

Comprehensive genome- and transcriptome-wide analyses of mutations associated with microsatellite instability in Korean gastric cancers

Kwiyeom Yoon,^{1,8} Sunghoon Lee,^{2,3,8} Tae-Su Han,^{4,8} So Yeon Moon,¹ Sun Mi Yun,^{1,5} Seong-Ho Kong,^{4,6} Sungwoong Jho,² Jinny Choe,¹ Jieun Yu,⁴ Hyuk-Joon Lee,^{4,6} Ji Hyun Park,¹ Hak-Min Kim,² So Yeun Lee,¹ Jongsun Park,² Woo-Ho Kim,^{4,7} Jong Bhak,^{2,3} Han-Kwang Yang,^{4,6} and Seong-Jin Kim^{1,5,9}

¹CHA Cancer Institute, CHA University, Seoul 135-081, Korea; ²Personal Genomics Institute, Genome Research Foundation, Suwon 443-270, Korea; ³TheragenEx Bio Institute Inc., Suwon 443-270, Korea; ⁴Cancer Research Institute, Seoul National University, Seoul 110-799, Korea; ⁵Department of Biomedical Science, College of Life Science, CHA University, Gangnam-gu, Seoul 135-081, Korea; ⁶Department of Surgery, Seoul National University College of Medicine, Seoul 110-799, Korea; ⁷Department of Pathology, Seoul National University College of Medicine, Seoul 110-799, Korea

Microsatellite instability (MSI) is a critical mechanism that drives genetic aberrations in cancer. To identify the entire MS mutation, we performed the first comprehensive genome- and transcriptome-wide analyses of mutations associated with MSI in Korean gastric cancer cell lines and primary tissues. We identified 18,377 MS mutations of five or more repeat nucleotides in coding sequences and untranslated regions of genes, and discovered 139 individual genes whose expression was down-regulated in association with UTR MS mutation. In addition, we found that 90.5% of MS mutations with deletions in gene regions occurred in UTRs. This analysis emphasizes the genetic diversity of MSI-H gastric tumors and provides clues to the mechanistic basis of instability in microsatellite unstable gastric cancers.

[Supplemental material is available for this article.]

Human cancer arises as a consequence of the accumulation of multiple genetic and epigenetic changes in the genome (Grady and Carethers 2008; Imai and Yamamoto 2008; Toyota and Suzuki 2010; Stratton 2011). The major forms of genomic instability observed in gastric cancer are chromosomal instability (CIN) and microsatellite (MS) instability (MSI) (Ottini et al. 2006). Deficiencies in the DNA mismatch repair (MMR) system cause MSI, which is characterized by nucleotide length abnormalities occurring in tandem repeat units of 1–6 bp (MS) (Vilar and Gruber 2010; Pino and Chung 2011). MSI is a key factor in several cancers, including colorectal, endometrial, and gastric cancers. MS mutations found in these cancers are expected to contribute to MSI-H (high levels of MSI) tumorigenesis (Menoyo et al. 2001; Duval et al. 2002; Mori et al. 2002; Woerner et al. 2003; Karamurzin and Rutgers 2009; Shin et al. 2011).

Recent advances in human genome analysis make it possible to consider listing all of the genetic lesions in specific cancers. As a result, appreciable changes in basic and clinical medicine are expected to lead to the development of new therapeutic and diagnostic approaches. Recently, high-throughput sequencing has emerged as a reliable alternative to microarrays to study genome and transcriptome profiles. Many studies have characterized one or several MSI target genes in human gastric cancers; however, no genome- and transcriptome-wide analyses have been conducted to date (Menoyo et al. 2001; Duval et al. 2002; Mori et al. 2002).

Here, we analyzed all the insertion/deletions (indels) in repetitive sequence tracts and the comprehensive pattern of MSI occurrence throughout the whole genome. We found that MSI has length specificity, and the instability was driven mostly by deletion rather than insertion. To explore gastric cancer-specific MS mutations, we developed a systematic identification by comparing genome results with transcriptome data. We identified 18,377 MS mutations in gene regions, including repeat tracts that affected the post-transcriptional regulation of mutated genes. Surprisingly, 3482 out of 14,895 deletion mutations were found only in transcriptome sequencing of Korean MSI-H cell lines and tissues. A great number of mutations may contribute to gastric carcinogenesis by functional inactivation or dysfunction of target genes. These findings further advance our understanding of gastric tumorigenesis in MSI-H cancers.

Results

Genome- and transcriptome-wide analyses of MS mutations through a high-throughput sequencing approach

To characterize MS mutations at the genome- and transcriptome-wide levels, we analyzed 18 human gastric cancer cell lines, including three Korean gastric cancer cell lines (SNU-1, SNU-520, and SNU-638) known to have MSI-H (Myeroff et al. 1995; Shin and Park 2000), as well as 16 pairs of primary Korean gastric cancer tissues and their adjacent normal tissues (eight MSI-H and eight MS stable [MSS] cancers) (Supplemental Fig. S1). Because the presence of stromal cells in tumor samples may confound the identification of genes with MS mutations, we first performed high-throughput RNA sequencing of the 18 gastric cancer cell lines and whole-genome

⁸These authors contributed equally to this work.

⁹Corresponding author
E-mail kimsj@cha.ac.kr

Article published online before print. Article, supplemental material, and publication date are at <http://www.genome.org/cgi/doi/10.1101/gr.145706.112>.

sequencing of six gastric cancer cell lines (three known MSI-H and three MSS), followed by RNA sequencing of 16 pairs of primary gastric cancer tissues.

The etiology, location, and carcinogenesis processes of gastric cancer are considered to be quite different in Korea as compared to Western countries. Therefore, we sequenced MSI-H gastric cancer cell lines derived from Korean gastric cancer patients because we performed transcriptome analysis and validation with primary Korean gastric cancer tissues. We sequenced an average of 100 Gbp of whole-genome sequencing with a mean depth of $34.9\times$ and 98.9% of bases covered and an average of 4.7 Gbp of RNA sequencing with an average of 47 million mapped reads and 87.4% of mapping rates for each sample (Supplemental Table S1). We used Sanger sequencing to determine if our method for analysis of indels was appropriate, since it was complex and difficult to define the accuracy of indel calling. The results validated 100% specificity in 123 MS mutations. The sensitivities were different in MSS and MSI-H samples. Overall, the sensitivities of MSI-H samples were much higher (87%–92%) than those of MSS samples (64%–78%) (Supplemental Tables S1, S2; Supplemental Notes).

We first analyzed the entire set of indels within MS in gastric cancer cell lines and primary tissues. We characterized MS mutations containing a length of five or more repeat nucleotides in the genomes and transcriptomes because the number of MSI target genes with less than six repeat sequences was very rare but was still reported (Ionov et al. 2004; Royrvik et al. 2007). The majority of indels (98.9%) were detected in nongene regions, including intergenic and intron (intragenic) regions, and in noncoding regions expressing noncoding RNAs and small RNAs. We found similar total numbers of indels in the transcriptomes of MSS and MSI-H gastric cancer cell lines, with a larger number of single-nucleotide insertions identified in the untranslated regions (UTRs) of MSS samples than those of MSI-H samples (Supplemental Tables S2–S4). The mechanisms underlying the relative increase in single-nucleotide insertions in the transcriptomes of MSS samples warrant further investigation. However, the total number of deletions in the MS mutations identified in the genomes and transcriptomes of MSI-H gastric cancer cell lines and tissues was considerably higher than that in the MSS samples, suggesting that instability was driven mostly by deletion, rather than by insertion (Supplemental Tables S3, S4). Therefore, we focused only on deletions connected to MS mutations in our further analyses.

