

Cell Genomics, Volume 3

Supplemental information

**Transposable elements are associated
with the variable response to influenza infection**

Xun Chen, Alain Pacis, Katherine A. Aracena, Saideep Gona, Tony Kwan, Cristian Groza, Yen Lung Lin, Renata Sindeaux, Vania Yotova, Alben Pramatarova, Marie-Michelle Simon, Tomi Pastinen, Luis B. Barreiro, and Guillaume Bourque

SUPPLEMENTARY FIGURES

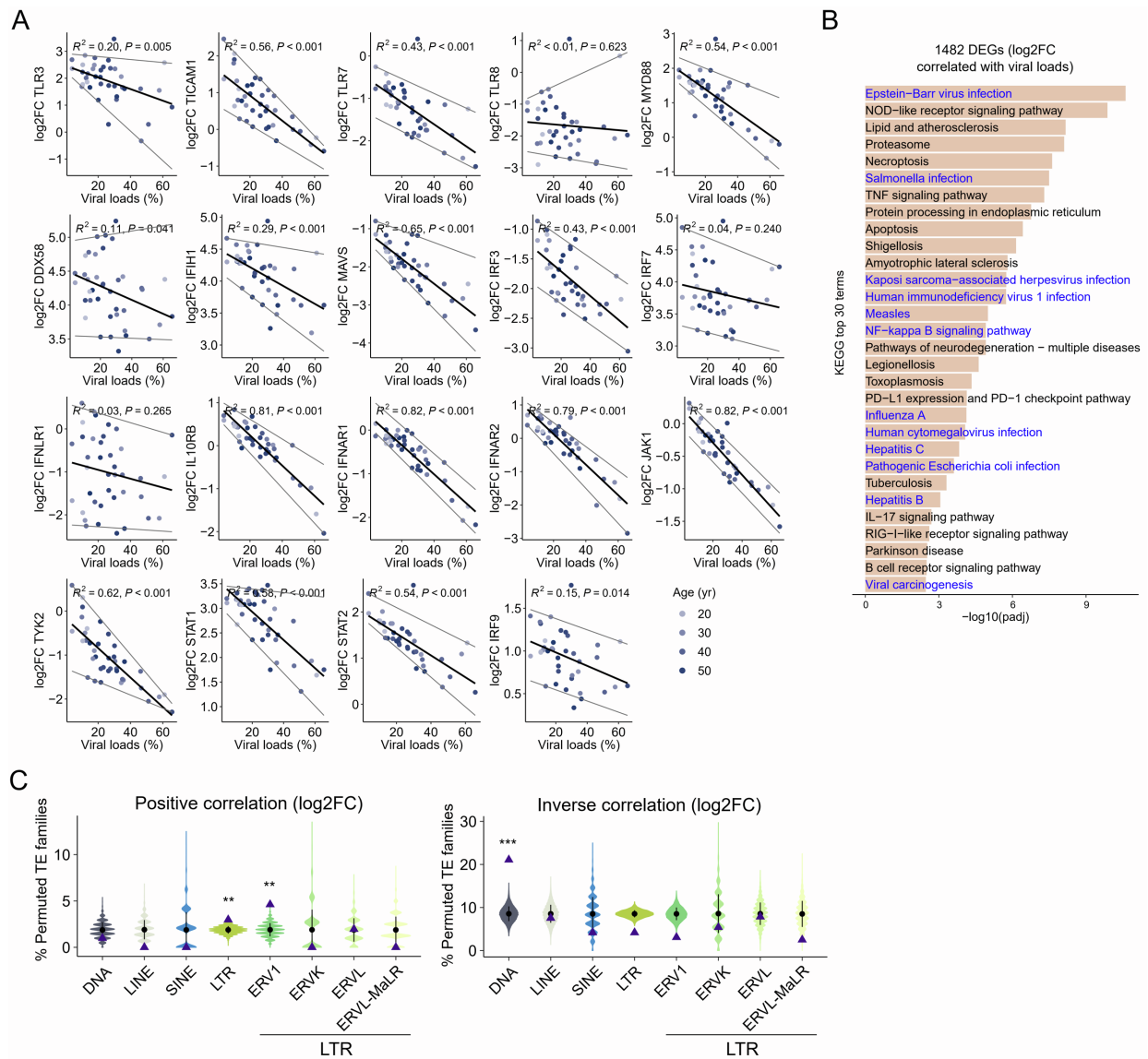


Figure S1. Immune genes and few TE families are prone to correlate with viral load post-infection, related to Figure 1

(A) Expression fold changes (log₂FCs) of most key immune regulators are inversely correlated with viral load post-infection. Many genes are shown to strongly correlate with viral load, including membrane-bound receptor genes sensing infection-induced interferons, *i.e.*, *IL10RB*, *IFNAR1*, *IFNAR2*, and *JAK1*. Linear regression model was used for the correlation analysis.

Black line represents the regression line and grey lines represent the 5% and 95% quantiles. **(B)** Expression fold change amongst differentially expressed genes (DEGs) correlated with viral load are enriched in multiple virus infection pathways. 1,482 DEGs with log₂FCs correlated with viral load ($R^2 \geq 0.3$, p value ≤ 0.05) were identified and used for downstream pathway enrichment analysis. **(C)** Enrichment of the proportion for log₂FC amongst TE families positively and inversely correlated with viral load post-infection. 17 positively and 77 inversely correlated TE families were analyzed. Purple triangle represents the actual proportion of correlated families among subclass or superfamily. Error bar represents the mean values and standard deviations of 10,000 randomized proportions. One-tailed student's t -test was used to compare the actual proportions with randomized proportions (* $p \leq 0.05$, ** $p \leq 0.01$, *** $p \leq 0.001$).

