# DNA sequencing of a cytogenetically normal acute myeloid leukemia genome

**Timothy J Ley**[1,2,3,4,*], **Elaine R Mardis**[2,3,*], **Li Ding**[2,3], **Bob Fulton**[3], **Michael D McLellan**[3], **Ken Chen**[3], **David Dooling**[3], **Brian H Dunford-Shore**[3], **Sean McGrath**[3], **Matthew Hickenbotham**[3], **Lisa Cook**[3], **Rachel Abbott**[3], **David E Larson**[3], **Dan C Koboldt**[3], **Craig Pohl**[3], **Scott Smith**[3], **Amy Hawkins**[3], **Scott Abbott**[3], **Devin Locke**[3], **LaDeana W Hillier**[5], **Tracie Miner**[3], **Lucinda Fulton**[3], **Vincent Magrini**[2,3], **Todd Wylie**[3], **Jarret Glasscock**[3], **Joshua Conyers**[3], **Nathan Sander**[3], **Xiaoqi Shi**[3], **John R Osborne**[3], **Patrick Minx**[3], **David Gordon**[5], **Asif Chinwalla**[3], **Yu Zhao**[1], **Rhonda E Ries**[1], **Jacqueline E Payton**[6], **Peter Westervelt**[1,4], **Michael H Tomasson**[1,4], **Mark Watson**[3,4,6], **Jack Baty**[7], **Jennifer Ivanovich**[4,8], **Sharon Heath**[1,4], **William D Shannon**[1,4], **Rakesh Nagarajan**[4,6], **Matthew J Walter**[1,4], **Daniel C Link**[1,4], **Timothy A Graubert**[1,4], **John F DiPersio**[1,4], and **Richard K Wilson**[2,3,4]

[1] Department of Medicine, Washington University School of Medicine, Washington University School of Medicine, St. Louis, MO

[2] Department of Genetics, Washington University School of Medicine, Washington University School of Medicine, St. Louis, MO

[3] The Genome Center at Washington University, Washington University School of Medicine, St. Louis, MO

[4] Siteman Cancer Center, Washington University School of Medicine, St. Louis, MO

[5] Department of Genome Sciences, University of Washington, Seattle WA

[6] Department of Pathology and Immunology, Washington University School of Medicine, St. Louis, MO

[7] Division of Biostatistics, Washington University School of Medicine, St. Louis, MO

[8] Department of Surgery, Washington University School of Medicine, St. Louis, MO

## Lay Summary

Acute myeloid leukemia is a highly malignant hematopoietic tumor that affects about 13,000 adults yearly in the United States. The treatment of this disease has changed little in the past two decades, since most of the genetic events that initiate the disease remain undiscovered. Whole genome sequencing is now possible at a reasonable cost and timeframe to utilize this approach for

unbiased discovery of tumor-specific somatic mutations that alter the protein-coding genes. Here we show the results obtained by sequencing a typical acute myeloid leukemia genome and its matched normal counterpart, obtained from the patient's skin. We discovered 10 genes with acquired mutations; two were previously described mutations thought to contribute to tumor progression, and 8 were novel mutations present in virtually all tumor cells at presentation and relapse, whose function is not yet known. Our study establishes whole genome sequencing as an unbiased method for discovering initiating mutations in cancer genomes, and for identifying novel genes that may respond to targeted therapies.

We used massively parallel sequencing technology to sequence the genomic DNA of tumor and normal skin cells obtained from a patient with a typical presentation of FAB M1 Acute Myeloid Leukemia (AML) with normal cytogenetics. 32.7-fold 'haploid' coverage (98 billion bases) was obtained for the tumor genome, and 13.9-fold coverage (41.8 billion bases) was obtained for the normal sample. Of 2,647,695 well-supported Single Nucleotide Variants (SNVs) found in the tumor genome, 2,588,486 (97.7%) also were detected in the patient's skin genome, limiting the number of variants that required further study. For the purposes of this initial study, we restricted our downstream analysis to the coding sequences of annotated genes: we found only eight heterozygous, non-synonymous somatic SNVs in the entire genome. All were novel, including mutations in protocadherin/cadherin family members (CDH24 and PCLKC), G-protein coupled receptors (GPR123 and EBI2), a protein phosphatase (PTPRT), a potential guanine nucleotide exchange factor (KNDC1), a peptide/drug transporter (SLC15A1), and a glutamate receptor gene (GRINL1B). We also detected previously described, recurrent somatic insertions in the *FLT3* and *NPM1* genes. Based on deep readcount data, we determined that all of these mutations (except FLT3) were present in nearly all tumor cells at presentation, and again at relapse 11 months later, suggesting that the patient had a single dominant clone containing all of the mutations. These results demonstrate the power of whole genome sequencing to discover novel cancer-associated mutations.

## INTRODUCTION

Acute Myeloid Leukemia refers to a group of clonal hematopoietic malignancies that predominantly affect middle-aged and elderly adults. An estimated 13,000 people will develop AML in the United States in 2008, and 8,800 will die from it1. Although the life expectancy from this disease has increased slowly over the past decade, the improvement is due predominantly to improvements in supportive care, not in the drugs or approaches used to treat patients.

