

LARVA: an integrative framework for Large-scale Analysis of Recurrent Variants in Annotations

PROT CODING

Lucas Lochovsky¹, Jing Zhang², Yao Fu¹, Ekta Khurana², and Mark Gerstein^{1,2,3,*}

¹Program in Computational Biology and Bioinformatics, Yale University, New Haven, Connecticut 06520, USA

²Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, Connecticut 06520, USA

³Department of Computer Science, Yale University, New Haven, Connecticut 06520, USA

* To whom correspondence should be addressed. Tel: +1 203 432 6105; Fax: +1 360 838 7861; Email: Mark.Gerstein@Yale.edu

Present Address: Mark Gerstein, Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, Connecticut 06520, USA

New 12/15/2014 12:02 AM
Deleted: GersteinYale

ABSTRACT

Noncoding variants are known to be associated with numerous diseases, but they are barely investigated due to their poor functional interpretation. Moreover, extensive overdispersion in the noncoding mutation count data, due to mutation rate heterogeneity or potential dependencies between neighboring nucleotides, further complicates the significance assessment in the mutation burden test. We address these issues with the development of a new computational framework called LARVA (Large-scale Analysis of Recurrent Variants in Annotations) for both germline and somatic variant analyses. LARVA first integrates a comprehensive set of noncoding elements from the ENCODE project, which span both proximal and distal gene regulators. Then, it models the mutation count in each regulatory element as a beta-binomial distributed random variable to handle the observed overdispersion. Furthermore, it incorporates regional genomic features like replication timing to achieve a more accurate local background estimate. Consequently, LARVA accurately evaluates the observed variant counts against this local null model to identify highly mutated regulatory elements in the human genome. We demonstrate LARVA by analyzing a set of 760 cancer whole genome sequences for recurrent mutations. LARVA is available at [url].

MORE

?

New 12/15/2014 12:02 AM
Deleted: difficulties in ...heir poor ... [1]

INTRODUCTION

Genomes of numerous patients have been sequenced (1-5), opening up opportunities to identify the underlying genetic causes for complex disease (6-9) and develop more effective therapies targeted at specific molecular disease subtypes (10). Most of these studies have so far focused on identifying mutations and defects in the protein coding regions, or exomes, of disease genomes (2, 11-14). These methods usually search for regions with higher than expected mutation frequencies in protein coding genes through rigorous background mutation rate control over a variety of genomic features (11). Such methods have been successfully used on numerous cancer genomes (15). However, the noncoding regions, which comprise more than 98% of the human genome, were rarely investigated, primarily due to the difficulty of functional interpretation of noncoding variants.

FINDING

New 12/15/2014 12:02 AM
Deleted: cancer ...atients have bec... [2]

The authors wish it to be known that, in their opinion, the first 2 authors should be regarded as joint First Authors

New 12/15/2014 12:02 AM
Formatted: Font color: Auto

PROT CODING REGIONS

Recent genome annotation analysis has revealed [a significant portion](#) of the human genome is functional in a certain tissue or development stage (16,17), and several noncoding variants has been implicated in disease (18). For example, [several](#) genome-wide association studies (GWASs) studies have discovered the phenotypic effect of common noncoding variants in regulatory regions (19,20). Other studies have reported that noncoding TERT mutations drive cancer progression in multiple tumor types, including melanomas and gliomas (21-23). Moreover, mutations in the promoter regions of PLEKHS1, WDR74 and SDHD were also identified as recurrent driver mutations in some cancer types (24). In another example, analysis of the miRNA-binding sites of BRCA1 and BRCA2, [the](#) established drivers of breast cancer, indicated that certain variants in these sites are associated with increased risk of early onset breast cancer (25). Histones also serve as important noncoding regulators, as demonstrated in an analysis of a histone H1 variant linked to oncogene expression in ovarian cancer (26). In light of these discoveries, and the growing availability of whole-genome sequencing data (2,27-32), a statistical framework facilitating the identification of [highly mutated](#) noncoding mutations would be useful.

New 12/15/2014 12:02 AM

Deleted: that up to 80%... significa ... [3]

More recently, ~~a pioneering~~ genome wide computational effort has been made to discover the noncoding [regions with higher mutation burden](#) in cancer genomes (24). The authors called whole genome somatic variants for 863 human tumor sequences from The Cancer Genome Atlas (TCGA) (28), and analyzed the variants that fall into noncoding annotations. A [p-value](#) was computed for each annotation reflecting the likelihood that the given annotation had more variants than expected from background mutation processes, which was modelled with a binomial distribution. [It successfully identified some known noncoding drivers, such as the TERT promoter, and reported some novel candidates that were not discovered previously.](#) The use of the binomial distribution is based on two assumptions: 1) the mutation rate is homogeneous; 2) [variants mutate independently.](#) However, [human genomes often violate these assumptions. First, studies on the coding variants already proved that the mutation rates in cancer genomes demonstrate](#) substantial cancer type, sample, and regional heterogeneity (11). [Second, in regions with high linkage disequilibrium \(LD\), it is observed that germline variants, and possibly somatic variants as well, tend to be mutated in a highly associated way.](#) Consistent with these statements, we observed that the somatic mutation counts in the noncoding elements exhibited substantially higher variance than expected, the so called overdispersion, indicating that binomial distribution might be potentially inadequate to handle such data, and the resultant [p-values might be](#) heavily inflated.

New 12/15/2014 12:02 AM

Deleted: the first... pioneering gen ... [4]

Here, we present a computational system, LARVA (Large-scale Analysis of Recurrent Variants in Annotations), that identifies [highly mutated noncoding regulatory elements](#) using whole genome sequencing (WGS) variant data from multiple genetic disease patients. [LARVA treats the mutations counts within a given regulatory element as a beta-binomial distributed random variable.](#) This design automatically accommodates the heterogeneous nature of mutation accumulation in cancer genomes and the potential dependency among neighboring loci by allowing [the mutation rate to be drawn from](#) a beta distribution. Furthermore, we also [divided the whole genome into several local bins and classified them using some known genomic confounders of the mutation rate, such as replication](#)

New 12/15/2014 12:02 AM

Formatted ... [5]

New 12/15/2014 12:02 AM

Deleted: recurrent patterns of disease mutations in various genome ... [6]

ON
2

timing, for a more accurate local background [mutation model](#) calculation. [Such integrative analysis](#) could potentially control the false positive rate in an effective manner. We demonstrate the usefulness of LARVA for finding noncoding elements [with higher mutation burdens](#) in a set of WGS cancer data that represents all the different types of whole genome sequenced cancers as of this writing (see Methods for details). [Besides the improved performance in cancer somatic mutation analysis, LARVA can also be immediately used for germline variant burden tests.](#) The following sections describe LARVA's concepts, their applications to the study of genetic disease, and our cancer findings.

MATERIAL AND METHODS

Whole genome cancer variant data

We collected whole genome cancer variant calls from a large number of previously sequenced cancer genomes. The majority of our data came from a set of 507 whole genome cancer samples published in Alexandrov *et al.* (27). This data spans breast cancer, lung cancer, leukemia, pancreatic cancer, pilocytic astrocytoma, medulloblastoma, liver cancer, and lymphoma (Fig 1 A and supplementary table 1). This was supplemented with a collection of 95 prostate cancer samples we obtained from publications (2,28-30), a set of 26 unpublished glioma tumor samples, 32 kidney cancer samples from the TCGA (28), a set of 100 stomach cancer samples from Wang *et al.* (31).

→ MAKE UP A TABLE

Quality control of the WGS variants

A number of genomic regions are known to have poor read mappability due to sequence phenomena that cause ambiguous mapping results, such as a large number of tandem repeats. These regions are known as [signal artifact blacklist regions](#) (33). Since it is likely that variant calls in this region are possibly inaccurate, we opted not to use these variants in our mutation rate calculations (details in Fig. S1). Blacklist regions were derived from (33), and downloaded from the UCSC Genome Browser. Variants intersecting these regions, as determined by BEDTools (34), were removed from the analysis.

