# Long-read whole genome sequencing identifies causal

# structural variation in a Mendelian disease

Jason D. Merker,[1,2,8] Aaron M. Wenger,[3,8] Tam Sneddon,[2] Megan Grove,[2] Daryl Waggott,[4,5] Sowmi Utiramerur,[2]

Yanli Hou,[1] Christine C. Lambert,[3] Kevin S. Eng,[3] Luke Hickey,[3] Jonas Korlach,[3] James Ford,[4,6,7] Euan A.

Ashley*[2,4,5,6]




[1]Department of Pathology, Stanford University, Stanford, California, 94305, USA

[2]Clinical Genomics Service, Stanford Health Care, Stanford, California, 94305, USA

[3]Pacific Biosciences, Menlo Park, California, 94025, USA

[4]Department of Medicine, Stanford University, Stanford, California, 94305, USA

[5]Stanford Center for Inherited Cardiovascular Disease, Stanford University, Stanford, California, 94305, USA

[6]Department of Genetics, Stanford University, Stanford, California, 94305, USA

[7]Stanford Cancer Institute, Stanford, California, 94305, USA

[8]These authors contributed equally to this work

*Correspondence: euan@stanford.edu

Abstract word count: 244

Main text word count: 1633

23     **Abstract**

24     Current clinical genomics assays primarily utilize short-read sequencing (SRS), which offers high throughput, high

25     base accuracy, and low cost per base. SRS has, however, limited ability to evaluate tandem repeats, regions with

26     high [GC] or [AT] content, highly polymorphic regions, highly paralogous regions, and large-scale structural

27     variants. Long-read sequencing (LRS) has complementary strengths and offers a means to discover overlooked

28     genetic variation in patients undiagnosed by SRS. To evaluate LRS, we selected a patient who presented with

29     multiple neoplasia and cardiac myxomata suggestive of Carney complex for whom targeted clinical gene testing and

30     whole genome SRS were negative. Low coverage whole genome LRS was performed on the PacBio Sequel system

31     and structural variants were called, yielding 6,971 deletions and 6,821 insertions > 50bp. Filtering for variants that

32     are absent in an unrelated control and that overlap a coding exon of a disease gene identified three deletions and

33     three insertions. One of these, a heterozygous 2,184 bp deletion, overlaps the first coding exon of *PRKAR1A*, which

34     is implicated in autosomal dominant Carney complex. This variant was confirmed by Sanger sequencing and was

35     classified as pathogenic using standard criteria for the interpretation of sequence variants. This first successful

36     application of whole genome LRS to identify a pathogenic variant suggests that LRS has significant potential to

37     identify disease-causing structural variation. We recommend larger studies to evaluate the diagnostic yield of LRS,

38     and the development of a comprehensive catalog of common human structural variation to support future studies.

39

40    Short-read sequencing (SRS) methods are currently favored in clinical medicine because of their cost effectiveness

41    and low per-base error rate. However, these methods are limited in their ability to capture the full range of genomic

42    variation.[1] Areas of low complexity such as repeats and areas of high polymorphism, such as the HLA region,

43    present challenges to SRS and reference-based genome assembly. Indeed, with 100 base pair (bp) read length, fully

44    5% of the genome cannot be uniquely mapped.[2] In addition, many diseases are caused by repeats that become

45    increasingly pathogenic in a range beyond the resolution of SRS. Another challenge comes in the form of structural

46    variation, and while SRS has been very successful in the genetic discovery of single nucleotide and small insertion-

47    deletion variation, recent findings suggest we have greatly underestimated the extent and complexity of structural

48    variation in the genome.[3,4]

49    Long-read sequencing (LRS), typified by PacBio® single molecule, real-time (SMRT®) sequencing, offers

50    complementary strengths to SRS. PacBio LRS produces reads of several thousand base pairs with uniform coverage

51    across sequence contexts.[5] Individual long reads have a lower accuracy (85%) than short reads, but errors are

52    random and are correctable with sufficient coverage, leading to extremely high consensus accuracy.[5,6] Further, long

53    reads are more accurately mapped to the genome and access regions that are beyond the reach of short reads.[1] Of

54    particular note, recent PacBio LRS *de novo* human genome assemblies have revealed tens of thousands of structural

55    variants per genome, many times more than previously observed with SRS.[3,7] These capabilities, together with

56    continuing progress in throughput and cost, have begun to make LRS an option for broader application in human

57    genomics.

