

1 **Title**

2 The Vertebrate Codex Gene Breaking Protein Trap Library For Genomic Discovery and Disease

3 Modeling Applications

4

5 **Authors**

6 Noriko Ichino<sup>1</sup>, MaKayla Serres<sup>1</sup>, Rhianna Urban<sup>1</sup>, Mark Urban<sup>1</sup>, Kyle Schaeftbauer<sup>1</sup>, Lauren Greif<sup>1</sup>,

7 Gaurav K. Varshney<sup>3</sup>, Kimberly J. Skuster<sup>1</sup>, Melissa McNulty<sup>1</sup>, Camden Daby<sup>1</sup>, Ying Wang<sup>4</sup>, Hsin-kai

8 Liao<sup>4</sup>, Suzan El-Rass<sup>5</sup>, Yonghe Ding<sup>1,2</sup>, Weibin Liu<sup>1,2</sup>, Lisa A. Schimmenti<sup>1</sup>, Sridhar Sivasubbu<sup>6</sup>, Darius

9 Balciunas<sup>7</sup>, Matthias Hammerschmidt<sup>8</sup>, Steven A. Farber<sup>9</sup>, Xiao-Yan Wen<sup>5</sup>, Xiaolei Xu<sup>1,2</sup>, Maura

10 McGrail<sup>4</sup>, Jeffrey J. Essner<sup>4</sup>, Shawn Burgess<sup>10</sup>, Karl J. Clark<sup>1\*</sup>, Stephen C. Ekker<sup>1\*</sup>

11 \*Corresponding authors

12

13 **Author contributions**

14 SCE and KJC conceived research; NI, SCE, KJC for basic experimental design including detailed

15 analyses of the collection; MS, CD, MU, SER, YD, RU, MM, JJE, XX, DB, SCE and KJC generated

16 GBT collection, NI, KS, MS, LG, CD, MU, RU, YD, SAF, WL and KJC conducted phenotype

17 screening of GBT mutant lines; NI, KS, MS, LG, CD, MU, RU, YD, WL and KJC conducted molecular

18 biology analyses; GV and SB conducted next generation sequencing; NI, KS, SCE and KJC conducted

19 bioinformatics-based analyses; LAS, NI, KJC and SCE conducted comparative genomics analyses, NI,

20 KS, MS, LG, KJS and SCE wrote the manuscript; SCE, KJC, JJE, XX, MH, SAF, XYW, SB, XX and

21 SL consulted to this research.

22

1 **Funding:** Supported by NIH grants (GM63904, DA14546, DK093399 and HG 006431), Natural  
2 Sciences and Engineering Research Council of Canada, grant RGPIN 05389-14 and the Mayo  
3 Foundation.

4  
5 **Affiliations**

6 1. Department of Biochemistry and Molecular Biology, Mayo Clinic, Minnesota, USA

7 2. Department of Cardiovascular Medicine, Mayo Clinic, Minnesota, USA

8 3. Functional & Chemical Genomics Program, Oklahoma Medical Research Foundation, Oklahoma City,  
9 OK, 73112, USA

10 4. Department of Genetics, Development and Cell Biology, Iowa State University, Iowa, USA

11 5. Zebrafish Centre for Advanced Drug Discovery & Keenan Research Centre for Biomedical Science,  
12 Li Ka Shing Knowledge Institute, St. Michael's Hospital, Toronto, Ontario, Canada.

13 6. CSIR–Institute of Genomics and Integrative Biology (CSIR-IGIB), Delhi, India.

14 7. Department of Biology Temple University, 435 Biology- Life Sciences Building, 1900 North 12th  
15 Street, Philadelphia, PA 19122, USA

16 8. Institute for Developmental Biology, Cologne University, Köln, Germany

17 9. Carnegie Institution, Baltimore, Maryland, USA

18 10. Translational and Functional Genomics Branch, National Human Genome Research Institute,  
19 National Institutes of Health, Bethesda, MD 20892-8004, USA.

20

21

22 **Keywords**

- 1 Gene-break transposon, zebrafish, disease model, human genetic disorders, gene ontology, Lightsheet
- 2 microscopy

1 **Abstract**

2

3 The zebrafish is a powerful model to explore the molecular genetics and expression of the vertebrate  
4 genome. The gene break transposon (GBT) is a unique insertional mutagen that reports the expression of  
5 the tagged member of the proteome while generating Cre-revertible genetic alleles. This 1000+ locus  
6 collection represents novel codex expression data from the illuminated mRFP protein trap, with 36%  
7 and 87% of the cloned lines showcasing to our knowledge the first described expression of these genes  
8 at day 2 and day 4 of development, respectively. Analyses of 183 molecularly characterized loci indicate  
9 a rich mix of genes involved in diverse cellular processes from cell signaling to DNA repair. The  
10 mutagenicity of the GBT cassette is very high as assessed using both forward and reverse genetic  
11 approaches. Sampling over 150 lines for visible phenotypes after 5dpf shows a similar rate of discovery  
12 of embryonic phenotypes as ENU and retroviral mutagenesis. Furthermore, five cloned insertions were  
13 in loci with previously described phenotypes; embryos homozygous for each of the corresponding GBT  
14 alleles displayed strong loss of function phenotypes comparable to published mutants using other  
15 mutagenesis strategies (*ryr1b*, *fras1*, *tnnt2a*, *edar* and *hmcn1*). Using molecular assessment after  
16 positional cloning, to date nearly all alleles cause at least a 99+% knockdown of the tagged gene.  
17 Interestingly, over 35% of the cloned loci represent 68 mutants in zebrafish orthologs of human disease  
18 loci, including nervous, cardiovascular, endocrine, digestive, musculoskeletal, immune and integument  
19 systems. The GBT protein trapping system enabled the construction of a comprehensive protein codex  
20 including novel expression annotation, identifying new functional roles of the vertebrate genome and  
21 generating a diverse collection of potential models of human disease.

22

23

## 1 **Introduction**

2 With the generation of more than 100 sequenced vertebrate genomes (Meadows & Lindblad-Toh, 2017),  
3 the current key question is how to determine the role(s) of uncharacterized gene products in specific  
4 biological and pathological processes. For example, genes associated with human disease are being  
5 discovered at a rapid rate. However, the biological functions underlying this linkage is often unclear  
6 (Kettleborough et al., 2013). Model system science using loss of function approaches has been essential  
7 to the annotation of the genome to date including the discovery of novel processes and the biological  
8 mechanisms underlying disease (Stoeger, Gerlach, Morimoto, & Nunes Amaral, 2018).

9 Among vertebrates, *Danio rerio* (zebrafish) has emerged as an outstanding model organism  
10 amenable to both forward and reverse genetic approaches. In addition, the natural transparency of the  
11 zebrafish embryo and larvae enables the unprecedented ability to non-invasively collect a rich set of  
12 expression data for the proteome and in the context of an entire living vertebrate. We describe here a  
13 1000+ collection of zebrafish lines made using the Protein Trap Gene- Breaking Transposon  
14 (GBT;(Clark, Balciunas, et al., 2011)to develop such a codex for the comparative vertebrate genomics  
15 field (Meadows & Lindblad-Toh, 2017), (Clark, Balciunas, et al., 2011).

16 The initial pGBT-RP 2.1 (RP2.1) vector has several features that efficiency cooperate to report  
17 gene sequence, expression and function (Clark, Balciunas, et al., 2011). Two main reporter components  
18 include a 5' protein trap and a 3' exon trap, with the entire cassette flanked by inverted terminal repeats  
19 (ITR) of the mini*Tol2* transposon to effectively deliver the transgene as single copy integrations into the  
20 zebrafish genome. In cases where RP2 integrates in the sense orientation of a transcription unit, the  
21 protein trap's splice acceptor overrides normal splicing of the transcription unit, creating a fusion  
22 between endogenous upstream exons and the monomeric RFP (mRFP) reporter sequences. The protein-  
23 trap domain in RP2.1 generates the expression profile, including subsequent protein localization and

1 accumulation when a functional in-frame fusion between the start codon-deficient mRFP reporter and  
2 the tagged protein. Mutagenesis is accomplished by the strong internal polyadenylation and putative  
3 border element, effectively truncating the endogenously tagged protein. The GBT mutagenesis system  
4 represented the first step toward a ‘codex’ of protein expression and functional annotation of the  
5 vertebrate genome (Clark, Balciunas, et al., 2011).

6 We report here the development of a series of GBT protein trap vectors including versions to trap  
7 expression in each of the three potential reading frames. In addition, we modified the 3’ exon trap to use  
8 a localized BFP rather than the more commonly used GFP to more effectively use these lines in  
9 conjunction with other transgenic fish. We deployed these vectors at scale, generating over 1000 protein  
10 trap lines with visible mRFP expression at either 2dpf (end of embryogenesis) or 4dpf (larval stage),  
11 with 36% and 87% of the cloned lines showcasing to our knowledge the first described expression of  
12 these genes at these stages, respectively. We used forward and reverse genetic tests to assess the  
13 mutagenicity of these vectors, noting similar rates of visible phenotypes at 5dpf as ENU and retroviral  
14 screening tools. We re-isolated five previously described loci, and embryos homozygous for each of the  
15 corresponding GBT alleles displayed strong loss of function phenotypes comparable to these previously  
16 published mutants generated using other mutagenesis strategies (*ryr1b*, *fras1*, *tnnt2a*, *edar* and *hmcn1*).  
17 Molecular assessment after positional cloning shows that nearly all alleles cause at least a 99+%  
18 knockdown of the tagged gene. Interestingly, over 35% of the cloned loci represent 68 mutants in  
19 zebrafish orthologs of human disease loci, including nervous, cardiovascular, endocrine, digestive,  
20 musculoskeletal, immune and integument systems. The GBT protein trapping system enabled the  
21 construction of a comprehensive protein codex including novel expression annotation, identifying new  
22 functional roles of the vertebrate genome and generating a diverse collection of potential models of  
23 human disease.

## 1 **Materials and Methods**

### 2 **Zebrafish husbandry**

3 All zebrafish (*Danio rerio*) was maintained according to the guidelines and the standard procedures  
4 established by the Mayo Clinic Institutional Animal Care and Use Committee (Mayo IACUC). The  
5 Mayo IACUC approved all protocols involving live vertebrate animals (A23107, A21710 and A34513).

### 6 **Generating GBT constructs, RP2 and RP8 series**

7 pGBT-RP8.2 and -RP8.3 were made by combining three restriction endonuclease fragments of pGBT-  
8 RP8.1, a 2.2 kb AflIII to AgeI, a 0.7 kb EcoRI to SpeI, and a 3.0kb SpeI to AflIII, with a short adapter to  
9 close the space between AgeI and EcoRI that effectively removed one or two thymine nucleotides just  
10 following the splice acceptor prior to the AUG-less mRFP cassette. For pGBT-RP8.2, Adapter-  
11 GBT(+2) was made by annealing oligos adapter-GBT(+2)-a  
12 [CCGGTTTTCTCATTCATTTACAGTCAGCCGG] and adapter-GBT (+2)-b  
13 [AATTCCGGCTGACTGTAAATGAATGAGAAAA]. For pGBT-RP8.3, Adapter-GBT(+3) was made  
14 by annealing oligos adapter-GBT (+3)-a [CCGGTTTTCTCATTCATTTACAGCAGCCGG] and  
15 adapter-GBT(+3)-b [AATTCCGGCTGCTGTAAATGAATGAGAAAA].

16  
17 pGBT-RP2.2 and -RP2.3 were made by combining three restriction endonuclease fragments of pGBT-  
18 RP2.1 (Clark, Balciunas, et al., 2011), a 3.6kb BlnI to AgeI, a 1.9kb EcoRI to AvrII, and a 3.55kb AvrII  
19 to BlnI, with a short adapter to close the space between AgeI and EcoRI that effectively removed one or  
20 two thymine nucleotides just following the splice acceptor prior to the AUG-less mRFP cassette. For  
21 pGBT-RP2.2, Adapter-GBT(+2) was made by annealing oligos adapter-GBT(+2)-a and adapter-GBT  
22 (+2)-b . For pGBT-RP2.3, Adapter-GBT(+3) was made by annealing oligos adapter-GBT (+3)-a and  
23 adapter-GBT(+3)-b.

1

2 pGBT-RP8.1 was made by cloning a mini-intron derived from carp beta actin intron 1 into pGBT-RP7.1.

3 The 234bp SalI to XhoI mini-intron fragment was isolated from pCR4-bactmIntron following digestion.

4 The pGBT-RP7.1 plasmid was digested with XhoI so that the SalI to XhoI fragment was cloned between  
5 the gamma-crystallin promoter and nls tagBFP.

6

7 pCR4-bactmIntron was made by removing a 1.1kb internal portion of the carp beta actin intron 1 by  
8 digestion of pCR4-bact\_I1 with BstBI and BssHII, followed by filling in 5' overhangs and ligating  
9 remaining vector fragment.

10

11 pCR4-bact\_I1 was cloning a PCR product containing the carp beta-actin intron into pCR4-TOPO  
12 (Invitrogen). The intron was amplified from pGBT-RP2.1 (Clark, Balciunas, et al., 2011) using MISC-  
13 bact\_exon-F1 [CAGCTAGTGCGGAATATCATCTGCC] and MISC-bact\_intron-R1  
14 [CTTCTCGAGGTGAATTCCGGCTGAACTGTA] primers.

15

16 pGBT-RP7.1 was made by replacing a 501bp PstI to PstI fragment of pGBT-RP6.1 with a 480bp PstI to  
17 PstI fragment of pRP2.1. This changed the nucleotide sequence between the carp beta-actin splice  
18 acceptor to replicate the sequences in pGBT-RP2.1. pGBT-RP7.1 was never directly tested in zebrafish.

19

20 pGBT-RP6.1 was made by flipping the internal trap cassette relative the Tol2 inverted terminal repeats  
21 in pGBT-RP5.1. To do this, pGBT-RP5.1 was cut with EcoRV and SmaI. The 2.27kb EcoRV to SmaI  
22 vector backbone fragment, which included the ITRs, was ligated to the 3.51kb EcoRV to SmaI trap



1 fragment. pGBT-RP6.1 was then selected based on the right ITR of Tol2 being in front of the RFP trap,  
2 which is the same orientation of pGBT-RP2.1.

3

4 pGBT-RP5.1 was made by cloning a PCR product with the AUG-less mRFP into pre(-1)GBT-RP5.1.  
5 The 698bp mRFP\* PCR product was obtained by amplification of pGBT-R15 (Clark, Balciunas, et al.,  
6 2011) with CDS-mRFP\*-F1 [AAGAATTCGAAGGTGCCTCCTCCGAGGATGTCATCAAGG] and  
7 CDS-mRFP-R1 [AAACTAGTCTTAGGCTCCGGTGGAGTGGCGG]. Prior to cloning the PCR  
8 mRFP\* product was digested with EcoRI and SpeI to prepare the ends for subcloning into pre(-1)GBT-  
9 RP5.1 that was opened between the carp beta actin splice acceptor and the ocean pout terminator.