Noncoding annotation summary

Our analysis covered a range of noncoding regulatory annotations. The GENCODE v16 main annotation file was parsed to derive the coordinates of regulatory annotations close to gene regions, including promoters and untranslated regions (UTRs)(35). Transcription factor (TF) binding sites were derived from the [Chip-seq](#) experiments conducted as part of the ENCODE project (36). We collected the full list of TF binding sites in all possible [tissues](#) and cell lines from ENCODE. Distal regulatory modules (DRM) enhancers, which regulate the expression of genes at non-adjacent sites, were derived from (37). Another class of regulators, the Dnase I hypersensitive (DHS) sites (38), were also derived from the ENCODE project. Additionally, we added a set of sites deemed "ultra-conserved" in (39), [due to their extremely high level of conservation across many species.](#) Furthermore, we used a set of "ultra-sensitive" sites from (40), so named because [they are noncoding regions under higher selective pressure from the population genetics point of view.](#) Finally, we studied transcription start sites (TSSes) by extracting the 100bp regions immediately upstream of GENCODE gene coding annotations (35).

Handwritten blue annotations: a large bracket on the right side of the text, and a signature-like mark at the bottom.

New 12/15/2014 12:02 AM

Deleted: rate

New 12/15/2014 12:02 AM

Deleted: , which

New 12/15/2014 12:02 AM

Deleted: recurrent mutated

New 12/15/2014 12:02 AM

Deleted: vast

New 12/15/2014 12:02 AM

Deleted: 25

New 12/15/2014 12:02 AM

Deleted: 26-

New 12/15/2014 12:02 AM

Deleted: 26

New 12/15/2014 12:02 AM

Deleted: 29

New 12/15/2014 12:02 AM

Deleted: , and a few cancer types from the PCAWG-50 dataset (30) not already represented in our other datasets, accounting for an additional 25 samples.

New 12/15/2014 12:02 AM

Deleted: -

... [7]

New 12/15/2014 12:02 AM

Deleted: 33

New 12/15/2014 12:02 AM

Deleted: Peak

New 12/15/2014 12:02 AM

Deleted: 34

New 12/15/2014 12:02 AM

Deleted: tissue

New 12/15/2014 12:02 AM

Deleted: 35

New 12/15/2014 12:02 AM

Deleted: 36

New 12/15/2014 12:02 AM

Deleted: 37

New 12/15/2014 12:02 AM

Deleted: due to their extremely high level of selection pressure imputed from large population genetic studies (16).

New 12/15/2014 12:02 AM

Deleted: 38

New 12/15/2014 12:02 AM

Deleted: of their high functional impact of variants in these

New 12/15/2014 12:02 AM

Deleted: 33

New 12/15/2014 12:02 AM

Deleted: . For comparison, we also created a collection of coding exon annotations from GENCODE.

Pseudogenes are known hot spots for artifacts due to their high context resemblance to their parent genes. In order to avoid potential variant calling bias, partially due to mapping difficulty, we removed the promoters, TSS, and UTR analyses for pseudogenes in the GENCODE annotation (details in Fig. S2 Text S1 section 1).

New 12/15/2014 12:02 AM
Deleted: the know...nown hot spot... [8]

Significance evaluation of mutation recurrence

We used our set of cancer variant calls to derive a "null model" representative of the mutation rate that would be expected due to background stochastic mutation processes. Some previous models (24,41) assumed that this rate is constant over the entire genome, and mutations occur in an independent way. Hence, the entire genome's expected mutation can be represented as a binomial distribution $B(n, p)$ where n is the number of basepairs (bp) of the regulatory element, and p is the probability that a single nucleotide is mutated. However, due to the heterogeneous nature of the cancer genomes, and the possible dependency among neighboring loci, the binomial distribution might be inadequate to describe the mutation count data. Also, taking into consideration the confounding influence of other genomic features, such as the replication timing, we proposed the following hierarchical model.

New 12/15/2014 12:02 AM
Deleted: whole genome ...utation r... [9]

New 12/15/2014 12:02 AM
Deleted: n

New 12/15/2014 12:02 AM
Deleted: in...f the human... [10]

New 12/15/2014 12:02 AM
Deleted: basepair...ucleotide is m... [11]

$$x_i | p_i \sim \text{Binomial}(n_i, p_i)$$

$$p_i \sim \text{Beta}(\alpha_i, \beta_i)$$

α_i, β_i : constant within the same r bin

where the mutation count is $x_i, i = 1, 2, \dots, k$, and the sample size, mutation probability, and replication timing can be expressed as n_i, p_i and r_i . Instead of the fixed mutation rate assumption, we provided p_i more flexibility by allowing it to follow a beta distribution $\text{Beta}(\alpha_i, \beta_i)$. The beta distribution can be expressed as

New 12/15/2014 12:02 AM

$$\pi(p | \alpha, \beta) = \text{Beta}(\alpha, \beta) = \frac{p^\alpha}{\Gamma(\alpha)\Gamma(\beta)} \Gamma(\alpha + \beta)$$

Deleted:

New 12/15/2014 12:02 AM
Deleted: $x_i, i = 1, 2, \dots, k$... EM... [13]

$$\pi(p_i | \alpha_i, \beta_i) = \text{Beta}(\alpha_i, \beta_i) = \frac{p_i^{\alpha_i-1} (1-p_i)^{\beta_i-1}}{\Gamma(\alpha_i)\Gamma(\beta_i)} \Gamma(\alpha_i + \beta_i)$$

To correct the effect of replication timing, we divided the regulatory elements into 10 bins according to their averaged replication timing signal, and assume that the beta distribution parameters are constant within the same bin. Then the marginal distribution of the total number of mutations within the bin with length n_i follows the beta binomial distribution.

New 12/15/2014 12:02 AM
Deleted: is a ...ccording to their... [14]

STUCK CONFUS

?

?

↑

$$\Pr\{X = x_i\} = \binom{n_i}{x_i} \frac{\Gamma(\alpha_i + \beta_i) \Gamma(\alpha_i + x_i) \Gamma(\alpha_i + n_i - x_i)}{\Gamma(\alpha_i) \Gamma(\beta_i) \Gamma(\alpha_i + \beta_i + n_i)}$$

The moment estimator mentioned in (42,43) was used to estimate the parameters in the beta-binomial distribution, and the p-values were calculated accordingly against the fitted null model (for details see section 2 in Text S1).

RESULTS

Overview of noncoding elements on various cancer genomes

We sought to study the whole genome somatic mutation patterns of as many different cancer patients as possible. To that end, we collected whole genome cancer variant call sets from a range of cancer data repositories (27,28) and publications (2,27,29-32). Our data spans 760 genomes, and includes 14 types of cancer (Fig 1A and Supplementary Table S1).

For our noncoding analysis, we conducted a literature survey of functional noncoding regions (Table 1). Our list consists of a number of proximal gene regulators, such as promoters, untranslated regions (UTRs), and transcription start sites (TSS) (35). We also collected a series of distal regulatory module (DRM) enhancers, which regulate expression for genes in other parts of the genome (37). We also collected the full set of transcription factor binding sites (TFBS) from ENCODE Chip-seq experiments (36) and Dnase I hypersensitive (DHS) sites from DNase-seq experiments (34). Our noncoding annotation set was further augmented by a list of "ultra-conserved" regions, a term defined in (39) that refers to sites with almost perfect conservation across many species. Our final annotation set consists of sites that are considered "sensitive" to mutation due to the high functional impact of variants in these regions, and therefore are called "ultra-sensitive" sites (40).

