

## Supplementary Material for LARVA

### 1. Pseudogene UTR, TSS, and promoter sites removal

Pseudogenes are known to be hotspots of artifacts in numerous genomics analyses. It is partially because read mapping in pseudogenes might be complicated due to their context similarity with their parent genes. In order to analyze the mutation events in the pseudogene regions, we extracted all the pseudogenes from the Gencode annotation (version 19) and calculated the average mutation counts from the pooled samples in gene and pseudogene regions, and also the up- and downstream 2kb region of all pseudogenes. Possibly due to the shorter length of pseudogenes, a larger variance of the mutation rate was observed in pseudogenes, although two-sided Wilcoxon test shows no significant difference ( $P = 0.453$ ). However, we observed a noticeable elevated mutation rate in the up- and downstream regions of pseudogenes (Fig. S2). In order to exclude artifacts, such as variant calling difficulties, we excluded the pseudogenes from the Gencode gene list in our analysis.

### 2. Details of model fittings

#### The constant mutation rate assumption and the resultant binomial distribution

The underlying assumption of the binomial model is that the mutation rate within the given region is constant. Suppose the target region has  $n$  bases in length, and the homogeneous mutation rate is  $p$ . Then mutation count  $x$  inside this region falls into a binomial distribution with the probability mass function as

$$\Pr(X = x) = \binom{n}{x} p^x (1-p)^{n-x} \quad (1.1).$$

Given the mutation count data, the maximum likelihood estimator of the mutation rate is just

$$\hat{p} = \frac{\sum_{i=1}^k x_i}{\sum_{i=1}^k n_i} \quad (1.2)$$

where  $k$  represents the total number of regions and  $i$  is the region index.

#### The beta-binomial distribution used in LARVA

Instead of the fixed mutation rate assumption, we provided more flexibility of the mutation rate by allowing it to follow a beta distribution

Unknown  
Field Code Changed

Unknown  
Field Code Changed

Unknown  
Field Code Changed

Jing Zhang 7/10/2015 4:31 PM  
Formatted: Space After: 6 pt

Jing Zhang 7/10/2015 4:32 PM  
Formatted: Space After: 6 pt

Unknown  
Field Code Changed

Jing Zhang 7/10/2015 4:32 PM  
Formatted: Space After: 6 pt

Unknown  
Field Code Changed

Jing Zhang 7/10/2015 4:32 PM  
Formatted: Space After: 6 pt

Jing Zhang 7/10/2015 4:32 PM  
Formatted: Space After: 6 pt

$$\begin{aligned}\pi(p|\alpha,\beta) &= \text{Beta}(\alpha,\beta) = \frac{p^{\alpha-1}(1-p)^{\beta-1}}{\text{Beta}(\alpha,\beta)} \\ &= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1}(1-p)^{\beta-1}\end{aligned}\quad (1.3).$$

Suppose the mutation count is  $x_i, i=1,2,\dots,k$ , and the sample size and binomial probability can be expressed as  $n_i$  and  $p_i$ . Instead of assuming the mutation in all the bins is a constant, we can set up a two-stage model

$$\begin{aligned}x_i|p_i &\sim \text{Binomial}(n_i, p_i) \\ p_i &\sim \text{Beta}(\alpha, \beta)\end{aligned}\quad (1.4).$$

Then the total number of mutations within the bin with length  $n$  follows the beta binomial distribution as in (1.5)

$$\Pr\{X = x_i\} = \binom{n_i}{x} \frac{\Gamma(\alpha+\beta)\Gamma(\alpha+x_i)\Gamma(\alpha+n_i-x_i)}{\Gamma(\alpha)\Gamma(\beta)\Gamma(\alpha+\beta+n_i)}\quad (1.5)$$

To estimate the parameters in beta-binomial distribution we used the scheme described in (1,2). When the target bin length is fixed, resulting in  $n_i = n, i=1,2,\dots,k$ , the mean and variance of mutation counts can be written as

$$\begin{aligned}E[X] &= n \frac{\alpha}{\alpha+\beta} = n\mu \\ \text{var}[X] &= n\mu(1-\mu)\sigma, \\ \sigma &= \frac{1}{\alpha+\beta+1}\end{aligned}\quad (1.6)$$

For simplicity, we directly estimate  $\mu$  and  $\sigma$  instead of  $\alpha$  and  $\beta$ . Hence the moment estimator can be immediately derived from equation (1.6).

