



# Development of a Transcriptional Regulation Bioinformatics Pipeline to Predict Co-Regulated Genes in Vascular Smooth Muscle Cell Phenotypic Transitions During Atherosclerosis

Mahima Reddy<sup>1</sup>, Vlad Serbulea<sup>1</sup>, Sohel Shamsuzzaman<sup>1</sup>, Anita Salamon<sup>1</sup>, Rupa Tripathi<sup>1</sup>, Clint Miller<sup>2</sup>, Giuseppe Mocci<sup>3</sup>, Johan Björkegren<sup>4</sup>, Gary Owens<sup>1</sup>

<sup>1</sup>Robert M. Berne Cardiovascular Research Center, Division of Cardiovascular Medicine, University of Virginia School of Medicine; <sup>2</sup>Center for Public Health Genomics, University of Virginia School of Medicine; <sup>3</sup>Department of Medicine Huddinge, Karolinska Institutet; <sup>4</sup>Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai

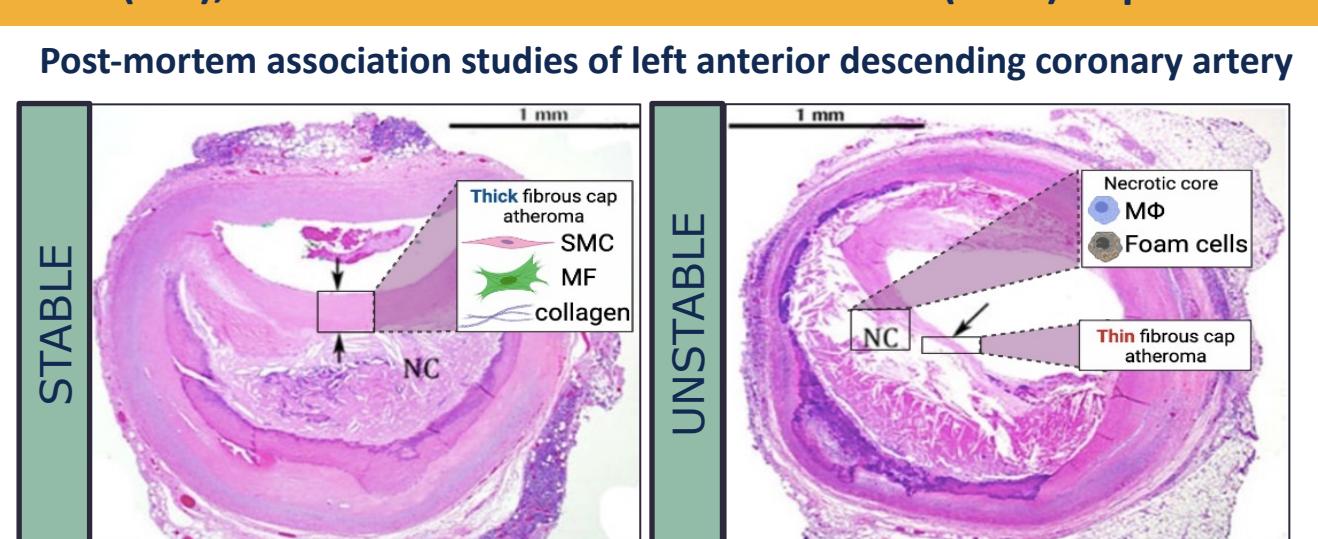


## Abstract

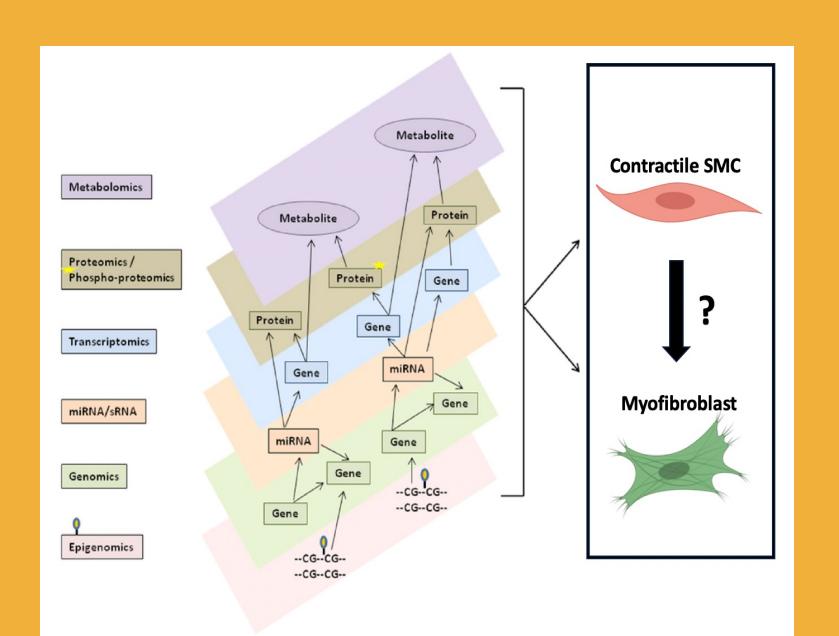
The rupture of vulnerable or advanced atherosclerotic lesions contributes to myocardial infarctions and strokes, the leading causes of death globally. During all stages of atherosclerosis, vascular smooth muscle cells (VSMCs) are recruited as a major source of plaque cells within lesions, where they undergo phenotypic modulation and give rise to both plaque-stabilizing and destabilizing phenotypes. Because most atherosclerosis studies test the functional role of single transcription factors, genes, and proteins in VSMCs, many aspects of how VSMC plasticity can be clinically exploited to prevent the rupture of vulnerable atherosclerotic lesions remain unclear. Systems biology approaches have emerged as a promising avenue to study gene regulation at the organismal level, but have remained underutilized in cardiovascular medicine. This study proposes a new systems biology pipeline that integrates experimental and computational tools using datasets measuring the transcriptomes and epigenomes of phenotypically transitioning VSMCs. Our pipeline allows users to identify transcription factors that co-regulate mouse or human gene sets controlling specific VSMC phenotypic transitions during atherosclerosis. The computational algorithms developed in this workflow integrate an array of publicly available databases, such as the Ensembl genome database and the JASPAR transcription factor database. To demonstrate the translational potential of our pipeline, we used it to test the hypothesis that the murine mesenchymal marker stem cell marker *Scf1* (*Ly6a*) has a human