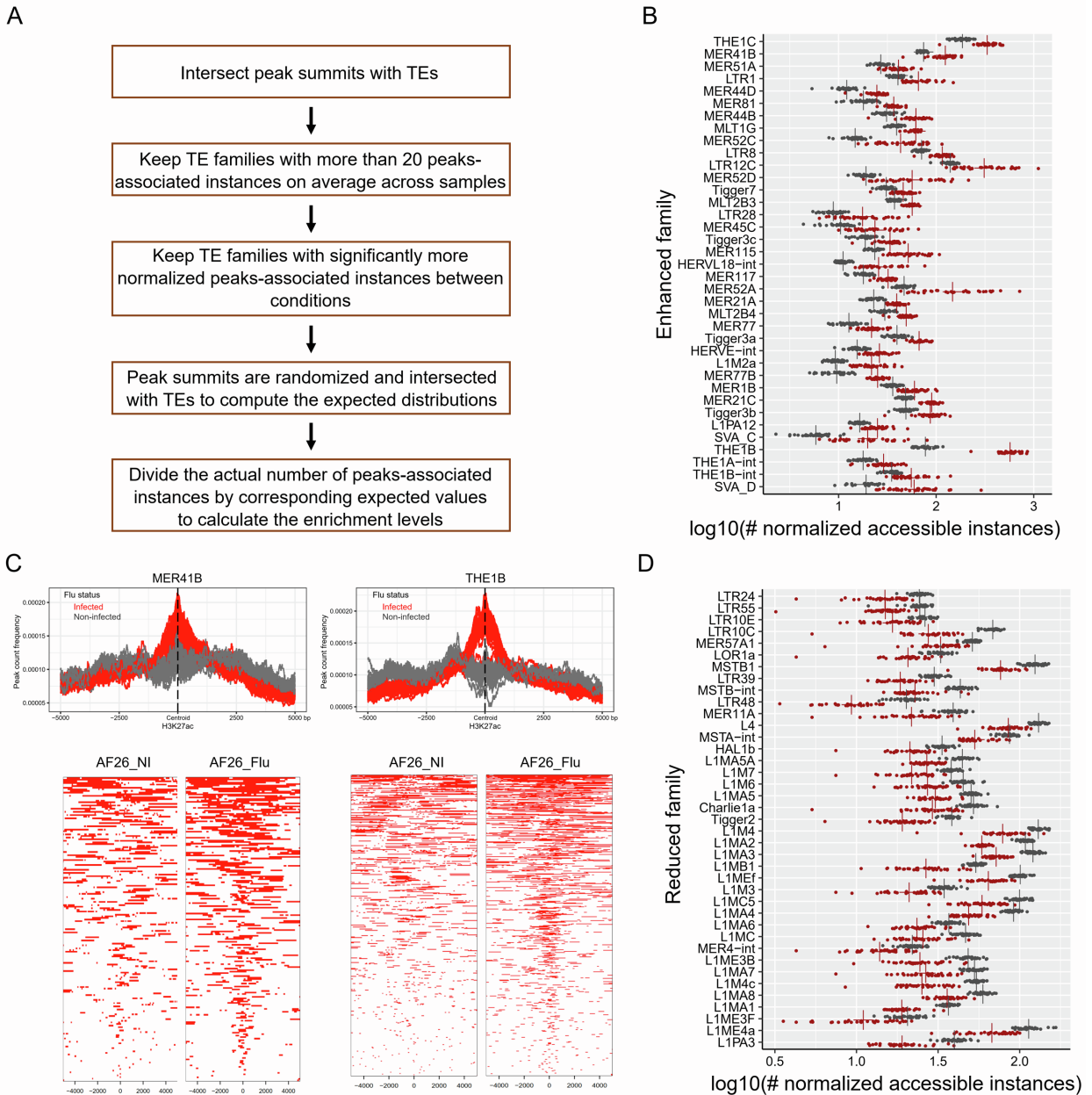


Figure S2. Detection of TE families with enhanced and reduced accessibility in response to IAV infection, related to Figure 2

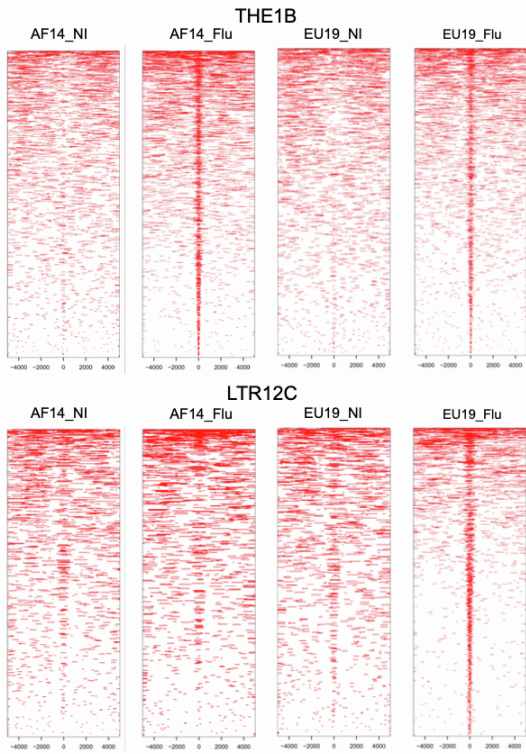
(A) Optimized TE enrichment analysis pipeline. Details were described in Methods. (B) Distributions of the normalized number of peaks-associated instances of TE families with enhanced accessibility. Each point represents one sample. Grey color represents the non-infected

sample and red color represents the infected sample. “+” indicates the mean value across non-infected (grey) and infected (red) samples. The number of accessible instances were normalized by the average number of peaks across infected and non-infected samples, respectively. **(C)**

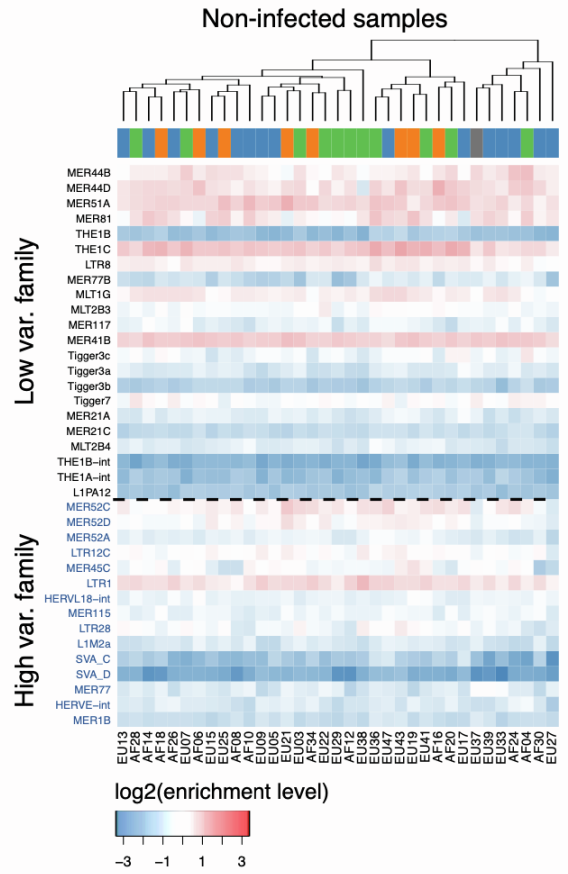
Average profiles (up) and heatmaps (bottom) of H3K27ac peaks at MER41B and THE1B (± 5 kb). H3K27ac peaks are centered at the median positions of peak summits across infected samples. Peak regions are shown as heatmaps at the bottom. AF26 non-infected and infected samples are shown as examples. Peak regions centered at each family are shown as red bars. **(D)**

Distributions of the normalized number of peaks-associated instances of TE families with reduced accessibility. Each point represents one non-infected (grey) and infected (red) sample. “+” indicates the mean value across non-infected (grey) and infected (red) samples. The number of accessible instances were normalized by the average number of peaks across infected and non-infected samples separately.