When the target region length is variable, estimation is a little bit more complicated. Define additional parameters for mathematical convenience as in (1.7).

$$\begin{aligned}\hat{p} &= \frac{\sum_{i=1}^k w_i \hat{p}_i}{w} \\ w_i &= \frac{n_i}{1+\sigma(n_i-1)} \\ w &= \sum_{i=1}^k w_i \\ S &= \sum_{i=1}^k w_i (\hat{p}_i - \hat{p})^2\end{aligned}\quad (1.7)$$

We can derive the moment estimator in (1.8)

$$\mu = \hat{p} = 1 - \hat{q}$$

$$\sigma = \frac{S - \hat{p}\hat{q} \left[ \sum_{i=1}^k \frac{w_i}{n_i} \left( 1 - \frac{w_i}{w} \right) \right]}{\hat{p}\hat{q} \left[ \sum_{i=1}^k w_i \left( 1 - \frac{w_i}{w} \right) - \sum_{i=1}^k \frac{w_i}{n_i} \left( 1 - \frac{w_i}{w} \right) \right]} \quad (1.8)$$

However, from (1.8),  $w_i$  is also a function of  $\sigma$  which is to be estimated, and there is no analytical solution to it. Hence as suggested in (2), we initially assigned the  $w_i$  proportional to  $n_i$  to get a rough estimate of  $\gamma$ . Then  $w_i$  was updated with this estimate to obtain a more accurate estimation of  $\sigma$ .

### 3. Coding Region Mutation Burden Analysis

LARVA is not designed for the coding region analysis due to the availability of synonymous sites, which serve as a natural and biologically meaningful background in these regions. For the sake of gaining additional insight from an exome performance calibration, we nevertheless evaluated LARVA's ability to identify statistically significant mutation burdens in genes. Exome variant data was obtained from The Cancer Genome Atlas (TCGA) Data Portal (3). The complete set of exome variant calls includes 20 cancer types and 5032 samples in total. A detailed graph of the collected data is provided in Fig S8.

Gene annotation data was derived from the GENCODE v19 annotation files (4). All complete protein-coding transcripts were extracted, and all the transcripts for each gene were merged, as demonstrated in Fig S9. This data spanned 19,822 genes, and a total of 252,356,877 nucleotides. We plotted the distribution of gene lengths in Fig S10. The total number of mutations falling into the merged gene regions is 3,547,350, and the average mutation rate is 0.0141 for the pooled samples. As with the noncoding regions, huge mutation heterogeneity was observed in the coding regions (Fig S11).

We removed the genes with length less than the bottom 5% of gene lengths for higher annotation confidence, and then compared the performance of LARVA and the binomial test. After p-value adjustment, LARVA found 7 genes that are potentially under higher mutation burden (Table S3). For each of these genes, we searched for literature supporting cancer association. Except for KRTAP4-11, we found all the remaining genes are clearly documented with some cancer association. Note that we reported only one Pubmed ID per gene, even if there are many more supporting references. Our findings effectively demonstrate that LARVA is capable of finding meaningful results on protein coding regions. On the other hand, the p-values for the binomial test method were heavily inflated. After p-value adjustment, there are 6759 out of 18,826 genes, roughly 35.90%, with p-value less than 0.05. It is very unlikely that all such genes are associated with cancer.

P-values given by LARVA and binomial test are given in Fig. S12. It is shown that the p-value distribution from the binomial test severely violates the uniform distribution assumption, which is consistent with its bad fitting of the data. On the

Unknown  
Field Code Changed

Jing Zhang 7/10/2015 4:32 PM  
Formatted: Space After: 6 pt

Jing Zhang 7/10/2015 4:31 PM  
Deleted: -

Jing Zhang 7/10/2015 4:31 PM  
Formatted: Heading 2, Space After: 0 pt, Line spacing: single

Jing Zhang 7/10/2015 4:31 PM  
Deleted: -

Lucas Lochovsky 7/10/2015 10:01 PM  
Deleted: because

Jing Zhang 7/10/2015 4:31 PM  
Formatted: Font:+Theme Body

Jing Zhang 7/10/2015 4:31 PM  
Formatted: Font:+Theme Body

Jing Zhang 7/10/2015 4:33 PM  
Formatted: Justified, Indent: Left: 0 cm, First line: 0.63 cm, Space After: 6 pt, Line spacing: at least 15 pt

Lucas Lochovsky 7/10/2015 10:02 PM  
Deleted: the

Jing Zhang 7/10/2015 4:31 PM  
Formatted: Font:+Theme Body

Lucas Lochovsky 7/10/2015 10:05 PM  
Deleted: However, more information is immediately available for the

Jing Zhang 7/10/2015 4:31 PM  
Formatted: Font:+Theme Body

Jing Zhang 7/10/2015 4:25 PM  
Deleted: We

Lucas Lochovsky 7/10/2015 10:05 PM  
Deleted: still

Jing Zhang 7/10/2015 4:31 PM  
Formatted: Font:+Theme Body

Jing Zhang 7/10/2015 4:28 PM  
Deleted: -

Jing Zhang 7/10/2015 4:33 PM  
Formatted: Normal, Left, Space After: 6 pt, Line spacing: single

other hand, the p-values from the LARVA method (Fig. S12, left hand side) roughly follow a uniform distribution. It is worth mentioning that after replication timing correction, the p-values from LARVA method have improved concordance with the theoretical distribution, indicating the importance of correction.