ortholog, which has not been previously validated as a target in clinical trials. We used our pipeline to predict novel sets of VSMC co-regulators, which will be tested in concert, to promote a distinct phenotypic state marked by *Ly6a*. We also integrated genomic and ATAC-seq profiles of VSMCs to identify mouse genes with human orthologs that bind *Ly6a*-specific transcription factors. Specifically, we predict that the transcription factors FOS, EGR1, JUN, *JunB* and *Atf3* co-regulate sets of human genes with *Ly6a*. While prior studies have associated these transcription factors with VSMC proliferation and phenotypic modulation, we specifically predict these *Ly6a*-specific transcription factor co-regulate the following human gene sets: *TNS1*, *THSD4*, *RCAN2*, *PTPRD*, *CNE4*, *COL15A1*, and *ABLIM1* (FOS and *Atf3*). To validate these predictions, we will determine if these transcription factors and their co-regulated genes are expressed in single cell RNA-seq analyses of human carotid artery, and can be detected with immunofluorescence staining of human coronary artery sections. We will also use siRNA to target these candidate transcription factors and genes in cultured VSMCs. The results of our pipeline will help identify novel transcriptional regulatory networks in VSMCs that can be therapeutically targeted for the treatment of advanced atherosclerosis.

## Introduction

- Rupture or erosion of unstable atherosclerotic plaques is the underlying cause of myocardial infarction and stroke-related mortality, the leading causes of death worldwide.<sup>1</sup>
- A stable plaque has a low ratio of macrophages (MΦ) to vascular smooth muscle cells (VSMCs) and myofibroblasts (MF), which increase extracellular matrix (ECM) deposition in the fibrous cap.<sup>2</sup>



- VSMCs de-differentiate during atherosclerosis to take on both plaque stabilizing phenotypes (collagen producing myofibroblast-like cells) and plaque destabilizing phenotypes (i.e. macrophage-like cells).<sup>3</sup>
- Over the decade, many mechanisms have been identified in mice that play roles in VSMC-derived plaque stability, but we do not understand the complete regulatory states of VSMCs.
- A lack of this understanding limits the ability of identified mechanisms to be directly exploited clinically to improve plaque stability in humans.



## Hypothesis

Conserved cis/trans elements coordinately regulate sets of genes that confer specific cellular (VSMC) phenotypes.

