

## **Genomic profiling of Collaborative Cross founder mice infected with respiratory viruses reveals novel transcripts and infection related strain-specific gene and isoform expression**

Hao Xiong<sup>\*,†</sup>, Juliet Morrison<sup>\*,†</sup>, Martin T. Ferris<sup>†,§</sup>, Lisa E. Gralinski<sup>†,‡</sup>, Alan C. Whitmore<sup>†,§</sup>, Richard Green<sup>\*,†</sup>, Matthew J. Thomas<sup>\*,†</sup>, Jennifer Tisoncik-Go<sup>\*,†</sup>, Gary P. Schroth<sup>\*\*</sup>, Fernando Pardo-Manuel de Villena<sup>§</sup>, Ralph S. Baric<sup>†,‡</sup>, Mark T. Heise<sup>†,§</sup>, Xinxia Peng<sup>\*,†</sup>, and Michael G. Katze<sup>\*,†,1</sup>

<sup>\*</sup>Department of Microbiology, School of Medicine, University of Washington, Seattle, Washington, United States of America

<sup>†</sup>Pacific Northwest Regional Center of Excellence for Biodefense and Emerging Infectious Diseases Research, Portland, Oregon, United States of America

<sup>‡</sup>Department of Epidemiology, University of North Carolina-Chapel Hill, Chapel Hill, North Carolina, United States of America

<sup>§</sup>Department of Genetics, University of North Carolina-Chapel Hill, Chapel Hill, North Carolina, United States of America

<sup>\*\*</sup>Illumina, Inc., San Diego, California, United States of America

<sup>1</sup>To whom correspondence should be addressed.