#### 4. Importance of covariate correction

It is well known that local mutation rate is affected by various factors, such as replication timing and GC content (3). Due to these confounding factors, the observed mutation count data distribution is actually a mixture of several different distributions, which further increases the overdispersion. We used some simulations to show this effect in Fig S14. We randomly simulated five binomial distribution of binomial( $n=100,000, p_i$ ), where  $p_i$  was uniformly sampled from  $[1e^{-6}, 5e^{-5}]$  to mimic the mutation rate difference from various replication timing regions. The empirical distribution of the pooled data was given in the pink line in Fig S14. It is shown that even if we ignore patient-specific heterogeneity, the observed data demonstrates much larger variation than expected simply by mixing several different binomial distributions. It is necessary to remove such effects to obtain a more reasonable P-value calculation.

#### 5. Factors that affect overdispersion in the mutation count data

##### 5.1 Heterogeneity in mutation rates in different patients/cancer types

Suppose that mutation rate per sample is constant across the whole genome, but varies from different patients/cancer types. Each patient-specific mutation rate can be considered a random sample from beta distribution, so after pooling all samples together, the mutation counts for each bin follows a beta-binomial distribution. The overdispersion under this condition depends on how different these patients are.

##### 5.2 Length of the target region to be analyzed

Assume that  $y$  is the number of somatic variants in  $n$  bases, the point mutation rate is  $\epsilon$ . Unlike the constant mutation assumption in the binomial distribution, we assume that  $\epsilon$  is a random variable with

$$E(\epsilon) = p$$

$$Var(\epsilon) = \phi p(1-p)$$

So the variance of  $y$  can be calculated as

Jing Zhang 7/10/2015 4:33 PM  
Formatted: Heading 2

Jing Zhang 7/10/2015 4:33 PM  
Formatted: Justified, Indent: First line: 0.63 cm, Space After: 6 pt, Line spacing: at least 15 pt

Lucas Lochovsky 7/10/2015 10:15 PM  
Deleted: {\cite: mutsig}

Jing Zhang 7/10/2015 4:33 PM  
Formatted

Lucas Lochovsky 7/10/2015 10:26 PM  
Deleted: Without distinguishing...u (... [1])

Jing Zhang 7/10/2015 4:33 PM  
Formatted

Lucas Lochovsky 7/10/2015 10:26 PM  
Deleted: and

Jing Zhang 7/10/2015 4:33 PM  
Formatted

Lucas Lochovsky 7/10/2015 10:27 PM  
Deleted: were

Jing Zhang 7/10/2015 4:33 PM  
Formatted (... [2])

Lucas Lochovsky 7/10/2015 10:27 PM  
Deleted: the ...atient-specific (... [3])

Jing Zhang 7/7/2015 4:08 PM  
Formatted: Heading 2

Jing Zhang 7/10/2015 4:19 PM  
Formatted: Justified, Indent: Left: 1.39 cm, Space After: 10 pt, Line spacing: multiple 1.15 li

Jing Zhang 7/10/2015 4:19 PM  
Formatted: Font:Italic, Underline

Lucas Lochovsky 7/10/2015 10:31 PM  
Deleted: Let's assume...oppose that (... [4])

Lucas Lochovsky 7/10/2015 10:32 PM  
Deleted: 1

Jing Zhang 7/10/2015 4:19 PM  
Formatted: Font:Italic, Underline

Unknown  
Field Code Changed (... [5])

Lucas Lochovsky 7/10/2015 10:32 PM  
Deleted: Different from...nlike the (... [6])

Unknown  
Field Code Changed

Unknown  
Field Code Changed

Unknown  
Field Code Changed

$$\begin{aligned}
\text{Var}(y) &= E_{\epsilon}(\text{Var}(y|\epsilon)) + \text{Var}(E(y|\epsilon)) \\
&= E_{\epsilon}[n\epsilon(1-\epsilon)] + \text{Var}(n\epsilon) \\
&= n[p - \phi p(1-p) - p^2] + n^2\phi p(1-p) \\
&= np(1-p)[1 + (n-1)\phi]
\end{aligned}$$

The variance in the mutation count data is scaled by a factor of  $1 + (n-1)\phi$ .

Biologically speaking, the difference of the point mutation rate within the analyzed noncoding region varies much more in longer noncoding elements (2.5kb promoters) as compared to smaller regions (200bp TSS), resulting in very different overdispersion parameters in the estimation stage. Hence, we do not recommend evaluation of regions of non-comparable length in the same run.