Comprehensive distribution of MS mutations with deletion through the genomes and the transcriptomes

We discovered that the number of deletions within MS in the SNU-1, SNU-520, and SNU-638 MSI-H gastric cancer cell lines was remarkably different from the SNU-16, SNU-668, and MKN-45 MSS gastric cancer cell lines, as shown in the Circos plots in Figure 1A. Because mutations in protein-coding genes are considered to be a causative factor in carcinogenesis (Ionov et al. 1993; Markowitz et al. 1995; Perucho 1996; Rampino et al. 1997; Alhopuro et al. 2010), we filtered MS mutations in gene regions (both coding sequences [CDSs] and UTRs) from the identified total deletions (Supplemental Table S2, S3). The number of deletion mutations in gene regions was four- to sixfold higher in MSI-H gastric cancer cell lines and primary tissues than MSS samples. We identified an average of 9554 deletions within MSs in the genomes of three MSI-H gastric cancer cell lines, whereas transcriptome analysis of the gastric cancer cell lines and primary tissues revealed averages of 3323 and 2786 MS mutations, respectively (Supplemental Table S3).

These mutations were compared with databases of genetic variations (dbSNP and 1000 Genomes) in order to highlight the novel deletion. Novel mutations found in the genome of three MSI-H gastric cancer cell lines, transcriptome of the gastric cancer cell lines, and transcriptome of primary tissues were ~41.7%, 40.2%, and 41.7% of total deletions in gene regions, respectively. In addition, we differentiated somatic deletion mutations from germline alterations in matched tumor and normal pairs of MSI-H gastric cancer primary tissues. Somatic or putative somatic mutations by the MMR deficiency consisted of 61% of total deletions in MSI-H gastric cancer tissues, whereas germline mutations that occurred in MSI-H or MSS tissues were 18.4% of MS mutations (Supplemental Table S2).

We computationally searched for all the repeat sequences of the human genome that contained ≥ 5 nucleotides (nt) in length in order to compare to our sequencing results. The search revealed the existence of 421,687 repeat sequences in gene regions of the human genome. However, our whole-genome sequencing showed only 14,895 deletion mutations in gene regions in MSI-H gastric cancer cell lines. Among these MS mutations, 3726 length alterations of genes were present in all three MSI-H gastric cancer cell lines but not in the three MSS gastric cancer cell lines (Fig. 1B). Interestingly, almost all of 421,687 repeat sequences were <10 nt in length, most of which were not identified as MS mutations by whole-genome sequencing (Fig. 1C). The percentage of mononucleotide repeats of length 10 or more in CDSs and UTRs is ~2.4% (10,292 of 421,687) and 4.6% (1,012,750 of 22,104,956) when intronic regions are included. The mutation ratios of the MSs in length of 10 and 11–21 mononucleotides were significantly higher, ~74% and 91%, respectively (Fig. 1C). However, compared with their longer repeats, only 1.3% of the repeat sequences ≤ 9 nt in length were mutated (Fig. 1C). The detection of mutations in repeat tracts of >21 mononucleotides in length is unreliable due to the limitations of current high-throughput sequencing technology, such as the generation of short read lengths. As consistent with the Figure 1C, the majority of deletions within MSs observed in all three MSI-H gastric cancer cell lines were 11–21 nt in length, suggesting that MSs of these lengths may be especially prone to genetic instability caused by MMR deficiency (Fig. 1D). However, the repeats of <10 nt in length were rarely subject to MSI in MSI-H gastric cancer cell lines. A similar pattern in the frequencies of MS mutations was also exhibited in whole-genome (including nongene region) analysis, which is available in Supplemental Figure S2. These results suggest that the majority of mononucleotide repeat tracts showing instability are a minimum of 10 nt in length.

Characterization of repeat sequence tract instabilities in gastric cancers

We explored the patterns of all deletion types, including A/T and G/C mononucleotide and CA/GT dinucleotide repeats that are prevalent in mammals. As shown in Figure 2, A through C, we did not detect substantial differences in the frequency of dinucleotide repeat tract deletions in gene regions between MSI-H and MSS gastric cancers, but mononucleotide repeat deletions were strikingly enriched in MSI-H gastric cancer cell lines and primary tissues (Supplemental Fig. S3). We further examined the number of MS mutations within mononucleotide repeats, in terms of deletion length, in MSI-H gastric cancers. Single-nucleotide deletions were most abundant. The number of deletions of >2 nt in length was considerably higher in MSI-H gastric cancers than in MSS gastric cancers (Fig. 2D–F; Supplemental Table S4).