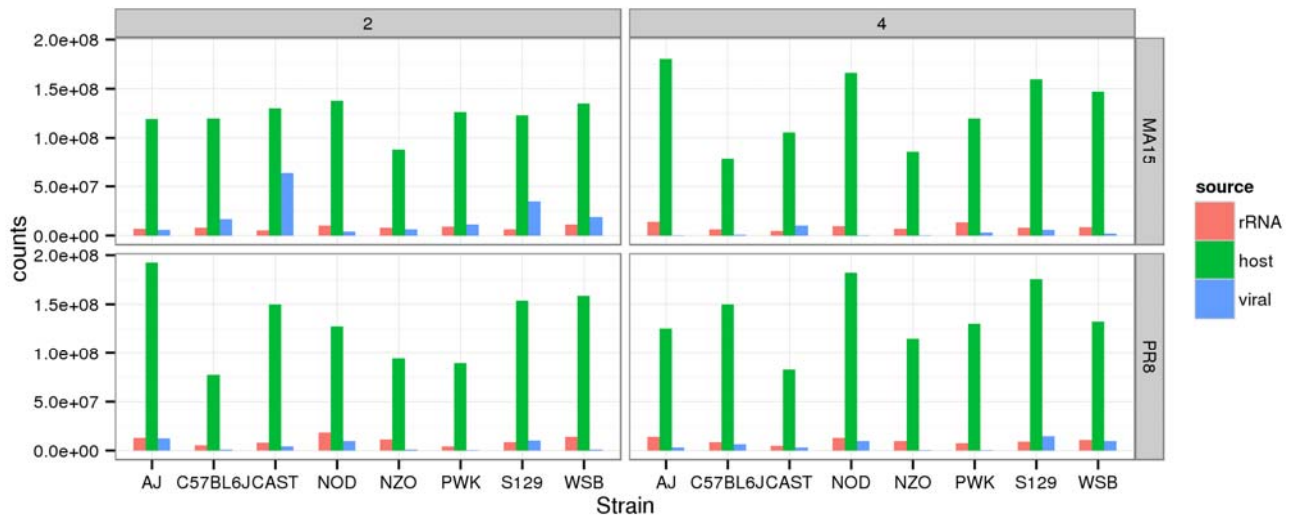
Data Access: GSE52405

**DOI: 10.1534/g3.114.011759**

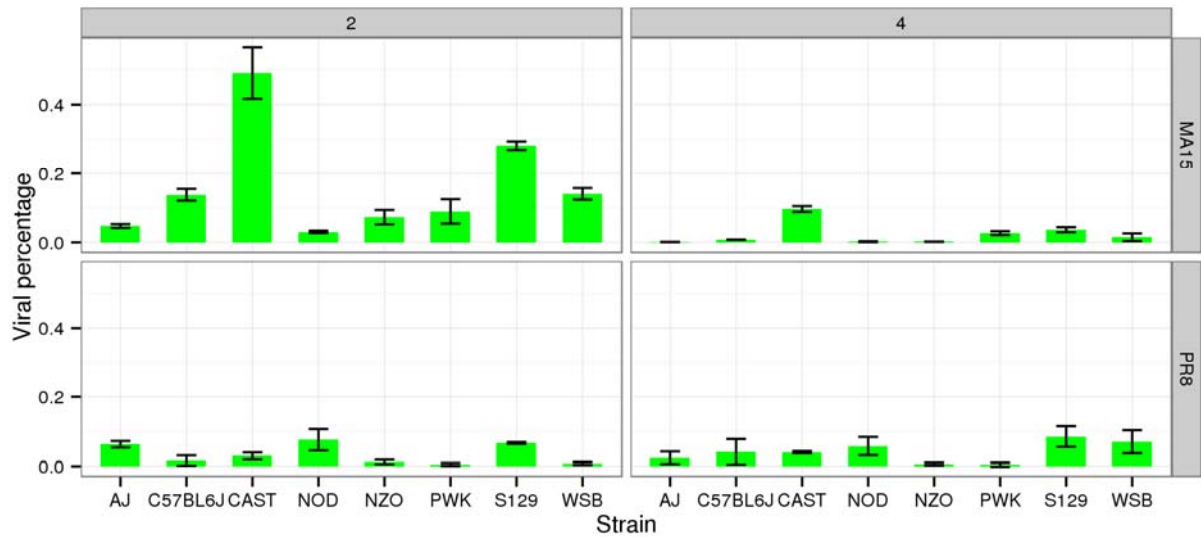
## File S1

### Transcript discovery pipeline

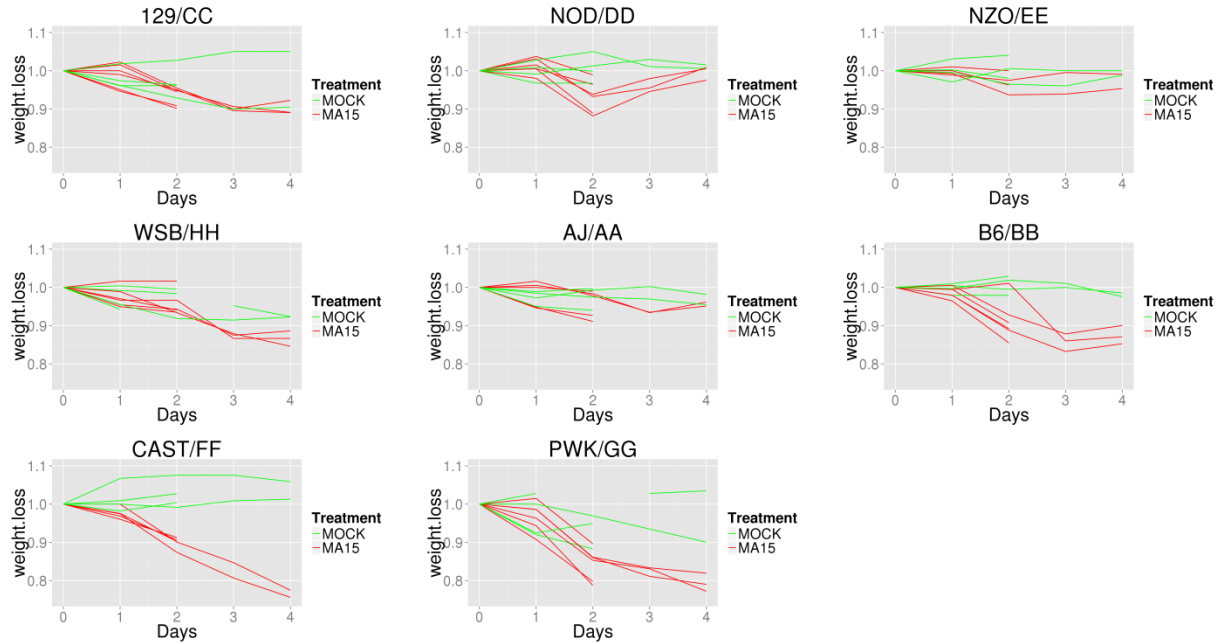
The transcript discovery pipeline accomplishes two tasks, adapting existing gene models to the different founder genomes and uncovering new transcripts not overlapping with any known transcripts. To use the most accurate genome for each founder we downloaded the pseudogenomes that had been generated by University of North Carolina based on DNA-seq data produced by the Sanger Institute. This incorporates single-nucleotide variations (SNV), indels, and other structural variations, which can number in the millions, and represents our most current description of the founder genomes. The reference annotation describes over thirty thousand genes, of which about twenty four thousand genes are coding genes. We mapped short reads to their respective pseudogenomes and used the Cufflinks program to predict new transcripts. To avoid confounding intronic reads, we retained only intergenic transcripts. We also took unmapped reads and used Trinity to discover *de novo* transcripts. After filtering *de novo* transcripts that bore resemblance to viral or host genomes, we combined the Trinity output that mapped to mouse genomes with the Cufflinks output, while the stand alone *de novo* transcripts were further filtered to identify those that were similar to known mouse, rat, or human sequences. The outline of our novel transcript discovery pipeline built is shown in Figure S5.



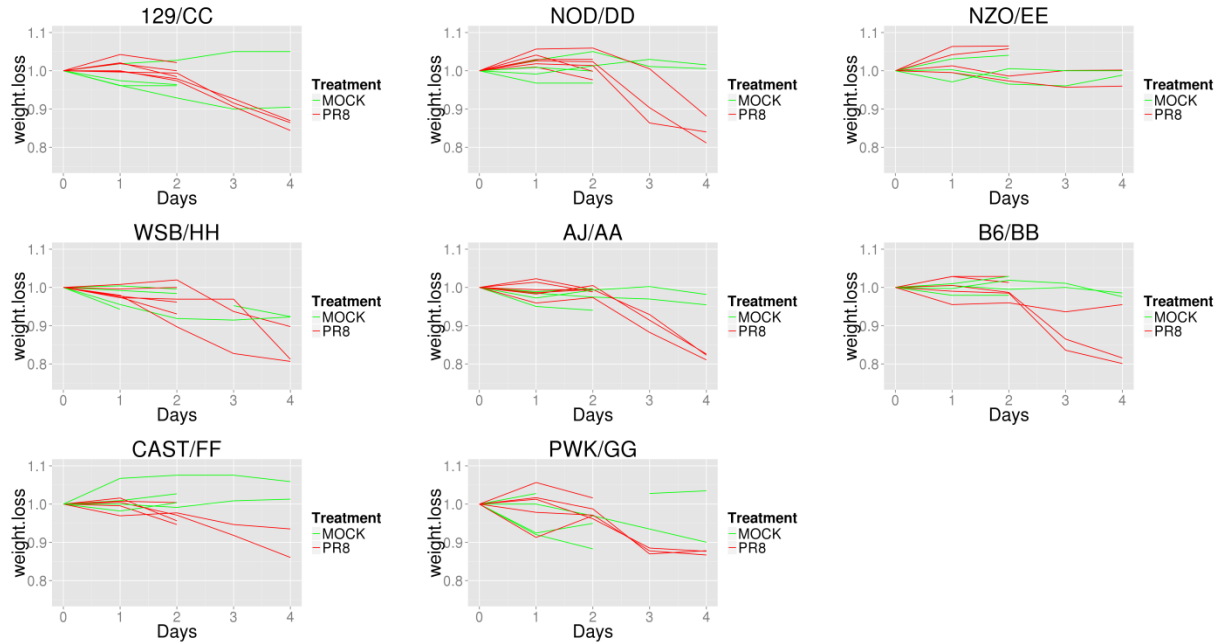
**Figure S1** The number of short RNA reads that map to host, ribosomes, and viruses. The read counts were summed by the unique combination of mouse strain, virus, and day post infection. We generally ensured each sample has at least 30 million host reads. The varying sequencing depth for ribosomal and host RNA among strains resulted from random sequencing variation. On the other hand, the viral read counts exhibited strain specificity. For example, MA15-infected CAST mice had 60M viral reads on day 2, while MA15-infected NZO mice had about 6M viral reads, which pointed to strain differences in viral replication.