Unknown  
Field Code Changed

Unknown  
Field Code Changed

Lucas Lochovsky 7/10/2015 10:43 PM

Deleted: From the b

Lucas Lochovsky 7/10/2015 10:43 PM

Deleted: side

Lucas Lochovsky 7/10/2015 10:33 PM

Deleted: -

Lucas Lochovsky 7/10/2015 10:42 PM

Deleted: quite

Lucas Lochovsky 7/10/2015 10:42 PM

Deleted: to

Lucas Lochovsky 7/10/2015 10:42 PM

Deleted: e

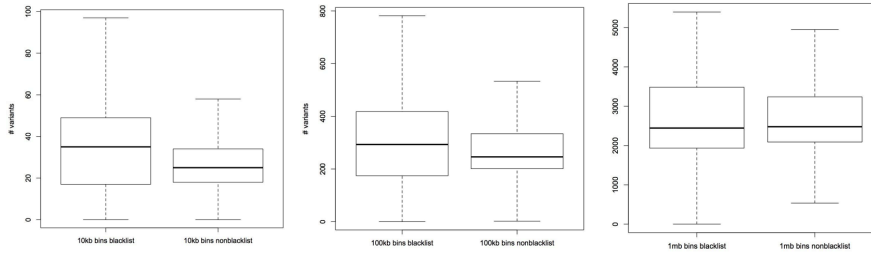
Lucas Lochovsky 7/10/2015 10:42 PM

Deleted: together

## Supplementary figures

**Figure S1**

Boxplot of mutations count in 10k, 100k, and 1mb regions with or without overlapping with the blacklist region. P-values were calculated from the two-sided Wilcoxon tests ( $P < 2.2 \times 10^{-16}$  for 10kb bin,  $P = 4.767 \times 10^{-5}$  and 0.473 for the 100k, and 1mb bins). In the smaller bin regions (10k and 100k), regions overlapped with blacklists demonstrates significantly higher mutation rate.



Jing Zhang 7/7/2015 4:44 PM

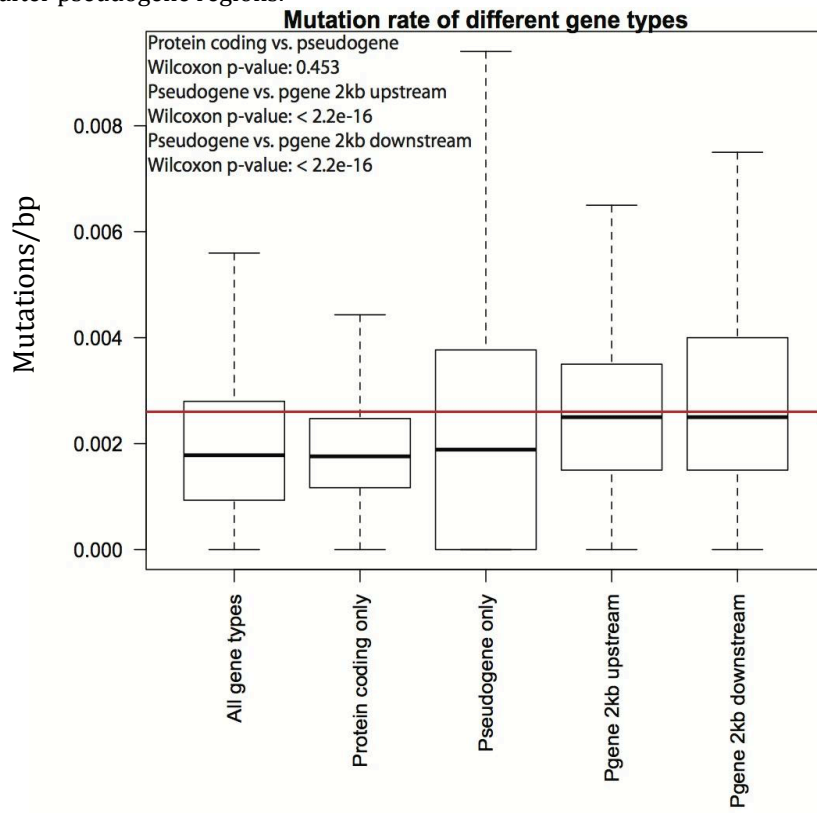
Deleted: [Page Break](#)

Jing Zhang 7/7/2015 4:10 PM

Formatted: Font:(Default) +Theme Body, 12 pt, Font color: Auto, (Asian) Chinese (PRC)

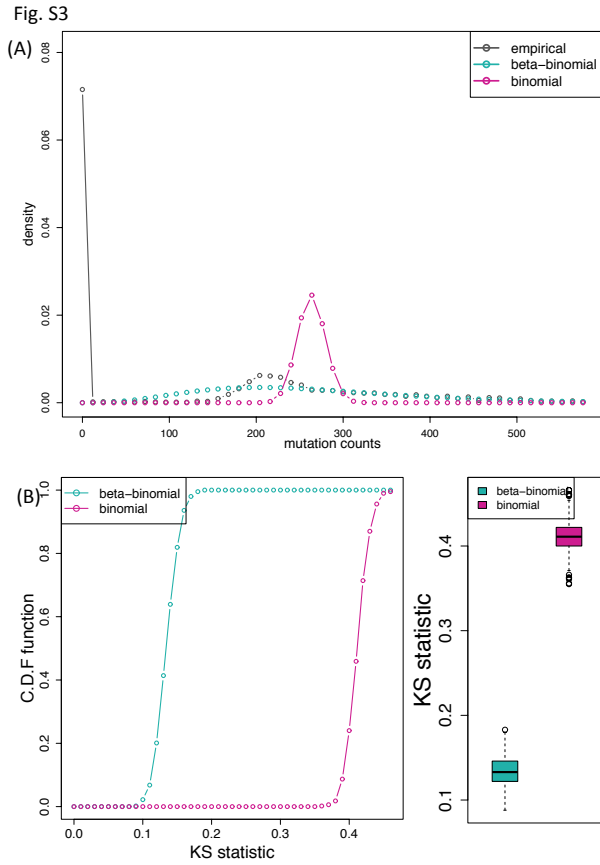
**Figure S2**

Average mutation rate estimation from gene, pseudogene, and regions before and after pseudogene regions.