## Approach

### Strategy for prediction of transcription factors (TFs) co-regulating gene sets

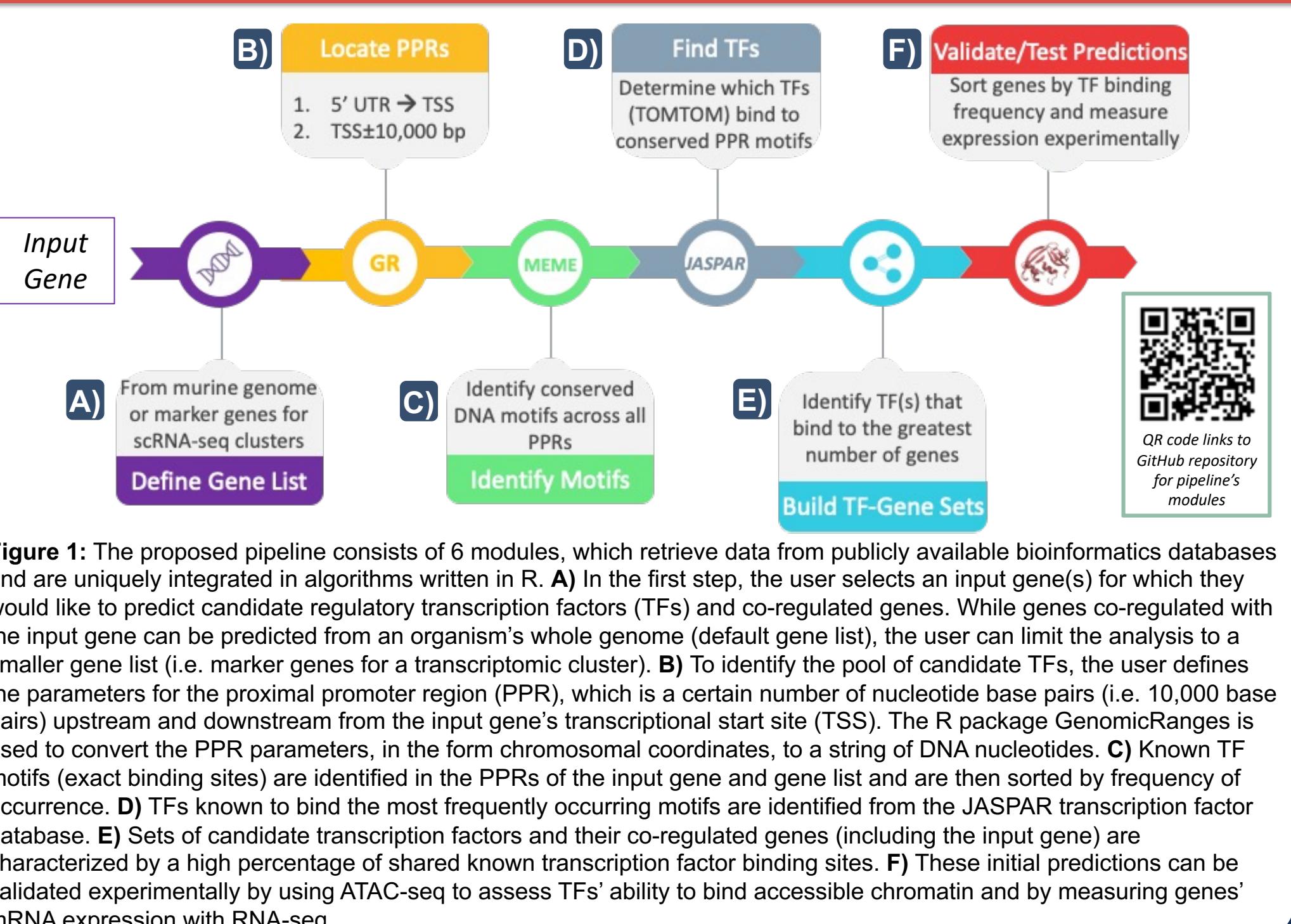
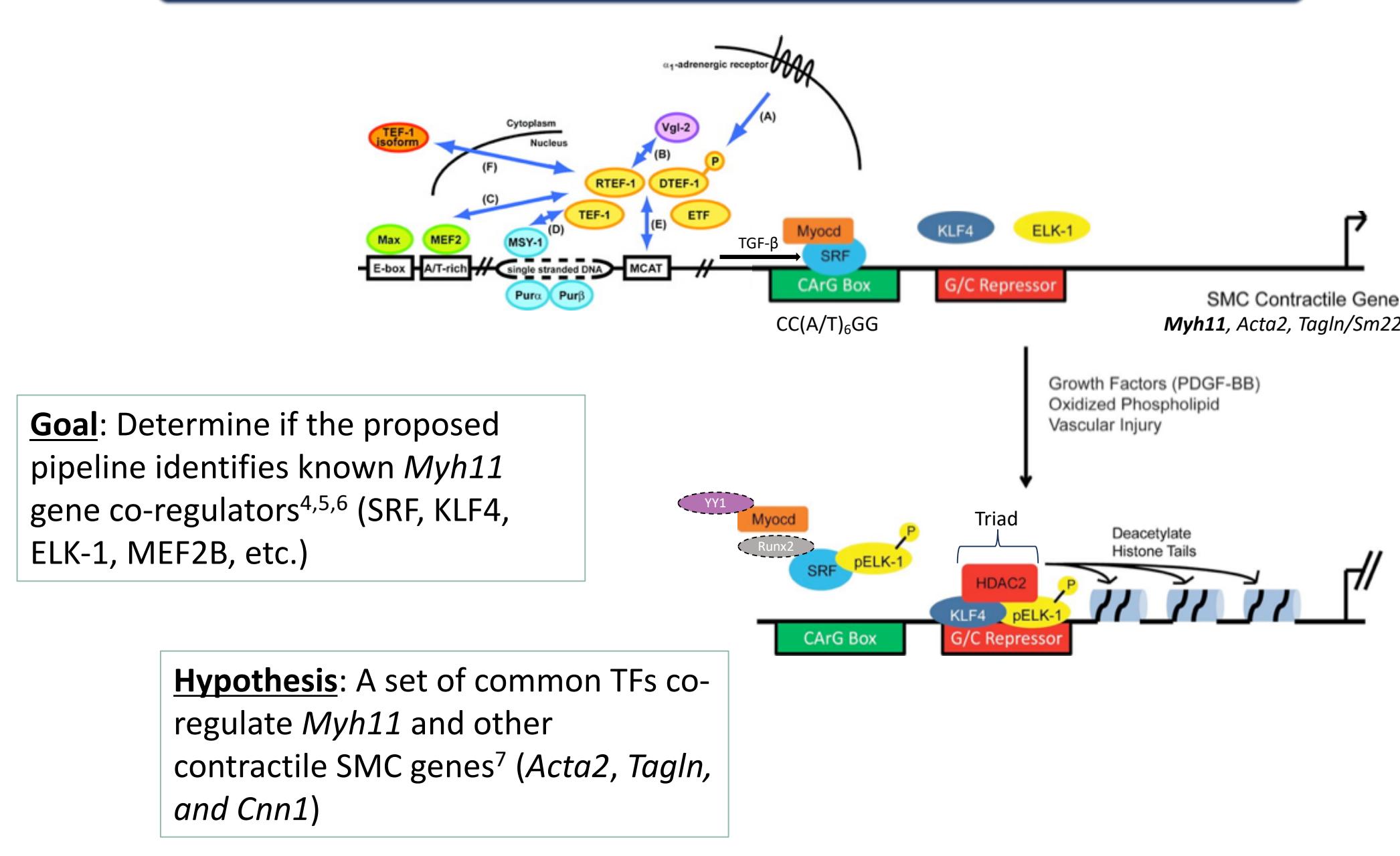