58    Here, we report the use of low coverage whole genome PacBio LRS to secure a diagnosis of Carney complex in a

59    patient unsolved by clinical single gene testing and whole genome SRS. The patient is an Asian/Hispanic male, the

60    product of an uncomplicated term pregnancy who was hospitalized for the first 10 days of life for cardiac and

61    respiratory issues (**Figure 1A**). He remained well until the age of 7 years, when, following the discovery of a heart

62    murmur, he was found to have a left atrial myxoma that was surgically removed. At 10 years, he was noted to have a

63    testicular mass that, at orchiectomy, was found to be a Sertoli-Leydig cell tumor. At 13 years, a pituitary tumor was

64    found and initial conservative management was adopted. Aged 16, he was noted to have both an adrenal

65    microadenoma and recurrence of the cardiac myxomata in the left ventricle and right atrium. Blue naevi were

66    reported. He underwent a second surgical resection of the myxomata with uncomplicated recovery. Aged 18,

3

67    recurrent cardiac myxomata including a right ventricular and two left ventricular tumors were once again resected

68    and a goretex patch was placed in the right ventricular wall. In the immediate post-operative period, he suffered

69    ventricular tachycardia (VT) and cardiac arrest with spontaneous return of circulation. At this time, a genetics

70    evaluation suggested the possibility of Carney complex but clinical genetic testing (sequencing of *PRKAR1A*) was

71    negative for disease causing variation. At age 19, multiple thyroid nodules were noted on ultrasound, and he was

72    also diagnosed with ACTH-independent Cushing's syndrome, secondary to the adrenal microadenoma. At 21, he

73    underwent trans-sphenoidal resection of the pituitary tumor. At this time, he was found to have recurrent myxomata

74    in the left ventricular outflow tract that have subsequently increased in size (**Figure 1B-C**). To date, these have been

75    treated conservatively with anti-coagulation to reduce the risk of stroke. As of 2016, he is under consideration for

76    heart transplantation and the transplant team judged a molecular diagnosis highly desirable prior to cardiac

77    transplant listing. As a result, whole genome SRS was performed. Genomic DNA was purified, and a library was

78    generated using the Illumina® TruSeq® DNA PCR-Free Library Prep Kit, and genome sequencing was performed

79    using the Illumina HiSeq® 2500 System with paired-end 2 x 100 bp reads to a 36-fold mean depth of coverage. The

80    data analysis and variant curation were performed by the Stanford Medicine Clinical Genomics Service. Single

81    nucleotide variants and small insertions and deletions were identified using MedGAP v2.0, a pipeline based on

82    GATK best practices for data pre-processing and variant discovery with GATK HaplotypeCaller v3.1.1.[8] This

83    analysis pipeline did not identify any variants that would explain the clinical findings in the patient.

84    To evaluate structural variation, low coverage whole genome LRS was performed on the PacBio Sequel™ system.

85    Following consent under a protocol approved by the Stanford University Institutional Review Board, DNA was

86    isolated from a peripheral blood specimen using the Gentra® Puregene® Blood Kit (Qiagen, Germantown, MD).

87    The DNA was sheared to 20 kb fragments on a Megaruptor® and size-selected to 10 kb using the Sage Science

88    BluePippin™ system. A SMRTbell™ library was prepared and sequenced on 10 Sequel SMRT Cells 1M with

89    chemistry S/P1-C1.2 and 6 hour collections. The sequencing yielded 26.7 Gb (8.6-fold coverage of human genome)

90    in 4.3 million reads with a read length N50 of 9,614 bp. Reads were mapped to the GRCh37/hg19 assembly of the

91    human genome using NGM-LR v0.1.4 with default parameters.[9] Structural variants were called using PBHoney

92    Spots with '-q 10 –m 10 –i 20 –e 1 –E 1 –spanMax 100000 –consensus None' for deletions and '-q 10 –m 70 –i 20 –

93    e 2 –E 2 –spanMax 10000 –consensus None' for insertions.[10] Variant calls were further refined to retain only those

94    larger than 50 bp, supported by at least 20% of local reads, and at least 100 bp from an assembly gap.

