

1 **Novel pedigree analysis implicates DNA repair and chromatin**  
2 **remodeling in Multiple Myeloma risk**

3 Rosalie G. Waller<sup>1¶</sup>, Todd M. Darlington<sup>1¶</sup>, Xiaomu Wei<sup>2</sup>, Myke Madsen<sup>1</sup>, Alun Thomas<sup>1</sup>, Karen  
4 Curtin<sup>1</sup>, Hilary Coon<sup>1</sup>, Venkatesh Rajamanickam<sup>1</sup>, Justin Musinsky<sup>3</sup>, David Jayabalan<sup>2</sup>, Djordje  
5 Atanackovic<sup>1</sup>, Vincent Rajkumar<sup>4</sup>, Shaji Kumar<sup>4</sup>, Susan Slager<sup>4</sup>, Mridu Middha<sup>5</sup>, Perrine Galia<sup>7</sup>,  
6 Delphine Demangel<sup>7</sup>, Mohamed Salama<sup>1</sup>, Vijai Joseph<sup>3</sup>, James McKay<sup>8</sup>, Kenneth Offit<sup>3</sup>, Robert  
7 J. Klein<sup>5</sup>, Steven M. Lipkin<sup>2</sup>, Charles Dumontet<sup>6</sup>, Celine M. Vachon<sup>4</sup>, Nicola J. Camp<sup>1\*</sup>

- 8 1. University of Utah School of Medicine, Utah, USA  
9 2. Weill Cornell Medical College, New York, USA  
10 3. Memorial Sloan Kettering Cancer Center, New York, USA.  
11 4. Mayo Clinic, Minnesota, USA  
12 5. Icahn School of Medicine at Mount Sinai, New York, USA  
13 6. INSERM 1052/CNRS 5286/UCBL  
14 7. ProfileXpert, Lyon, France  
15 8. International Agency for Research on Cancer, Lyon, France

16 \*Corresponding author

17 E-mail: [nicola.camp@utah.edu](mailto:nicola.camp@utah.edu) (NJC)

18 <sup>¶</sup>Equal contribution

## 1 **ABSTRACT**

2           The high-risk pedigree (HRP) design is an established strategy to discover rare, highly-  
3 penetrant, Mendelian-like causal variants. Its success, however, in complex traits has been  
4 modest, largely due to challenges of genetic heterogeneity and complex inheritance models. We  
5 describe a HRP strategy that addresses intra-familial heterogeneity, and identifies inherited  
6 segments important for mapping regulatory risk. We illustrate this new Shared Genomic  
7 Segment (SGS) method in 11 extended, Utah, multiple myeloma (MM) HRPs, and subsequent  
8 exome sequencing in SGS regions of interest in 1064 MM/MGUS (monoclonal gammopathy of  
9 undetermined significance – a precursor to MM) cases and 964 controls from a jointly-called  
10 collaborative resource, including cases from the initial 11 HRPs. One genome-wide significant  
11 1.8 Mb shared segment was found at 6q16. Exome sequencing in this region revealed predicted  
12 deleterious variants in *USP45* (p.Gln691\*, p.Gln621Glu), a gene known to influence DNA repair  
13 through endonuclease regulation. Additionally, a 1.2 Mb segment at 1p36.11 is inherited in two  
14 Utah HRPs, with coding variants identified in *ARID1A* (p.Ser90Gly, p.Met890Val), a key gene in  
15 the SWI/SNF chromatin remodeling complex. Our results provide compelling statistical and  
16 genetic evidence for segregating risk variants for MM. In addition, we demonstrate a novel  
17 strategy to use large HRPs for risk-variant discovery more generally in complex traits.

## 18 **AUTHOR SUMMARY**

19           Although family-based studies demonstrate inherited variants play a role in many  
20 common and complex diseases, finding the genes responsible remains a challenge. High-risk  
21 pedigrees, or families with more disease than expected by chance, have been helpful in the  
22 discovery of variants responsible for less complex diseases, but have not reached their potential  
23 in complex diseases. Here, we describe a method to utilize high-risk pedigrees to discover risk-

1 genes in complex diseases. Our method is appropriate for complex diseases because it allows  
2 for genetic-heterogeneity, or multiple causes of disease, within a pedigree. This method allows  
3 us to identify shared segments that likely harbor disease-causing variants in a family. We  
4 demonstrate our method in Multiple Myeloma, a heritable and complex cancer of plasma cells.  
5 We identified two genes *USP45* and *ARID1A* that fall within shared segments with compelling  
6 statistical evidence. Exome sequencing of these genes revealed likely-damaging variants  
7 inherited in Myeloma high-risk families, suggesting these genes likely play a role in development  
8 of Myeloma. Our Myeloma findings demonstrate our high-risk pedigree method can identify  
9 genetic regions of interest in large high-risk pedigrees that are also relevant to smaller nuclear  
10 families and overall disease risk. In sum, we offer a strategy, applicable across phenotypes, to  
11 revitalize high-risk pedigrees in the discovery of the genetic basis of common and complex  
12 disease.

## 13 **INTRODUCTION**

14 Rare risk variants have been suggested as a source of missing heritability in the majority  
15 of complex traits [1–3]. High-risk pedigrees (HRPs) are a mainstay for identifying rare, highly  
16 penetrant, Mendelian-like causal variants [4–11]. However, while successful for relatively  
17 simple traits, genetic heterogeneity remains a major obstacle that reduces the effectiveness of  
18 HRPs for gene mapping in complex traits [12,13]. Also challenging is mapping regulatory  
19 variants, likely to be important for complex traits, necessitating interrogation outside the well-  
20 annotated coding regions of the genome [14,15]. Localizing chromosomal regions to target the  
21 search for rare risk variants will be instrumental in mapping them.

22 Here we develop a HRP strategy based on our previous Shared Genomic Segment  
23 (SGS) approach [16] that focuses on pedigrees sufficiently large to singularly identify

1 segregating chromosomal segments of statistical merit. The method addresses genetic  
2 heterogeneity by optimizing over all possible subsets of studied cases in a HRP. Key to the  
3 utility of the method is the derivation of significance thresholds for interpretation. These  
4 thresholds address the genome-wide search and the multiple testing, inherent from the  
5 optimization, through use of distribution fitting and the Theory of Large Deviations.

6 We apply this novel method to 11 MM HRPs, and use exome sequencing from a  
7 collaborative resource of 68 multiplex MM pedigrees to perform subsequent targeted searches  
8 at the variant level. MM is a complex cancer of the plasma cells with 30,330 new cases annually  
9 (incidence 6.5/100,000 per year) [17]. Despite survival dramatically increasing from 25.8% in  
10 1980 to 48.5% in 2012, MM remains a cancer with one of the lowest 5-year survival rates in  
11 adult hematological malignancies [17]. MM is preceded by a condition referred to as monoclonal  
12 gammopathy of undetermined significance (MGUS). Evidence for the familial clustering of MM is  
13 consistently replicated [18–21], as is its clustering with MGUS [22–25]. Genetic pedigree studies  
14 in MM are scarce as it remains a challenge to acquire samples in pedigrees due to rarity and  
15 low survival rates. The Utah MM HRPs are one of only a few pedigree resources worldwide and  
16 contains unparalleled multi-generational high-risk pedigrees. Thus far, no segregating risk  
17 variants have been identified for MM.

