

RESEARCH ARTICLE

# Discovery of gene regulatory elements through a new bioinformatics analysis of haploid genetic screens

Bhaven B. Patel<sup>1,2</sup>, Andres M. Lebensohn<sup>1,2</sup>, Ganesh V. Pusapati<sup>1</sup>, Jan E. Carette<sup>3</sup>, Julia Salzman<sup>1,4</sup>, Rajat Rohatgi<sup>1,2</sup>

**1** Department of Biochemistry, Stanford University School of Medicine, Stanford, California, United States of America, **2** Department of Medicine, Stanford University School of Medicine, Stanford, California, United States of America, **3** Department of Microbiology and Immunology, Stanford University School of Medicine, Stanford, California, United States of America, **4** Department of Biomedical Data Science, Stanford University School of Medicine, Stanford, California, United States of America

☞ These authors contributed equally to this work.

✉ Current address: Laboratory of Cellular and Molecular Biology, Center for Cancer Research, National Cancer Institute, National Institutes of Health, Bethesda, Maryland, United States of America

\* [andres.lebensohn@nih.gov](mailto:andres.lebensohn@nih.gov) (AML); [julia.salzman@stanford.edu](mailto:julia.salzman@stanford.edu) (JS); [rrohatchi@stanford.edu](mailto:rrohatchi@stanford.edu) (RR)



**OPEN ACCESS**

**Citation:** Patel BB, Lebensohn AM, Pusapati GV, Carette JE, Salzman J, Rohatgi R (2019) Discovery of gene regulatory elements through a new bioinformatics analysis of haploid genetic screens. *PLoS ONE* 14(1): e0198463. <https://doi.org/10.1371/journal.pone.0198463>

**Editor:** Bart O. Williams, Van Andel Institute, UNITED STATES

**Received:** May 16, 2018

**Accepted:** December 17, 2018

**Published:** January 29, 2019

**Copyright:** © 2019 Patel et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** All FASTQ files containing the sequencing data for unsorted (control) and sorted cells from the screens analyzed in this study have been deposited in the National Center for Biotechnology Information (NCBI) Sequence Read Archive (SRA) with Study accession number SRP094861. All BAIMS pipeline software files are available through Github (<https://github.com/RohatgiLab/BAIMS-Pipeline>).

**Funding:** This work was funded in part by a National Cancer Institute (<https://www.cancer.gov>)

## Abstract

The systematic identification of regulatory elements that control gene expression remains a challenge. Genetic screens that use untargeted mutagenesis have the potential to identify protein-coding genes, non-coding RNAs and regulatory elements, but their analysis has mainly focused on identifying the former two. To identify regulatory elements, we conducted a new bioinformatics analysis of insertional mutagenesis screens interrogating WNT signaling in haploid human cells. We searched for specific patterns of retroviral gene trap integrations (used as mutagens in haploid screens) in short genomic intervals overlapping with introns and regions upstream of genes. We uncovered atypical patterns of gene trap insertions that were not predicted to disrupt coding sequences, but caused changes in the expression of two key regulators of WNT signaling, suggesting the presence of cis-regulatory elements. Our methodology extends the scope of haploid genetic screens by enabling the identification of regulatory elements that control gene expression.

## Introduction

An outstanding challenge in genomics is the identification of functional regulatory elements that control spatial and temporal expression of protein-coding genes and non-coding RNAs. The Encyclopedia of DNA Elements (ENCODE) project has the ambitious goal of generating a candidate list of all functional elements in the human genome using sequence features, such as evolutionary conservation, and biochemical features, such as chromatin accessibility and chromatin modifications [1]. Functional approaches to identify regulatory elements have thus far focused on specific regions of the genome and include massively parallel reporter assays or

grant R00 CA168987-03, a National Institute of General Medical Sciences (<https://www.nigms.nih.gov>) grant R01 GM116847, a Joint Initiative for Metrology in Biology (<http://jimb.stanford.edu>) seed grant, a National Science Foundation (<https://www.nsf.gov>) CAREER Award, a McCormick-Gabilan Fellowship and a Baxter Family Fellowship (all to JS), by National Institutes of Health (<https://www.nih.gov>) grants DP2 AI104557 (JEC), DP2 GM105448 (RR), R35 GM118082 (RR) and by startup funds from the Stanford Cancer Institute (RR). BBP was supported by the Stanford Bio-X Undergraduate Summer Research Program and a Major Grant from the Stanford University Office of Undergraduate Advising and Research. AML was supported by the Stanford Dean's Postdoctoral Fellowship, the Stanford Cancer Biology Program Training Grant and the Novartis sponsored Fellowship from the Helen Hay Whitney Foundation (<http://hhwf.org>). JEC is a David and Lucile Packard Foundation (<https://www.packard.org>) fellow, JS is an Alfred P. Sloan Foundation (<https://sloan.org>) fellow in Computational & Evolutionary Molecular Biology, and RR is a Josephine Q. Berry Faculty Scholar in Cancer Research at Stanford. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

dense clustered regularly-interspaced short palindromic repeats (CRISPR)-mediated mutagenesis of <1 megabase pair segments around a locus of interest (reviewed in [2]).

In work published recently [3], we conducted a comprehensive set of forward genetic screens in haploid human cells to uncover genes required for signaling through the WNT pathway, which plays central roles in development, stem cell function, and cancer. The power of these screens, which used a quantitative transcriptional reporter as the basis for phenotypic selection, was highlighted by the identification of genes encoding both known and novel components that function at most levels of the WNT pathway, from the cell surface to the nucleus. Our previous analysis focused primarily on annotated protein-coding genes and non-coding RNAs. Since the mutant cell libraries used in these screens were generated through untargeted insertional mutagenesis of the genome with a gene trap (GT)-bearing retrovirus, we wondered whether we could use the datasets generated by these screens to uncover gene regulatory mechanisms that modulate the WNT signaling pathway. One advantage of using genome-wide insertional mutagenesis data to identify gene regulatory regions is that the analysis does not have to be targeted to a specific region of the genome, in contrast to the currently used CRISPR methods or massively parallel reporter assays [2] mentioned above. The insertional bias of the retroviral mutagen used in haploid screens limits the regions of the genome that can be searched because retroviruses have a propensity to insert around transcriptional start sites (TSS), promoters, and enhancers [4]. For this reason, we focused our analysis of retroviral GT insertions on non-coding regions in genes and immediately upstream of them. We note that the use of other insertional mutagens, such as other viruses or transposons, that have distinct insertional biases could allow this strategy to be used to search for regulatory elements in other regions of the genome.

Here we present a new bioinformatics pipeline designed to uncover gene regulatory elements and we provide evidence for regulatory regions in the first intron of the gene encoding the transcription factor AP4 (TFAP4), a positive regulator of WNT signaling [3], and in the genomic region upstream of the promoter for the gene encoding the WNT co-receptor LRP6.

## Materials and methods

### Reagent providers

Reagents were obtained from the following companies: Thermo Fisher Scientific, Waltham, MA; Sigma-Aldrich, St. Louis, MO; Bio-Rad, Hercules, CA; Promega, Madison, WI; GE Healthcare Life Sciences, Logan, UT; BD Biosciences, San Jose, CA; EMD Millipore, Billerica, MA; Cell Signaling Technology, Danvers, MA; Li-Cor, Lincoln, NE; Jackson ImmunoResearch Laboratories, West Grove, PA; Atlanta Biologicals, Flowery Branch, GA; QIAGEN Sciences, Hilden, Germany; New England Biolabs (NEB), Ipswich, MA.

### Antibodies

**For immunoblotting.** Primary antibodies: rabbit anti-AP4 (TFAP4) serum (1:2000, a gift from Takeshi Egawa [5]); rabbit anti-LRP6 (C5C7) (1:500, Cell Signaling Technologies Cat. # 2560); mouse anti-ACTIN (clone C4) (1:500, EMD Millipore Cat. # MAB1501).

Secondary antibodies: peroxidase AffiniPure goat anti-rabbit IgG (H+L) (1:7500, Jackson ImmunoResearch Laboratories Cat. # 111-035-003); IRDye 800CW donkey anti-rabbit IgG (H+L) (1:10,000, Li-Cor Cat. # 925-32213); IRDye 800CW donkey anti-mouse IgG (H+L) (1:10,000, Li-Cor Cat. # 926-32212).

Primary and secondary antibodies used for detection with the Li-Cor Odyssey imaging system were diluted in a 1 to 1 mixture of Odyssey Blocking Buffer (Li-Cor Cat. # 927-40000) and TBST (Tris buffered saline (TBS) + 0.1% Tween-20), and those used for detection by

chemiluminescence were diluted in TBST + 5% skim milk. Primary antibody incubations were done overnight at 4°C, and secondary antibody incubations were done for 1 hr at room temperature (RT).

**For immunostaining.** Primary antibodies: mouse anti-LRP6 (clone A59) (5μg/mL, EMD Millipore Cat. # MABS274).

Secondary antibodies: donkey anti-mouse IgG (H+L) Alexa Fluor 647 conjugate (1:200, Thermo Fisher Scientific Cat. # A-31571).

## Cell lines and growth conditions

WT HAP1-7TGP cells (S1 Fig) and genetically modified clonal derivatives were grown at 37°C and 5% CO<sub>2</sub> in CGM 2: Iscove's Modified Dulbecco's Medium (IMDM) with L-glutamine, with HEPES, without Alpha-Thioglycerol (GE Healthcare Life Sciences Cat. # SH30228.01); 1X GlutaMAX-I (Thermo Fisher Scientific Cat. # 35050079); 40 Units/ml Penicillin, 40 μg/ml Streptomycin (Thermo Fisher Scientific Cat. # 15140122); 10% Fetal Bovine Serum (FBS) (Atlanta Biologicals Cat. # S11150) [3].

## Bioinformatics analysis

**Bin-based Analysis of Insertional Mutagenesis Screens (BAIMS).** Genetic screens were conducted as described in the "Reporter-based forward genetic screens" section of Materials and methods in [3], except that GT integrations were mapped as follows. FASTQ files containing 36 base pair (bp) sequencing reads corresponding to genomic sequences flanking retroviral integration sites in both the sorted and unsorted control cells were obtained for the various genetic screens described (National Center for Biotechnology Information (NCBI) Sequence Read Archive (SRA) Study accession number SRP094861). Reads were aligned to the human genome version "GRCh38" using Bowtie alignment software, version 1.0.1 [6], allowing up to 3 base pair mismatches, and only reads that aligned to a single locus of the human genome were considered for downstream analysis. The orientation of the reads relative to the "+" or "-" strand of the chromosome, as defined in human genome version GRCh38, was noted.