Fig 1B illustrates our procedure for analyzing cancer variant data for significant recurrently mutated annotations. The cancer variants in VCF format pass through a quality control filter that include removing those variants that fall into blacklist regions. The preprocessed variants, along with our collected set of noncoding annotations that do not overlap blacklist regions, are used in the main computation. The main processing step includes counting all variant intersections with the noncoding annotations, and adjusting the expected background mutation rate with regional mutation rate corrections, such as DNA replication timing. A DNA replication timing correction is useful because a positive correlation has been observed between genome regions' mutation rate and time of replication in the cell's S phase (43). Our software then fits the beta binomial distribution described in our methods to the observed cancer variant data, and produces a list of the top noncoding annotations with significant mutations.

As it is shown in table 1, our noncoding annotation list spans approximately 30% of the human genome. We observed different cancer types demonstrate distinct mutational preferences over these noncoding regions. To illustrate this phenomenon, we used 11 types of cancer from our overall

INTUITION

DIAGNOSTIC

New 12/15/2014 12:02 AM

Deleted: $\Pr\{X = x_i\} = \binom{n_i}{x} \frac{\Gamma(\alpha + \dots [15]$

New 12/15/2014 12:02 AM

Deleted: were

New 12/15/2014 12:02 AM

Deleted: P

New 12/15/2014 12:02 AM

Deleted: 25,26

New 12/15/2014 12:02 AM

Deleted: 25,

New 12/15/2014 12:02 AM

Deleted: -30

New 12/15/2014 12:02 AM

Deleted: 785

New 12/15/2014 12:02 AM

Deleted: 27

New 12/15/2014 12:02 AM

Deleted: 33

New 12/15/2014 12:02 AM

Deleted: 35

New 12/15/2014 12:02 AM

Deleted: Clip

New 12/15/2014 12:02 AM

Deleted: 34

New 12/15/2014 12:02 AM

Deleted: (40).

New 12/15/2014 12:02 AM

Deleted: 37

New 12/15/2014 12:02 AM

Deleted: 38

New 12/15/2014 12:02 AM

Deleted: (41).

New 12/15/2014 12:02 AM

Deleted: recurrent

dataset for which there are at least 20 samples and calculated the fraction of WGS mutations within each noncoding element category (boxplots of different color in Fig. 2). The overall nucleotide percentage of each annotation over the genome was used as the background (black dash lines in Fig.2). In one instance representative of the large differences observed between cancer types, variants in kidney cancer was found to be preferentially located in the TF binding site while lung adenocarcinoma is mutation depleted in this region (0.140 average vs. 0.098 average, light green color vs. aquamarine color in Fig. 2). Huge sample effect was also observed in several cancer types. For instance, within Pilocytic Astrocytoma, there are samples that have a TF binding peak mutation fraction as high as 0.252 and as low as 0.011, which represents a ~23-fold difference. Hence, it is important to understand the mutation patterns in these annotations, and take their unique characteristics into consideration.

New 12/15/2014 12:02 AM
Deleted: plots
New 12/15/2014 12:02 AM
Deleted: compared to

SPR Use

Large cancer type, sample, regional heterogeneity of cancer genomes, and the potential dependency among neighboring regions violate the binomial assumption

New 12/15/2014 12:02 AM
Deleted: and

In (24), the mutation burden tests are performed based on the binomial distribution, which inherently assumes a constant mutation rate and completely independent mutation events. However, these assumptions might not be appropriate for either somatic or germline variant analysis.

New 12/15/2014 12:02 AM
Deleted: In (22), the mutation recurrence p-values are calculated based on a binomial distribution, which is based on the assumption of a constant mutation rate. However, in our analysis of hundreds of WGS mutation In (22), the mutation recurrence p-values can be calculated from binomial distribution, which is based on the assumption of a constant mutation rate. However, in our analysis of hundreds of WGS mutation signatures, we observed huge cancer type, sample, and regional somatic mutation rate heterogeneity.

First, in our analysis of hundreds of WGS somatic mutation signatures, we observed huge cancer type, sample, and regional somatic mutation rate heterogeneity. To demonstrate cancer type and sample mutation rate heterogeneity, we selected all cancer types with more than 20 samples in it. We split the human genome into 1 mega basepair (Mbp) size bins, and intersected the individual sample variants from our data set to calculate the mutation rate of each sample. Consistent with the analysis in coding regions (11), we observed huge mutation rate differences between cancer types. For instance, the average whole genome mutation rate in stomach cancer is as high as 11.389 mutations/Mbp (violet red colors in Fig 3.A), which is ~800 times of the mutation rate in medulloblastoma (0.0142, blue colors in Fig 3.A). Furthermore, the whole genome mutation rate also fluctuates wildly across samples, and such changes may go up to 100 times within the same cancer type (0.359 VS. 21.8 in breast cancer for example). Additionally, to illustrate regional mutation rate heterogeneity, we randomly selected 50 one-megabase-length regions to calculate the mean and standard deviation (SD) of the local mutation rate across samples in lung cancer and prostate cancer (Fig 3. B, dashed lines show the SD). As shown in Fig.3 B, the average local mutation may vary from 0 to 50.8 mutations/Mbp across the randomly selected bins, and the SD range is unusually huge for each bin. Similar results were also observed in prostate cancer (red dots and lines in Fig 3. B).

New 12/15/2014 12:02 AM
Deleted: mb
New 12/15/2014 12:02 AM
Deleted: 9
New 12/15/2014 12:02 AM
Deleted: of 0.0142
New 12/15/2014 12:02 AM
Deleted: .

Several biological signatures could partially explain the observed mutation rate heterogeneity. Replication timing is one of the well-known genomic features that affects both germline and somatic mutation rate. The later replicating regions accumulate DNA damage, such as oxidation and deamination, making them prone to mutations of all sources (45). Besides, methylated cytosines in CpG sites are often unstable and undergo deamination to thymine, which yields a C to T transition

New 12/15/2014 12:02 AM
Deleted:

New 12/15/2014 12:02 AM
Deleted: Mutation

SPR Use

(46). Hence, there is a noticeable mutation rate difference at CpG and non CpG sites. Several other hypothesis were also proposed and summarized in the excellent review paper (46).

Second, mutation events might not be independent of each other. For example, in germline mutation analysis, mutations with high LD are prone to mutate together. Additionally, some passenger mutations are generated by other driver mutations. The driver mutation might be a mutation in a DNA replication or repair gene. Moreover, some structural variations, such as long insertions or deletions, might cause problems in pairing during meiosis and thus to generate additional point mutations in neighboring regions (47). Consistent with this hypothesis, the mutation rates of the surrounding structure variations are elevated in several eukaryotic-species (47-49).

Improved mutation count fitting with a beta-binomial distribution

As shown in Fig.3, cancer genomes display large mutational heterogeneity due to various factors, therefore the theoretical binomial distribution model used in (24), which assumes a constant mutation rate, could be problematic in more practical data analysis applications. Consistent with this assumption, we did observe a much higher than expected variance in the mutation count data. For example, at a 10kb bin resolution, the observed mutation count variance is 7.679 times of the expected valued under the binomial assumption. Hence, it is necessary to introduce other statistical models to handle such overdispersion in the mutation count data.