For most patients with a 'sporadic' presentation of AML, it is not yet clear whether inherited susceptibility alleles play a role in pathogenesis2. Further, the nature of the initiating or progression mutations is for the most part unknown3. Recent attempts to identify additional progression mutations by extensively re-sequencing tyrosine kinase genes yielded very few new mutations, and most were not recurrent4,5. Expression profiling studies have yielded signatures that correlate with specific cytogenetic subtypes of AML, but have not yet suggested new initiating mutations6–8. Recent studies using array-based comparative genomic hybridization and/or SNP arrays, while identifying important gene mutations in acute lymphoblastic leukemia9,10 have revealed very few recurrent submicroscopic somatic

copy number variants in AML (Walter, MJ, manuscript in preparation, and references 11–13). Together, these studies suggest that *we have not yet discovered most of the relevant mutations that contribute to the pathogenesis of AML.* We therefore believe that unbiased whole genome sequencing will be required to identify most of these mutations. Until recently, this approach has not been feasible, due to the high cost of conventional capillary-based approaches, and the large numbers of primary tumor cells required to yield the necessary genomic DNA. "Next-generation" sequencing approaches, however, have changed this landscape.

Our group has pioneered the use of whole genome re-sequencing and variant discovery approaches using the Illumina/Solexa technology with the genome of the nematode worm *C. elegans* as a proof-of-principle14. This approach has distinct advantages in reduced cost, dramatically increased data production rate, and a low input requirement of DNA for library construction. In the present study, we used a similar approach to sequence the tumor genome of a single AML patient, and the matched normal genome (derived from a skin biopsy) of the same patient. Following alignment to the human reference genome, sequence variants were discovered in the tumor genome and compared to the patient's normal sequence, to dbSNP, and to variants recently reported for two additional human genomes15,16; revealing novel single nucleotide and small insertion/deletion ("indel") variants genome-wide. Novel somatic mutations were detected in genes not previously implicated in AML pathogenesis, demonstrating the need for unbiased whole genome approaches to discover all mutations associated with cancer pathogenesis.

## RESULTS

### Rationale for selecting the FAB M1 AML subtype for sequencing

Of the 8 FAB subtypes of AML, M1 AML is one of the most common (~20% of all cases). No specific cytogenetic abnormalities or somatic initiating mutations have been identified for this subtype; in fact, about half of the patients with *de novo* M1 AML have normal cytogenetics17–19. The frequency of well-described progression mutations (e.g. activating alleles of FLT3, c-Kit, and Ras) is similar to that of other common FAB subtypes5. We therefore decided to sequence the genome of tumor cells derived from a patient with M1 AML, since so little is known about the molecular pathogenesis of this common subtype. The criteria used to select the sample are outlined in Supplementary Materials.

### Case presentation of UPN 933124

The case presentation is described in detail in the Supplementary Materials. Briefly, a previously healthy woman in her mid-50s presented suddenly with fatigue and easily bruisability, and was found to have a peripheral WBC count of 105,000 cells per microliter, with 85% myeloblasts. A bone marrow examination revealed 100% myeloblasts with morphologic features and cell surface markers consistent with FAB M1 AML (Supplementary Figure 1). Cytogenetic analysis of tumor cells revealed a normal 46 XX karyotype. Although the patient experienced a complete remission with conventional therapies, she relapsed at 11 months and expired 24 months after her initial diagnosis was

made. Informed consent for whole genome sequencing was subsequently obtained from her next of kin.

### An essentially diploid tumor genome with an expression profile typical of M1 AML

The tumor sample from patient 933124 contained no somatic copy number changes at a resolution of ~5 kb (further confirmed on the NimbleGen 2.1M array platform, data not shown), and no evidence of copy number neutral LOH, indicating that the genome was essentially diploid at this level of resolution (see Supplementary Figure 2). Further analysis of the 933124-derived tumor and skin samples revealed 26 inherited copy number variants (i.e. detected in both the tumor and skin samples). All but two of these previously had been reported in the Database of Genomic Variants (see Supplementary Table 1). All of the CNVs detected in this genome were found in at least one other AML patient (89 additional cases, mostly Caucasian, have been queried using the same SNP array platform), and all but one were found in at least one of the 160 Caucasian HapMap and Coriell samples that were studied on the same array platform (Supplementary Table 1).

To determine whether the tumor cells of 933124 were typical of M1 AML, we compared the expression signatures of 111 *de novo* AML cases using unsupervised clustering (Ward's method, see Materials and Methods). 933124 clustered with multiple other M1 (and M2) AML cases with normal cytogenetics, suggesting that the genetic events underlying the pathogenesis of this case are similar to those of other cases exhibiting normal cytogenetics (Supplementary Fig. 3).

### Coverage depth of the tumor and skin genomes

Since most of the acquired mutations in cancer genomes have been shown to be heterozygous, the complete sequencing of a cancer genome requires the detection of *both* alleles at most positions in the genome[20]. We therefore designed sequence coverage metrics to define the point at which 90% diploid coverage had been reached. To minimize errors associated with any single platform or measurement, diploid coverage for this genome was assessed using a set of High-Quality (HQ) SNPs derived from two different SNP array platforms, Affymetrix 6.0 and Illumina Infinium 550K. For a SNP to be included in the HQ set, the following criteria had to be satisfied: (1) identical genotypes were called from both assays at the same genomic positions and (2) the resulting genotype was heterozygous. For the 933124 tumor genome, 46,494 heterozygous SNPs passed the above criteria and were defined as HQ SNPs. 46,572 HQ SNPs were defined for the skin sample.