Next, the genome was divided into contiguous, non-overlapping intervals of arbitrary length (250–1000 bp as indicated in the Results and figure legends), which are referred to as "bins", regardless of the location of genes and other genetic elements. Each bin was annotated with any overlapping genes and corresponding features (5'UTR, CDS, intron, and 3'UTR), according to the RefSeq annotations from the University of California, Santa Cruz Table Browser [7] for the GRCh38 assembly of the human genome. An additional genetic feature that we defined as "promoter," encompassing the 2000 bp directly upstream of the TSS of every gene, was also used to annotate any overlapping bins. The orientation of each genetic feature with respect to the chromosome (whether it resides on the "+" or "-" strand of the chromosome, as specified by the RefSeq annotation) was also noted.

Each GT insertion considered for downstream analysis was mapped to the bin that encompassed its location in the genome. For each bin, we tallied the number of insertions that mapped to the "+" and to the "-" chromosome strand. This enabled us to determine the number of sense and antisense insertions relative to any genetic feature. For example, a GT insertion that aligned to the "+" chromosome strand was considered to be in the sense orientation with respect to a genetic feature that resided on the "+" chromosome strand, whereas an insertion that aligned to the "-" chromosome strand was considered to be in the antisense orientation with respect to the same genetic feature. Histograms depicting the orientation of

insertions across genomic regions or genes of interest could then be generated using insertion counts from the bins contained within the region of interest.

**Gene-based insertion enrichment analysis.** To determine which genes were enriched for total GT insertions in the sorted versus the unsorted cells, all insertions in bins annotated with a given gene and its associated promoter as defined above were aggregated separately for the sorted and unsorted cell populations. Thus, the sum of insertions for a specific gene included both sense and antisense insertions that overlapped with the gene's features, including the promoter. For each gene, a  $p$ -value for the significance of enrichment was calculated using a one-sided Fisher's exact test run using the "scipy" package (version 0.7.2) in Python 2.7.5 by comparing the frequency of insertions in the gene in the sorted cells to the frequency of insertions in the gene in the unsorted cells; this  $p$ -value was then corrected for false-discovery rate. Genes were ranked in ascending order based on FDR-corrected  $p$ -value.

**Antisense intronic insertion enrichment analysis.** This analysis included bins annotated exclusively as intron and containing at least one GT insertion in the antisense orientation with respect to the gene in the sorted cells. An FDR-corrected  $p$ -value for the significance of antisense insertion enrichment in each of these bins was determined using a one-sided Fisher's exact test (from the "scipy" package for Python) comparing the frequency of antisense insertions in the bin for the sorted versus the unsorted cells. Bins were then ranked in ascending order based on FDR-corrected  $p$ -value (S1 File).

**Upstream insertion enrichment analysis.** This analysis included bins annotated exclusively as promoter and containing at least one GT insertion regardless of orientation in the sorted cells. An FDR-corrected  $p$ -value for the significance of insertion enrichment in each of these bins was determined using a one-sided Fisher's exact test (from "scipy" package for Python) comparing the frequency of insertions in the bin for the sorted versus the unsorted cells. Bins were then ranked in ascending order based on FDR-corrected  $p$ -value (S1 File).

**Inactivating insertion enrichment analysis.** This analysis included bins annotated with any exonic feature (5'UTR, CDS, 3'UTR) and containing at least one GT insertion regardless of orientation in the sorted cells, as well as bins annotated exclusively with intron and containing at least one GT insertion in the sense orientation with respect to the gene in the sorted cells. An FDR-corrected  $p$ -value for the significance of inactivating insertion (all insertions in bins annotated with 5'UTR, CDS, or 3'UTR and only sense insertions in bins annotated exclusively with intron) enrichment in the bin was determined using a one-sided Fisher's exact test (from "scipy" package for Python) comparing the frequency of insertions in the bin for the sorted versus the unsorted cells. Bins were ranked in ascending order based on FDR-corrected  $p$ -value (S1 File).

**BAIMS pipeline code.** The BAIMS pipeline code used for the bioinformatics analysis is available through Github (<https://github.com/RohatgiLab/BAIMS-Pipeline>).

## Isolation of cell lines containing GT insertions

All clonal cell lines containing specified GT insertions were isolated as described in the "Isolation of APC<sup>KO-2</sup> mutant cell line containing a GT insertion" section of Materials and methods in [3]. Briefly, following the WNT positive regulator high stringency screen, the same FACS gate used during the screen was used to sort single cells into 96-well plates. Colonies were harvested after 16 days, 1/10<sup>th</sup> of each clone was passaged for continued growth, and the remainder of the cells were collected and centrifuged. Genomic DNA was prepared from these cells using the QIAamp DNA mini kit (QIAGEN Sciences Cat. # 51304), and a nested PCR strategy was used to identify clones containing GTs in either *TFAP4* or *LRP6*. A genomic region of *TFAP4* or *LRP6* enriched for GT insertions was amplified by PCR using a forward primer

complementary to a unique sequence in the GT (pGT-Puro4: 5' -TCTCCAAATCTCGGTG GAAC-3') and a reverse primer complementary to a unique genomic sequence adjacent to the GT-enriched region in *TFAP4* or *LRP6*. 400 ng of genomic DNA was used as input for PCR amplification in 25  $\mu$ l reactions containing 1X LongAmp *Taq* reaction buffer, 300  $\mu$ M of each dNTP, 400 nM of each primer and 0.1 units/ $\mu$ l of LongAmp *Taq* DNA polymerase (NEB Cat. # M0323L). The presence of clones containing a GT insertion was evident as discrete bands when the PCR products were analyzed on a 1% agarose gel.

The *TFAP4*<sup>GT</sup> cell line containing an antisense GT insertion in the first intron of *TFAP4* was isolated from the WNT positive regulator high stringency screen using the reverse primer Wntlow *TFAP4* AS II (5' -GCTGCACACGTGTAGACACTC-3').

*LRP6*<sup>GT</sup>-1(Up) and *LRP6*<sup>GT</sup>-2(Up) cell lines, containing antisense GT insertions upstream of the *LRP6* TSS, and the *LRP6*<sup>GT</sup>-3(Int) cell line, containing a sense GT insertion in the first intron of *LRP6*, were isolated from the WNT positive regulator high stringency screen using the reverse primers *LRP6*UP-ASGT-Loc-2 (5' -GCAGTGTGTAATATCTCATTCCC-3'), *LRP6*UP-ASGT-Loc-1 (5' -GGAGACTCCCATTACTCTCTGTT-3') and Wntlow *LRP6* (5' -TGTGGGAAAACCTTTGTAATATGC-3'), respectively.

The genomic location of the GT insertion in each isolated cell line is indicated in [S2 File](#).

### Analysis of WNT reporter fluorescence

WNT reporter fluorescence was measured in WT HAP1-7TGP cells or derivatives thereof as described in the "Analysis of WNT reporter fluorescence" section of Materials and methods of [3].

### Quantitative RT-PCR (qPCR) analysis

All mRNA measurements were made as described in the "Quantitative RT-PCR analysis" section of Materials and methods in [3] using the *AXIN2* and *HPRT1* primer pairs described therein, the following forward and reverse primers for *TFAP4*: h*TFAP4*-RT-PCR-1-FOR (5' -GAGGGCTCTGTAGCCTTGC-3') and h*TFAP4*-RT-PCR-1-REV (5' -GAATCCCGCT TGATGCTCT-3'), and the following forward and reverse primers spanning two pairs of contiguous exons for *LRP6*: qPCR-*LRP6*-Exons-1-2-For (5' -GCTTCTGTGTGCTCCTGAG-3'), qPCR-*LRP6*-Exons-1-2-Rev (5' -TCCAAGCCTCCAACACTACAATC-3'), qPCR-*LRP6*-Exons-7-8-For (5' - GGAGATGCCAAAACAGACAAG -3'), and qPCR-*LRP6*-Exons-7-8-Rev (5' - CAGTCCAGTAAACATAGTCACCC -3').

### Immunoblot analysis of WT HAP1-7TGP and mutant cell lines

**Immunoblot analysis of *TFAP4*.** This analysis was performed as described in the "Immunoblot analysis of HAP1-7TGP and mutant cell lines—Immunoblot analysis of total AXIN1 and AXIN2" section of Materials and methods in [3] with some modifications. Cell pellets harvested from confluent 6 cm dishes were resuspended in 100  $\mu$ l of ice-cold RIPA lysis buffer (50 mM Tris-HCl pH 8.0, 150 mM NaCl, 2% NP-40, 0.25% deoxycholate, 0.1% SDS, 1X SIGMA-FAST protease inhibitors (Sigma-Aldrich Cat. # S8820), 10% glycerol), sonicated in a Bioruptor 300 (Diagenode) 2 x 30 sec in the high setting, centrifuged 10 min at 20,000 x g, and the supernatant was recovered.

The protein concentration in the supernatant was quantified using the Pierce BCA Protein Assay Kit. Samples were normalized by dilution with RIPA lysis buffer, further diluted with 4X LDS sample buffer supplemented with 50 mM TCEP, heated for 5 min at 95°C, and 40  $\mu$ g of total protein were electrophoresed alongside Precision Plus Protein All Blue Prestained Protein Standards (Bio-Rad Cat. # 1610373) in NuPAGE 4–12% Bis-Tris (Thermo Fisher



Scientific) gels using 1X NuPAGE MOPS SDS running buffer (Thermo Fisher Scientific Cat. # NP0001).

Proteins were transferred to nitrocellulose membranes using 1X NuPAGE transfer buffer (Thermo Fisher Scientific Cat. # NP0006) + 10% methanol. Membranes were stained with Ponceau S to assess loading, washed and blocked with TBST + 5% skim milk. Blots were incubated with rabbit anti-AP4 (TFAP4), washed with TBST, incubated with Peroxidase AffiniPure anti-rabbit secondary antibody, washed with TBST followed by TBS, and developed with SuperSignal West Pico Chemiluminescent Substrate and SuperSignal West Femto Maximum Sensitivity Substrate (Thermo Fisher Scientific Cat. # 34080 and 34095).

**Immunoblot analysis of total LRP6.** This analysis was performed as described in the previous section with the following modifications. 75  $\mu$ g of total protein were loaded in duplicate and electrophoresed in a NuPAGE 4–12% Bis-Tris gel. Following the transfer step, the nitrocellulose membrane was cut between the 50 and 75 kDa molecular weight standards and blocked for 1 hour with Odyssey Blocking Buffer. The top membrane was incubated with rabbit anti-LRP6 primary antibody, and the bottom membrane was incubated with mouse anti-ACTIN primary antibody. Membranes were washed with TBST and incubated with IRDye 800CW donkey anti-rabbit IgG and IRDye 800CW donkey anti-mouse IgG secondary antibodies, respectively. Membranes were washed with TBST followed by TBS, and imaged using the Li-Cor Odyssey imaging system. Acquisition parameters in the manufacturer's Li-Cor Odyssey Image Studio software were set so as to avoid saturated pixels in the bands of interest, and bands were quantified using manual background subtraction. The integrated intensity for LRP6 was normalized to that for ACTIN in the same lane and the average  $\pm$  SD from duplicate lanes was used to represent the data.