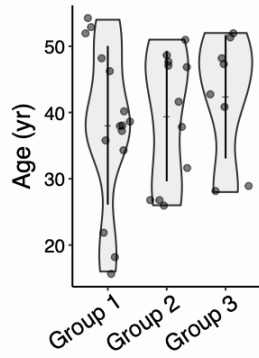
A



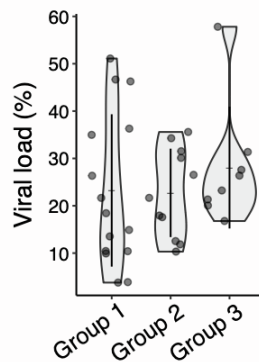
B



C



D



E

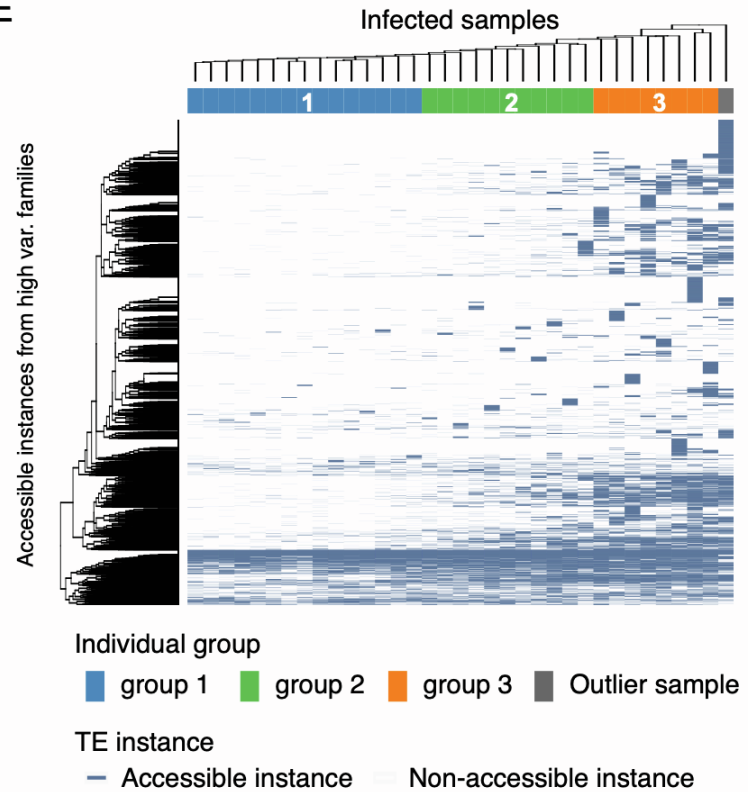
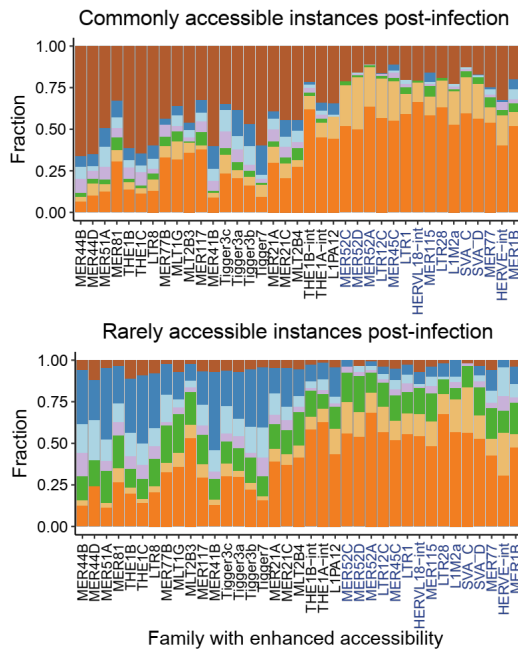


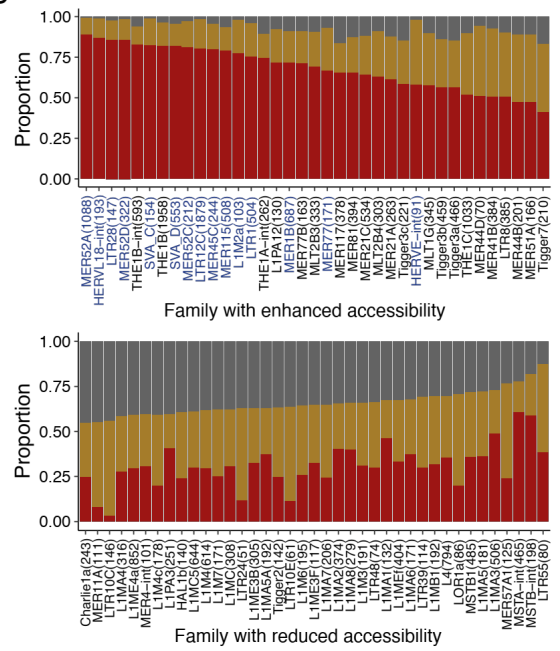
Figure S3. High variable families display high variability in chromatin accessibility post-infection, related to Figure 3

(A) Heatmaps of open chromatin regions at THE1B and LTR12C (± 5 kb). Two infected and non-infected samples are shown as examples. LTR12C shows a high variability of accessibility between randomly-selected samples post-infection. ATAC-seq peaks are centered at the summits in TEs. 5 kb upstream and downstream regions are shown. Compared to AF14, EU19 displays a higher enrichment at LTR12C but comparable enrichment at THE1B. (B) Heatmap of log₂ enrichment levels of 37 enhanced families in 35 non-infected samples. Semi-clustering analysis was performed. Three individual groups observed at infected samples are not clustered together. High variable families are highlighted in blue color. (C,D) Violin plots of age of macrophage donors and viral load of the three individual groups. Group 3 individuals have relatively older ages and higher viral load compared to group 1 individuals. The dot represents each individual and the error bar represents the mean value and standard deviation. (E) Heatmap of the chromatin state of accessible instances from high variable families in 35 infected samples. The state of open chromatin is in blue color and closed chromatin in white color. Unsupervised clustering analysis was performed, and the three individual groups are clustered together. A fraction of instances shows an enrichment in group 3 compared to group 1 individuals.

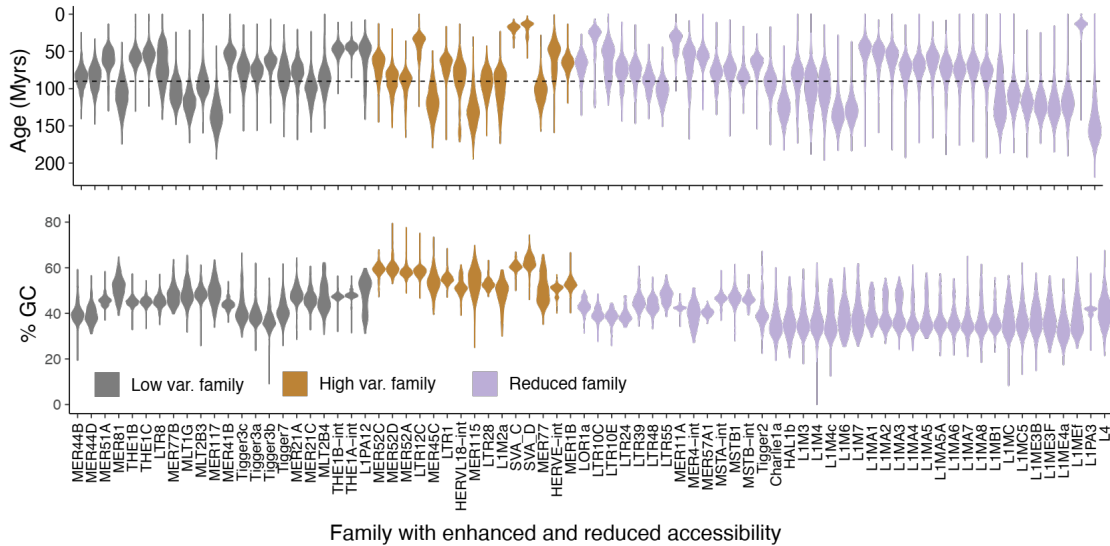
A



B



C



D

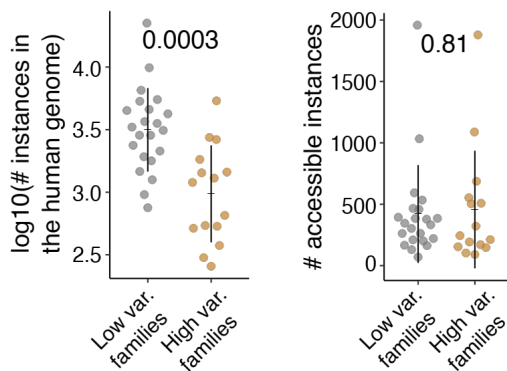


Figure S4. Characteristics of TE families display high variability in chromatin accessibility post-infection, related to Figure 3