## 18 **RESULTS**

### 19 **Pedigree analysis strategy**

20 We developed a gene mapping strategy, based on the SGS method [16,26], that  
21 accounts for intra-familial heterogeneity and multiple testing. The basic SGS method identifies  
22 genomic segments shared identical-by-state (sharing without regard to inheritance) between a  
23 defined set of cases using a dense genome-wide map of common single nucleotide

1 polymorphisms (SNPs), either from a genotyping platform or extracted from sequence data. If  
2 the length of a shared segment is significantly longer than by chance, inherited sharing is  
3 implied; theoretically, chance inherited sharing in distant relatives is extremely improbable.  
4 Nominal chance occurrence (nominal p-value) for shared segments is assessed empirically  
5 using gene-drop simulations to create a null distribution, as follows. Null genotype  
6 configurations are generated by assigning haplotypes to pedigree founders according to a  
7 publicly available linkage disequilibrium (LD) map, followed by segregation of these through the  
8 pedigree structure to the case set via simulated Mendelian inheritance according to a genetic  
9 (recombination) map. Gene-drops are performed independent of disease status and the  
10 resulting genotype data in the case set are representative of chance sharing. This basic method  
11 was shown to have excellent power in homogeneous pedigrees [16].

12 In our new strategy, we iterate over all non-trivial combinations of the cases (subsets) in  
13 each pedigree to address heterogeneity in a “brute-force” fashion. For each subset, shared  
14 segments throughout the genome are identified and nominal p-values assigned. Across  
15 subsets, an optimization procedure is performed at every marker across the genome to identify  
16 the segment with the most significant sharing evidence. All shared segments selected by the  
17 optimization procedure, and their respective p-values, comprise the final optimized SGS results.

18 To perform significance testing and identify segments that are unexpected by chance  
19 (hypothesized to harbor risk loci), we derive significance thresholds to account for the genome-  
20 wide optimization. Acknowledging that the vast majority of observed sharing across a genome is  
21 under the null (true risk loci are a very small minority of the genome), we use the observed  
22 optimized results ( $Y = -\log_{10}(p)$ ) to model the distribution for optimized SGS results. We note  
23 that this approach may be slightly conservative because signals for true risk loci are also  
24 included. We identified the gamma distribution as adequate to represent the distribution (Fig. 1).  
25 Based on the fitted distribution,  $Y \sim \Gamma(k, \sigma)$ , where  $k$  and  $\sigma$  and the shape and scale parameters,

1 we apply the Theory of Large Deviations; previously applied to successfully model genome-wide  
2 fluctuations in linkage analysis [27]. The significance threshold,  $T$ , accounts for multiple testing  
3 of optimized segments across the genome, and is found by solving Eq. 1:

$$4 \quad \mu(X) = [C + 2GX]\alpha(X) \quad (1)$$

5 where  $T = 10^{-X\sigma/2}$ ,  $X = 2Y/\sigma \sim \chi_{2k}^2$ ,  $\mu(X)$  is the genome-wide false positive rate required,  $\alpha(X)$   
6 is nominal probability of exceeding  $X$ ,  $C$  is the number of chromosomes considered, and  $G$  is the  
7 genome length in Morgans. A criterion of  $\mu(X) = 0.05$  is typically used to define the genome-  
8 wide significant threshold (false positive rate of 0.05 per genome), and  $\mu(X) = 1$  to define the  
9 genome-wide suggestive threshold (false positive rate of 1 per genome).

10 In general, we found that the fitted distributions were sufficiently stable to produce robust  
11 significance thresholds after 100,000-300,000 simulations (Table 1). Typically, threshold  
12 determination requires 1,000-3,000 CPU hours per pedigree, increasing with the number of  
13 subsets and separating meioses between pedigree cases. For example, in pedigree UT-  
14 571744, 300k simulations genome-wide (2,513,408 segments) took 1,275 CPU hours on  
15 tangent nodes featuring Intel Xeon E5-2650 processors. Once significance thresholds are  
16 established, subset/segment combinations of potential interest are identified and additional  
17 simulations are restricted to those combinations to gain the required p-value resolution. For  
18 these subsequent targeted simulations, we use a marginalized LD map specific for the segment  
19 of interest, dramatically reducing the analysis time. For example, in pedigree UT-571744, 600M  
20 simulations on one segment took 325 CPU hours on tangent nodes featuring Intel Xeon E5-  
21 2650 processors. See S1 Fig. for an overview of the strategy pipeline.

**Table 1. Genome-wide Significance Thresholds.** Fitted distributions are stable enough for threshold determination after 100,000 to 300,000 simulations.

Pedigree	100k	200k	300k	1M
260	$6.36 \times 10^{-6}$	$6.35 \times 10^{-6}$	$6.28 \times 10^{-6}$	$6.25 \times 10^{-6}$
576834	$3.50 \times 10^{-6}$	$3.53 \times 10^{-6}$	$3.53 \times 10^{-6}$	$3.51 \times 10^{-6}$
571744	$3.80 \times 10^{-6}$	$3.83 \times 10^{-6}$	$3.75 \times 10^{-6}$	$3.80 \times 10^{-6}$
34955	$5.67 \times 10^{-6}$	$5.60 \times 10^{-6}$	$5.61 \times 10^{-6}$	$5.61 \times 10^{-6}$

## 1 Application to Utah, MM HRP

2 We applied our new pedigree analysis strategy to 11 Utah MM HRP using high-density  
3 OMNI Express SNP array genotype data. Each pedigree was selected to contain excess MM (4-  
4 37 MM total per pedigree), had 2-4 sampled MM cases with genotype data, and 8-23 meioses  
5 per pedigree between the sampled cases. After quality control, a consistent set of 678,447  
6 SNPs were used for all SGS analyses. The total number of shared segments for each pedigree  
7 across all subsets ranged from 638,525 to 6,765,500 (larger pedigrees with more subsets  
8 producing larger numbers of segments). After optimization,  $Y = -\log_{10}(p)$  for 6,697 to 10,369  
9 segments were fit to gamma distributions for each pedigree, and used to determine genome-  
10 wide significant and suggestive thresholds (Eq. 1). The genome-wide significant thresholds  
11 ranged from  $6.2 \times 10^{-5}$  to  $7.8 \times 10^{-7}$  and genome-wide suggestive from  $8.2 \times 10^{-4}$  to  $2.1 \times 10^{-5}$  (S1  
12 Table).

13 A genome-wide significant, 1.8 Mb shared segment ( $p = 3.3 \times 10^{-6}$ ) was observed in  
14 pedigree UT-571744. All three genotyped MM cases, separated by 20 meioses, share the  
15 segment (Fig. 2a and Table 2). The segment is located at chromosome 6q16 (98.49-100.24 Mb;  
16 hg19) and includes 9 genes: *POU3F2*, *FBXL4*, *FAXC*, *COQ3*, *PNISR*, *USP45*, *TSTD3*, *CCNC*,  
17 and *PRDM13* (Figure 2b).