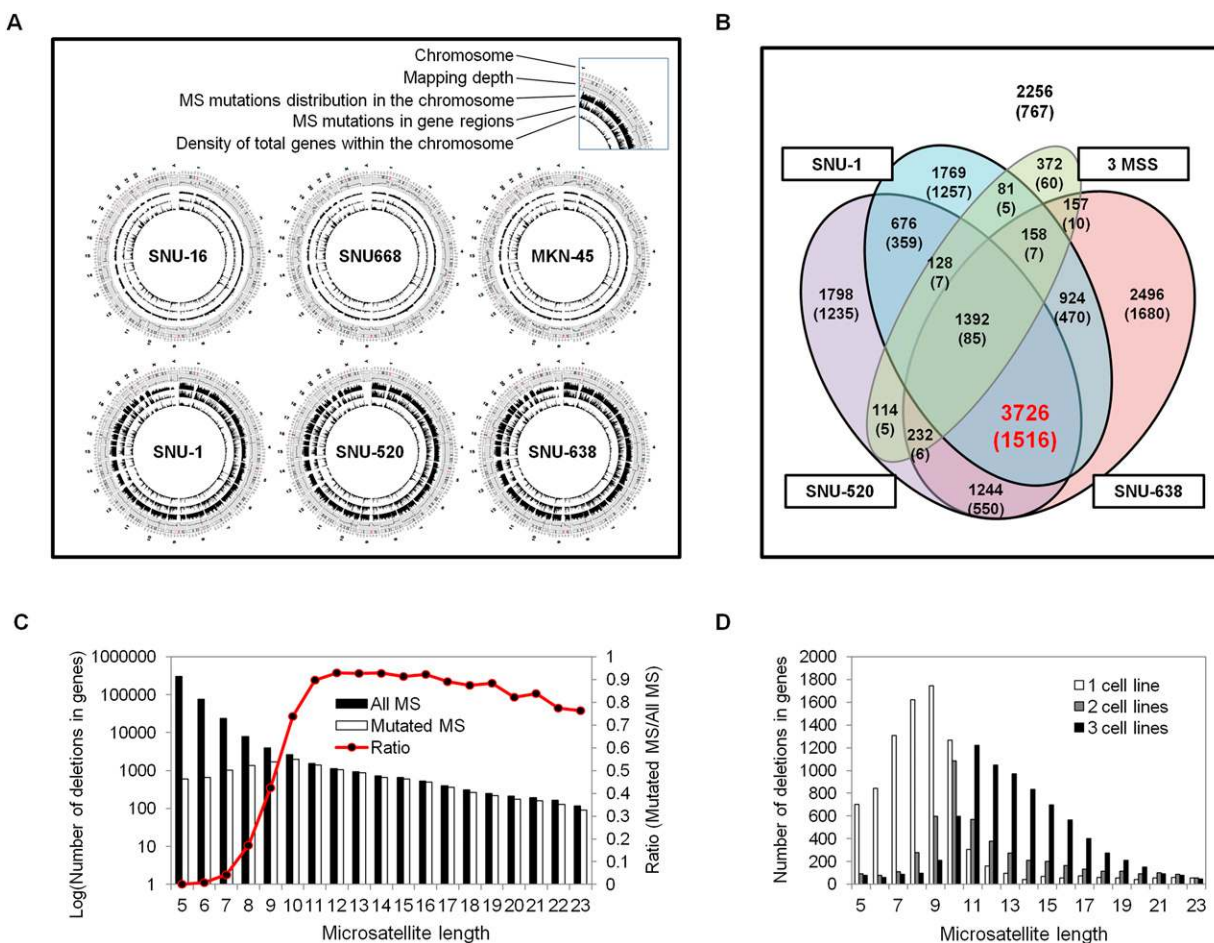


Figure 1. The landscape of MS mutations in human gastric cancer. (A) Graphical representation of six gastric cancer genomes in a Circos plot (Krzywinski et al. 2009). Remarkable differences in the number of MS mutations were identified in three MSI-H cell lines, with mapping depth and gene density displaying similar patterns. (B) Venn diagram depicting the dispersion of MS mutations in CDS and UTRs of genes between three MSI-H and sum of three MSS gastric cancer cell lines from whole-genome sequencing. The overlapping regions indicate the number of length alterations in gene regions that the cell lines have in common. The number of novel genes is displayed in parentheses. (C) Comparisons between mutated repeat tracts (mutated MS) and all existed repeat tracts (all MS) in gene regions in which both contain five or more mononucleotides in length (A/T and G/C) are displayed as ratios. (D) Length distribution of mononucleotide repeat instabilities in gene regions found in one, two, or all three MSI-H cell lines.

The number of mutations in G or C repeat tracts was higher in MSI-H gastric cancer cell lines and primary tissues than in MSS cell lines and primary tissues, but they were much less common than mutations involving A/T repeats (Fig. 2A–C). Importantly, we found 1103 different genes with deletion mutations in their CDS regions in MSI-H gastric cancer cell lines. Of these genes, 92 were detected in all three MSI-H gastric cancer cell lines (Fig. 3A–C). We extended our analysis from the CDS region to the UTRs, which contained longer MSs that are highly susceptible to mutation. Figure 3, A–C shows the numbers of deletion mutations in their UTRs. MSs with longer repeat lengths exhibited a higher instability rate, which led to multiple nucleotide deletions (Suraweera et al. 2001). Taken together, these results demonstrate that the number of candidate genes with MS mutations is markedly greater than had been previously known (Suraweera et al. 2001; Kim et al. 2002; Woerner et al. 2003; Royrvik et al. 2007; Shin et al. 2011).

Identification of frameshift mutations in gastric cancers

We identified 1103 different genes with deletion mutations in their CDS regions and used *SelTarbase*, a comprehensive mononucleotide

repeat mutation database of MSI-H tumors, to extract candidate genes that have not been identified in gastric cancers (Woerner et al. 2010). We found that 956 genes with MS mutations in their CDS regions have not been validated in gastric cancer tissues. Through pathway analysis and a driver gene score, we selected 24 candidate driver genes with nine or more repeat tracts for validation (Supplemental Notes). These genes are known to be involved in either chromosomal instability or cancer development. To study the mutation frequency of these genes in primary tumors, we carried out a mutational analysis of 17 uncharacterized genes associated with tumorigenesis and seven previously reported genes, including six genes that have not previously been verified in gastric cancer (Table 1; Supplemental Table S4; Kim et al. 2002; Tougeron et al. 2009; Williams et al. 2010; Woerner et al. 2010). We performed PCR and Sanger sequencing analysis with genomic DNAs of gastric cancer patient tissue to determine whether mutations in the 24 candidate genes were present in 22 MSI-H and 13 MSS primary Korean gastric tumors (Supplemental Fig. S4; Supplemental Table S6). As in colorectal cancers (Kim et al. 2002; Tougeron et al. 2009), mutations within mononucleotide tracts in *TGFBR2*, *MIS18BP1*, and *RNPC3* were detected in >60% of the MSI-H primary gastric

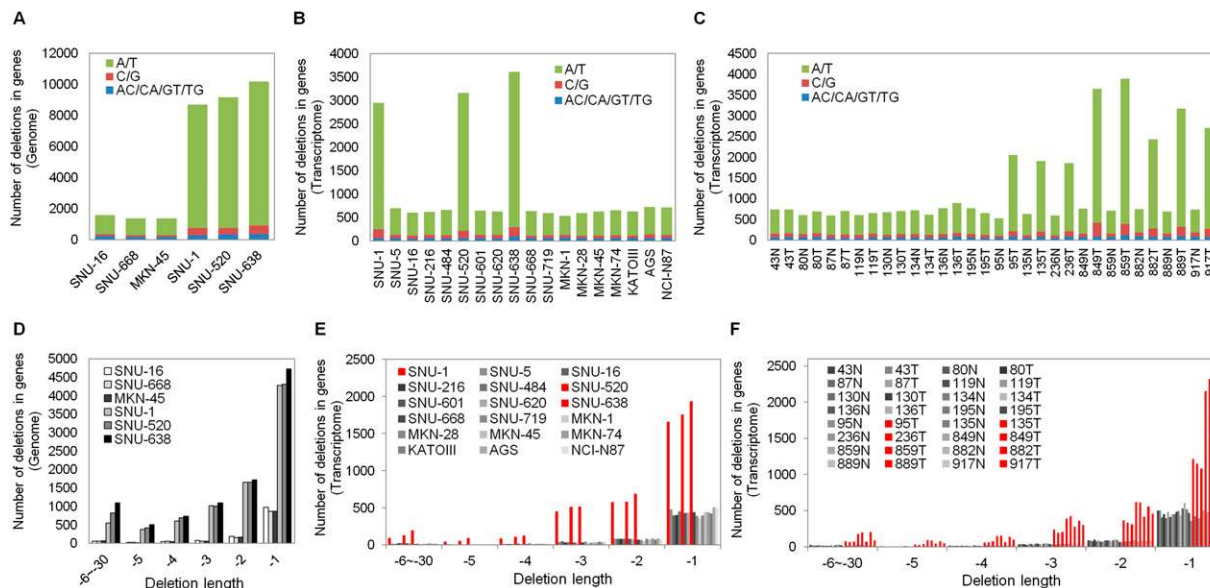


Figure 2. Refinement of repeat sequence tract instabilities within gene regions in gastric cancer. (A–C) Total Deletion counts in gene regions from whole-genome sequencing of six gastric cancer cell lines (A), RNA sequencing of 18 gastric cancer cell lines (B), and RNA sequencing of 16 pairs of gastric cancer and normal matched control tissues (C). (D–F) Deletion frequency of MSIs that contain mononucleotide repeats in human gastric cancer, identified by whole-genome sequencing (D) and RNA sequencing (E, F). Red bars indicate MSI-H cell lines (E) or tissues (F).

cancers. Interestingly, although *CEP164* mutations are relatively uncommon in colorectal cancers, 16 of the 22 (72.7%) MSI-H gastric cancers were demonstrated to contain mutations within the *CEP164* polyadenine repeats. Among the novel candidate driver genes, *KIAA2018* (77.3%), *CNOT1* (77.3%), and *CCDC150* (63.6%) showed relatively high mutational frequencies in MSI-H gastric cancer tissues (Table 1; Supplemental Table S6).