We performed 98 full runs on the Illumina Genome Analyzer to achieve the targeted level of 90% diploid coverage as determined by coverage of the HQ SNP set. Maq[21] was used to perform alignment, determine consensus, and identify SNVs within the 98 billion bases generated from the tumor genome (see Table 1). Maq predicted a total of 3.81M SNVs (Maq SNP Quality    15) in the tumor genome, including matching heterozygous genotypes for 91.2% of the 46,494 HQ SNPs. When we lowered the Maq SNP Quality cutoff to 0, 94.06% HQ SNPs were predicted. Further investigation of Maq alignments revealed coverage for both alleles at an additional 5.38% of the HQ SNPs, but Maq did not predict a SNP or matching heterozygous genotype due to insufficient depth or quality of coverage. Additional

analysis revealed coverage at 46,484 of 46,494 HQ SNPs for at least one allele (i.e., 99.98% haploid coverage for the tumor genome).

We sequenced the genome of normal skin cells from the same patient to enable the identification of inherited sequence variants in the tumor genome. Our targeted diploid coverage goal for the skin-derived genome was 80%. We achieved this goal with only 34 Solexa runs (41.8 billion bases), utilizing improved reagents and longer read lengths to attain 82.6% diploid and 84.2% haploid coverage (Table 1).

To begin evaluating the quantity and quality of the detected sequence variants in the tumor and skin genomes, we compared the overlap and uniqueness of this genome's variants with respect to the Watson and Venter genomes, and to dbSNP (v127) (Figure 1). Of the 3.68 M single nucleotide variants (SNVs, Maq SNP Quality   15, excluding SNVs found on chromosome X) predicted by Maq in the tumor genome, 2.36 M were present in dbSNP, 2.36 M were detected in the skin genome (Fig. 1A), 1.50 M were detected in the Venter genome, and 1.58 M were found in the Watson genome (Fig. 1B). Ultimately, 1.70 M SNVs were unique to the 933124 tumor genome. Upon filtering the 933124 SNVs at different Maq quality values to determine the stability of results, we observed that the proportion of 933124 SNVs that also are in dbSNP increases from 63.9% to 69.48% when the Maq quality threshold score increases from 15 to 30, as expected.

### Refining variant detection to detect potential somatic mutations

Because the number of sequence variants initially detected by Maq was high, we developed improved filtering tools to effectively separate true variants from false positives. To this end, we generated an experimental dataset by re-sequencing Maq-predicted SNVs, randomly selecting a training subset and a test dataset, whose annotations and features were submitted to Decision Tree C4.522. This approach identified parameters that separated true variants from false positives, revealing that SNV-supporting read counts (unique based on read start position, and base position in supporting reads), base quality, and Maq quality scores are major determinants for identifying false positives. Implementing rules obtained from the Decision Tree analysis resulted in 91.9% sensitivity and 83.5% specificity for validated SNVs.

### Identification of somatic mutations in coding sequences

The patient had 3,813,205 sequence variants in her tumor genome, as defined by Maq scores of >15 (Table 1). Of these, 2,647,695 were supported by the Decision Tree analysis in the tumor genome, of which 2,584,418 (97.7%) also were detected in the skin genome (Figure 2). The detailed algorithm for selecting putative somatic variants is described in Supplemental Materials. Most of the 63,277 tumor-specific variants we detected were either present in dbSNP or were previously described in the Watson or Venter genomes (31,645), or occurred in non-genic regions (20,440). A total of 11,192 variants were located within the boundaries of annotated genes; 216 of these variants were in untranslated regions, and 10,735 were in introns (but not involving splice junctions) and were not explored further in our analysis. Of the coding sequence variants, 60 were synonymous, and not further evaluated. The remaining 181 variants were either non-synonymous, or were predicted to

alter splice site function. By sequencing PCR-generated amplicons from the tumor and skin samples (and also from the relapse tumor sample obtained 11 months after the original presentation), we determined that 152 of these variants were false positive (i.e. wild type) calls, 14 were inherited SNPs, and eight were somatic mutations in both the original tumor and the relapse sample (Table 2). Seven variants could not be validated, either because the regions involved were repetitive, or because all attempts to obtain PCR amplicons failed. All of the PCR-amplified exons from the eight genes containing validated somatic mutations were sequenced in 187 additional cases of AML using samples from our discovery and validation sets23; no additional somatic mutations were detected in these genes (data not shown). A description of how we estimated the false negative (12.45%) and false positive (0.06%) rates for SNVs over the entire genome is presented in Supplementary Materials. Using these estimates, we can predict that very few somatic, non-synonymous variants were missed by our analysis of this deeply covered genome.

### Defining mutation frequencies in the tumor sample

To better define the percentage of tumor cells that contained each of the discovered somatic mutations, we amplified each mutation-containing locus from non-amplified genomic DNA derived from the de novo and relapse tumor samples, and from the skin biopsy obtained at presentation. The resulting amplicons were sequenced using the Roche/454 FLX platform, and the frequency of reads containing the reference and variant alleles were defined (Figure 3, and Table 3). Control amplicons containing a known heterozygous SNP in BRCA2 (encoding N372H) and a homozygous SNP in TP53 (encoding P72R) were analyzed similarly. The BRCA2 SNP yielded ~50% variant frequencies in the tumor and skin samples, while nearly 100% of the TP53 alleles were variant in all three samples, as expected. Remarkably, all eight somatic SNVs were detected at ~50% frequencies in the primary tumor sample (100% blasts), and at ~40% frequencies in the relapse sample (78% blasts; if the variant frequencies are corrected for blast counts [i.e. multiplied by 1.28], the frequencies at relapse also were ~50%). The NPMc mutation also was detected at a frequency of ~50%, but the FLT3 ITD allele was detected in only 35.1% of the 454 reads at diagnosis, and 31.3% at relapse, suggesting that the mutation was not present in all tumor cells at diagnosis or relapse.