**Table 2. Significant or overlapping SGSs and segregating SNVs.**

Family	Cases	Me	Position	Len	p-value	Gene	Conseq	Impact	AAF
UT 571744	3	20	6:98,489,655— 100,243,996	1.8	3.3x10 <sup>-6*</sup>				
PET-Nice 0909	3(2)	3	6:99,891,443			<i>USP45</i>	p.Gln691*	SG	None
Mayo 458	2(1)	2	6:99,893,787			<i>USP45</i>	p.Gln621Glu	MS	None
UT 576834	3	12	1:24,389,214— 33,298,821	8.9	3.0x10 <sup>-4</sup>				
UT 260	3	16	1:26,224,634— 27,384,988	1.2	2.1x10 <sup>-4</sup>				
UT 576834	3	12	1:27,023,162			<i>ARID1A</i>	p.Ser090Gly	MS	0.0000
Cornell MM12	2	4	1:27,089,712			<i>ARID1A</i>	p.Met890Val	MS	0.0001

**Legend:** UT – Utah, Cases – total MM and MGUS cases (number of MGUS), Me – meioses, Position – build HG19, Len – length in mega-bases, p-value – SGS p-value, \*genome-wide significant, Conseq – exome-variant consequence, SG – stop gain variant, MS – missense variant, AAF – alternate allele frequency based on the non-TCGA, non-Finnish, European ExAC individuals.

1 We also identified two HRPs, UT-576834 and UT-260, with overlapping shared  
2 segments at 1p36.11 (Fig. 3). A 8.9 Mb (24.39-33.30 Mb,  $p = 3.0 \times 10^{-4}$ ) segment was observed  
3 in 3 of the 4 genotyped MM cases in UT-576834, shared across 12 meioses (Fig. 3b and Table  
4 2). A nested 1.2 Mb shared segment (26.22-27.38 Mb;  $p = 2.1 \times 10^{-4}$ ) segregated to 3 MM cases  
5 separated by 16 meioses in UT-260 (Fig. 3a and Table 2). The overlapping segment contains  
6 30 genes (Fig. 3d).

### 7 Exome follow-up of shared segments in HRPs

8 Whole-exome sequencing (WES) data was interrogated, targeted to the shared  
9 segments, to identify potential risk variants in the pedigree sharers in the HRP and in a broader  
10 set of 57 pedigrees. WES data was available for: 28 cases from the 11 extended Utah HRPs;  
11 and 162 exomes from 57 densely clustered MM/MGUS families from Mayo Clinic Rochester,



1 Weill Cornell, Memorial Sloan Kettering Cancer Center, International Agency for Research on  
2 Cancer, and INSERM France (S2 Table). Prioritization was used to identify variants that were:  
3 in the target segment; rare (alternate allele frequency, AAF<0.001), potentially deleterious  
4 (variant impact predicted to be high or moderate); and observed recurrently in the appropriate  
5 segment sharers (if observed in the segment discovery pedigree).

6 At 6q16, no rare, potentially deleterious coding risk variants were shared by the 3 UT-  
7 571744 MM cases in the 1.8 Mb genome-wide significant segment, indicating non-coding  
8 regulatory variants may be responsible for MM risk in this pedigree. However, two, rare coding  
9 and potentially deleterious single nucleotide variants (SNVs) were identified in two MM/MGUS  
10 families (Fig. 2c-e and Table 2). Both SNVs are in the hydrolase domain of *USP45*: a stop gain  
11 (p.Gln691\*) shared by 3 sibling cases (1 MM and 2 MGUS) in an INSERM family (PET-Nice  
12 0909) and a missense SNV (p.Gln621Glu) shared by 2 siblings (1 MM and 1 MGUS) but not  
13 their 2 screened unaffected siblings in Mayo family 485. Coverage of these positions in ExAC  
14 sequence data is high (> 99% of the 60,706 ExAC samples had at least 10x read coverage) and  
15 neither variant was observed.

16 Pedigree exomes in the 1.2 Mb segment at 1p36.11 revealed two, rare and potentially  
17 deleterious SNVs. The first in discovery pedigree UT-576834: a missense SNV (rs752026201,  
18 p.Ser90Gly, AAF = 0.0016 in ExAC) in *ARID1A* (Fig. 3e) shared by 3 of the 4 Utah MM cases,  
19 concordant with the segment sharing pattern. A second rare, missense SNV in *ARID1A*  
20 (rs140664170, p.Met890Val, AAF < 0.0001 in ExAC) was found to be carried by a pair of MM  
21 cousins in Weill-Cornell family 12 (Fig. 3c and e, and Table 2). Based on the ExAC data,  
22 *ARID1A* is extremely intolerant to missense variants ( $Z = 4.1$ ) and loss of function (LoF) SNVs  
23 (pLI = 1) [28].

## 1 **Pathway follow-up of candidate genes**

2 Our SGS findings and pedigree WES identify *USP45* and *ARID1A* as candidate genes  
3 for inherited MM risk. We further investigated shared segments and WES for evidence  
4 supporting the complexes *USP45* and *ARID1A* are involved in. Here we further expanded our  
5 WES to: 154 early-onset MM/MGUS cases from our collaborative group, 733 sporadic MM  
6 cases from dbGaP [29], and 964 controls [30].

7 *USP45* is an essential DNA repair regulator, de-ubiquitylating *ERCC1* to allow for DNA  
8 translocation of the *ERCC1-ERCC4* endonuclease [31,32]. This endonuclease is a part of the  
9 global genome nucleotide-excision repair (GG-NER) incision complex, a 22 protein complex  
10 essential to removing lesions from DNA and cancer prevention [33–36] (see S3 Table). We  
11 reviewed SGS results in the Utah HRP at the location of these 22 genes and identified a  
12 genome-wide suggestive segment in pedigree UT-34955 (S2 Fig.). This HRP identified a 0.8 Mb  
13 segment at 19q13 (45.71-46.51 Mb; hg19), containing 31 genes including *ERCC1* and *ERCC2*  
14 (S2 Fig. and S4 Table). The segment is shared by 3 MM cases separated by 16 meioses ( $p =$   
15  $6.6 \times 10^{-5}$ ). No rare, coding variants were identified from the WES in the 3 MM cases in UT-  
16 34955, nor in the remaining 67 pedigrees/families. We interrogated the 22 GG-NER incision  
17 complex genes in our 885 early-onset and sporadic MM exomes. This identified a ClinVar-  
18 annotated pathogenic, missense SNV in *ERCC4* (p.Arg799Trp) in one early-onset MM case and  
19 one sporadic MM case, and a stop-gain SNV in *ERCC3* (p.Arg574Ter), in the same domain as a  
20 ClinVar-annotated pathogenic variant, in a second early-onset MM case (S4 Table). Further,  
21 burden testing in all MM cases vs controls was significant in 7 of the 22 GG-NER genes:  
22 *ERCC3*, *PARP1*, *GTF2H1*, *GTF2H2*, *DDB1*, *RPA2*, and *CHD1L* (S3 Table).