Length alterations in MS tracts in the coding regions of genes are known to affect gene expression because they can generate premature termination codons that produce C-terminally truncated proteins or transcripts that are actively degraded by nonsense-mediated mRNA decay (Rebbapragada and Lykke-Andersen 2009; Silva and Romao 2009). We extracted a list of genes with frameshift mutations that were detected by the whole-genome analysis but not by transcriptome analysis (i.e., because their expression levels were low or undetectable). Differential expression analysis of transcriptome was performed systematically by using the variance analysis algorithm DESeq with the *P*-value of 0.05 as a threshold (Anders and Huber 2010). The expressions of 23 genes with frameshift mutations, including *TGFBR2*, were low to undetectable in transcriptomes (Fig. 4A; Supplemental Table S7). *TGFBR2* contained A repetitive sequences in its coding region, and its expression was not detected in MSI-H samples. Notably, ~25% of the genes identified solely from whole-genome sequencing were not expressed in any of the MSS or MSI-H gastric cancer cell lines examined, suggesting that these genes may not be expressed in gastric tissues or that their expression may be epigenetically suppressed in gastric cancers.

Post-transcriptional dysregulation of genes by instabilities of MSs in UTRs

We found that 90.5% of MS mutations with deletions in gene regions occurred in UTRs (Fig. 3A–C). UTRs are emerging as essential factors in mRNA stability, translation, and localization

through *cis*-elements, which include iron-responsive elements (IREs) and AU-rich elements (AREs) and often their secondary structures, or *trans*-elements, which include noncoding RNAs, miRNAs, and RNA-binding proteins. Such elements are critical in influencing the overall translation rate and stability of transcripts (Pickering and Willis 2005; Lopez de Silanes et al. 2007; Chatterjee and Pal 2009). Because recent studies have shown that *CEACAM1*, *RB1CC1*, and *EGFR* 3' UTR MSs control the post-transcriptional activity of their mRNA (Ruggiero et al. 2003; Paun et al. 2009; Yuan et al. 2009), we also examined whether the mutations in UTRs influenced expression levels by comparing the mutation status and transcriptional levels of specific genes.

We first identified genes that were significantly up-regulated (*P* < 0.05) in three MSI-H gastric cancer cell lines in comparison with 15 MSS gastric cancer cell lines (Supplemental Table S7). We found 137 genes with 210 MS mutations showed up-regulated mRNA expression; 96.2% of MS mutations in total up-regulated candidate genes were found in UTRs. In addition, we also investigated 139 genes with 199 MS mutations in their UTRs that were down-regulated due to the destabilization of mRNA caused by mutations in MSs. Surprisingly, the expression levels of the 139 genes with UTR MSIs were extremely low compared with genes without UTR mutations. Among these candidate genes, the expression levels of 48 genes (34.5%) were low to undetectable in all three MSI-H gastric cancer cell lines (Fig. 4A; Supplemental Table S7). We assessed five genes whose expressions were more significantly down-regulated in all three MSI-H against MSS cell lines and confirmed the expression levels of these genes by Q-PCR (Supplemental Table S8). The Q-PCR results for five genes—*MGLL*, *SORL1*, *C20orf194*, *WWC3*, and *PXDC1*—correlated to FPKM (fragments per kilobase of exon model per million mapped reads) values estimated by RNA sequencing (Fig. 4B; Supplemental Fig. S5). We next investigated whether MSI within the 3' UTR leads to post-transcriptional dysregulation of the transcripts. We examined the 3' UTR of *MGLL* because its expression level was most significantly

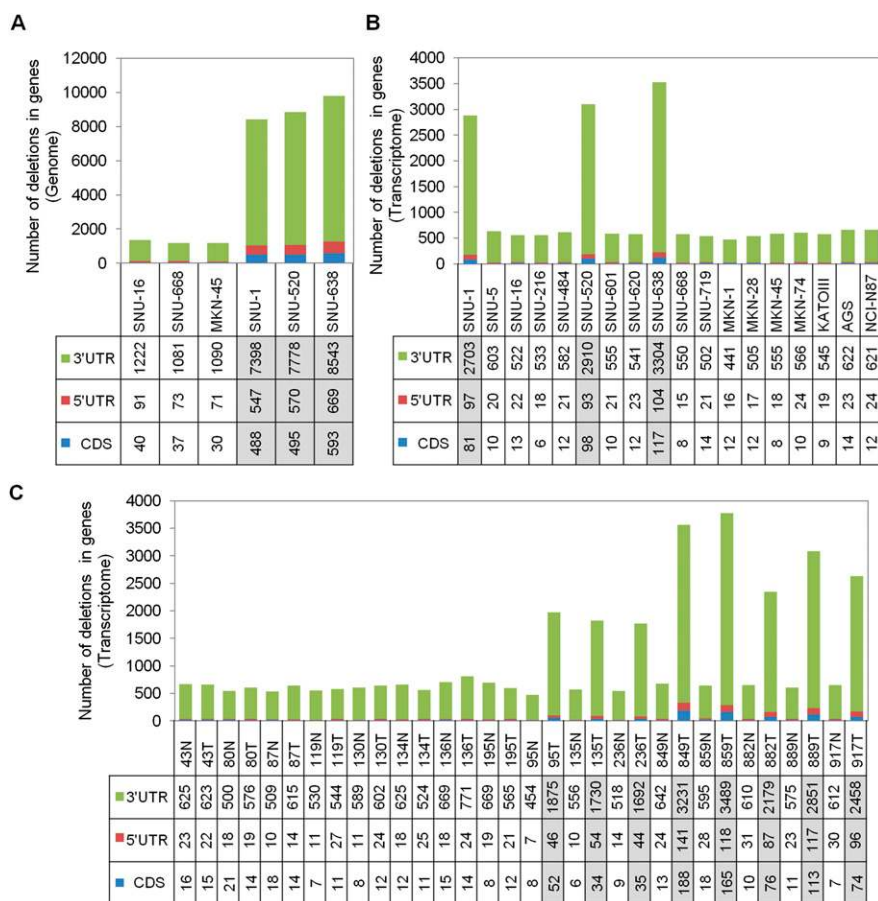


Figure 3. Analysis of mononucleotide MS mutations associated-genes in gastric cancer. (A–C) Deletions of specific genes carrying mononucleotide MSs were analyzed by whole-genome sequencing (A), and RNA sequencing in gastric cancer cell lines (B) and primary tissues (C).

down-regulated ($P < 3.58 \times 10^{-8}$) among the 139 genes. *MGLL*, also known as *MAGL*, serves as a critical regulator in endocannabinoid signaling, which is responsible for pain control (Jhaveri et al. 2007; Chanda et al. 2010; Lichtman et al. 2010). We examined MSI frequencies in *MGLL* by PCR and Sanger sequencing analysis in MSI-H and MSS primary gastric tissues (Supplemental Fig. S4C). Notably, mutations in the 3' UTR of *MGLL* were absent in all but one of the MSS samples (only 1 nt deletion), whereas deletions of >4 nt in the *MGLL* 3' UTR were found in 50% of MSI-H tumors (Table 1; Supplemental Table S6). To assess the functional impact of 3' UTR MSI on mRNA stability, we conducted a reporter assay using the luciferase gene fused to either the wild-type *MGLL* 3' UTR or the mutant *MGLL* 3' UTRs that were identified in MSI-H gastric cancer cell lines. The level of luciferase gene expression was markedly decreased (up to 42.6%) compared with wild type when the luciferase gene was fused to the mutant 3' UTR of *MGLL* (Fig. 4C). These results suggest that MSI within the 3' UTR may contribute to aberrant gene expression in MSI-H tumors.