(A) Proportions of accessible instances per enhanced family are variable between three individual groups post infection. Commonly accessible instances represent instances that are accessible in more than 25% samples from at least one group (left); rarely accessible instances represent instances that are accessible in less than 25% samples from any groups (right). Enrichment in one individual group refers to instances that are accessible in more than 25% samples for commonly accessible instances and one or more samples for rarely accessible instances. High variable families are highlighted in blue color. (B) Proportions of flu-specific, shared, and NI-specific instances of each family with enhanced and reduced accessibility. Flu-specific instances represent instances that are accessible in ≥ 1 infected and no non-infected sample; NI-specific instances represent instances that are accessible in ≥ 1 non-infected and no infected sample; Shared instances represent instances that are accessible in ≥ 1 non-infected and ≥ 1 infected samples. High variable families are in blue color. (C) Violin plot of the estimated TE evolutionary ages (up) and GC contents (bottom). Dotted lines indicate the evolution time when primates diverged from other mammals (~90 million years ago). No distinct patterns are observed between high variable, low variable families, and reduced families. The estimation of ages was described in Methods. High variable families show a higher GC content compared to others. Group 3 individuals have comparable or lower GC content than other individuals, supporting that the higher accessibility in high variable families for group 3 individuals are not derived from sequencing artifacts. (D) Number of instances and accessible instances from high variable and low variable families. *P* values computed by two-tailed student's *t*-test are shown above the dot plots.

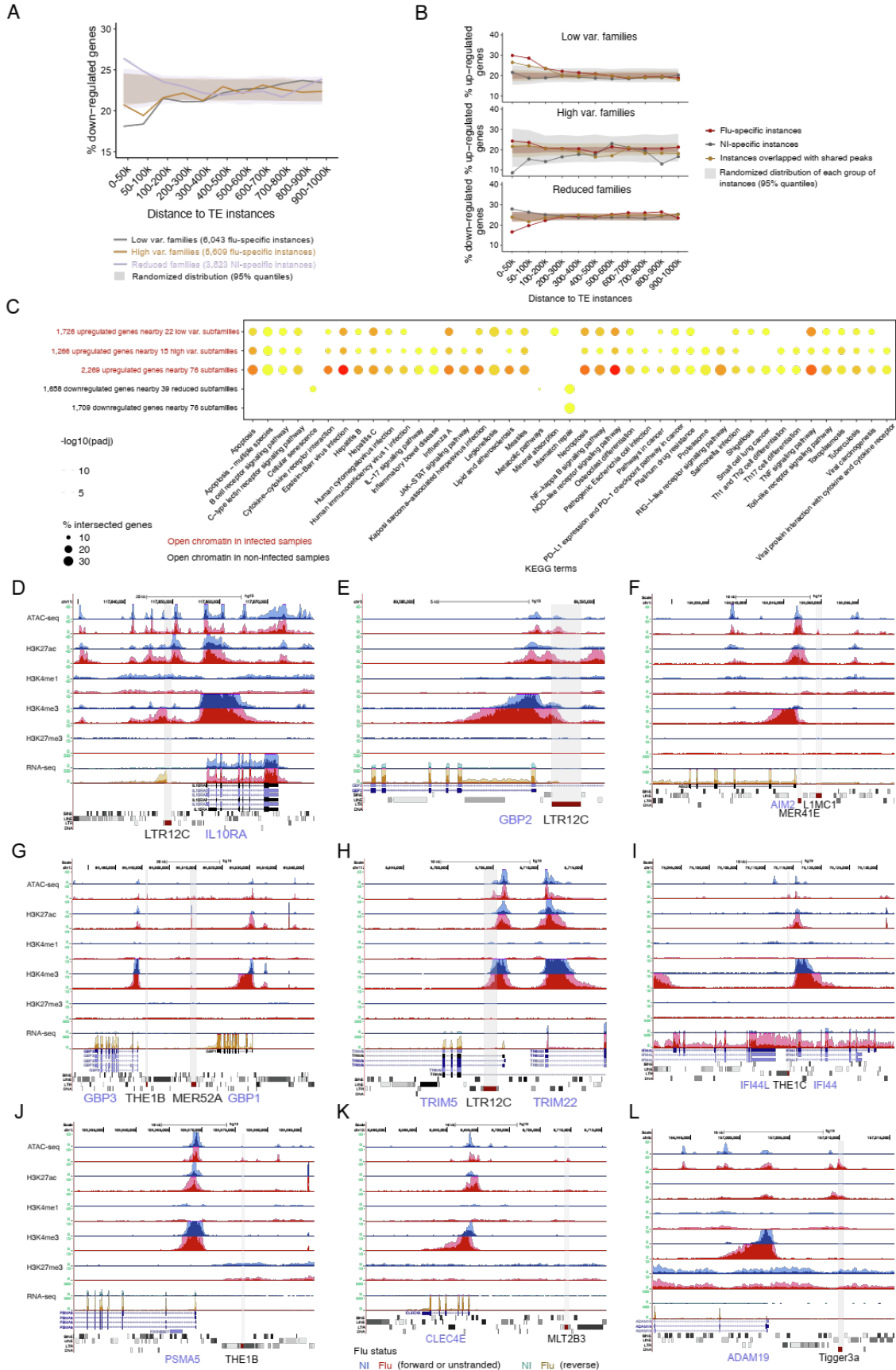


Figure S5. TE families with accessibility changes may co-opt in the immune response to IAV infection, related to Figure 4

(A) Fractions of down-regulated genes near accessible TEs relative to the random distributions. Proportions of down-regulated genes are shown within each of the genomic intervals relative to nearby accessible TEs. More details were described in **Figure 4A**. (B) Proportions of up/down-regulated genes near flu-specific, NI-specific, and shared instances. Proportions were computed within each of the genomic intervals relative to nearby accessible instances. Upregulated genes were analyzed for high variable and low variable families, and downregulated genes were analyzed for reduced families. Expected distributions were computed as we described in Methods (shaded regions, 95% confidence intervals). The proportions were compared with corresponding expected distributions. Flu-specific instances from low variable and high variable families and NI-specific instances from reduced families display the highest proportions of up/down-regulated genes within 100 kb relative to nearby accessible instances, particularly within 50 kb. (C) Pathway enrichment analysis of DEGs adjacent (≤ 50 kb) to each category of families. It shows the enrichment of multiple immune-related pathways near families with enhanced accessibility. Some pathways are differentially enriched between high variable and low variable families, including RIG-I-like receptor signaling pathway. (D) Genomic view of an accessible LTR12C with the expression was upregulated and initiated at the open chromatin region post-infection. The LTR12C instance highlighted as the shaded area shows an upregulated accessibility, expression, and H3K4me3 activity. *IL10RA* gene located near the LTR12C instance is also significantly upregulated post-infection. (E-L) Example genomic views of instances with enhanced accessibility post-infection. Instances are highlighted as the shaded areas. Eight TE immune-related gene pairs are shown, i.e., LTR12C-*GBP2*, MER41-*AIM2*, MER52A/*THE1B*-

GBP1/3, *LTR12C-TRIM22*, *THE1C-IFI44*, *THE1B-PSMA5*, and *MLT2B3-CLEC4E*, and *Tigger3a-ADAM19*. *GBP2* has been validated to be regulated by the upstream *LTR12C* instance.¹ *AIM2* has also been validated to be regulated by a *MER41* instance;² interestingly, it may also be regulated by another TE instance. Other three TE instances reported by Chuong et al.² that potentially regulate *APOL1*, *IFI6*, and *SECTMI* did not show chromatin change in macrophages (<https://computationalgenomics.ca/tools/epivar>). The dark shaded area denotes the distribution of the average RPM values and the light shaded area denotes the standard deviation. Signals of various epigenetic marks are shown in blue color for non-infected samples and red color for infected samples. For RNA-seq, forward and reverse transcripts are shown in blue and green color separately for non-infected samples; while forward and reverse transcripts are shown in red and brown color separately for infected samples.