Figure 1: The proposed pipeline consists of 6 modules, which retrieve data from publicly available bioinformatics databases and are uniquely integrated in algorithms written in R. A) In the first step, the user selects an input gene(s) for which they would like to predict candidate regulatory transcription factors (TFs) and co-regulated genes. While genes co-regulated with the input gene can be predicted from an organism's whole genome (default gene list), the user can limit the analysis to a smaller gene list (i.e. marker genes for a transcriptomic cluster). B) To identify the pool of candidate TFs, the user defines the parameters for the proximal promoter region (PPR), which is a certain number of nucleotide pairs (i.e. 10,000 base pairs) upstream and downstream from the input gene's transcriptional start site (TSS). The R package GenomicRanges is used to convert the PPR parameters, in the form chromosomal coordinates, to a string of DNA nucleotides. C) Known TF motifs (exact binding sites) are identified in the PPRs of the input gene and gene list and are then sorted by frequency of occurrence. D) TFs known to bind the most frequently occurring motifs are identified from the JASPAR transcription factor database. E) Sets of candidate transcription factors and their co-regulated genes (including the input gene) are characterized by a high percentage of shared known transcription factor binding sites. F) These initial predictions can be validated experimentally by using ATAC-seq to assess TFs' ability to bind accessible chromatin and by measuring genes' mRNA expression with RNA-seq.