Discussion

We have examined, for the first time, the entire MS mutations in MS unstable human gastric cancers through systematic whole-genome and whole-transcriptome approaches. Through comparison of whole-genome and whole-transcriptome data, we identified

frequently mutated genes and aberrantly expressed genes in Korean gastric cancers. We detected only 14,895 deletion mutations in genes, although the human genome has 421,687 repeat sequences in gene regions. Thus, our findings suggest that comprehensive MS mutation analysis using next-generation sequencing technology is a valid strategy to identify genes mutated in human gastric cancers. Because gastric cancer shows distinctive patterns related to ethnics, further studies pertaining to different racial/ethnic groups or cancer types may reinforce our investigations. Our study shows that mutations at MSs may not occur randomly while previously unrecognized features of the sequences, genomic structure, or surrounding nucleotide sequences may contribute to mutability. It is important to note that 27.2% (1392/5118) of MS mutations identified in all three MSI-H cell lines were also detected in MSS cell lines (Fig. 1B). Thus, these target sequences are particularly prone to mutation during the formation of gastric cancer, suggesting that mutations at mononucleotide repeats occur not only by virtue of MMR deficiency but also by the presently unknown mechanisms in gastric cancers.

Previous studies demonstrated that the frequency of mutation within MS tracts is dependent on the number of the repeats in the tract (Chen et al. 1995; Woerner et al. 2003; Paun et al. 2009). We showed that the mutation frequency is very high when repeat lengths are 10 or more units (Fig. 1C). These genes may present high mutational frequencies in a majority of MSI-H gastric cancers and may influence gastric cancer development. Although several genes with less than 10 repeat tracts were previously identified as MSI target genes, we found that repeat tracts with fewer than 10 mononucleotides are less susceptible to mutations in MSI-H tumors, which are likely to occur randomly (Mori et al. 2002; Royrvik et al. 2007). This indicates large mutation variability in short repeat tracts under MMR deficiency and emphasizes the genetic diversity of MSI-H gastric tumors.

In this study, the majority of mutations (99.9%) were detected within the UTRs, as well as the intronic and intergenic regions in MSI-H tumors. Evaluating the functional significance of mutations within these regions is a somewhat difficult task. Even though we have not validated all mutations in 3' UTRs, our results suggest that the 3' UTRs have relatively longer MSs, which are expected to be more unstable in MSI-H tumors. Recent exome sequencing study of gastric cancers reported that *ARID1A* showed mutations related to MSI status (Wang et al. 2011; Zang et al. 2012). We also found *ARID1A* deletion mutations, which were located in CDS and in two long A repeat tracts in 3' UTR, in the three MSI-H but not in the MSS gastric cancer cell lines (Supplemental Tables S2). We discovered 139 candidate genes whose expression was down-regulated in association with the UTR MSI-related alteration. Previous studies suggested that mutations in the 3' UTRs may have

Table 1. Frequencies of genes with MS mutations in primary human gastric cancers

Gene symbol	TGFB2	CEP164	MIST1BP1	RNPC3	SREK1P1	TMBIM4	PRRC2C	KIAA2018	CNOT1	CCDC150	RNF145	CCDC178	VCP	TVP23A	ULK4	PDSSB	UPF3A	PRRT2	RBM43	CD3G	INO80E	IPH4	SLAMF1	GIN1	MGLL
Repeat Location	A, 10	A, 11	A, 11	A, 12	A, 10	T, 10	A, 10	A, 11	T, 13	A, 11	A, 11	A, 10	A, 9	T, 9	A, 10	A, 9	A, 9	C, 9	A, 10	A, 9	A, 9	G, 9	A, 9	A, 9	A, 21
Frequency of mutations	86.4	72.7	63.6	68.2	40.9	31.8	9.1	77.3	77.3	63.6	59.1	59.1	50.0	50.0	40.9	40.9	40.9	27.3	27.3	22.7	22.7	13.6	13.6	13.6	50.0
% of MSI tumors	(19/22)	(16/22)	(14/22)	(15/22)	(9/22)	(7/22)	(2/22)	(17/22)	(17/22)	(14/22)	(13/22)	(13/22)	(11/22)	(11/22)	(9/22)	(9/22)	(9/22)	(6/22)	(6/22)	(5/22)	(5/22)	(3/22)	(3/22)	(3/22)	(11/22)
% of MSS tumors	0	0	0	0	0	0	0	7.7	7.7	0	0	0	0	0	0	0	0	0	0	0	0	7.7	0	0	7.7
Driver gene score	3.19	1.74	0.48	1.49	1.52	1.31	0.48	1.37	0.52	1.39	—	0.96	2.28	2.14	0.79	1.19	2.57	2.07	1.02	2.47	1.58	2.95	2.86	—	—

Frameshift mutations in *CEP164*, *MIST1BP1*, *RNPC3*, *SREK1P1*, *TMBIM4*, and *PRRC2C*: reported in colorectal cancer but have not been validated in gastric cancer. Frameshift mutations in *KIAA2018*, *CNOT1*, *CCDC150*, *RNF145*, *CCDC178*, *VCP*, *TVP23A*, *ULK4*, *PDSSB*, *UPF3A*, *PRRT2*, *RBM43*, *CD3G*, *INO80E*, *IPH4*, *SLAMF1*, and *GIN1*: 17 novel candidate genes identified in this study. Mutation in the *MGLL* 3' UTR MS: identified in this study. *RNF145* and *PDSSB* have frameshift mutations in their second A11 and second A9 tracts, respectively.

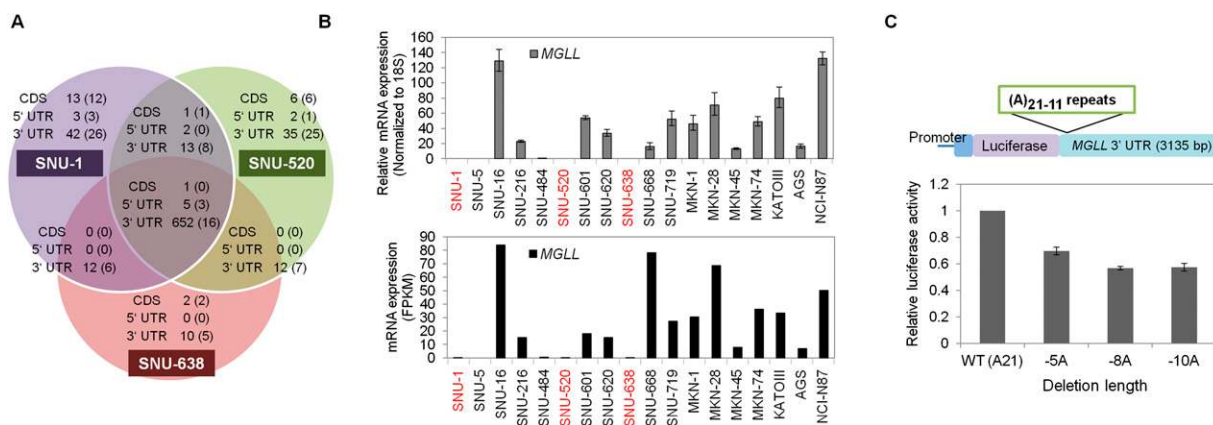


Figure 4. mRNA expression is dysregulated by the deletion of MSs in UTRs. (A) The number of mononucleotide repeat tracts in genes that correlate with down-regulated mRNA expression were identified by comparing the results of whole-genome sequencing with those of RNA sequencing in three MSI-H cell lines. (B) mRNA transcript levels of *MGLL* with MSs within their 3' UTR in MSI-H cell lines were measured by Q-PCR. The levels of mRNA expression quantified by RNA sequencing are shown. (C) A schematic of the luciferase gene construct is displayed above. Significant deletion of polyA in the 3' UTR causing decreased mRNA stability was confirmed via a reporter assay in MKN-1 cell line.

little impact on gene function (Zhang et al. 2001; Hienonen et al. 2005); however, our result speculated that mutations in the 3' UTRs influence gene expression in MSI-H tumors. Several *cis*- or *trans*-elements of the UTRs might fail to regulate their gene functions such as stability or activity if the mutations alter RNA sequences or structure. Thus, our findings suggest that aberrant expression of genes may create a growth or survival advantage for MSI-H gastric cancer.