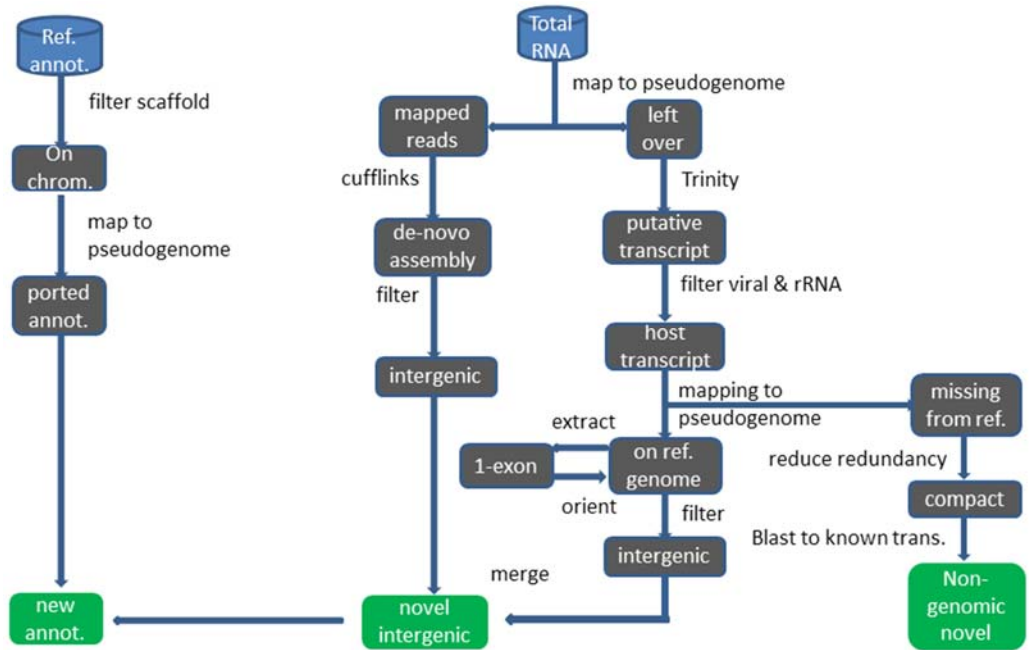
**Figure S2 The percentage of RNA reads that belong to MA15 or PR8 genomes.** The percentages were calculated by summing all viral reads from animals infected with one virus type for each founder strain at different time points (2 and 4 days post infection) and divided by the total read counts. Each strain had at least two replicates and most have three replicates. There are clear strain-specific differences in the viral read percentages, ranging from almost zero percent of RNA reads in MA15-infected AJ mice to close to half of reads in MA15-infected CAST mice on day 2 post infection.



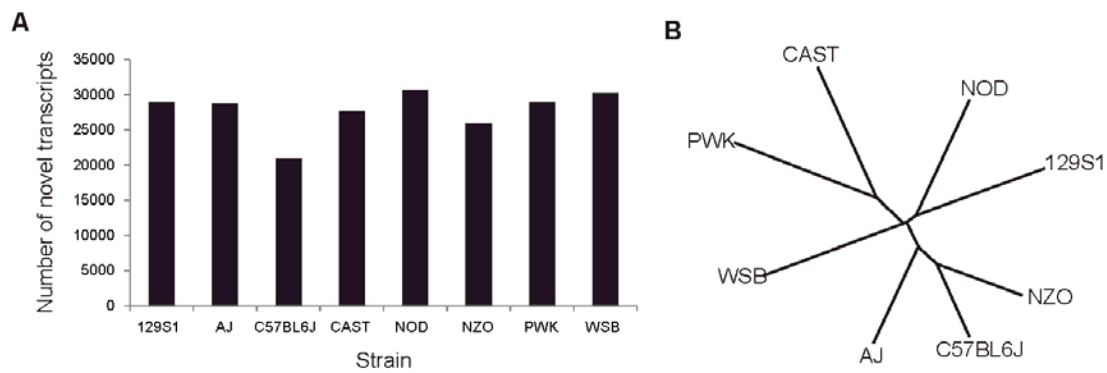
**Figure S3 Weight loss of founder mice infected with MA15.** Mice were infected with  $10^5$  PFU of MA15 and monitored for weight loss for 4 days. Half of the infected animals were sacrificed for expression profiles at day 2 post- nfection, while the second half were sacrificed at day 4 post infection. In general, there were observable strain differences in weight loss. For example, NZO lost less than 5% weight, while NOD and AJ mice were able to regain their lost weight by day 4. In contrast, the remaining mouse strains sustained about 10% weight loss, with the CAST and PWK strains showing close to 20% weight loss.