## Proof of Concept



### Step 1: Define input gene of interest (Myh11)

Because the regulatory mechanisms controlling the expression of the smooth muscle contractile marker Myosin Heavy Chain 11 (*Myh11*) have been extensively studied in mice, murine *Myh11* was chosen as the input gene of interest (positive control) to test whether the proposed pipeline identifies transcription factors (TEF-1, MAX, MEF2, MSY1, SRF, KLF4, ELK-1, etc.) that are known to co-regulate *Myh11* and other contractile genes like *Acta2* and *Tagln*.

### Step 2: Identify all TFs that bind the PPR of Myh11

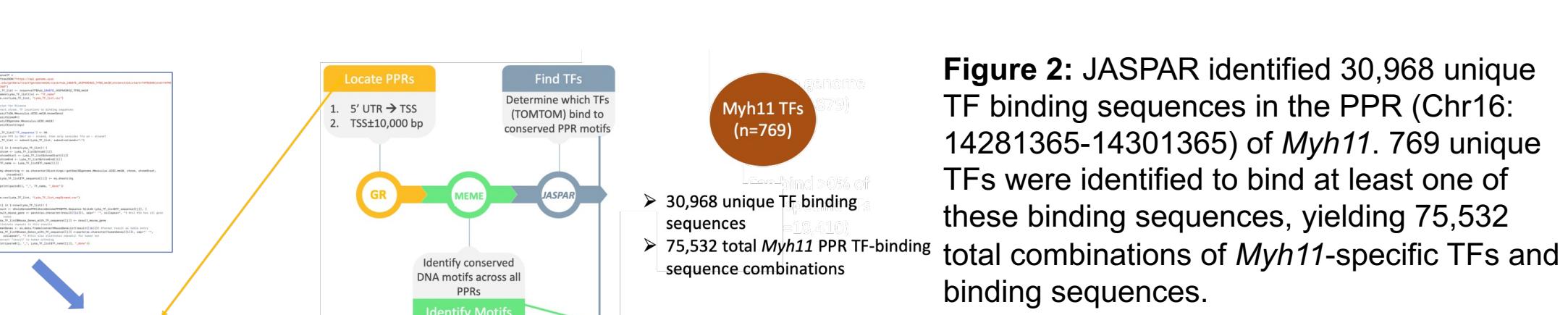


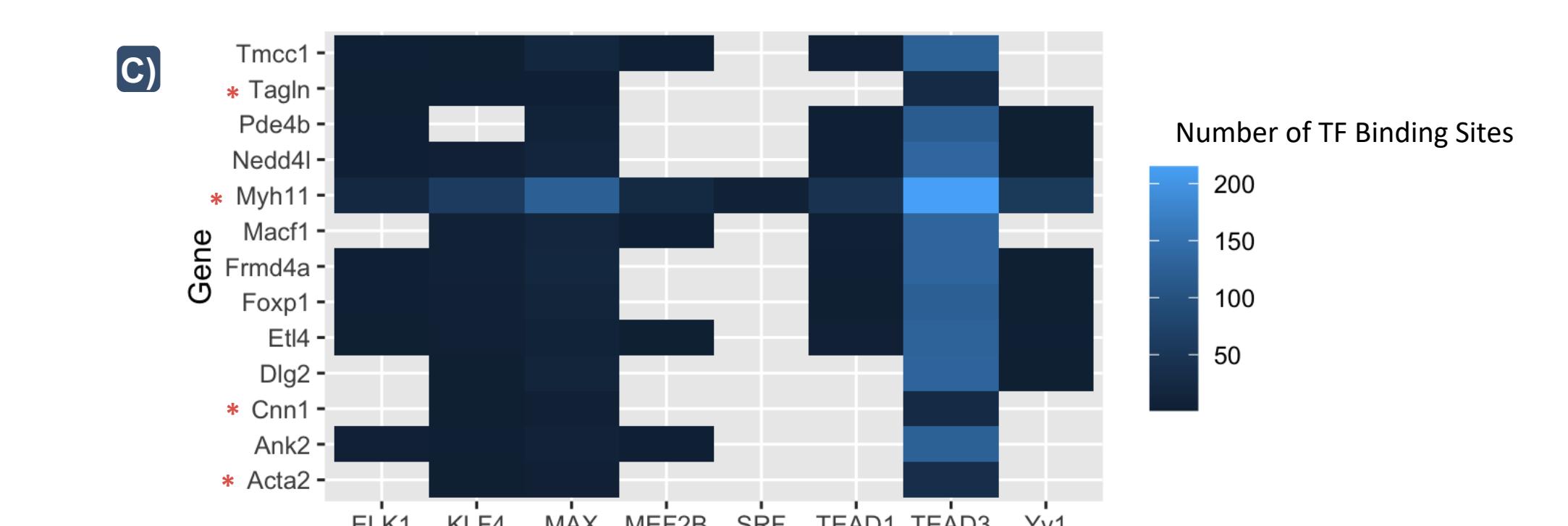
Figure 2: JASPAR identified 30,968 unique TF binding sequences in the PPR (Chr16: 14281365-14301365) of *Myh11*. 769 unique TFs were identified to bind at least one of these binding sequences, yielding 75,532 total *Myh11* PPR TF-binding sequence combinations.