23 *ARID1A* is a member of the SWI/SNF chromatin remodeling complex, a 15 gene  
24 complex involved in DNA transcription regulation [37] (see S5 Table). Members of this complex

1 are mutated in >20% of malignancies [38–40], but are extremely intolerant to LoF and missense  
2 variation [41] (S5 Table). We reviewed SGS results in the Utah HRP at the location of these 15  
3 genes and identified a marginal, genome-wide suggestive segment in pedigree UT-549917  
4 shared by 4 MM cases across 21 meioses ( $p = 2.17 \times 10^{-5}$ , S3 Fig. and S6 Table). This 1.5 Mb  
5 segment at chr3p21.1-p21.2 (52.01-53.56 Mb; hg19) contains 32 genes including *PBRM1* from  
6 the SWI/SNF complex. No coding variants were identified in this gene in UT-549917, nor in the  
7 remaining 67 pedigrees/families. Burden testing was significant for 8 of the 15 genes in the  
8 complex: *ARID1A*, *ARID1B*, *ARID2*, *SMARCA4*, *ACTL6A*, *SMARCD3*, *SMARCC2*, and  
9 *SMARCE1* (S5 Table).

## 10 DISCUSSION

11 We developed a novel strategy to identify segregating chromosomal segments shared  
12 by subsets of cases in HRP. It focuses on extended HRP that are singularly powerful to  
13 identify significant genetic segregation. Our strategy allows for genetic heterogeneity within such  
14 pedigrees and provides formal significance thresholds for valid interpretation. Previously,  
15 extended HRP have not delivered on their potential in complex traits because in common,  
16 complex traits, HRP are likely enriched for multiple susceptibility variants and may capture both  
17 familial and sporadic cases in their branches. Our optimization strategy over subsets is  
18 attractive because it allows for heterogeneity without prior knowledge of genetic similarities or  
19 deep phenotyping. Application of the method to extended MM pedigrees demonstrated the utility  
20 of this new method and illustrated that the segments identified were used successfully to narrow  
21 the search for risk variants in smaller pedigrees, allowing for an overall strategy that can utilize  
22 both large pedigrees and smaller families together for discovery (see Table 2, Fig. 1 and Fig. 2).  
23 Post-hoc, additional value can be gained from demographic and/or clinical data on the sharing

1 subsets shedding light on other shared characteristics that may aid future mapping. Also, we  
2 note that in the absence of any significant findings, genome-wide SGS results can be used as  
3 genomic annotations of segregation evidence for more heuristic approaches.

4 While we identified several rare, potentially deleterious coding variants of interest,  
5 several of the SGS discovery pedigrees had no coding variants that satisfied prioritization  
6 criteria. We believe this will be characteristic of complex traits and that regulatory variants will  
7 also play a substantial role. Mutations with strong causal likelihood found in other disease  
8 cohorts may focus the search for regulatory variation to particular genes within a shared  
9 segment, as with *USP45* in MM. In the absence of such compelling evidence, a return to  
10 pedigree segregation methods will provide identification of statistically compelling regions which  
11 can concentrate efforts to identify and characterize regulatory risk variants. Our proposed  
12 method is a new analytic tool with the potential to reinvigorate the use of extended HRP in the  
13 identification of risk variants that contribute to common, complex disease.

14 Multiple myeloma is a malignancy of the plasma cells that has been shown to be familial  
15 [42]. Consistent with a role for genetics, case-control studies have been successful in identifying  
16 association signals for 17 low-risk variants [43–47]. However, despite consistent evidence for  
17 familial clustering, our study is the first to explore high-risk MM pedigrees. Using the unique  
18 genealogical database available in Utah, we identified and studied extended MM HRP. We  
19 identified a genome-wide significant segment containing *USP45*, an important regulator of DNA  
20 repair (see Fig. 1 and Table 2), and a genome-wide suggestive segment harboring other genes  
21 in the endonuclease regulation pathway (*ERCC1* and *ERCC2*). Exome sequencing in a  
22 collaborative resource of high-risk families and early-onset cases revealed four rare, potentially  
23 deleterious coding variants; two novel variants in *USP45* segregating in two pedigrees and two  
24 variants in early-onset cases in *ERCC3* and *ERCC4*, the latter annotated as pathogenic in  
25 ClinVar. Burden testing including sporadic MM, and comparing to controls, identified significant

1 enrichment for variants in MM cases in 7 of the 22 GG-NER genes in the protein endonuclease  
2 regulation complex.

3 In particular, the functional literature supports *USP45* as a candidate cancer risk gene.  
4 *USP45* has been shown to deubiquitylate ERCC1, a catalytic subunit of the ERCC1-ERCC4  
5 DNA repair endonuclease (ERCC4 also known as XPF) [31]. This endonuclease is a critical  
6 regulator of DNA repair processes [34]. The complex repairs recombination, double strand  
7 break, and inter-strand crosslink by cutting DNA overhangs around a lesion, degrades 3' G-rich  
8 overhangs in telomere maintenance, and plays a role in cancer prevention and in tumor  
9 resistance to chemotherapy [31,34]. Mouse models have shown *USP45* knockout cells have  
10 higher levels of ubiquitylated ERCC1 and that cells are hypersensitive to UV radiation and DNA  
11 inter-strand cross-links, repair of UV-induced DNA damage, and ERCC1 translocation to DNA  
12 damage is impaired [31]. Hence, the deubiquitylase activity of *USP45* is important for  
13 maintaining the DNA repair ability of ERCC1-ERCC4. In total, these observations implicate the  
14 GG-NER pathway and specifically the interaction of *USP45* and the disruption of the ERCC1-  
15 ERCC4 role in DNA repair as a mechanism of potential importance in MM risk.

16 Our strategy also identified shared segments overlapping at chr1p36.11 in two Utah  
17 pedigrees containing *ARID1A* (Fig. 2) and a genome-wide suggestive segment in a third  
18 pedigree harboring another gene in the SWI/SNF complex (*PBRM1*). For the SWI/SNF  
19 complex, exome sequencing revealed two rare, potentially deleterious variants in *ARID1A*  
20 segregating in two pedigrees. Burden testing provided further evidence for enrichment of  
21 variants in *ARID1A* specifically, and in 7 of the 15 genes in the complex. As a component of the  
22 SWI/SNF chromatin remodeling complex, *ARID1A* facilitates gene activation by assisting  
23 transcription machinery gain access to gene targets [48]. Based on the patterns of mutations in  
24 tumor cells, *ARID1A* likely functions as a tumor-suppressor [49]. Members of the SWI/SNF  
25 chromatin remodeling complexes are mutated in 20% of malignancies [38], but are extremely

1 intolerant to LoF and missense variation [41] (see S5 Table). Blockage of chromatin remodeling  
2 may sustain cancer development [39]. Aberrant chromatin remodeling contributes to the  
3 pathogenesis of ovarian clear-cell carcinoma [49]. It has previously been shown that *ARID1A* is  
4 intolerant to variation (LoF and missense mutations) [28], consistent with its prominent somatic  
5 role in multiple tumors [38,49,50], including hematological malignancies [51–53]. These  
6 observations implicate the SWI/SNF chromatin remodeling complex, and specifically *ARID1A* in  
7 MM risk.

