https://onlinelibrary.wiley.com/doi/10.1002/glia.23785>, doi: 10.1002/GLIA.23785.

494 **Ji RR**, Chamesian A, Zhang YQ, Pain regulation by non-neuronal cells and inflammation. *American Association for the Advancement of*  
495 *Science*; 2016. doi: 10.1126/science.aaf8924.

496 **Ji RR**, IV RWG, Malcangio M, Strichartz GR, MAP kinase and pain; 2009. doi: 10.1016/j.brainresrev.2008.12.011.

497 **Kim HW**, Wang S, Davies AJ, Oh SB, The therapeutic potential of natural killer cells in neuropathic pain; 2023. doi: 10.1016/j.tins.2023.05.008.

498 **Kondo M**, Shibuta I, Extracellular signal-regulated kinases (ERK) 1 and 2 as a key molecule in pain research; 2020. doi: 10.2334/josnurd.19-  
499 0470.

500 **Kumar Y**, Koul A, Singla R, Ijaz MF. Artificial intelligence in disease diagnosis: a systematic literature review, synthesizing framework and  
501 future research agenda. *Journal of Ambient Intelligence and Humanized Computing*. 2023; 14. doi: 10.1007/s12652-021-03612-z.

502 **LaCroix-Fralish ML**, Ledoux JB, Mogil JS. The Pain Genes Database: An interactive web browser of pain-related transgenic knockout studies.  
503 Pain. 2007; 131. doi: 10.1016/j.pain.2007.04.041.

504 **Lundberg SM**, Lee SI. A unified approach to interpreting model predictions. In: *Advances in Neural Information Processing Systems*, vol.  
505 2017-December; 2017. .

506 **Maratou K**, Wallace VCJ, Hasnie FS, Okuse K, Hosseini R, Jina N, Blackbeard J, Pheby T, Orengo C, Dickenson AH, McMahon SB, Rice ASC.  
507 Comparison of dorsal root ganglion gene expression in rat models of traumatic and HIV-associated neuropathic pain. *European Journal*  
508 *of Pain* (London, England). 2008 7; 13:387–398. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC18606552/?tool=EBIhttps://europepmc.org/article/MED/18606552>, doi: 10.1016/J.EJPAIN.2008.05.011.

510 **McDermott LA**, Weir GA, Themistocleous AC, Segerdahl AR, Blesneac I, Baskozos G, Clark AJ, Millar V, Peck LJ, Ebner D, Tracey I, Serra J,  
511 Bennett DL. Defining the Functional Role of Na V 1.7 in Human Nociception. *Neuron*. 2019; 101:905–919.e8. <https://doi.org/10.1016/j.neuron.2019.01.047>, doi: 10.1016/j.neuron.2019.01.047.

513 **Meloto CB**, Benavides R, Lichtenwalter RN, Wen X, Tugarinov N, Zorina-Lichtenwalter K, Chabot-Doré AJ, Piltonen MH, Cattaneo S, Verma  
514 V, Klares R, Khoury S, Parisien M, Diatchenko L. Human pain genetics database: A resource dedicated to human pain genetics research.  
515 Pain. 2018; 159. doi: 10.1097/j.pain.0000000000001135.

516 **Middleton SJ**, Barry AM, Comini M, Li Y, Ray PR, Shiers S, Themistocleous AC, Uhelski ML, Yang X, Dougherty PM, Price TJ, Bennett DL.  
517 Studying human nociceptors: From fundamentals to clinic. *Brain*. 2021; doi: 10.1093/brain/awab048.

518 **Ochoa D**, Hercules A, Carmona M, Suveges D, Baker J, Malangone C, Lopez I, Miranda A, Cruz-Castillo C, Fumis L, Bernal-Llinares M, Tsukanov  
519 K, Cornu H, Tsigiris K, Razuvayevskaya O, Buniello A, Schwartzentruber J, Karim M, Ariano B, Osorio REM, et al. The next-generation Open  
520 Targets Platform: reimaged, redesigned, rebuilt. *Nucleic Acids Research*. 2023; 51. doi: 10.1093/nar/gkac1046.

521 **Paige C**, Plasencia-Fernandez I, Kume M, Papalampropoulou-Tsiridou M, Lorenzo LE, David ET, He L, Mejia GL, Driskill C, Ferrini F, Feldhaus  
522 AL, Garcia-Martinez LF, Akopian AN, Koninck YD, Dussor G, Price TJ. A Female-Specific Role for Calcitonin Gene-Related Peptide (CGRP) in  
523 Rodent Pain Models. *Journal of Neuroscience*. 2022 3; 42:1930–1944. <https://www.jneurosci.org/content/42/10/1930https://www.jneurosci.org/content/42/10/1930.abstract>, doi: 10.1523/JNEUROSCI.1137-21.2022.