Surprisingly, the variant alleles also were detected at frequencies of ~5–10% in the skin sample (the PCLKC gene was the single outlier, with 23% variant reads in the skin sample for unclear reasons). In retrospect, it is clear that the skin sample contained contaminating leukemic cells, since the patient's WBC count at presentation was 105,000 per microliter, with 85% blasts. This information was used to inform the Decision Tree analysis described above: we allowed high quality tumor variants to move forward in the discovery pipeline if they were detected at a low frequency (two or fewer reads) in the skin sample (as defined by a binomial test).

### Detecting insertions and deletions (indels)

To discover small indels (<6 bp) from sequence reads (32–35 bp long), we started with a set of 236 million reads that were not confidently aligned by Maq to the reference genome. We applied Cross_Match and BLAT to identify gapped alignments that are unique in the

genome. To detect indels longer than 6 bp, we developed a novel "split reads" algorithm (see Supplementary Materials) that aligns sub-segments of reads independently to the genome, and computes a mapping quality for the derived gapped alignment based on the number of hits and the quality of the bases. These efforts resulted in the identification of 726 putative small indels (1 to 30 bp in size) that occur in coding exons, 393 of which (54.2%) were found in dbSNP. After manual review, we selected a set of 28 putative somatic coding indels for validation using PCR-based dye terminator sequencing. Of these putative indels, 22 were validated, but were found present in both tumor and skin (15 of these were in dbSNP), 2 were false positive calls, 2 had no coverage, and two were previously validated somatic insertions in *NPM1* (4 bp) and *FLT3* (30 bp).

## DISCUSSION

In this report, we describe the sequencing and analysis of a primary human cancer genome using next-generation sequencing technology. Our patient's tumor genome was essentially diploid, and contained ten non-synonymous somatic mutations that may be relevant for her disease. These mutations affect genes participating in several well-described pathways that are known to contribute to cancer pathogenesis, but most of these genes would not have been candidates for directed re-sequencing based on our current understanding of cancer. Hence, these results justify the use of next-generation whole genome sequencing approaches to reveal somatic mutations in cancer genomes.

As we demonstrated in our re-sequencing of the genome of the *C. elegans* N2 Bristol strain14, and again in this study, massively-parallel short-read sequencing provides an effective method for examining single nucleotide and short indel variants by comparison of the aligned reads to a reference genome sequence. By sequencing our patient's tumor genome to a depth of >30-fold coverage, and gauging our ability to detect known heterozygous positions across the genome, we have produced a sufficient depth and breadth of sequence coverage to comprehensively discover somatic genome variants. A slightly lower coverage of the normal genome from this individual helped to identify nearly 98% of potential variants as being inherited, a critical filter that allowed us to more readily identify the true somatic mutations in this tumor. Our results strongly support the notion that hypothesis-driven (e.g. candidate gene-based) examination of tumor genomes by PCR-directed or capture-based methods is inherently limited, and will miss key mutations. An additional and important consideration is the demand for large amounts of genomic DNA by these techniques; this is a serious limitation when precious clinical samples are being studied. The Illumina/Solexa technology requires only ~1 ug of DNA per library, enabling the study of primary tumor DNA rather than requiring the use of tumor cell lines, which may contain genetic changes and adaptations required for immortalization and maintenance in tissue culture conditions.

A total of 10 non-synonymous somatic mutations were identified in this patient's tumor genome. Two are well known AML-associated mutations, including an internal tandem duplication of the FLT3 receptor tyrosine kinase gene, which constitutively activates kinase signaling, and portends a poor prognosis5,24,25, and a four base insertion in exon 12 of the NPM1 gene (NPMc)26–28. Both of these mutations are common (25–30%) in AML tumors,

and both are thought to contribute to progression of the disease, rather than to cause it directly29. Interestingly, the frequency of the mutant FLT3 allele in the primary and relapse tumor samples (35.08% and 31.30%, respectively) was significantly less than that of the other 9 mutations (p<0.000001 for both the primary and relapse samples). These data suggest that the FLT3 ITD may not have been present in all tumor cells, and further, that it may have been the last mutation acquired.

The other eight somatic mutations that we detected are all single base changes, and none has previously been detected in an AML genome. Four of the genes affected, however, are in gene families that are strongly associated with cancer pathogenesis (including PTPRT, CDH24, PCLKC, and SLC15A1). The other four somatic mutations occurred in genes not previously implicated in cancer pathogenesis, but whose potential functions in metabolic pathways suggest mechanisms by which they could act to promote cancer (including KNDC1, GPR123, EBI2, and GRINL1B). We speculate regarding the roles of these mutations for the pathogenesis of this patient's disease in Supplementary Materials.

The importance of the eight newly defined somatic mutations for AML pathogenesis is not yet known, and will require functional validation studies in tissue culture cells and mouse models to assess their relevance. Even though we could not detect recurrent mutations in the limited AML sample set we surveyed, several lines of evidence suggest that these mutations may not be random, "passenger" mutations, as follows:

1.  Somatic mutations in this genome are extremely rare. The rarity of somatic variants, and the normal diploid structure of the tumor genome, argues strongly against genetic instability or DNA repair defects in this tumor. Conceptually, this result is further supported by the very small number of novel somatic mutations discovered in the expressed tyrosine kinases of AML samples4,5; genetic instability does not appear to be a general feature of AML genomes.