95    The resulting call set consisted of 6,971 deletions and 6,821 insertions. To prioritize candidate pathogenic variants,

96    the call set was filtered to exclude variants within a segmental duplication or present in the unrelated control

97    individual NA12878 (A.W., unpublished data). This left 2,368 deletions and 3,174 insertions. Focusing on variants

98    that overlap a RefSeq coding exon resulted in 20 deletions and 16 insertions, with 3 deletions and 3 insertions in

99    genes tied to a genetic disease in OMIM (**Table 1**). Manual review of the 6 candidate variants and correlation with

100    phenotype identified a heterozygous deletion that removes the first coding exon of *PRKAR1A* (NM_212472.2).

101    Germline variants in *PRKAR1A* cause Carney complex, type 1 (MIM #160980), an autosomal dominant multiple

102    neoplasia syndrome.[11] Two of four reads at the locus unambiguously support the presence of a deletion variant

103    (**Figure 2A**). Because of the random errors in LRS, individual reads from the same allele can have slight

104    disagreements, and two reads can be insufficient to define exact deletion breakpoints with full confidence. Here, the

105    higher quality read supports a 2,184 bp deletion of GRCh37/hg19 chr17:66,510,475-66,512,658

106    (NC_000017.10:g.66510475_66512658del). This heterozygous deletion variant was validated by Sanger

107    sequencing, which in this case confirmed the precise breakpoints identified by LRS (**Figure 2B**).

108    It is difficult to call structural variants in SRS data with simultaneously high sensitivity and specificity that is

109    necessary for clinical laboratory testing. Nevertheless, once a small candidate gene list or approximate breakpoints

110    are known, many variants can be identified retrospectively.[5] In such cases, SRS often provides exact breakpoints to

111    refine the variant discovered by LRS.[12] Manual inspection of SRS data from the *PRKAR1A* locus shows support for

112    the heterozygous deletion through a drop in read depth and alignment clipping at the deletion breakpoints (**Figure**

113    **2C**). Multiple short-read structural variant callers, including Pindel, Lumpy, BreakDancer, Manta, CNVKit, and

114    CNVnator, were retrospectively used to identify structural variants.[13–18] All tools were run with default parameters.

115    Pindel, Lumpy, BreakDancer, and Manta all identify a deletion in the locus. Pindel and Manta approximately match

116    the breakpoints identified from LRS and Sanger sequencing.

117    This case demonstrates the ability of whole genome LRS to detect causal structural variation in a rare disease, and to

118    our knowledge, this is the first reported application of whole genome LRS to identify a pathogenic variant in a

119    patient. Although manual inspection of the aligned read data and short-read structural variant callers are able to

120    identify this 2,184 bp deletion, these approaches are not practical to apply genome wide due to limited throughput

121    and high false-positive call rates, respectively. Looking forward, clinical-grade genomics demands strong precision

5

122    and recall across the full spectrum of genetic variation. SRS has limited sensitivity for variants larger than a few

123    base pairs, and it can miss up to 80% of the structural variants in an individual genome.[3] LRS appears to be capable

124    of identifying much of the missed variation, and manifests high recall of structural variants even at low depths of

125    coverage.[12]

126    To accelerate the adoption of sequencing-based structural variant analysis into clinical practice, it will be important

127    for the community to develop and expand an ecosystem of tools and databases similar to that which has arisen

128    around smaller variants. We advocate further development and continued evaluation of tools and best practices for

129    calling structural variants from SRS, LRS and orthogonal data. Additionally, we recommend that the community

130    prioritize creation of a catalog of structural variation derived from these data sources. Databases of common single

131    nucleotide variation, such as ExAC, have proven incredibly valuable.[19] We expect that a comparable database of

132    structural variants would be similarly valuable and that building the database from LRS would greatly expand

133    current catalogs such as DGV and dbVar.[20,21]

134    LRS has seen limited adoption in clinical genomics laboratories, in large part due to the per base error rate and cost.

135    Although the individual read error rate requires higher coverage to provide clinical-grade identification of single

136    nucleotide variants, high coverage is not necessarily required for sensitive and specific detection of larger structural

137    variants. Cost effectiveness will ultimately be judged not on cost per base, but on cost per diagnosis. Larger studies

138    on the diagnostic yield of various approaches using LRS will be required to answer the question of the most cost

139    effective technologies for clinical genomics moving forward.

140    **Conflict of interest statement**

141    AW, CL, KE, LH, and JK are employees and shareholders of Pacific Biosciences, a company commercializing DNA

142    sequencing technologies.