**Figure S3**

Fitting comparison between beta-binomial and binomial distribution. (A) Density plot of the beta-binomial, binomial, and empirical distribution of read count data in 100kb bins; (B) C.D.F curve of the KS statistics of beta-binomial and binomial generated counts VS. random samplings in the observed counts; (C) Boxplots of the KS statistics.

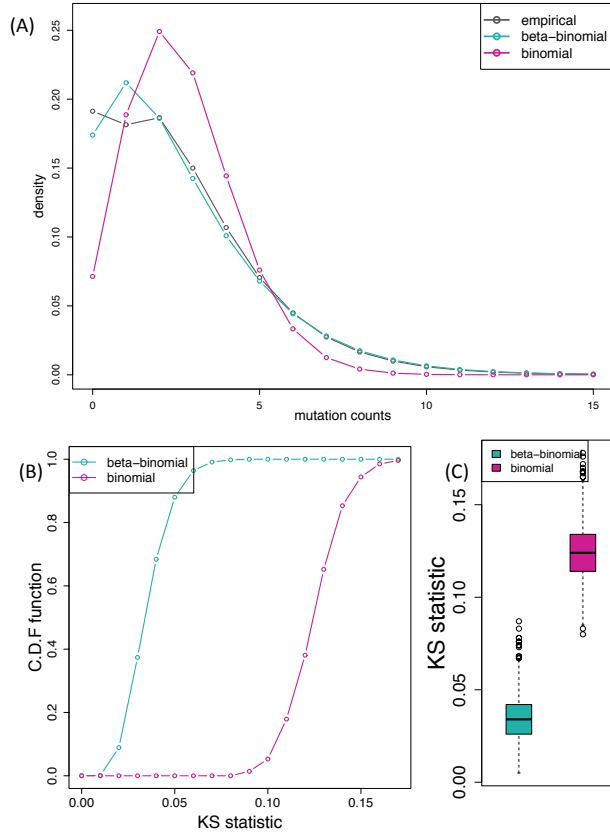




**Figure S4**

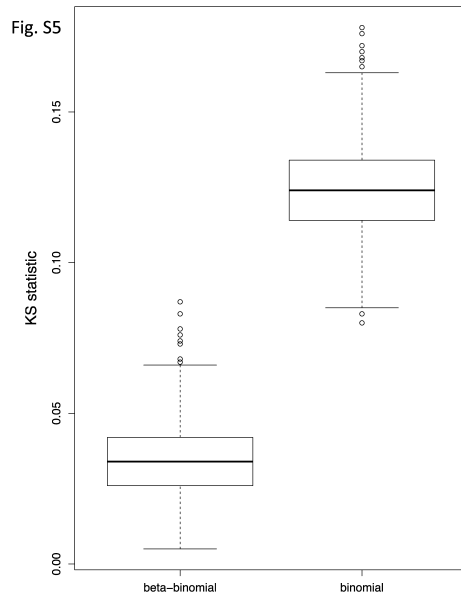
Fitting comparison between beta-binomial and binomial distribution. (A) Density plot of the beta-binomial, binomial, and empirical distribution of read count data in 1kb bins; (B) C.D.F curve of the KS statistics of beta-binomial and binomial generated counts VS. random samplings in the observed counts; (C) Boxplots of the KS statistics.

Fig. S4



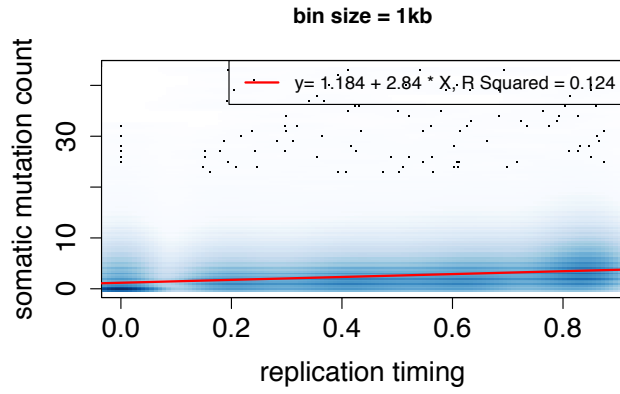
**Figure S5**

Half of the observed data is used for model fitting of both beta-binomial and binomial distribution, and the remaining half was used to calculate the KS statistics with generalizations from the fitted distributions. Boxplots of 100 repeats were given below.