Colorectal cancer patients with tumors exhibiting MSS or low levels of MSI (MSI-L) are reported to have a favorable response to fluorouracil-based adjuvant chemotherapy; however, such chemotherapy did not benefit patients with MSI-H tumors and may, in fact, have led to worse outcomes among such patients (Ribic et al. 2003). Even though the mutational spectra for genes with high mutational frequencies were known to be quite different between MSI-H gastric cancers and colorectal cancers, our analysis of mutations associated with MSI in human gastric cancers will provide further information about discrete molecular pathways, which may explain the difference.

In summary, we analyzed all the mutations in repetitive sequence tracts in MSI-H Korean gastric cancers. The number of candidate genes with deletion mutations is significantly greater than previously reported, implying that the accumulation of MSs contributes to the genetic complexities of gastric cancer. This study will enhance our understanding of gastric tumorigenesis in MSI-H cancers.

Methods

Preparation of cell lines and primary gastric cancer tissues

Total 18 gastric cancer cell lines, including 15 MSS cell lines (SNU-5, SNU-16, SNU-216, SNU-484, SNU-601, SNU-620, SNU-668, SNU-719, MKN-1, MKN-28, MKN-45, MKN-74, KATOIII, AGS, and NCI-N87) and three MSI-H cell lines (SNU-1, SNU-520, and SNU-638), were purchased from the Korean Cell Line Bank. The cells were maintained under standard conditions (RPMI-1640 containing 25 mM HEPES, 10% fetal bovine serum, 100 unit/mL streptomycin, and 100 units/mL penicillin at 37°C, 5% CO₂).

Korean gastric cancer tissues were from the gastric cancer tissue depository of the Gastrointestinal Division in the Department of

Surgery at Seoul National University Hospital. The gastric cancer tissues and normal mucosa controls were obtained from the removed stomach specimens immediately after gastrectomy in an aseptic condition. Sixteen pairs of gastric cancer and normal matched control tissues for RNA sequencing were randomly selected from the tissue depository samples. Twenty-two MSI-H and 13 MSS tissue samples for validation were selected from MSI information. The management of the tissue depository and the usage of the tissues and associated clinicopathologic data were approved by the institutional review board of the Seoul National University Hospital (IRB no. H-0806-072-248). The clinical information associated to the gastric tumors is provided in Supplemental Table S9.

RNA sequencing and whole-genome sequencing

TRIzol Reagent (Invitrogen) was used to isolate total RNA for RNA sequencing following the manufacturer's instructions. The total RNA was treated with DNase I and then was purified with an miRNeasy Mini Kit according to the manufacturer's instructions (Qiagen). The quality of the RNA was checked with the Agilent 2100 Bioanalyzer (Agilent) prior to sequencing. Genomic DNAs for whole-genome sequencing were prepared with DNeasy blood and tissue kit (Qiagen) according to their specific protocol.

We used the Illumina platform for analyzing indels of gastric cancer genomes and transcriptomes with a 90-bp paired-end library according to the manufacturer's instructions (Illumina). Libraries were constructed following the Illumina Paired-End Sequencing Library Preparation Protocol. Library quality and concentration were determined using an Agilent 2100 BioAnalyzer (Agilent). Each sample was paired-end sequenced with the Illumina Genome Analyzer II or with the Illumina HiSeq 2000 using HiSeq Sequencing kits. A base-calling pipeline (Sequencing Control Software [SCS], Illumina) was used to process the raw fluorescent images and the called sequences.

Reads alignment and small indel detection from whole-genome sequencing data

Paired-end sequence reads were aligned to hg19 human reference genome (NCBI build 37) with a BWA algorithm (Li and Durbin 2010) version 0.5.9. We permitted two mismatches within 45-bp seed sequence when aligning sequence reads. Aligned read files were

converted to SAM and BAM files using Samtools (Li et al. 2009) version 0.1.17. To remove PCR duplicates of sequence reads, which could be generated during the library construction process, we used the `rmDup` command from Samtools. To enhance read alignment accuracy, aligned reads were realigned at putative indel positions with the GATK (McKenna et al. 2010) IndelRealigner algorithm. Base quality scores were recalibrated using the GATK TableRecalibration algorithm.

Putative small indels were called and filtered using the GATK UnifiedGenotyper algorithm. The options used for small indel calling were a minimum of five to a maximum of 200 read mapping depth with a consensus quality of 30 and the prior likelihood for heterozygosity value of 0.001.

Reads alignment and small indel detection from RNA sequencing data.

For comprehensive detection of small indels, two independent analysis procedures were used, and the resulting small indels were merged as shown in Supplemental Figure S1.

1. All 90-bp paired-end sequence reads were aligned to hg19 human reference genome (NCBI build 37) with a TopHat algorithm (Trapnell et al. 2009) version 2.0.4. by using splice-junction information of GRCh37.58 (<http://asia.ensembl.org/info/data/ftp/index.html>). The resulting two outputs, deletions.bed and insertions.bed, describing small indels, were filtered by two thresholds, minimum mapping depth of five and supporting read ratio larger than 0.2, which were calculated by PILEUP files generated by the `mpileup` algorithm of Samtools.
2. All 90-bp paired-end sequence reads were aligned to hg19 reference genome with BWA, processed with Samtools and GATK as the same procedure to whole-genome sequence analysis except removing PCR duplicates. The minimum mapping depth is five, and the maximum mapping depth threshold was not applied to include small indels in genes that were expressed highly. Indels within range of 30 bp downstream from and 30 bp upstream of splicing sites were filtered out to reduce false positives.

Annotation of MSs and known small indel information

MSs, which are >4 bp in length with one or two base repeating units, in hg19 reference genome were mapped on the GRCh37.58 gene table by genomic features such as CDS, 5' UTR, and 3' UTR. The positions of MSs were compared with NCBI dbSNP 135 (<http://www.ncbi.nlm.nih.gov/projects/SNP/>) and 1000 Genomes (<http://www.1000genomes.org/>) to annotate known indel information. *SeITarbase* (Woerner et al. 2010) was used to annotate the observed frequency of indel occurrences of MSs in four types of cancers, such as colon, colorectal, gastric, and endometrial cancer. When the end of an indel overlapped with an MS, it is annotated as an MS indel. When more than two indels were mapped in one MS, these indels were annotated as the same MS indel. Small Indels found in our analysis were also mapped to MS and shared the annotation information.

