

201220055A

厚生労働科学研究費補助金  
第3次対がん総合戦略研究事業

骨髄異形成症候群における  
エピゲノム修飾分子異常の解明

平成24年度 総括研究報告書

研究代表者 真田 昌

平成25(2013)年 4月

目 次

I. 平成24年度総括研究報告 骨髄異形成症候群におけるエピゲノム修飾分子異常の解明 真田 昌	-----	1
II. 研究成果の刊行に関する一覧表	-----	7
III. 研究成果の刊行物・別刷	-----	11

# I . 総括研究報告

## 骨髄異形成症候群におけるエピゲノム修飾分子異常の解明

研究代表者 真田 昌 東京大学がんゲノミクスプロジェクト 特任助教

### 研究要旨

骨髄異形成症候群(Myelodysplastic syndrome, MDS)は高齢者に好発する難治性造血器腫瘍である。MDSにおいてエピゲノム修飾分子にゲノム異常が生じていることが、近年明らかとなった。MDSで観察されるエピゲノム異常が、エピゲノム修飾分子の遺伝子変異に起因しているのか、明らかではなく、本研究ではゲノム・エピゲノム解析を統合し、両異常の関連を明らかとすることを目指した。本年度は、変異の有無が明らかとなった症例群について網羅的なメチル化解析を行い、前年度のゲノム異常との比較検討を行った。

すなわち、エピゲノム関連遺伝子の変異解析を行った192例のMDS（低リスクMDS 145例、高リスクMDS 47例）についてイルミナ社のHumanMethylation 450 BeadChipを用いて、DNAのシトシンのメチル化状態を解析した。MDSにおいて、正常末梢血に比しメチル化を受けている傾向にある4000遺伝子について教師なしクラスタリング解析を行い、メチル化パターンの異なる、いくつかのサブクラスに階層化された。192例において変異頻度が高いエピゲノム修飾関連遺伝子、*TET2* (35%)、*ASXL1* (20%)、*DNMT3A* (12%)、*IDH1/2* (6%) および *EZH2* (5%) について変異の有無（内は192例での変異頻度）、および臨床情報とクラスタリング結果を比較検討した。骨髄芽球の増加を伴うRAEB-1および2と診断されている症例は、増加を伴わない病型に比し、解析した遺伝子群においてはメチル化されている傾向が高かった。DNAメチルトランスフェラーゼとしてシトシンのメチル化に直接的に関わる*DNMT3A*遺伝子においては、変異の有無とクラスタリング結果には関連は認められなかったが、脱メチル化過程で重要な働きを有することが明らかとなりつつある*TET2*および*IDH*変異例は、特徴的なメチル化パターンを示すクラスター群に集中する傾向が認められ、ゲノムレベルでの変異がエピゲノム異常に関わっていることが示唆された。

### A. 研究目的

骨髄異形成症候群(MDS)は高齢者に好発する難治性造血器腫瘍であるが、高齢者に適した根治的治療がなく、急速な少子高齢化による患者数の増加も危惧される。MDSにおけるDNAのメチル化などのエピゲノム異常が生じていることは1990年代から報告され、近年脱メチル化剤やヒストン修飾酵素阻害剤などが

欧米で臨床応用され、一定の臨床効果が得られ、従来の抗腫瘍剤では効果が期待できないMDSにおいて重要な治療薬として認識され、本邦でも臨床応用が開始された。一方で、最近になり、MDSにおいて*EZH2*、*TET2*、*ASXL1*などのエピゲノム関連分子の後天的変異が報告されたが、他にも多くの分子がエピゲノム修飾に関与しており、それらの分子の変異の



有無は明らかではない。しかしMDSにおけるエピゲノム修飾異常が、同修飾に関わる分子のゲノム異常に起因しているのか、多数存在する他のエピゲノム修飾分子の変異の有無など不明な点は多く、MDSにおけるエピゲノム関連分子異常の全体像は明らかではない。本研究では、最新のゲノム解析技術を駆使し、MDS 検体におけるエピゲノム修飾関連分子の変異プロファイルを明らかとし(H23)、網羅的なメチル化プロファイル(H24)との関連を明らかとする。

## B. 研究方法

### (1) 実施経過

エピゲノム関連遺伝子の変異解析を行った192例のMDS（低リスクMDS 145例、高リスクMDS 47例）についてイルミナ社のHumanMethylation 450 BeadChipを用いて、DNAのシトシンのメチル化状態を解析した。本アレイ上には、480万箇所、17,990遺伝子のプローブが搭載をされているが、今回は、遺伝子のプロモーター領域に存在し、遺伝子の発現調整に関わりがあると推測されるプローブに限って解析を行った。正常末梢血における各プローブのシグナル強度と比較したメチル化の程度を3段階にスコア化し、MDSにおいて、正常末梢血に比しメチル化を強く受けている遺伝子について教師なしクラスタリング解析を行い、昨年度明らかとした変異プロファイル結果と比較検討を行った。

### (2) 倫理面の配慮

本研究で実施される患者検体を用いた遺伝子解析研究は、原則としてMDS細胞の体細胞突然変異を扱うものであるが、平成16年（平成20年改訂）文部科学省、厚生労働省および経済産業省告示第1号「ヒトゲノム・遺伝子研究に関する倫理指針」を遵守し、事前に検体提供施設の倫理委員会の承認を得ている。東京大学における遺伝子解析研究については、学内のヒトゲノム遺伝子解析研究倫理委員会の審査・承認済み（「造血器腫瘍における遺伝

子異常の網羅的解析」948-7）である。なお、承認済みの研究計画書に基づき、研究対象（検体提供）者から文書による同意を得た上で検体を採取し、匿名化作業を行った上で、遺伝子解析研究に用いた。

## C. 研究結果

昨年度、変異解析を行った192例のMDSについてイルミナ社のHumanMethylation 450 BeadChipを用いて、DNAのシトシンのメチル化状態を解析した。正常末梢血に比しメチル化を受けている傾向にある4000遺伝子に絞り込んだ後に、教師なし階層クラスタリング解析を行い、メチル化パターンの異なる、いくつかのサブクラスに階層化された。192例において変異頻度が高いエピゲノム修飾関連遺伝子、TET2（35%）、ASXL1（20%）、DNMT3A（12%）、IDH1/2（6%）およびEZH2（5%）について変異の有無（内は192例での変異頻度）、および臨床情報とクラスタリング結果を比較検討した。骨髄芽球の増加を伴うRAEB-1および2と診断されている症例は、増加を伴わない病型に比し、解析した遺伝子群においてはメチル化されている傾向が高かった。DNAメチルトランスフェラーゼとしてシトシンのメチル化に直接的に関わるDNMT3A遺伝子においては、変異の有無とクラスタリング結果には関連は認められなかったが、脱メチル化過程で重要な働きを有することが明らかとなりつつあるTET2およびIDH変異例は、特徴的なメチル化パターンを示すクラスター群に集中する傾向が認められ、ゲノムレベルでの変異がエピゲノム異常に関わっていることが示唆された。

## D. 考察

エピゲノム修飾の異常は様々ながん種で観察され、がん抑制遺伝子のプロモーター領域の異常メチル化を介した発現低下など、重要な発がんメカニズムの一つとして認識されてきた。また、近年の遺伝子解析技術の進歩に伴い、エピゲノム修飾分子に遺伝子変異が

広範ながんにおいて生じていることが明らかとなっている。しかし、ゲノムの異常である遺伝子変異が、患者細胞において、エピゲノム修飾に関わっているのか、さらには如何にして腫瘍化に寄与しているかなど不明な点が多い。

MDS は造血幹細胞に由来する腫瘍性疾患であるが、メチル化異常が生じていることが数多く報告をされ、また、脱メチル化剤などエピゲノム修飾を標的とした治療薬剤の臨床応用もされ、エピゲノム修飾遺伝子の変異が高頻度に生じていることが近年明らかとなっており、エピゲノム異常が病態に大きく関与していると推測される代表的な腫瘍である。本研究では、MDS においてエピゲノム修飾関連遺伝子における遺伝子変異は高頻度に観察され、エピゲノム修飾異常を招いていると考えられることを、TET2 および IDH1/2 変異とメチル化異常について臨床検体の解析を通じて明らかとした。近年の研究により、TET2 はメチル化シトシンの脱メチル化過程で重要なメチル化シトシンヒドロキシラーゼ活性を有していることが明らかとなり、また IDH 変異体は  $\alpha$  ケトグルタル酸から 2 ヒドロキシグルタル酸への変換を介して、TET2 の酵素活性を阻害することが示されている。すなわち、TET2 の不活化変異および IDH 変異は、ともに脱メチル化過程が障害を受け、メチル化状態が維持されることが予測され、我々のメチル化解析結果とも合致する。TET2 および IDH1/2 変異と関連の高いクラスターを特徴づけるメチル化されている領域・遺伝子を検索することを現在進めており、本変異に導かれるエピゲノム修飾異常のメカニズムならびに標的遺伝子を通じた MDS の分子病態の解明が進むことが期待される。

脱メチル化剤を用いた治療は、日本でも一昨年より高リスク MDS 例を中心に行われているが、奏効率は必ずしも高くはなく、作用機序も不明である。今後、脱メチル化剤投与前後の検体の解析を行うことで、脱メチル化剤投

与によるメチル化の変化、有効例と無効例における変化の違いを解析することで、脱メチル化剤の作用機序を明らかとし、治療反応性を事前に予測する臨床上有用なバイオマーカーの確立も望まれる。

MDS における RNA スプライシング関連分子の変異は排他的に生じているのに対し、エピゲノム修飾関連遺伝子の変異は、TET2 と IDH1/2 の変異の重複例は既報の通りに少ないものの、TET2 変異と ASXL1 変異など、しばしば重複して観察をされ、アレル頻度からも、同一の細胞に変異が生じていると推測される。すなわち、遺伝子異常が、エピゲノム修飾全体そして遺伝子発現に与える影響は単純ではないことが推測をされる。今後、遺伝子発現解析やヒストン修飾の解析も含めた、より多層的な解析が必要であると思われる。