8         This study has limitations. First, the method is applicable only to extended HRP that are  
9 singularly effective for identifying segregating segments (15 meioses between cases is optimal  
10 [16]). The method is not directly applicable to the many smaller family-based resources that  
11 have been gathered in the complex trait field and may therefore result in findings from single  
12 large pedigrees that are private and difficult to replicate. However, as illustrated in our example,  
13 in a collaborative setting containing both extended HRP and smaller families, the approach can  
14 be mutually beneficial. Second, our observation of two borderline genome-wide suggestive  
15 overlapping segments at 1p36 led to our identification of *ARID1A* as a potential candidate risk  
16 gene and illustrates the potential for discoveries using overlapping subthreshold evidence.  
17 However, it raises analytical questions of how to systematically identify such segments. This  
18 segment would have been ignored based on strict individual-pedigree thresholds and highlights  
19 an important area for further methodological development. Finally, this study is observational  
20 and cannot describe causation. We have identified two complexes, several genes and specific  
21 variants as compelling candidates involved in MM risk, but further functional studies will be  
22 required to determine and characterize the mechanisms involved in risk.

23         In conclusion, we have developed a strategy for gene mapping in complex traits that  
24 accounts for heterogeneity within HRP and formally corrects for multiple testing to allow for  
25 statistically rigorous discovery. We applied this strategy to MM, a complex cancer of plasma

1 cells, and identified multiple shared segments containing genes in nucleotide excision repair  
2 and SWI/SNF chromatin remodeling. Exome follow-up supported these segments in both the  
3 Utah large HRP and smaller families from other sites. Our study offers a novel technique for  
4 HRP gene mapping and demonstrates its utility to narrow the search for risk-variants in complex  
5 traits.

## 6 **METHODS**

### 7 **SGS Analysis in Utah, Myeloma HRPs**

8 **HRPs and genotyping.** All participants were studied with informed consent under protocols  
9 approved by the University of Utah IRB. Using the statewide Utah Cancer Registry (UCR), all  
10 living individuals with MM in Utah were invited to participate and peripheral blood was collected  
11 for DNA extraction. Participants were linked in the Utah Population Database (UPDB), a unique  
12 resource that integrates UCR records with a 5M person genealogy. HRPs were defined as  
13 pedigrees containing statistical excess of MM ( $p < 0.05$ ), based on sex and cohort-specific rates  
14 in Utah. Eleven of the HRPs identified in the UPDB contained 3-4 MM cases with DNA (total  
15 MM cases per pedigree ranged from 4 to 37) with 8-23 meioses between studied MM cases.  
16 DNA from the 28 cases was genotyped on the Illumina Omni Express high-density SNP array.

17 **Genotype quality control.** Only bi-allelic SNPs were considered. Genotypes and individual  
18 call-rates were used to ensure high quality data. PLINK was used to remove SNPs with  $< 95\%$   
19 call rate across individuals [54]. The final SNP set contained 678,447 single nucleotide variants.  
20 After SNP removal for low call rates, individuals were removed based on  $< 90\%$  call rate across  
21 the genome, or if they failed the PLINK sex check. One MM case was removed. The QC'ed  
22 SNP data were transformed to match strand orientation of the 1000Genomes.

23

1 **Probability of sharing a segment.** SGS analysis identifies contiguous SNPs that are  
2 shared identical-by-state (IBS) by cases in a HRP and assigns an empirical probability of  
3 chance ancestral sharing [26]. First, a set of cases in a HRP are defined and all segments of  
4 contiguous SNPs shared IBS are identified. Only shared segments > 20 SNPs are considered  
5 as lengths shorter than 20 are commonly shared between unrelated individuals. Second,  
6 population-based data (here we used CEU and GBR data from the 1000Genomes Project [55])  
7 are used to estimate a graphical model for linkage disequilibrium (LD) [56], providing a  
8 probability distribution of chromosome-wide haplotypes in the population. Third, pairs of  
9 haplotypes are randomly assigned to pedigree founders according to the haplotype distribution.  
10 Founders are individuals whose parents are not specified in the pedigree. For chromosome-  
11 wide haplotype simulations the full chromosome LD model is used. Fourth, Mendelian  
12 segregation and recombination are simulated to generate genotypes for all pedigree members.  
13 The Rutgers genetic map [57] is used for a genetic map for recombination, with interpolation  
14 based on physical base pair position for SNPs not represented. Steps two through four create  
15 one simulated data set, a random sample from the null hypothesis. This process is repeated  
16 hundreds of thousands to millions of times.

17 Each shared segment in the real data (step one) is compared to the simulated segments  
18 at the precise genomic location. The number of times the null segment equals or encompasses  
19 the observed segment is counted and divided by the total number of simulations to generate the  
20 empirical nominal p-value for the observed shared segment. The simulations continue until a p-  
21 value has been estimated to a required resolution, or until it surpasses a defined significance  
22 threshold. To facilitate this in an efficient manner, we follow-up specific segments using  
23 marginal distributions from the LD model, established using standard graphical modeling  
24 methods [58]. The marginalized LD model encompassing only the region of interest, but



1 capturing relevant LD to accurately simulate genotypes from this region alone. This reduction in  
2 markers vastly increases the speed in which simulations are generated. The graphical model  
3 estimation, marginalization, and simulation processes are computationally efficient requiring  
4 time and storage that is linear with the number of SNPs being considered.

5 **Heterogeneity optimization.** We systematically perform SGS analysis on each subset of  
6 cases in a HRP. At each marker position across the genome, the optimized segment is the one  
7 minimizing the p-value across all subsets considered. All segments selected by the optimization  
8 procedure, and their respective p-values, comprise the final optimized SGS results.

9 **Significance threshold determination.** A transformation,  $Y = -\log_{10}(p)$  is performed to  
10 the optimized genome-wide SGS p-value vector. The results are fit to a gamma distribution  
11 using the MLE method.  $Y \sim \Gamma(k, \sigma)$  ( $k$  shape,  $\sigma$  scale parameterization). The Theory of Large  
12 Deviations has previously been used in pedigree studies to model extreme values in a genome  
13 wide genetic setting [27], and it has been shown that for a statistic following a Gaussian  
14 distribution, the number of segments where the statistic exceeds a threshold  $W$  has mean:

$$15 \quad \mu(W) = [C + 2\rho GW^2]\alpha(W) \quad (2),$$

16 where  $\alpha(W)$  is the pointwise significance level of exceeding  $W$ ,  $C$  is the number of  
17 chromosomes considered,  $\rho$  reflects the recombination rate ( $\rho = 1$  for general pedigrees), and  
18  $G$  is genetic length in Morgans. Lander & Kruglyak demonstrated that the same equation  
19 extends a statistic following the chi-squared distribution:

$$20 \quad \mu(X) = [C + 2\rho GX]\alpha(X) \quad (3),$$

21 based on the distributional relationship between the chi-squared and Normal distributions  $W^2 =$   
22  $X$ . Here, we use the distributional relationship between the gamma and chi-square distributions,

1 our estimated  $k$  and  $\sigma$  gamma parameters, where  $T = 10^{-X\sigma/2}$ ,  $X = 2Y/\sigma \sim \chi_{2k}^2$ , and the genetic  
2 length of the genome (matched to that used in the gene-drop) to utilize Eq. 3 and derive  $\mu(X)$   
3 thresholds. Solving for  $\mu(X) = 0.05$  and  $\mu(X) = 1$  produced significance and suggestive  
4 thresholds, respectively. These thresholds are remarkably stable after a few hundred thousand  
5 simulations. For pedigrees with very large numbers of meioses (>50) between the full case-set  
6 a larger number of simulations may be required.