525 **Pedregosa F**, Gramfort NA, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau  
526 D, Pedregosa F, Varoquaux G, Gramfort A, Thirion B, Prettenhofer P, Vanderplas J, Brucher M, an Edouard Duchesnay MP, et al. Scikit-learn:  
527 Machine Learning in {P}ython. *Journal of Machine Learning Research*. 2011; 12.

528 **Perkins JR**, Lees J, Antunes-Martins A, Diboun I, McMahon SB, Bennett DLH, Orengo C. PainNetworks: A web-based resource for  
529 the visualisation of pain-related genes in the context of their network associations. Pain. 2013; 154:2586.e1–2586.e12. doi:  
530 10.1016/j.pain.2013.09.003.

531 **Raja SN**, Carr DB, Cohen M, Finnerup NB, Flor H, Gibson S, Keefe FJ, Mogil JS, Ringkamp M, Sluka KA, Song XJ, Stevens B, Sullivan MD,  
532 Tutelman PR, Ushida T, Vader K, The revised International Association for the Study of Pain definition of pain: concepts, challenges, and  
533 compromises; 2020. doi: 10.1097/j.pain.0000000000001939.

534 **Ray PR**, Shiers S, Caruso JP, Tavares-Ferreira D, Sankaranarayanan I, Uhelski ML, Li Y, North RY, Tatsui C, Dussor G, Burton MD, Dougherty  
535 PM, Price TJ. RNA profiling of human dorsal root ganglia reveals sex differences in mechanisms promoting neuropathic pain. *Brain*. 2022  
536 2; doi: 10.1093/brain/awac266.

537 **Reel PS**, Reel S, Pearson E, Trucco E, Jefferson E, Using machine learning approaches for multi-omics data analysis: A review; 2021. doi:  
538 10.1016/j.biotechadv.2021.107739.

539 **Ren D**, Sodium leak channels in neuronal excitability and rhythmic behaviors; 2011. doi: 10.1016/j.neuron.2011.12.007.

540 **Schapire RE**. In: Explaining adaboost; 2013. doi: 10.1007/978-3-642-41136-6\_5.

541 **Schou WS**, Ashina S, Amin FM, Goadsby PJ, Ashina M, Calcitonin gene-related peptide and pain: a systematic review; 2017. doi:  
542 10.1186/s10194-017-0741-2.

543 **Shannon P**, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B, Ideker T. Cytoscape: A software Environment for  
544 integrated models of biomolecular interaction networks. *Genome Research*. 2003; 13. doi: 10.1101/gr.1239303.

545 **Singhmar P**, Huo X, Eijkelkamp N, Berciano SR, Baameur F, Mei FC, Zhu Y, Cheng X, Hawke D, Mayor F, Murga C, Heijnen CJ, Kavelaars A.  
546 Critical role for Epac1 in inflammatory pain controlled by GRK2-mediated phosphorylation of Epac1. *Proceedings of the National Academy*  
547 *of Sciences of the United States of America*. 2016; 113. doi: 10.1073/pnas.1516036113.

548 **Szklarczyk D**, Gable AL, Nastou KC, Lyon D, Kirsch R, Pyysalo S, Doncheva NT, Legeay M, Fang T, Bork P, Jensen LJ, von Mering C. The STRING  
549 database in 2021: customizable protein–protein networks, and functional characterization of user-uploaded gene/measurement sets.  
550 *Nucleic Acids Research*. 2021 1; 49:D605. /pmc/articles/PMC7779004//pmc/articles/PMC7779004/?report=abstracthttps://www.ncbi.nlm.  
551 nih.gov/pmc/articles/PMC7779004/, doi: 10.1093/NAR/GKAA1074.

552 **Tao F**, Li Q, Liu S, Wu H, Skinner J, Hurtado A, Belegu V, Furmanski O, Yang Y, McDonald JW, Johns RA. Role of neuregulin-1/ErbB signaling  
553 in stem cell therapy for spinal cord injury-induced chronic neuropathic pain. *Stem Cells*. 2013; 31. doi: 10.1002/stem.1258.

554 **Tavares-Ferreira D**, Shiers S, Ray PR, Wangzhou A, Jeevakumar V, Sankaranarayanan I, Cervantes AM, Reese JC, Chamesian A, Copits BA,  
555 Dougherty PM, Gereau RW, Burton MD, Dussor G, Price TJ. Spatial transcriptomics of dorsal root ganglia identifies molecular signatures  
556 of human nociceptors. *Science translational medicine*. 2022 2; 14:eabj8186. http://www.ncbi.nlm.nih.gov/pubmed/35171654, doi: DOI:  
557 10.1126/scitranslmed.abj8186.