143

144    **Acknowledgements**

145    The authors thank the research subject and clinical care teams for their participation in this research study; Chen-

146    Shan (Jason) Chin for helpful discussions; and Primo Baybayan and Matt Boitano for PacBio library preparation and

147    sequencing.

148    **Web Resources**

149    OMIM, http://www.omim.org/

150    **References**

151    1. Ashley, E.A. (2016). Towards precision medicine. Nat. Rev. Genet. *17*, 507–522.

152    2. Goldfeder, R.L., Priest, J.R., Zook, J.M., Grove, M.E., Waggott, D., Wheeler, M.T., Salit, M., and

153    Ashley, E.A. (2016). Medical implications of technical accuracy in genome sequencing. Genome Med. *8*,

154    1–12.

155    3. Chaisson, M.J.P., Huddleston, J., Dennis, M.Y., Sudmant, P.H., Malig, M., Hormozdiari, F.,

156    Antonacci, F., Surti, U., Sandstrom, R., Boitano, M., et al. (2015). Resolving the complexity of the human

157    genome using single-molecule sequencing. Nature *517*, 608–611.

158    4. Huddleston, J., Chaisson, M.J., Meltz Steinberg, K., Warren, W., Hoekzema, K., Gordon, D.S., Graves-

159    Lindsay, T.A., Munson, K.M., Kronenberg, Z.N., Vives, L., et al. (2016). Discovery and genotyping of

160    structural variation from long-read haploid genome sequence data. Genome Res. Advance online

161    publication. doi:10.1101/gr.214007.116

162    5. Chaisson, M.J.P., Wilson, R.K., and Eichler, E.E. (2015). Genetic variation and the de novo assembly

163    of human genomes. Nat. Rev. Genet. *16*, 627–640.

164    6. Chin, C.-S., Alexander, D.H., Marks, P., Klammer, A.A., Drake, J., Heiner, C., Clum, A., Copeland,

165    A., Huddleston, J., Eichler, E.E., et al. (2013). Nonhybrid, finished microbial genome assemblies from

166     long-read SMRT sequencing data. Nat. Methods *10*, 563–569.

167     7. Seo, J.-S., Rhie, A., Kim, J., Lee, S., Sohn, M.-H., Kim, C.-U., Hastie, A., Cao, H., Yun, J.-Y., Kim, J.,

168     et al. (2016). De novo assembly and phasing of a Korean human genome. Nature *538*, 243–247.

169     8. Van der Auwera, G.A., Carneiro, M.O., Hartl, C., Poplin, R., Del Angel, G., Levy-Moonshine, A.,

170     Jordan, T., Shakir, K., Roazen, D., Thibault, J., et al. (2013). From FastQ data to high confidence variant

171     calls: the Genome Analysis Toolkit best practices pipeline. Curr. Protoc. Bioinformatics *43*, 11.10.1–33.

172     9. Rescheneder, P., Sedlazeck, F.J., von Haeseler, A., and Schatz, M.C. (2016). NextGenMap-LR

173     computer program, https://github.com/philres/nextgenmap-lr (Cold Spring Harbor Laboratory).

174     10. English, A.C., Salerno, W.J., and Reid, J.G. (2014). PBHoney: identifying genomic variants via long-

175     read discordance and interrupted mapping. BMC Bioinformatics *15*, 180.

176     11. Kirschner, L.S., Carney, J.A., Pack, S.D., Taymans, S.E., Giatzakis, C., Cho, Y.S., Cho-Chung, Y.S.,

177     and Stratakis, C.A. (2000). Mutations of the gene encoding the protein kinase A type I-alpha regulatory

178     subunit in patients with the Carney complex. Nat. Genet. *26*, 89–92.

179     12. English, A.C., Salerno, W.J., Hampton, O.A., Gonzaga-Jauregui, C., Ambreth, S., Ritter, D.I., Beck,

180     C.R., Davis, C.F., Dahdouli, M., Ma, S., et al. (2015). Assessing structural variation in a personal

181     genome-towards a human reference diploid genome. BMC Genomics *16*, 286.

182     13. Ye, K., Schulz, M.H., Long, Q., Apweiler, R., and Ning, Z. (2009). Pindel: a pattern growth approach

183     to detect break points of large deletions and medium sized insertions from paired-end short reads.