Instead of assuming a homogenous mutation rate, we utilized a beta distribution to describe the mutation rate more flexibly. This two-parameter distribution conveniently provides the underlying mutation rate with desired mean and variance properties, and could also easily model the mutation count data within a specific region as a beta-binomial distribution (details in methods). We fitted the mutation count data at a 10kb bin resolution of the 760 WGS cancer genomes under the fixed (binomial) and variable (beta-binomial) mutation rate assumptions in Fig. 4. We calculated the frequency of the observed mutation count in each bin and compared it with the binomial and beta-binomial fittings respectively. It is shown in Fig. 4 A that the observed data (black dots) demonstrates much heavier tails than the binomial distribution (purple dots), while the beta-binomial distribution (green dots) fits very well at the right tail. In order to quantitatively exhibit the improved performance of beta-binomial fitting, we utilized Kolmogorov-Smirnov (KS) statistics to compare the two distributions with the observed data in a nonparametric way. A larger KS statistic indicates a higher level of deviation between the two distributions. Specifically, 1000 bins were simulated from beta-binomial and binomial fitted distributions separately to calculate the KS statistic against the randomly sampled 1000 mutation counts from the observed data. This scheme was repeated 1000 times and the cumulative distribution function (C.D.F) of the KS statistics were given in Fig. 4B. The median KS statistic value for the beta-binomial distribution was 0.087, significantly smaller than 0.218 of the binomial distribution (p-value for two-sided Wilcoxon test $< 2.2 \times 10^{-16}$, boxplots given in Fig. 4C). Different bin sizes were analyzed using the sample method and results were similar (Fig. S3-Fig. S4). In order to avoid overfitting, we utilized half of the data for distribution fitting, and the remaining half as the input to calculate the KS statistic for evaluation. This scheme was repeated for 100 times. The

New 12/15/2014 12:02 AM
Deleted: improves cancer genome heterogeneity modelling

New 12/15/2014 12:02 AM
Deleted: (22)

New 12/15/2014 12:02 AM
Deleted: 78.901

New 12/15/2014 12:02 AM
Deleted: larger than

New 12/15/2014 12:02 AM
Deleted: mutation heterogeneity and the resultant

New 12/15/2014 12:02 AM
Deleted: allow

New 12/15/2014 12:02 AM
Deleted: to change over regions in our model.

New 12/15/2014 12:02 AM
Deleted: flexibly

New 12/15/2014 12:02 AM
Deleted: conveniently

New 12/15/2014 12:02 AM
Deleted: 785

New 12/15/2014 12:02 AM
Deleted: the empirical distribution from

New 12/15/2014 12:02 AM
Deleted: counts

New 12/15/2014 12:02 AM
Deleted: 220

New 12/15/2014 12:02 AM
Deleted: the problem of

beta-binomial distribution still significantly outperforms the binomial distribution (0.0821 vs. 0.216, p-value for two sided Wilcoxon test $< 2.2 \times 10^{-16}$, Fig. S5). Hence, the improved performance of the beta-binomial distribution is due to its enhanced flexibility to handle the overdispersed mutation count data instead of overfitting.

In the significance analysis of recurrent mutations, p-values were usually calculated from the right tails of the null distribution. But the huge deviation of the binomial distribution from the observed one could potentially introduce huge p-value inflation, and consequently result in numerous false positives. We defined the p-values for the observed distribution as the percentage of bins with equal or larger mutation counts. However, the improved fitting of the beta-binomial distribution could solve this problem and provide more accurate p-value assessment.

Replication timing correction helps to control both false positives and negatives

Recently, several computational efforts have been made to link the somatic mutation rates with several genomic features in the protein-coding regions (11,46). A particularly well-known example is DNA replication timing. During replication, the single stranded DNA usually suffers from endogenous DNA damage, such as oxidation and deamination (45). Hence, DNA that is replicated in a later stage would be susceptible to the effects of accumulative damage, and would be prone to all classes of substitutions. Consistent with this assumption, scientists observed that the later replicating regions demonstrate remarkably higher mutation rate (45). Although replication timing has been used successfully in the coding regions, little work has been done in the noncoding regions in cancer genomics. Hence, we explored the effect of replication timing on the mutation rate calculation, and the consequential effect on the p-value evaluation.

Using 1kb bins, we counted the average replication timing value within the bin, and then separated the top and bottom 10% of replication timing bins for mutation rate calculation. As shown in Fig. 5 A, we observed noticeable differences in the mutation rate vis-a-vis the replication timing signal. The average mutation count of the 760 samples was 1,200 for the bottom 10% replicating timing bins, as compared to 4,028 for the top 10 percent counterparts (p-value for two-sided Wilcoxon test $< 2.2 \times 10^{-16}$). A KS test was performed to determine whether these two sets of mutation counts data follow the same distribution, and the p-value is less than 2.2×10^{-16} , indicating that the two distributions are significantly different.

Moreover, we observed that the mutation counts data for bins with similar replication timing values still shows extensive overdispersions. For example, for the bottom 10% replication timing bins, the observed variance of mutation counts was 4,168, which is 3,477 times that under the binomial assumption. Consistently, we observed poor fitting of binomial distribution against the observed distribution, especially in the right tails (dark black bars vs. dark purple bars in Fig. 5A). The huge deviation in the right tails would result in huge p-value calculation inflation as shown in Fig. 5B. The p-value for 16 mutations in the bottom replication timing 1kb region from the empirical distribution shows only marginal significance (3.994×10^{-4}), but the binomial distribution could inflate it to 2.585×10^{-13} due

New 12/15/2014 12:02 AM

Deleted: 0824

New 12/15/2014 12:02 AM

Deleted: 218

New 12/15/2014 12:02 AM

Deleted: empirical

New 12/15/2014 12:02 AM

Deleted: Covariant

New 12/15/2014 12:02 AM

Deleted: 9,42

New 12/15/2014 12:02 AM

Deleted: The later replicating regions were reported to demonstrate remarkably increased mutation rate

New 12/15/2014 12:02 AM

Deleted: (41).

New 12/15/2014 12:02 AM

Deleted: vast majority of

New 12/15/2014 12:02 AM

Deleted: median

New 12/15/2014 12:02 AM

Deleted: counts

New 12/15/2014 12:02 AM

Deleted: 785

New 12/15/2014 12:02 AM

Deleted: 241

New 12/15/2014 12:02 AM

Deleted: 136

New 12/15/2014 12:02 AM

Deleted: (p-value inflation due to replication timing)

New 12/15/2014 12:02 AM

Deleted: 422

New 12/15/2014 12:02 AM

Deleted: 568

New 12/15/2014 12:02 AM

Deleted: larger than expected

New 12/15/2014 12:02 AM

Deleted: empirical

New 12/15/2014 12:02 AM

Deleted: 15

New 12/15/2014 12:02 AM

Deleted: 0.0005097721

New 12/15/2014 12:02 AM

Deleted: 4.250e-

New 12/15/2014 12:02 AM

Formatted: Superscript

to its [pad fitting of the heavy tails on the right side](#). But our beta-binomial distribution rigorously controls the p-values through the flexible mutation rate assumption (p-value = 1.002×10^{-3}). We demonstrated the better p-value curve of the beta-binomial distribution in a variety of data points and replication timings, indicating the robustness of our method (Fig. 5B).

Additionally, the replication timing effect correction further improves the p-value calculation to avoid potential false positives and false negatives. For instance, [for a region with the top replication timing regions, 8 mutations in 1kb bin would give a p-value at 0.094 after replication correction from beta-binomial model, but might be reported as positive when ignoring replication timing effect \(p-value = 0.038 from beta-binomial by mixing the top and bottom 10% replication timing points\)](#). Similarly, [a p-value of 0.064 would refuse 7 mutations within 1kb bin as significant without correction](#). However, [if this point comes from the bottom 10% replication timing region, the true p-value should be 0.030 due to its relatively lower local mutation rate](#). Hence, it is important to perform covariate correction before calculating p-values.

LARVA discovered a list of highly recurrent noncoding regulatory regions from WGS data

We first applied LARVA to the [760](#) genomes' variants, intersecting them with the noncoding regions listed in Table 1. In total, LARVA reported [3964](#) and [3776](#) highly [mutated](#) regions before and after replication timing corrections, respectively (as shown in Table 2). On the other hand, the binomial distribution models reported at least 30 times more regions as significant because of the aforementioned p-value inflation, giving rise to a high false positive rate. We also tested the immediate 100bp upstream of every possible transcription start sites (see methods for details), the results of which are depicted in Fig. 6 B. Forty-five TSSs passed the 0.05 p-value thresholds after p-value adjustment (BH method, [\(50\)](#)). Consistent with previous studies, we observed that the TSS for TERT came up in the top regions (Fig. 6 B), and the oncogene TP53 also ranked second among all sites. LMO3, which ranked third after replication timing correction, is a protein-coding oncogene that is predominantly expressed in brain tissue. It has been reported to be involved in a variety of cancer types, such as lung cancer [\(51\)](#) and neuroblastomas [\(52\)](#). PRRC2B's TSS was reported as the most significantly recurrent region among all TSSes. It is a protein-coding gene that is extensively expressed in brain tissue, but to our best of knowledge, there is no study to show the link of PRRC2B to cancer. Further investigations should be performed for the purpose of validation. [Similar](#) results were given for promoters and UTR regions as well. We selected all the genes with highly mutated TSSes, promoters, or UTRs (adjusted p-values after corrections ≤ 0.05) and performed GO analysis (<http://amigo.geneontology.org>, [\(53\)](#) results in [Table S2](#)). The top three enriched GO terms are: "negative regulation of fibroblast proliferation", "regulation of extrinsic apoptotic signaling pathway in absence of ligand", and "regulation of cell growth".