## E. 結論

MDS 症例においては、エピゲノム修飾分子にゲノムレベルでの異常が既知の遺伝子のみならず、高頻度に生じていることが本研究を通じて明らかとなった。更にエピゲノム解析を行うことにより、MDS におけるエピゲノム修飾分子のゲノム異常とエピゲノム異常の関わりが明らかになることが期待される。

## F. 健康危険情報

なし

## G. 研究発表

### 1 論文発表

- 1) Ueda T, Sanada M, Matsui H, Yamasaki N, Honda ZI, Shih LY, Mori H, Inaba T, Ogawa S, Honda H. EED mutants impair polycomb repressive complex 2 in myelodysplastic syndrome and related neoplasms. **Leukemia**. 2012;26(12):2557-60.
- 2) Hosokawa K, Katagiri T, Sugimori N, Ishiyama K, Sasaki Y, Seiki Y, Sato-Otsubo A, Sanada M, Ogawa

- S, Nakao S. Favorable outcome of patients who have 13q deletion: a suggestion for revision of the WHO 'MDS-U' designation. **Haematologica**. 2012; 97(12): 1845-9.
- 3) Takita J, Yoshida K, Sanada M, Nishimura R, Okubo J, Motomura A, Hiwatari M, Oki K, Igarashi T, Hayashi Y, Ogawa S. Novel splicing-factor mutations in juvenile myelomonocytic leukemia. **Leukemia**. 2012; 26: 1879-81.
- 4) Sanada M, Ogawa S. Genome-wide analysis of myelodysplastic syndromes. **Curr Pharm Des**. 2012;18:3163-9.
- 5) Shiraishi Y, Sato Y, Chiba K, Okuno Y, Nagata Y, Yoshida K, Shiba N, Hayashi Y, Kume H, Homma Y, Sanada M, Ogawa S, Miyano S. An empirical Bayesian framework for somatic mutation detection from cancer genome sequencing data. **Nucleic Acids Res**. 2013 in press
- 6) Kon S, Minegishi N, Tanabe K, Watanabe T, Funaki T, Wong WF, Sakamoto D, Higuchi Y, Kiyonari H, Asano K, Iwakura Y, Fukumoto M, Osato M, Sanada M, Ogawa S, Nakamura T, Satake M. Smad1 deficiency perturbs receptor trafficking and predisposes mice to myelodysplasia. **J Clin Invest**. 2013;123(3):1123-37.
- 7) Kunishima S, Okuno Y, Yoshida K, Shiraishi Y, Sanada M, Muramatsu H, Chiba K, Tanaka H, Miyazaki K, Sakai M, Ohtake M, Kobayashi R, Iguchi A, Niimi G, Otsu M, Takahashi Y, Miyano S, Saito H, Kojima S, Ogawa S. ACTN1 Mutations Cause Congenital Macrothrombocytopenia. **Am J Hum Genet**. 2013; 92(3):431-8.
- 2 学会発表
- 真田 昌 Novel pathway mutations in myelodysplasia revealed by high-throughput sequencing technology. 第 71 回日本癌学会学術総会 International session (招待講演) 2012 年 9 月 20 日 札幌
- H. 知的財産権の出願・登録状況
1. 特許出願中  
なし
  2. 実用新案登録  
なし
  3. その他  
なし

## Ⅱ．研究成果の刊行に関する一覧表



研究成果の刊行に関する一覧表

書籍

著者氏名	論文タイトル名	書籍全体の編集者名	書 籍 名	出版社名	出版地	出版年	ページ
真田 昌	正常核型急性骨髄性白血病のゲノム異常	高久史麿 小澤敬也 坂田洋一 金倉 讓 小島勢二	Annual Review w 2013 血液	中外医学 社	東京都 新宿区	2012年	80-86ページ

雑誌

発表者氏名	論文タイトル名	発表誌名	巻号	ページ	出版年
Sanada M, Ogawa S	Genome-wide analysis of myelodysplastic syndromes.	Curr Pharm D es.	18	3163-9.	2012
Shiraishi Y, Sato Y, Chiba K, Okuno Y, Nagata Y, Yoshida K, Shiba N, Hayashi Y, Kume H, Homma Y, Sanada M, Ogawa S, Miyano S.	An empirical Bayesian framework for somatic mutation detection from cancer genome sequencing data.	Nucleic Acid Research		In press	2013
Kon S, Minegishi N, Tanabe K, Watanabe T, Funaki T, Wong WF, Sakamoto D, Higuchi Y, Kiyonari H, Asano K, Iwakura Y, Fukumoto M, Osato M, Sanada M, Ogawa S, Nakamura T, Satake M.	Smad1 deficiency perturbs receptor trafficking and predisposes mice to myelodysplasia.	J Clin Invest.	123(3)	1123-37	2013
Ueda T, Sanada M, Matsui H, Yamasaki N, Honda ZI, Shih LY, Moris. H, Inaba T, Ogawa S, Honda H.	EED mutants impair polycomb repressive complex 2 in myelodysplastic syndrome and related neoplasms.	Leukemia	26(12)	2557-60	2012
Takita J, Yoshida K, Sanada M, Nishimura R, Okubo J, Motomura A, Hiwatari M, Oki K, Igarashi T, Hayashi Y, Ogawa S.	Novel splicing-factor mutations in juvenile myelomonocytic leukemia.	Leukemia	26(8)	1879-81	2012
Hosokawa K, Katagiri T, Sugimori N, Ishiyama K, Sasaki Y, Seiki Y, Sato-Otsubo A, Sanada M, Ogawa S, Nakao S.	Favorable outcome of patients who have 13q deletion: a suggestion for revision of the WHO 'MDS-U' designation.	Haematologic	97	1845-9	2012

### Ⅲ. 研究成果の刊行物・別刷

# Genome-wide Analysis of Myelodysplastic Syndromes

Masashi Sanada\* and Seishi Ogawa

Cancer Genomics Project, Graduate School of Medicine, The University of Tokyo, Tokyo, Japan

**Abstract:** Myelodysplastic syndromes (MDS) are heterogeneous hematopoietic neoplasms characterized by ineffective hematopoiesis and a risk for progression to acute myeloid leukemia. A number of cytogenetic changes have been described that are characteristic to MDS and of clinical relevance; the specific gene targets of these alterations were largely unknown. On the other hand, over the past decade, technologies have been dramatically improved to enable high-throughput analysis of entire MDS genomes, leading to identification of frequent copy number neutral events and a number of novel gene targets implicated in the pathogenesis of MDS. In this review, we briefly overview the recent progress in the genetics of MDS, focusing on the newly identified gene targets in MDS.

**Keywords:** Microarray, SNP array, CNN-LOH, somatic mutation, high-throughput parallel sequencing.

## INTRODUCTION

Myelodysplastic syndromes (MDS) are intractable clonal disorders of hematopoietic systems characterized by bone marrow dysplasia, peripheral blood cytopenia due to ineffective hematopoiesis, and a high propensity to acute myeloid leukemia (AML) [1, 2]. One of the prominent features of MDS is the high frequency of unbalanced chromosomal changes that accompany copy number alterations of chromosomal segments. Gains and losses of one or more chromosomal segments are found in approximately 50% of MDS patients in conventional cytogenetics and represent major determinants of the prognosis of MDS [3-5], indicating that these changes could be closely related to the pathogenesis of MDS. Unfortunately, however, most of the common changes typically involve large chromosomal segments, and with the lack of specific positional markers that pinpointed the critical genetic loci, the gene targets of these chromosomal lesions have not been determined until recently. This shows a stark contrast to *de novo* AML, where the breakpoints of disease type-specific translocations provided reliable positional markers to identify the major gene fusions that are relevant to molecular classification and characterization of AML [6,7]. The breakthrough for this situation has been brought about over the past decade, during which there have been dramatic improvements in genome technologies that allowed high-throughput/ resolution analysis of genomes [8], particularly with the development of single nucleotide polymorphism (SNP) array-based technology for copy number analysis. The SNP array-based copy number detection technologies enabled detection of copy-number (CN) alterations as well as allelic imbalances or loss of heterozygosity (LOH) in cancer genomes [9-13] and successfully applied to the analysis of MDS genomes, leading to the identification of a number of novel gene targets, frequently mutated in MDS as well as other myeloid cancers [14-18]. Interestingly, many of the newly identified mutational targets are those involved in epigenetic regulation, such as DNA methylation and chromatin modifications, which is in accordance with the clinical observation that demethylating agents (azacitidine and decitabine) have been demonstrated to be effective in the treatment of high-risk MDS patients [19-21]. Thus, the frequent mutations of epigenesis-regulating genes support the possibility that the epigenetic alterations in MDS could be at least partly explained by the primary genetic alterations.

## CYTOGENETICS IN MDS

Conventional cytogenetics provides an invaluable clue to the management of MDS, since the types and numbers of chromosomal lesions have been tightly linked to the prognosis of MDS cases.

Cytogenetic findings are among the key parameters for the prediction of prognosis in the International Prognostic Scoring System (IPSS), and also adopted for the World Health Organization (WHO) classification-based Prognostic Scoring System (WPSS) [22]. Hasse *et al.* and other researchers also demonstrated that rare but recurrent cytogenetic alterations and specific karyotypic combinations could be used as beneficial markers for determining the prognosis of MDS [4, 23-25]. On the other hand, a potential caveat in conventional cytogenetics is that it absolutely depends on viable cells to obtain metaphases for analysis. Conventional cytogenetics fails to detect any abnormalities in approximately half of the patients with MDS. In fact, using interphase fluorescent in situ hybridization (FISH) analysis with 4 FISH probes, Rigolin *et al.* reported occult cytogenetic alterations in 17.8% of MDS patients with normal karyotype, including deletions of 5q31, 7q31 and 17p13, as well as trisomy8 [26]. Although providing a sensitive method for detecting submicroscopic alterations of known targets that are present in a small fraction of tumor samples without depending on cell divisions, interphase FISH analysis cannot be applied to genome-wide detection of genetic lesions.