184     Bioinformatics *25*, 2865–2871.

185     14. Layer, R.M., Chiang, C., Quinlan, A.R., and Hall, I.M. (2014). LUMPY: a probabilistic framework

186     for structural variant discovery. Genome Biol. *15*, R84.

187     15. Chen, K., Wallis, J.W., McLellan, M.D., Larson, D.E., Kalicki, J.M., Pohl, C.S., McGrath, S.D.,

188    Wendl, M.C., Zhang, Q., Locke, D.P., et al. (2009). BreakDancer: an algorithm for high-resolution

189    mapping of genomic structural variation. Nat. Methods *6*, 677–681.

190    16. Chen, X., Schulz-Trieglaff, O., Shaw, R., Barnes, B., Schlesinger, F., Källberg, M., Cox, A.J.,

191    Kruglyak, S., and Saunders, C.T. (2016). Manta: rapid detection of structural variants and indels for

192    germline and cancer sequencing applications. Bioinformatics *32*, 1220–1222.

193    17. Talevich, E., Shain, A.H., Botton, T., and Bastian, B.C. (2016). CNVkit: Genome-Wide Copy

194    Number Detection and Visualization from Targeted DNA Sequencing. PLoS Comput. Biol. *12*,

195    e1004873.

196    18. Abyzov, A., Urban, A.E., Snyder, M., and Gerstein, M. (2011). CNVnator: an approach to discover,

197    genotype, and characterize typical and atypical CNVs from family and population genome sequencing.

198    Genome Res. *21*, 974–984.

199    19. Lek, M., Karczewski, K.J., Minikel, E.V., Samocha, K.E., Banks, E., Fennell, T., O'Donnell-Luria,

200    A.H., Ware, J.S., Hill, A.J., Cummings, B.B., et al. (2016). Analysis of protein-coding genetic variation in

201    60,706 humans. Nature *536*, 285–291.

202    20. MacDonald, J.R., Ziman, R., Yuen, R.K.C., Feuk, L., and Scherer, S.W. (2014). The Database of

203    Genomic Variants: a curated collection of structural variation in the human genome. Nucleic Acids Res.

204    *42*, D986–D992.

205    21. Lappalainen, I., Lopez, J., Skipper, L., Hefferon, T., Spalding, J.D., Garner, J., Chen, C., Maguire, M.,

206    Corbett, M., Zhou, G., et al. (2013). DbVar and DGVa: public archives for genomic structural variation.

207    Nucleic Acids Res. *41*, D936–D941.

208

209

210 **Figure Titles and Legends**

211 **Figure 1**. **Clinical history and three-dimensional transthoracic echocardiography of patient with multiple**

212 **neoplasia including cardiac myxomata**. (A) Patient narrative.  VT= ventricular tachycardia (B) A 2 x 3 cm

213 myxoma is seen in the left ventricular outflow tract (white arrow).  (C) The 2 x 3 cm myxoma is seen from another

214 perspective (lower left, white arrow).  A 5 x 4 cm myxoma is seen in the right atrium (lower right, white arrow).
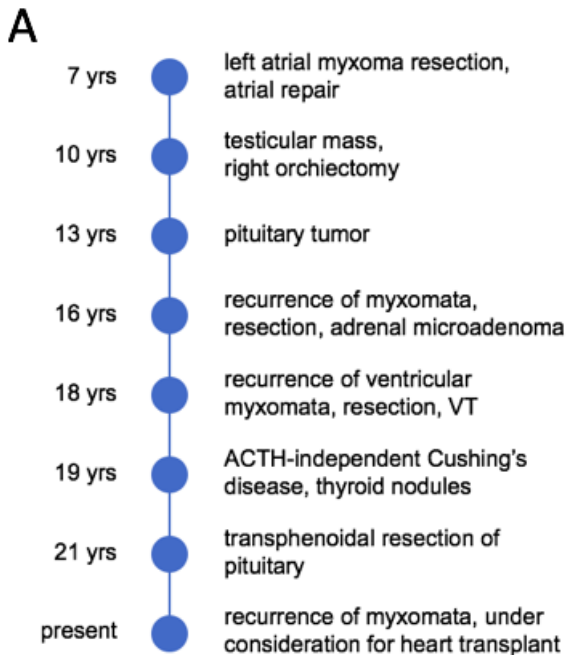
215

216 **Figure 2. Heterozygous deletion in *PRKAR1A*.** (A) PacBio long reads identify a heterozygous 2,184 bp deletion