### Luciferase reporter assay

A portion of the *TFAP4* intronic region enriched for GT insertions (chr16:4,270,498–4,271,890, hg38) and a control intronic region lacking GT insertions (chr16:4,264,430–4,265,871, hg38) were amplified from HAP1 genomic DNA and cloned in an antisense orientation in front of a minimal promoter driving a firefly luciferase reporter gene (pGL4.23; Promega Cat. # E8411). To analyze the effect of each fragment on luciferase transcription, HAP1-7TGP cells were seeded in a 96-well plate at a density of 5000 per well, and the following day each well was transfected with a mixture containing 40 ng of firefly luciferase reporter plasmid and 10 ng of a control plasmid encoding a renilla luciferase gene driven by a constitutive Thymidine Kinase (TK) promoter to account for differences in transfection. After 72 hours, cells were washed in PBS, lysed, and firefly and renilla enzymatic activities were measured using the Dual Luciferase System (Promega Cat. # E1910). Reporter activity for each well was calculated as the firefly/renilla ratio.

### LRP6 cell-surface staining of WT HAP1-7TGP and mutant cell lines

Approximately 24 hr before staining, cells were seeded in a 6-well plate at a density of  $2 \times 10^5$  per well and grown in 2.5 ml of CGM 2. On the following day the cells were washed once in 3 ml of phosphate buffered saline (PBS), harvested in 0.5 ml of Accutase Cell Detachment Reagent (Thermo Fisher Scientific Cat. # NC9839010), resuspended in 1.5 ml of ice-cold CGM 2 and centrifuged 4 min at  $400 \times g$  at  $4^\circ\text{C}$  (all subsequent centrifugation steps were done in the same way). Cells were washed with 2 ml of ice-cold Iscove's Modified Dulbecco's Medium (IMDM) with L-glutamine, with HEPES, without Alpha-Thioglycerol (GE Healthcare Life Sciences Cat. # SH30228.01) + 1% Fetal Bovine Serum (FBS) (Atlanta Biologicals Cat. # S11150), centrifuged and resuspended in 150  $\mu$ l of mouse anti-LRP6 primary antibody in IMDM + 1%

FBS. Following a 1 hr incubation on ice, cells were washed with 1.8 ml of ice-cold IMDM + 1% FBS, centrifuged, washed with 2 ml of ice-cold IMDM + 1% FBS and centrifuged again. Cells were resuspended in 150  $\mu$ l of secondary antibody in ice-cold IMDM + 1% FBS and incubated on ice for 30 minutes. Cells were washed with 1.8 ml of ice-cold IMDM + 1% FBS, centrifuged, washed with 2 ml of ice-cold IMDM + 1% FBS and centrifuged again. Cells were resuspended in 200  $\mu$ l of PBS + 2% FBS and LRP6 cell-surface fluorescence was analyzed on a BD Accuri RUO Special Order System (BD Biosciences).

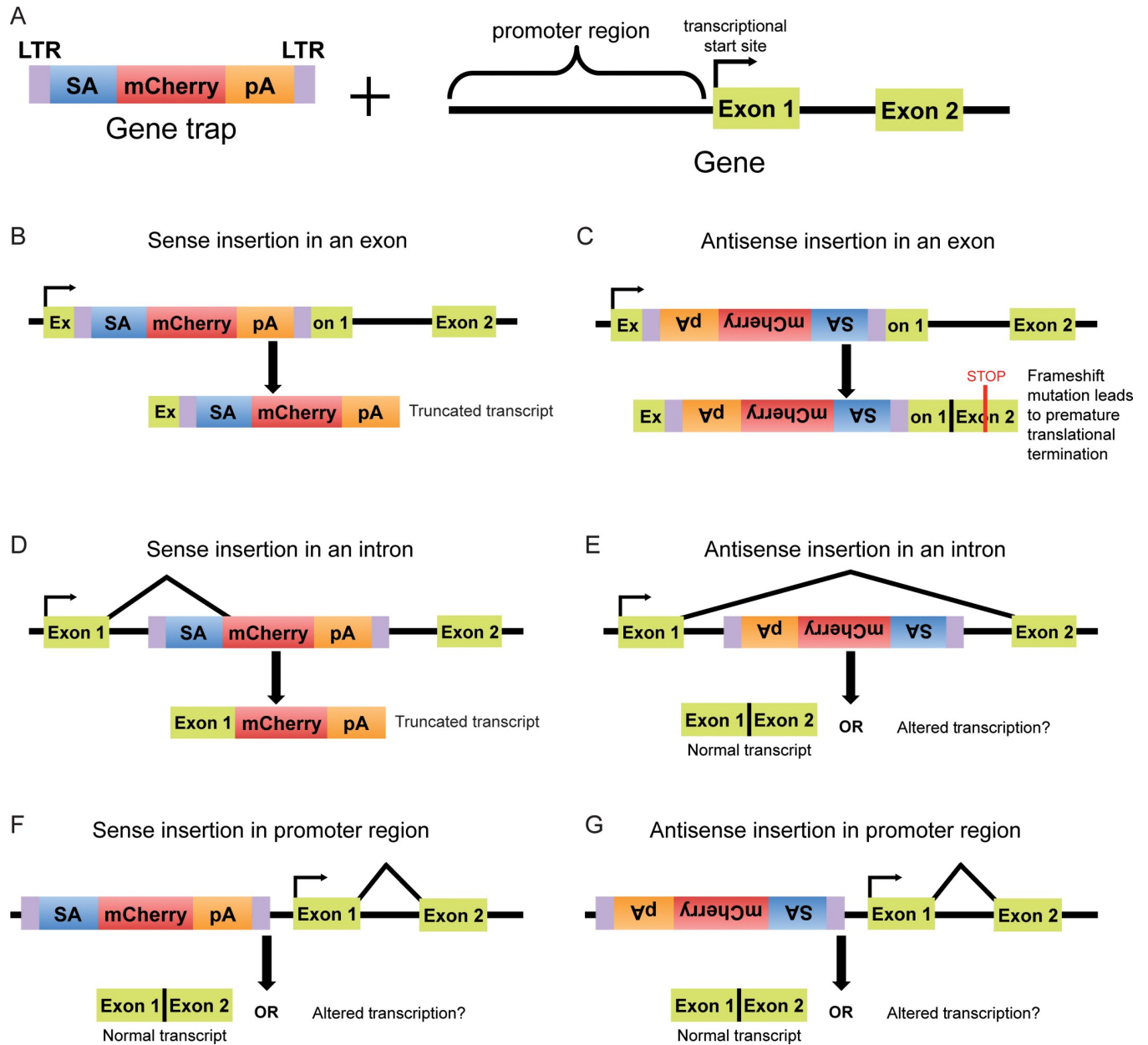
## Results

### Bin-based Analysis of Insertional Mutagenesis Screens (BAIMS)

Haploid genetic screens rely on the phenotypic selection of a population of cells mutagenized by integration of a GT-bearing retrovirus. GTs, which contain a splice acceptor (SA) and a transcriptional termination polyadenylation signal (pA), can disrupt protein-coding genes in two ways: (1) by inserting into an exon in either the sense or antisense orientation relative to the coding sequence of the gene or (2) by inserting into an intron in the sense orientation, such that the directional SA causes the GT to be spliced to the 3'-end of the preceding exon, resulting in a transcript that undergoes premature termination (Fig 1A–1D). Indeed, most hits in haploid genetic screens exhibit a bias towards such inactivating sense insertions in introns [8]. In contrast, antisense GT insertions in introns, and sense or antisense GT insertions in the promoter region of genes typically will not disrupt protein-coding transcripts (Fig 1E–1G), but such GT insertions could theoretically perturb gene regulation by directly interrupting a regulatory protein-binding site on DNA, by terminating a regulatory transcript, or by altering chromatin structure. Therefore, in principle it should be possible to find GT insertion patterns indicative of such regulatory mechanisms.

In order to map GT insertions in a way that would enable us to identify regulatory elements, we devised a bioinformatics pipeline that was completely agnostic to the boundaries of annotated genes and simply tracked the number and orientation of GT insertions in short genomic intervals of arbitrary size, defined as “bins” (Fig 2A). We refer to this approach as “Bin-based Analysis of Insertional Mutagenesis Screens”, or BAIMS. Sequencing reads adjacent to the location of GT insertions found in sorted (phenotypically selected) and unsorted (control) cells from haploid genetic screens were aligned to the human genome and assigned to the bin that encompassed the insertion (Fig 2B). The orientation of each insertion on the chromosome was defined according to the GRCh38 assembly of the human genome. Each bin was also annotated with any relevant genetic features it overlapped with—5' untranslated region (5'UTR), coding domain sequence (CDS), intron and 3' untranslated region (3'UTR)—using the RefSeq annotations from the University of California, Santa Cruz Table Browser [7] for the GRCh38 assembly of the human genome. We also defined an additional genetic feature, designated “promoter”, as the 2000 base pairs (bp) upstream of the TSS of each gene. This region typically includes the minimal promoter but may also contain other cis-regulatory elements. We annotated bins overlapping with this feature accordingly. The relative orientation of any insertion with respect to any feature can therefore be determined, allowing us to observe patterns of sense and antisense GT insertions across features of interest (Fig 2C). This information can be displayed in a histogram depicting insertions over any genomic region of interest (Fig 2C), providing a high-resolution picture of the insertional landscape. Thus, BAIMS enables us to identify atypical patterns of GT insertions in specific genetic features that could be indicative of regulatory elements.