**Figure S4 Weight loss of founder mice infected with PR8.** Mice were infected with  $5 \times 10^2$  PFU of PR8 and monitored for weight loss for 4 days. Half of the infected animals were sacrificed for expression profiles at day 2 post infection, while the second half were sacrificed at day 4 post infection. There was little weight loss and no discernable strain difference on day 2 post infection. On day 4, NZO had little to no weight loss while all other strains have appreciable weight loss, with most samples showing 10% loss. Half of infected animals were sacrificed for expression profiles at day 2 post-infect while the second half was sacrificed at day 4.

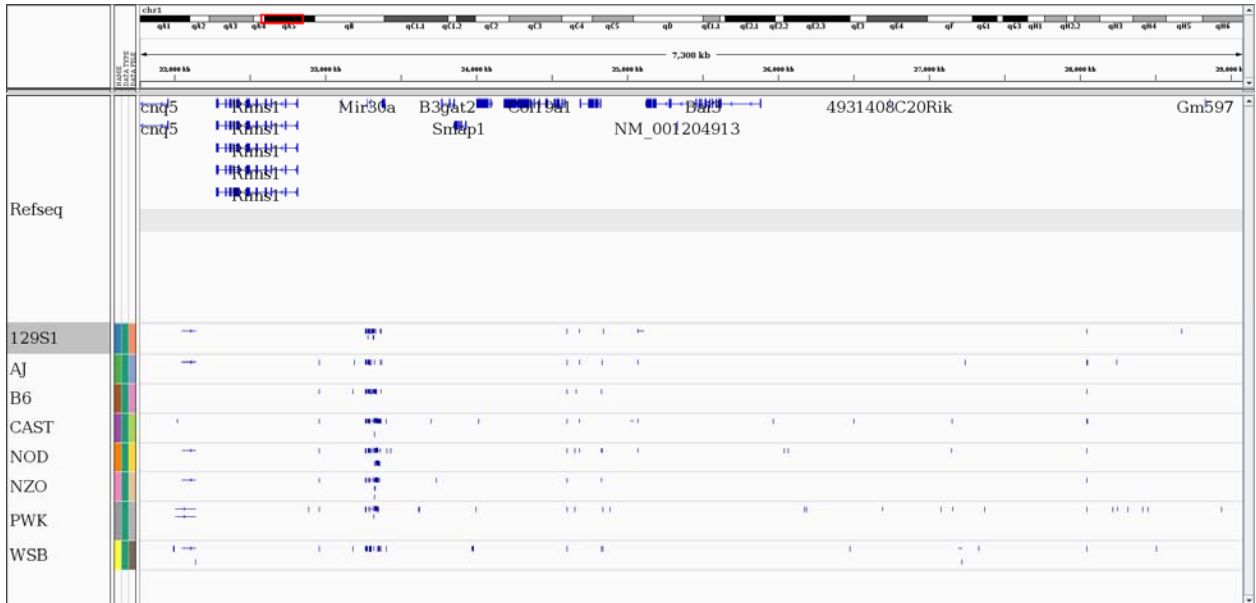


**Figure S5 Schematic of the annotation pipeline.** The left branch adapts the reference annotation to eight founders' genomes, since the reference annotation is based on C57BL/6J mouse's genome. Before adapting the reference annotation, genes not on chromosomes but on scaffolds are filtered out because by definition they cannot be placed on chromosomes. The right branch is for novel transcript discovery. There are two complementary approaches to discovering transcripts: genome-based and non-genome-based. Cufflinks, a *de novo* assembler, was used for the genome-based approach, and Trinity, another *de novo* assembler, does not require a genome for novel transcripts discovery. After several post-processing steps (see Methods) are conducted with Trinity, output is merged with Cufflinks' output to obtain novel transcripts. An extra step of filtering viral and ribosomal sequences was added to remove any residual viral and ribosomal sequences due to the weaker mappability of short reads compared to the assembled contigs. Using this approach some *de novo* transcripts could not be placed on existing mouse genomes. After reducing redundancy between *de novo* transcripts, we filter out transcripts similar to sequences on scaffolds and organisms other than human and mouse. There is also an additional stage for orienting single-exon transcripts. The final product of this pipeline was a new annotation with novel transcripts.