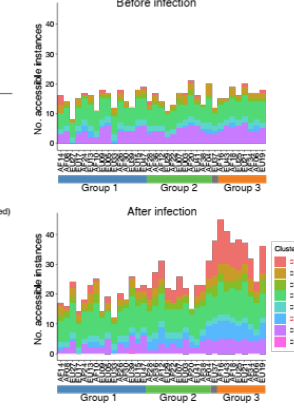
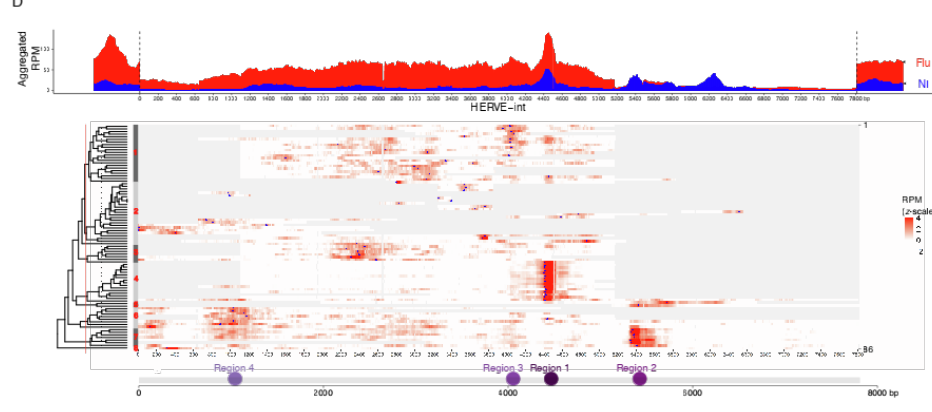
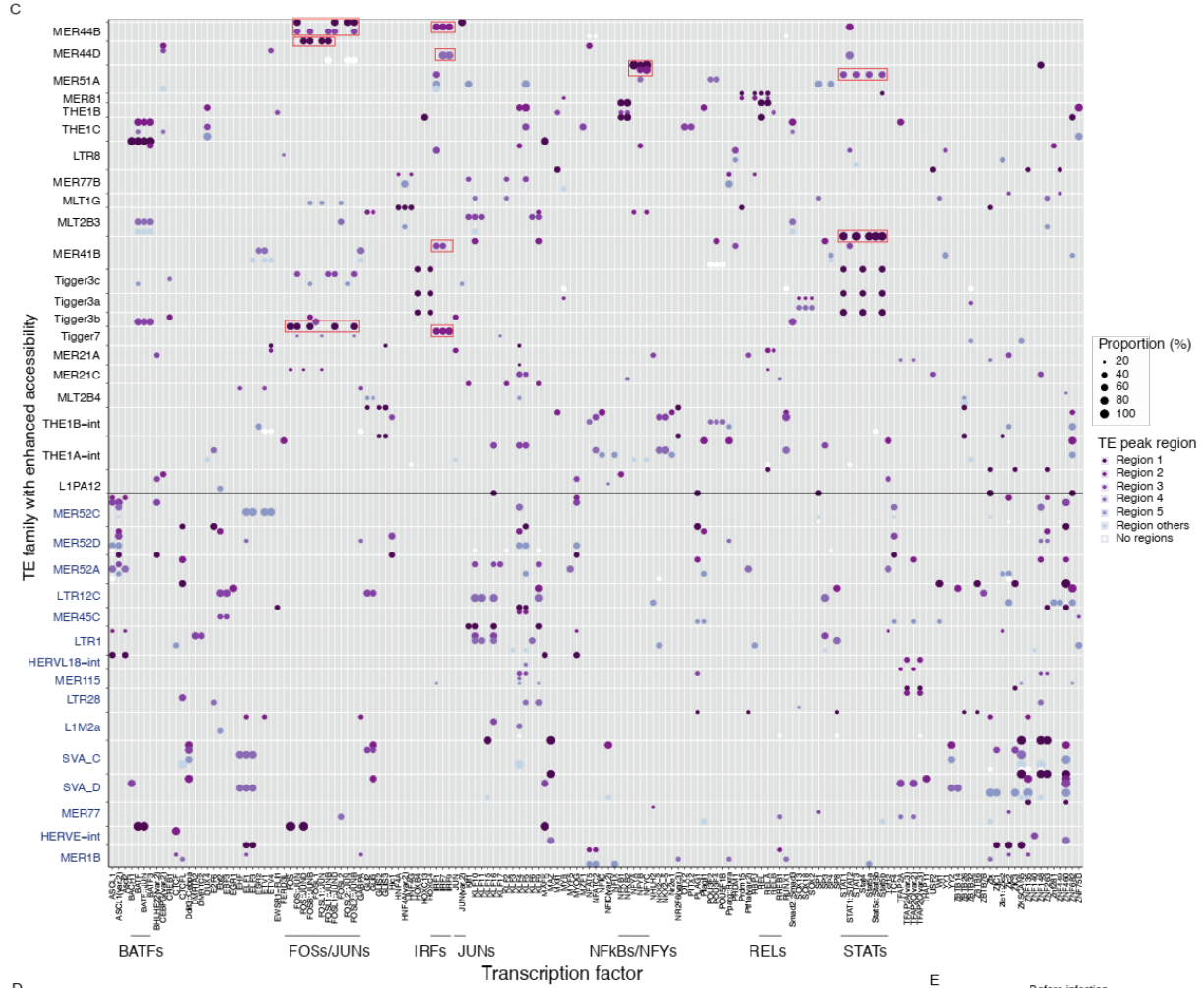
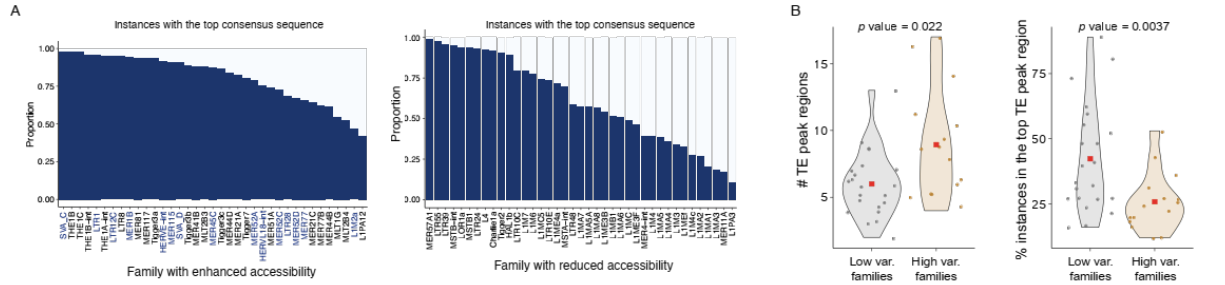
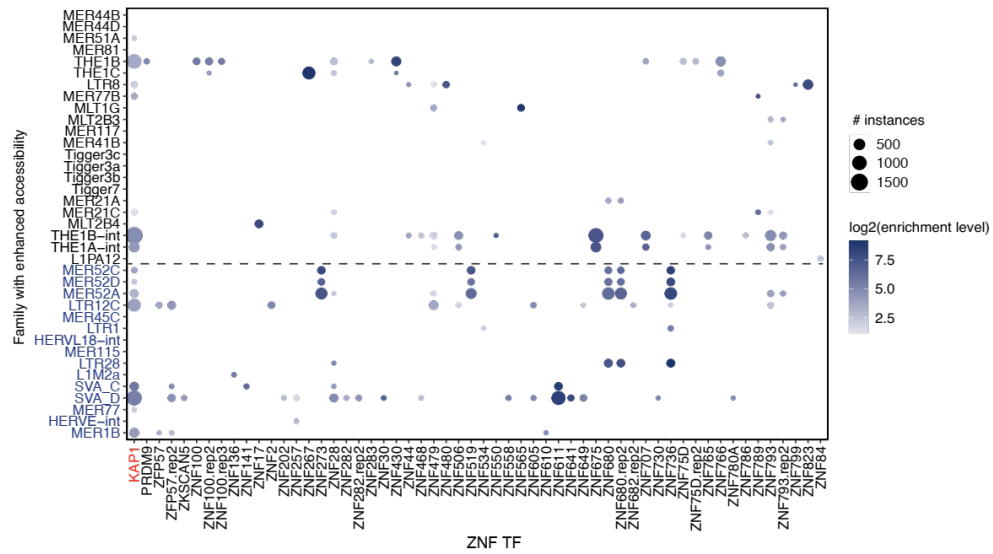