### 7 **Interrogating target regions defined by shared segments in WES**

8 **Participants.** WES data were interrogated in the regions defined by the shared segments of  
9 interest. WES data was available on 964 controls [30] and 1,064 MM or MGUS cases including:  
10 28 MM from the 11 Utah HRPs; 70 MM and 79 MGUS from 57 densely clustered families (each  
11 containing at least 2 MM/MGUS); 154 sporadic or early-onset MM/MGUS cases; and 733  
12 sporadic MM cases from dbGaP [29]. Of the 57 densely clustered MM/MGUS families, 37 were  
13 ascertained by INSERM, France (36 MM, 69 MGUS, 2 controls), 10 by Mayo Clinic, Minnesota  
14 (10 MM, 10 MGUS, 11 controls), 6 by Memorial Sloan Kettering Cancer Center, New York (14  
15 MM, 0 MGUS, 0 controls), 3 by International Agency for Research on Cancer, France (8 MM, 0  
16 MGUS, 0 controls), and 1 by Weill Cornell, New York (2 MM, 0 MGUS, 0 controls). Most of the  
17 families had both MM and MGUS cases (32 families total), while 13 families only had MGUS  
18 and 12 families only had MM cases sequenced. Nine families had at least one unaffected  
19 relative sequenced as a control. (See S2 Table.)

1 **Joint calling analysis.** To perform joint calling of all of the exome sequences, we utilized the  
2 calling pipeline developed at the Icahn School of Medicine at Mt. Sinai, based on GATK Best  
3 Practices [59]. Briefly, fastq files were aligned to genome build 37 using bwa version 0.7.8,  
4 indels were realigned using GATK, duplicates were removed using Picard MarkDuplicates, and  
5 base quality scores were recalibrated using GATK. HaplotypeCaller was then used to generate  
6 individual GVCF files for each individual, and GenotypeGVCFs was used to generate the final  
7 joint calling. The jointly-called VCF was annotated with SNPEff and loaded into a GEMINI  
8 database for easy of querying. Additional functional annotations available in the GEMINI suite  
9 include CADD, ANNOVAR, conservation, location, and if the variant was listed in OMIM.

10 **Variant prioritization.** A GEMINI query was developed to select variants with high or medium  
11 impact, AAF < 0.001, in the shared segments of interest, and shared by the MM cases is the  
12 discovery pedigree, or identify other sequence variants in candidate gene across collaborating  
13 and publicly available exomes.

14 **Burden testing.** Burden testing was performed on jointly called and processed WES from  
15 1,064 MM/MGUS cases and 964 unaffected controls for the 22 genes in the GG-NER incision  
16 complex and 15 genes in the SWI/SNF chromatin remodeling complex. The GEMINI software  
17 [60] was used to perform a c-alpha test [61] with 1000 permutations. Only variants with AAF <  
18 0.05 and high or moderate predicted impact were included in the analysis.

## 19 **ACKNOWLEDGMENTS**

20 This work was supported in part by the DNA Sequencing Core Facility and Genomics Core  
21 Facility at the University of Utah, and through the computational resources and staff expertise

1 provided by Scientific Computing at the Icahn School of Medicine at Mount Sinai. Data  
2 collection was made possible, in part, by the Utah Population Database and the Utah Cancer  
3 Registry. We thank the participants and their families who make this research possible.

#### 4 **REFERENCES**

- 5 1. Maher B. Personal genomes: The case of the missing heritability. *Nature*. 2008;456: 18–  
6 21. doi:10.1038/456018a
- 7 2. Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorff LA, Hunter DJ, et al. Finding the  
8 missing heritability of complex diseases. *Nature*. 2009;461: 747–53.  
9 doi:10.1038/nature08494
- 10 3. Eichler EE, Flint J, Gibson G, Kong A, Leal SM, Moore JH, et al. Missing heritability and  
11 strategies for finding the underlying causes of complex disease. *Nat Rev Genet*.  
12 2010;11: 446–50. doi:10.1038/nrg2809
- 13 4. Miki Y, Swensen J, Shattuck-Eidens D, Futreal PA, Harshman K, Tavtigian S, et al. A  
14 strong candidate for the breast and ovarian cancer susceptibility gene BRCA1. *Science*.  
15 1994;266: 66–71. Available: <http://www.ncbi.nlm.nih.gov/pubmed/7545954>
- 16 5. Wooster R, Neuhausen SL, Mangion J, Quirk Y, Ford D, Collins N, et al. Localization of  
17 a breast cancer susceptibility gene, BRCA2, to chromosome 13q12-13. *Science*.  
18 1994;265: 2088–90. Available: <http://www.ncbi.nlm.nih.gov/pubmed/8091231>
- 19 6. Vance JM, Pericak-Vance MA, Yamaoka LH, Speer MC, Rosenwasser GO, Small K, et  
20 al. Genetic linkage mapping of chromosome 17 markers and neurofibromatosis type I.  
21 *Am J Hum Genet*. 1989;44: 25–9. Available:  
22 <http://www.ncbi.nlm.nih.gov/pubmed/2491777>

- 1 7. Cannon-Albright LA, Goldgar DE, Meyer LJ, Lewis CM, Anderson DE, Fountain JW, et  
2 al. Assignment of a locus for familial melanoma, MLM, to chromosome 9p13-p22.  
3 Science. 1992;258: 1148–52. Available: <http://www.ncbi.nlm.nih.gov/pubmed/1439824>
- 4 8. Leppert M, Dobbs M, Scambler P, O’Connell P, Nakamura Y, Stauffer D, et al. The gene  
5 for familial polyposis coli maps to the long arm of chromosome 5. Science. 1987;238:  
6 1411–3. Available: <http://www.ncbi.nlm.nih.gov/pubmed/3479843>
- 7 9. Nishisho I, Nakamura Y, Miyoshi Y, Miki Y, Ando H, Horii A, et al. Mutations of  
8 chromosome 5q21 genes in FAP and colorectal cancer patients. Science. 1991;253:  
9 665–9. Available: <http://www.ncbi.nlm.nih.gov/pubmed/1651563>
- 10 10. Ng SB, Buckingham KJ, Lee C, Bigham AW, Tabor HK, Dent KM, et al. Exome  
11 sequencing identifies the cause of a mendelian disorder. Nat Genet. 2010;42: 30–5.  
12 doi:10.1038/ng.499
- 13 11. Ng SB, Bigham AW, Buckingham KJ, Hannibal MC, McMillin MJ, Gildersleeve HI, et al.  
14 Exome sequencing identifies MLL2 mutations as a cause of Kabuki syndrome. Nat  
15 Genet. 2010;42: 790–3. doi:10.1038/ng.646
- 16 12. McClellan J, King M-C. Genetic Heterogeneity in Human Disease. Cell. 2010;141: 210–  
17 217. doi:10.1016/j.cell.2010.03.032
- 18 13. Mitchell KJ. What is complex about complex disorders? Genome Biol. 2012;13: 237.  
19 doi:10.1186/gb-2012-13-1-237
- 20 14. Li X, Montgomery SB. Detection and Impact of Rare Regulatory Variants in Human  
21 Disease. Front Genet. 2013;4. doi:10.3389/fgene.2013.00067
- 22 15. Albert FW, Kruglyak L. The role of regulatory variation in complex traits and disease. Nat  
23 Rev Genet. 2015;16: 197–212. doi:10.1038/nrg3891