The overall enrichment of GT insertions for any given gene in the selected versus the control cells from a haploid genetic screen can also be assessed by aggregating the insertions

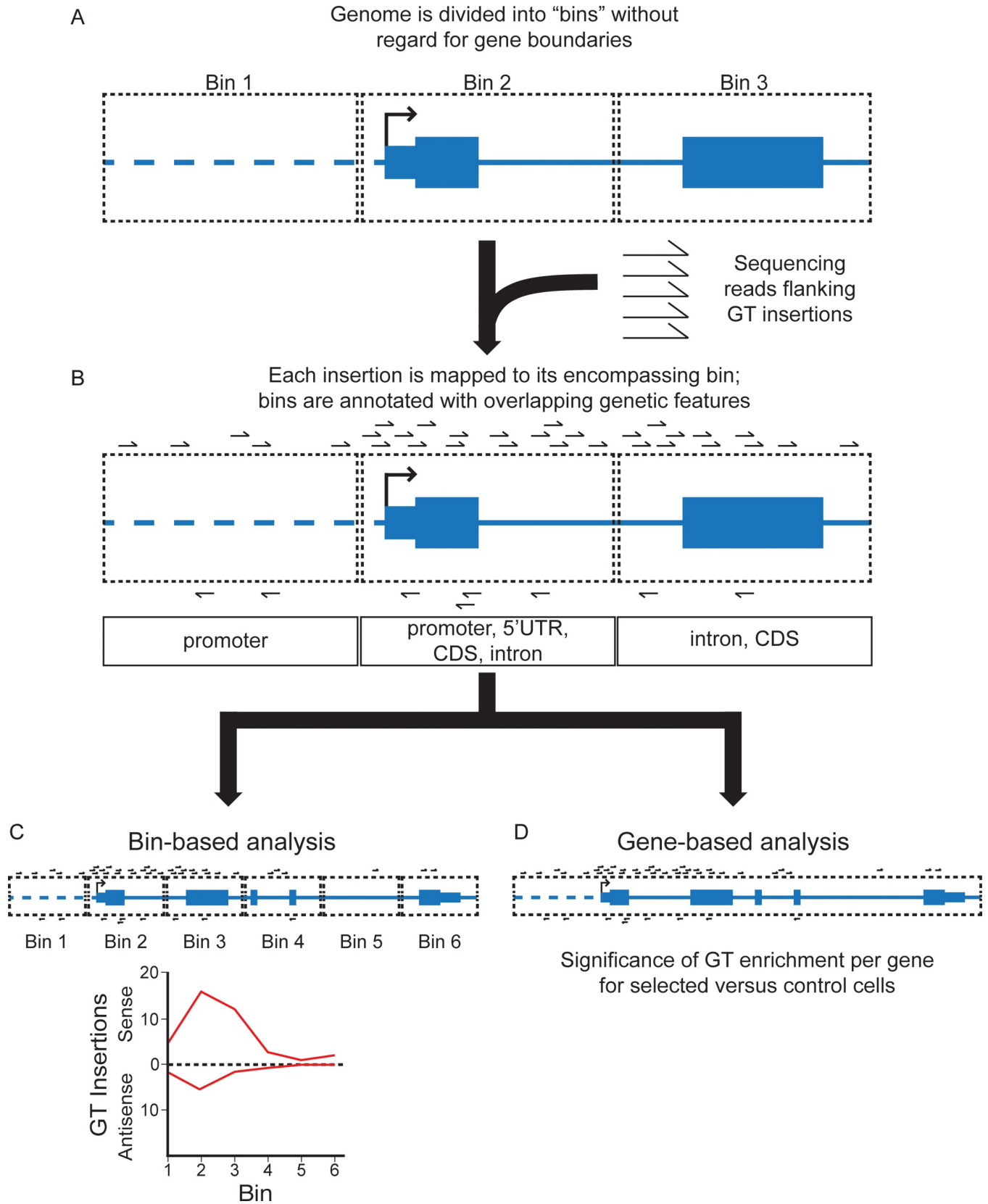


**Fig 1. Possible outcomes of GT insertions in different genetic features.** (A) A GT consists of direct long terminal repeats (LTRs), a strong splice acceptor (SA), a reporter gene (mCherry) and a poly-adenylation (pA) sequence. A schematic of the 5' end of a gene, including the promoter region, is also shown. (B) A GT can disrupt a gene by inserting into an exon in the sense orientation (with respect to the coding sequence of the gene), interrupting the coding sequence and causing premature transcriptional termination due to the pA sequence. (C) An antisense GT insertion into an exon interrupts the coding sequence of the gene and typically causes a frameshift mutation that leads to premature translational termination, producing a truncated protein. (D) When a GT integrates into an intron in the sense orientation, the SA causes the reporter gene and pA sequence to be spliced to the preceding exon, inevitably leading to premature transcriptional termination due to the pA sequence. (E) An antisense GT insertion in an intron will typically not disrupt a gene due to the directionality of the SA; however, it could interfere with regulatory elements or with transcripts present on the antisense strand. (F, G) GT insertions in the promoter region of a gene in either the sense or antisense orientation generally do not affect the downstream transcript; however, they could potentially disrupt regulatory elements and alter transcription.

<https://doi.org/10.1371/journal.pone.0198463.g001>

found in all bins that overlap with the gene (Fig 2D; see Materials and methods). We refer to this analysis, which produces a significance score for GT enrichment comparable to that of previous analyses [3], as “gene-based insertion enrichment analysis”.





**Fig 2. Schematic of Bin-based Analysis of Insertional Mutagenesis Screens (BAIMS).** (A) The human genome is computationally divided into “bins” (pictured as rectangles with black dotted lines), which encompass contiguous segments of DNA of an equal arbitrary length. Throughout this study, we used bins of 250 bp or 1000 bp in length, depending on the resolution required for the analysis. The boundaries of annotated genetic features, including genes and regulatory elements, are ignored. The depicted fictitious gene is modeled after a RefSeq gene track following the University of California, Santa Cruz (UCSC) genome browser display conventions: coding exons are represented by tall blocks, UTRs by shorter blocks, and introns by horizontal lines connecting the blocks. The arrow indicates the gene’s TSS. (B) Sequencing reads flanking the location of individual GT insertions in the control and selected cell populations from a haploid genetic screen are mapped to the human genome and assigned to the bin that encompasses the location of the insertion. The orientation of each insertion relative to the chromosome is noted. Bins are also annotated with any overlapping genetic features. These include promoter (defined as the 2000 bp upstream of the TSS, indicated by a horizontal dotted line), 5’UTR, CDS, intron, and 3’UTR. The orientation of the feature relative to the chromosome is also noted. (C) For the bin-based analysis, the number and orientation of GT insertions in consecutive bins along any defined portion of the genome (including but not limited to genes) is determined and can be depicted in a histogram (the number of sense GT insertions per bin is arbitrarily shown above the horizontal line labeled “0”, and the number of antisense insertions below), enabling the visualization of insertion patterns at sub-gene resolution. (D) For the gene-based analysis, GT insertions in bins that overlap with genes can be summed to obtain a total insertion count for each gene. The significance of GT enrichment for every gene is calculated by comparing the total number of insertions per gene found in the selected versus the control cell populations (see [Materials and methods](#) for details).

<https://doi.org/10.1371/journal.pone.0198463.g002>

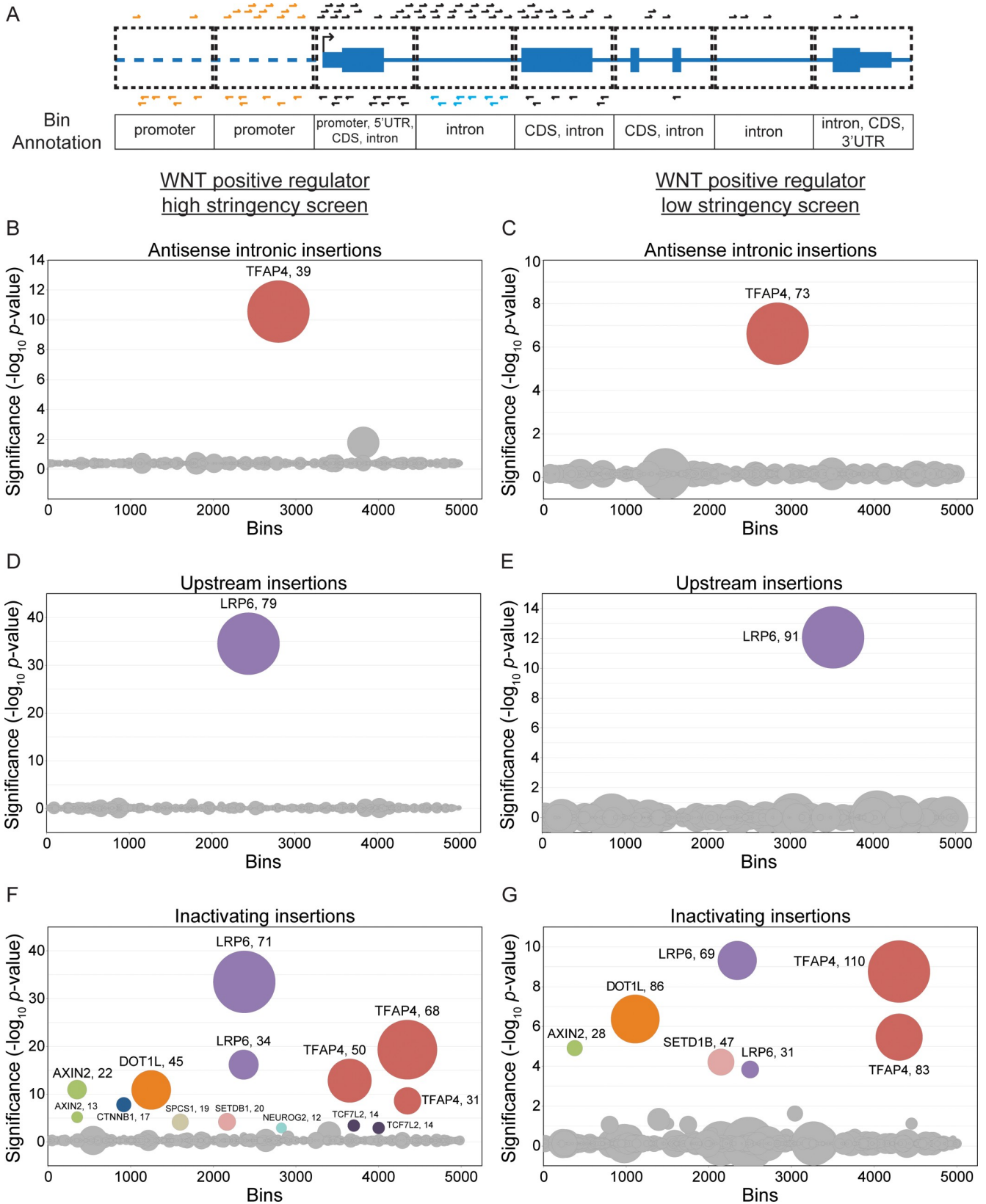
## Identification of regulatory elements through the analysis of bins with atypical GT insertion patterns

Our previous analysis [3] focused on GT insertions predicted to inactivate protein-coding genes and non-coding RNAs as outlined above: sense and antisense insertions in exons, and sense insertions in introns (Fig 1B–1D). To identify regulatory elements, we searched for GT insertion patterns distinct from these. Because the GT retrovirus has a strong propensity to integrate near TSSs, promoters and enhancers, we limited our analysis to non-coding regions within and adjacent to genes. We used BAIMS to look for enrichment of antisense insertions in introns, which would not be predicted to interrupt protein-coding transcripts (Fig 1E), and for enrichment of insertions in either orientation in the regions upstream of the TSS of genes (Fig 1F and 1G). Since each bin is annotated with the genetic features it overlaps with (Fig 3A), we could readily identify these distinct patterns of GT insertions.