10

11 pre(-1)GBT-RP5.1 was made by cloning 1.2kb SpeI to AvrII fragment from pGBT-PX (Sivasubbu et al.,  
12 2006) that contained the ocean pout terminator into the SpeI site of pre(-2)GBT-RP5.1. The resulting  
13 products were screened for the proper orientation of the ocean pout terminator relative to the carp beta  
14 actin splice acceptor.

15

16 pre(-2)GBT-RP5.1 was made by inserting an expression cassette to make a 3' poly(A) trap that makes  
17 blue lenses. A 1.15kb SpeI to BglIII fragment from pKTol2gC-nlsTagBFP was cloned into pre(-3)GBT-  
18 RP5.1 that had been cut with AvrII and BglIII. This moved the *Xenopus* gamma crystallin promoter  
19 driving a nuclear-localized TagBFP in front of the carp beta actin splice donor within pre(-3)GBT-RP5.1  
20 to create a localized BFP poly(A) trap signal replacing the ubiquitous GFP signal that was in pGBT-  
21 RP2.1.

22

1 pre(-3)GBT-RP5.1 was made by cloning a 492bp XmaI to NheI scaffold fragment from pUC57-I-  
2 SceI\_loxP\_splice into pKTol2-SE (Clark, Balciunas, et al., 2011) opened with XmaI and NheI.  
3 pUC57-I-SceI\_LoxP\_Splice contains a synthetic sequence (see below) cloned into pUC57 (Genscript).  
4 The scaffold contains an I-SceI site; loxP site; carp beta actin splice acceptor; cloning sites for mRFP,  
5 ocean pout terminator, and BFP lens cassettes; carp beta actin splice donor; loxP site; and an I-SceI site.  
6 [cccgggatagggataacagggtaataataacttcgtatagcatacattatacgaagtatcgttaccaccactagcggtcagactgcagattgcagcac  
7 gaaacaggaagctgactccacatggtcacatgctcactgaagtgtgacttcctgacagctgtgcactttctaaaccggttttctcattcatttacagttca  
8 gcctgttacctgactcaccgacaagctgttaccctggaattcgtttaaacactagtcaccggcgttctaggtataagatctacctaaggtgagttgatct  
9 ttaagcttttacatttcagctcgcataatcaattcgaacgtttaattagaatgtttaataaagctagattaaatgattaggctcagttaccggctttttttct  
10 catttactgagctcaagacgtctgataacttcgtatagcatacattatacgaagttattaccctgttatccctatggctagc]

11

## 12 **Generating GBT collection**

13 Generation of the GBT collection was based on the prior described protocols (Clark, Balciunas, et al.,  
14 2011; Clark, Urban, Skuster, & Ekker, 2011; J. Ni et al., 2016).

15

## 16 **Fluorescent microscopy of mRFP reporter protein expression**

17 Larvae were treated with 0.2 mM phenylthiocarbamide at 1 dpf to inhibit pigment formation. The  
18 anesthetized fish were mounted in 1.5% agarose (Fisher Scientific BP1360) prepared with 0.017mg/ml  
19 tricaine solution in an agarose column in the imaging chamber. The protocol of ApoTome microscopy  
20 was described in previous publication. (Clark, Balciunas, et al., 2011) For Lightsheet microscopy, larval  
21 zebrafish were anesthetized with 0.017g/ml tricaine (Ethyl 3-aminobenzoate methanesulfonate salt) in  
22 embryo water during imaging procedure. To capture RFP expression patterns of 2 dpf and 4 dpf larval  
23 zebrafish, LP 560 nm filter as excitation and LP 585nm as emission was used for Lightsheet microscopy.

1 The sagittal-, dorsal-, and ventral- oriented z-stacks of the mRFP expression were captured at either 50x  
2 magnification using an ApoTome microscope (Zeiss) with a 5x/0.25 NA dry objective (Zeiss) or 50x  
3 magnification using a Lightsheet Z.1 microscope (Zeiss) 5x/0.16 NA dry objective. Each set of images  
4 were obtained from the same larva and the images shown are composites of the maximum image  
5 projections of the z-stacks obtained from each direction.

6

### 7 **Sperm Cryopreservation**

8 Sperm collection and cryopreservation was initially based on the protocol described in (Draper & Moens,  
9 2009) and moved to the ZIRC protocol described in (Matthews et al., 2018).

10

### 11 **Genomic DNA isolation**

12 Genomic DNA was isolated from F1 fish tail biopsies to conduct next generation sequencing and from  
13 both WT and heterozygous larva to manually perform the PCR-based mRFP linkage analysis. Zebrafish  
14 larva were individually placed to 0.2 ml PCR tubes and sacrificed to extract genomic DNA in 50 mM  
15 NaOH for 20 min at 95 C°.

16

### 17 **PCR-based linkage analysis of GBT insertions loci**

18 TALE-PCR: The protocol used was designed to amplify and clone junction fragments from Tol2-based  
19 gene-break transposons (GBT) in the zebrafish genome. Although modified, it is based on a protocols  
20 received from Alexi Parnov, Vladimir Korsch, and Karuna Sampath. The following primer mixtures  
21 (containing 0.4  $\mu$ M GBT specific primer and 2  $\mu$ M DP primer) were prepared: for primary PCR: 5R-  
22 mRFP-P1/DP1, 5R-mRFP-P1/DP2, 5R-mRFP-P1/DP3, 5R-mRFP-P1/DP4, 3R-GM2-P1/DP1, 3R-  
23 GM2-P1/DP2, 3R-GM2-P1/DP3, 3R-GM2-P1/DP4, 3R-tagBFP-P1/DP1, 3R-tagBFP-P1/DP2, 3R-

1 tagBFP-P1/DP3, 3R-tagBFP-P1/DP4 ; for secondary PCR: 5R-mRFP-P2/DP1, 5R-mRFP-P2/DP2, 5R-  
2 mRFP-P2/DP3, 5R-mRFP-P2/DP4, 3R-GM2-P2/DP1, 3R-GM2-P2/DP2, 3R-GM2-P2/DP3, 3R-GM2-  
3 P2/DP4, 3R-tagBFP-P2/DP1, 3R-tagBFP-P2/DP2, 3R-tagBFP-P2/DP3, 3R-tagBFP-P2/DP4; for tertiary  
4 PCR: TAIL-bA-SA/DP1, TAIL-bA-SA/DP2, TAIL-bA-SA/DP3, TAIL-bA-SA/DP4, Tol2-ITR(L)-  
5 O1/DP1, Tol2-ITR(L)-O1/DP2, Tol2-ITR(L)-O1/DP3, Tol2-ITR(L)-O1/DP4, Tol2-ITR(L)-O3/DP1,  
6 Tol2-ITR(L)-O3/DP2, Tol2-ITR(L)-O3/DP3, Tol2-ITR(L)-O3/DP4. A total of 1  $\mu$ l of primer mixtures  
7 were added to PCR reaction (total volume 25  $\mu$ l). Cycle settings were as follows. Primary: (1) 95°C, 3  
8 min; (2) 95°C, 20 sec; (3) 61°C, 30 sec; (4) 70°C, 3 min; (5) go to “cycle 2” 5 times; (6) 95°C, 20 sec;  
9 (7) 25°C, 3 min; (8) ramping 0.3°/sec to 70°C; (9) 70°C, 3 min; (10) 95°C, 20 sec; (11) 61°C, 30 sec;  
10 (12) 70°C, 3 min; (13) 95°C, 20 sec; (14) 61°C, 30 sec; (15) 70°C, 3 min; (16) 95°C, 20 sec; (17) 44°C,  
11 1 min; (18) 70°C, 3 min; (19) go to “cycle 10” 15 times; (20) 70°C, 5 min; Soak at 12 °C . A total of 5  
12  $\mu$ l of the primary reaction was diluted with 95  $\mu$ l of 10mM Tris-Cl or TE buffers and 1 $\mu$ l of the mixture  
13 was added to the secondary reaction. Secondary: (1) 95°C, 2 min (2) 95°C, 20 sec; (3)61°C, 30 sec; (4)  
14 70°C, 3 min; (5) 95°C, 20 sec; (6) 61°C, 30 sec ; (7) 70°C, 3 min; (8) 95°C, 20 sec; (9) 44°C, 1 min;  
15 (10) ramping 1.5°/sec to 70°C; (11) 70°C, 3 min; (12) go to “cycle 2” 15times; (13) 70°C, 5 min; Soak  
16 at 12°C. A total of 5  $\mu$ l of the primary reaction was diluted with 95  $\mu$ l of 10mM Tris-Cl or TE buffers  
17 and 1 $\mu$ l of the mixture was added to the tertiary reaction. Tertiary: (1) 95°C, 2 min; (2) 95°C, 20 sec; (3)  
18 44°C, 1 min; (3) ramping 1.5°/sec to 70°C; (4) 70°C, 3 min; (5) go to “cycle 2” 32 times; (6) 70°C, 5  
19 min; Soak at 12°C. Products of the secondary and tertiary reactions were separated by using 1-1.5%  
20 agarose gel. The individual bands from the “band shift” pairs were cut from the gel and purified by  
21 using QIAquick Gel Extraction Kit (QIAGEN, Germany), and sequenced by using ABI Cycle  
22 Sequencing chemistry (PE Applied Biosystems, CA) and an ABI Prism 310 Genetic Analyzer with Data  
23 Collection Software (PE Applied Biosystems, Foster City, CA) supplied by the producer.

1 5' and 3' RACE PCRs: The protocol used was designed to use cDNA to amplify and clone junction  
2 fragments of Tol2-based gene-break transposons (GBT) in the zebrafish genome. 5' RACE allows PCR  
3 amplification of unknown sequence at the 5' end of a cDNA as long as there is enough known sequence  
4 within the cDNA to design two antisense primers. Although modified, it is based on the protocol in  
5 described in (Clark, Balciunas, et al., 2011). The following primer mixtures were prepared: for primary  
6 PCR: 0.20  $\mu$ M GBT specific primer (5R-mRFP-P1), and a mix of universal 5' RACE primers 2.5  $\mu$ M  
7 5R-UP-S and 0.5  $\mu$ M 5R-UP-L . Secondary reaction: 25  $\mu$ M GSP (5R-mRFP-P2), and 25  $\mu$ M universal  
8 primer 5R-N1. The reaction mix used was as follows: Primary: (25  $\mu$ l reaction) Template (RR-cDNA) 2  
9  $\mu$ l, Biorline 5X myTaq buffer 5  $\mu$ l, myTaq 0.25  $\mu$ l, GSP 0.5  $\mu$ l, Universal 5' RACE primer mix (URS  
10 mix) 2.0  $\mu$ l, Water 15.25  $\mu$ l. Secondary: (50  $\mu$ l reaction) Template (2:100 dilution 1<sup>o</sup> PCR Reaction),  
11 Biorline 5X myTaq buffer 10  $\mu$ l, myTaq 0.3  $\mu$ l, GSP-P2 0.9  $\mu$ l, URS-P2 0.9  $\mu$ l, Water 35.9  $\mu$ l. Cycle  
12 Settings are as follows. Primary: (1) 95°, 3'; (2) 95°, 30"; (3) 65°, 30" -0.5°/cycle; (4) 70°, 2'; Go To (2)  
13 x 15 cycles; (5) 95°, 30"; (6) 57°, 30"; (7) 70°, 2'; Go To (5) x 20 cycles, (8) 70°, 10'; (9) Soak at 12°.  
14 Dilute 2 $\mu$ L of the primary PCR reaction with 198 $\mu$ L of 10mM Tris-Cl; 1mM EDTA pH8.0. Secondary:  
15 (1) 95° for 3' (2) 95°, 30"; (3) 63°, 30" -0.5°/cycle; (4) 70°, 2'; Go To (2) x 10 cycles; (5) 95°, 30"; (6)  
16 58°, 30"; (7) 70°, 2'; Go To (5) x 25 cycles; (8) 70°, 10'; Soak at 12°. After the PCR reactions finish  
17 20 $\mu$ L of each sample were run on a 1.2% agarose gel. You may also run 10 $\mu$ L of undiluted primary  
18 PCR reactions on the same gel; however, most of the time the bands of interest do not all appear until  
19 after the nested PCR reactions have been run. The individual bands from the gel were excised and  
20 purified by using QIAquick Gel Extraction Kit (QIAGEN, Germany), and sequenced by using ABI  
21 Cycle Sequencing chemistry (PE Applied Biosystems, CA) and an ABI Prism 310 Genetic Analyzer  
22 with Data Collection Software (PE Applied Biosystems, Foster City, CA) supplied by the producer.  
23

## 1 **Forward genetic screening with next-generation sequencing**

2 Isolated genomic DNA (300-500ng) was digested with MseI, and BfaI in parallel for 3h at 37°C and  
3 heat inactivated for 10 min at 80°C. The digested samples from each enzyme were pooled with  
4 prealiquoted barcoded linker in individual wells. The T4 DNA ligase (New England Biolabs, Inc.) was  
5 added, and the reaction mix was incubated for 2 h at 16°C. The linker-mediated PCR was performed in  
6 two steps. In the first step, PCR was done with one primer specific to the 3' - ITR (5' -  
7 GACTTGTGGTCTCGCTGTTCCCTTGG-3') and the other primer specific to linker sequences (5' -  
8 GTAATACGACTCACTATAGGGC- 3') using the following conditions: 2 min at 95°C, 25 cycles of 15  
9 sec at 95°C, 30 seconds at 55°C and 30 seconds at 72°C. The PCR products were diluted to 1:50 in  
10 dH<sub>2</sub>O, and a second round of PCR was performed using ITR (5' -  
11 TCACTTGAGTAAAATTTTTGAGTACTTTTTACACCTC-3') and linker specific (5' -  
12 GCGTGGTTCGACTGCGCAT-3') nested primers to increase sensitivity and avoid non- specific  
13 amplification using the following conditions: 2 min at 95°C, 20 cycles of 15 sec at 95°C, 30 seconds at  
14 58°C and 30 seconds at 72°C. The nested PCR products from each 96-well plate are pooled and  
15 processed for Illumina library preparation as per manufacturer's instructions.

16

## 17 **Protein classification of the cloned zebrafish genes**

18 The 183 cloned zebrafish genes are classified by using PANTHER (Mi, Muruganujan, & Thomas, 2013).  
19 PANTHER provided protein classes of the molecule coded by the cloned zebrafish genes.

20

## 21 **Annotating human orthologues of GBT-tagged genes and disease-causing genes**

22 The human orthologues of 192 cloned zebrafish genes were mainly collected by using a data mining tool,  
23 ZebrafishMine (Van Slyke et al., 2018) supported by the ZFIN database. In some cases, the candidates

1 of human orthologues unlisted in ZFIN database were manually searched by using both Ensembl  
2 (<https://useast.ensembl.org/index.html>) and InParanoid8 (<http://inparanoid.sbc.su.se/cgi-bin/index.cgi>)  
3 databases. In parallel, the candidates were manually identified by the result of BLASTP assembled with  
4 human proteins and by the result of an online synteny analysis tool, SynFind  
5 (<https://genomeevolution.org/CoGe/SynFind.pl>). If the candidate multiply hit in those manual  
6 assessments, it was annotated as a human orthologue. The human phenotype data caused by mutations  
7 of 68 human orthologues were collected by using another data mining tool, BioMart  
8 (<http://useast.ensembl.org/biomart/martview/cfe15ead83199a0b7c7997f5a4ce9e6b>) supported by  
9 Ensembl database.