2.  Based on the equivalent frequencies of the variant and wild type alleles for the mutations in the tumor genome (except for FLT3 ITD), it is highly likely that all the mutations are heterozygous, and are present in virtually all of the tumor cells (Figure 3). The latter suggests that these mutations all may have been selected for and retained because they are important for disease pathogenesis in this patient. Alternatively, all may have occurred simultaneously in the same leukemia-initiating cell, but only a subset of the mutations (or an as-yet undetected mutation) is truly important for pathogenesis (i.e. disease "drivers" vs. passengers). Although we suggest that the latter scenario is very unlikely based on our current understanding of tumor progression, many more AML genomes will need to be sequenced to resolve this issue.

3.  The same mutations were detected in the tumor cells in the relapse sample at approximately the same frequencies as in the primary sample. All of these mutations were therefore present in the resistant tumor cells that contributed to the patient's relapse, further suggesting that a single clone contains all 10 mutations.

4.  Seven of the 10 genes containing somatic mutations were detectably expressed in the tumor sample. FLT3 and NPM1 mRNAs were highly expressed in this tumor

sample, as they are in virtually all AML samples. We detected mRNA from the CDH24, SLC15A1, and EBI2 genes on the Affymetrix expression array, while expression of GRINL1B and PCLKC were detected by RT-PCR (data not shown). Expression of KNDC1, PTPRT, and GPR123 was not detected by either approach, but we cannot rule out expression of these genes in a small subset of tumor cells (e.g. leukemia initiating cells).

5.  For the five point mutations where data are available, the mutated base is highly conserved across multiple species (Table 2).

Although we performed whole genome sequencing on this cancer sample, we restricted our initial validation studies to the 1–2% of the genome that encodes genes. This raises the issue of whether sequencing the cDNA transcriptome of this tumor would have been a faster, cheaper, and more efficient way of finding the mutations. While this approach will undoubtedly be an important adjunct to whole genome sequencing, there are several advantages to the approach we used:

1.  Coverage models for whole genome libraries are currently better understood than for cDNA libraries, where transcript abundance can vary over many orders of magnitude.

2.  Even if the transcriptome had been sequenced, extensive characterization of the normal genome would have been required to distinguish inherited variants from somatic mutations.

3.  Relevant non-synonymous mutations could be missed by cDNA sequencing, including mutations that result in RNA instability (splice variants, nonsense mutations, etc.), and/or mutations in genes expressed at low levels, or in only a small subset of tumor cells.

The additional non-coding and non-genic somatic variants in this genome (which we currently estimate at 500–1000, based on our assessment of false positive and negative rates for non-synonymous mutations), which will be fully described later, will provide a rich source of potentially relevant sequence changes that will be better understood as more cancer genomes are sequenced.

In summary, we have successfully used a next-generation whole genome sequencing approach to identify new candidate genes that may be relevant for AML pathogenesis. We cannot overemphasize the importance of parallel sequencing of the patient's normal genome to determine which variants were inherited; the identification of the true somatic mutations in this tumor genome would not have been feasible without this approach. Furthermore, until hundreds (or perhaps thousands) of normal genomes and additional AML tumors are sequenced, the contextual relevance of the mutations found in this genome will be unknown. Regardless, the somatic mutations that we did find were neither predicted by the curation of previously defined cancer genes, nor by the study of this tumor using unbiased, high-resolution array-based genomic approaches. For AML and other types of cancer, whole genome sequencing may therefore be the only effective means for discovering all of the mutations that are relevant for pathogenesis.

## Methods Summary

Sequence end reads (average length for tumor genome, 32 bp, and for skin, 35 bp) were generated from Illumina/Solexa fragment libraries derived from the tumor or skin cells of patient 933124, using the Illumina Genome Analyzer. The analyzed reads were aligned to the human reference (NCBI Build 36) using Maq21. Coverage of the tumor and normal genomes was ascertained by comparison to the patient's heterozygous SNPs, established by compiling shared SNP calls monitored on the Affymetrix 6.0 and Illumina Infinium 550K genotyping platforms. We examined the Maq alignments by Decision Tree analysis to discover single nucleotide variants, as well as to identify copy number variants. Nonaligned reads were further analyzed for indel discovery. For all putative variants, we attempted validation using custom PCR and capillary sequencing on the ABI 3730 platform. All validated somatic mutations were further analyzed by Roche/454 sequencing of PCR-generated amplicons made from primary genomic DNA to compare readcounts of wild-type and mutant alleles in the primary tumor, skin, and relapse tumor samples. A complete description of the AML case sequenced, and the Materials and Methods used to generate this dataset, are provided in the Supplementary Materials.