In terms of transcription factor binding sites, LARVA claimed [2054](#) out of the 5710954 binding sites as highly recurrent (0.036%). The transcription factor CTCF had [852](#) binding sites reported as significant (Table 3). CTCF is a multifunctional protein that is linked with multiple cancer types [\(54\)](#).

New 12/15/2014 12:02 AM
Deleted: failure to properly model
New 12/15/2014 12:02 AM
Deleted: 0.001188101

New 12/15/2014 12:02 AM
Deleted: the empirical p-values

New 12/15/2014 12:02 AM
Deleted: 9 mutation counts

New 12/15/2014 12:02 AM
Deleted: a

New 12/15/2014 12:02 AM
Deleted: in

New 12/15/2014 12:02 AM
Deleted: regions are 0.101 and 0.019, respectively. A false positive can be generated by mistakenly calculating the top replication timing

New 12/15/2014 12:02 AM
Deleted: significance against

New 12/15/2014 12:02 AM
Deleted: in the bottom replication timing regions. Similarly, false negatives could be generated by using the inverse calculation

New 12/15/2014 12:02 AM
Deleted: 785

New 12/15/2014 12:02 AM
Deleted: 3988

New 12/15/2014 12:02 AM
Deleted: 3816

New 12/15/2014 12:02 AM
Deleted: recurrent

New 12/15/2014 12:02 AM
Deleted: (43).

New 12/15/2014 12:02 AM
Deleted: 44

New 12/15/2014 12:02 AM
Deleted: 45

New 12/15/2014 12:02 AM
Deleted: Simiar

New 12/15/2014 12:02 AM
Deleted: <=

New 12/15/2014 12:02 AM
Deleted: 46

New 12/15/2014 12:02 AM
Formatted: Not Highlight

New 12/15/2014 12:02 AM
Deleted: 2063

New 12/15/2014 12:02 AM
Deleted: 870

New 12/15/2014 12:02 AM
Deleted: (47).

Specifically, several studies have reported that disruption of CTCF binding sites through mutations or abnormal methylation sites is closely associated with cancer (55,56). Moreover, we found that the oncogene BCL3 has a noticeably higher significant percentage with respect to the average (7.721 times of the average, p-value for two-sided binomial test = 6.762×10^{-13}). Interestingly, BCL3 is a proto-oncogene candidate which is closely associated with progression of diverse solid tumors (57). For example, BCL3 is aberrantly up- and down-regulated in breast cancer and nasopharyngeal carcinomas respectively, and is also reported to be strongly associated with survival in colorectal cancer. However, it is not a highly mutation gene according to our data: BCL3's mutation rate is 1.22 mutations/Mbp while the gene average is 2.52 mutations/Mbp. Our analysis suggests another possibility that the misregulation of BCL3 is possibly due to binding site disruption instead of the changes in the protein itself. Further computational and experimental effort should be made to clarify the mechanism of BCL3 regulation in different cancer types.

Whole genome recurrent events evaluation

Despite great efforts to annotate noncoding regions, there are still many regions with as yet unknown regulatory roles. In order to evaluate the recurrent events in these regions, LARVA provides all possible p-values, whether before or after adjustment, and with or without replication timing corrections, for high confident bins on the genome (see methods for details) of variable length. We also compared the results from our beta-binomial model with the binomial models. For example, we randomly sampled 5000 10kb bins from the whole genome and made a Manhattan plot of p-values from both methods. It is obvious that the p-values from the binomial distribution were noticeably inflated (Fig. 7 B), while our beta-binomial model effectively controls the p-values (Fig. 7 A). The p-values obtained from using different bin sizes are provided in Table S4.

DISCUSSION

Due to the rapid decline in time and money involved to perform whole genome sequencing, data is now available for thousands of genomes where previously only a handful were available (58). However, the analyses necessary for finding useful patterns in this data, and making sense of it for clinical benefit, have not kept pace with this sudden increase. Therefore, it is important that new algorithms are developed that can efficiently mine relevant patterns from genome sequence data, and that user interfaces for finding and understanding that data are optimized so that clinicians and biologists, who may not have extensive technical expertise, can use these results effectively in their work.

Compared with the extensive computational and experimental efforts on the mutation patterns in the protein coding regions in the past decade (59), the noncoding regions, which was viewed as 'dark matter', and comprises up to 98% of the human genome, are barely investigated in cancer research studies, partially due to the limited knowledge of noncoding function. However, recently several examples clearly pinpointed the phenotypic effect of mutations in noncoding regulatory regions in a variety of cancer types. For instance, TERT promoter, a well-known example, has been associated

New 12/15/2014 12:02 AM

Deleted: 48,49

New 12/15/2014 12:02 AM

Deleted: 802

New 12/15/2014 12:02 AM

Deleted: 4.744e-

New 12/15/2014 12:02 AM

Formatted: Superscript

New 12/15/2014 12:02 AM

Deleted: (50).

New 12/15/2014 12:02 AM

Deleted: 18

New 12/15/2014 12:02 AM

Formatted: Indent: First line: 0 cm

New 12/15/2014 12:02 AM

Deleted: time

New 12/15/2014 12:02 AM

Deleted: 100kb

New 12/15/2014 12:02 AM

Deleted: 6

New 12/15/2014 12:02 AM

Deleted: 6

New 12/15/2014 12:02 AM

Deleted: (51).

New 12/15/2014 12:02 AM

Deleted: (52)

New 12/15/2014 12:02 AM

Deleted: our

with several cancer types (21-23). Fusions of the 5' UTR of TMPRSS2 with ETS genes frequently observed in prostate cancer, as well as mutations in certain miRNA binding sites (60), can influence the binding affinity at these sites, and thus affect androgen receptor regulation in prostate cancer. Hence, it is important to explore the mutation landscapes of such noncoding regions.

In this paper, we have introduced a new computational framework for exploring patterns of mutation across either somatic or rare germline variants, especially in the noncoding regions of human genomes. We took advantage of complete genome annotation effort of the ENCODE project (16) to extract the most extensive catalog of non-coding regulatory regions to date. We included the TF binding sites and DHS sites from all ENCODE experiments, promoters, UTRs, predicted enhancers, conserved and sensitive noncoding regions from our previous efforts (18). We then explored 760 cancer genomes on this comprehensive list of noncoding annotations to search for the highly mutated regulatory regions as potential noncoding driver candidates.

Moreover, consistent with the highly heterogeneous protein coding regions (11), we observed larger than expected mutation variation across cancer types, samples, and genomic regions (Fig. 3). Therefore, the recently proposed binomial models, which assume a constant mutation rate and independence of mutation events, might be inadequate for the observed data (Fig. 4, Fig. S3-S4). Instead, we set up a hierarchical model to handle mutation count overdispersion. First, we flexibly modeled the mutation rate in the regulatory elements as a two-parameter beta distribution, hence the corresponding mutation count could be conveniently described as a beta-binomial distribution. It provided significant improvement over the binomial model. In addition, we found that other genomic features, such as replication timing, would largely affect the background mutation rate (Fig. S6) and consequently generate both false positives and negatives. We corrected the replication timing effect by estimating the local mutation parameters in the beta-binomial distribution for better p-value assessment.