217 that includes the first coding exon of *PRKAR1A*. Two of four reads at the locus support the deletion. (B) Sanger

218 sequencing confirms the deletion. The forward (YH_479426-1073) and reverse (YH_479426-1074) sequences from

219 a representative amplicon agree to the base pair with the higher quality PacBio read, PacBio_53019216. (C)

220 Illumina short reads support the heterozygous deletion variant through a drop in read coverage and clipped reads at
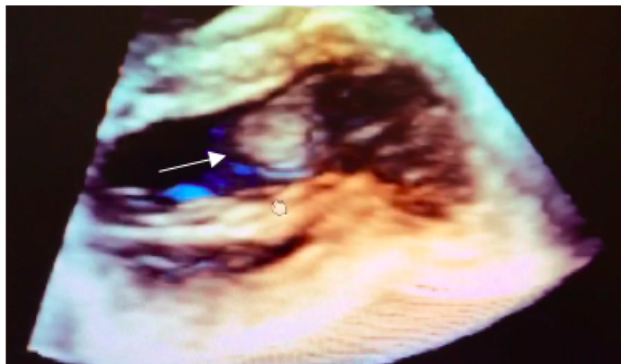
221 the deletion breakpoints.

222

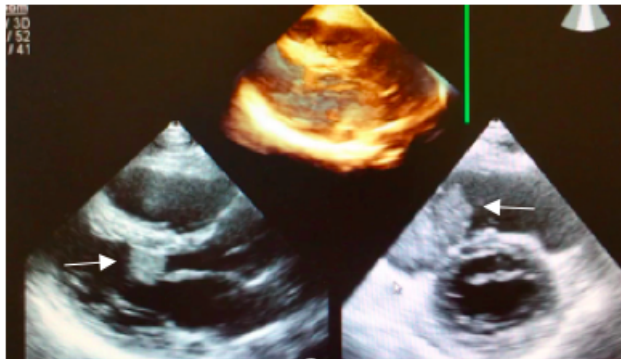|                               | Deletions > 50 bp | Insertions > 50 bp |
|-------------------------------|-------------------|--------------------|
| Initial call set              | 6,971             | 6,821              |
| Not in segmental duplication  | 5,818             | 6,254              |
| Not in NA12878 control        | 2,368             | 3,174              |
| Overlaps RefSeq coding exon   | 20                | 16                 |
| Gene linked to disease in OMIM| 3                 | 3                  |

223

224 **Table 1. Prioritizing candidate pathogenic variants.** The initial call set of 6,971 deletions and 6,821 insertions

225 was filtered to remove variants in segmental duplications or the NA12878 control and to focus on variants that

226 overlap coding exons of genes with a known link to genetic disease.
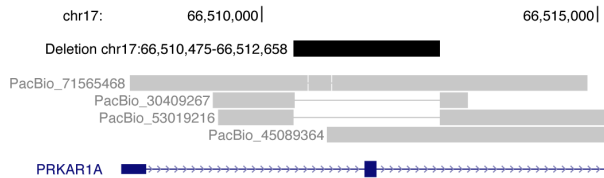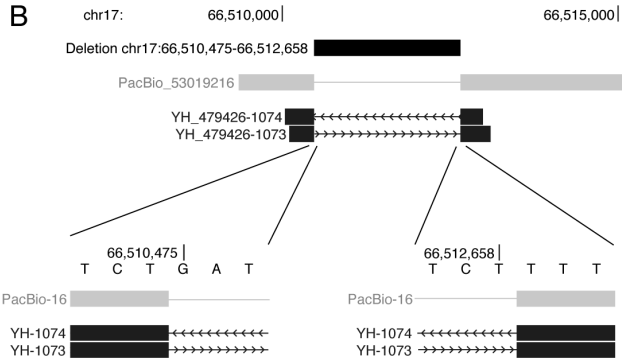
11

**A**

| | |
|---|---|
| 7 yrs | left atrial myxoma resection, atrial repair |
| 10 yrs | testicular mass, right orchiectomy |
| 13 yrs | pituitary tumor |
| 16 yrs | recurrence of myxomata, resection, adrenal microadenoma |
| 18 yrs | recurrence of ventricular myxomata, resection, VT |
| 19 yrs | ACTH-independent Cushing's disease, thyroid nodules |
| 21 yrs | transphenoidal resection of pituitary |
| present | recurrence of myxomata, under consideration for heart transplant |