10

## 11 **Finding Disease Models in Vertebrates**

12 Mouse models were found by using both descriptions of animal models in Online Mendelian Inheritance  
13 in Man (OMIM; <https://www.omim.org/>) and in Mouse Genome Informatics (MGI;  
14 <http://www.informatics.jax.org/humanDisease.shtml>). MGI provided the details of mouse models of  
15 human disease, such as the number of models have been established. Zebrafish model were also found  
16 by using both OMIM and the Zebrafish Information Network (ZFIN; <https://zfin.org/>). ZFIN provided  
17 all data of fish strains listed in this database.

18

## 19 **Gene expression profiling of the cloned zebrafish genes**

20 The cloned genes with unpublished expression data were isolated by using “Gene Expression” tool of  
21 ZebrafishMine (Van Slyke et al., 2018). In parallel, some published expression data were also manually  
22 searched from ZFIN database or in some references. To isolate the gene with the expression localized in  
23 the tissues or the organs shows abnormalities in the causing diseases of the human orthologues, the

- 1 mRFP reporter expression patterns of the cloned genes 2 and 4 dpf are manually analyzed using
- 2 zfishbook database (Clark, Argue, Petzold, & Ekker, 2012).
- 3



## 1 **Results**

### 2 **The features of GBT constructs RP2 and RP8 – capturing all three proteomic reading frames**

3 In our previous study, we reported the intronic-based gene-breaking transposons (GBTs) as effective and  
4 revertible loss-of function tools for zebrafish (Clark, Balciunas, et al., 2011). The main features of the  
5 RP2.1 vector system are as follows (Fig1A): 1) Genetically engineered cargo is flanked by miniTol2  
6 sequences necessary and sufficient for Tol2 transposase-mediated transposition, an efficient transgenesis  
7 vector in zebrafish (Kawakami et al., 2004); (Balciunas et al., 2006); (Urasaki, Morvan, & Kawakami,  
8 2006); 2) protein trap that enables *in vivo* expression selection of the vertebrate proteome. The addition  
9 of the AUG-free mRFP reporter has yielded an effective protein trap for both organ-specific and  
10 subcellular localization of the tagged locus (Petzold et al., 2009); (Clark, Balciunas, et al., 2011); (Liao  
11 et al., 2012);(Xu et al., 2012);(Ding et al., 2013); (Westcot et al., 2015); 3) Mutagenic transcriptional  
12 terminator. The 5' cassette is a combination of a strong splice-acceptor (SA), poly adenylation signal  
13 (pA), and putative border element (red octagon) in conjunction with a start codon (AUG)-free  
14 monomeric red fluorescent protein (mRFP) reporter. These elements have been shown to be effective as  
15 a transcriptional stop in zebrafish by hijacking endogenous splicing (Sivasubbu et al., 2006). These  
16 elements are very effective at inducing a quantitative knockdown in all 26 lines assessed to date using  
17 qRT-PCR; 97% or higher knockdown in all lines (Clark, Balciunas, et al., 2011); (Ding et al., 2013);  
18 (Ding et al., 2016), this manuscript). The GBT mutagenesis system is thus an effective first step to  
19 creating a gene codex simultaneously combining expression and loss of function genetics.  
20 However, some limitations were noted with the initial RP2.1 vector – notably the effective trapping of  
21 transcripts without detectable expression of the mRFP reporter. Molecular cloning of these GFP+/RFP-  
22 lines demonstrated the fidelity of expression requiring the capture of an appropriate reading frame.  
23 RP2.1 was designed to use one main reading frame, and some lines with expression were noted to

1 include the use of a secondary splice acceptor (data not shown). To maximize genome coverage of this  
2 insertional mutagen, we created all three reading frames of the RP2 and RP8 vector series (Fig. 1). 4) 3'  
3 exon trap. These vectors also encode a 3' exon trap with preferential expression following intragenic  
4 insertions in zebrafish (Sivasubbu et al., 2006); (Petzold et al., 2009); (Clark, Balciunas, et al., 2011).  
5 The function of this cassette complements the obligate, in-frame protein trapping effect and is used for  
6 both quality control during mutagenesis and for genotyping of more weakly expressing protein trap  
7 alleles (Fig. 1A). In the RP2 vector series, the nearly ubiquitous b-actin promoter drives expression of  
8 GFP. Expression of integrated GFP becomes detectable between early developmental stages such as  
9 seven- to eight-somite-stage and provides bright expression with a good signal-to-noise ratio at 25 hpf  
10 (Davidson AE et al., 2003). However, the ubiquitous GFP expression from the 3' exon trap cassettes  
11 could interfere with another fluorescent marker system based on GFP labeling in further studies. The  
12 RP8 vector series includes all reading frames for the AUG-free mRFP reporter and a new 3' exon trap  
13 cassette with expression of tagBFP driven by the lens-specific gamma-crystalline promoter (Fig. 1B).  
14 Using the tissue-specific reporter system with BFP is helpful to easily detect F1 founder with weak  
15 mRFP expression and to avoid interference with GFP-based multi-labeling purposes when crossed with  
16 other GFP-labeled transgenic fish lines. 5) Reversible mutagenic cassette. The flanking loxP sites enable  
17 reversion of the tagged locus by Cre – mediated recombination. This facilitates both somatic (Clark,  
18 Balciunas, et al., 2011); (Ding et al., 2013) and germline approaches (Petzold et al., 2009).  
19 We generated more than eleven hundred independent lines by using all six constructs of the GBT system  
20 (Supplemental Table 1). We conducted an initial screening expression of the mRFP fusion protein and  
21 showed that RP2 and RP8 vector series with all reading frames of mRFP reporter protein readily detects  
22 the distribution of the fusion proteins expressed from their own promoter in zebrafish (Supplemental  
23 Figure 1).

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

## **Annotation of protein localization and trafficking of the GBT strain collection**

The ability to non-invasively obtain temporal and spatial expression pattern information is a key feature of these protein trap strains. Pilot data from our first RFP lines rapidly demonstrated this step was going to be a major bottleneck for our pipeline if we used standard documentation methods. Consequently, we established a capillary-based confinement and imaging protocol (SCORE imaging; (Petzold et al., 2010) for quickly capturing and holding living zebrafish for rapid, high quality fluorescent expression in precise longitudinal imaging angles. The animals are easily inserted into a capillary of the correct diameter for the particular fish stage, then placed on the microscope and covered with a matched solution to remove the optical distortion from the capillary housing. The capillary is simply rotated for precise 0 degree (dorsal), 90 degree (lateral), and 180 degree (ventral) images on standard fluorescent microscopes, such as the Apotome. The technical bottleneck of standard microscopy to scan a whole embryo and larva is also labor-intensive and time-consuming. For instance, fluorescent scanning of a quarter of a 2dpf or 4dpf larvae required about 15 min exposure using a Zeiss Apotome microscope, one imaging modality deployed for this collection. To accelerate the expression profiling of these GBT lines, we subsequently utilized a Zeiss Lightsheet Z.1 SPIM microscope. The Lightsheet enabled high speed scanning of a whole embryo or larva, resulting in a nearly 20x faster image acquisition rate than the primary imaging process utilizing the Apotome. We prioritized and cataloged lines with robust expression at 2 and 4 days post fertilization (dpf). All zebrafish lines are freely available now through zfishbook (Clark et al., 2012) and are partially accessible from the Zebrafish International Resource Center. Consequently, we have generated 1138 lines by using each vector system showing in Supplemental Table 1), and updated results are posted at zfishbook.

1 Our throughput for cloning the GBT lines using traditional molecular methods was clearly an initial  
2 bottleneck. To help address this, we deployed a rapid cloning process based on methods used to isolate  
3 retroviral integrations (Varshney, Huang, et al., 2013; Varshney, Lu, et al., 2013). This method  
4 leverages the massive parallel sequencing technology of the Illumina MiSeq, yielding 101 bp  
5 sequencing reads, followed by a custom bioinformatics pipeline that involves both mapping and  
6 annotation. Fin-clips from four male animals per GBT locus are obtained during sperm cryopreservation  
7 and are used as a source of DNA. Shared inserts in multiple individuals from a single GBT line are  
8 considered candidate loci. This information is subsequently used to generate gene-specific primers to  
9 manually confirm linkage using 5' or 3'RACE, inverse PCR, or TAIL PCR to generate locus-specific  
10 primers for downstream molecular genotyping applications. 212 of the lines met the highest stringency  
11 of confirmed expression linkage and are classified as confirmed GBT integrations. An additional 143  
12 lines have been initially tagged using this high throughput method, yielding candidate integration  
13 annotation as listed in zfishbook. The ability to complete the annotation status from candidate to  
14 confirmed for any given GBT locus with a desired expression profile is enhanced by the continued  
15 refinement of the zebrafish genome.

16

### 17 **mRFP Expression profiling reveals overlap with known annotation at 2dpf and substantive new** 18 **expression data at both 2dpf and 4dpf**

19 GBT lines are currently imaged to capture mRFP expression at both 2 and 4 dpf, including dorsal,  
20 sagittal and ventral views using the SCORE imaging method. (Petzold et al., 2010) In an openly  
21 accessible database of GBT lines, zfishbook, the images of mRFP expression pattern were stored within  
22 the media gallery associated with each line (Clark et al., 2012). We summarized published expression  
23 data of the tagged genes in both zfishbook and ZFIN in Table 1. Imaging the localization of transcripts

1 and proteins at 4dpf is more difficult than those at 2 dpf, because accessibility of antisense RNA probes  
2 and antibodies into the larva's body is technically limited for in the methods of both *in situ* hybridization  
3 and immunohistochemistry. Compared with the published data of gene expression in ZFIN, zfishbook  
4 currently provides almost the double number of genes with expression data at 2 dpf and 14 times the  
5 number of genes with expression data at 4 dpf (Table 1). In addition, zfishbook also provides novel  
6 expression data for 61 genes at any developmental stage (Fig.2).

7

### 8 **High knockdown efficiency of endogenous transcripts induced by RP2**

9 We directly compared published transposon insertional mutant vector systems (Fig. 3). The range and  
10 average knockdown levels in the FlipTrap system (Trinh le et al., 2011) produced a range of 4-30% (70-  
11 96% knockdown capacity) in six tested fish alleles, a similar range to our initial R-series protein trap  
12 vectors (R14-R15) that used a simple transcriptional terminator (Liao et al., 2012; Petzold et al., 2009)  
13 The pFT1 appears to be an improvement over these systems, in which the overall range and average  
14 read-through is reduced to 6-11% (89-94% knockdown) from four tested fish alleles (T. T. Ni et al.,  
15 2012). In contrast, the RP2.1 vector (Clark, Balciunas, et al., 2011), (Ding et al., 2013), (Ding et al.,  
16 2016) and this manuscript) maintains a strong knockdown (1% or less read-through) in 26 lines tested.  
17 Though deployed here using a nearly random, transposon-based delivery platform, the GBT vector  
18 system is an effective insertional mutagen suitable for an array of other – including targeted integration -  
19 genome-wide applications.

20

21 **Phenotypic Appearance Rate is Similar to Other Mutagenic Technologies for Forward Genetics**  
22 **Screening through 5 dpf**

1 We conducted an initial forward genetic screen on embryos and early larvae of 179 RFP-positive GBT  
2 lines, identifying 12 recessive phenotypes, such as *ryr1b*, *fras1*, *tnnt2a*, *edar* and *hmcn1*, (Clark,  
3 Balciunas, et al., 2011; Westcot et al., 2015) visible during the first five days of development including  
4 lethality, heart, muscle, skin and other phenotypes. This 7% recovery of visible early developmental  
5 mutants is very similar to the 5% recovered visible mutants from the Sanger TILLING consortium  
6 analysis of truncated zebrafish genes (Kettleborough et al., 2013). This 7% is also comparable to prior  
7 retroviral (Amsterdam & Hopkins, 2004) and ENU (Haffter et al., 1996) zebrafish mutagenesis work that  
8 estimated between 1400 and 2400 genes (~5-9% of the genome) would result in a visible embryonic  
9 phenotype when mutated.

10

#### 11 **GBT alleles phenocopy known embryonic mutations**

12 We tested the first five GBT lines in genes with known loss of function mutant phenotypes (*ryr1b*;  
13 *fras1*; *tnnt2a*; *edar*; *hmcn1*). All five of these alleles in genes with described loss of function defects are  
14 phenocopied by these GBT insertional alleles (Clark, Balciunas, et al., 2011); (Westcot et al., 2015); this  
15 manuscript). These loci represent a critical internal methods reference further validating the  
16 mutagenicity of these novel insertional vectors.

17

#### 18 **Gene ontology analysis of GBT-tagged loci**

19 To assess the diversity of GBT loci molecularly characterized to date, we utilized the PANTHER  
20 classification system (v.14.0, pantherdb.org, (Mi et al., 2013) and generated a table of protein class  
21 ontology tags in the molecularly isolated GBT lines. The PANTHER Protein Class ontology was  
22 adapted from the PANTHER INDEX (PANTHER/X) ontology that comprises two types of  
23 classifications: molecular function and biological process and includes commonly used classes of

1 protein families. The molecular function schema classifies a protein based on its biochemical properties,  
2 such as receptor, cell adhesion molecule, or kinase. The biological process schema classifies a protein  
3 based on the cellular role or process in which it is involved, for example, carbohydrate metabolism  
4 (cellular role), TCA cycle (pathway), neuronal activities (process), or developmental processes (process)  
5 (Thomas PD et al, 2013). As of April 2018, almost half of known zebrafish genes (10626/25289 genes)  
6 were tagged in the PANTHER Protein Class ontology. 168 of our cloned GBT alleles mapped in the  
7 PANTHER system with 21 types of Protein Classes (Table 2). 18% and 16 % of the mapped GBT  
8 alleles are classified to nucleic acid binding (PC00171) and transcription factor (PC00218), respectively  
9 (Fig. 4). This result reveals that a quarter of the mapped genes possibly play a role in regulatory  
10 processes. Overall, however, the rich diversity of protein classes observed in our cloned traps suggests a  
11 large diversity will be represented by the overall collection and consistent with the random nature of  
12 genome integration events by the Tol2 transposon (Clark, Balciunas, et al., 2011).

13

#### 14 **Disease-causing human orthologues of GBT-tagged genes**

15 Of 183 GBT-tagged genes, 171 human orthologues were annotated in at least one public database such  
16 as ZFIN, Ensembl, Homologene, and InParanoid (Table 3). Several human orthologues were  
17 provisionally annotated using BLASTP and a synteny analysis tool, SynFind. In a previous study  
18 comparing the list of human genes possessing at least one zebrafish orthologue with the 3,176 genes  
19 bearing morbidity descriptions that are listed in the OMIM database, 82 % morbid genes (2,601 genes)  
20 can be related to at least one zebrafish orthologue (Howe et al., 2013). Surprisingly, 67 genes (about  
21 37%) of 183 annotated human orthologues are associated with human disease involved in multi-organ  
22 system including nervous, circulatory, endocrine, metabolic, digestive, musculoskeletal, immune, and  
23 integument systems (Fig. 5 and Table 3) and many are not established in rodents and zebrafish (Table

1 5). The GBT protein-trap system provides a variety of potential human disease models which have a  
2 revertible allele that can interchange between disease and healthy cellular, organ and physiological states.