## ARRAY COMPARATIVE GENOMIC HYBRIDIZATION

Array-based comparative genomic hybridization (aCGH) enables comprehensive genome-wide analysis of genetic aberrations in cancers [8], in which differentially labeled DNAs from both tumor and normal reference samples are comparatively hybridized to a large number of probes on microarray. The ratio of the signal intensity of the test to that of the reference DNA is then calculated for the measurements of genomic copy numbers. The density of probes on microarray has been increased up to 4.2 million (NimbleGen), allowing for detection of smaller, more focal amplifications and deletions [27,28]. In the previous studies of MDS, a number of small, cryptic chromosomal abnormalities were identified using a CGH that could otherwise escape conventional cytogenetic analysis [29-32].

## SNP ARRAY ANALYSIS

High density SNP arrays were originally developed for large-scale genotyping that is required for genome-wide association studies (GWAS) [33, 34]. However, the quantitative nature of the preparative whole-genome amplification and array hybridization thereafter allows for accurate estimation of genomic copy numbers at high resolution [35-37]. Furthermore, SNP array analysis also enables genome-wide LOH detection using genotyping data. With these desirable features, SNP arrays are now widely used for genome-wide copy number and LOH analyses in cancer research and the diagnosis of rare congenital disorders [10, 12-14,38,39]. Currently, two SNP array platforms are commercially available, Affymetrix GeneChip SNP Genotyping array [33] and Illumina beads array [40]. A number of software are developed for the analysis of

Address correspondence to this author at the Cancer Genomics Project, Graduate School of Medicine, The University of Tokyo, Tokyo, Japan; Tel: +81-3-3815-5411; Ext: 35615; Fax: +81-3-5800-9047; E-mail: sanada-ky@umin.ac.jp

genomic copy numbers [35, 37, 41, 42], among which CNAG/AsCNAR software [36, 43], is one of the most widely used for this purpose. CNAG/AsCNAR implements with a series of data compensation algorithms to accurately estimate copy numbers. In addition, by detecting subtle distortions in allele-specific signals caused by allelic imbalance, CNAG/AsCNAR enables sensitive detection of LOH with accurate determination of allele-specific copy numbers even in the face of up to 80% normal cell contamination [43].

Using AffymetrixGeneChip50k or 250k array, we analyzed a total of 222 MDS and myelodysplastic/myeloproliferative neoplasms (MDS/MPN) specimens, 87 of the 137 MDS cases (63.5%) had one or more regions showing allelic imbalances [14] Fig. (1). In accordance with previous cytogenetic studies, MDS genomes showed high frequencies of unbalanced genetic changes, including  $-5/5q-$ ,  $-7/7q-$ ,  $+8$ ,  $9p+$ ,  $12p-$ ,  $17p-$ ,  $18q+$ ,  $19p+$ ,  $19q+$ ,  $20q-$ , and  $21q+$ , which were detected with higher sensitivity using SNP arrays. For example, hidden copy number alterations were successfully detected by SNP array-based copy number analysis in 14 out of 55 cases of normal karyotype MDS in our series [14]. However, the major advantage of SNP array analysis is the ability to detect genome-wide copy-number neutral (CNN)-LOH, which is undetectable by conventional cytogenetics, FISH or array CGH.

### CNN-LOH IN MDS

CNN-LOH or uniparental disomy (UPD) is a common genetic alteration in cancer genome, majority of LOH in cancer being due to CNN-LOH rather than simple allelic deletion. Although CNN-LOH has been considered to be a common mechanism of inactivation of tumor suppressor genes, the discovery of a gain-of-function mutation of *JAK2* kinase associated with  $9pUPD$  in myeloproliferative neoplasms (MPN) lead to a concept that CNN-LOH could also

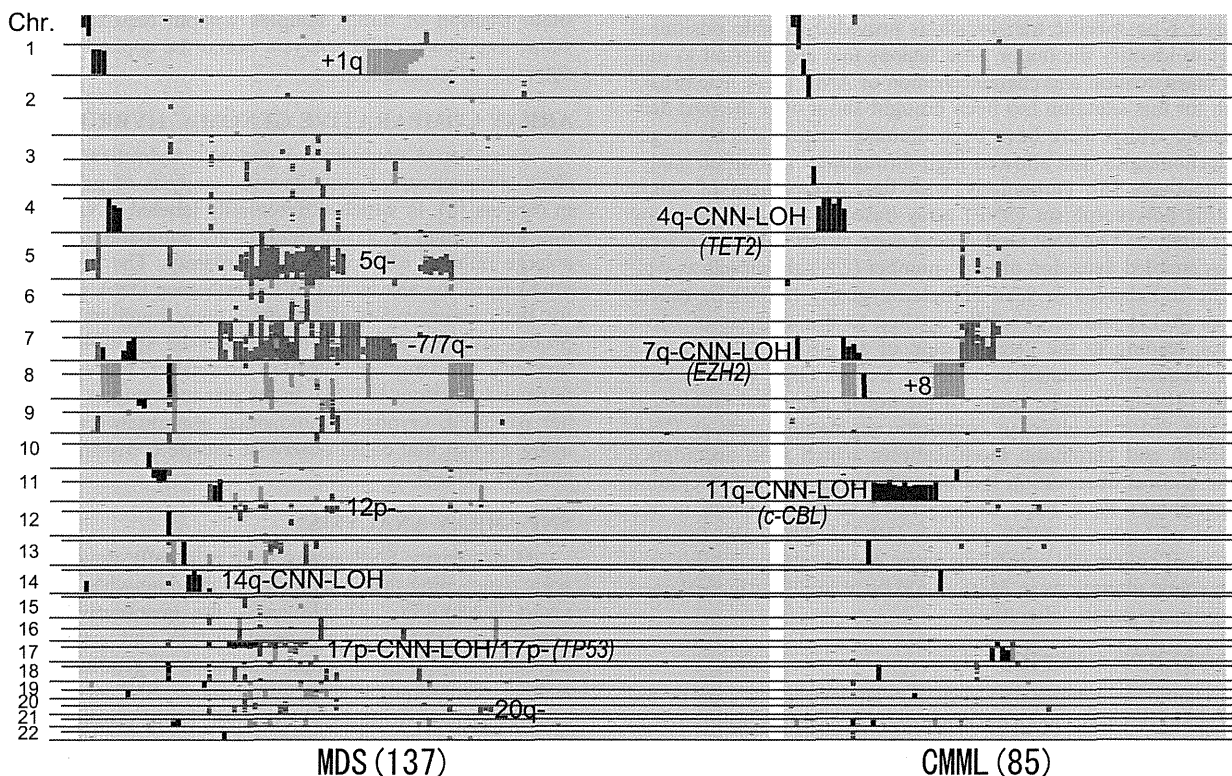
provide the genetic mechanism for clonal selection of a gain-of-function mutation [44]. CNN-LOH has been documented in 10-25% of MDS cases [14, 45, 46], 10-20% of *de novo* AML [47-52], and over 35% of chronic myelomonocytic leukemia (CMML) cases [14, 45].

Similar to other allelic imbalances, CNN-LOH was not randomly distributed throughout the MDS genomes, but tended to involve particular chromosomal arms in a relatively mutually exclusive manner, including 1p, 1q, 4q, 7q, 11p, 11q, 14q, 17p, and 21q Fig. (1). Among these, 7q, 17p, and 21q are also affected by deletions, while LOH in other arms were largely caused by UPD. In contrast, 5q and 20q are frequent targets of deletion in MDS cases, but rarely show CNN-LOH. CNN-LOH in 11p, 13q, 17p and 21q were also seen in *de novo* AML cases, whereas 11q CNN-LOH was typically found in cases with MDS/MPN. A significant finding about these recurrent CNN-LOH is that they are frequently associated with homozygous mutations of known gene targets of myeloid neoplasms, including *c-MPL* or *N-RAS* in 1pCNN-LOH [14, 53], *JAK2* in 9pCNN-LOH [43, 44], *FLT3* in 13qCNN-LOH [54], *TP53* in 17pCNN-LOH [14], and *RUNX1* in 21qCNN-LOH [14, 54] (Table 1). CNN-LOH could result in the duplication of mutated oncogenes after the loss of the normal allele or by inducing deletion of tumor suppressor genes.

### MUTATED GENE TARGETS IN MDS (FIG. 2)

#### 1) TET2

The long arm of chromosome 4 has not been reported as a common target of chromosomal abnormalities in myeloid malignancies in conventional cytogenetics [4], but recently turned out to be a recurrent target of CNN-LOH in MDS and CMML in SNP array analysis. Delhommeau *et al.* and Langeimer *et al.* identified



**Fig. (1).** The genome profile of 222 cases of MDS and related myeloid neoplasms detected by SNP array analysis.

The genetic alterations, including CN gains, losses and CNN-LOH, are color-coded, light gray, gray, and dark gray, respectively. These lesions are plotted vertically in chromosomal order for each sample. Vertical positions of each lesion are proportional to the genetic length and thus the size of the color-coded corresponds to the length of alterations. CNN-LOH, in particular chromosomal arms tends to be found in mutually exclusive cases, enabling clustering based on the site of CNN-LOH, except for 17pLOH, which was frequently accompanied by loss of 5q, loss of chromosome 7 or 7q, and loss of 12p. Common genetic alterations and their target genes are indicated.

loss of function mutations of *TET2* as the target of 4qLOH [15,16], and also mutated frequently in other cases without having 4qLOH. In fact, *TET2* is now shown to represent one of the most frequently mutated genes in MDS (~20%) as well as other myeloid neoplasms [55], including MPN with or without *JAK2-V617F* mutations (~10%), CMML (30-50%), and part of AML (13%) [15, 16, 56, 57]. *TET2* mutations frequently occur during progression of MPN or MDS to secondary AML. The impact of *TET2* mutations on clinical outcomes is still controversial. Some reports demonstrated significantly shorter overall survival in patients with *TET2* mutations [56-58], while others reported favorable or no prognostic impact of *TET2* mutations [16, 55, 59].

TET family proteins (Tet1, Tet2, and Tet3) catalyze the conversion of 5-methyl-cytosine to 5-hydroxymethyl-cytosine (5hmC) [60, 61]. In ES cells, TET1 plays a functional role in maintaining the pluripotent state [61-63]. A recent study demonstrated that 5hmC generated by TET activity is an intermediate during the process of DNA demethylation [64]. In addition, TET1 directly interacts with Sin3A, a co-repressor protein essential for inhibiting the transcription of a subset of genes [65]. *Tet2* deficiency in mice lead to the progressive enlargement of the hematopoietic stem and progenitor compartment, and also results in abnormalities in mature myeloid and lymphoid cells, and leading to fatal hematopoietic malignancies [66]. Quivoron *et al.* also found that *TET2* mutations were not only seen in myeloid neoplasms but also in various types of B- and T-cell lymphoid tumors in humans.