To identify regulatory elements in introns, we looked for enrichment of antisense insertions in bins exclusively annotated as intron (Fig 3A); we refer to this analysis as “antisense intronic insertion enrichment analysis.” To identify regulatory elements in regions immediately upstream of genes, we looked for enrichment of both sense and antisense GT insertions in bins exclusively annotated as promoter (Fig 3A); we refer to this analysis as “upstream insertion enrichment analysis.” To distinguish features identified in these two new analyses from the more typical disruption of protein-coding genes or non-coding RNAs by GT insertions, we looked for enrichment of gene-inactivating insertions, as defined above (sense and antisense insertions in bins annotated with 5’UTR, CDS or 3’UTR, and sense insertions in bins annotated exclusively as intron; see Fig 3A); we refer to this analysis as “inactivating insertion enrichment analysis.”

These three analyses were applied to the data from two screens for positive regulators of signaling following stimulation with WNT3A, henceforth referred to as the WNT positive regulator high and low stringency screens, which differed only in the stringency of selection [3]. In these screens, HAP1 cells harboring a WNT-responsive fluorescent reporter (described in [9] and S1 Fig), hereafter referred to as WT HAP1-7TGP, were mutagenized with GT retrovirus, treated with WNT3A and sorted by fluorescence activated cell sorting (FACS) for cells that exhibited the lowest 2% (high stringency screen) or the lowest 10% (low stringency screen) signaling activity. These screens enabled us to identify known and new regulators in the WNT pathway [3].

Antisense intronic insertion enrichment analysis of the WNT positive regulator high and low stringency screens produced only one significant (FDR-corrected  $p$ -value < 0.01) bin (Fig 3B and 3C, S1 File), which mapped to the gene *TFAP4*, one of the top hits from these screens



**Fig 3. BAIMS identifies atypical GT insertion patterns in screens for regulators of WNT signaling.** (A) Schematic depicting various patterns of GT insertions relative to genetic features in the bins, used for the antisense intronic, upstream, and inactivating insertion enrichment analyses (see text for details). A fictitious gene modeled after a RefSeq gene track, with GT insertions in the sense orientation relative to the gene depicted above the track and in the antisense orientation depicted below it. The antisense intronic insertion enrichment analysis accounts for antisense GT insertions in bins annotated exclusively as intron (depicted in blue) and the upstream insertion enrichment analysis accounts for both sense and antisense insertions in bins annotated exclusively as promoter (depicted in orange). These two classes of insertions had been ignored in previous gene-based analyses of haploid genetic screens [3]. The inactivating insertion enrichment analysis accounts for both sense and antisense insertions in bins annotated as 5'UTR, CDS, or 3'UTR, as well as sense insertions in bins annotated exclusively as intron; these insertions (depicted in black) include all the gene-inactivating insertions used in previous analyses. (B-G) Circle plots depicting the results of antisense intronic (B, C), upstream (D, E), and inactivating (F, G) insertion enrichment analyses for the WNT positive regulator high stringency (B, D, and F) and low stringency (C, E, and G) screens. Circles represent individual 1000 bp bins. The y-axis indicates the significance of GT insertion enrichment in the selected versus the control cells, expressed in units of  $-\log_{10}$ (FDR-corrected  $p$ -value), and the x-axis indicates the 5000 bins with the smallest FDR-corrected  $p$ -values, arranged in random order. Circles representing bins with an FDR-corrected  $p$ -value  $< 0.01$  are colored and labeled with the name of the gene with which the bin overlaps. Circles representing bins corresponding to the same gene are depicted in the same color. The diameter of each circle is proportional to the number of independent GT insertions mapped to the corresponding bin in the selected cells, which is also indicated next to the gene name for enriched bins.

<https://doi.org/10.1371/journal.pone.0198463.g003>

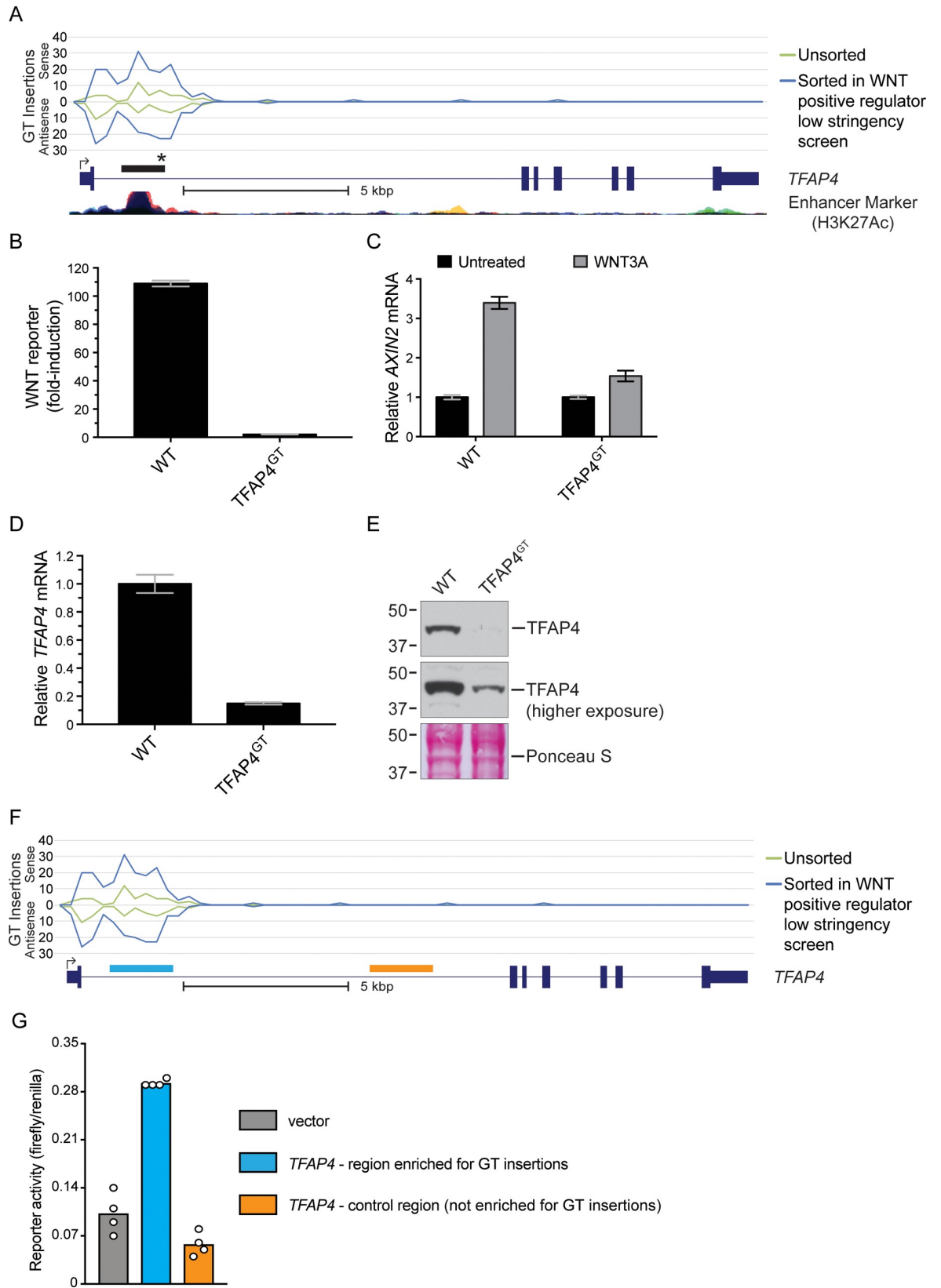
[3]. Upstream insertion enrichment analysis of the same screens produced only one significant bin upstream of *LRP6* (Fig 3D and 3E, S1 File), which was the top hit of both of these screens [3]. These results are markedly different from those of the inactivating insertion enrichment analysis of the same screens (Fig 3F and 3G, S1 File), which revealed bins in many of the same genes identified as significant hits in these screens [3].

In the sections that follow, we tested if the GT insertion patterns identified in *TFAP4* and *LRP6* by the antisense intronic and upstream insertion enrichment analyses, respectively, reflected regulatory effects on gene expression.

### Antisense GT insertions in the first intron of *TFAP4* disrupt a transcriptional enhancer element

The second top hit in the WNT positive regulator high and low stringency screens was *TFAP4*, encoding the transcription factor TFAP4, which we have shown to be a positive regulator of the WNT pathway acting downstream of the key transcriptional co-activator  $\beta$ -catenin (CTNNB1) [3]. As is common for top hits of haploid genetic screens, the 5' end of *TFAP4* was significantly enriched for inactivating GT insertions, including many sense and antisense insertions in the first exon as well as sense insertions in the first intron, which are all expected to disrupt the *TFAP4* coding sequence (Fig 4A and A in S2 Fig). However, the single bin identified in the antisense intronic insertion enrichment analysis (Fig 3B and 3C, S1 File) was also located in the first intron and it contained a comparable number of sense and antisense GT insertions (Fig 4A and A in S2 Fig). The enrichment of anti-sense insertions in the first intron was unexpected since these would not be expected to disrupt the *TFAP4* coding sequence. This pattern of GT insertion enrichment was not seen for *TFAP4* in the mutagenized but unsorted cells used as a control for the WNT positive regulator screens (Fig 4A) or for other top hits, such as *DOT1L*, in the sorted cells from these same screens (B in S2 Fig). These results suggested that antisense GT insertions in the first intron of *TFAP4* (which would not be predicted to terminate the *TFAP4* transcript) reduced WNT signaling.