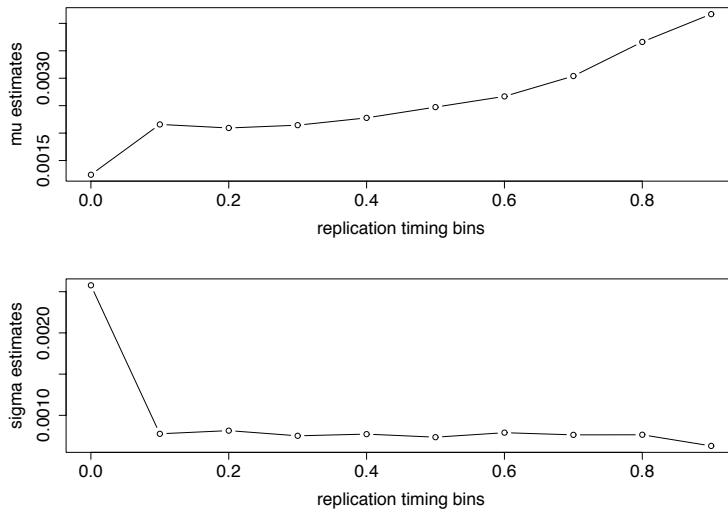
**Figure S6**

The smooth scatter plots of the mutations count in all tumor samples within 1kb bins vs. its averaged replication timing value. A linear regression was fitted and the R-squared values are up to 0.124.



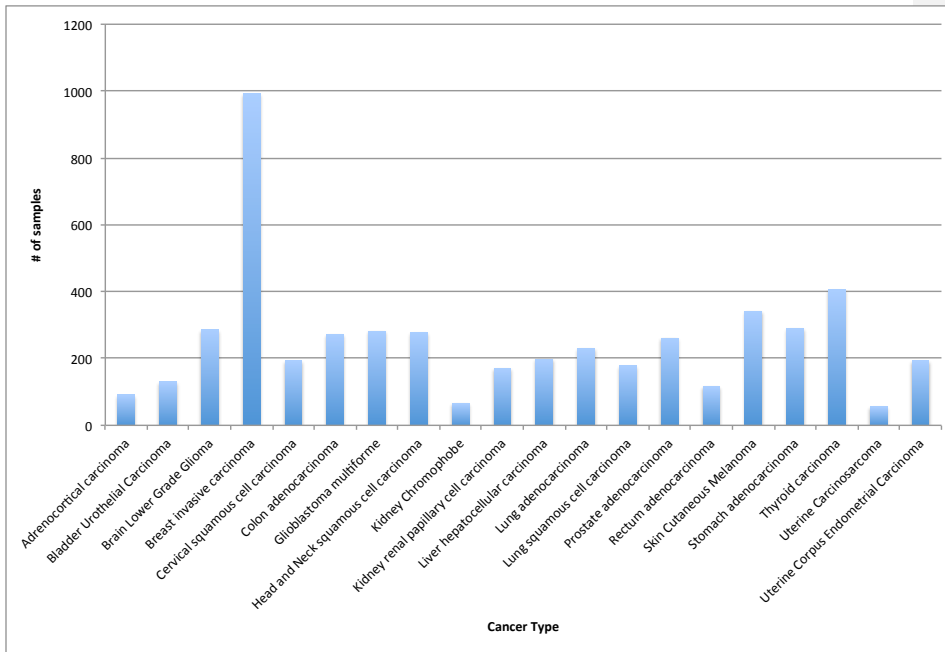
**Figure S7**

The fitted  $\mu$  and  $\sigma$  were plotted for each the 10 used replication timing bins.



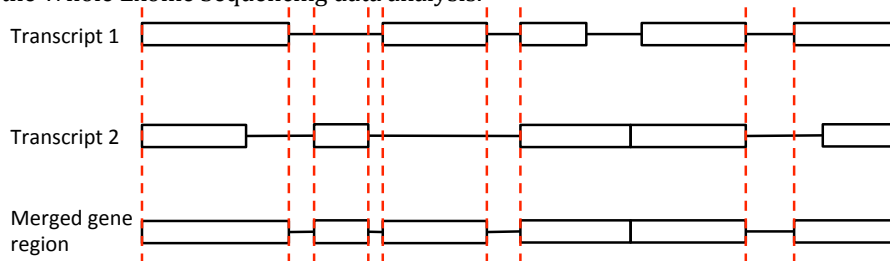
**Figure S8**

TCGA Whole Exome Sequencing samples by cancer types.

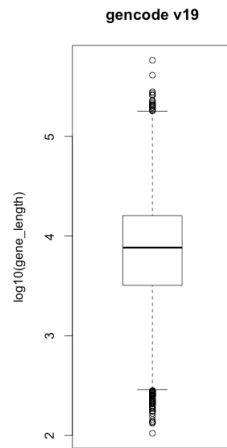


**Figure S9**

Details of gene regions definition. Note that only coding transcripts were used for the Whole Exome Sequencing data analysis.