558 **Themistocleous AC**, Baskozos G, Blesneac I, Comini M, Megy K, Chong S, Deevi SVV, Ginsberg L, Gosal D, Hadden RDM, Horvath R,  
559 Mahdi-Rogers M, Manzur A, Mapeta R, Marshall A, Matthews E, McCarthy MI, Reilly MM, Renton T, Rice ASC, et al. Investigating  
560 genotype–phenotype relationship of extreme neuropathic pain disorders in a UK national cohort. *Brain Communications*. 2023 3; 5.  
561 https://academic.oup.com/braincomms/article/doi/10.1093/braincomms/fcad037/7048922, doi: 10.1093/braincomms/fcad037.

562 **Ulgen E**, Ozisik O, Sezerman OU. pathfindR: An R Package for Comprehensive Identification of Enriched Pathways in Omics Data Through  
563 Active Subnetworks. *Frontiers in genetics*. 2019 9; 10. https://pubmed.ncbi.nlm.nih.gov/31608109/, doi: 10.3389/FGENE.2019.00858.

564 **Vega-Avelaira D**, Géranton SM, Fitzgerald M. Differential regulation of immune responses and macrophage/neuron interactions in the  
565 dorsal root ganglion in young and adult rats following nerve injury. *Molecular Pain*. 2009 12; 5. doi: 10.1186/1744-8069-5-70.

566 **Wahlman C**, Doyle TM, Little JW, Luongo L, Janes K, Chen Z, Esposito E, Tosh DK, Cuzzocre S, Jacobson KA, Salvemini D. Chemotherapy-  
567 induced pain is promoted by enhanced spinal adenosine kinase levels through astrocyte-dependent mechanisms. *Pain*. 2018; 159. doi:  
568 10.1097/j.pain.0000000000001177.

569 **Wu I**, Wang X. A novel approach to topological network analysis for the identification of metrics and signatures in non-small cell lung cancer.  
570 *Scientific Reports*. 2023; 13. doi: 10.1038/s41598-023-35165-w.

571 **Wu T**, Hu E, Xu S, Chen M, Guo P, Dai Z, Feng T, Zhou L, Tang W, Zhan L, Fu X, Liu S, Bo X, Yu G. clusterProfiler 4.0: A universal enrichment  
572 tool for interpreting omics data. *The Innovation*. 2021 8; 2:100141. doi: 10.1016/J.XINN.2021.100141.

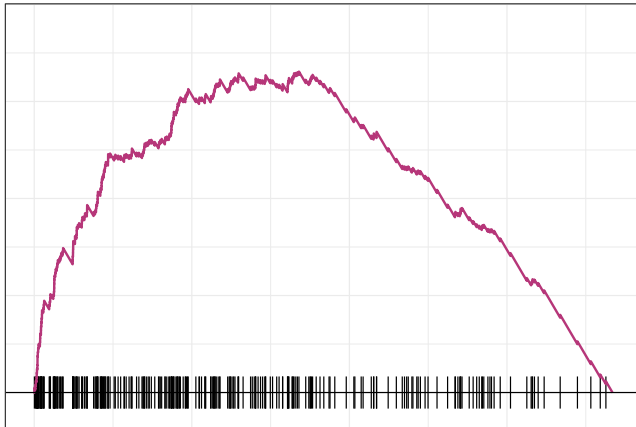
573 **Young MD**, Wakefield MJ, Smyth GK, Oshlack A. Gene ontology analysis for RNA-seq: accounting for selection bias. *Genome Biology*. 2010  
574 2; 11:1–12. https://genomebiology.biomedcentral.com/articles/10.1186/gb-2010-11-2-r14, doi: 10.1186/GB-2010-11-2-R14/TABLES/4.

575 **Yu X**, Lai S, Chen H, Chen M. Protein-protein interaction network with machine learning models and multiomics data reveal potential  
576 neurodegenerative disease-related proteins. *Human Molecular Genetics*. 2020; 29. doi: 10.1093/hmg/ddaa065.

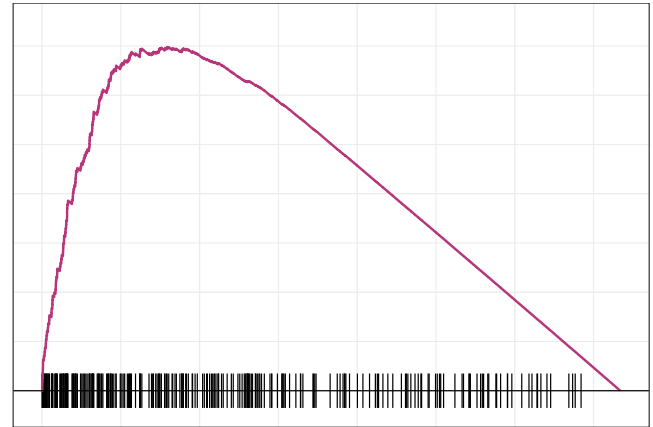
577 **Zhang D**, Zhao W, Liu J, Ou M, Liang P, Li J, Chen Y, Liao D, Bai S, Shen J, Chen X, Huang H, Zhou C. Sodium leak channel contributes to  
578 neuronal sensitization in neuropathic pain. *Progress in Neurobiology*. 2021; 202. doi: 10.1016/j.pneurobio.2021.102041.