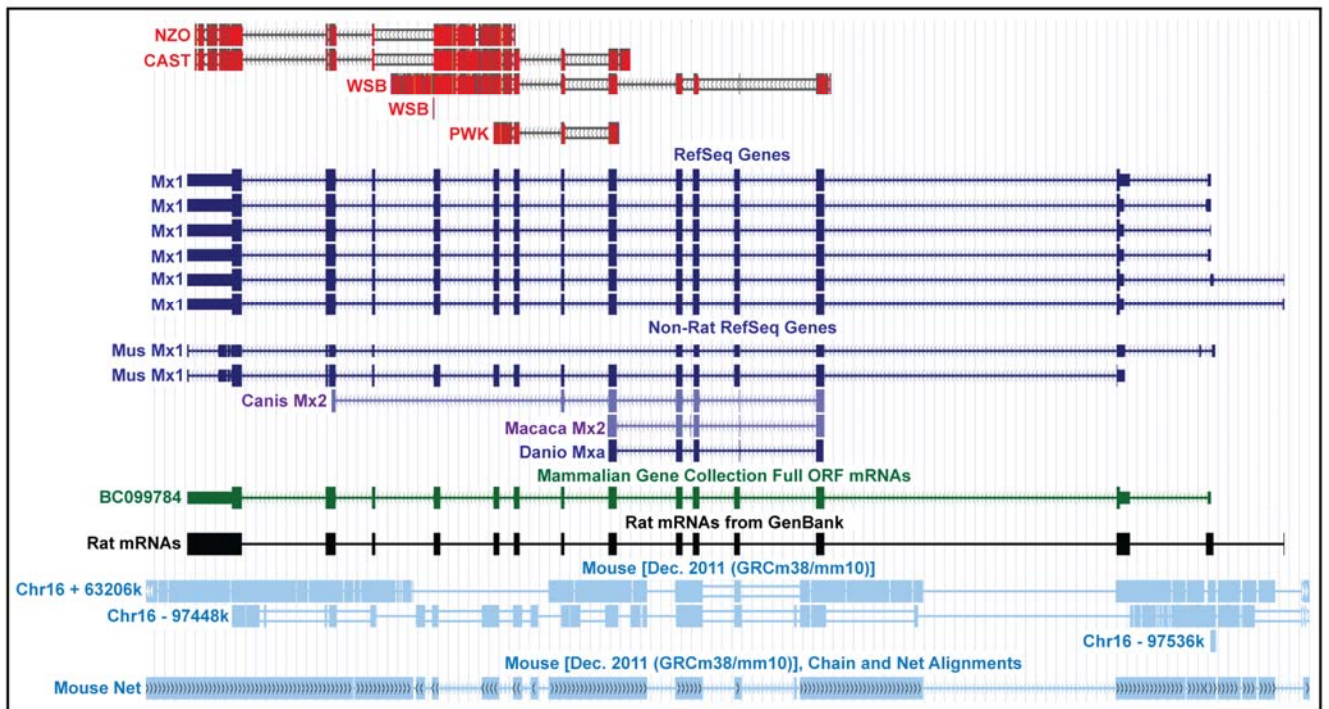


**Figure S6 Summary of novel transcripts in eight founder strains and their relationships.** (A) the number of novel intergenic novel transcripts in the eight founders. C57BL/6J had the fewest new transcripts because the reference annotation is based on data on the same strain and therefore is most complete for this strain. Other strains had more new intergenic transcripts but wild-derived strains did not appear to possess more than laboratory strains (except C57BL/6J). (B) clustering of founder strains by chromosome 2 novel-transcript density. The clustering results are to the phylogeny of the eight mouse strains. PWK and CAST are separate from all other six strains while WSB is the wild-derived strain closest to the classical laboratory strains.