**Figure S10**  
Distribution of gene length



**Figure S11**  
Distribution of the pooled mutation rates

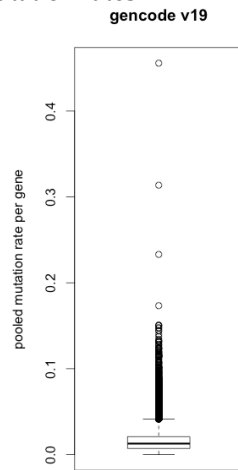
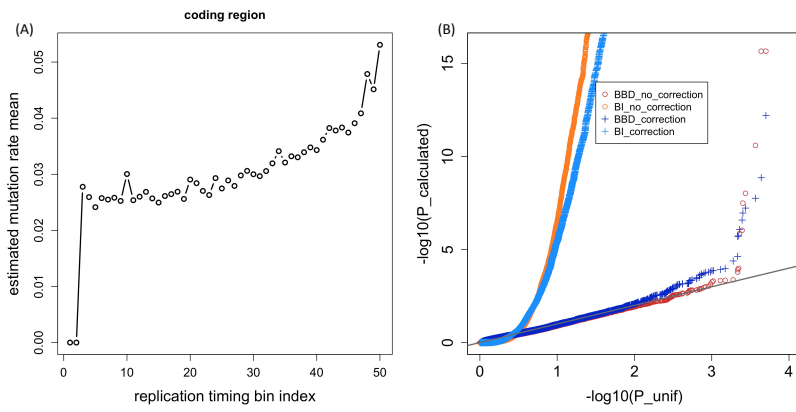


Figure S12

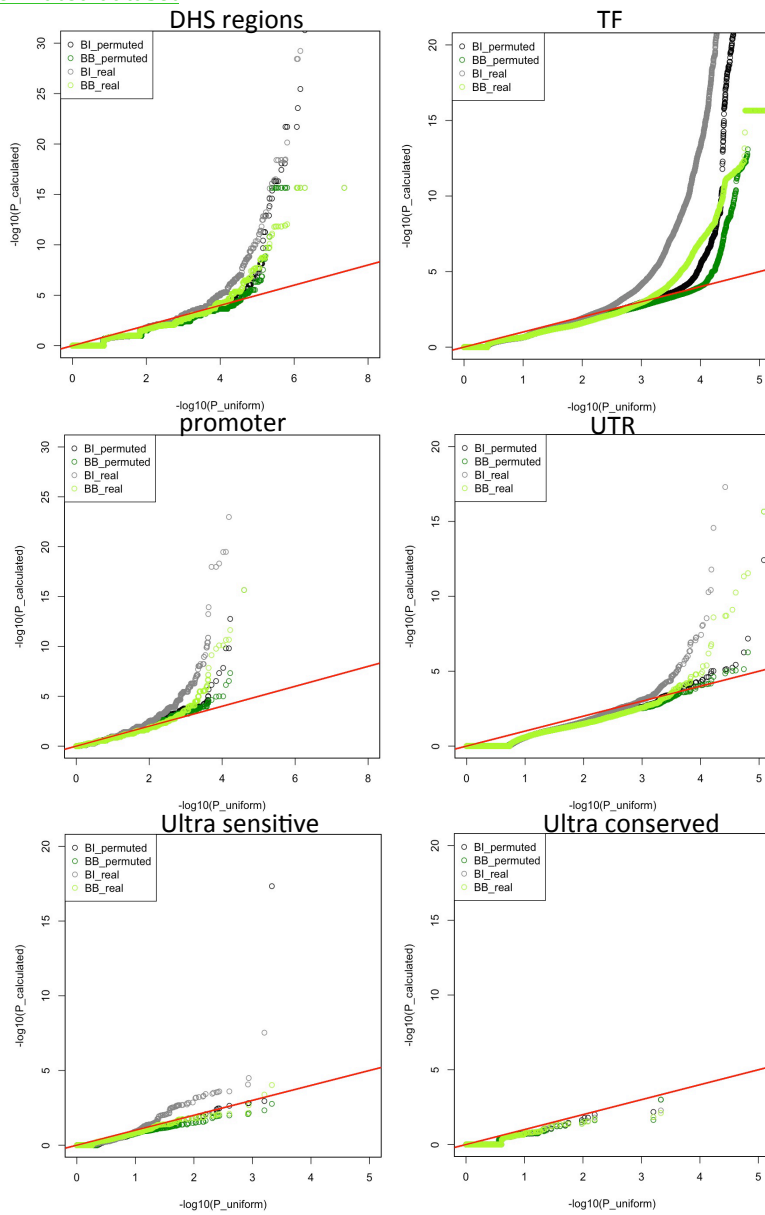
In the coding regions: (A). the average mutation rate vs. replication timing; (B). QQ plots of calculated P values VS. uniform theoretical ones of the coding region analysis.