579 **Zheng Y**, Liu P, Bai L, Trimmer JS, Bean BP, Ginty DD. Deep Sequencing of Somatosensory Neurons Reveals Molecular Deter-  
580 minants of Intrinsic Physiological Properties. *Neuron*. 2019; 103:598–616.e7. https://doi.org/10.1016/j.neuron.2019.05.039, doi:  
581 10.1016/j.neuron.2019.05.039.

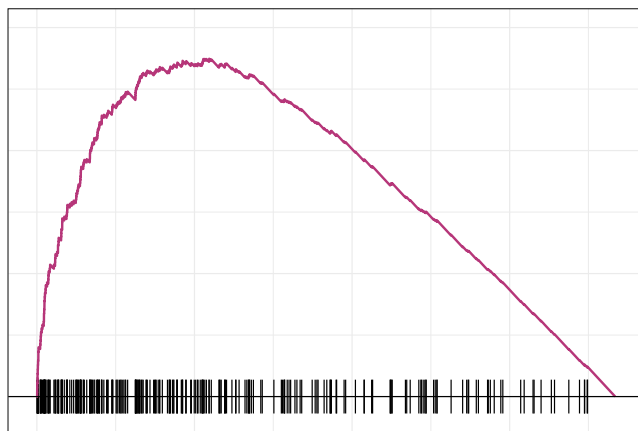
AdaBoost



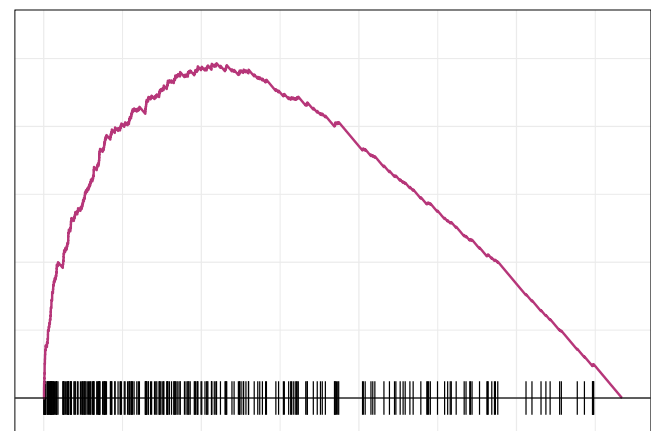
GradientBoost



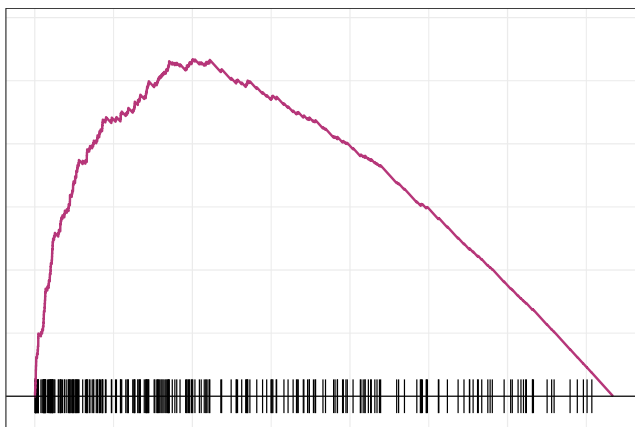
Stacking



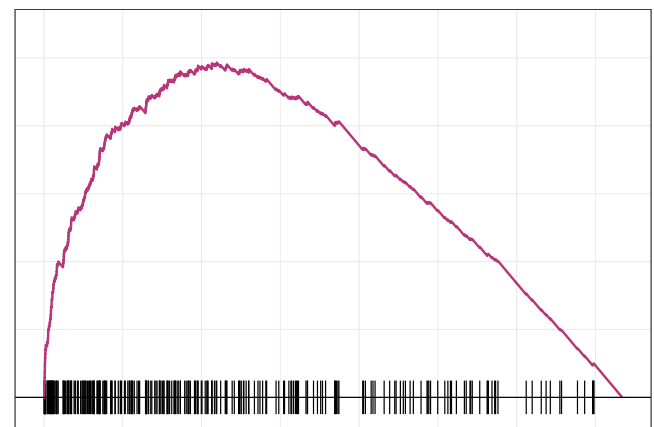
Voting



Random Forest



XGBoost



**Figure supplement 1.** GSEA validation analyses against HPGDB without "pain" labels