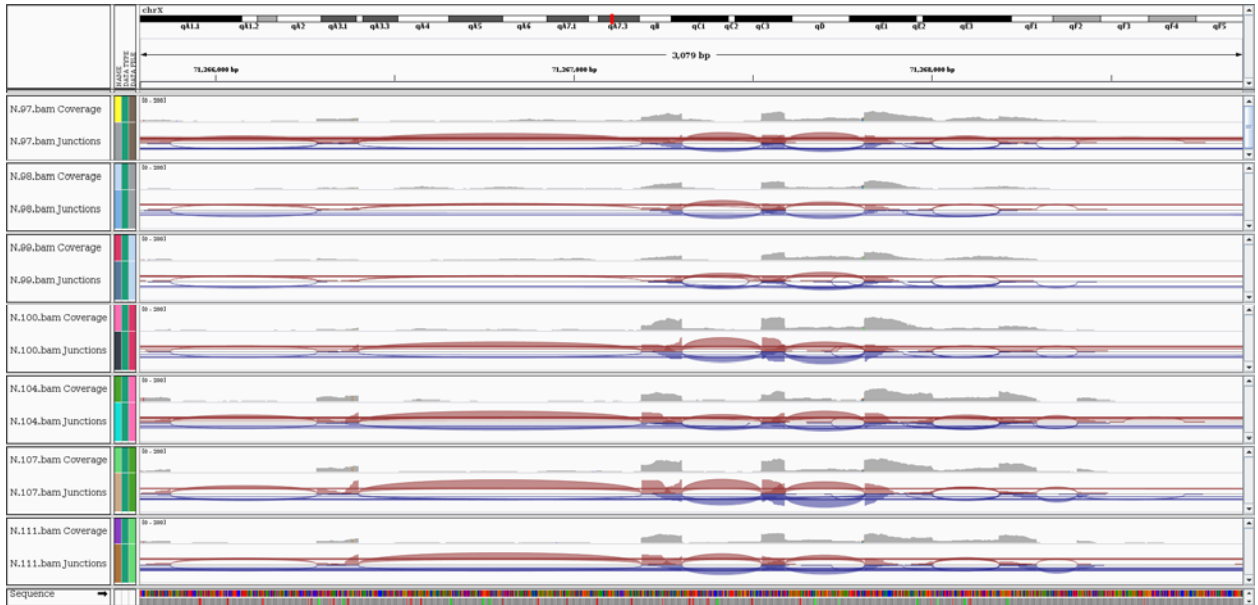




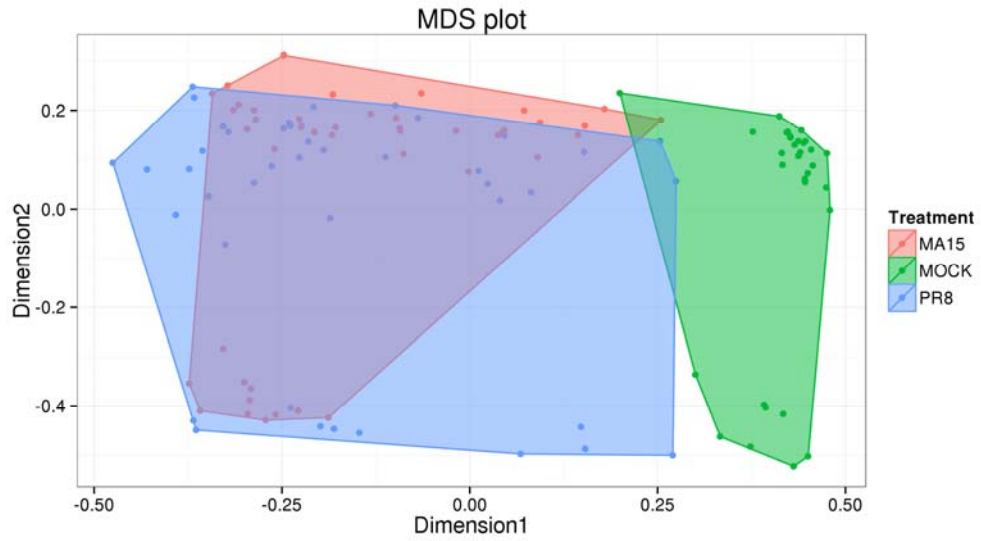
**Figure S7** An IGV view of annotated and novel transcripts within Hrl3 (chr1:21,767,867–29,085,401) which was found to be associated with pulmonary edema in mice infected with PR8 (Ferris et al. 2013). The top half are the annotated genes from the reference annotation while the rest shows novel transcripts in eight founders. The cluster of novel transcripts near Mir30a could herald previously unknown genes or considerable extension of nearby genes. It is also easy to similarity between the eight founders as well as divergence.



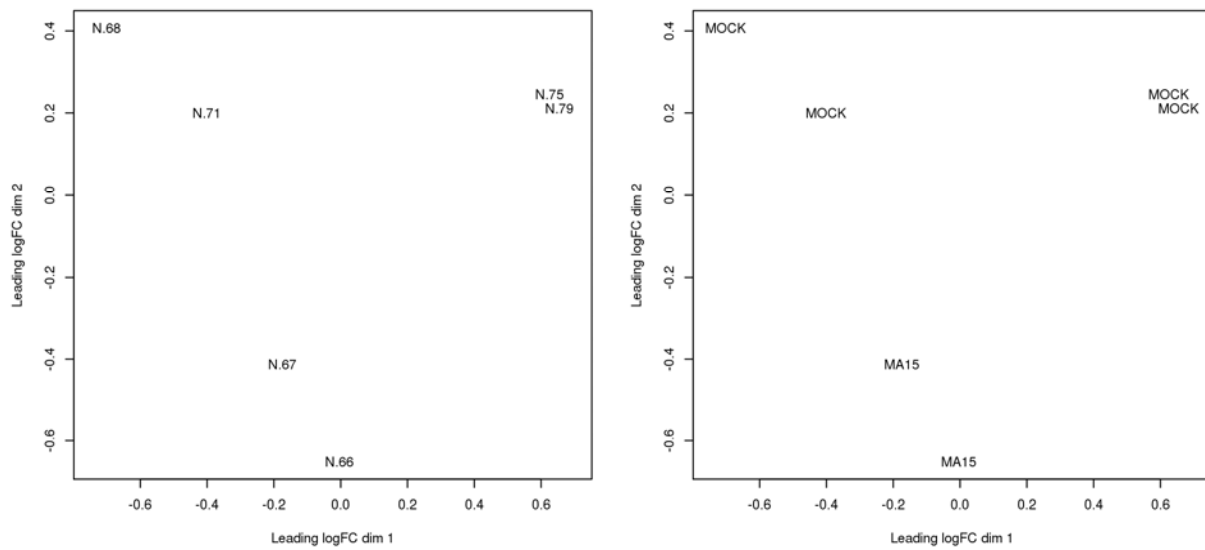
**Figure S8 Alignment of de novo transcript with affinity to four extra *Mx1* exons.** Ferris et al reported that four founder strains have extra exons in *Mx1* gene compared to the reference annotation. We therefore mapped our novel transcripts against these extra exons and found four novel transcripts containing all of parts of these exons in the four previously reported strains. Mapping the four novel transcripts to rat genome found was successful and shown here. The four transcripts were found in the four strains that Ferris et al originally discovered the extra exons. The other four strains didn't have these missing exons and therefore had no novel transcript that could be mapped here.



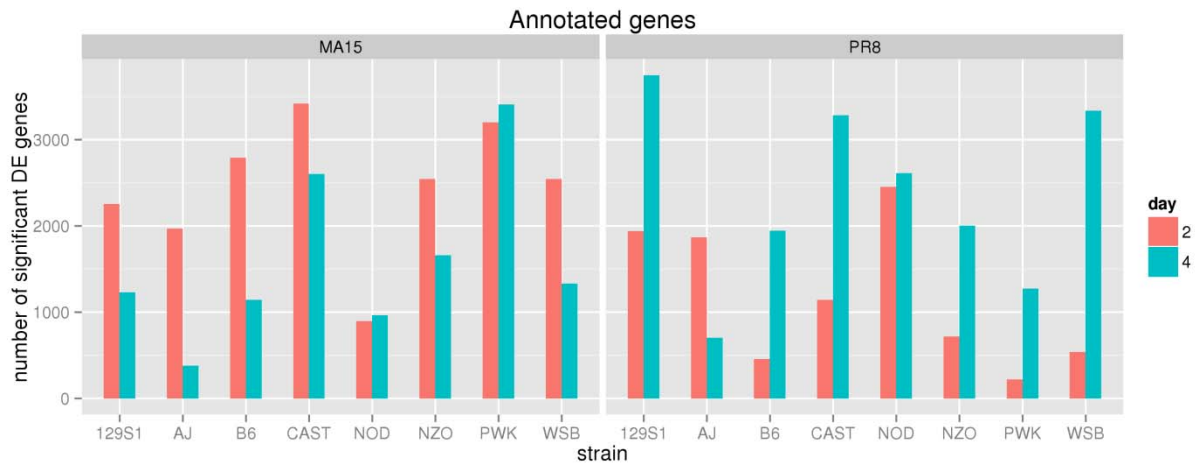
**Figure S9** Differential isoform expression of *Irak1* gene in MA15-infected CAST mice at day 2 post- infection by MA15. The rightmost splicing junction was differentially expressed: the infected samples (top three) had much lower level of expression than mock samples (bottom four). The gene overall was not differentially expressed as can be seen in the coverage graphs.