## Supplementary Material

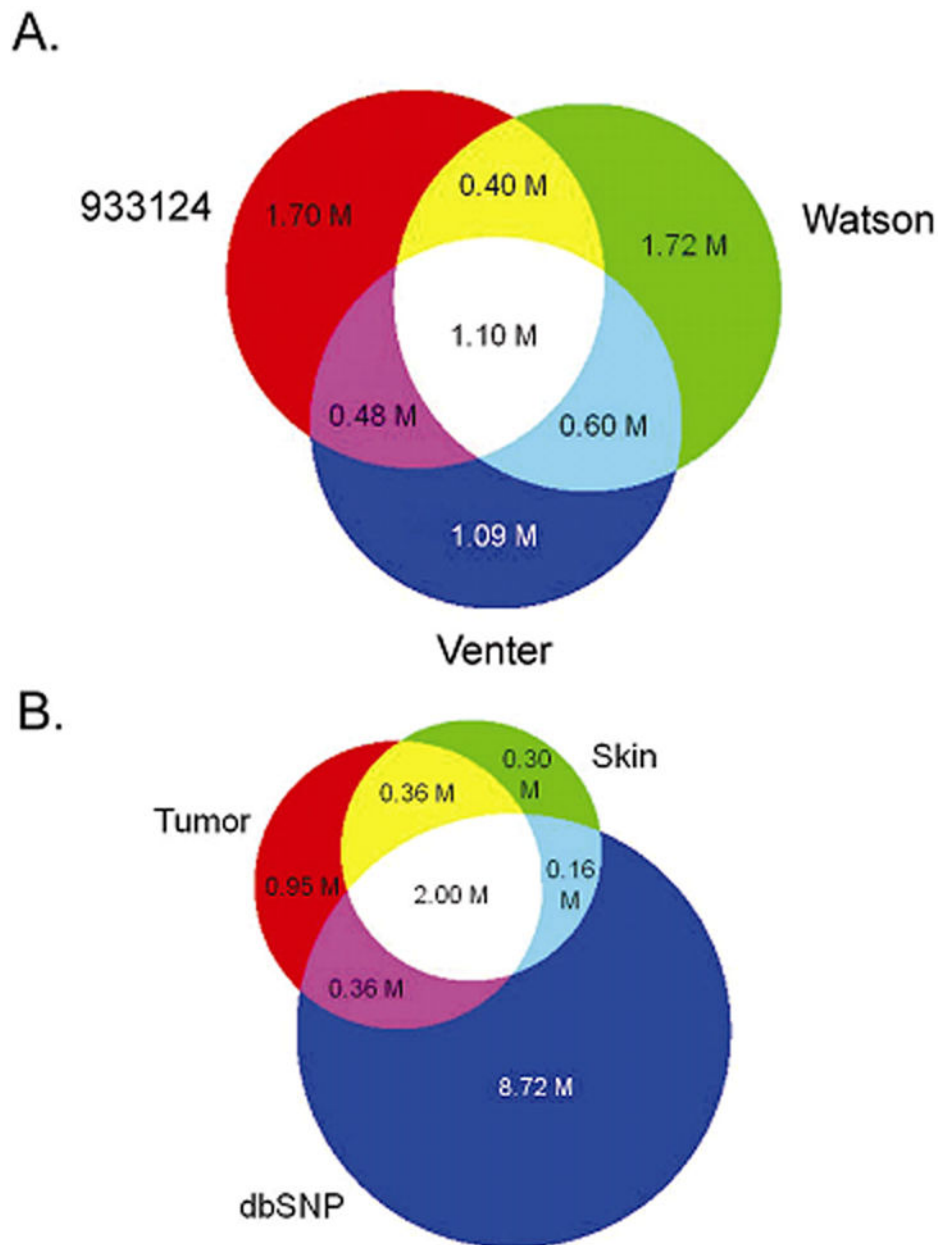Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

## References

1. Jemal A, et al. Cancer statistics, 2008. CA Cancer J Clin. 2008; 58:71–96. [PubMed: 18287387]

2. Owen C, Barnett M, Fitzgibbon J. Familial myelodysplasia and acute myeloid leukaemia--a review. Br J Haematol. 2008; 140:123–32. [PubMed: 18173751]

3. Mrozek K, Marcucci G, Paschka P, Whitman SP, Bloomfield CD. Clinical relevance of mutations and gene-expression changes in adult acute myeloid leukemia with normal cytogenetics: are we ready for a prognostically prioritized molecular classification? Blood. 2007; 109:431–48. [PubMed: 16960150]

4. Loriaux MM, et al. High-throughput sequence analysis of the tyrosine kinome in acute myeloid leukemia. Blood. 2008; 111:4788–96. [PubMed: 18252861]

5. Tomasson MH, et al. Somatic mutations and germline sequence variants in the expressed tyrosine kinase genes of patients with de novo acute myeloid leukemia. Blood. 2008; 111:4797–808. [PubMed: 18270328]

6. Schoch C, et al. Acute myeloid leukemias with reciprocal rearrangements can be distinguished by specific gene expression profiles. Proc Natl Acad Sci U S A. 2002; 99:10008–13. [PubMed: 12105272]

7. Bullinger L, et al. Use of gene-expression profiling to identify prognostic subclasses in adult acute myeloid leukemia. N Engl J Med. 2004; 350:1605–16. [PubMed: 15084693]

8. Valk PJ, et al. Prognostically useful gene-expression profiles in acute myeloid leukemia. N Engl J Med. 2004; 350:1617–28. [PubMed: 15084694]

9. Mullighan CG, et al. Genome-wide analysis of genetic alterations in acute lymphoblastic leukaemia. Nature. 2007; 446:758–64. [PubMed: 17344859]

10. Mullighan CG, et al. BCR-ABL1 lymphoblastic leukaemia is characterized by the deletion of Ikaros. Nature. 2008; 453:110–4. [PubMed: 18408710]

11. Raghavan M, et al. Genome-wide single nucleotide polymorphism analysis reveals frequent partial uniparental disomy due to somatic recombination in acute myeloid leukemias. Cancer Res. 2005; 65:375–8. [PubMed: 15695375]