Unknown  
Formatted: Font:(Default) Helvetica  
Jing Zhang 7/7/2015 6:23 PM  
Deleted: <sp>

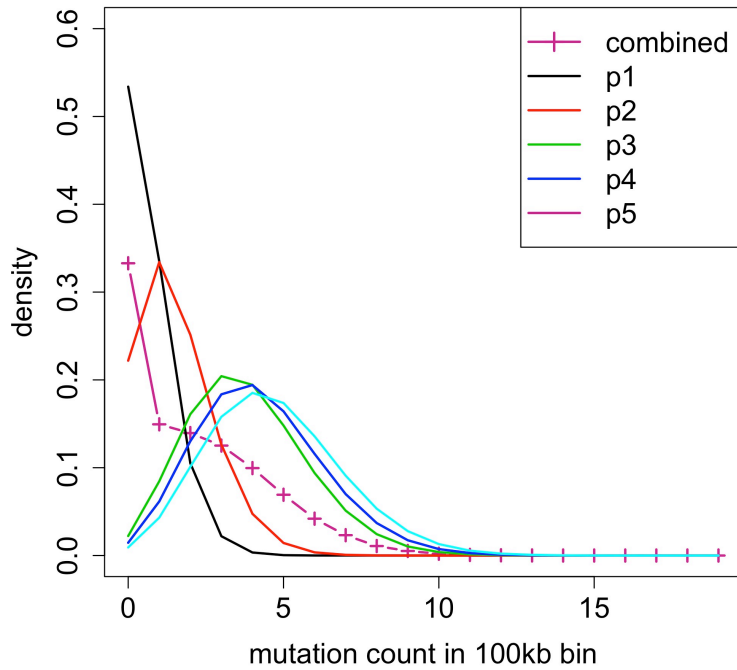
Jing Zhang 7/7/2015 1:26 PM  
Formatted: Heading 3

Figure S13: QQ plots of calculated P values VS. uniform theoretical ones of the real and permuted dataset.



Jing Zhang 7/8/2015 4:54 PM  
 Deleted: <sp>  
 Unknown  
 Formatted: Font:(Default) Helvetica  
 Jing Zhang 7/11/2015 11:52 AM  
 Deleted: of the whole genome 10kb bins analysis. The discrete dots in the top two figures are due to the limited number of unique mutation counts  
 Jing Zhang 7/11/2015 11:51 AM  
 Formatted: Normal  
 Jing Zhang 7/11/2015 11:51 AM  
 Formatted: Font:+Theme Body

Figure S14: Importance of covariant correction



Jing Zhang 7/7/2015 1:27 PM  
Formatted: Heading 3, Left

Jing Zhang 7/7/2015 1:27 PM  
Formatted: Font: Cambria, Not Bold, Font color: Text 1

Jing Zhang 7/7/2015 1:27 PM  
Formatted: Font: Theme Body

Jing Zhang 7/7/2015 1:27 PM  
Formatted: Justified

### Supplementary tables

Table S1

Summary of the whole genome sequencing cancer data used in this study

Cancer Type	# of Samples
Acute Lymphoblastic Leukemia	1
Acute Myeloid Leukemia	7
Breast Cancer	119
Chronic Lymphocytic Leukemia	28
Glial Tumor	26
Kidney Carcinoma	32
Liver Cancer	88
Lung Adenocarcinoma	24
Lymphoma B-cell	24



Medulloblastoma	100
Pancreatic Cancer	15
Pilocytic Astrocytoma	101
Prostate Cancer	95
Stomach Cancer	100
Sum	760

**Table S2**

Percentage of coding mutations in all samples (attached in supplementary file)

**Table S3**

Genes with significant mutation burden, according to LARVA's exome analysis.

Gene	Simple annotation	Supporting Reference
<b>TP53</b>	Well-known oncogene	PMID:20182602
<b>BRAF</b>	B-Raf proto-oncogene	<a href="http://ghr.nlm.nih.gov/gene/BRAF">http://ghr.nlm.nih.gov/gene/BRAF</a>
<b>KRTAP4-11</b>	Unknown	
<b>IDH1</b>	Glioblastomas, astrocytomas, oligodendroglial tumors	PMID:19435942
<b>FRG1B</b>	lineage-specific mutation patterns in many cancer types	PMID: 24465236
<b>CDKN2A</b>	pancreatic cancer	PMID: 21150883
<b>PRSS1</b>	pancreatic cancer	PMID:22379635

## References

1. Young-Xu, Y. and Chan, K.A. (2008) Pooling overdispersed binomial data to estimate event rate. *BMC medical research methodology*, **8**, 58.

2. Kleinman, J.C. (1975) Proportions with extraneous variance: two dependent samples. *Biometrics*, **31**, 737-743.
3. Lawrence, M.S., Stojanov, P., Polak, P., Kryukov, G.V., Cibulskis, K., Sivachenko, A., Carter, S.L., Stewart, C., Mermel, C.H., Roberts, S.A. *et al.* (2013) Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature*, **499**, 214-218.