Statistical analysis of gene expressions using RNA sequencing data

The mapping results of RNA-seq reads by the TopHat algorithm to hg19 reference were used to count the number of reads per gene with an `htseq-count` algorithm (<http://www-huber.embl.de/users/anders/HTSeq/doc/count.html>). Then, DESeq algorithm (Anders and Huber 2010) version 1.10.1 was used to analyze differentially expressed genes between MSI-H and MSS samples with the *P*-value of 0.05 as a threshold.

Validation of candidate genes with MS mutations

DNA was extracted by QIAmp DNA mini kit (Qiagen) according to the manufacturer's instructions for isolation of genomic DNA from human gastric tissues. Buffer ATL and proteinase K were mixed with 25 mg of tissue and 4 h incubation at 56°C in an upside-down position. Selected genes with MS mutations were amplified and sequenced in cell lines as well as MSI-H and MSS tissues genomic DNA. mRNA expressions for comparing to their transcript levels in FPKM to the transcriptomes were verified by Q-PCR with the cDNA of 18 gastric cancer cell lines.

Plasmid construction and Luciferase assay

To elucidate regulation of mRNA expression by UTRs, the corresponding *MGLL* 3' UTR DNA fragments from SNU-1 and *MKN-74* cDNA were cloned into pGL3-promoter vector (Promega) after luciferase gene position without a polyA signal of the vector. *MKN-1* cells (1.2×10^5 cells/well) were plated on a 24-well plate 1 d prior to transfection. Lipofectamine 2000 (Invitrogen) was used as the transfection reagent following the manufacturer's instructions. After 36 h, the whole-cell extracts were then prepared for the luciferase assay. The luciferase activity was measured by using the Luciferase Assay System (Promega) and normalized relative to β -galactosidase activity (Sigma) assessed by ELISA. The data were obtained from three independent experiments were performed each time in triplicate and presented as the fold increase in luciferase activities (mean \pm SD) relative to the control.

Data access

Sequencing reads are available in the NCBI Sequence Read Archive (SRA) (<http://www.ncbi.nlm.nih.gov/sra>) under accession number SRP014574.

Competing interest statement

The authors declare competing financial interests. S.L., J.B., and S-J.K. have personal financial interests as shareholders in TheragenEtex.

Acknowledgments

We thank Drs. W. Grady, P. Tan, and C. Lee for reviewing the manuscript and providing comments; Dr. H. Lee for the pathway analysis; and S.Y. Hwang for editing the manuscript. This work was supported in part by the National Research Foundation of Korea (NRF) grants (2009-0081756 and 2012M3A9C4048736) and by the Korea Health Industry Development Institute (KHIDI; A090726). J.B. and S.L. were supported by TheragenEtex and Genome Research Foundation internal funds.

Author contributions: K.Y. performed research, analyzed data, and wrote the manuscript; S.L. analyzed and interpreted sequencing data and drafted the manuscript; S.Y.M, S.M.Y., J.C., J.H.P., and S.Y.L. performed research, collected data, and contributed to the manuscript; T.S.H., S-H.K., J.Y., H-J.L., W-H.K., and H-K.Y. collected clinical information and prepared clinical samples; S.J., H-M.K., J.P., and J.B. analyzed and collected sequencing data; and S-J.K. designed research, interpreted data, and drafted the manuscript.

References

- Alhopuro P, Bjorklund M, Sammalkorpi H, Turunen M, Tuupainen S, Bistrom M, Niittymaki I, Lehtonen HJ, Kivioja T, Launonen V, et al. 2010. Mutations in the circadian gene *CLOCK* in colorectal cancer. *Mol Cancer Res* **8**: 952–960.