3

4 **The GBT Protein Trap Reveals Protein-Coding Regions Not Predicted by the Zebrafish Genome**  
5 **Project**

6 16 of 211 molecularly confirmed lines by manual, PCR-based mRFP linkage analysis using TALE,  
7 inverse, and RACE PCRs (Table 4) do not match any expressed sequence tag (EST) or predicted genes.  
8 However, in each case we were able to confirm transcription at the locus in wild-type animals, yielding  
9 new annotation for these loci in the zebrafish genome.

10



## 1 **Discussions**

### 2 **Novel Expression Annotation Revealed By Protein Trapping of Endogenous Genes**

3 The RP2 and RP8 vectors of GBT system were assembled to capture all three reading frames of fusion  
4 protein of the trapped gene with the mRFP reporter. This GBT system reveals that mRFP-truncated  
5 fusion proteins exhibit distinct subcellular localization. This is particularly noteworthy for the 4dpf  
6 stages because published gene expression data at late developmental stages have been limited by the  
7 technical difficulties in conducting such analyses using traditional techniques such as whole mount *in*  
8 *situ* hybridization (WISH) at these larval and later time points. We note that the mRFP-truncated fusion  
9 protein may localize ectopically in cases where the protein localization signal is contained in the C-  
10 terminal domain (Clark, Balciunas, et al., 2011; Trinh le & Fraser, 2013). Although the extent that  
11 subcellular localization recapitulates the endogenous protein is dependent on each insertion locus  
12 specifics, the visualizing and illuminating spatiotemporal expression patterns of trapped protein may  
13 facilitate dynamic studies of specific cell types and molecular functions in a living vertebrate. The novel  
14 expression description at 4dpf in nearly five of six cloned loci demonstrates the dearth of 3D expression  
15 annotation for the overwhelming number of genes, even in one of the most studied model system such as  
16 the zebrafish. With ever-improving microscope-based imaging tools (Liu et al., 2018), these lines have  
17 the potential to help annotate at diverse developmental and adult stages, while also potentially imaging  
18 subcellular expression for a subset of protein trap fusions.

19

### 20 **Gene-Break Protein Trap in an Effective Insertional Mutagen with High Knockdown Efficiency** 21 **and Cre-Reversion to WT allele**

1 Transposons offer several unique features over and above traditional static mutational approaches,  
2 including high quality expression tools and new regulated mutagenesis methodologies. From the  
3 perspective of genome engineering development, GBT technology was the first method for revertible  
4 allele generation of vertebrates outside of the mouse (Clark, Balciunas, et al., 2011); (Ding et al., 2013).  
5 We know two major potential biases that may yield non-random trapping coverage of the genome. First,  
6 the RP2.1 protein trap was initially designed around a single reading frame. Upon molecular analysis of  
7 our first lines, however, we discovered that RP2.1 encodes a second, alternative splice acceptor yielding  
8 protein trap expression from a second reading frame due to this alternative splicing event in a significant  
9 number of our lines. The deliberate development of RP2 and RP8 vectors for each reading frame  
10 obviates this potential limitation when used in other delivery modalities besides the Tol2 transposon. For  
11 example, these vectors would be suitable for gene editing-based targeted knockin methods.

12 The GBT system deployed here has been joined by two new and complementary transposon  
13 mutagenesis systems. The FlipTrap system by Dr. Fraser's group (Trinh le et al., 2011), which can be  
14 mutagenic when provided Cre recombinase, is primarily focused on imaging fusion proteins in vivo and  
15 addressing cellular dynamics and related questions. The FT1 system by the Chen lab is a complementary  
16 flipping trap that can use either Cre or Flp recombinases to regulate alleles depend on the original  
17 orientation of the insertion (T. T. Ni et al., 2012). GBT-based zebrafish alleles are highly  
18 complementary and non-redundant to those generated by other mutagenesis methods, including these  
19 other transposons and demonstrated higher knockdown efficiency of the WT transcripts than those  
20 mutagens because of the use of an enhanced polyadenylation signal and a putative boundary element  
21 between 5' protein trap and 3' exon trap cassettes. Importantly, since the initiation of this project to  
22 generate a collection of Cre-revertible mutant alleles, several groups have now reported collections of  
23 tissue-specific Cre driver lines including the Brand lab (Jungke et al., 2013; <http://crezoo.crt->

1 dresden.de/crezoo/), the Zcre consortium (<http://zcre.org.uk/>) and the Wen lab in the PTC Consortium .  
2 The GBT system described here is a two-component, molecularly regulatable mutagenesis approach that  
3 offers the ability to test for the sufficiency of protein-encoding loci in regulated, tissue- and cell- specific  
4 applications.

5

## 6 **Functional Diversity of the Trapped Proteins by the GBT system**

7 To analyze distribution of protein functions of the trapped genes by GBT system, we performed GO  
8 analysis using the PANTHER protein classification. Although the protein functions related in  
9 transcriptional regulatory process, such as nucleic acid binding and transcription factors represented one  
10 relatively common class of isolated genes, the protein GO analysis indicated that the GBT protein trap  
11 was a useful tool for capturing a wide range of protein functions in addition to cell fate regulators and  
12 related nuclear genes.

13

## 14 **Cloned GBT Loci Represent a Rich Collection of Potential Human Disease Models**

15 More than 7,000 human diseases have already been described, and 80% of those are thought to have a  
16 genetic origin (Varga et al., 2018). Model organism studies can be a pivotal resource for understanding  
17 gene function which possibly provides additional insight into the cause of particular disease, thereby  
18 contributing to understanding of the pathogenic process and discovery of the therapeutic strategy  
19 (Wangler et al., 2017). Annotation of human orthologues of 160 tagged genes revealed that GBT  
20 technology yield a high frequency potential human disease models with loss of gene function, tracking  
21 expression of the truncated protein and Cre-revertible mutated allele to rescue the phenotype.

1 Furthermore, the modeling human genetic diseases using the GBT system has advantages compared  
2 with reverse genetic approaches such as TALEN and CRISPR-Cas9 systems. In initial phenotype  
3 screening, PR2.1 mutagenesis demonstrated 7% phenotype appearance as much as the other ENU- and  
4 retroviral forward genetic screenings (Amsterdam & Hopkins, 2004; Haffter et al., 1996; Kettleborough  
5 et al., 2013). For example, the annotation of disease-causing human orthologues of the tagged genes also  
6 revealed that the GBT system comprehensively developed mutants in zebrafish orthologues of human  
7 disease loci, including nervous, cardiovascular, endocrine, digestive, musculoskeletal, immune, and  
8 integument systems. Surprisingly, this system generated 68 of pioneering mutants in orthologs of human  
9 disease loci (Supplemental Table 2).

10

### 11 **Discovery of Novel Transcripts by Trapping Unpredicted Genes**

12 Since the completion of the zebrafish reference genome sequencing, it has enabled many new  
13 discoveries to be made, in particular the positional cloning of hundreds genes from mutation affecting  
14 embryogenesis behavior, physiology, and health and disease. However, a few poorly assembled regions  
15 remain (Howe et al., 2013). In molecular cloning of GBT lines generated, we found that a surprising  
16 proportion of the sequenced insertions does not correspond to any predicted genes. Although we have  
17 not formally excluded that mRFP expression might, in some case, be an artifact, the data of gene  
18 prediction provided in genome databases reveals some prediction errors. These results suggest that the  
19 algorithms used to predict genes from genome databases have missed a significant number of genes. The  
20 protein trapping by using GBT system may useful in identifying unsuspected novel genes, expressions  
21 and functions *in vivo* in real time.

22

## 1 **New Genomic Insights Using the GBT Random Insertion Mutagen**

2 The ready ease of these mRFP-based protein traps for basic expression analyses demonstrates how much  
3 we do not know about our overall proteome and the codex that is our genome. Nearly forty percent and  
4 five of six of the cloned genes show no expression at 2dpf and 4dpf, respectively, in ZFIN. At the  
5 subcellular level, such protein traps in conjunction with new microscopy techniques represent just  
6 another method for cellular and mechanistic analyses in an *in vivo* cellular context. Although these were  
7 made using random insertional approaches, new targeted integration tools using gene editing such as  
8 GeneWeld(Wierson et al., 2018) should readily empower labs to build their custom GBT lines in the  
9 future for genes not in this collection. Together, this initial 1100+ GBT collection is a new contribution  
10 to the use of the zebrafish to annotate the vertebrate genome.

11

1 **Competing interests**

2 The authors declare no competing interests.

1 **Figure legends**

2 **Table 1 Comparison availability of expression data**

3 **Table 2 Protein classification**

4 **Table 3 Disease-causing human orthologues**

5 **Table 4 Potential novel human disease models**

6 **Table 5 Novel transcripts**

7 **Figure 1 Schematic of the RP2 and RP8 gene-break transposon system with all three reading**  
8 **frames of AUG-less mRFP reporter**

9 A. Schematic of the RP2 system fused with 3 reading frames of AUG-less mRFP reporter (RP2.1, RP2.2  
10 and RP2.3). B. Schematic of the RP2 system fused with 3 reading frames of AUG-less mRFP reporter  
11 (RP8.1, RP8.2 and RP8.3). ITR, inverted terminal repeat; SA, loxP; Cre recombinase recognition  
12 sequence, splice acceptor; \*mRFP' AUG-less mRFP sequence; poly (A)+, polyadenylation signal; red  
13 octagon, extra transcriptional terminator and putative border element; *β-act*, carp beta-actin enhancer,  
14 SD, splice donor; E, enhancer; P, promoter; and WT, wild type.

15 **Figure 2 Novel protein expression**

16 Representative novel expression data of trapped proteins fused to the mRFP reporter in this GBT  
17 collection. A. unkl localized in olfactory pit, cerebrum and spinal cord at 4dpf. B. nusap1 localized in  
18 retina and the top layer of both forebrain and midbrain at 4 dpf. C. zgc:194659 strongly expressed in the  
19 brain and spinal cord at 4 dpf. D. marcksl1a expressed in the lens, skin and notochord at 2 dpf. E. pipp2a  
20 specifically localized to otoliths at 2 dpf. F. ahnak specifically expressed in skin at 4 dpf. G. dph1  
21 ubiquitously expressed showing granulized pattern in somites at 4 dpf. H. nfatc3a expressed in skeletal  
22 muscle and skin at 4 dpf . I. pard3bb localized to the pronephros and gut at 4 dpf. White arrowhead

1 shows an artificial expression in lens driven by the promoter of 3' exon trap of RP8.1. J.

2 LOC100537272 expressed in circulatory cells in the blood stream at 4 dpf.

3 **Figure 3 High knockdown efficiency of RP2.1 compared with other previous gene-trap systems**

4 Black dots of bar graphs shows percentage of remaining endogenous transcripts in homozygous larvae  
5 with mean and 95% confidence interval indicated by individual lines. The data of previous protein trap  
6 systems were also converted from the data in the original articles, R14-R15, our initial R-series protein  
7 trap vectors (n= 6),(Clark, Balciunas, et al., 2011) ; FlipTrap, FlipTrap vectors (n= 6), (Trinh le et al.,  
8 2011); FT1, FT1 vector (n=4),(T. T. Ni et al., 2012); RP2.1 (n=26), (Clark, Balciunas, et al., 2011; Ding  
9 et al., 2013; El-Rass et al., 2017; Westcot et al., 2015) and unpublished data),

10 **Figure 4 Summary of protein classes categorized the trapped proteins using PANTHER algorithm**

11 77 trapped genes were successfully categorized at least one of 21 protein classes by using PANTHER  
12 gene ontology algorithm. The details of the trapped genes classified in each protein class are listed in  
13 Table 2.

14 **Figure 5 Disease-causing human orthologues of the trapped genes involved in human genetic  
15 disorders in multi-organ systems.**

16 **Supplemental Table 1 Generation of 11 hundred independent lines by all three reading frames of  
17 mRFP reporter in both RP2 and RP8 cassettes.**

18 **Supplemental Table 2 Disease-causing human orthologues of the trapped protein**

19 **Supplemental Figure 1 Representative expression patterns of mRFP fusion protein by RP2 and  
20 RP8 integration**

21



1 **Acknowledgements**

2 This work is supported by grants from the National Institutes of Health (GM63904; DA14546;  
3 HG006431) and the Mayo Foundation. We thank Zoltan Varga for sharing of the ZIRC sperm  
4 cryopreservation protocol prior to publication. Appreciation is also extended to the Mayo Clinic  
5 Zebrafish Facility staff for their excellent support. This research was in part funded by the Intramural  
6 Research Program of the National Human Genome Research Institute; National Institutes of Health  
7 (S.M.B.: 1ZIAHG000183)