Amongst the 769 candidate TFs, known *Myh11* transcriptional regulators, such as SRF and RUNX2, were identified as having the ability to bind the *Myh11* PPR.

### Step 3: Identify all mouse genes with PPRs that can bind *Myh11*-specific TFs

TF Name	TF binding sequence	Mouse genes with PPRs that have <i>Myh11</i> TF sequence
SRF	TTGCTTATAAGAGAT	<i>Myh11</i> , <i>Efcab12</i> , <i>Oas1e</i>
KLF4	AGTCCACACACAT	<i>Btn1</i> , <i>Gm10486</i> , <i>Gm16405</i> , <i>Gm10488</i> , <i>Gm10230</i> , <i>Sloc15a5</i> , <i>Csq1</i> , <i>Dpp6</i> , <i>Gk2</i> , <i>Mdfc</i> , <i>Myh11</i> , <i>Pearl15a</i> , <i>Pecam1</i> , <i>Dusp9</i> , <i>Oncut2</i> , <i>Sts15a6</i> , <i>Mlr1</i> , <i>C1qtrn12</i> , <i>Ergic1</i> , <i>C1qtnf2</i> , <i>Skl1</i> , <i>Cnnm4</i>

Figure 3: The 20 kb DNA nucleotide sequences for the PPRs of all genes in the mouse genome were identified in the pipeline's second module. A) Murine genes able to bind *Myh11*-specific TFs (i.e. SRF or KLF4) were determined by an exact nucleotide sequence match between the TF's binding sequence and the genes' PPR sequences. B) To identify candidate genes co-regulated with *Myh11*, all murine genes were ranked by the percentage of *Myh11*-specific TFs able to bind their PPRs. Here, the top ten genes whose PPRs share the greatest number of *Myh11*-specific TFs and binding sites are shown. Despite containing <10% of all possible *Myh11*-specific TF binding sites in their PPRs, previously validated contractile genes (*Acta2*, *Cnn1*, *Tagln*) were also identified as sharing *Myh11* regulatory TFs. C) This heatmap shows the diversity in the number of binding sites for *Myh11*-specific TFs in both the novel and previously validated contractile genes' PPRs\*.