## 2) IDH1/IDH2

Mutations of isocitrate dehydrogenase (IDH) 1 and IDH2 are initially identified through comprehensive mutation studies in glioblastoma as well as *de novo* AML in high frequencies [67, 68], and also reported in other myeloid malignancies including secondary AML, MDS and MPN [69-73]. IDH1 and IDH2 are components of TCA enzymes that catalyze isocitrate to  $\alpha$ -ketoglutarate conversion in cytoplasm and mitochondria, respectively. Mutations of IDH1 and IDH2 exclusively involved in amino acid positions of R132 in IDH1 and R140 and R172 in IDH2, respectively, indicating they represent gain-of-function, rather than loss of function mutations. In fact, these mutations were shown to cause dramatic alteration of substrate specificity. As a result, the mutated enzymes show severely compromised activity of the intrinsic isocitrate to  $\alpha$ -ketoglutarate conversion, but in turn acquire a *de novo* activity to catalyze  $\alpha$ -ketoglutarate to 2-hydroxyglutarate (2HG) conversion. The 2HG represents the first example of oncogenic metabolite in human cancers. Intriguingly, 2HG competitively inhibits TET2 function, which absolutely depends on  $\alpha$ -ketoglutarate as a substrate [74]. In fact, the *IDH1/2* mutations were always heterozygous and tend to occur in a mutually exclusive manner with *TET2* mutations.

## 3) C-CBL

11qUPD is one of the most common targets of UPD found in myelodysplasia, particularly in CMML with normal karyotypes. We and other groups identified *C-CBL* mutations as the critical gene affected by 11qCNN-LOH [14, 45, 75, 76]. *C-CBL* is the cellular homolog of the *v-Cbl* transforming gene of Cas NS-1 murine leukemia virus, and is thought to negatively regulate tyrosine kinase signaling, mainly through the down-regulation of activated tyrosine kinases by E3 ubiquitin ligase activity [77]. *C-CBL* mutations are frequently seen in MDS/MPN cases with a tight association with 11q-CNN-LOH. *C-CBL* mutations and other *RAS* pathway mutations (*NRAS*, *KRAS*, *PTPN11*, and *NF1*) occur in a mutually exclusive manner in CMML and juvenile myelomonocytic leukemia (JMML) [76, 78, 79]. Interestingly in this regard, similar to other mutations of *RAS* pathway genes, heterozygous germ-line *C-CBL* mutations may predispose the development of JMML with a Noonan Syndrome-like phenotype [80, 81]. Most *C-CBL* mutations

in myeloid malignancies are found in the linker and RING finger domains, which are central to the E3 ubiquitin ligase activity [82]. *C-CBL* mutants show compromised E3 ubiquitin ligase activity, and also inhibit wild type *C-CBL* and *CBLB*, leading to prolonged activation of tyrosine kinases following cytokine stimulation [14, 83, 84], leading to hypersensitivity to a wide spectrum of cytokines that underlies the pathogenesis of the myeloproliferative phenotype commonly found in CMML and JMML [82, 84].

## 4) EZH2

Loss of chromosomes 7 and 7q are one of the most frequent genetic alterations in MDS and known as a reliable predictor of adverse prognosis. Approximately 10% of the patients with MDS carry an abnormality of chromosome 7, either alone or as part of a complex karyotype. This frequency is higher in therapy-related MDS associated with a prior history of treatment with alkylating agents. SNP array analysis has revealed that not only copy number loss but also CNN-LOH is the cause of 7qLOH in MDS and related myeloid neoplasms. Recently, Ernst *et al.* and Nikoloski *et al.* have shown that *EZH2* is mutated in some cases with 7q-LOH [17,18], indicating that *EZH2* is one of the gene targets in 7qLOH. *EZH2* encodes a histone methyltransferase that is the catalytic component of the polycomb repressive complex-2 (PRC2), a highly conserved histone H3 at lysine-27 methyl transferase, which functions to initiate epigenetic silencing of genes involved in cell fate decisions [85]. Loss of PRC2 function increases hematopoietic stem cell activity and expansion, which may explain how loss of function mutations of *EZH2* leads myeloid neoplasms [86]. On the other hand, at least three common deleted regions (CDRs) on 7q (7q22, 7q32-33, and 7q35-36) have been identified in myeloid malignancies [87-89], and therefore, *EZH2*(7q36) does not seem to be the sole target for the deletions of chromosome 7q.

## 5) Ribosomal Protein

Deletion of chromosome 5q is also a common cytogenetic alteration in MDS, and isolated 5q- is associated with a favorable prognosis and a favorable response to lenalidomide [90, 91]. Many studies attempted to narrow the region of recurrent somatic deletion to identify the critical gene in this region, but no somatic mutations have been identified among genes located within the CDR of 5q [92, 93]. SNP array analysis did not contribute to narrow the 5qCDR, which is rarely affected by CNN-LOH in MDS. It has been suggested that haplo-insufficiency in one or more genes may explain 5q- pathogenesis, instead of bi-allelic inactivation of a tumor suppressor gene. Ebert *et al.* performed an RNA interference screen against all 40 genes located within the 5qCDR and implicated haplo-insufficiency of the *RPS14* gene as a major contributor to the hematologic manifestations of 5q-[94]. Barlow *et al.* generated deletions of portions of syntenic lesion (containing *RPS14*) with the human 5q region in mouse, haplo-insufficiency of this loci caused macrocytic anemia, increased apoptosis and the morphologic abnormalities found in the erythroid compartment [95]. Loss-of-function mutations involving other ribosomal components (e.g., *RPS19* and *RPS24*) have also been implicated in rare congenital bone marrow failure syndromes, Diamond-Blackfan anemia [96, 97]. Nevertheless, haploinsufficiency of *RPS14* does not seem to explain several other features of the 5q-syndrome, which also shows thrombocytosis associated with megakaryocytic dysplasia, neutropenia, and clonal dominance [98, 99]. Interestingly, a recent study has demonstrated that haplo-insufficiency of two micro RNAs within CDR, *miR-145* and *miR-146*, could also contribute to the pathogenesis of 5q- syndrome, supporting a model of haploinsufficiency of multiple gene targets in this syndrome [100].

## CLINICAL APPLICATION

Given that cytogenetic information provides a valuable clue to the management of MDS as prognostic makers, a more accurate prognosis could be established based on SNP array or other CGH

Table 1. Recurrent Gene Mutations in Myeloid Malignancies

Mutated Gene	Diseases	frequency in MDS	frequency in de novo AML	Associated chromosomal alterations	pathway
<i>TET2</i>	MDS, CMML, MPN	20.0%	13.2%	4qUPD	epigenetic modification
<i>EZH2</i>	MDS, CMML	6.0%	rare	7qUPD	epigenetic modification
<i>ASXL1</i>	AML, MDS, CMML	10-15%	10.8%		epigenetic modification
<i>DNMT3A</i>	AML, MDS	8.0%	22.1%		epigenetic modification
<i>IDH1</i>	AML, MDS	rare-5.2%	6.6-8.5%	normal cytogenetics	epigenetic modification
<i>IDH2</i>	AML, MDS, CMML	4.2%	11-15.4%		epigenetic modification
<i>TP53</i>	AML, MDS	5-10%	<10%	17ploss/UPD, complex karyotype	cell cycle, apoptosis
<i>Nras</i>	MDS, AML, MDS/MPN	3.6-6.3%	10-15%	1pUPD	signal transduction
<i>Kras</i>	MDS, AML	rare	5.0%		signal transduction
<i>cMPL</i>	MPN, RARS <sup>t</sup>	rare-5%	rare	1pUPD	signal transduction
<i>JAK2</i>	MPN, RARS <sup>t</sup>	rare-50%	rare	9pUPD	signal transduction
<i>c-CBL</i>	CMML, JMML	rare	rare	11qUPD	signal transduction
<i>FLT3</i>	AML	rare	28-33%(ITD), 5-10%	13qUPD	signal transduction
<i>NF1</i>	JMML	rare	rare	17qUPD	signal transduction
<i>PTPN11</i>	JMML	rare	rare		signal transduction
<i>c-KIT</i>	AML	rare	6-10%		signal transduction
<i>RUNX1</i>	AML, MDS	15-20%	8.6%	21qloss/UPD	transcriptional factor
<i>WT1</i>	AML	rare	10.0%	11pUPD	transcriptional factor
<i>CEBPA</i>	AML	rare	4-9%	19pUPD	transcriptional factor
<i>U2AF35</i>	MDS	11.6%	rare		RNA splicing
<i>SRSF2</i>	MDS, CMML	11.6%	rare		RNA splicing
<i>SF3B1</i>	RARS, MDS	6.5-75.3%	rare		RNA splicing
<i>ZRSR2</i>	MDS	7.7%	rare		RNA splicing
<i>NPM1</i>	AML	rare	25-35%	normal cytogenetics	other

rare, mutations present in <3% of patients

MDS, myelodysplastic syndrome; RARS, refractory anemia with ringed sideroblasts; RARS<sup>t</sup>, RARS and thrombocytosis

MPN, myeloproliferative neoplasm; AML, acute myeloid leukemia; CMML, chronic myelomonocytic leukemia; JMML, juvenile myelomonocytic leukemia

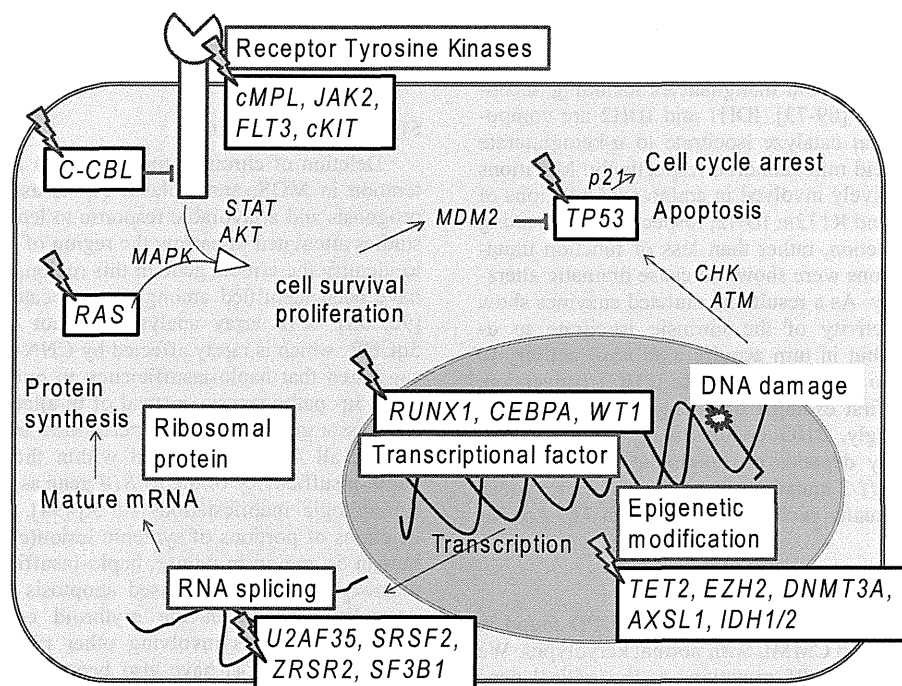


Fig. (2). Molecular pathways of genes affected in MDS.