In the 760 cancer whole genomes in our analysis, we discovered 3776 noncoding regulatory regions that have significantly higher mutations than expected and provided the mutation enrichment significance of bins with variable length on the whole genome (Table 2). A list of known noncoding hypomutated regions, such as TERT and TP53 TSS, were also reported by our analysis, which convincingly proved the effectiveness of LARVA in discovering functionally relevant results. We also observed some relatively novel results such as PRRC2B TSS, CTCF and BCL3 binding sites. BCL3 is a known oncogene that is highly associated with several solid tumors (57,61), but this gene itself is not enriched in somatic mutations. Our results advocate an alternate possibility that its involvement in cancer cells is actually in the disruption of its binding sites, rather than the disabling of the protein itself.

In summary, LARVA is a powerful computational method to explore a broad range of genome annotations to uncover the ones that are mutated across many samples. It is worthwhile to point out that although we demonstrated the effectiveness of our method by analyzing somatic genomes, LARVA is also potentially suitable to discover noncoding regions under higher germline mutation

New 12/15/2014 12:02 AM

Deleted: 19-

New 12/15/2014 12:02 AM

Deleted: ,

New 12/15/2014 12:02 AM

Deleted: recurrent

New 12/15/2014 12:02 AM

Deleted: and

New 12/15/2014 12:02 AM

Deleted: 14

New 12/15/2014 12:02 AM

Deleted: (16). We then explored the 785

New 12/15/2014 12:02 AM

Deleted: for

New 12/15/2014 12:02 AM

Deleted: 9

New 12/15/2014 12:02 AM

Deleted: and annotation categories

New 12/15/2014 12:02 AM

Deleted: across all the genomes, are a demonstrably poor fit

New 12/15/2014 12:02 AM

Deleted: two-layer

New 12/15/2014 12:02 AM

Deleted: genome

New 12/15/2014 12:02 AM

Deleted: heterogeneity

New 12/15/2014 12:02 AM

Deleted: target region

New 12/15/2014 12:02 AM

Deleted: 785

New 12/15/2014 12:02 AM

Deleted: 3816

New 12/15/2014 12:02 AM

Deleted: know

New 12/15/2014 12:02 AM

Deleted: 50,53

New 12/15/2014 12:02 AM

Deleted: , making

burden. LARVA makes it possible to predict putative noncoding drivers of genetic disease, and prioritize these predicted drivers for more rigorous downstream analysis. This may lead to faster identification of important targets that may be used to suppress disease with therapies and drugs.

FUNDING

This work was supported by the National Institutes of Health [5R01CA152057-02]. Funding for open access charge: National Institutes of Health.

REFERENCES

TABLE AND FIGURES LEGENDS

Figure 1. (A) A pie chart representing the distribution of samples in our dataset of collected whole genome sequenced (WGS) cancers. (B) A flowchart of LARVA's procedure for identifying significant highly mutated noncoding elements. Cancer variants in VCF format are passed through quality control filters, and then intersected with our noncoding annotation corpus. After factoring in regional mutation rate corrections, a beta-binomial distribution is fitted to the observed data, which allows the identification of elements with a significant mutational burden.

Figure 2. Mutational heterogeneity between different types of cancer within several prominent classes of noncoding annotations. The percentage of mutations varies widely between noncoding element types, between cancer types, and between samples of the same cancer type.

Figure 3. (A) Between samples of the same cancer type, there is huge mutation rate heterogeneity. For most cancers, the mutation rate spans several orders of magnitude. (B) Variation in the mutation rate across chromosome 1 in lung cancer (top) and prostate cancer (bottom).

Figure 4. (A) The beta-binomial distribution (pink line) provides better fitting to the observed mutation counts at 10kb resolution (black line) of 760 cancer genomes, especially at the right tail as compared to the binomial distribution (turquoise line). (B) A comparison of the cumulative distribution function (CDF) of the binomial distribution and the beta-binomial distribution from part A. (C) Boxplots of the Kolmogorov-Smirnov (KS) statistics.

Figure 5. The 1 kb genome bins representing the top 10% and bottom 10% of the DNA replication timing were used to derive an observed distribution of mutation counts, demonstrating the influence of replication timing. The fitted binomial and beta-binomial distributions are plotted as bar plots (A). P-values at different mutation counts were given by the observed, beta-binomial, and binomial distribution.

Figure 6. (A) The number of significant p-values implied by beta-binomial distribution and binomial distribution (with and without DNA replication timing correction). (B) A sorted p-value plot of the top significant TSSes derived from the LARVA analysis.

New 12/15/2014 12:02 AM

Deleted: in

New 12/15/2014 12:02 AM

Deleted: ACKNOWLEDGEMENT . << Details of all funding sources for the work in question, and for the Open Access charge, should be given. An example is shown below ... [16]

New 12/15/2014 12:02 AM

Deleted: << Details of all funding sources for the work in question, and for the Open Access charge, should be given. An example is shown below ... [17]

New 12/15/2014 12:02 AM

Deleted: AA123456 to A.B., BB123456 to C.D.]; and the Alcohol & Education Research Council [abcde123456

New 12/15/2014 12:02 AM

Deleted: <<Figures and tables should be submitted Tagged Image File Format (.tif), Encapsulated PostScript (.eps), Joint Photographic Experts Group (.jpg), Graphics Interchange Format (.gif), Adobe Illustrator (.ai) (please save your files in Illustrator's EPS format), Portable Network Graphics (.png), Microsoft Word (. ... [18]

New 12/15/2014 12:02 AM

Deleted: recurrently

New 12/15/2014 12:02 AM

Deleted: (A) We fit the mutation c ... [19]

New 12/15/2014 12:02 AM

Deleted: The

New 12/15/2014 12:02 AM

Deleted: test was performed betw ... [20]

New 12/15/2014 12:02 AM

Deleted: 100

New 12/15/2014 12:02 AM

Deleted: distribution.

New 12/15/2014 12:02 AM

Deleted: empirical

New 12/15/2014 12:02 AM

Deleted: best fit

New 12/15/2014 12:02 AM

Deleted:) and as line plots (B) for ... [21]