- 1 16. Knight S, Abo RP, Abel HJ, Neklason DW, Tuohy TM, Burt RW, et al. Shared Genomic  
2 Segment Analysis: The Power to Find Rare Disease Variants. *Ann Hum Genet.* 2012;76:  
3 500–509. doi:10.1111/j.1469-1809.2012.00728.x
- 4 17. Myeloma - SEER Stat Fact Sheets [Internet]. Available:  
5 <https://seer.cancer.gov/statfacts/html/mulmy.html>
- 6 18. Cannon-Albright LA, Thomas A, Goldgar DE. Familiality of cancer in Utah. *Cancer Res.*  
7 1994;54: 2378–2385.
- 8 19. Landgren O, Linet MS, McMaster ML, Gridley G, Hemminki K, Goldin LR. Familial  
9 characteristics of autoimmune and hematologic disorders in 8,406 multiple myeloma  
10 patients: A population-based case-control study. *Int J Cancer.* 2006;118: 3095–3098.  
11 doi:10.1002/ijc.21745
- 12 20. Albright F, Teerlink C, Werner TL, Cannon-Albright L a. Significant evidence for a  
13 heritable contribution to cancer predisposition: a review of cancer familiality by site. *BMC*  
14 *Cancer.* BioMed Central Ltd; 2012;12: 138. doi:10.1186/1471-2407-12-138
- 15 21. Schinasi LH, Brown EE, Camp NJ, Wang SS, Hofmann JN, Chiu BC, et al. Multiple  
16 myeloma and family history of lymphohaematopoietic cancers: Results from the  
17 International Multiple Myeloma Consortium. *Br J Haematol.* England; 2016;175: 87–101.  
18 doi:10.1111/bjh.14199
- 19 22. Landgren O, Kristinsson SY, Goldin LR, Caporaso NE, Blimark C, Mellqvist U-H, et al.  
20 Risk of plasma cell and lymphoproliferative disorders among 14621 first-degree relatives  
21 of 4458 patients with monoclonal gammopathy of undetermined significance in Sweden.  
22 *Blood.* 2009;114: 791–5. doi:10.1182/blood-2008-12-191676
- 23 23. Greenberg AJ, Rajkumar SV, Vachon CM. Familial monoclonal gammopathy of  
24 undetermined significance and multiple myeloma: epidemiology, risk factors, and

- 1 biological characteristics. *Blood*. 2012;119: 5359–66. doi:10.1182/blood-2011-11-  
2 387324
- 3 24. Greenberg AJ, Rajkumar SV, Larson DR, Dispenzieri A, Therneau TM, Colby CL, et al.  
4 Increased prevalence of light chain monoclonal gammopathy of undetermined  
5 significance (LC-MGUS) in first-degree relatives of individuals with multiple myeloma. *Br*  
6 *J Haematol*. 2012;157: 472–5. Available: <http://www.ncbi.nlm.nih.gov/pubmed/22629552>
- 7 25. Vachon CM, Kyle RA, Therneau TM, Foreman BJ, Larson DR, Colby CL, et al.  
8 Increased risk of monoclonal gammopathy in first-degree relatives of patients with  
9 multiple myeloma or monoclonal gammopathy of undetermined significance. *Blood*.  
10 2009;114: 785–90. doi:10.1182/blood-2008-12-192575
- 11 26. Thomas A, Camp NJ, Farnham JM, Allen-Brady K, Cannon-Albright LA. Shared  
12 Genomic Segment Analysis. Mapping Disease Predisposition Genes in Extended  
13 Pedigrees Using SNP Genotype Assays. *Ann Hum Genet*. 2008;72: 279–287.  
14 doi:10.1111/j.1469-1809.2007.00406.x
- 15 27. Lander E, Kruglyak L. Genetic dissection of complex traits: guidelines for interpreting  
16 and reporting linkage results. *Nat Genet*. 1995;11: 141–147.
- 17 28. Lek M, Karczewski KJ, Minikel E V., Samocha KE, Banks E, Fennell T, et al. Analysis of  
18 protein-coding genetic variation in 60,706 humans. *Nature*. *Nature Research*; 2016;536:  
19 285–291. doi:10.1038/nature19057
- 20 29. Myeloma data downloaded from the dbGaP web site under accessions:  
21 phs000348.v2.p1 and phs000748.v4.p3. [Internet].
- 22 30. Control data downloaded from the dbGaP web site under accessions:  
23 phs000209.v13.p3, phs000276.v2.p1, phs000179.v5.p2, phs000298.v3.p2,  
24 phs000424.v6.p1, phs000653.v2.p1, phs000687.v1.p1, phs000814.v1.p1, and  
25 phs000806.v1.p1.

- 1 31. Perez-Oliva AB, Lachaud C, Szyniarowski P, Muñoz I, Macartney T, Hickson I, et al.  
2 USP45 deubiquitylase controls ERCC1-XPF endonuclease-mediated DNA damage  
3 responses. *EMBO J.* 2015;34: 326–43. doi:10.15252/embj.201489184
- 4 32. USP45 in the GG-NER Pathway [Internet]. Available:  
5 <http://www.reactome.org/PathwayBrowser/#/R-HSA-5696398&SEL=R-HSA->  
6 [5696465&PATH=R-HSA-73894](http://www.reactome.org/PathwayBrowser/#/R-HSA-5696398&SEL=R-HSA-5696465&PATH=R-HSA-73894)
- 7 33. Marteiijn JA, Lans H, Vermeulen W, Hoeijmakers JHJ. Understanding nucleotide excision  
8 repair and its roles in cancer and ageing. *Nat Rev Mol Cell Biol.* 2014;15: 465–81.  
9 doi:10.1038/nrm3822
- 10 34. Kirschner K, Melton DW. Multiple roles of the ERCC1-XPF endonuclease in DNA repair  
11 and resistance to anticancer drugs. *Anticancer Res.* 2010;30: 3223–3232. doi:30/9/3223  
12 [pii]
- 13 35. Friedberg EC. How nucleotide excision repair protects against cancer. *Nat Rev Cancer.*  
14 2001;1: 22–33. doi:10.1038/35094000
- 15 36. Christmann M, Tomicic MT, Roos WP, Kaina B. Mechanisms of human DNA repair: an  
16 update. *Toxicology.* 2003;193: 3–34. Available:  
17 <http://www.ncbi.nlm.nih.gov/pubmed/14599765>
- 18 37. SWI/SNF Chromatin Remodeling Complex [Internet]. Available:  
19 <http://www.reactome.org/PathwayBrowser/#/R-HSA-5696398&PATH=R-HSA-73894>
- 20 38. Biegel JA, Busse TM, Weissman BE. SWI/SNF chromatin remodeling complexes and  
21 cancer. *Am J Med Genet C Semin Med Genet.* 2014;166C: 350–66.  
22 doi:10.1002/ajmg.c.31410
- 23 39. Romero O a, Sanchez-Cespedes M. The SWI/SNF genetic blockade: effects in cell  
24 differentiation, cancer and developmental diseases. *Oncogene.* Nature Publishing  
25 Group; 2014;33: 2681–9. doi:10.1038/onc.2013.227