- Anders S, Huber W. 2010. Differential expression analysis for sequence count data. *Genome Biol* **11**: R106.
- Chanda PK, Gao Y, Mark L, Btesh J, Strassle BW, Lu P, Piesla MJ, Zhang MY, Bingham B, Uveges A, et al. 2010. Monoacylglycerol lipase activity is a critical modulator of the tone and integrity of the endocannabinoid system. *Mol Pharmacol* **78**: 996–1003.
- Chatterjee S, Pal JK. 2009. Role of 5'- and 3'-untranslated regions of mRNAs in human diseases. *Biol Cell* **101**: 251–262.
- Chen J, Heerdt BG, Augenlicht LH. 1995. Presence and instability of repetitive elements in sequences the altered expression of which characterizes risk for colonic cancer. *Cancer Res* **55**: 174–180.
- Duval A, Reperant M, Compoin A, Seruca R, Ranzani GN, Iacopetta B, Hamelin R. 2002. Target gene mutation profile differs between gastrointestinal and endometrial tumors with mismatch repair deficiency. *Cancer Res* **62**: 1609–1612.
- Grady WM, Carethers JM. 2008. Genomic and epigenetic instability in colorectal cancer pathogenesis. *Gastroenterology* **135**: 1079–1099.
- Hienonen T, Sammalkorpi H, Enholm S, Alhopuro P, Barber TD, Lehtonen R, Nupponen NN, Lehtonen H, Salovaara R, Mecklin JP, et al. 2005. Mutations in two short noncoding mononucleotide repeats in most microsatellite-unstable colorectal cancers. *Cancer Res* **65**: 4607–4613.
- Imai K, Yamamoto H. 2008. Carcinogenesis and microsatellite instability: The interrelationship between genetics and epigenetics. *Carcinogenesis* **29**: 673–680.
- Ionov Y, Peinado MA, Malkhosyan S, Shibata D, Perucho M. 1993. Ubiquitous somatic mutations in simple repeated sequences reveal a new mechanism for colonic carcinogenesis. *Nature* **363**: 558–561.
- Ionov Y, Matsui S, Cowell JK. 2004. A role for p300/CREB binding protein genes in promoting cancer progression in colon cancer cell lines with microsatellite instability. *Proc Natl Acad Sci* **101**: 1273–1278.
- Jhaveri MD, Richardson D, Chapman V. 2007. Endocannabinoid metabolism and uptake: Novel targets for neuropathic and inflammatory pain. *Br J Pharmacol* **152**: 624–632.
- Karamurzin Y, Rutgers JK. 2009. DNA mismatch repair deficiency in endometrial carcinoma. *Int J Gynecol Pathol* **28**: 239–255.
- Kim NG, Rhee H, Li LS, Kim H, Lee JS, Kim JH, Kim NK. 2002. Identification of MARCKS, FLJ11383 and TAF1B as putative novel target genes in colorectal carcinomas with microsatellite instability. *Oncogene* **21**: 5081–5087.
- Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, Jones SJ, Marra MA. 2009. Circos: An information aesthetic for comparative genomics. *Genome Res* **19**: 1639–1645.
- Li H, Durbin R. 2010. Fast and accurate long-read alignment with Burrows-Wheeler transform. *Bioinformatics* **26**: 589–595.
- Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078–2079.
- Lichtman AH, Blankman JL, Cravatt BF. 2010. Endocannabinoid overload. *Mol Pharmacol* **78**: 993–995.
- Lopez de Silanes I, Quesada MP, Esteller M. 2007. Aberrant regulation of messenger RNA 3'-untranslated region in human cancer. *Cell Oncol* **29**: 1–17.
- Markowitz S, Wang J, Myeroff L, Parsons R, Sun L, Lutterbaugh J, Fan RS, Zborowska E, Kinzler KW, Vogelstein B, et al. 1995. Inactivation of the type II TGF- β receptor in colon cancer cells with microsatellite instability. *Science* **268**: 1336–1338.
- McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, Garimella K, Altshuler D, Gabriel S, Daly M, et al. 2010. The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**: 1297–1303.
- Menoyo A, Alazzouzi H, Espin E, Armengol M, Yamamoto H, Schwartz S Jr. 2001. Somatic mutations in the DNA damage-response genes ATR and CHK1 in sporadic stomach tumors with microsatellite instability. *Cancer Res* **61**: 7727–7730.
- Mori Y, Sato F, Selaru FM, Oлару A, Perry K, Kimos MC, Tamura G, Matsubara N, Wang S, Xu Y, et al. 2002. Instability typing reveals unique mutational spectra in microsatellite-unstable gastric cancers. *Cancer Res* **62**: 3641–3645.
- Myeroff LL, Parsons R, Kim SJ, Hedrick L, Cho KR, Orth K, Mathis M, Kinzler KW, Lutterbaugh J, Park K, et al. 1995. A transforming growth factor β receptor type II gene mutation common in colon and gastric but rare in endometrial cancers with microsatellite instability. *Cancer Res* **55**: 5545–5547.
- Ottini L, Falchetti M, Lupi R, Rizzolo P, Agnese V, Colucci G, Bazan V, Russo A. 2006. Patterns of genomic instability in gastric cancer: Clinical implications and perspectives. *Ann Oncol* **17** (Suppl 7): vii97–vii102.
- Paun BC, Cheng Y, Leggett BA, Young J, Meltzer SJ, Mori Y. 2009. Screening for microsatellite instability identifies frequent 3'-untranslated region mutation of the RB1-inducible coiled-coil 1 gene in colon tumors. *PLoS ONE* **4**: e7715.
- Perucho M. 1996. Microsatellite instability: The mutator that mutates the other mutator. *Nat Med* **2**: 630–631.
- Pickering BM, Willis AE. 2005. The implications of structured 5' untranslated regions on translation and disease. *Semin Cell Dev Biol* **16**: 39–47.
- Pino MS, Chung DC. 2011. Microsatellite instability in the management of colorectal cancer. *Expert Rev Gastroenterol Hepatol* **5**: 385–399.
- Rampino N, Yamamoto H, Ionov Y, Li Y, Sawai H, Reed JC, Perucho M. 1997. Somatic frameshift mutations in the BAX gene in colon cancers of the microsatellite mutator phenotype. *Science* **275**: 967–969.
- Rebbapragada I, Lykke-Andersen J. 2009. Execution of nonsense-mediated mRNA decay: What defines a substrate? *Curr Opin Cell Biol* **21**: 394–402.
- Ribic CM, Sargent DJ, Moore MJ, Thibodeau SN, French AJ, Goldberg RM, Hamilton SR, Laurent-Puig P, Gryfe R, Shepherd LE, et al. 2003. Tumor microsatellite-instability status as a predictor of benefit from fluorouracil-based adjuvant chemotherapy for colon cancer. *N Engl J Med* **349**: 247–257.
- Royrvik EC, Ahlquist T, Rognes T, Lothe RA. 2007. Slip slidin' away: A duodecennial review of targeted genes in mismatch repair deficient colorectal cancer. *Crit Rev Oncog* **13**: 229–257.
- Ruggiero T, Olivero M, Follenzi A, Naldini L, Calogero R, Di Renzo MF. 2003. Deletion in a (T)₈ microsatellite abrogates expression regulation by 3'-UTR. *Nucleic Acids Res* **31**: 6561–6569.
- Shin KH, Park JG. 2000. Microsatellite instability is associated with genetic alteration but not with low levels of expression of the human mismatch repair proteins hMSH2 and hMLH1. *Eur J Cancer* **36**: 925–931.
- Shin N, You KT, Lee H, Kim WK, Song M, Choi HJ, Rhee H, Nam SW, Kim H. 2011. Identification of frequently mutated genes with relevance to nonsense mediated mRNA decay in the high microsatellite instability cancers. *Int J Cancer* **128**: 2872–2880.
- Silva AL, Romao L. 2009. The mammalian nonsense-mediated mRNA decay pathway: To decay or not to decay! Which players make the decision? *FEBS Lett* **583**: 499–505.
- Stratton MR. 2011. Exploring the genomes of cancer cells: Progress and promise. *Science* **331**: 1553–1558.
- Suraweera N, Iacopetta B, Duval A, Compoin A, Tubacher E, Hamelin R. 2001. Conservation of mononucleotide repeats within 3' and 5' untranslated regions and their instability in MSI-H colorectal cancer. *Oncogene* **20**: 7472–7477.
- Tougeron D, Fauquemburgue E, Rouquette A, Le Pessot F, Sesboue R, Laurent M, Berthet P, Mauillon J, Di Fiore F, Sabourin JC, et al. 2009. Tumor-infiltrating lymphocytes in colorectal cancers with microsatellite instability are correlated with the number and spectrum of frameshift mutations. *Mod Pathol* **22**: 1186–1195.
- Toyota M, Suzuki H. 2010. Epigenetic drivers of genetic alterations. *Adv Genet* **70**: 309–323.
- Trapnell C, Pachter L, Salzberg SL. 2009. TopHat: Discovering splice junctions with RNA-Seq. *Bioinformatics* **25**: 1105–1111.
- Vilar E, Gruber SB. 2010. Microsatellite instability in colorectal cancer—the stable evidence. *Nat Rev Clin Oncol* **7**: 153–162.
- Wang K, Kan J, Yuen ST, Shi ST, Chu KM, Law S, Chan TL, Kan Z, Chan AS, Tsui WY, et al. 2011. Exome sequencing identifies frequent mutation of ARID1A in molecular subtypes of gastric cancer. *Nat Genet* **43**: 1219–1223.
- Williams DS, Bird MJ, Jorissen RN, Yu YL, Walker F, Zhang HH, Nice EC, Burgess AW. 2010. Nonsense mediated decay resistant mutations are a source of expressed mutant proteins in colon cancer cell lines with microsatellite instability. *PLoS ONE* **5**: e16012.
- Woerner SM, Benner A, Sutter C, Schiller M, Yuan YP, Keller G, Bork P, Doberitz MK, Gebert JF. 2003. Pathogenesis of DNA repair-deficient cancers: A statistical meta-analysis of putative real common target genes. *Oncogene* **22**: 2226–2235.
- Woerner SM, Yuan YP, Benner A, Korff S, von Knebel Doeberitz M, Bork P. 2010. SelTarbase, a database of human mononucleotide-microsatellite mutations and their potential impact to tumorigenesis and immunology. *Nucleic Acids Res* **38**: D682–D689.
- Yuan Z, Shin J, Wilson A, Goel S, Ling YH, Ahmed N, Dopeso H, Hawer M, Nasser S, Montagna C, et al. 2009. An A13 repeat within the 3'-untranslated region of epidermal growth factor receptor (EGFR) is frequently mutated in microsatellite instability colon cancers and is associated with increased EGFR expression. *Cancer Res* **69**: 7811–7818.
- Zang ZJ, Cutcutache I, Poon SL, Zhang SL, McPherson JR, Tao J, Rajasegaran V, Heng HL, Deng N, Gan A, et al. 2012. Exome sequencing of gastric adenocarcinoma identifies recurrent somatic mutations in cell adhesion and chromatin remodeling genes. *Nat Genet* **44**: 570–574.
- Zhang L, Yu J, Willson JK, Markowitz SD, Kinzler KW, Vogelstein B. 2001. Short mononucleotide repeat sequence variability in mismatch repair-deficient cancers. *Cancer Res* **61**: 3801–3805.

Received July 10, 2012; accepted in revised form April 3, 2013.