Figure S6. TE peak regions and TF binding motifs reveal the association of high variability in chromatin accessibility with KRAB-ZNFs, related to Figure 5

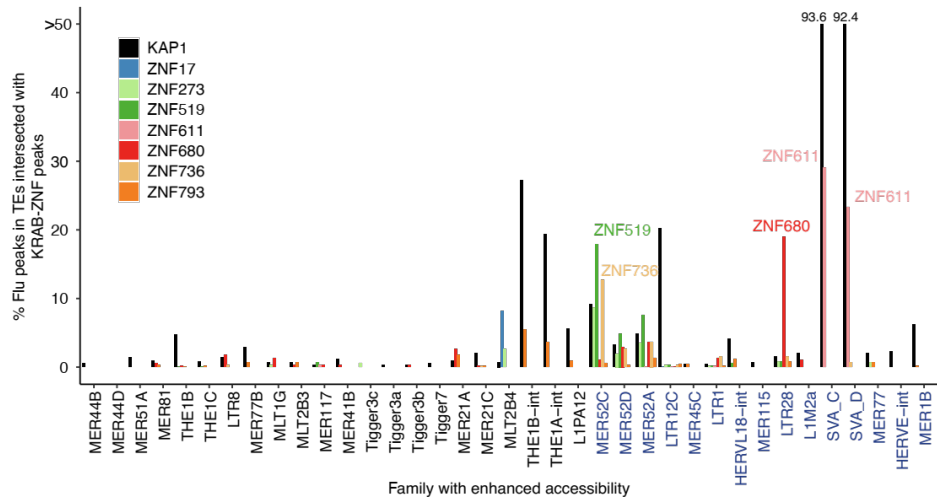
(A) Proportions of accessible instances with top candidate consensus sequences. Consensus sequence information was achieved from the “.align” files generated by RepeatMasker. It shows a high proportion for TE families with enhanced accessibility. High variable families are highlighted in blue color. (B) Number of TE peak regions detected for high variable and low variable families. Compared to low variable families, high variable families have significantly more TE peak regions and more instances in the top TE peak region. *P* values were computed by the two-tailed student’s *t*-test. (C) TF binding motifs enriched in each TE peak region of families with enhanced accessibility. Top five TE peak regions are shown. Instances in other regions and instances not within TE peak regions were analyzed separately. TE peak regions of each family are shown as separate rows. Dot size refers to the proportion of instances in each region containing each motif. Dotted line separates high variable and low variable families and high variable families are also highlighted in blue color. (D) Aggregated RPM (up) and RPM (bottom) values on each instance along the HERVE-int consensus sequence. Infected (red) and non-infected (blue) samples are shown separately including the upstream and downstream regions ($\pm 20\%$ of the consensus sequence length). Computed RPM values were *z*-scaled in the heatmap while values below zero are in white color. Deletions relative to the consensus sequence are shown in grey color. Unsupervised clustering analysis was performed with the scaled RPM values to determine the main clusters. Blue triangles indicate the peak centroids referring to the highest RPM values. TE peak regions and positions are shown in the bottom (Same as **Figure 5B**). More details were described in Methods. (E) Number of accessible instances from each cluster. Infected and non-infected samples are shown separately. Individuals are ordered based

on the clustering obtained in **Figure 3C**. Instances from Cluster 1 and 6, most of which contain TE peak region 3 and 4, are more abundant in group 3 individuals than group 1 individuals.

A



B



C

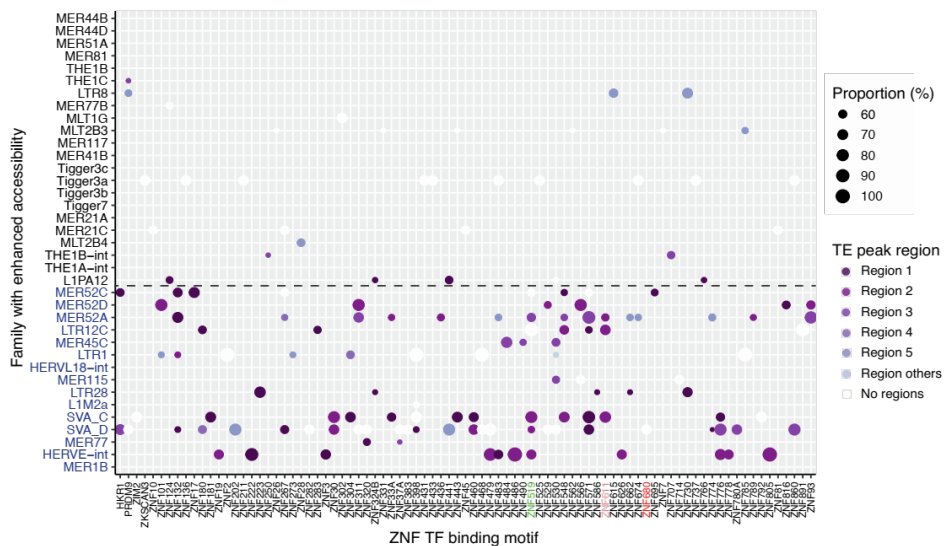
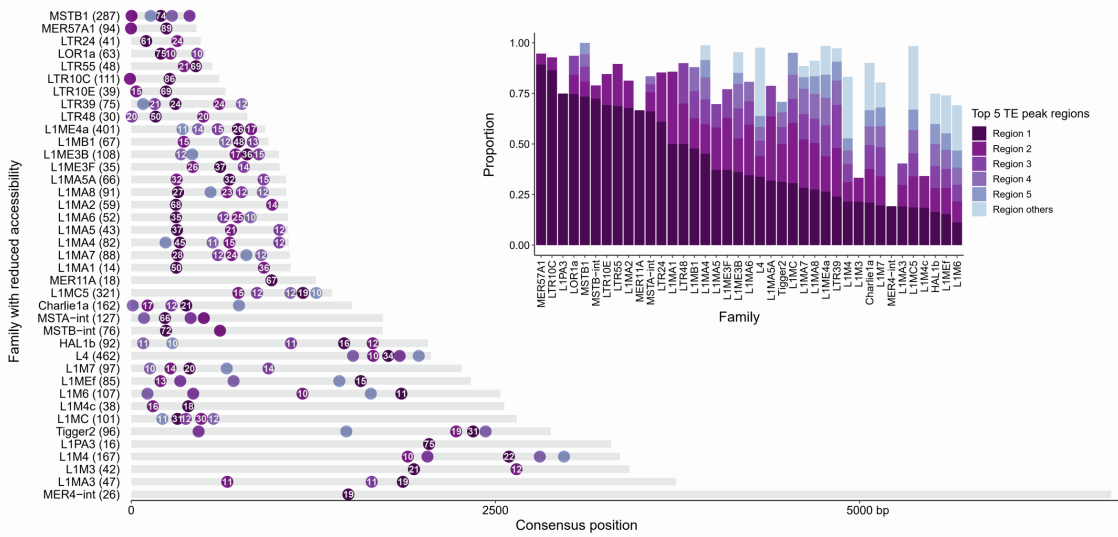