## Proof of Concept

### Step 4: Validate *Myh11*-associated TFs that co-regulate genes with *Myh11*

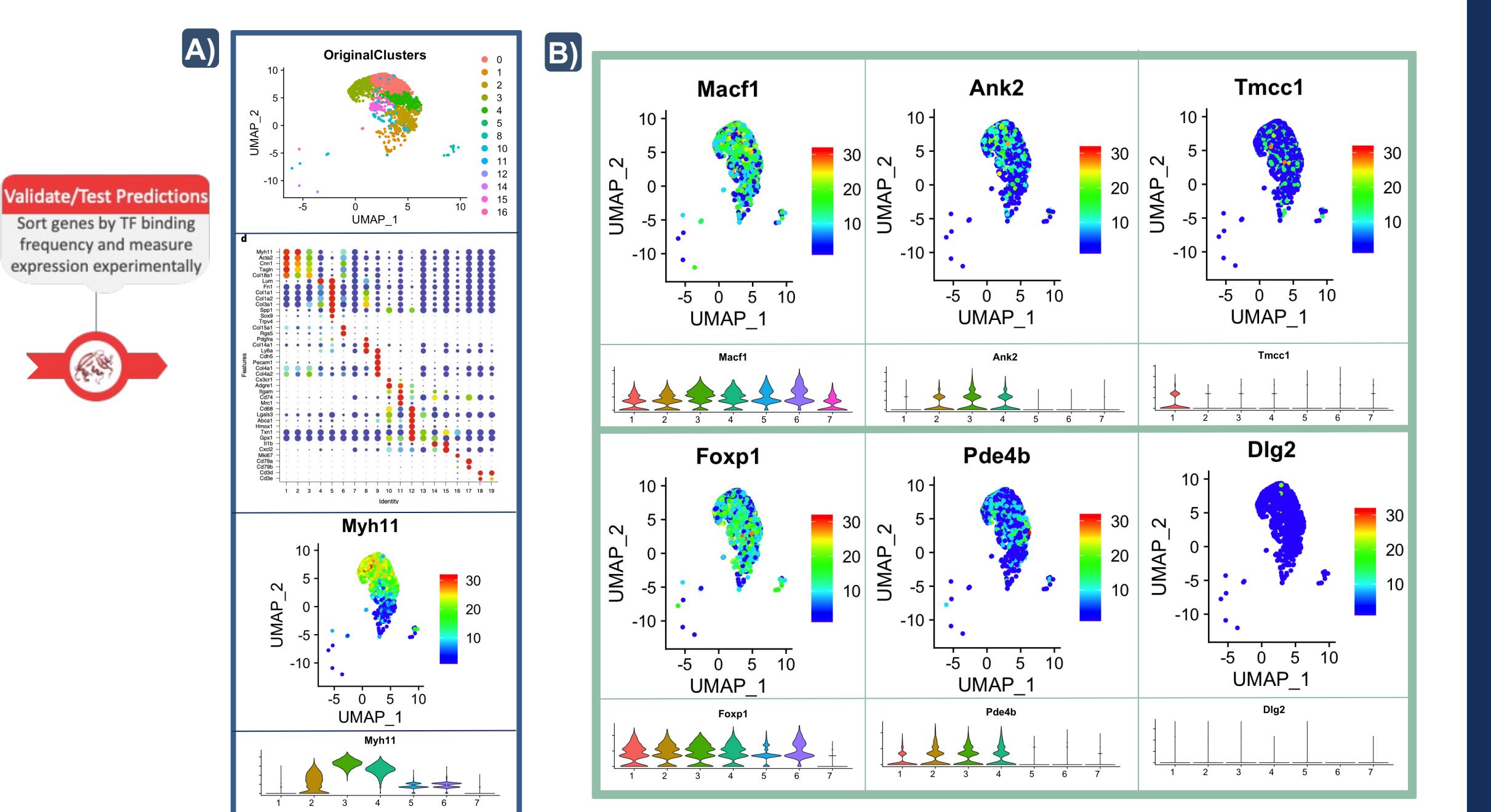


Figure 4: Candidate genes' co-expression is depicted in this UMAP presentation of scRNA-seq data from *Pdgfbrcm*-WT/WT mice.<sup>8</sup> A) scRNA-seq analysis of advanced BC lesions from *Pdgfbrcm*-WT/WT identifies 19 distinct cell clusters. Clusters 1-7 are *eyfp* positive (SMC-derived). UMAP representation of *Myh11* expression is seen in clusters 1-2. B) UMAP representation of the expression of 6 candidate genes (*Macf1*, *Foxp1*, *Ank2*, *Pde4b*, *Tmcc1*, *Dig2*) correlate with *Myh11* expression, supporting these genes' co-regulation with *Myh11*. At the bottom of each UMAP, violin plots show the corresponding gene's expression in SMC-clusters (clusters 1-7)

## Summary

- The proposed pipeline (TRBP) identifies both known genes (*Acta2*, *Tagln*, and *Cnn1*) and novel genes (*Macf1*, *Ank2*, *Foxp1*, and *Pde4b*) that are co-regulated with the contractile SMC marker *Myh11* and its associated TFs (TEAD1/TEF-1, TEAD3, KLF4, MAX, and Yy1).
  - Novel genes *Macf1* and *Ank2* are associated with the actin cytoskeleton, supporting their co-expression with *Myh11*.
- The murine *Ly6a*-specific TFs *FOS* and *Atf3* are predicted to co-regulate *Tns1*, *Thsd4*, *Rcan2*, *Ptpd*, *Cne4*, *Col15a1*, and *Ablim1*, which have known orthologs that are expressed in human SMCs.
- A complex network of TFs co-regulate multiple genes conferring VSMC contractile and plastic phenotypes.