New 12/15/2014 12:02 AM

Deleted: is a better model for this data

New 12/15/2014 12:02 AM

Deleted: the best fit

New 12/15/2014 12:02 AM

Deleted: ,

New 12/15/2014 12:02 AM

Deleted: the best fit

New 12/15/2014 12:02 AM

Deleted: graph

New 12/15/2014 12:02 AM

Deleted: noncoding annotations

<p>12. Rudd, M.L., Mohamed, H., Price, J.C., AJ, O.H., Le Gallo, M., Urick, M.E., Cruz, P., Zhang, S., Hansen, N.F., Godwin, A.K. <i>et al.</i> (2014) Mutational analysis of the tyrosine kinome in serous and clear cell endometrial cancer uncovers rare somatic mutations in TNK2 and DDR1. <i>BMC cancer</i>, 14, 884.</p>	<p>New 12/15/2014 12:02 AM Deleted: 10</p>
<p>13. Long, G.V., Fung, C., Menzies, A.M., Pupo, G.M., Carlino, M.S., Hyman, J., Shahheydari, H., Tembe, V., Thompson, J.F., Saw, R.P. <i>et al.</i> (2014) Increased MAPK reactivation in early resistance to dabrafenib/trametinib combination therapy of BRAF-mutant metastatic melanoma. <i>Nature communications</i>, 5, 5694.</p>	<p>New 12/15/2014 12:02 AM Deleted: 11</p>
<p>14. Yadav, M., Jhunjhunwala, S., Phung, Q.T., Lupardus, P., Tanguay, J., Bumbaca, S., Franci, C., Cheung, T.K., Fritsche, J., Weinschenk, T. <i>et al.</i> (2014) Predicting immunogenic tumour mutations by combining mass spectrometry and exome sequencing. <i>Nature</i>, 515, 572-576.</p>	<p>New 12/15/2014 12:02 AM Deleted: 12</p>
<p>15. Youn, A. and Simon, R. (2011) Identifying cancer driver genes in tumor genome sequencing studies. <i>Bioinformatics</i>, 27, 175-181.</p>	<p>New 12/15/2014 12:02 AM Deleted: 13</p>
<p>16. Consortium, E.P. (2012) An integrated encyclopedia of DNA elements in the human genome. <i>Nature</i>, 489, 57-74.</p>	<p>New 12/15/2014 12:02 AM Deleted: 14</p>
<p>17. Gerstein, M.B., Rozowsky, J., Yan, K.K., Wang, D., Cheng, C., Brown, J.B., Davis, C.A., Hillier, L., Sisu, C., Li, J.J. <i>et al.</i> (2014) Comparative analysis of the transcriptome across distant species. <i>Nature</i>, 512, 445-448.</p>	<p>New 12/15/2014 12:02 AM Deleted: 15</p>
<p>18. Fu, Y., Liu, Z., Lou, S., Bedford, J., Mu, X., Yip, K.Y., Khurana, E. and Gerstein, M. (2014) FunSeq2: A framework for prioritizing noncoding regulatory variants in cancer. <i>Genome biology</i>, 15, 480.</p>	<p>New 12/15/2014 12:02 AM Deleted: 16</p>
<p>19. Futreal, P.A., Coin, L., Marshall, M., Down, T., Hubbard, T., Wooster, R., Rahman, N. and Stratton, M.R. (2004) A census of human cancer genes. <i>Nature reviews. Cancer</i>, 4, 177-183.</p>	<p>New 12/15/2014 12:02 AM Deleted: 17</p>
<p>20. Dees, N.D., Zhang, Q., Kandoth, C., Wendl, M.C., Schierding, W., Koboldt, D.C., Mooney, T.B., Callaway, M.B., Dooling, D., Mardis, E.R. <i>et al.</i> (2012) MuSiC: identifying mutational significance in cancer genomes. <i>Genome research</i>, 22, 1589-1598.</p>	<p>New 12/15/2014 12:02 AM Deleted: 18</p>
<p>21. Vinagre, J., Almeida, A., Populo, H., Batista, R., Lyra, J., Pinto, V., Coelho, R., Celestino, R., Prazeres, H., Lima, L. <i>et al.</i> (2013) Frequency of TERT promoter mutations in human cancers. <i>Nature communications</i>, 4, 2185.</p>	<p>New 12/15/2014 12:02 AM Deleted: 19</p>
<p>22. Maurano, M.T., Humbert, R., Rynes, E., Thurman, R.E., Haugen, E., Wang, H., Reynolds, A.P., Sandstrom, R., Qu, H., Brody, J. <i>et al.</i> (2012) Systematic localization of common disease-associated variation in regulatory DNA. <i>Science</i>, 337, 1190-1195.</p>	<p>New 12/15/2014 12:02 AM Deleted: 20</p>
<p>23. Grossman, S.R., Andersen, K.G., Shlyakhter, I., Tabrizi, S., Winnicki, S., Yen, A., Park, D.J., Griesemer, D., Karlsson, E.K., Wong, S.H. <i>et al.</i> (2013) Identifying recent adaptations in large-scale genomic data. <i>Cell</i>, 152, 703-713.</p>	<p>New 12/15/2014 12:02 AM Deleted: 21</p>
<p>24. Weinhold, N., Jacobsen, A., Schultz, N., Sander, C. and Lee, W. (2014) Genome-wide analysis of noncoding regulatory mutations in cancer. <i>Nature genetics</i>, 46, 1160-1165.</p>	<p>New 12/15/2014 12:02 AM Deleted: 22</p>
<p>25. Erturk, E., Cecener, G., Polatkan, V., Gokgoz, S., Egeli, U., Tunca, B., Tezcan, G., Demirdogen, E., Ak, S. and Tasdelen, I. (2014) Evaluation of Genetic Variations in miRNA-Binding Sites of BRCA1 and BRCA2 Genes as Risk Factors for the Development of Early-Onset and/or Familial Breast Cancer. <i>Asian Pacific journal of cancer prevention : APJCP</i>, 15, 8319-8324.</p>	<p>New 12/15/2014 12:02 AM Deleted: 23</p>
<p>26. Medrzycki, M., Zhang, Y., Zhang, W., Cao, K., Pan, C., Lailier, N., McDonald, J.F., Bouhassira, E.E. and Fan, Y. (2014) Histone h1.3 suppresses h19 noncoding RNA expression and cell growth of ovarian cancer cells. <i>Cancer research</i>, 74, 6463-6473.</p>	<p>New 12/15/2014 12:02 AM Deleted: 24</p>
<p>27. Alexandrov, L.B., Nik-Zainal, S., Wedge, D.C., Aparicio, S.A., Behjati, S., Biankin, A.V., Bignell, G.R., Bolli, N., Borg, A., Borresen-Dale, A.L. <i>et al.</i> (2013) Signatures of mutational processes in human cancer. <i>Nature</i>, 500, 415-421.</p>	<p>New 12/15/2014 12:02 AM Deleted: 25</p>
<p>28. Cancer Genome Atlas Research, N. (2008) Comprehensive genomic characterization defines human glioblastoma genes and core pathways. <i>Nature</i>, 455, 1061-1068.</p>	<p>New 12/15/2014 12:02 AM Deleted: 26</p>