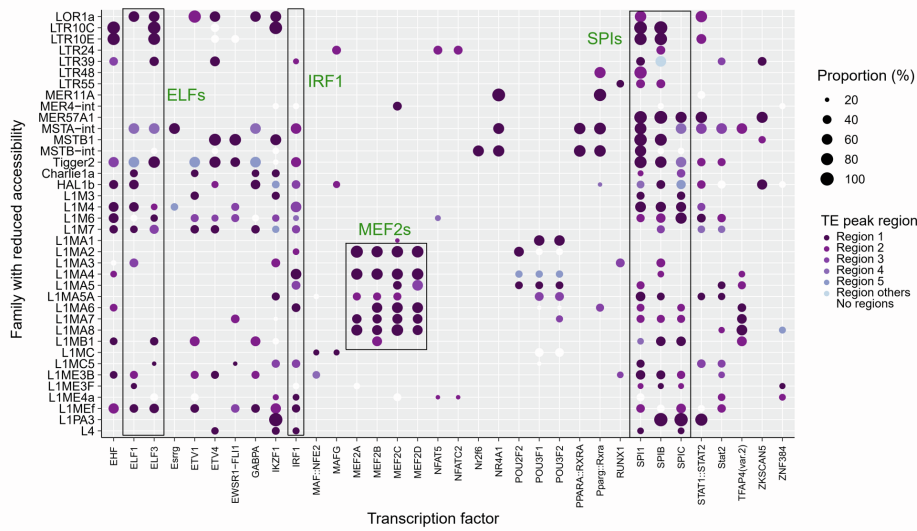
Figure S7. KAP1 and KRAB-ZNFs are associated with high variability in chromatin accessibility in high variable families, related to Figure 5

(A) Enrichment levels of KRAB-ZNF binding sites in high variable and low variable families in 257 HEK293T cell lines.³ High variable families like MER52s, SVAs, LTR12C, and LTR28 are shown to be enriched for KRAB-ZNF binding sites. Color intensity refers to the fold enrichment relative to the random distribution (see Methods). (B) Proportion of KAP1 and KRAB-ZNF binding sites that overlap with accessible regions in TEs post-infection. A 100-bp of genomic region centered at the ATAC-seq peak centroids was used for this analysis. KRAB-ZNFs with a minimum of 5% across enhanced families were visualized. (C) KRAB-ZNF binding motifs enriched in enhanced families. Motifs were obtained from Barazandeh et al.⁴ Same motifs enriched across TE peak regions are aggregated. TE peak regions with the most instances are shown as representatives. KRAB-ZNFs with their enrichment of binding sites in high variable families are highlighted.

A



B



C

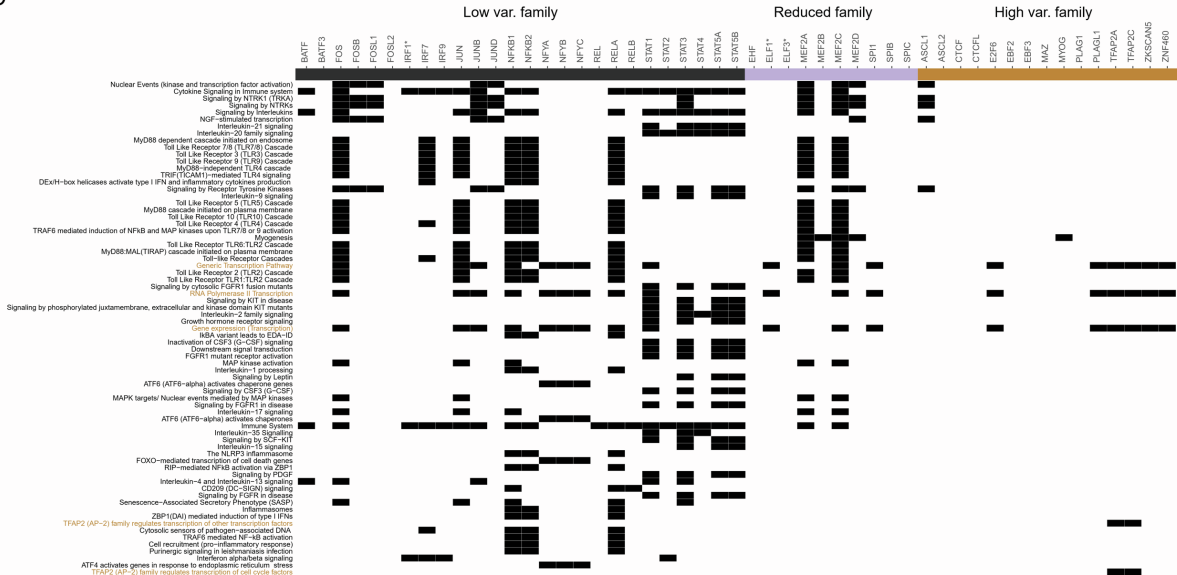


Figure S8. TE Peak regions and TF binding motifs enriched in reduced families and different sets of TE families may co-opt in distinct regulatory pathways in response to IAV infection, related to Figure 5

(A) Distribution and TE peak regions on reduced families. The inset barplot shows the proportion of instances in each TE peak region. The locations and proportions (%) of the top-five TE peaks regions along the consensus sequence are shown. The number in each dot refers to the proportion (above 10%) among accessible instances in each TE peak region. It shows that most instances with reduced accessibility (in Region 1) from L1MA families are consistently located at similar locations of the 3' end of consensus sequence. Y-axis shows the family name and the number of accessible instances mapped to the top consensus sequence. TE peak regions were detected as we described in Methods. (B) TF binding motifs enriched in reduced families. Same motifs enriched across TE peak regions were aggregated and the proportions of instances containing each motif are shown. TE peak regions with more accessible instances are shown as representatives. Black boxes highlighted the SPI and MEF2 related motifs (green color) enriched in reduced families. It shows that Region 1s (located at around 300 bp, **Figure S6B**) of L1MA families are consistently enriched for MEF2 related motifs. MEF2 related TFs were previously reported to regulate anti-microbial genes.⁵ (C) TFs potentially bound to different categories of families are involved in distinct pathways. Apart from AP-2 related pathways, TFs bound to high variable families are mainly involved in transcription-related pathways. TFs bound to low variable families and reduced families are mainly involved in cytokine signaling and other immune-related pathways. Bars in different colors represent different categories of families they are enriched. * indicates motifs that are enriched in different categories of families. IRF1 motif is

enriched in both low variable families and reduced families. ELF1/3 motifs are enriched in both high variable families and reduced families.