Mutations of multiple pathways have been indicated in the pathogenesis of MDS. The mutated components are indicated by arrows.

based genomic analysis of MDS. Array-based genome-wide copy number analysis can provide much information on genetic alterations, especially on CNN-LOH, although array-based analysis cannot detect the balanced translocations that are relevant to the management of a large number of hematopoietic malignancies.

Some studies showed that the presence of newly detected alterations by microarray were useful as novel predictors of prognosis [101]. Heinrichs *et al.* and Godek *et al.* showed that 7q-CNN-LOH is a possible marker for poor prognosis [45, 46], although the evi-

dence for the value of each alteration identified with SNP array or aCGH has so far been still incomplete. Clearly, further studies are required to establish the clinical values of array-based karyotyping technologies in MDS. Recently, Bejar *et al.* examined whether the mutation profile of known target genes was associated with the clinical phenotype, and found that mutations in *TP53*, *EZH2*, *ETV6*, *RUNX1* and *ASXL1* are independent predictors of poor prognosis [55]. However, most reported mutations occur infrequently in MDS cases and are also found in the case of AML and other myeloid



neoplasms (Table 1, Fig. (2)). These mutations may explain the limited aspect of pathogenesis of MDS.

## CONCLUSION AND RECENT PROGRESS

One of the best targets of SNP-array based genome-wide allelkaryotyping would be MDS and related disorders in which CNN-LOH and unbalanced genetic changes are predominant. Using SNP array, several novel gene mutations, *C-CBL*, *TET2*, and *EZH2*, have been identified in MDS and related myeloid neoplasms. However, as many as 20-30% of primary MDS cases do not show any genetic changes even with SNP array karyotyping or mutation analysis of previously known targets. More problematic is that no gene mutations are specific to MDS but also found in other myeloid cancers, indicating that we still have incomplete knowledge about the molecular pathogenesis of MDS. In this regard, the development of high-throughput parallel sequencing technologies has provided an opportunity to characterize genetic changes across the genome-wide sequences at single nucleotide level [102], and is expected to be successfully applied to the genetic analysis of MDS to reveal more aspects of their pathogenesis in near future. In fact, our recent study using whole exome sequencing has revealed high frequencies (45-85% depending on subtypes of MDS) of pathway mutations involving multiple components of the splicing machinery that are highly specific to myeloid neoplasms showing features of myelodysplasia [103], although more studies are required to elucidate their roles in the pathogenesis of MDS.

## REFERENCES

- Tefferi A, Vardiman JW. Myelodysplastic syndromes. *N Engl J Med* 2009; 361(19): 1872-85.
- Bejar R, Levine R, Ebert BL. Unraveling the molecular pathophysiology of myelodysplastic syndromes. *J Clin Oncol* 2011 Feb 10; 29(5): 504-15.
- Greenberg P, Cox C, LeBeau MM, *et al.* International scoring system for evaluating prognosis in myelodysplastic syndromes. *Blood* 1997; 89(6): 2079-88.
- Haase D, Germing U, Schanz J, *et al.* New insights into the prognostic impact of the karyotype in MDS and correlation with subtypes: evidence from a core dataset of 2124 patients. *Blood* 2007; 110(13): 4385-95.
- Garcia-Manero G. Prognosis of myelodysplastic syndromes. *Hematology Am Soc Hematol Educ Program* 2010; 2010: 330-7.
- Schoch C, Haferlach T. Cytogenetics in acute myeloid leukemia. *Curr Oncol Rep* 2002; 4(5): 390-7.
- Swerdlow S, Campo E, Harris N, *et al.* World Health Organization Classification of Tumours: Pathology and Genetics of Tumours of Haematopoietic and Lymphoid Tissues. Lyon: IARC; 2008.
- Pinkel D, Albertson DG. Array comparative genomic hybridization and its applications in cancer. *Nat Genet* 2005; 37 Suppl: S11-7.
- Rauch A, Ruschendorf F, Huang J, *et al.* Molecular karyotyping using an SNP array for genomewide genotyping. *J Med Genet* 2004; 41(12): 916-22.
- Mullighan CG, Goorha S, Radtke I, *et al.* Genome-wide analysis of genetic alterations in acute lymphoblastic leukaemia. *Nature* 2007; 446(7137): 758-64.
- Kawamata N, Ogawa S, Zimmermann M, *et al.* Molecular allelkaryotyping of pediatric acute lymphoblastic leukemias by high-resolution single nucleotide polymorphism oligonucleotide genomic microarray. *Blood* 2008; 111(2): 776-84.
- Chen Y, Takita J, Choi YL, *et al.* Oncogenic mutations of ALK kinase in neuroblastoma. *Nature* 2008; 455(7215): 971-4.
- Kato M, Sanada M, Kato I, *et al.* Frequent inactivation of A20 in B-cell lymphomas. *Nature* 2009; 459(7247): 712-6.
- Sanada M, Suzuki T, Shih LY, *et al.* Gain-of-function of mutated C-CBL tumour suppressor in myeloid neoplasms. *Nature* 2009; 460(7257): 904-8.
- Delhommeau F, Dupont S, Della Valle V, *et al.* Mutation in TET2 in myeloid cancers. *N Engl J Med* 2009; 360(22): 2289-301.
- Langemeijer SM, Kuiper RP, Berends M, *et al.* Acquired mutations in TET2 are common in myelodysplastic syndromes. *Nat Genet* 2009; 41(7): 838-42.
- Ernst T, Chase AJ, Score J, *et al.* Inactivating mutations of the histone methyltransferase gene *EZH2* in myeloid disorders. *Nat Genet* 2010; 42(8): 722-6.
- Nikoloski G, Langemeijer SM, Kuiper RP, *et al.* Somatic mutations of the histone methyltransferase gene *EZH2* in myelodysplastic syndromes. *Nat Genet* 2010; 42(8): 665-7.
- Silverman LR, Demakos EP, Peterson BL, *et al.* Randomized controlled trial of azacitidine in patients with the myelodysplastic syndrome: a study of the cancer and leukemia group B. *J Clin Oncol* 2002; 20(10): 2429-40.
- Silverman LR, McKenzie DR, Peterson BL, *et al.* Further analysis of trials with azacitidine in patients with myelodysplastic syndrome: studies 8421, 8921, and 9221 by the Cancer and Leukemia Group B. *J Clin Oncol* 2006; 24(24): 3895-903.
- Kantarjian H, Oki Y, Garcia-Manero G, *et al.* Results of a randomized study of 3 schedules of low-dose decitabine in higher-risk myelodysplastic syndrome and chronic myelomonocytic leukemia. *Blood* 2007; 109(1): 52-7.
- Malcovati L, Germing U, Kuendgen A, *et al.* Time-dependent prognostic scoring system for predicting survival and leukemic evolution in myelodysplastic syndromes. *J Clin Oncol* 2007; 25(23): 3503-10.
- Sole F, Espinet B, Sanz GF, *et al.* Incidence, characterization and prognostic significance of chromosomal abnormalities in 640 patients with primary myelodysplastic syndromes. *Grupo Cooperativo Espanol de Citogenetica Hematologica. Br J Haematol* 2000; 108(2): 346-56.
- Sole F, Luno E, Sanzo C, *et al.* Identification of novel cytogenetic markers with prognostic significance in a series of 968 patients with primary myelodysplastic syndromes. *Haematologica* 2005; 90(9): 1168-78.
- Bernasconi P, Klersy C, Boni M, *et al.* World Health Organization classification in combination with cytogenetic markers improves the prognostic stratification of patients with de novo primary myelodysplastic syndromes. *Br J Haematol* 2007; 137(3): 193-205.
- Rigolin GM, Bigoni R, Milani R, *et al.* Clinical importance of interphase cytogenetics detecting occult chromosome lesions in myelodysplastic syndromes with normal karyotype. *Leukemia* 2001; 15(12): 1841-7.
- Barrett MT, Scheffer A, Ben-Dor A, *et al.* Comparative genomic hybridization using oligonucleotide microarrays and total genomic DNA. *Proc Natl Acad Sci USA* 2004; 101(51): 17765-70.
- Coe BP, Ylstra B, Carvalho B, Meijer GA, Macaulay C, Lam WL. Resolving the resolution of array CGH. *Genomics* 2007; 89(5): 647-53.
- Paulsson K, Heidenblad M, Strombeck B, *et al.* High-resolution genome-wide array-based comparative genome hybridization reveals cryptic chromosome changes in AML and MDS cases with trisomy 8 as the sole cytogenetic aberration. *Leukemia* 2006; 20(5): 840-6.
- O'Keefe CL, Tiu R, Gondek LP, *et al.* High-resolution genomic arrays facilitate detection of novel cryptic chromosomal lesions in myelodysplastic syndromes. *Exp Hematol* 2007; 35(2): 240-51.
- Starczynowski DT, Vercauteren S, Telenius A, *et al.* High-resolution whole genome tiling path array CGH analysis of CD34+ cells from patients with low-risk myelodysplastic syndromes reveals cryptic copy number alterations and predicts overall and leukemia-free survival. *Blood* 2008; 112(8): 3412-24.
- Thiel A, Beier M, Ingenhag D, *et al.* Comprehensive array CGH of normal karyotype myelodysplastic syndromes reveals hidden recurrent and individual genomic copy number alterations with prognostic relevance. *Leukemia* 2011; 25(3): 387-99.
- Matsuzaki H, Dong S, Loi H, *et al.* Genotyping over 100,000 SNPs on a pair of oligonucleotide arrays. *Nat Methods* 2004; 1(2): 109-11.
- Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature* 2007; 447(7145): 661-78.
- Zhao X, Li C, Paez JG, *et al.* An integrated view of copy number and allelic alterations in the cancer genome using single nucleotide polymorphism arrays. *Cancer Res* 2004; 64(9): 3060-71.
- Nannya Y, Sanada M, Nakazaki K, *et al.* A robust algorithm for copy number detection using high-density oligonucleotide single nucleotide polymorphism genotyping arrays. *Cancer Res* 2005; 65(14): 6071-9.