29.	Berger, M.F., Lawrence, M.S., Demichelis, F., Drier, Y., Cibulskis, K., Sivachenko, A.Y., Sboner, A., Esqueva, R., Pflueger, D., Sougnez, C. <i>et al.</i> (2011) The genomic complexity of primary human prostate cancer. <i>Nature</i> , 470 , 214-220.	New 12/15/2014 12:02 AM Deleted: 27
30.	Weischenfeldt, J., Simon, R., Feuerbach, L., Schlagen, K., Weichenhan, D., Minner, S., Wuttig, D., Warnatz, H.J., Stehr, H., Rausch, T. <i>et al.</i> (2013) Integrative genomic analyses reveal an androgen-driven somatic alteration landscape in early-onset prostate cancer. <i>Cancer cell</i> , 23 , 159-170.	New 12/15/2014 12:02 AM Deleted: 28
31.	Wang, K., Yuen, S.T., Xu, J., Lee, S.P., Yan, H.H., Shi, S.T., Siu, H.C., Deng, S., Chu, K.M., Law, S. <i>et al.</i> (2014) Whole-genome sequencing and comprehensive molecular profiling identify new driver mutations in gastric cancer. <i>Nature genetics</i> , 46 , 573-582.	New 12/15/2014 12:02 AM Deleted: 29
32.	Cancer Genome Atlas Research, N., Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw, K.R., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C. and Stuart, J.M. (2013) The Cancer Genome Atlas Pan-Cancer analysis project. <i>Nature genetics</i> , 45 , 1113-1120.	New 12/15/2014 12:02 AM Deleted: 30
33.	Derrien, T., Estelle, J., Marco Sola, S., Knowles, D.G., Raineri, E., Guigo, R. and Ribeca, P. (2012) Fast computation and applications of genome mappability. <i>PLoS one</i> , 7 , e30377.	New 12/15/2014 12:02 AM Deleted: 31
34.	Quinlan, A.R. and Hall, I.M. (2010) BEDTools: a flexible suite of utilities for comparing genomic features. <i>Bioinformatics</i> , 26 , 841-842.	New 12/15/2014 12:02 AM Deleted: 32
35.	Harrow, J., Frankish, A., Gonzalez, J.M., Tapanari, E., Diekhans, M., Kokocinski, F., Aken, B.L., Barrell, D., Zadissa, A., Searle, S. <i>et al.</i> (2012) GENCODE: the reference human genome annotation for The ENCODE Project. <i>Genome research</i> , 22 , 1760-1774.	New 12/15/2014 12:02 AM Deleted: 33
36.	Rozowsky, J., Euskirchen, G., Auerbach, R.K., Zhang, Z.D., Gibson, T., Bjornson, R., Carriero, N., Snyder, M. and Gerstein, M.B. (2009) PeakSeq enables systematic scoring of ChIP-seq experiments relative to controls. <i>Nature biotechnology</i> , 27 , 66-75.	New 12/15/2014 12:02 AM Deleted: 34
37.	Yip, K.Y., Cheng, C., Bhardwaj, N., Brown, J.B., Leng, J., Kundaje, A., Rozowsky, J., Birney, E., Bickel, P., Snyder, M. <i>et al.</i> (2012) Classification of human genomic regions based on experimentally determined binding sites of more than 100 transcription-related factors. <i>Genome biology</i> , 13 , R48.	New 12/15/2014 12:02 AM Deleted: 35
38.	Thurman, R.E., Rynes, E., Humbert, R., Vierstra, J., Maurano, M.T., Haugen, E., Sheffield, N.C., Stergachis, A.B., Wang, H., Vernot, B. <i>et al.</i> (2012) The accessible chromatin landscape of the human genome. <i>Nature</i> , 489 , 75-82.	New 12/15/2014 12:02 AM Deleted: 36
39.	Bejerano, G., Pheasant, M., Makunin, I., Stephen, S., Kent, W.J., Mattick, J.S. and Haussler, D. (2004) Ultraconserved elements in the human genome. <i>Science</i> , 304 , 1321-1325.	New 12/15/2014 12:02 AM Deleted: 37
40.	Khurana, E., Fu, Y., Colonna, V., Mu, X.J., Kang, H.M., Lappalainen, T., Sboner, A., Lochovsky, L., Chen, J., Harman, A. <i>et al.</i> (2013) Integrative annotation of variants from 1092 humans: application to cancer genomics. <i>Science</i> , 342 , 1235-1238.	New 12/15/2014 12:02 AM Deleted: 38
41.	Ding, L., Wendl, M.C., Koboldt, D.C. and Mardis, E.R. (2010) Analysis of next-generation genomic data in cancer: accomplishments and challenges. <i>Human molecular genetics</i> , 19 , R188-196.	New 12/15/2014 12:02 AM Deleted: 39
42.	Young-Xu, Y. and Chan, K.A. (2008) Pooling overdispersed binomial data to estimate event rate. <i>BMC medical research methodology</i>, 8, 58.	New 12/15/2014 12:02 AM Deleted: 40
43.	Kleinman, J.C. (1975) Proportions with extraneous variance: two dependent samples. <i>Biometrics</i>, 31, 737-743.	
44.	Song, L. and Crawford, G.E. (2010) DNase-seq: a high-resolution technique for mapping active gene regulatory elements across the genome from mammalian cells. <i>Cold Spring Harbor protocols</i> , 2010 , pdb prot5384.	
45.	Stamatoyannopoulos, J.A., Adzhubei, I., Thurman, R.E., Kryukov, G.V., Mirkin, S.M. and Sunyaev, S.R. (2009) Human mutation rate associated with DNA replication timing. <i>Nature genetics</i> , 41 , 393-395.	New 12/15/2014 12:02 AM Deleted: 41
46.	Hodgkinson, A. and Eyre-Walker, A. (2011) Variation in the mutation rate across mammalian genomes. <i>Nature reviews. Genetics</i> , 12 , 756-766.	New 12/15/2014 12:02 AM Deleted: 42

47.	Tian, D., Wang, Q., Zhang, P., Araki, H., Yang, S., Kreitman, M., Nagylaki, T., Hudson, R., Bergelson, J. and Chen, J.Q. (2008) Single-nucleotide mutation rate increases close to insertions/deletions in eukaryotes. <i>Nature</i>, 455, 105-108.	New 12/15/2014 12:02 AM Deleted: 43
48.	Hollister, J.D., Ross-Ibarra, J. and Gaut, B.S. (2010) Indel-associated mutation rate varies with mating system in flowering plants. <i>Molecular biology and evolution</i>, 27, 409-416.	
49.	McDonald, M.J., Wang, W.C., Huang, H.D. and Leu, J.Y. (2011) Clusters of nucleotide substitutions and insertion/deletion mutations are associated with repeat sequences. <i>PLoS biology</i>, 9, e1000622.	
50.	Benjamini, Y. and Hochberg, Y. (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. <i>Journal of the Royal Statistical Society. Series B (Methodological)</i> . 57 , 289-300.	
51.	Kwon, Y.J., Lee, S.J., Koh, J.S., Kim, S.H., Lee, H.W., Kang, M.C., Bae, J.B., Kim, Y.J. and Park, J.H. (2012) Genome-wide analysis of DNA methylation and the gene expression change in lung cancer. <i>Journal of thoracic oncology : official publication of the International Association for the Study of Lung Cancer</i>, 7, 20-33.	New 12/15/2014 12:02 AM Deleted: 44
52.	Isogai, E., Ohira, M., Ozaki, T., Oba, S., Nakamura, Y. and Nakagawara, A. (2011) Oncogenic LMO3 collaborates with HEN2 to enhance neuroblastoma cell growth through transactivation of Mash1. <i>PloS one</i>, 6, e19297.	New 12/15/2014 12:02 AM Deleted: 45
53.	Ashburner, M., Ball, C.A., Blake, J.A., Botstein, D., Butler, H., Cherry, J.M., Davis, A.P., Dolinski, K., Dwight, S.S., Eppig, J.T. et al. (2000) Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. <i>Nature genetics</i>, 25, 25-29.	New 12/15/2014 12:02 AM Deleted: 46
54.	Filippova, G.N. (2008) Genetics and epigenetics of the multifunctional protein CTCF. <i>Current topics in developmental biology</i>, 80, 337-360.	New 12/15/2014 12:02 AM Deleted: 47
55.	Ohlsson, R., Renkawitz, R. and Lobanenkov, V. (2001) CTCF is a uniquely versatile transcription regulator linked to epigenetics and disease. <i>Trends in genetics : TIG</i>, 17, 520-527.	New 12/15/2014 12:02 AM Deleted: 48
56.	Takai, D., Gonzales, F.A., Tsai, Y.C., Thayer, M.J. and Jones, P.A. (2001) Large scale mapping of methylcytosines in CTCF-binding sites in the human H19 promoter and aberrant hypomethylation in human bladder cancer. <i>Human molecular genetics</i>, 10, 2619-2626.	New 12/15/2014 12:02 AM Deleted: 49
57.	Maldonado, V. and Melendez-Zajgla, J. (2011) Role of Bcl-3 in solid tumors. <i>Molecular cancer</i>, 10, 152.	New 12/15/2014 12:02 AM Deleted: 50
58.	Shendure, J. and Ji, H. (2008) Next-generation DNA sequencing. <i>Nature biotechnology</i>, 26, 1135-1145.	New 12/15/2014 12:02 AM Deleted: 51
59.	Koch, L. (2014) Cancer genomics: Non-coding mutations in the driver seat. <i>Nature reviews. Genetics</i> , 15 , 574-575.	New 12/15/2014 12:02 AM Deleted: 52
60.	Lin, P.C., Chiu, Y.L., Banerjee, S., Park, K., Mosquera, J.M., Giannopoulou, E., Alves, P., Tewari, A.K., Gerstein, M.B., Beltran, H. et al. (2013) Epigenetic repression of miR-31 disrupts androgen receptor homeostasis and contributes to prostate cancer progression. <i>Cancer research</i>, 73, 1232-1244.	New 12/15/2014 12:02 AM Deleted: 53
61.	Kim, Y.M., Sharma, N. and Nyborg, J.K. (2008) The proto-oncogene Bcl3, induced by Tax, represses Tax-mediated transcription via p300 displacement from the human T-cell leukemia virus type 1 promoter. <i>Journal of virology</i> , 82 , 11939-11947.	