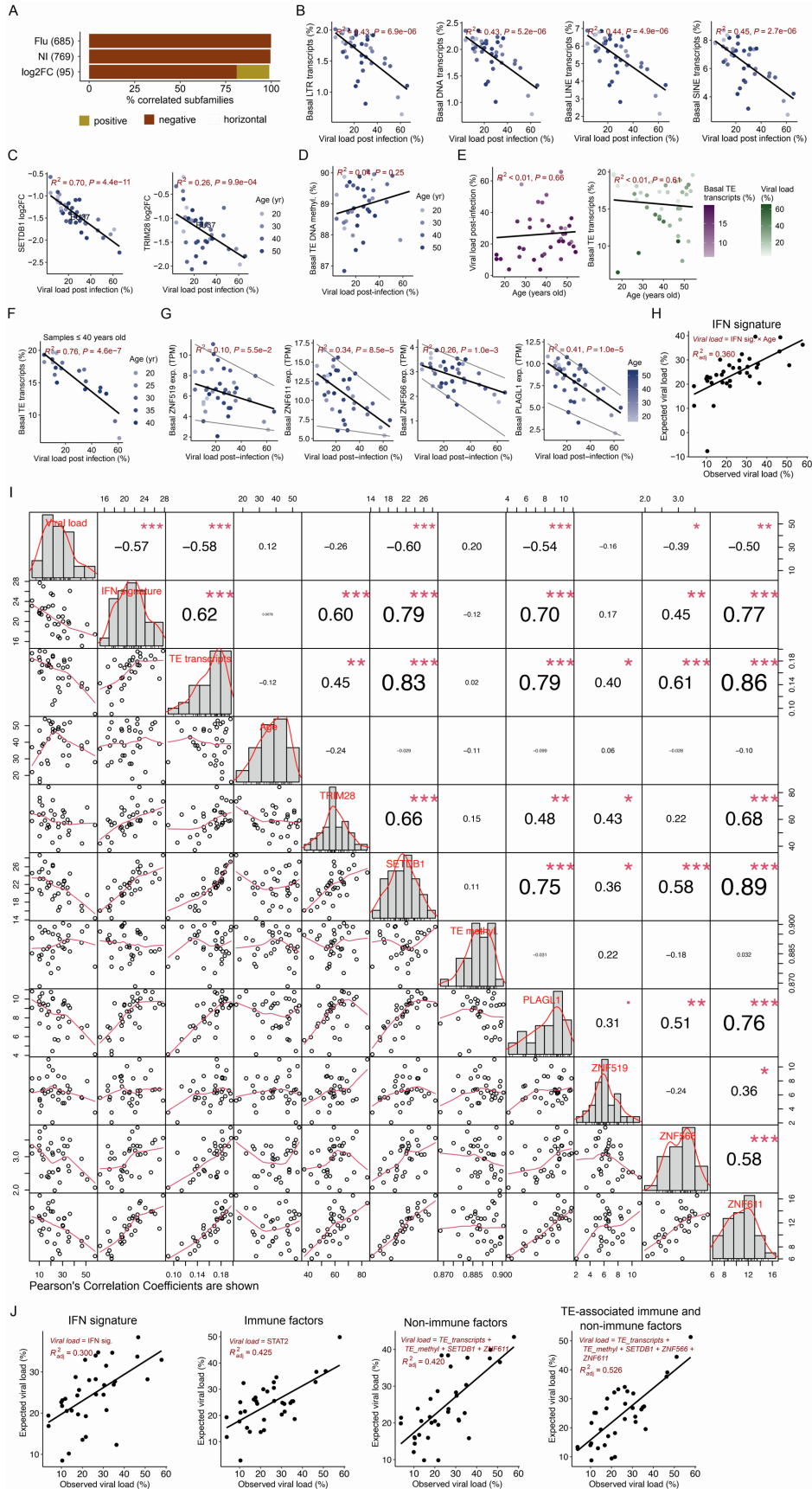


Figure S9. Correlations between TE-associated host factors and viral load post-infection, related to Figure 6

(A) Correlation directions among TE families that are correlated ($R^2 \geq 0.3$ and p value ≤ 0.05) with viral load post-infection (**Figure 6A**). More details were described in Methods. (B) Consistent inverse correlations between the basal transcripts of four main TE subclasses and viral load, including LTR, DNA, LINE, and SINE. The basal amount of transcript refers to the proportion of aggregated normalized read counts in each subclass among the global transcripts. Black line represents the regression line. R^2 and p values computed by the linear regression model are shown. (C) Correlations between the expression fold changes (log₂FCs) of both *SETDB1* (left) and *TRIM28* (right) and viral load. *SETDB1* rather than *TRIM28* is shown to be inversely correlated with viral load. (D) Correlation between the basal DNA methylation level in TEs and viral load. Individuals with low viral load are likely to have lower TE DNA methylation levels. (E) Correlations between both viral load post-infection (left) and basal TE transcripts (right) and ages. Samples from younger individuals have a relatively higher amount of basal TE transcripts with a smaller deviation compared to older samples in macrophages. Two samples were identified as outliers having among the lowest amount of basal TE transcripts; strikingly, they also preserved among the highest viral load. (F) Correlation between the basal TE transcripts and viral load for individuals 40 years old or younger. It shows an increased correlation compared to the correlation among 39 samples (**Figure 6B**). (G) Correlations between the basal expression of candidate host factors and viral load post-infection. Basal *ZNF519*, *ZNF566*, *ZNF611* and *PLAGL1* expressions are both inversely correlated with viral load post-infection. (H) Multivariable regression model developed for the predictive of viral load post-infection using type I interferon (IFN) signature and age only. Details were described in

Methods. **(I)** Correlation matrix chart of viral load and variables that are potentially associated with IAV infection and TEs. Histograms and kernel density overlays of each variable are shown. Scatterplot matrix and absolute correlations between each of two variables and viral load are also shown. Red lines represent the distribution. R *chart.Correlation* function was used for the analysis. Pearson's correlation coefficients are shown (* $p \leq 0.05$, ** $p \leq 0.01$, *** $p \leq 0.001$).

(J) Multivariable regression models developed for the predictive of viral load post-infection while age was included as an independent variable. We used the same sets of variables for the models included in **Figure S9H** and **Figure 6H-J**. The adjusted R^2 is significantly lower than the model developed by the inclusion of age as an interaction term variable. Details were described in Methods.

REFERENCES

1. Srinivasachar Badarinarayan, S., Shcherbakova, I., Langer, S., Koepke, L., Preising, A., Hotter, D., Kirchhoff, F., Sparrer, K.M.J., Schotta, G., and Sauter, D. (2020). HIV-1 infection activates endogenous retroviral promoters regulating antiviral gene expression. *Nucleic Acids Research* 48, 10890–10908. 10.1093/nar/gkaa832.
2. Chuong, E.B., Elde, N.C., and Feschotte, C. (2016). Regulatory evolution of innate immunity through co-option of endogenous retroviruses. *Science* 351, 1083–1087. 10.1126/science.aad5497.
3. Imbeault, M., Helleboid, P.-Y., and Trono, D. (2017). KRAB zinc-finger proteins contribute to the evolution of gene regulatory networks. *Nature* 543, 550–554. 10.1038/nature21683.
4. Barazandeh, M., Lambert, S.A., Albu, M., and Hughes, T.R. (2018). Comparison of ChIP-Seq Data and a Reference Motif Set for Human KRAB C2H2 Zinc Finger Proteins. *G3 Genes|Genomes|Genetics* 8, 219–229. 10.1534/g3.117.300296.
5. Clark, R.I., Tan, S.W.S., Péan, C.B., Roostalu, U., Vivancos, V., Bronda, K., Pilátová, M., Fu, J., Walker, D.W., Berdeaux, R., et al. (2013). MEF2 Is an In Vivo Immune-Metabolic Switch. *Cell* 155, 435–447. 10.1016/j.cell.2013.09.007.