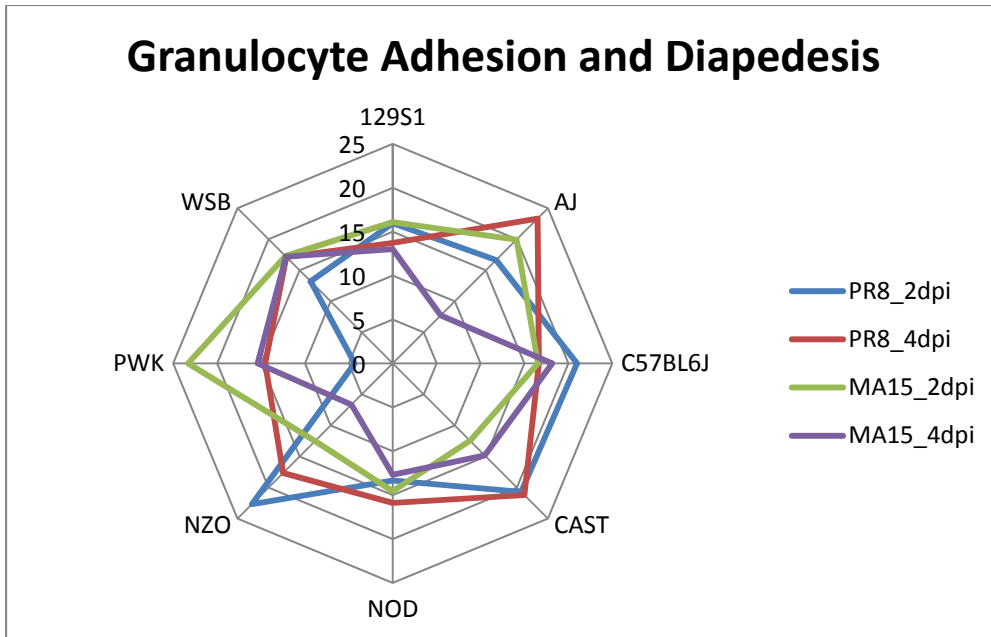
**Figure S10** MDS plot of all 119 lung samples that were analyzed for differential expression. Mock samples are colored green, MA15 infection samples are red, and PR8 infection samples are blue. The same colors are also used to define the space spanned by the respective samples.



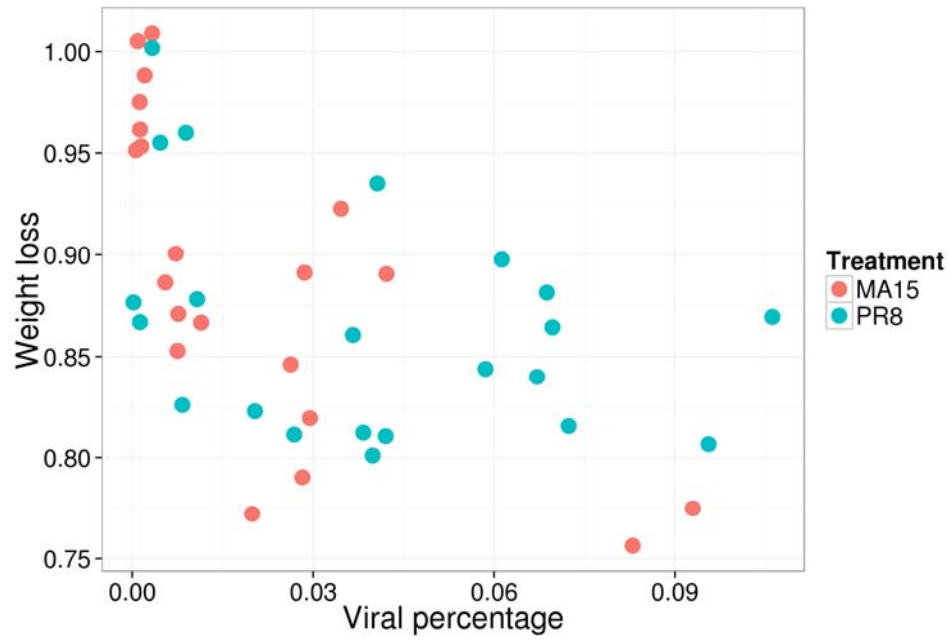
**Figure S11** The MDS plot of six samples for AJ strain infected with MA15 at day 2 post infection. The two subplots have different text but are based on the same data. MDS plot is used to gauge possible batch-effect and separation between mocks and infected samples, which are well separated.



**Figure S12 Differentially expressed genes among the eight founder mouse strains following either influenza or SARS-CoV infection.** The counts for two viruses were separated into two panels, the left panel for MA15 and the right panel for PR8. Red bars depict DE genes for day 2 post-infection and blue bars depict DE genes for day 4 post-infection.



**Figure S13 Radial plot showing functional enrichment of granulocyte adhesion and diapedesis pathway genes.** This pathway is related to a vascular cuffing phenotype. The strain difference was quite pronounced for both viruses and time points except PR8 on day 4 post-infection. At day 2 post influenza infection the biggest contrast is between PWK and NZO strain, between no enrichment and extreme enrichment. For SARs infection, the difference isn't so big but they were clearly visible. At day 2 post infection, NZO and PWK strains showed mild and strong enrichment, respectively; and at day 4 NZO and A/J strains had less enrichment than other strains.