8

9

## 1 References

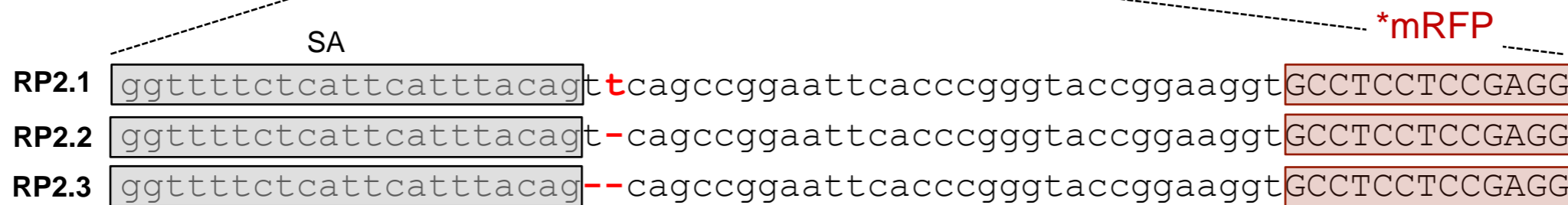
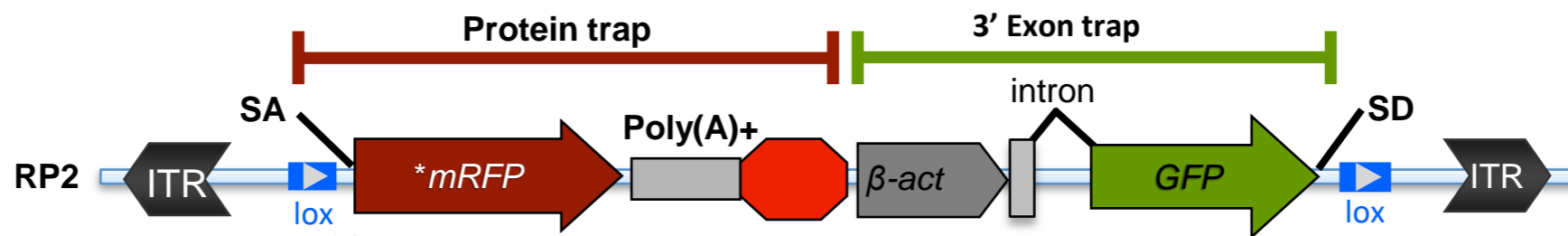
- 2 Amsterdam, A., & Hopkins, N. (2004). Retroviral-mediated insertional mutagenesis in zebrafish. *Methods Cell*  
3 *Biol*, 77, 3-20.
- 4 Balciunas, D., Wangenstein, K. J., Wilber, A., Bell, J., Geurts, A., Sivasubbu, S., . . . Ekker, S. C. (2006). Harnessing  
5 a high cargo-capacity transposon for genetic applications in vertebrates. *PLoS Genet*, 2(11), e169.  
6 doi:10.1371/journal.pgen.0020169
- 7 Clark, K. J., Argue, D. P., Petzold, A. M., & Ekker, S. C. (2012). zfishbook: connecting you to a world of zebrafish  
8 revertible mutants. *Nucleic Acids Res*, 40(Database issue), D907-911. doi:10.1093/nar/gkr957
- 9 Clark, K. J., Balciunas, D., Pogoda, H. M., Ding, Y., Westcot, S. E., Bedell, V. M., . . . Ekker, S. C. (2011). In vivo  
10 protein trapping produces a functional expression codex of the vertebrate proteome. *Nat Methods*, 8(6),  
11 506-515. doi:10.1038/nmeth.1606
- 12 Clark, K. J., Urban, M. D., Skuster, K. J., & Ekker, S. C. (2011). Transgenic zebrafish using transposable elements.  
13 *Methods Cell Biol*, 104, 137-149. doi:10.1016/B978-0-12-374814-0.00008-2
- 14 Ding, Y., Liu, W., Deng, Y., Jomok, B., Yang, J., Huang, W., . . . Xu, X. (2013). Trapping cardiac recessive mutants  
15 via expression-based insertional mutagenesis screening. *Circ Res*, 112(4), 606-617.  
16 doi:10.1161/CIRCRESAHA.112.300603
- 17 Ding, Y., Long, P. A., Bos, J. M., Shih, Y. H., Ma, X., Sundsbak, R. S., . . . Xu, X. (2016). A modifier screen identifies  
18 DNAJB6 as a cardiomyopathy susceptibility gene. *JCI Insight*, 1(14). doi:10.1172/jci.insight.88797
- 19 Draper, B. W., & Moens, C. B. (2009). A high-throughput method for zebrafish sperm cryopreservation and in  
20 vitro fertilization. *J Vis Exp*(29). doi:10.3791/1395
- 21 El-Rass, S., Eisa-Beygi, S., Khong, E., Brand-Arzamendi, K., Mauro, A., Zhang, H., . . . Wen, X. Y. (2017). Disruption  
22 of pdgfra alters endocardial and myocardial fusion during zebrafish cardiac assembly. *Biol Open*, 6(3),  
23 348-357. doi:10.1242/bio.021212
- 24 Haffter, P., Granato, M., Brand, M., Mullins, M. C., Hammerschmidt, M., Kane, D. A., . . . Nusslein-Volhard, C.  
25 (1996). The identification of genes with unique and essential functions in the development of the  
26 zebrafish, *Danio rerio*. *Development*, 123, 1-36.
- 27 Harris, M. P., Rohner, N., Schwarz, H., Perathoner, S., Konstantinidis, P., & Nusslein-Volhard, C. (2008). Zebrafish  
28 eda and edar mutants reveal conserved and ancestral roles of ectodysplasin signaling in vertebrates.  
29 *PLoS Genet*, 4(10), e1000206. doi:10.1371/journal.pgen.1000206
- 30 Howe, K., Clark, M. D., Torroja, C. F., Torrance, J., Berthelot, C., Muffato, M., . . . Stemple, D. L. (2013). The  
31 zebrafish reference genome sequence and its relationship to the human genome. *Nature*, 496(7446),  
32 498-503. doi:10.1038/nature12111
- 33 Kawakami, K., Takeda, H., Kawakami, N., Kobayashi, M., Matsuda, N., & Mishina, M. (2004). A transposon-  
34 mediated gene trap approach identifies developmentally regulated genes in zebrafish. *Dev Cell*, 7(1),  
35 133-144. doi:10.1016/j.devcel.2004.06.005
- 36 Kettleborough, R. N., Busch-Nentwich, E. M., Harvey, S. A., Dooley, C. M., de Bruijn, E., van Eeden, F., . . .  
37 Stemple, D. L. (2013). A systematic genome-wide analysis of zebrafish protein-coding gene function.  
38 *Nature*, 496(7446), 494-497. doi:10.1038/nature11992
- 39 Liao, H. K., Wang, Y., Noack Watt, K. E., Wen, Q., Breitbart, J., Kemmet, C. K., . . . McGrail, M. (2012). Tol2 gene  
40 trap integrations in the zebrafish amyloid precursor protein genes appa and apl2 reveal accumulation  
41 of secreted APP at the embryonic veins. *Dev Dyn*, 241(2), 415-425. doi:10.1002/dvdy.23725
- 42 Liu, T. L., Upadhyayula, S., Milkie, D. E., Singh, V., Wang, K., Swinburne, I. A., . . . Betzig, E. (2018). Observing the  
43 cell in its native state: Imaging subcellular dynamics in multicellular organisms. *Science*, 360(6386).  
44 doi:10.1126/science.aaq1392
- 45 Matthews, J. L., Murphy, J. M., Carmichael, C., Yang, H., Tiersch, T., Westerfield, M., & Varga, Z. M. (2018).  
46 Changes to Extender, Cryoprotective Medium, and In Vitro Fertilization Improve Zebrafish Sperm  
47 Cryopreservation. *Zebrafish*, 15(3), 279-290. doi:10.1089/zeb.2017.1521

- 1 Meadows, J. R. S., & Lindblad-Toh, K. (2017). Dissecting evolution and disease using comparative vertebrate  
2 genomics. *Nat Rev Genet*, 18(10), 624-636. doi:10.1038/nrg.2017.51
- 3 Mi, H., Muruganujan, A., & Thomas, P. D. (2013). PANTHER in 2013: modeling the evolution of gene function,  
4 and other gene attributes, in the context of phylogenetic trees. *Nucleic Acids Res*, 41(Database issue),  
5 D377-386. doi:10.1093/nar/gks1118
- 6 Ni, J., Wangenstein, K. J., Nelsen, D., Balciunas, D., Skuster, K. J., Urban, M. D., & Ekker, S. C. (2016). Active  
7 recombinant Tol2 transposase for gene transfer and gene discovery applications. *Mob DNA*, 7, 6.  
8 doi:10.1186/s13100-016-0062-z
- 9 Ni, T. T., Lu, J., Zhu, M., Maddison, L. A., Boyd, K. L., Huskey, L., . . . Chen, W. (2012). Conditional control of gene  
10 function by an invertible gene trap in zebrafish. *Proc Natl Acad Sci U S A*, 109(38), 15389-15394.  
11 doi:10.1073/pnas.1206131109
- 12 Petzold, A. M., Balciunas, D., Sivasubbu, S., Clark, K. J., Bedell, V. M., Westcot, S. E., . . . Ekker, S. C. (2009).  
13 Nicotine response genetics in the zebrafish. *Proc Natl Acad Sci U S A*, 106(44), 18662-18667.  
14 doi:10.1073/pnas.0908247106
- 15 Petzold, A. M., Bedell, V. M., Boczek, N. J., Essner, J. J., Balciunas, D., Clark, K. J., & Ekker, S. C. (2010). SCORE  
16 imaging: specimen in a corrected optical rotational enclosure. *Zebrafish*, 7(2), 149-154.  
17 doi:10.1089/zeb.2010.0660
- 18 Sivasubbu, S., Balciunas, D., Davidson, A. E., Pickart, M. A., Hermanson, S. B., Wangenstein, K. J., . . . Ekker, S. C.  
19 (2006). Gene-breaking transposon mutagenesis reveals an essential role for histone H2afza in zebrafish  
20 larval development. *Mech Dev*, 123(7), 513-529. doi:10.1016/j.mod.2006.06.002
- 21 Stoeger, T., Gerlach, M., Morimoto, R. I., & Nunes Amaral, L. A. (2018). Large-scale investigation of the reasons  
22 why potentially important genes are ignored. *PLoS Biol*, 16(9), e2006643.  
23 doi:10.1371/journal.pbio.2006643
- 24 Trinh le, A., & Fraser, S. E. (2013). Enhancer and gene traps for molecular imaging and genetic analysis in  
25 zebrafish. *Dev Growth Differ*, 55(4), 434-445. doi:10.1111/dgd.12055
- 26 Trinh le, A., Hochgreb, T., Graham, M., Wu, D., Ruf-Zamojski, F., Jayasena, C. S., . . . Fraser, S. E. (2011). A  
27 versatile gene trap to visualize and interrogate the function of the vertebrate proteome. *Genes Dev*,  
28 25(21), 2306-2320. doi:10.1101/gad.174037.111
- 29 Urasaki, A., Morvan, G., & Kawakami, K. (2006). Functional dissection of the Tol2 transposable element  
30 identified the minimal cis-sequence and a highly repetitive sequence in the subterminal region essential  
31 for transposition. *Genetics*, 174(2), 639-649. doi:10.1534/genetics.106.060244
- 32 Van Slyke, C. E., Bradford, Y. M., Howe, D. G., Fashena, D. S., Ramachandran, S., Ruzicka, L., & Staff\*, Z. (2018).  
33 Using ZFIN: Data Types, Organization, and Retrieval. *Methods Mol Biol*, 1757, 307-347. doi:10.1007/978-  
34 1-4939-7737-6\_11
- 35 Varga, M., Ralbovszki, D., Balogh, E., Hamar, R., Keszthelyi, M., & Tory, K. (2018). Zebrafish Models of Rare  
36 Hereditary Pediatric Diseases. *Diseases*, 6(2). doi:10.3390/diseases6020043
- 37 Varshney, G. K., Huang, H., Zhang, S., Lu, J., Gildea, D. E., Yang, Z., . . . Burgess, S. M. (2013). The Zebrafish  
38 Insertion Collection (ZInC): a web based, searchable collection of zebrafish mutations generated by DNA  
39 insertion. *Nucleic Acids Res*, 41(Database issue), D861-864. doi:10.1093/nar/gks946
- 40 Varshney, G. K., Lu, J., Gildea, D. E., Huang, H., Pei, W., Yang, Z., . . . Lin, S. (2013). A large-scale zebrafish gene  
41 knockout resource for the genome-wide study of gene function. *Genome Res*, 23(4), 727-735.  
42 doi:10.1101/gr.151464.112
- 43 Wangler, M. F., Yamamoto, S., Chao, H. T., Posey, J. E., Westerfield, M., Postlethwait, J., . . . Bellen, H. J. (2017).  
44 Model Organisms Facilitate Rare Disease Diagnosis and Therapeutic Research. *Genetics*, 207(1), 9-27.  
45 doi:10.1534/genetics.117.203067
- 46 Westcot, S. E., Hatzold, J., Urban, M. D., Richetti, S. K., Skuster, K. J., Harm, R. M., . . . Ekker, S. C. (2015). Protein-  
47 Trap Insertional Mutagenesis Uncovers New Genes Involved in Zebrafish Skin Development, Including a

- 1            Neuregulin 2a-Based ErbB Signaling Pathway Required during Median Fin Fold Morphogenesis. *PLoS One*,  
2            10(6), e0130688. doi:10.1371/journal.pone.0130688
- 3            Wierson, W. A., Welker, J. M., Almeida, M. P., Mann, C. M., Webster, D. A., Torrie, M. E., . . . Essner, J. J. (2018).  
4            GeneWeld: a method for efficient targeted integration directed by short homology. *bioRxiv*, 431627.  
5            doi:10.1101/431627
- 6            Xu, J., Gao, J., Li, J., Xue, L., Clark, K. J., Ekker, S. C., & Du, S. J. (2012). Functional analysis of slow myosin heavy  
7            chain 1 and myomesin-3 in sarcomere organization in zebrafish embryonic slow muscles. *J Genet*  
8            *Genomics*, 39(2), 69-80. doi:10.1016/j.jgg.2012.01.005

9

A



B

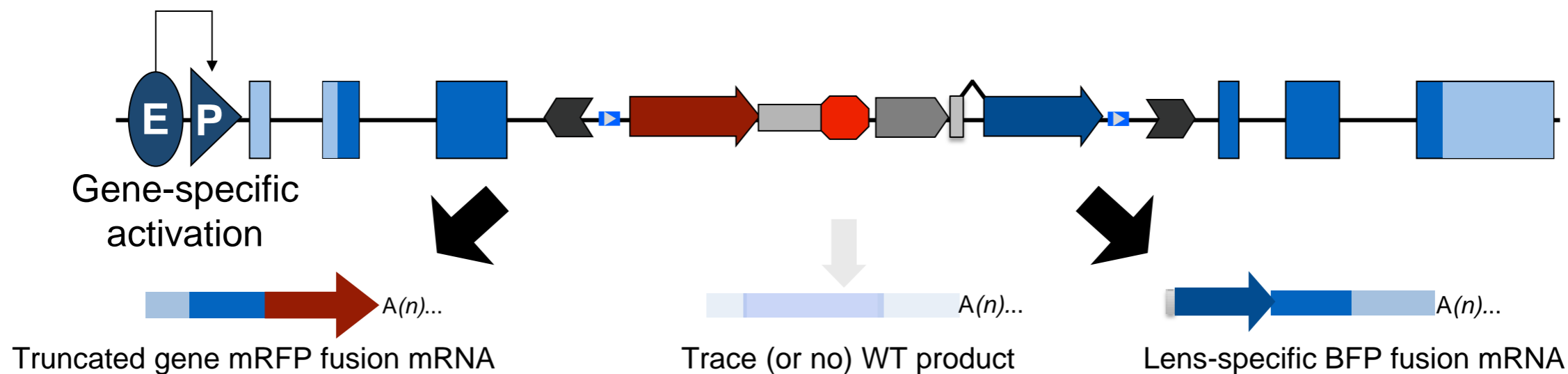
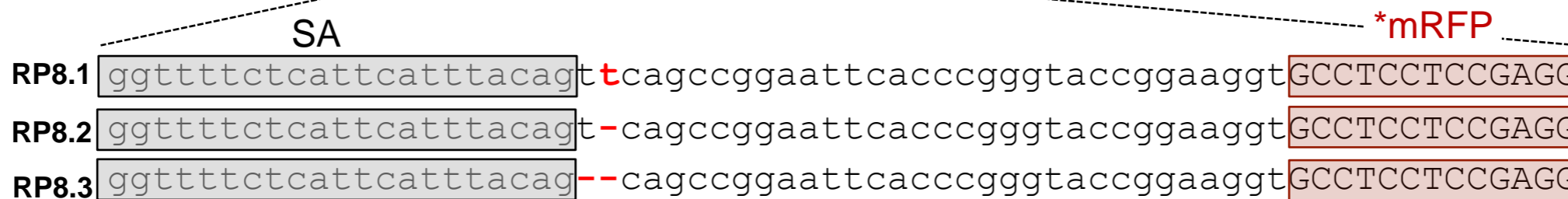
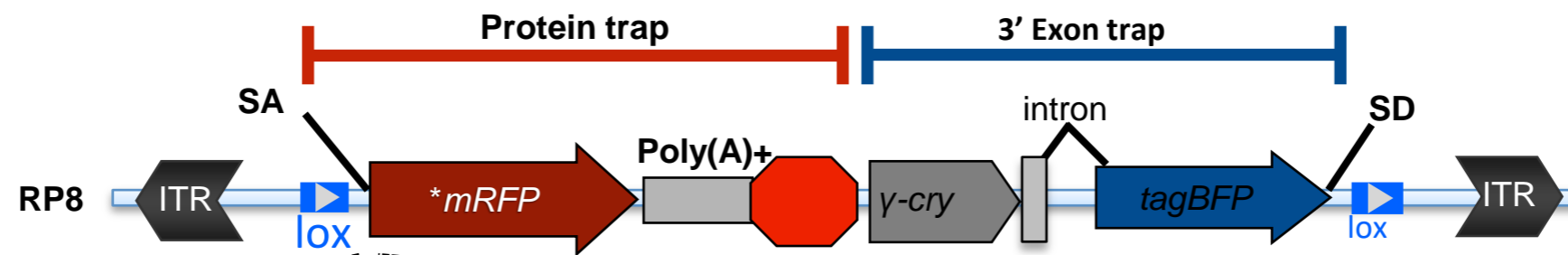