12. Paulsson K, 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; 20:840–6. [PubMed: 16498392]

13. Rucker FG, et al. Disclosure of candidate genes in acute myeloid leukemia with complex karyotypes using microarray-based molecular characterization. J Clin Oncol. 2006; 24:3887–94. [PubMed: 16864856]

14. Hillier LW, et al. Whole-genome sequencing and variant discovery in C. elegans. Nat Methods. 2008; 5:183–8. [PubMed: 18204455]

15. Wheeler DA, et al. The complete genome of an individual by massively parallel DNA sequencing. Nature. 2008; 452:872–6. [PubMed: 18421352]

16. Levy S, et al. The diploid genome sequence of an individual human. PLoS Biol. 2007; 5:e254. [PubMed: 17803354]

17. Byrd JC, et al. Pretreatment cytogenetic abnormalities are predictive of induction success, cumulative incidence of relapse, and overall survival in adult patients with de novo acute myeloid leukemia: results from Cancer and Leukemia Group B (CALGB 8461). Blood. 2002; 100:4325–36. [PubMed: 12393746]

18. Grimwade D, et al. The importance of diagnostic cytogenetics on outcome in AML: analysis of 1,612 patients entered into the MRC AML 10 trial. The Medical Research Council Adult and Children's Leukaemia Working Parties. Blood. 1998; 92:2322–33. [PubMed: 9746770]

19. Mrozek K, Heerema NA, Bloomfield CD. Cytogenetics in acute leukemia. Blood Rev. 2004; 18:115–36. [PubMed: 15010150]

20. Wendl MC, Wilson RK. Aspects of coverage in medical DNA sequencing. BMC Bioinformatics. 2008; 9:239. [PubMed: 18485222]

21. Li H, Ruan J, Durbin R. Mapping short DNA sequencing reads and calling variants using mapping quality scores. Genome Res. 2008

22. Quinlan, JR. C4.5: Programs for Machine Learning. Vol. 302. Morgan Kaufmann Publishers; San Mateo, Calif: 1993.

23. Link DC, et al. Distinct patterns of mutations occurring in de novo AML versus AML arising in the setting of severe congenital neutropenia. Blood. 2007; 110:1648–55. [PubMed: 17494858]

24. Frohling S, et al. Identification of driver and passenger mutations of FLT3 by high-throughput DNA sequence analysis and functional assessment of candidate alleles. Cancer Cell. 2007; 12:501–13. [PubMed: 18068628]

25. Levis M, Small D. FLT3: ITDoes matter in leukemia. Leukemia. 2003; 17:1738–52. [PubMed: 12970773]

26. Falini B, et al. Cytoplasmic Nucleophosmin in Acute Myelogenous Leukemia with a Normal Karyotype. N Engl J Med. 2005; 352:254–266. [PubMed: 15659725]

27. Thiede C, et al. Prevalence and prognostic impact of NPM1 mutations in 1485 adult patients with acute myeloid leukemia (AML). Blood. 2006

28. den Besten W, Kuo ML, Williams RT, Sherr CJ. Myeloid leukemia-associated nucleophosmin mutants perturb p53-dependent and independent activities of the Arf tumor suppressor protein. Cell Cycle. 2005; 4:1593–8. [PubMed: 16205118]

29. Kelly LM, et al. PML/RARalpha and FLT3-ITD induce an APL-like disease in a mouse model. Proc Natl Acad Sci U S A. 2002; 99:8283–8. [PubMed: 12060771]

A.



B.



**Figure 1. Overlap of SNPs detected in 933124 and other genomes**
(A) Venn diagram of overlap between SNPs detected in the 933124 tumor genome and the genomes of Watson and Venter. (B) Venn Diagram of overlap among 933124 tumor genome, skin genome, and dbSNP (ver. 127). Single nucleotide variants were defined with a MAQ SNP quality     15.

3,813,205 tumor single nucleotide variants
(SNVs) (Maq15)

2,647,695 well supported SNVs (Decision Tree)

2,584,418 present in skin (SNPs)

63,277 tumor-specific SNVs

31,645 in dbSNP/Watson/Venter

31,632 novel SNVs

20,440 in non-genic regions

11,192 SNVs in genic regions

216 in UTRs

10,735 intronic

241 SNVs in coding sequence

60 synonymous

181 SNVs predicted to alter gene function
(non-synonymous and splice junctions)

7 unable to be validated
(technical failures)

14 validated as germline SNVs (SNPs)

8 validated as somatic SNVs
(acquired mutations)

152 validated as wild type
(false positives)

**Figure 2. Filters used to identify somatic point mutations in the tumor genome**
See text for details.

**Figure 3. Summary of Roche/454 FLX readcount data obtained for 10 somatic mutations and 2 validated SNPs in the primary tumor, relapse tumor, and skin specimens**

The readcount data for the variant alleles in the primary tumor sample and relapse tumor sample are statistically different than that of the skin sample for all mutations (p<0.000001 for all mutations, Fisher's exact test, denoted by a single asterisk in all cases). Note that the normal skin sample was contaminated with leukemic cells containing the somatic mutations. The patient's WBC count was 105,000 (85% blasts) when the skin punch biopsy was obtained.