**Figure S14 Percent weight loss correlates with viral read percentages.** The viral read percentages were plotted against day 4 weight loss. All infected samples were included, with the two viral infections represented by different colored points.



**Table S1** The number of splice junctions that were discovered in one strain but not in another strain. One strain (in the row heading) is compared to another strain (in the column heading). For example, there are 3515 junctions in 129S1 that are absent from AJ.

	<b>129S1</b>	<b>AJ</b>	<b>B6</b>	<b>CAST</b>	<b>NOD</b>	<b>NZO</b>	<b>PWK</b>	<b>WSB</b>
<b>129S1</b>	0	3515	10824	14581	3279	7901	8250	7576
<b>AJ</b>	6507	0	13508	17474	4980	9888	9294	9925
<b>B6</b>	1421	1113	0	7877	974	2526	4363	2976
<b>CAST</b>	3156	3057	5855	0	2957	4556	4746	4101
<b>NOD</b>	5309	4018	12407	16412	0	8954	9255	9108
<b>NZO</b>	3187	2182	7215	11267	2210	0	7140	5255
<b>PWK</b>	11641	9693	17157	19562	10616	15245	0	14535
<b>WSB</b>	3154	2511	7957	11104	2656	5547	6722	0

**Table S2** The number of splice junctions that are observed in infected samples but not in mock samples. Thresholds were set at 10 reads for infected samples and 5 reads for mock samples.

	<b>MA15</b>		<b>PR8</b>	
	<b>Day 2</b>	<b>Day4</b>	<b>Day 2</b>	<b>Day 4</b>
<b>129S1</b>	2114	1868	2854	4504
<b>AJ</b>	4615	1818	2674	788
<b>B6</b>	2344	757	371	7877
<b>CAST</b>	1487	2666	2078	3934
<b>NOD</b>	2294	3173	1825	2496
<b>NZO</b>	1844	2827	4353	3562
<b>PWK</b>	4675	4093	429	2455
<b>WSB</b>	2842	2377	6359	4616

**Table S3 The number of differentially expressed splicing junctions.** Here a DE junction must be outside DE genes to disambiguate isoform DE from gene DE. There is one order of magnitude of difference between different strains or time points. We took a decidedly conservative approach. There likely are more differentially expressed splicing junctions than we found.

Strain	MA15		PR8	
	Day 2	Day 4	Day 2	Day 4
<b>129S1</b>	453	326	219	720
<b>AJ</b>	25	237	306	174
<b>B6</b>	817	183	42	628
<b>CAST</b>	855	581	192	783
<b>NOD</b>	206	282	521	438
<b>NZO</b>	612	184	101	268
<b>PWK</b>	692	976	64	235
<b>WSB</b>	964	276	75	822

**Table S4** The replicate size for each combination of mouse strain, virus, and days-post-infection. Note that mocks are pooled for all DE analysis.

<b>Strain</b>	<b>Virus</b>	<b>Day</b>	<b>Sample size</b>
<b>129S1</b>	MA15	2	3
<b>129S1</b>	MA15	4	3
<b>129S1</b>	MOCK	2	2
<b>129S1</b>	MOCK	4	2
<b>129S1</b>	PR8	2	3
<b>129S1</b>	PR8	4	3
<b>AJ</b>	MA15	2	3
<b>AJ</b>	MA15	4	2
<b>AJ</b>	MOCK	2	2
<b>AJ</b>	MOCK	4	2
<b>AJ</b>	PR8	2	3
<b>AJ</b>	PR8	4	3
<b>B6</b>	MA15	2	3
<b>B6</b>	MA15	4	3
<b>B6</b>	MOCK	2	2
<b>B6</b>	MOCK	4	2
<b>B6</b>	PR8	2	2
<b>B6</b>	PR8	4	3
<b>CAST</b>	MA15	2	3
<b>CAST</b>	MA15	4	2
<b>CAST</b>	MOCK	2	2
<b>CAST</b>	MOCK	4	2
<b>CAST</b>	PR8	2	3
<b>CAST</b>	PR8	4	2
<b>NOD</b>	MA15	2	3
<b>NOD</b>	MA15	4	3
<b>NOD</b>	MOCK	2	2
<b>NOD</b>	MOCK	4	2
<b>NOD</b>	PR8	2	3
<b>NOD</b>	PR8	4	3
<b>NZO</b>	MA15	2	2
<b>NZO</b>	MA15	4	2
<b>NZO</b>	MOCK	2	2
<b>NZO</b>	MOCK	4	2
<b>NZO</b>	PR8	2	2
<b>NZO</b>	PR8	4	2
<b>PWK</b>	MA15	2	3
<b>PWK</b>	MA15	4	3
<b>PWK</b>	MOCK	2	2

---

<b>PWK</b>	MOCK	4	2
<b>PWK</b>	PR8	2	2
<b>PWK</b>	PR8	4	3
<b>WSB</b>	MA15	2	3
<b>WSB</b>	MA15	4	3
<b>WSB</b>	MOCK	2	2
<b>WSB</b>	MOCK	4	2
<b>WSB</b>	PR8	2	3
<b>WSB</b>	PR8	4	3

---

**Table S5 Viral read counts in the mock-infected founders.** Each strain's total comes from four mice. The counts are low, several orders of magnitude lower than infected samples. The fact we still get non-zero reads may be an artifact of mapping and select strains may have partial sequence similarity to viruses.

<b>Strain</b>	<b>MA15</b>	<b>PR8</b>
<b>129S1</b>	5084	11394
<b>AJ</b>	9549	2659
<b>B6</b>	521	1656
<b>CAST</b>	10490	3060
<b>NOD</b>	608	857
<b>NZO</b>	629	636
<b>PWK</b>	957	634
<b>WSB</b>	144	7491