- [37] Van Loo P, Nordgard SH, Lingjaerde OC, *et al.* Allele-specific copy number analysis of tumors. *Proc Natl Acad Sci USA* 2010; 107(39): 16910-5.
- [38] Thoenissen NH, Krug UO, Lee DH, *et al.* Prevalence and prognostic impact of allelic imbalances associated with leukemic transformation of Philadelphia chromosome-negative myeloproliferative neoplasms. *Blood* 2010; 115(14): 2882-90.
- [39] Greenway SC, Pereira AC, Lin JC, *et al.* De novo copy number variants identify new genes and loci in isolated sporadic tetraploidy of Fallopian tube. *Nat Genet* 2009; 41(8): 931-5.
- [40] Murray SS, Oliphant A, Shen R, *et al.* A highly informative SNP linkage panel for human genetic studies. *Nat Methods* 2004; 1(2): 113-7.
- [41] Beroukhi R, Lin M, Park Y, *et al.* Inferring loss-of-heterozygosity from unpaired tumors using high-density oligonucleotide SNP arrays. *PLoS Comput Biol* 2006; 2(5): e41.
- [42] Laframboise T, Harrington D, Weir BA. PLASQ: a generalized linear model-based procedure to determine allelic dosage in cancer cells from SNP array data. *Biostatistics* 2007; 8(2): 323-36.
- [43] Yamamoto G, Nannya Y, Kato M, *et al.* Highly sensitive method for genomewide detection of allelic composition in nonpaired, primary tumor specimens by use of affymetrix single-nucleotide-polymorphism genotyping microarrays. *Am J Hum Genet* 2007; 81(1): 114-26.
- [44] Kralovics R, Passamonti F, Buser AS, *et al.* A gain-of-function mutation of JAK2 in myeloproliferative disorders. *N Engl J Med* 2005; 352(17): 1779-90.
- [45] Gondek LP, Tiu R, O'Keefe CL, Sekeres MA, Theil KS, Maciejewski JP. Chromosomal lesions and uniparental disomy detected by SNP arrays in MDS, MDS/MPD, and MDS-derived AML. *Blood* 2008; 111(3): 1534-42.
- [46] Heinrichs S, Kulkarni RV, Bueso-Ramos CE, *et al.* Accurate detection of uniparental disomy and microdeletions by SNP array analysis in myelodysplastic syndromes with normal cytogenetics. *Leukemia* 2009; 23(9): 1605-13.
- [47] Raghavan M, Smith LL, Lillington DM, *et al.* Segmental uniparental disomy is a commonly acquired genetic event in relapsed acute myeloid leukemia. *Blood* 2008; 112(3): 814-21.
- [48] Akagi T, Shih LY, Kato M, *et al.* Hidden abnormalities and novel classification of t(15; 17) acute promyelocytic leukemia (APL) based on genomic alterations. *Blood* 2009; 113(8): 1741-8.
- [49] Akagi T, Shih LY, Ogawa S, *et al.* Single nucleotide polymorphism genomic arrays analysis of t(8; 21) acute myeloid leukemia cells. *Haematologica* 2009; 94(9): 1301-6.
- [50] Walter MJ, Payton JE, Ries RE, *et al.* Acquired copy number alterations in adult acute myeloid leukemia genomes. *Proc Natl Acad Sci USA* 2009; 106(31): 12950-5.
- [51] Tiu RV, Gondek LP, O'Keefe CL, *et al.* New lesions detected by single nucleotide polymorphism array-based chromosomal analysis have important clinical impact in acute myeloid leukemia. *J Clin Oncol* 2009; 27(31): 5219-26.
- [52] Bullinger L, Kronke J, Schon C, *et al.* Identification of acquired copy number alterations and uniparental disomies in cytogenetically normal acute myeloid leukemia using high-resolution single-nucleotide polymorphism analysis. *Leukemia* 2010; 24(2): 438-49.
- [53] Dunbar AJ, Gondek LP, O'Keefe CL, *et al.* 250K single nucleotide polymorphism array karyotyping identifies acquired uniparental disomy and homozygous mutations, including novel missense substitutions of c-Cbl, in myeloid malignancies. *Cancer Res* 2008; 68(24): 10349-57.
- [54] Fitzgibbon J, Smith LL, Raghavan M, *et al.* Association between acquired uniparental disomy and homozygous gene mutation in acute myeloid leukemias. *Cancer Res* 2005; 65(20): 9152-4.
- [55] Bejar R, Stevenson K, Abdel-Wahab O, *et al.* Clinical effect of point mutations in myelodysplastic syndromes. *N Engl J Med* 2011; 364(26): 2496-506.
- [56] Kosmider O, Gelsi-Boyer V, Ciudad M, *et al.* TET2 gene mutation is a frequent and adverse event in chronic myelomonocytic leukemia. *Haematologica* 2009; 94(12): 1676-81.
- [57] Abdel-Wahab O, Mullally A, Hedvat C, *et al.* Genetic characterization of TET1, TET2, and TET3 alterations in myeloid malignancies. *Blood* 2009; 114(1): 144-7.
- [58] Metzeler KH, Maharry K, Radmacher MD, *et al.* TET2 mutations improve the new European LeukemiaNet risk classification of acute myeloid leukemia: a Cancer and Leukemia Group B study. *J Clin Oncol* 2011; 29(10): 1373-81.
- [59] Smith AE, Mohamedali AM, Kulasekararaj A, *et al.* Next-generation sequencing of the TET2 gene in 355 MDS and CMML patients reveals low-abundance mutant clones with early origins, but indicates no definite prognostic value. *Blood* 2010; 116(19): 3923-32.
- [60] Tahiliani M, Koh KP, Shen Y, *et al.* Conversion of 5-methylcytosine to 5-hydroxymethylcytosine in mammalian DNA by MLL partner TET1. *Science* 2009; 324(5929): 930-5.
- [61] Ito S, D'Alessio AC, Taranova OV, *et al.* Hong K, Sowers LC, Zhang Y. Role of Tet proteins in 5mC to 5hmC conversion, ES-cell self-renewal and inner cell mass specification. *Nature* 2010; 466(7310): 1129-33.
- [62] Wu H, D'Alessio AC, Ito S, *et al.* Dual functions of Tet1 in transcriptional regulation in mouse embryonic stem cells. *Nature* 2011; 473(7347): 389-93.
- [63] Pastor WA, Pape UJ, Huang Y, *et al.* Genome-wide mapping of 5-hydroxymethylcytosine in embryonic stem cells. *Nature* 2011; 473(7347): 394-7.
- [64] Guo JU, Su Y, Zhong C, *et al.* Hydroxylation of 5-methylcytosine by TET1 promotes active DNA demethylation in the adult brain. *Cell* 2011; 145(3): 423-34.
- [65] Williams K, Christensen J, Pedersen MT, *et al.* TET1 and hydroxymethylcytosine in transcription and DNA methylation fidelity. *Nature* 2011; 473(7347): 343-8.
- [66] Quivoron C, Couronne L, Della Valle V, *et al.* TET2 inactivation results in pleiotropic hematopoietic abnormalities in mouse and is a recurrent event during human lymphomagenesis. *Cancer Cell* 2011; 20(1): 25-38.
- [67] Parsons DW, Jones S, Zhang X, *et al.* An integrated genomic analysis of human glioblastoma multiforme. *Science* 2008; 321(5897): 1807-12.
- [68] Mardis ER, Ding L, Dooling DJ, *et al.* Recurring mutations found by sequencing an acute myeloid leukemia genome. *N Engl J Med* 2009; 361(11): 1058-66.
- [69] Marcucci G, Maharry K, Wu YZ, *et al.* IDH1 and IDH2 gene mutations identify novel molecular subsets within de novo cytogenetically normal acute myeloid leukemia: a Cancer and Leukemia Group B study. *J Clin Oncol* 2010; 28(14): 2348-55.
- [70] Paschka P, Schlenk RF, Gaidzik VI, *et al.* IDH1 and IDH2 mutations are frequent genetic alterations in acute myeloid leukemia and confer adverse prognosis in cytogenetically normal acute myeloid leukemia with NPM1 mutation without FLT3 internal tandem duplication. *J Clin Oncol* 2010; 28(22): 3636-43.
- [71] Yoshida K, Sanada M, Kato M, *et al.* A nonsense mutation of IDH1 in myelodysplastic syndromes and related disorders. *Leukemia* 2011; 25(1): 184-6.
- [72] Tefferi A, Lasho TL, Abdel-Wahab O, *et al.* IDH1 and IDH2 mutation studies in 1473 patients with chronic-, fibrotic- or blast-phase essential thrombocythemia, polycythemia vera or myelofibrosis. *Leukemia* 2010; 24(7): 1302-9.
- [73] Pardanani A, Lasho TL, Finke CM, *et al.* IDH1 and IDH2 mutation analysis in chronic- and blast-phase myeloproliferative neoplasms. *Leukemia* 2010; 24(6): 1146-51.
- [74] Figueroa ME, Abdel-Wahab O, Lu C, *et al.* Leukemic IDH1 and IDH2 mutations result in a hypermethylation phenotype, disrupt TET2 function, and impair hematopoietic differentiation. *Cancer Cell* 2010; 18(6): 553-67.
- [75] Grand FH, Hidalgo-Curtis CE, Ernst T, *et al.* Frequent CBL mutations associated with 11q acquired uniparental disomy in myeloproliferative neoplasms. *Blood* 2009; 113(24): 6182-92.
- [76] Loh ML, Sakai DS, Flotho C, *et al.* Mutations in CBL occur frequently in juvenile myelomonocytic leukemia. *Blood* 2009; 114(9): 1859-63.
- [77] Thien CB, Langdon WY. Cbl: many adaptations to regulate protein tyrosine kinases. *Nat Rev Mol Cell Biol* 2001; 2(4): 294-307.
- [78] Perez B, Kosmider O, Cassinat B, *et al.* Genetic typing of CBL, ASXL1, RUNX1, TET2 and JAK2 in juvenile myelomonocytic leukaemia reveals a genetic profile distinct from chronic myelomonocytic leukaemia. *Br J Haematol* 2010; 151(5): 460-8.
- [79] Kohlmann A, Grossmann V, Klein HU, *et al.* Next-generation sequencing technology reveals a characteristic pattern of molecular mutations in 72.8% of chronic myelomonocytic leukemia by detecting frequent alterations in TET2, CBL, RAS, and RUNX1. *J Clin Oncol* 2010; 28(24): 3858-65.