**Table 1**

Assessments of haploid and diploid coverage of the tumor and skin genomes from patient 933124.

|  | Tumor | | Skin | |
| --- | --- | --- | --- | --- |
| Libraries | 4 | | 3 | |
| Runs | 98 | | 34 | |
| Reads obtained | 5,858,992,064 | | 2,122,836,148 | |
| Reads passing quality filter | 3,025,923,365 | | 1,228,177,690 | |
| Bases passing quality filter | 98,184,511,523 | | 41,783,794,834 | |
| Reads aligned by Maq | 2,729,957,053 | | 1,080,576,680 | |
| Reads unaligned by Maq | 295,966,312 | | 138,276,594 | |
| SNVs detected with respect to hg18 (no Y) | 3,811,115 | | 2,918,446 | |
| SNVs (chr 1–22) detected with respect to hg18 | 3,681,968 | (100.0%) | 2,830,292 | (100.0%) |
| SNVs also present in dbSNP | 2,368,458 | (64.3%) | 2,161,695 | (76.4%) |
| SNVs also present in Venter genome | 1,499,010 | (40.7%) | 1,383,431 | (48.9%) |
| SNVs also present in Watson genome | 1,573,435 | (42.7%) | 1,456,822 | (51.5%) |
| SNVs not in dbSNP/Venter/Watson | 1,223,830 | (33.2%) | 591,131 | (20.9%) |
| SNVs not in dbSNP/Venter/Watson/skin | 925,200 | (25.1%) | – | |
| HQ SNPs | 46,494 | (100.0%) | 46,572 | (100.0%) |
| HQ SNPs where reference allele is detected | 42,419 | (91.2%) | 38,454 | (82.6%) |
| HQ SNPs where variant allele is detected | 43,164 | (92.9%) | 39,220 | (84.2%) |
| HQ SNPs where both alleles are detected | 42,415 | (91.2%) | 38,454 | (82.6%) |

**Table 2**

Somatic mutations detected in the de novo and relapse AML samples.

| Gene | Consequence | Type | Solexa Tumor Reads WT:Variant | Solexa Skin Reads WT:Variant | Conservation Score of mutant base | Mutations in other AML cases[*] |
|------|-------------|------|-------------------------------|------------------------------|-----------------------------------|-------------------------------|
| CDH24 | Y590X | nonsense | 9:9 | 16:0 | 0.998 | 0/187 |
| SLC15A1 | W77X | nonsense | 15:12 | 19:0 | 1.000 | 0/187 |
| KNDC1 | L799F | missense | 7:8 | 20:0 | NA | 0/187 |
| PTPRT | P1235L | missense | 9:13 | 16:0 | 1.000 | 0/187 |
| GRINL1B | R176H | missense | 15:10 | 14:0 | NA | 0/187 |
| GPR123 | T38I | Missense | 11:11 | 13:0 | NA | 0/187 |
| EBI2 | A338V | Missense | 7:12 | 18:2 | 1.000 | 0/187 |
| PCLKC | P1004L | missense | 19:9 | 15:1 | 0.98 | 0/187 |
| FLT3 | ITD | indel | 18:12 | 8:0 | NA | 51/185 |
| NPM1 | CATG insertion | indel | 36:6 | 33:0 | NA | 43/180 |

[*]
patient cohort defined in Link, et al23.

**Table 3**

454 Readcount data for Somatic Mutations and known SNPs in primary tumor, skin, and relapse tumor samples.

| Gene | Variant | Primary AML (100% blasts) | | | Skin | | | Relapse (78% blasts) | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Variant | Ref | % Variant | Variant | Ref | % Variant | Variant | Ref | % Variant |
| CDH24 | Y590X | 5672 | 4890 | 53.70 | 564 | 10358 | 5.16 | 3108 | 4599 | 40.33 |
| SLC15A1 | W77X | 3817 | 4962 | 43.48 | 875 | 10773 | 7.51 | 4714 | 7173 | 39.66 |
| KNDC1 | L799F | 4640 | 4848 | 48.90 | 770 | 8972 | 7.90 | 3883 | 6342 | 37.98 |
| PTPRT | P1235L | 998 | 1058 | 48.54 | 126 | 1489 | 7.80 | 350 | 493 | 41.52 |
| GRINL1B | R176H | 2211 | 2674 | 45.26 | 318 | 4461 | 6.65 | 1447 | 2070 | 41.14 |
| GPR123 | T38I | 4618 | 4569 | 50.27 | 850 | 9751 | 8.02 | 3660 | 6057 | 37.67 |
| EBI2 | A338V | 12750 | 15453 | 45.21 | 458 | 10088 | 4.34 | 2646 | 3627 | 42.18 |
| PCLKC | P1004L | 9216 | 8815 | 51.11 | 6617 | 21786 | 23.29 | 8600 | 8822 | 49.36 |
| FLT3 | ITD | 4220 | 7810 | 35.08 | 3475 | 23159 | 13.05 | 3870 | 8495 | 31.30 |
| NPM1 | CATG ins | 1550 | 1974 | 43.98 | 143 | 2390 | 5.65 | 2303 | 3910 | 37.07 |
| BRCA2 | N372H | 778 | 752 | 50.85 | 763 | 876 | 46.55 | 285 | 303 | 48.47 |
| TP53 | P72R | 8989 | 1 | 99.99 | 8161 | 0 | 100.00 | 7914 | 6 | 99.92 |

The differences between variant frequencies in primary or relapse tumor samples and skin were highly significant for all somatic mutations (p<0.000001, Fisher's exact test, one tailed). The BRCA2 variant is a known heterozygous SNP in this genome, and the TP53 variant is a known homozygous SNP.