Fig. 1

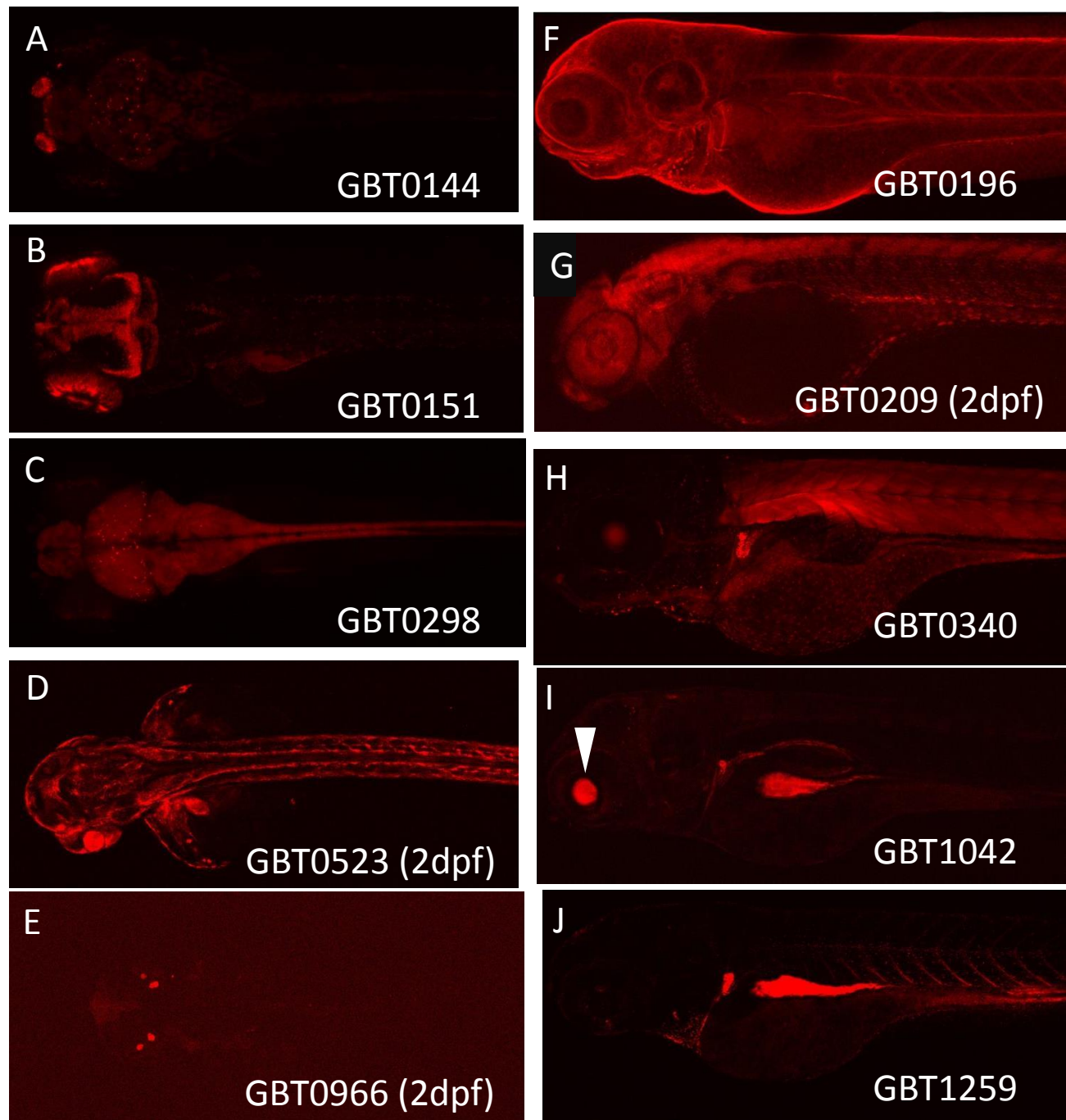


Fig.2

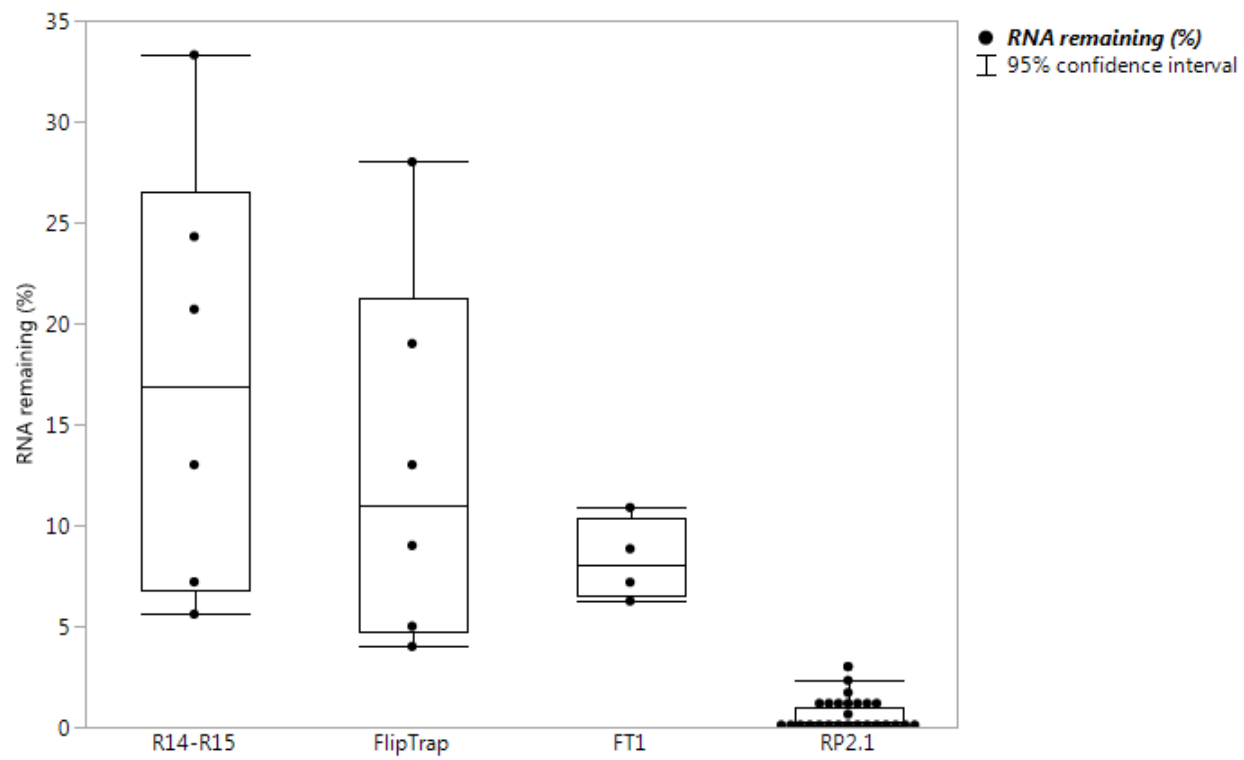


Fig. 3

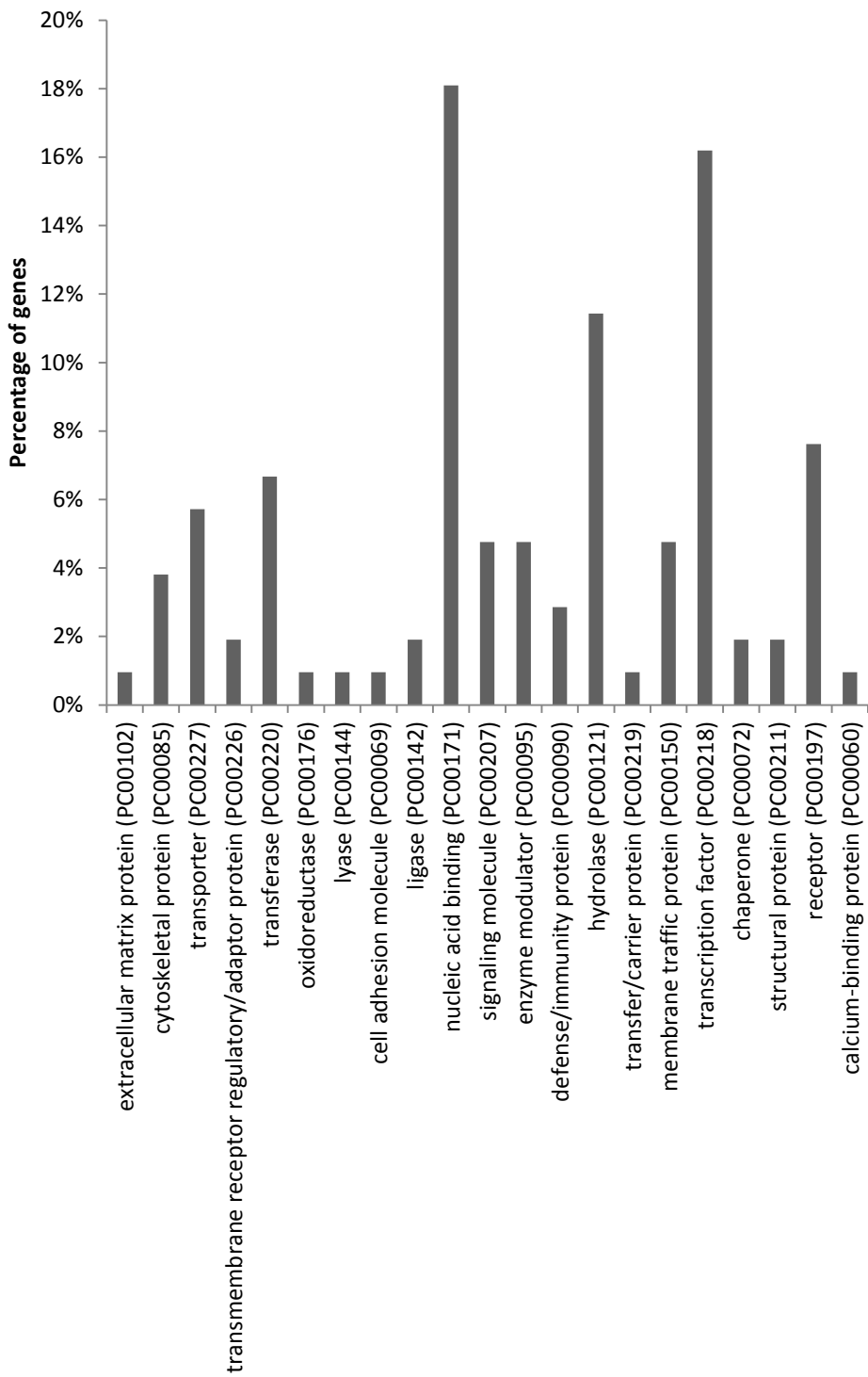


Figure 4



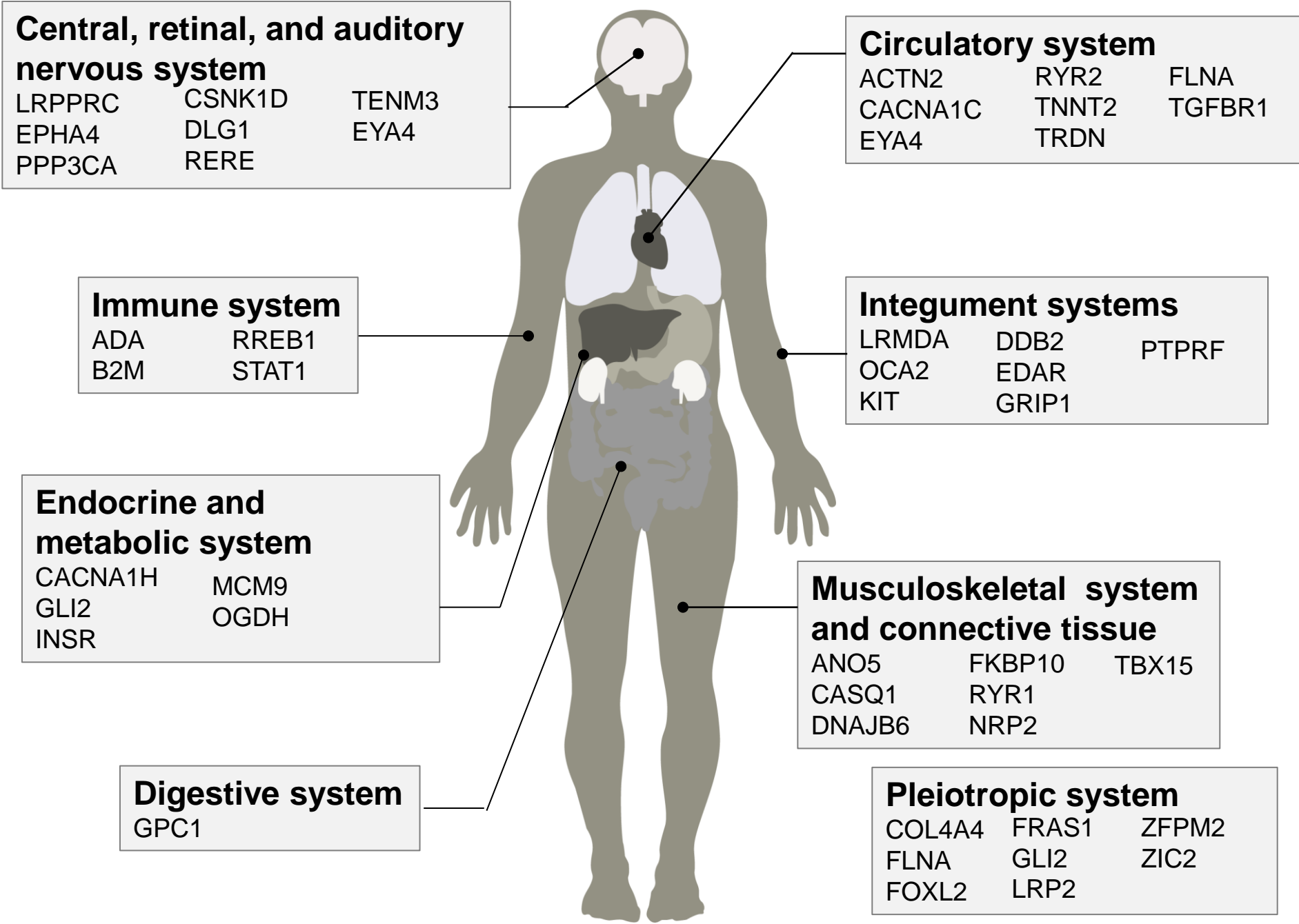


Fig. 5

Expression data exist	The number of genes	
	2 dpf	4 dpf
zfishbook and ZFIN	45	10
zfishbook	75	155
ZFIN	20	2
Not available	36	9