- [80] Niemeyer CM, Kang MW, Shin DH, *et al.* Germline CBL mutations cause developmental abnormalities and predispose to juvenile myelomonocytic leukemia. *Nat Genet* 2010; 42(9): 794-800.
- [81] Martinelli S, De Luca A, Stellacci E, *et al.* Heterozygous germline mutations in the CBL tumor-suppressor gene cause a Noonan syndrome-like phenotype. *Am J Hum Genet* 2010; 87(2): 250-7.
- [82] Ogawa S, Sanada M, Shih LY, *et al.* Gain-of-function c-CBL mutations associated with uniparental disomy of 11q in myeloid neoplasms. *Cell Cycle* 2010; 9(6): 1051-6.
- [83] Sargin B, Choudhary C, Crosetto N, *et al.* Flt3-dependent transformation by inactivating c-Cbl mutations in AML. *Blood* 2007; 110(3): 1004-12.
- [84] Rathinam C, Thien CB, Flavell RA, Langdon WY. Myeloid leukemia development in c-Cbl RING finger mutant mice is dependent on FLT3 signaling. *Cancer Cell* 2010; 18(4): 341-52.
- [85] Valk-Lingbeek ME, Bruggeman SW, van Lohuizen M. Stem cells and cancer; the polycomb connection. *Cell* 2004; 118(4): 409-18.
- [86] Majewski IJ, Ritchie ME, Phipson B, *et al.* Opposing roles of polycomb repressive complexes in hematopoietic stem and progenitor cells. *Blood* 2010; 116(5): 731-9.
- [87] Le Beau MM, Espinosa R, 3rd, Davis EM, *et al.* Eisenbart JD, Larson RA, Green ED. Cytogenetic and molecular delineation of a region of chromosome 7 commonly deleted in malignant myeloid diseases. *Blood* 1996; 88(6): 1930-5.
- [88] Dohner K, Brown J, Lehmann U, *et al.* Molecular cytogenetic characterization of a critical region in bands 7q35-q36 commonly deleted in malignant myeloid disorders. *Blood* 1998; 92(11): 4031-5.
- [89] Wong JC, Zhang Y, Lieu KH, *et al.* Use of chromosome engineering to model a segmental deletion of chromosome band 7q22 found in myeloid malignancies. *Blood* 2010; 115(22): 4524-32.
- [90] List A, Dewald G, Bennett J, *et al.* Lenalidomide in the myelodysplastic syndrome with chromosome 5q deletion. *N Engl J Med* 2006; 355(14): 1456-65.
- [91] Nimer SD. Clinical management of myelodysplastic syndromes with interstitial deletion of chromosome 5q. *J Clin Oncol* 2006; 24(16): 2576-82.
- [92] Le Beau MM, Espinosa R, 3rd, Neuman WL, *et al.* Cytogenetic and molecular delineation of the smallest commonly deleted region of chromosome 5 in malignant myeloid diseases. *Proc Natl Acad Sci USA* 1993; 90(12): 5484-8.
- [93] Boulwood J, Fidler C, Strickson AJ, *et al.* Narrowing and genomic annotation of the commonly deleted region of the 5q- syndrome. *Blood* 2002; 99(12): 4638-41.
- [94] Ebert BL, Pretz J, Bosco J, *et al.* Identification of RPS14 as a 5q-syndrome gene by RNA interference screen. *Nature* 2008; 451(7176): 335-9.
- [95] Barlow JL, Drynan LF, Hewett DR, *et al.* A p53-dependent mechanism underlies macrocytic anemia in a mouse model of human 5q- syndrome. *Nat Med* 2010; 16(1): 59-66.
- [96] Drapchinskaia N, Gustavsson P, Andersson B, *et al.* The gene encoding ribosomal protein S19 is mutated in Diamond-Blackfan anaemia. *Nat Genet* 1999; 21(2): 169-75.
- [97] Boria I, Garelli E, Gazda HT, *et al.* The ribosomal basis of Diamond-Blackfan Anemia: mutation and database update. *Hum Mutat* 2010; 31(12): 1269-79.
- [98] Van den Berghe H, Cassiman JJ, David G, *et al.* Fryns JP, Michaux JL, Sokal G. Distinct haematological disorder with deletion of long arm of no. 5 chromosome. *Nature* 1974; 251(5474): 437-8.
- [99] Tinegate H, Gaunt L, Hamilton PJ. The 5q-syndrome: an underdiagnosed form of macrocytic anaemia. *Br J Haematol* 1983; 54(1): 103-10.
- [100] Starczynowski DT, Kuchenbauer F, Argiropoulos B, *et al.* Identification of miR-145 and miR-146a as mediators of the 5q-syndrome phenotype. *Nat Med* 2010; 16(1): 49-58.
- [101] Tiu RV, Gondek LP, O'Keefe CL, *et al.* Prognostic impact of SNP array karyotyping in myelodysplastic syndromes and related myeloid malignancies. *Blood* 2011; 117(17): 4552-60.
- [102] Shendure J, Ji H. Next-generation DNA sequencing. *Nat Biotechnol* 2008; 26(10): 1135-45.
- [103] Yoshida K, Sanada M, Shiraishi Y, *et al.* Frequent pathway mutations of splicing machinery in myelodysplasia. *Nature* 2011; 478(7367): 64-9.

# An empirical Bayesian framework for somatic mutation detection from cancer genome sequencing data

Yuichi Shiraishi<sup>1,\*</sup>, Yusuke Sato<sup>2,3</sup>, Kenichi Chiba<sup>1</sup>, Yusuke Okuno<sup>2</sup>, Yasunobu Nagata<sup>2</sup>, Kenichi Yoshida<sup>2</sup>, Norio Shiba<sup>2,4</sup>, Yasuhide Hayashi<sup>4</sup>, Haruki Kume<sup>3</sup>, Yukio Homma<sup>3</sup>, Masashi Sanada<sup>2</sup>, Seishi Ogawa<sup>2,\*</sup> and Satoru Miyano<sup>1,\*</sup>

<sup>1</sup>Laboratory of DNA Information Analysis, Human Genome Center, Institute of Medical Science, The University of Tokyo, 4-6-1, Shirokanedai, Minato-ku, Tokyo 108-8639, Japan, <sup>2</sup>Cancer Genomics Project, Graduate School of Medicine, The University of Tokyo, 7-3-1, Hongo, Bunkyo-ku, Tokyo 113-8655, Japan, <sup>3</sup>Department of Urology, Graduate School of Medicine, The University of Tokyo, 7-3-1, Hongo, Bunkyo-ku, Tokyo 113-8655, Japan and <sup>4</sup>Department of Hematology/Oncology, Gunma Children's Medical Center, 779, Shimohakoda, Hokkutsu-machi, Shibukawa, Gunma 377-0061, Japan

Received October 14, 2012; Revised January 25, 2013; Accepted February 10, 2013

## ABSTRACT

Recent advances in high-throughput sequencing technologies have enabled a comprehensive dissection of the cancer genome clarifying a large number of somatic mutations in a wide variety of cancer types. A number of methods have been proposed for mutation calling based on a large amount of sequencing data, which is accomplished in most cases by statistically evaluating the difference in the observed allele frequencies of possible single nucleotide variants between tumours and paired normal samples. However, an accurate detection of mutations remains a challenge under low sequencing depths or tumour contents. To overcome this problem, we propose a novel method, Empirical Bayesian mutation Calling (<https://github.com/friend1ws/EBCall>), for detecting somatic mutations. Unlike previous methods, the proposed method discriminates somatic mutations from sequencing errors based on an empirical Bayesian framework, where the model parameters are estimated using sequencing data from multiple non-paired normal samples. Using 13 whole-exome sequencing data with 87.5–206.3 mean sequencing depths, we demonstrate that our method not only outperforms several existing methods in the calling of mutations with moderate allele frequencies but also enables accurate calling of mutations with

low allele frequencies ( $\leq 10\%$ ) harboured within a minor tumour subpopulation, thus allowing for the deciphering of fine substructures within a tumour specimen.

## INTRODUCTION

Cancer is caused by genetic alterations in which acquired or somatic gene mutations, together with germline factors, play definitive roles in cancer development. As such, comprehensive knowledge regarding somatic mutations in the cancer genome is indispensable for the ultimate understanding of cancer pathogenesis. In this regard, the recent advances in massively parallel sequencing technologies have provided an unprecedented opportunity to decipher a full registry of somatic events in the cancer genome at a single nucleotide resolution (1). However, accurate detection of somatic mutations from high-throughput sequencing data may not always be a straightforward task because ambiguities in short read alignment and sequencing errors are inevitably introduced during sample preparation and signal processing, making it difficult to discriminate true somatic mutations from sequencing errors, especially for those mutations with low sequencing depths or allele frequencies. The detection of low allele frequency mutations is not only required for specimens with low tumour contents but is also important for capturing minor tumour subclones to understand the heterogeneity of cancer (2–5) and the underlying causes of tumour recurrence and therapeutic resistance.