To confirm this prediction, we isolated a clonal cell line harboring an antisense GT insertion in the first intron of *TFAP4* (we designate this cell line TFAP4<sup>GT</sup>; see Fig 4A and Materials and methods). WNT3A-induced reporter activation was nearly eliminated in TFAP4<sup>GT</sup> cells when compared to WT HAP1-7TGP cells (Fig 4B). Expression of *AXIN2* mRNA, a universal target gene of the pathway, following treatment with WNT3A was also reduced substantially in TFAP4<sup>GT</sup> cells (Fig 4C). Given the insertion's location within the boundaries of the *TFAP4* gene, we tested whether the antisense GT insertion affected expression of *TFAP4* itself. Both





**Fig 4. Antisense GT insertions in the first intron of *TFAP4* disrupt a transcriptional enhancer element and impair WNT signaling.** (A) The histogram indicates the number and orientation of GT insertions mapped to *TFAP4* in unsorted cells and in the sorted cells from the WNT positive regulator low stringency screen. Values above the horizontal line labeled “0” indicate sense insertions relative to the coding sequence of the gene, and values below it indicate antisense insertions. The x-axis represents contiguous 250 bp bins to which insertions were mapped (Chromosome 16, 4257249–4273000 bp). Insertions mapped for the different cell populations indicated in the legend are depicted by traces of different colors. A RefSeq gene track for *TFAP4* (following UCSC genome browser display conventions, described in the legend of Fig 2A) and an ENCODE track for histone3-lysine27-acetylation, a marker for enhancer activity (taken from the UCSC genome browser), are shown underneath the graph. The black rectangle above the gene track indicates the location of the bin identified in the antisense intronic insertion enrichment analyses of both the WNT positive regulator low stringency and high stringency screens. The black star denotes the position of the antisense GT insertion (located at NC\_000016.13:g.4271036\_4271037insGenetrap (Dec.2013: hg38, GRCh38) [10]; see S2 File) in the *TFAP4*<sup>GT</sup> clonal cell line used for further characterization. A scale bar is provided beneath the gene track for reference. (B) Fold-induction in WNT reporter (median  $\pm$  standard error of the median (SEM) EGFP fluorescence from 10,000 cells) following treatment with 50% WNT3A conditioned media (CM). (C) *AXIN2* mRNA (average  $\pm$  standard deviation (SD) of *AXIN2* mRNA normalized to *HPRT1* mRNA, each measured in triplicate qPCR reactions) relative to untreated cells. Where indicated, cells were treated with 50% WNT3A CM. (D) *TFAP4* mRNA (average  $\pm$  SD of *TFAP4* mRNA normalized to *HPRT1* mRNA, each measured in triplicate qPCR reactions) relative to WT HAP1-7TGP cells. (E) Immunoblot of *TFAP4*. The middle panel shows a higher exposure of the same blot shown in the top panel, and the bottom panel displays Ponceau S staining of the same blot as a loading control. Molecular weight standards in kiloDaltons (kDa) are indicated to the left of each blot. (F) Histogram of GT insertions mapped to *TFAP4* as in (A), with blue and orange boxes depicting the regions within the first intron tested in the transcriptional reporter assays shown in (G). GT insertions were enriched in the genomic region marked in blue (Chromosome 16, 4270498–4271890 bp) but were not enriched in a nearby control region marked in orange (Chromosome 16, 4264430–4265871 bp). (G) Luciferase reporter activity (ratio of firefly to renilla luciferase) in extracts of WT HAP1-7TGP cells transfected with a firefly luciferase gene driven by a minimal promoter alone (vector control, grey bar) or by the same minimal promoter with either of the two regions of *TFAP4* shown in (F) cloned upstream (blue and orange bars). Renilla luciferase was driven by a constitutive promoter and serves as a control to normalize for differences in transfection. Bars show the average firefly to renilla luciferase ratio from 4 replicate wells, and circles indicate the ratio for each replicate.

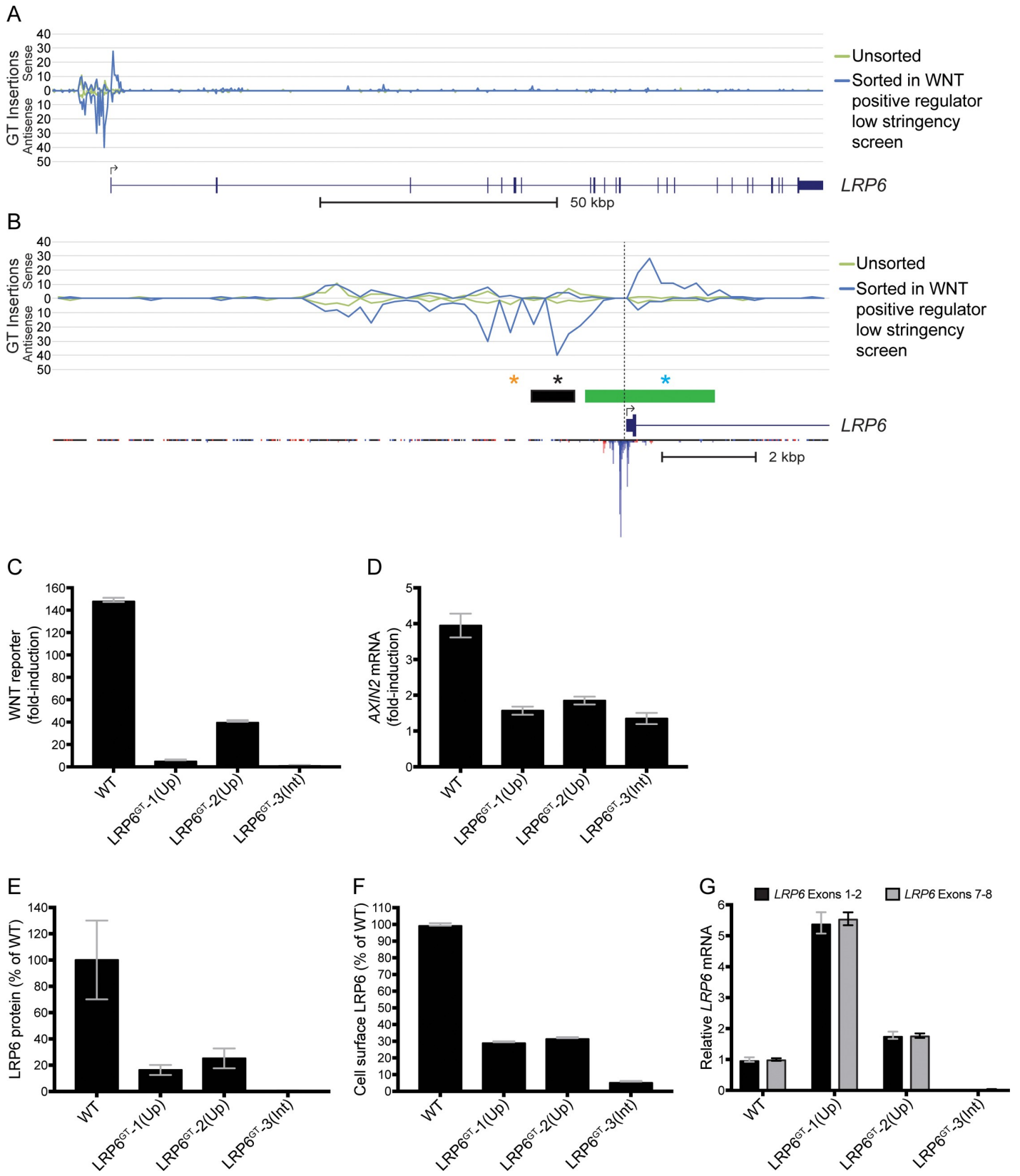
<https://doi.org/10.1371/journal.pone.0198463.g004>

*TFAP4* mRNA and protein levels were severely reduced in *TFAP4*<sup>GT</sup> cells, explaining the observed defect in pathway activity (Fig 4D and 4E). A higher exposure of the *TFAP4* immunoblot from *TFAP4*<sup>GT</sup> cells revealed a faint band corresponding to *TFAP4* (Fig 4E), indicating that the antisense GT insertion in the first intron of *TFAP4* reduced expression of a full-length transcript and protein as opposed to disrupting the coding sequence.

These results suggested the possibility that the antisense GT insertions disrupted an intronic regulatory element that enhances the transcription of *TFAP4*. To test this hypothesis, we cloned the intronic region encompassing most of the GT insertions upstream of a luciferase reporter gene driven by a minimal promoter (Fig 4F). As a control, we also cloned a nearby region in the same *TFAP4* intron that was not enriched for GT insertions. Only the intronic region enriched for GT insertions, but not the control region, enhanced expression of the reporter gene (Fig 4G). We conclude from these experiments that the BAIMS antisense intronic insertion enrichment analysis uncovered a transcriptional enhancer element that controls *TFAP4* expression and, consequently, WNT signaling output.

### Antisense GT insertions upstream of *LRP6* reduce *LRP6* protein abundance independently of mRNA levels

*LRP6* encodes a required co-receptor for WNT ligands and was the top hit of the WNT positive regulator high and low stringency screens [3]. As expected, most GT insertions in the *LRP6* gene proper (downstream of the TSS) were in the sense orientation with respect to the coding sequence (Fig 5A and 5B, and A and B in S3 Fig). However, our upstream insertion enrichment analysis also revealed a bin enriched for GT insertions located upstream of the TSS (Fig 3D and 3E, S1 File). A closer inspection of the region surrounding this bin revealed a pronounced enrichment of antisense insertions extending from about 1 to 3.5 kilobase pairs (kbp) upstream of the TSS (Fig 5B and 5B in S3 Fig). Importantly, this region was located upstream of the annotated *LRP6* promoter in *Ensembl* (Fig 5B). These GT insertion patterns were not observed in the mutagenized but unsorted cells used as a control for the WNT



**Fig 5. Antisense GT insertions upstream of *LRP6* reduce *LRP6* protein expression and impair WNT signaling.** (A) The histogram indicates the number and orientation of GT insertions mapped to *LRP6* and to the region ~12.5 kbp upstream of the TSS in unsorted cells and in the sorted cells from the WNT positive regulator low stringency screen. See legend to Fig 4A for details. The x-axis represents contiguous 250 bp bins to which insertions were mapped (Chromosome 12, 12116000–12279249 bp). (B) The histogram shows an expanded view of the 5' end of *LRP6* and the region ~12.5 kbp upstream of the TSS (left of the vertical dotted line), with traces for GT insertions mapped in unsorted cells and in the sorted cells from the WNT positive regulator low stringency screen. The x-axis represents contiguous 250 bp bins to which insertions were mapped (Chromosome 12, 12262500–12279249 bp). The green rectangle above the gene track indicates the location of the *LRP6* promoter according to *Ensembl* and the black rectangle indicates the location of the bin identified in the upstream insertion enrichment analyses of the WNT positive regulator low stringency and high stringency screens. The black and orange stars denote the positions of the antisense GT insertions (located at NC\_000012.13:g.12268371\_12268372insGenetrap (Dec.2013: hg38, GRCh38) and NC\_000012.13:g.12269383\_12269384insGenetrap (Dec.2013: hg38, GRCh38) respectively [10]; see S2 File) in the *LRP6*<sup>GT</sup>-1(Up) and *LRP6*<sup>GT</sup>-2(Up) clonal cell lines, respectively, and the blue star denotes the position of the sense GT insertion (located at NC\_000012.13:g.12266072\_12266073insGenetrap (Dec.2013: hg38, GRCh38); see S2 File) in the *LRP6*<sup>GT</sup>-3(Int) cell line. The inverted histogram below the RefSeq gene track for *LRP6* indicates the maximum CAGE read count found in any tissue sample from the FANTOM5 database [11]. (C) Fold-induction in WNT reporter (median +/- SEM EGFP fluorescence from 20,000 cells) following treatment with 50% WNT3A CM. (D) Fold-induction in *AXIN2* mRNA (average +/- SD of *AXIN2* mRNA normalized to *HPRT1* mRNA, each measured in triplicate qPCR reactions) following treatment with 50% WNT3A CM. (E) Quantification of immunoblot analysis of total *LRP6* protein (average +/- SD *LRP6* intensity normalized to *ACTIN* intensity from samples run in duplicate) shown as percentage of WT HAP1-7TGP. The blot used for quantification is shown in C in S3 Fig. (F) Cell surface *LRP6* protein (median +/- SEM cell surface *LRP6* immunofluorescence from 20,000 cells) shown as percentage of WT HAP1-7TGP. (G) *LRP6* mRNA (average +/- SD of *LRP6* mRNA, measured using two different primer pairs, normalized to *HPRT1* mRNA, each measured in triplicate qPCR reactions) shown relative to WT HAP1-7TGP cells.