# Table 1

Protein Class	Mapped ID	Gene Symbol	Gene Name
<b>nucleic acid binding (PC00171)</b>			
	ZDB-GENE-040426-893	gar1	H/ACA ribonucleoprotein complex subunit 1
	ZDB-GENE-060503-160	adarb2	Adenosine deaminase, RNA-specific, B2 (non-functional) (Fragment)
	ZDB-GENE-000405-1	pbx1a	Pbx1a homeodomain protein
	ZDB-GENE-060130-4	zfpm2a	Zinc finger protein, FOG family member 2a
	ZDB-GENE-050417-327	eef1a1b	Elongation factor 1-alpha
	ZDB-GENE-081104-328	enox1	Ecto-NOX disulfide-thiol exchanger 1
	ZDB-GENE-050419-169	ddb2	DNA damage-binding protein 2
	ZDB-GENE-980526-306	msxc	Homeobox protein MSH-C
	ZDB-GENE-020529-1	tbx15	T-box 15
	ZDB-GENE-070314-2	rxraa	Retinoic acid receptor RXR-alpha-A
	ZDB-GENE-030131-6117	dido1	Death inducer-obliterator 1
	ZDB-GENE-050913-153	barhl2	BarH-class homeodomain transcription factor
	ZDB-GENE-040426-1272	ddb1	Damage-specific DNA-binding protein 1
	ZDB-GENE-980526-499	stat1a	Signal transducer and activator of transcription
	ZDB-GENE-050517-31	abcf1	ATP-binding cassette, sub-family F (GCN20), member 1
	ZDB-GENE-080403-11	kat2a	Histone acetyltransferase KAT2A
	ZDB-GENE-030131-4505	dhx37	DEAH (Asp-Glu-Ala-His) box polypeptide 37
	ZDB-GENE-060512-241	foxl2a	Forkhead box L2
	ZDB-GENE-041111-41	nfatc3a	Nuclear factor of-activated T cells 3a
<b>transcription factor (PC00218)</b>			
	ZDB-GENE-050419-261	wtip	Wilms tumor protein 1-interacting protein homolog
	ZDB-GENE-050913-139	znf1015	Zinc finger protein 1015
	ZDB-GENE-000405-1	pbx1a	Pbx1a homeodomain protein
	ZDB-GENE-060130-4	zfpm2a	Zinc finger protein, FOG family member 2a
	ZDB-GENE-000710-4	zic2a	ZIC family member 2 (Odd-paired homolog, Drosophila), A
	ZDB-GENE-061207-62	bcl11ba	B cell CLL/lymphoma 11Ba
	ZDB-GENE-050419-146	bhlhe41	BHLH protein DEC2
	ZDB-GENE-030131-6789	taf6l	TAF6-like RNA polymerase II, p300/CBP-associated factor (PCAF)-associated factor
	ZDB-GENE-980526-306	msxc	Homeobox protein MSH-C
	ZDB-GENE-020529-1	tbx15	T-box 15
	ZDB-GENE-070314-2	rxraa	Retinoic acid receptor RXR-alpha-A
	ZDB-GENE-030131-6117	dido1	Death inducer-obliterator 1
	ZDB-GENE-050913-153	barhl2	BarH-class homeodomain transcription factor
	ZDB-GENE-980526-499	stat1a	Signal transducer and activator of transcription
	ZDB-GENE-060130-108	casz1	Castor zinc finger 1
	ZDB-GENE-060512-241	foxl2a	Forkhead box L2
	ZDB-GENE-041111-41	nfatc3a	Nuclear factor of-activated T cells 3a
<b>hydrolase (PC00121)</b>			
	ZDB-GENE-040611-3	nrp2b	Neuropilin
	ZDB-GENE-060503-160	adarb2	Adenosine deaminase, RNA-specific, B2 (non-functional) (Fragment)
	ZDB-GENE-070112-1812	pnpla7a	Patatin-like phospholipase domain-containing 7a
	ZDB-GENE-050417-327	eef1a1b	Elongation factor 1-alpha
	ZDB-GENE-040718-393	ada	Adenosine deaminase
	ZDB-GENE-140703-2	plpp2a	Phospholipid phosphatase 2a
	ZDB-GENE-080818-1	ca16b	Carbonic anhydrase XVI b
	ZDB-GENE-050517-31	abcf1	ATP-binding cassette, sub-family F (GCN20), member 1
	ZDB-GENE-060503-530	ptprf	Receptor-type tyrosine-protein phosphatase F
	ZDB-GENE-070117-757	parga	Poly(ADP-ribose) glycohydrolase
	ENSDARG0000010758	capn12	Calpain 12
	ZDB-GENE-041014-310	mcm9	DNA helicase MCM9
<b>receptor (PC00197)</b>			
	ZDB-GENE-050419-65	v2rl1	Vomer nasal 2 receptor, l1
	ZDB-GENE-060503-5	gabbr1b	Gamma-aminobutyric acid (GABA) B receptor, 1b
	ZDB-GENE-070314-2	rxraa	Retinoic acid receptor RXR-alpha-A
	ZDB-GENE-030131-2427	col7a1	Collagen, type VII, alpha 1
	ZDB-GENE-030909-10	itgb1b	Integrin beta
	ZDB-GENE-080818-1	ca16b	Carbonic anhydrase XVI b
	ZDB-GENE-091027-1	tgfbr1b	Serine/threonine-protein kinase receptor
	ZDB-GENE-070905-3	crfb4	Cytokine receptor family member b4
<b>transferase (PC00220)</b>			
	ZDB-GENE-030131-2049	mat2aa	S-adenosylmethionine synthase
	ZDB-GENE-040426-1516	mboat7	Lysophospholipid acyltransferase 7
	ZDB-GENE-121129-1	csgalnact1a	Hexosyltransferase
	ZDB-GENE-060616-238	mgat5	Mannosyl (alpha-1,6-)-glycoprotein beta-1,6-N-acetyl-glucosaminyltransferase
	ZDB-GENE-080403-11	kat2a	Histone acetyltransferase KAT2A
	ZDB-GENE-030131-825	csnk1da	Casein kinase I isoform delta-A
	ZDB-GENE-091027-1	tgfbr1b	Serine/threonine-protein kinase receptor

bioRxiv preprint doi: <https://doi.org/10.1101/630236>; this version posted May 7, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

Table 2\_Protein Classification.xlsx

<b>transporter (PC00227)</b>			
ZDB-GENE-130103-1	cacna1ha		Voltage-dependent T-type calcium channel subunit alpha
ZDB-GENE-050419-65	v2rl1		Vomeronal 2 receptor, l1
ZDB-GENE-070705-417	ryr1b		Ryanodine receptor 1b (skeletal)
ZDB-GENE-001127-2	atp1b2a		Sodium/potassium-transporting ATPase subunit beta
ZDB-GENE-050517-31	abcf1		ATP-binding cassette, sub-family F (GCN20), member 1
ZDB-GENE-070718-4	oca2		Oculocutaneous albinism II
<b>enzyme modulator (PC00095)</b>			
ZDB-GENE-060503-160	adarb2		Adenosine deaminase, RNA-specific, B2 (non-functional) (Fragment)
ZDB-GENE-050417-327	eef1a1b		Elongation factor 1-alpha
ZDB-GENE-030131-9805	ppp1r13ba		Protein phosphatase 1, regulatory subunit 13Ba
ZDB-GENE-110411-39	si:dkey-222l13.1		Si:dkey-222l13.1
ZDB-GENE-040426-1131	phactr4b		Phosphatase and actin regulator 4B
<b>membrane traffic protein (PC00150)</b>			
ZDB-GENE-030131-5290	napab		N-ethylmaleimide sensitive fusion protein attachment protein alpha
ZDB-GENE-050522-134	syt5b		Synaptotagmin Vb
ZDB-GENE-040718-281	vapal		VAMP (vesicle-associated membrane protein)-associated protein A,-like
ZDB-GENE-050522-235	sncgb		Synuclein, gamma b (breast cancer-specific protein 1)
ZDB-GENE-110411-39	si:dkey-222l13.1		Si:dkey-222l13.1
<b>signaling molecule (PC00207)</b>			
ZDB-GENE-080225-17	arngel25b		Rho guanine nucleotide exchange factor (GEF) 25b
ZDB-GENE-070615-10	nrg2a		Neuregulin 2a
ZDB-GENE-041114-101	fgf13a		Fibroblast growth factor
ZDB-GENE-040426-1793	fgf13b		Fibroblast growth factor
<b>cytoskeletal protein (PC00085)</b>			
ZDB-GENE-070112-1512	wasf3b		WAS protein family, member 3b
ZDB-GENE-080215-23	map7d1b		MAP7 domain-containing 1b
ZDB-GENE-060503-22	map7d1a		MAP7 domain-containing 1a
ZDB-GENE-000626-1	tnnt2a		Troponin T type 2a (cardiac)
<b>defense/immunity protein (PC00090)</b>			
ZDB-GENE-060503-160	adarb2		Adenosine deaminase, RNA-specific, B2 (non-functional) (Fragment)
ZDB-GENE-040426-2136	b2ml		Beta-2-microglobulin
ZDB-GENE-070905-3	crfb4		Cytokine receptor family member b4
<b>chaperone (PC00072)</b>			
ZDB-GENE-040426-876	cct8		Chaperonin-containing TCP1, subunit 8 (theta)
ZDB-GENE-050522-235	sncgb		Synuclein, gamma b (breast cancer-specific protein 1)
<b>ligase (PC00142)</b>			
ZDB-GENE-090313-285	zmiz2		Zinc finger, MIZ-type-containing 2
ZDB-GENE-070912-399	si:dkey-181m9.8		Si:dkey-181m9.8
<b>structural protein (PC00211)</b>			
ZDB-GENE-030710-8	gpm6ab		Glycoprotein M6Ab (Fragment)
ZDB-GENE-030710-9	gpm6ba		DMgamma1
<b>transmembrane receptor regulatory/adaptor protein (PC00226)</b>			
ZDB-GENE-081030-10	cntn3a.1		Contactin 3a, tandem duplicate 1
ZDB-GENE-010724-8	dlg1		Disks large homolog 1
<b>cell adhesion molecule (PC00069)</b>			
ZDB-GENE-030909-10	itgb1b		Integrin beta
<b>extracellular matrix protein (PC00102)</b>			
ZDB-GENE-101112-3	megf6b		Multiple EGF-like-domains 6b
<b>lyase (PC00144)</b>			
ZDB-GENE-141030-2	npr2		Guanylate cyclase
<b>oxidoreductase (PC00176)</b>			
ZDB-GENE-081104-328	enox1		Ecto-NOX disulfide-thiol exchanger 1
<b>transfer/carrier protein (PC00219)</b>			
ZDB-GENE-060628-2	xpo7		Exportin 7
<b>calcium-binding protein (PC00060)</b>			
ENSDARG00000010758	capn12		Calpain 12

## Table 2

Line	ZF NCBI gene ID	ZF gene symbol	Human Orthologue NCBI Gene ID	Human Orthologue Gene symbol
GBT0001	559053	casz1	54897	CASZ1
GBT0002	563408	sorbs2b	8470	SORBS2
GBT0005	570216	itgb1b	3688	ITGB1
GBT0007	794348	nrde2	55051	NRDE2
GBT0010	30461	cdh11	1009	CDH11
GBT0016	58138	pbx1a	5087	PBX1
GBT0019	565629	kcnk10b	54207	KCNK10
GBT0020	569437	cntn3a.1	5067	CNTN3
GBT0021	571891	cntnap5b	129684	CNTNAP5
GBT0025	553277	dido1	11083	DIDO1
GBT0028	101882228	CABZ01057928.1	1729	DIAPH1
GBT0031	58071	tnnt2a	7139	TNNT2
GBT0033	100003333	lrch4	4034	LRCH4
GBT0034	567642	si:dkey-181m9.8		
GBT0035	559134	parga	8505	PARG
GBT0038	541513	srpx	8406	SRPX
GBT0039	559579	gabbr1,2	2550	GABBR1
GBT0040	58049	hoxa3a	3201	HOXA4
GBT0043	323266	cd99l2	83692	CD99L2
GBT0046	64270	epha4b	2043	EPHA4
GBT0060	797527	crfb4	3588	IL10RB
GBT0067	559614	myom3	127294	MYOM3
GBT0070	555610	mcm9	254394	MCM9
GBT0073	406467	abcf1	23	ABCF1
GBT0077	406467	abcf1	23	ABCF1
GBT0078	558006	grip1	23426	GRIP1
GBT0082	394037	cct8	10694	CCT8
GBT0091	100002220	enox1	55068	ENOX1
GBT0094	403003	fgf13b	2258	FGF13
GBT0096	394094	sh3glb2b	56904	SH3GLB2
GBT0101	553277	dido1	11083	DIDO1
GBT0103	100002190	cyth3a	9265	CYTH3
GBT0111	569081	rerea	473	RERE
GBT0113	569437	cntn3a.1	5067	CNTN3
GBT0125	100333685	cd302	9936	CD302
GBT0126	405902	nrp2b	8828	NRP2
GBT0131	100004133	ddb2	1643	DDB2
GBT0133	58077	zic2a	7546	ZIC2
GBT0135	563771	bhlhe41	79365	BHLHE41
GBT0137	100004503	eef1a1b	1915	EEF1A1
GBT0141	368225	gpm6ba	2824	GPM6B
GBT0142	100004850	aatka	9625	AATK
GBT0143	566046	wtip	126374	WTIP
GBT0144	565405	unkl	64718	UNKL