\*To whom correspondence should be addressed. Tel: +81 3 5449 5615; Fax: +81 3 5449 5442; Email: yshira@hgc.jp  
Correspondence may also be addressed to Seishi Ogawa. Tel: +81 3 5800 9045; Fax: +81 3 5800 9047; Email: sogawa-tyk@umin.ac.jp  
Correspondence may also be addressed to Satoru Miyano. Tel: +81 3 5449 5615; Fax: +81 3 5449 5442; Email: miyano@hgc.jp

For calling somatic mutations, each candidate has to be discriminated from germline variants and artifacts appearing from sequencing errors. Although germline variants can be effectively detected by relying on the base calls in paired normal samples, the elimination of sequencing errors may be a more complex task because of uncertain allele frequencies and tumour contents. Most existing approaches have adopted variants whose allele frequencies in tumour samples are significantly higher than those in normal samples, excluding variants whose allele frequencies are high enough to indicate that they are putative germline variants. Sequencing errors can be eliminated to some extent by testing the differences in allele frequencies, as they are expected to occur with equal probability between tumour and normal samples. To measure the significance of the difference in allele frequencies, *SomaticSniper* (6) and *jointSNVmix* (7) estimate the Bayesian posterior probability that tumour and normal samples have different genotypes, whereas our previous approach (8) and *VarScan 2* (9) both rely on the *P*-values from Fisher's exact test.

Although a direct comparison between tumour and normal samples has achieved a measure of success, a more efficient approach to discriminate between sequencing errors and genuine somatic mutations is possible when prior information on sequencing errors is given. In fact, the susceptibility to sequencing errors in each genomic position is not uniform, but there are many common sequencing error-prone sites across different experiments, as shown by several previous studies (10–12) as well as our current study. This implies that, by inferring the susceptibility to sequencing errors at each genomic site, we can achieve greater sensitivity in the detection of somatic mutations at sites with no sequencing errors while efficiently filtering false positives at sequencing error-prone sites (Figure 1).

In this article, we propose a novel statistical approach for the detection of somatic mutations, which explicitly takes into account prior information of sequencing errors. By introducing a Bayesian statistical model, we propose a framework for empirically estimating the distribution of sequencing errors by using a set of non-paired normal samples. Using this approach, we can directly evaluate the discrepancy between the observed allele frequencies and the expected scope of sequencing errors. The proposed approach, which we call Empirical Bayesian mutation Calling (*EBCall*), is superior to several existing methods in calling somatic mutations with moderate allele frequencies. In addition, we demonstrate that *EBCall* can effectively detect a series of somatic mutations that have allele frequencies of <10% with a high degree of accuracy, thereby identifying sub-clonal structures of cancer cells that cannot otherwise be found.

## MATERIALS AND METHODS

### Patient samples and sequencing procedures

After receiving informed consent, paired tumour-normal samples were obtained from 20 patients with clear cell

renal cell carcinoma (ccRCC) by sampling their specimens during surgical operations. Of the samples obtained, 13 paired tumour-normal samples were used for a performance evaluation of the mutation detection, and all 20 of the normal samples were used for estimating the sequencing errors as non-paired normal reference samples. In addition, to compare the choice of normal reference samples, 20 normal samples collected from patients with paediatric acute myeloid leukemia (ped-AML) were also used; the informed consent for these sample collections were obtained from the patients' parents. This study was approved by the ethics committees of the University of Tokyo and Gunma Children's Medical Center.

Genomic DNA and total RNA were extracted from the samples using QIAamp DNA Investigator kit (Qiagen) and the RNeasy Total RNA kit (Qiagen) with DNase treatment, respectively, according to the manufacturers' protocols. For whole-exome sequencing, SureSelect-enriched exon fragments were subjected to sequencing using HiSeq 2000, as previously described (8). The ccRCC samples were sequenced from October 2011 to February 2012, whereas the ped-AML samples were sequenced from April 2012 to June 2012. For 10 ccRCC samples, whole-genome sequencing and RNA sequencing were performed using HiSeq 2000, according to standard protocols recommended by Illumina. The mean sequencing depth for each sample was 65.9–223.0 (Supplementary Table S1 and S2).

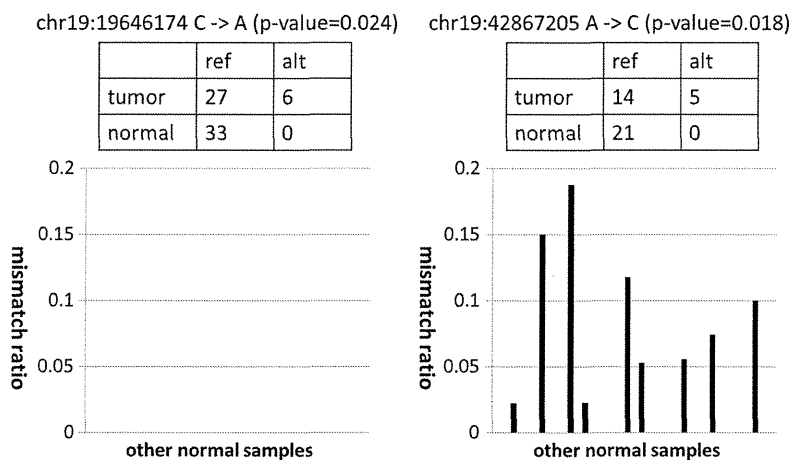
### Outline of the mutation calling method

The outline of *EBCall* is shown in Figure 2. The key concept in *EBCall* is that sequencing data of multiple non-paired normal samples are used to estimate possible sequencing errors at each genomic site. For this purpose, we modelled the sequencing errors that follow a Beta-binomial distribution, the parameters of which were estimated using the sequencing data from multiple non-paired normal samples (Figure 3). The allele frequencies of the observed variants in the tumour DNA were then compared with the inferred sequencing error distribution at the corresponding genomic positions to exclude sequencing errors. Germline Single Nucleotide Polymorphism (SNPs) were eliminated using sequencing data from the paired normal DNA.

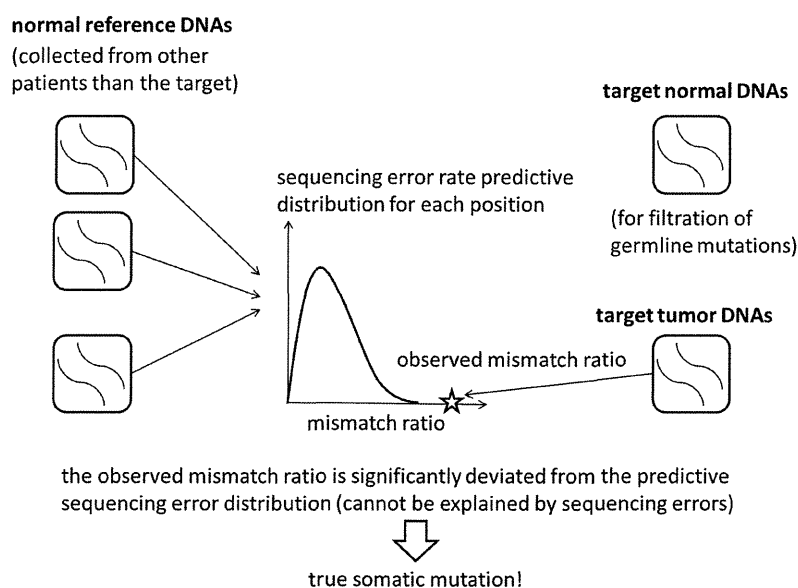
### Alignment of sequencing data

The sequencing reads were aligned to NCBI Human Reference Genome Build 37 using Burrows-Wheeler Aligner, version 0.5.8 (13) with the default parameter settings. Polymerase chain reaction (PCR) duplications were eliminated using Picard (<http://picard.sourceforge.net/>). Low-quality reads showing >5 mismatches with the reference genome or those whose mapping quality was <30 were excluded from further analysis as we did in (8).

For RNA sequencing data, a two-step alignment strategy adopted in *Genomon-fusion* (under submission) was used, in which all sequence reads were first aligned to the known transcript sequences (UCSC known genes)



**Figure 1.** Examples of mismatch ratios of other normal samples for mutation candidates with moderate  $P$ -values. In both cases, although the mismatch ratios of the target tumour sample were relatively high, the numbers of corresponding supporting variant reads were small. For the candidate on the left, the frequencies of non-reference alleles for other normal samples were consistently zero. Therefore, this supports the prediction that the observed variant reads in the target tumour sample came from a true somatic mutation and not from sequencing errors. On the other hand, for the candidate on the right, we often observed high frequencies of non-reference alleles for several different normal samples. Therefore, the observed variant reads in the target tumour sample likely came from sequencing errors, and it was just by chance that there was no variant read in the target normal sample.



**Figure 2.** An illustrative description of the proposed method. For each genomic site, the distribution of sequencing errors is estimated using non-paired normal samples from patients other than the target. The mismatch ratio of the target tumour sample is then compared with the distribution. If the mismatch ratio deviates significantly from the distribution, the corresponding variant is then extracted as a somatic mutation candidate. The target normal sample is used for filtering germline mutations.

using bowtie (14), and the non-aligned reads were then aligned to the genome sequences using blat (15). For the whole-genome sequencing data, all reads were aligned using blat.

#### Definition of variables

Let  $\Omega$  be an entire set of possible nucleotide variations consisting of combinations of genomic positions and

types of nucleotide changes (e.g. chr1:5, C > A or chr20:10 000, A > AAG). Because sequencing errors are often biased to one strand (6,9,16), the number of total ( $d$ ) and variant reads ( $x$ ) for a given variant,  $v \in \Omega$ , were enumerated for each strand separately to distinguish between short reads aligned with the positive ( $x_{a,v,+}$ ,  $d_{a,v,+}$ ) and negative ( $x_{a,v,-}$ ,  $d_{a,v,-}$ ) strands, respectively, where  $a$  denotes the type of sample, which is either