<https://doi.org/10.1371/journal.pone.0198463.g005>

positive regulator screens (Fig 5A and 5B). These results suggested that antisense insertions upstream of *LRP6* impaired WNT signaling.

To test this possibility, we isolated two clonal cell lines containing antisense GT insertions in the region upstream of *LRP6* (we designate these cell lines *LRP6*<sup>GT</sup>-1(Up) and *LRP6*<sup>GT</sup>-2(Up); see Fig 5B and Materials and methods) and as a control we isolated a clonal cell line with a sense GT insertion in the first intron of *LRP6* that is predicted to disrupt the *LRP6* coding sequence (we designate this cell line *LRP6*<sup>GT</sup>-3(Int); see Fig 5B and Materials and methods). Both *LRP6*<sup>GT</sup>-1(Up) and *LRP6*<sup>GT</sup>-2(Up) cells demonstrated significantly reduced WNT reporter activation and *AXIN2* mRNA accumulation following treatment with WNT3A when compared to WT HAP1-7TGP cells (Fig 5C and 5D). The most plausible explanation for how the GT insertions reduced WNT signaling would be down-regulation of *LRP6*, which is indeed what we observed when we measured total and cell-surface levels of *LRP6* protein. *LRP6*<sup>GT</sup>-1(Up) and *LRP6*<sup>GT</sup>-2(Up) cells exhibited a 75–84% reduction in total *LRP6* protein and a 68–71% reduction in cell-surface *LRP6* compared to WT cells (Fig 5E and 5F). *LRP6*<sup>GT</sup>-3(Int) cells exhibited >99% and 94% reductions in total and cell-surface *LRP6* respectively, compared to WT cells, as would be expected from the disruption of the *LRP6* coding sequence caused by the sense GT insertion in the first intron (Fig 5E and 5F).

Unexpectedly, despite the reduction in *LRP6* protein observed in *LRP6*<sup>GT</sup>-1(Up) and *LRP6*<sup>GT</sup>-2(Up) cells harboring antisense GT insertions upstream of the *LRP6* promoter, we did not observe a corresponding decrease in *LRP6* mRNA (Fig 5G). In an important control, *LRP6* mRNA levels were indeed markedly reduced in *LRP6*<sup>GT</sup>-3(Int) cells carrying a sense intronic GT insertion that disrupts the coding sequence (Fig 5G).

We considered the possibility that the antisense GT insertions interfere with an unannotated TSS located upstream of the annotated TSS for *LRP6*. Cap analysis of gene expression (CAGE) has been demonstrated to be the best genome-scale method to identify TSSs from hundreds of human tissues [12]. However, the FANTOM5 database [11], which aggregates CAGE data from hundreds of human tissues, did not reveal significant reads upstream of the annotated TSS for *LRP6* (Fig 5B), indicating that the antisense GT insertions are unlikely to disrupt an upstream TSS for *LRP6*. In summary, these results suggest that antisense GT insertions upstream of *LRP6* diminished signaling by an enigmatic mechanism that reduced *LRP6* protein levels without altering mRNA levels, rather than by simply disrupting the *LRP6* promoter. Interestingly, sequence elements with similar properties have been described upstream of promoter elements for heat shock target genes in yeast [13].

## Discussion

We developed a new bioinformatics tool to analyze haploid genetic screens with the explicit goal of uncovering regulatory elements. We analyzed screen data in a way that discerned GT insertion patterns distinct from those predicted to disrupt the coding sequence of genes, and found that atypical insertions in introns and regions upstream of the TSS can cause down-regulation of genes, leading to the phenotype selected for during the screens. When we applied this new analysis to haploid genetic screens interrogating the WNT pathway, we found that antisense GT insertions in the first intron of *TFAP4* and upstream of the *LRP6* promoter resulted in marked changes in the expression of these genes. These types of insertions had not been accounted for in previous analyses of haploid genetic screens.

The identified GT insertions could disrupt regulatory elements such as promoters, enhancers, antisense transcripts or splicing sequences. In the case of *TFAP4*, most of the insertions were located in the first intron and overlapped with a strong enhancer signal (Fig 4A), suggesting they may disrupt an enhancer. We confirmed that this region contains an enhancer by showing that it could increase transcription of a reporter gene when transplanted in front of a minimal promoter (Fig 4F and 4G). Previous studies have shown that *TFAP4* is directly regulated by c-MYC and that the first intron of *TFAP4* in fact contains four c-MYC binding sites [14, 15], three of which are encompassed by the bin identified in our antisense intronic insertion enrichment analysis (Fig 3B and 3C). In future studies, it will be important to test whether the antisense insertions found in the first intron of *TFAP4* down-regulate *TFAP4* mRNA (Fig 4D) by disrupting c-MYC binding or through an alternative mechanism.

Similarly, LRP6 protein was down-regulated in the LRP6<sup>GT</sup>-1(Up) and LRP6<sup>GT</sup>-2(Up) cell lines containing antisense GT insertions upstream of the *LRP6* promoter (Fig 5E and 5F). Surprisingly, *LRP6* mRNA levels were not reduced in these same cell lines (Fig 5G), suggesting a mechanism regulating LRP6 protein levels. In yeast, genomic sequences upstream of genes that have no effect on mRNA levels can instead regulate protein levels, perhaps by regulating mRNA translation or mRNA localization, although the precise mechanisms remain unknown [13].

The selective enrichment of antisense versus sense GT insertions in the region upstream of the *LRP6* promoter in cells sorted for low WNT reporter fluorescence (Fig 5A and 5B) suggests that such insertions are not merely disrupting an enhancer or promoter. A prior study narrowed down the location of the minimal promoter for *LRP6* to the region 242 to 352 bp upstream of the annotated TSS [16], whereas the region enriched for antisense GT insertions that we identified is located 1124 to 2123 bp upstream of the annotated TSS, supporting our conclusion that these antisense GT insertions do not disrupt the activity of the *LRP6* promoter. Furthermore, that same study observed no change in *LRP6* transcription when the region from 352 to 2523 bp upstream of the transcriptional start site was deleted. This deleted region overlaps with the upstream region enriched for antisense GT insertions in our study. Thus, we speculate that these GT insertions may disrupt an antisense transcript or another directional element residing on the antisense strand that positively regulates *LRP6* expression. Since no such elements have been described, it will be important to elucidate the nature of this regulatory mechanism in future studies.

Unlike other more focused efforts to identify regulatory regions associated with a given gene or set of genes [17–23], our untargeted approach enables the genome-wide identification of cis-regulatory elements involved in any phenotype that can be probed through a haploid genetic screen. Identification of such elements may not be feasible with RNA interference-based screens because they require that the target genomic sequences be transcribed. CRISPR-based technologies to screen for regulatory modules on a genome scale are still limited by the

focused mutagenesis or transcriptional modulation of predetermined sequences in the genome [24–27]. However, focused CRISPR-based approaches would be powerful tools to precisely delineate any regulatory element found through our analysis.

While we found new regulatory elements in two central regulators of WNT signaling, there are a few reasons why our current study may be under-powered to comprehensively detect all regulatory elements in the genome affecting the WNT pathway. First, we used deep sequencing datasets from previous screens [3] that were designed to uncover protein coding genes involved in WNT signaling. The sequencing depth used to map insertions in these screens was sufficient to saturate the protein-coding genome, but is probably insufficient to interrogate the much larger non-coding genome. Second, the propensity of the retroviral-based mutagen used in the screens analyzed in this study to insert near TSSs, promoters, and enhancers limited our search for regulatory elements to regions within and adjacent to genes. Our methodology could in principle be extended to identify regulatory elements located anywhere in the genome by using agents that integrate in a truly unbiased manner and then exhaustively mapping insertions in both the selected and unselected cell populations by sequencing at greater depth. Finally, because we assigned bins disregarding gene boundaries, our analysis may have missed regulatory elements in bins that overlapped with both an exon and an intron (such bins would have been excluded from the antisense intronic insertion enrichment analysis), and elements in bins that overlapped with features spanning regions located both upstream and downstream of the TSS (such bins would have been excluded from the upstream insertion enrichment analysis). Reducing the size of the bins could ameliorate this problem, but at the expense of statistical power to determine the significance of GT insertion enrichment due to a reduction in GT insertions per bin and an increase in the multiple testing correction for a larger number of bins. Alternatively, computing GT insertions in intervals defined by the boundaries of genetic features such as introns or promoters (rather than bins of a predetermined size) could also help this issue, but would limit the analysis to annotated regions of the genome.

The analysis described in this work provides an untargeted and genome-scale method to identify both genes and regulatory elements involved in any biological process that can be probed by a haploid genetic screen. We hope that this bioinformatics analysis, available through Github (<https://github.com/RohatgiLab/BAIMS-Pipeline>), empowers other researchers to extract new insights about gene regulation from the growing body of insertional mutagenesis screen data.