Table 3\_Human orthologues of GBT confirmed genes updated Apr8-2019.xlsx

GBT0145	560542	e pn2	22905	EPN2
GBT0151	567446	nusap1	51203	NUSAP1
GBT0154	449761	si:ch211-163l21.8	57535	KIAA1324
GBT0156	563428	fras1	80144	FRAS1
GBT0157	724006	mgat5	4249	MGAT5
GBT0166	64269	atp1b2a	482	ATP1B2
GBT0168	492758	fgf13a	2258	FGF13
GBT0170	378997	glcci1a	113263	GLCCI1
GBT0172	337397	tmem30aa	55754	TMEM30A
GBT0175	569561	arhgef25b	115557	ARHGEF25
GBT0178	436919	ada	100	ADA
GBT0181	30155	tenm3	55714	TENM3
GBT0186	170581	cacna1c	775	CACNA1C
GBT0187	100007836	si:dkey-253d23.3		
GBT0189	562459	znf414	84330	ZNF414
GBT0190	100001699	barhl2	343472	BARHL2
GBT0195	562459	znf414	84330	ZNF414
GBT0196	559276	ahnak	79026	AHNAK
GBT0200	541428	znrd1	30834	ZNRD1
GBT0201	556341	mhc1zca	3140	MR1
GBT0202	569081	rerea	473	RERE
GBT0203	678534	ntm	150084	IGSF5
GBT0204	492787	cadm4	199731	CADM4
GBT0205	100151756	fam117ab	81558	FAM117A
GBT0208	641575	actn2b	88	ACTN2
GBT0209	550559	dph1	1801	DPH1
GBT0230	797651	jam3a	83700	JAM3
GBT0231	266983	neo1a	4756	NEO1
GBT0232	553567	syt5b	6861	SYT5
GBT0235	557195	lrpprc	10128	LRPPRC
GBT0237	100006951	nrg2a	9542	NRG2
GBT0238	792928	tgfbr1b	7046	TGFBR1
GBT0239	796270	map7d1b	55700	MAP7D1
GBT0240	100093707	bcl11ba	64919	BCL11B
GBT0241	561900	fip1l1a	81608	FIP1L1
GBT0242	554131	tex261	113419	TEX261
GBT0243	563577	ppp1r13ba	23368	PPP1R13B
GBT0245	30526	msx3		
GBT0249	791742	b2ml	567	B2M
GBT0250	386769	ptprma	5797	PTPRM
GBT0251	692279	foxl2a	668	FOXL2
GBT0255	114441	ncam2	4685	NCAM2
GBT0256	550554	ppp3cb	5530	PPP3CA
GBT0261	799247	capn12	147968	CAPN12
GBT0263	559150	hmcn1	83872	HMCN1
GBT0268	100534737	ano5a	203859	ANO5
GBT0270	568996	zfpm2a	23414	ZFPM2

Table 3\_Human orthologues of GBT confirmed genes updated Apr8-2019.xlsx

GBT0271	565172	map7d1a	55700	MAP7D1
GBT0275	554270	col4a4	1286	COL4A4
GBT0281	797715	ogdhb	4967	OGDH
GBT0283	793623	sh3kbp1	30011	SH3KBP1
GBT0286	797252	emid1	129080	EMID1
GBT0292	541428	znrd1	30834	ZNRD1
GBT0298	100007373	zgc:194659		
GBT0312	799867	ptprfb	5792	PTPRF
GBT0313	565460	csmd2	114784	CSMD2
GBT0316	324381	fkbp10b	60681	FKBP10
GBT0319	791987	ino80c	125476	INO80C
GBT0321	791987	ino80c	125476	INO80C
GBT0322	101886266	si:ch211-160j14.3	91522	COL23A1
GBT0323	558149	adarb2	105	ADARB2
GBT0325	557764	megf6b	1953	MEGF6
GBT0329	793671	csgalnact1a	55790	CSGALNACT1
GBT0340	561772	nfatc3a	4775	NFATC3
GBT0346	553730	lrmda	83938	LRMDA
GBT0348	570245	ryr1b	6261	RYR1
GBT0357	569183	ca16b	5793	PTPRG
GBT0364	323329	mat2aa	4144	MAT2A
GBT0365	553679	sncgb	6623	SNCG
GBT0369	407734	gpm6ab	2823	GPM6A
GBT0380	795234	znf1140		
GBT0383	554230	kdr	3791	KDR
GBT0389	767723	mosmob	730094	MOSMO
GBT0396	562557	smg9	56006	SMG9
GBT0397	393643	phactr4b	65979	PHACTR4
GBT0398	791173	pnpla7a	375775	PNPLA7
GBT0399	571887	kirrel3a	84623	KIRREL3
GBT0401	560003	ERC1	23085	ERC1
GBT0402	566484	scaf11	9169	SCAF11
GBT0404	564731	CABZ01045212.1	ENSGGOG00000007420	
GBT0409	100331745	npr2	4882	NPR2
GBT0410	436819	vapal	9218	VAPA
GBT0411	393275	dnajb6b	10049	DNAJB6
GBT0412	407710	xpo7	23039	XPO7
GBT0415	570029	arrdc1b	92714	ARRDC1
GBT0416	322795	csrnp1b	64651	CSRNP1
GBT0419	555578	rxraa	6256	RXRA
GBT0422	245700	insrb	3643	INSR
GBT0424	566039	v2r11		
GBT0425	550455	mrps18b	28973	MRPS18B
GBT0433	393599	ddb1	1642	DDB1
GBT0434	393950	gar1	54433	GAR1
GBT0435	558326	nrxn2a	9379	NRXN2
GBT0437	445226	casq1a	844	CASQ1

Table 3\_Human orthologues of GBT confirmed genes updated Apr8-2019.xlsx



GBT0503	30768	stat1a	6772	STAT1
GBT0505	572369	kirrel3b	84623	KIRREL3
GBT0510	553343	si:ch211-266g18.10	10345	TRDN
GBT0511	567419	oca2	4948	OCA2
GBT0513	65239	map2k6	5608	MAP2K6
GBT0520	553367	gpc1a	2817	GPC1
GBT0522	30708	fabp2	2169	FABP2
GBT0523	406407	marcksl1a	65108	MARCKSL1
GBT0525	100333948	mapk8ip1b	9479	MAPK8IP1
GBT0527	386629	tob1a	10140	TOB1
GBT0528	323269	zbtb16a	7704	ZBTB16
GBT0534	327082	napab	8775	NAPA
GBT0545	563632	wasf3b	10810	WASF3
GBT0552	562615	znf1015	163033	ZNF579
GBT0554	492500	gsto2	119391	GSTO2
GBT0570	30154	gli2a	2736	GLI2
GBT0572	568184	lrp2a	4036	LRP2
GBT0585	30256	kita	3815	KIT
GBT0591	100009635	dhx37	57647	DHX37
GBT0593	393716	skp1	6500	SKP1
GBT0599	768182	mcu	90550	MCU
GBT0700	558875	ttc23	64927	TTC23
GBT0710	564112	magi2a	9863	MAGI2
GBT0717	557363	taf6l	10629	TAF6L
GBT0722	569986	dip2cb	22982	DIP2C
GBT0726	100003611	radx	55086	RADX
GBT0734	100149334	si:dkey-15h8.15		
GBT0750	557081	col7a1	1294	COL7A1
GBT0757	555517	kat2a	2648	KAT2A
GBT0760	110437953	adgrl2b.1	23266	ADGRL2
GBT0776	796370	edar	10913	EDAR
GBT0785	393509	mboat7	79143	MBOAT7
GBT0795	562459	znf414	84330	ZNF414
GBT0906	393638	timmm50	92609	TIMM50
GBT0926	790941	dele1	9812	DELE1
GBT0936	541552	eya4	2070	EYA4
GBT0941	100333571	zmiz2	83637	ZMIZ2
GBT0951	393521	macrod2	140733	MACROD2
GBT0959	562529	flna	2316	FLNA
GBT0966	563806	plpp2a	8612	PLPP2
GBT0972	114446	dlg1	1739	DLG1
GBT0978	568996	zfpm2a	23414	ZFPM2
GBT0980	559150	hmcn1	83872	HMCN1
GBT0993	557073	cdip1	29965	CDIP1
GBT1023	322106	csnk1da	1453	CSNK1D
GBT1027	557315	kirrel3l	84063	KIRREL2
GBT1033	555701	rreb1b	6239	RREB1

Table 3\_Human orthologues of GBT confirmed genes updated Apr8-2019.xlsx



GBT1042	562146	pard3bb	117583	PARD3B
GBT1093	100126126	ryr2a	6262	RYR2
GBT1105	246222	tbx15	6913	TBX15
GBT1129	560875	cacna1ha	8912	CACNA1H
GBT1248	100148041	tnk2a	10188	TNK2
GBT1259	100537272	LOC100537272		
GBT1278	735249	selenos	55829	SELENOS
GBT1300	386856	pdgfra	5156	PDGFRA

## Table 3

GBT	Tagged Gene	Human Orthologue	Human Disease	Disease Models (Number of models)	Reference (PMID)
GBT0016	pbx1a	PBX1	CONGENITAL ANOMALIES OF KIDNEY AND URINARY TRACT SYNDROME WITH ORWITHOUT HEARING LOSS ABNORMAL EARS OR DEVELOPMENTAL DELAY	mouse (1)	12591246
GBT0031	tnnt2a	TNNT2	CARDIOMYOPATHY DILATED 1D	mouse (5)	18606313, 27936050, 17556660, 18349139
GBT0031	tnnt2a	TNNT2	CARDIOMYOPATHY FAMILIAL HYPERTROPHIC 2	mouse (9)	16326803, 9788962, 11171784, 18349139, 10449439, 9637714, 23532597
GBT0078	grip1	GRIP1	FRASER SYNDROME 3	mouse (2)	14730302, 16880404
GBT0131	ddb2	DDB2	XERODERMA PIGMENTOSUM COMPLEMENTATION GROUP E	mouse (4)	14769931, 15558025
GBT0133	zic2a	ZIC2	HOLOPROSENCEPHALY 5	mouse (3)	27466203, 18617531, 10677508
GBT0135	bhlhe41	BHLHE41	SHORT SLEEPER	mouse (1)	19679812
GBT0156	fras1	FRAS1	FRASER SYNDROME 1	mouse (5)	12766769, 12766770, 24143185, 15623520, 15838507
GBT0178	ada	ADA	SEVERE COMBINED IMMUNODEFICIENCY AUTOSOMAL RECESSIVE TCELL-NEGATIVE B CELL-NEGATIVE NK CELL- NEGATIVE DUE TO ADENOSINEDEAMINASE DEFICIENCY	mouse (1)	9478961
GBT0186	cacna1c	CACNA1C	TIMOTHY SYNDROME	mouse (1)	21878566
GBT0240	bcl11ba	BCL11B	IMMUNODEFICIENCY 49	zebrafish (1)	
GBT0251	foxl2a	FOXL2	Blepharophimosis ptosis and epicanthus inversus	mouse (3)	15056605, 14736745, 24565867
GBT0268	ano5a	ANO5	MUSCULAR DYSTROPHY LIMB-GIRDLE AUTOSOMAL RECESSIVE 12	mouse (1)	26911675
GBT0270	zfpm2a	ZFPM2	Tetralogy of Fallot	mouse (1)	10892744, 12223418
GBT0270	zfpm2a	ZFPM2	DIAPHRAGMATIC HERNIA 3	mouse (1)	16103912
GBT0275	col4a4	COL4A4	ALPORT SYNDROME AUTOSOMAL RECESSIVE	mouse (4)	21196518, 24522496
GBT0348	ryr1b	RYR1	CENTRAL CORE DISEASE OF MUSCLE	mouse (3)	25564733, 19959667, 7515481
GBT0396	smg5	SMG5	HEART AND BRAIN MALFORMATION SYNDROME	mouse (1)	27018474
GBT0409	npr2	NPR2	Acromesomelic dysplasia Maroteaux type	mouse (2)	23065701, 17728275
GBT0411	dnajb6b	DNAJB6	autosomal dominant limb-girdle muscular dystrophy type 1	zebrafish (1-NOT), mouse (1)	26362252
GBT0422	insrb	INSR	DONOHUE SYNDROME	mouse (2-NOT)	
GBT0511	oca2	OCA2	ALBINISM OCULOCUTANEOUS TYPE II	rabbit (1)	Magnussen, K. (1952).
GBT0572	lrp2a	LRP2	DONNAI-BARROW SYNDROME	mouse (1)	20653565
GBT0585	kita	KIT	MASTOCYTOSIS CUTANEOUS	mouse (5)	24788138, 21148330
GBT0585	kita	KIT	Piebald trait	mouse (1)	20095975
GBT0585	kita	KIT	Gastrointestinal stromal tumor	mouse (5)	16061643, 22652566, 18098338, 12754375
GBT0710	magi2a	MAGI2	NEPHROTIC SYNDROME TYPE 15	mouse (1)	25271328
GBT0750	col7a1	COL7A1	recessive dystrophic epidermolysis bullosa	mouse (4)	18382769, 19893033, 10523500
GBT0776	edar	EDAR	hypohidrotic ectodermal dysplasia	mouse (1)	10431242, 17148670, 9799835
GBT0776	edar	EDAR	HAIR MORPHOLOGY 1	mouse (1)	18561327
GBT0785	mboat7	MBOAT7	MENTAL RETARDATION AUTOSOMAL RECESSIVE 57	mouse (1)	23097495
GBT0959	flna	FLNA	Periventricular nodular heterotopia 1	mouse (1)	16825286
GBT1023	csnk1da	CSNK1D	ADVANCED SLEEP PHASE SYNDROME FAMILIAL 2	mouse (2)	15800623
GBT1093	ryr2a	RYR2	Ventricular tachycardia catecholaminergic polymorphic 1 with or without atrial dysfunction and/or dilated cardiomyopathy	mouse (10)	27482086, 25775566, 26121139, 24755079, 23152493, 22828895, 20224043, 18419777, 16873551, 15890976
GBT1093	ryr2a	RYR2	ARRHYTHMOGENIC RIGHT VENTRICULAR DYSPLASIA FAMILIAL 2	mouse (1)	16873551
GBT1105	tbx15	TBX15	COUSIN SYNDROME	mouse (1)	19068278

Table 4

Line	Vector	Genome location	Genome version
GBT0115	RP2	chr8:12923125-12923133	ZV9
GBT0129	RP2	chr12:45451589-45451597	ZV9
GBT0148	RP2	chr16:55188302-55188310	ZV9
GBT0264	RP2	chr5:49960612-49960620	ZV9
GBT0506	RP2	chr18:46186350-46186358	ZV9
GBT0573	RP2	chr6:33630487-33630495	GRCz10
GBT0586	RP2	chr13:9112806-9112814	GRCz10
GBT0702	RP2	chr13:9112806-9112814	GRCz10
GBT0724	RP8	chr7:7633071-7633079	GRCz10
GBT0994	RP2	chr3:39627361-39627369	GRCz10
GBT1024	RP2	chr25:5175242-5175250	GRCz10
GBT1071	RP8	chr9:21818261-21818269	GRCz10
GBT1100	RP2	chr5:25509243-25509251	GRCz10
GBT1116	RP2	chr14:14724098-14724106	GRCz10
GBT1168	RP2	chr5:39664405-39664412	GRCz10

## Table 5