## Future Directions

- Approach 1: Knock down predicted *Ly6a*-associated TFs and co-regulated genes in murine and human SMCs
  - Simultaneously treat cultured *eyfp*+ murine SMCs, in addition to HCAsMCs, with multiple siRNAs targeting Atf3, FOS, *Tns1*, *Thsd4*, *Rcan2*, *Ptpd*, etc. in *Ly6a*-stimulating conditions (i.e. PDGF-BB supplementation)
  - Current Expectation: Reduction in proliferative/migratory capacities in both human and murine SMCs, in addition to reduced *Ly6a* gene expression in murine SMCs
- Approach 2: Identify clinically approved drugs that target predicted co-regulatory TF-gene sets
  - Integrate DrugBank pharmacogenomics and proteomics data as a module in the current TRBP
  - Couple "outside-in" approach with current "inside-out" pipeline
  - Administer candidate drug(s) to cultured SMCs (and eventually *Myh11*-ERT2<sup>Cre</sup> *eyfp* ApoE<sup>-/-</sup> mice), measuring ability of SMCs to participate in injury-repair processes during atherosclerosis
  - Current Expectation: Loss of *Ly6a* phenotype in SMCs after treatment with a known drug (i.e. valproic acid)
- Approach 3: Account for nucleotide sequence variability in TF binding sites
  - Compute a "conservation score" that accounts for the frequencies of nucleotides in TF binding sequences (instead of using exact PPR-TF binding sequence matches)
  - Current Expectation: More TFs are predicted to co-regulate common gene sets

## Acknowledgements

Lab Management: Rebecca Deaton, Laura Shankman  
Computational Support: UVA Research Computing (Rivanna)

## Contact

Mahima Reddy msr3nf@virginia.edu  
twitter.com/MahimaReddy11  
linkedin.com/in/mahima-reddy

## References

- Benjamin, E. J. et al. Heart Disease and Stroke Statistics-2019 Update: A Report From the American Heart Association. *Circulation*. 2019.
- Tavora, F., Cresswell, N., Li, L., Fowler, D. & Burke, A. Frequency of Acute Plaque Ruptures and Thin Cap Atheromas at Sites of Maximal Stenosis. *Arg Bras Cardiol*. 2010.
- Shankman L., Gomez D., Cherepanova O.A., Salmon M., Alencar G.F., Haskins R.M., Switwielska P., Newman A.C., Greene E.S., Straub A.C., Isackson B.E., Randolph G.J., Owens G.K. KLF4-dependent phenotypic modulation of smooth muscle cells has a key role in atherosclerotic plaque pathogenesis. *Nature Medicine*. 2015.
- Yoshida, T. MCAT elements and the TEF-1 family of transcription factors in muscle development and disease. *Arteriosclerosis, thrombosis, and vascular biology*. 2008.
- Katoh, Y., Molentkin, J.D., Dave, V., Olson, E.N. and Periasamy, M., MEF2B is a component of a smooth muscle-specific complex that binds an A/T-rich element important for smooth muscle myosin heavy chain gene expression. *Journal of Biological Chemistry*. 1998.
- Huang, J. and Parmacek, M.S., Modulation of smooth muscle cell phenotype: The other side of the story. *Circulation research*. 2012.
- Sun, Q., Chen, G., Streb, J.W., Long, X., Yang, Y., Stoeckert, C.J. and Miano, J.M., Defining the mammalian CarGome. *Genome research*. 2006.
- Newman, A.A., Serbulea, V., Baylis, R.A., Shankman, L.S., Bradley, X., Alencar, G.F... & Owens, G. K. Multiple cell types contribute to the atherosclerotic lesion fibrous cap by PDGF $\beta$  and bioenergetic mechanisms. *Nature metabolism*. 2021.