## Supporting information

**S1 Fig. HAP1-7TGP reporter cell line.** HAP1-7TGP cells harbor an enhanced green fluorescent protein (EGFP) reporter driven by an established WNT-responsive element containing seven TCF/LEF-binding sites upstream of a minimal promoter.  
(PDF)

**S2 Fig. GT insertion patterns found in *TFAP4* and *DOT1L* in the WNT positive regulator low stringency and high stringency screens.** (A) The histogram indicates the number and orientation of insertions mapped to *TFAP4* in the sorted cell populations from the WNT positive regulator low stringency and high stringency screens. See legend to Fig 4A for details. (B) The histogram indicates the number and orientation of insertions mapped to *DOT1L* (Chromosome 19, 2163750–2232749 bp) in unsorted cells and in the sorted cell populations from the WNT positive regulator low stringency and high stringency screens. The pattern of GT insertions seen in *DOT1L*, predominantly enriched for sense insertions in the first intron, differs from the observed enrichment for both sense and antisense insertions seen in the first intron of *TFAP4*.  
(PDF)



**S3 Fig. GT insertion patterns found in *LRP6* in the WNT positive regulator low stringency and high stringency screens, and immunoblot analysis of *LRP6*.** (A) The histogram indicates the number and orientation of insertions mapped to *LRP6* and to the region ~12.5 kbp upstream of the TSS in the sorted cell populations from the WNT positive regulator low stringency and high stringency screens. See legend to Fig 5A for details. (B) The histogram shows an expanded view of the 5' end of *LRP6* and the region ~12.5 kbp upstream of the TSS (left of the vertical dotted line), with traces for GT insertions mapped in the sorted cell populations from the WNT positive regulator low stringency and high stringency screens. See legend to Fig 5B for details. (C) Immunoblot analysis of *LRP6*. The top and bottom parts of the same membrane were probed for *LRP6* and *ACTIN* (loading control), respectively. The cell lines from which the samples were prepared and loaded in duplicate are indicated above the blots. Molecular weight standards in kDa are indicated to the left of each panel. (PDF)

**S1 File. BAIMS output.** Ranked lists of bins from the bin-based analyses. Each sheet of the Excel file contains a ranked list of bins determined by either the antisense intronic, upstream, or inactivating insertion enrichment analysis applied to either the WNT positive regulator low stringency or high stringency screen. The screen and type of bin-based analysis is indicated at the top of every sheet. The location of each bin in the human genome, the genes overlapping with the bin, and the FDR-corrected *p*-values generated by the bin-based analysis are specified. For each bin, the number of antisense intronic insertions, upstream insertions (sum of sense and antisense insertions), or inactivating insertions (sum of sense and antisense insertions for bins overlapping with a 5'UTR, CDS, or 3'UTR, or sense insertions only for bins overlapping exclusively with an intron) found within the bin in unsorted (control) and sorted cells are also indicated. The total number of insertions mapped in the unsorted cells and in the sorted cells are also shown. (XLSX)

**S2 File. List of clonal cell lines containing GT insertions.** The genomic sequences flanking GT insertion sites in the clonal cell lines used in this study are described. The first column ("Clone Name") indicates the names of the clonal cell lines and the second column ("Genomic sequence flanking GT") indicates the genomic sequences 5' and 3' from the GT insertion site (relative to the sense orientation of the GT as described in Fig 1). (XLSX)

## Acknowledgments

We thank members of the Rohatgi, Salzman and Carette labs for input on the project.

## Author Contributions

**Conceptualization:** Bhaven B. Patel, Andres M. Lebensohn, Rajat Rohatgi.

**Data curation:** Bhaven B. Patel, Ganesh V. Pusapati.

**Formal analysis:** Bhaven B. Patel, Andres M. Lebensohn, Ganesh V. Pusapati, Jan E. Carette, Julia Salzman, Rajat Rohatgi.

**Funding acquisition:** Jan E. Carette, Julia Salzman, Rajat Rohatgi.

**Investigation:** Bhaven B. Patel, Andres M. Lebensohn, Ganesh V. Pusapati.

**Methodology:** Bhaven B. Patel, Andres M. Lebensohn, Jan E. Carette, Julia Salzman, Rajat Rohatgi.

**Project administration:** Rajat Rohatgi.

**Software:** Bhaven B. Patel, Julia Salzman.

**Supervision:** Andres M. Lebensohn, Julia Salzman, Rajat Rohatgi.

**Validation:** Bhaven B. Patel, Andres M. Lebensohn.

**Visualization:** Bhaven B. Patel, Andres M. Lebensohn.

**Writing – original draft:** Bhaven B. Patel, Andres M. Lebensohn, Rajat Rohatgi.

**Writing – review & editing:** Bhaven B. Patel, Andres M. Lebensohn, Rajat Rohatgi.

## References

1. Consortium EP. An integrated encyclopedia of DNA elements in the human genome. *Nature*. 2012; 489(7414):57–74. <https://doi.org/10.1038/nature11247> PMID: 22955616
2. Wright JB, Sanjana NE. CRISPR Screens to Discover Functional Noncoding Elements. *Trends Genet*. 2016; 32(9):526–9. <https://doi.org/10.1016/j.tig.2016.06.004> PMID: 27423542
3. Lebensohn AM, Dubey R, Neitzel LR, Tacchelly-Benites O, Yang E, Marceau CD, et al. Comparative genetic screens in human cells reveal new regulatory mechanisms in WNT signaling. *Elife*. 2016; 5.
4. Vrljicak P, Tao S, Varshney GK, Quach HN, Joshi A, LaFave MC, et al. Genome-Wide Analysis of Transposon and Retroviral Insertions Reveals Preferential Integrations in Regions of DNA Flexibility. *G3 (Bethesda)*. 2016; 6(4):805–17.
5. Egawa T, Littman DR. Transcription factor AP4 modulates reversible and epigenetic silencing of the Cd4 gene. *Proc Natl Acad Sci U S A*. 2011; 108(36):14873–8. <https://doi.org/10.1073/pnas.1112293108> PMID: 21873191
6. Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol*. 2009; 10(3):R25. <https://doi.org/10.1186/gb-2009-10-3-r25> PMID: 19261174
7. Karolchik D, Hinrichs AS, Furey TS, Roskin KM, Sugnet CW, Haussler D, et al. The UCSC Table Browser data retrieval tool. *Nucleic Acids Res*. 2004; 32(Database issue):D493–6. <https://doi.org/10.1093/nar/gkh103> PMID: 14681465
8. Carette JE, Guimaraes CP, Wuethrich I, Blomen VA, Varadarajan M, Sun C, et al. Global gene disruption in human cells to assign genes to phenotypes by deep sequencing. *Nat Biotechnol*. 2011; 29(6):542–6. <https://doi.org/10.1038/nbt.1857> PMID: 21623355
9. Fuerer C, Nusse R. Lentiviral vectors to probe and manipulate the Wnt signaling pathway. *PLoS One*. 2010; 5(2):e9370. <https://doi.org/10.1371/journal.pone.0009370> PMID: 20186325
10. den Dunnen JT, Dalgleish R, Maglott DR, Hart RK, Greenblatt MS, McGowan-Jordan J, et al. HGVS Recommendations for the Description of Sequence Variants: 2016 Update. *Hum Mutat*. 2016; 37(6):564–9. <https://doi.org/10.1002/humu.22981> PMID: 26931183
11. Lizio M, Harshbarger J, Shimoji H, Severin J, Kasukawa T, Sahin S, et al. Gateways to the FANTOM5 promoter level mammalian expression atlas. *Genome Biol*. 2015; 16:22. <https://doi.org/10.1186/s13059-014-0560-6> PMID: 25723102
12. Adiconis X, Haber AL, Simmons SK, Levy Moonshine A, Ji Z, Busby MA, et al. Comprehensive comparative analysis of 5'-end RNA-sequencing methods. *Nat Methods*. 2018; 15(7):505–11. <https://doi.org/10.1038/s41592-018-0014-2> PMID: 29867192
13. Zid BM, O'Shea EK. Promoter sequences direct cytoplasmic localization and translation of mRNAs during starvation in yeast. *Nature*. 2014; 514(7520):117–21. <https://doi.org/10.1038/nature13578> PMID: 25119046
14. Jung P, Hermeking H. The c-MYC-AP4-p21 cascade. *Cell Cycle*. 2009; 8(7):982–9. <https://doi.org/10.4161/cc.8.7.7949> PMID: 19270520
15. Jung P, Menssen A, Mayr D, Hermeking H. AP4 encodes a c-MYC-inducible repressor of p21. *Proc Natl Acad Sci U S A*. 2008; 105(39):15046–51. <https://doi.org/10.1073/pnas.0801773105> PMID: 18818310
16. Pellicelli M, Hariri H, Miller JA, St-Arnaud R. Lrp6 is a target of the PTH-activated alphaNAC transcriptional coregulator. *Biochim Biophys Acta*. 2018; 1861(2):61–71.

17. Fulco CP, Munschauer M, Anyoha R, Munson G, Grossman SR, Perez EM, et al. Systematic mapping of functional enhancer-promoter connections with CRISPR interference. *Science*. 2016; 354(6313):769–73. <https://doi.org/10.1126/science.aag2445> PMID: 27708057
18. Korkmaz G, Lopes R, Ugalde AP, Nevedomskaya E, Han R, Myacheva K, et al. Functional genetic screens for enhancer elements in the human genome using CRISPR-Cas9. *Nat Biotechnol*. 2016; 34(2):192–8. <https://doi.org/10.1038/nbt.3450> PMID: 26751173
19. Smeenk L, van Heeringen SJ, Koeppel M, van Driel MA, Bartels SJ, Akkers RC, et al. Characterization of genome-wide p53-binding sites upon stress response. *Nucleic Acids Res*. 2008; 36(11):3639–54. <https://doi.org/10.1093/nar/gkn232> PMID: 18474530
20. Canver MC, Lessard S, Pinello L, Wu Y, Ilboudo Y, Stern EN, et al. Variant-aware saturating mutagenesis using multiple Cas9 nucleases identifies regulatory elements at trait-associated loci. *Nat Genet*. 2017.
21. Canver MC, Smith EC, Sher F, Pinello L, Sanjana NE, Shalem O, et al. BCL11A enhancer dissection by Cas9-mediated in situ saturating mutagenesis. *Nature*. 2015; 527(7577):192–7. <https://doi.org/10.1038/nature15521> PMID: 26375006
22. Rajagopal N, Srinivasan S, Kooshesh K, Guo Y, Edwards MD, Banerjee B, et al. High-throughput mapping of regulatory DNA. *Nat Biotechnol*. 2016; 34(2):167–74. <https://doi.org/10.1038/nbt.3468> PMID: 26807528
23. Sanjana NE, Wright J, Zheng K, Shalem O, Fontanillas P, Joung J, et al. High-resolution interrogation of functional elements in the noncoding genome. *Science*. 2016; 353(6307):1545–9. <https://doi.org/10.1126/science.aaf7613> PMID: 27708104
24. Larson MH, Gilbert LA, Wang X, Lim WA, Weissman JS, Qi LS. CRISPR interference (CRISPRi) for sequence-specific control of gene expression. *Nat Protoc*. 2013; 8(11):2180–96. <https://doi.org/10.1038/nprot.2013.132> PMID: 24136345
25. Qi LS, Larson MH, Gilbert LA, Doudna JA, Weissman JS, Arkin AP, et al. Repurposing CRISPR as an RNA-guided platform for sequence-specific control of gene expression. *Cell*. 2013; 152(5):1173–83. <https://doi.org/10.1016/j.cell.2013.02.022> PMID: 23452860
26. Perez-Pinera P, Kocak DD, Vockley CM, Adler AF, Kabadi AM, Polstein LR, et al. RNA-guided gene activation by CRISPR-Cas9-based transcription factors. *Nat Methods*. 2013; 10(10):973–6. <https://doi.org/10.1038/nmeth.2600> PMID: 23892895
27. Kearns NA, Genga RM, Enuameh MS, Garber M, Wolfe SA, Maehr R. Cas9 effector-mediated regulation of transcription and differentiation in human pluripotent stem cells. *Development*. 2014; 141(1):219–23. <https://doi.org/10.1242/dev.103341> PMID: 24346702