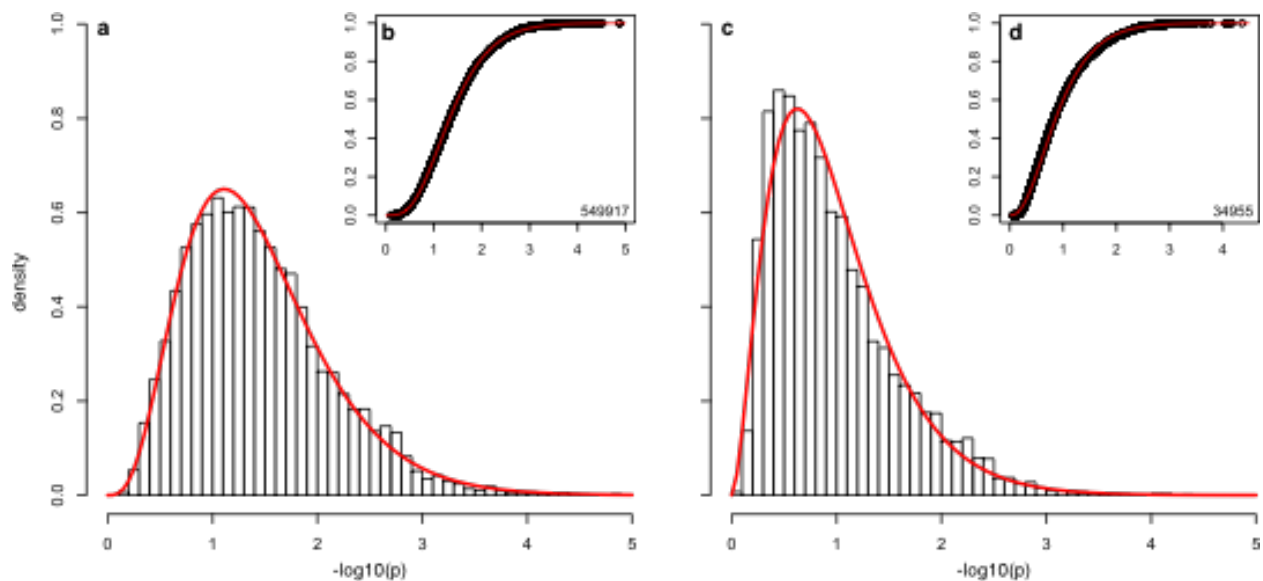
- 1 40. Roberts CWM, Orkin SH. The SWI/SNF complex - chromatin and cancer. Nat Rev  
2 Cancer. Nature Publishing Group; 2004;4: 133–142. Available:  
3 <http://dx.doi.org/10.1038/nrc1273>
- 4 41. Lek M. Analysis of protein-coding genetic variation in 60,706 humans. 2015;  
5 doi:10.1101/030338
- 6 42. Morgan GJ, Johnson DC, Weinhold N, Goldschmidt H, Landgren O, Lynch HT, et al.  
7 Inherited genetic susceptibility to multiple myeloma. Leukemia. 2014;28: 518–24.  
8 doi:10.1038/leu.2013.344
- 9 43. Broderick P, Chubb D, Johnson DC, Weinhold N, Försti A, Lloyd A, et al. Common  
10 variation at 3p22.1 and 7p15.3 influences multiple myeloma risk. Nat Genet. Nature  
11 Publishing Group; 2011;44: 58–61. doi:10.1038/ng.993.Common
- 12 44. Chubb D, Weinhold N, Broderick P, Chen B, Johnson DC, Försti A, et al. Common  
13 variation at 3q26.2, 6p21.33, 17p11.2 and 22q13.1 influences multiple myeloma risk. Nat  
14 Genet. Nature Publishing Group; 2013;45: 1221–1225. doi:10.1038/ng.2733
- 15 45. Weinhold N, Johnson DC, Chubb D, Chen B, Försti A, Hosking FJ, et al. The CCND1  
16 c.870G>A polymorphism is a risk factor for t(11;14)(q13;q32) multiple myeloma. Nat  
17 Genet. 2013;45: 522–5. doi:10.1038/ng.2583
- 18 46. Swaminathan B, Thorleifsson G, Jöud M, Ali M, Johnsson E, Ajore R, et al. Variants in  
19 ELL2 influencing immunoglobulin levels associate with multiple myeloma. Nat Commun.  
20 2015;6: 7213. doi:10.1038/ncomms8213
- 21 47. Mitchell JS, Li N, Weinhold N, Försti A, Ali M, Duin M Van, et al. Genome-wide  
22 association study identifies multiple susceptibility loci for multiple myeloma. Nat  
23 Commun. 2016;7: 12050. doi:10.1038/ncomms12050

- 1 48. Nie Z, Xue Y, Yang D, Zhou S, Deroo BJ, Archer TK, et al. A specificity and targeting  
2 subunit of a human SWI/SNF family-related chromatin-remodeling complex. *Mol Cell*  
3 *Biol.* 2000;20: 8879–88. Available: <http://www.ncbi.nlm.nih.gov/pubmed/11073988>
- 4 49. Jones S, Wang T-L, Shih I-M, Mao T-L, Nakayama K, Roden R, et al. Frequent  
5 mutations of chromatin remodeling gene ARID1A in ovarian clear cell carcinoma.  
6 *Science.* 2010;330: 228–31. doi:10.1126/science.1196333
- 7 50. Hodges C, Kirkland JG, Crabtree GR. The Many Roles of BAF (mSWI/SNF) and PBAF  
8 Complexes in Cancer. *Cold Spring Harb Perspect Med.* 2016;6.  
9 doi:10.1101/cshperspect.a026930
- 10 51. Puente XS, Beà S, Valdés-Mas R, Villamor N, Gutiérrez-Abril J, Martín-Subero JI, et al.  
11 Non-coding recurrent mutations in chronic lymphocytic leukaemia. *Nature.* 2015;526:  
12 519–524. doi:10.1038/nature14666
- 13 52. Lunning MA, Green MR. Mutation of chromatin modifiers; an emerging hallmark of  
14 germinal center B-cell lymphomas. *Blood Cancer J.* 2015;5: e361.  
15 doi:10.1038/bcj.2015.89
- 16 53. Choi J, Goh G, Walradt T, Hong BS, Bunick CG, Chen K, et al. Genomic landscape of  
17 cutaneous T cell lymphoma. *Nat Genet.* 2015;47: 1–11. doi:10.1038/ng.3356
- 18 54. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: A  
19 Tool Set for Whole-Genome Association and Population-Based Linkage Analyses. *Am J*  
20 *Hum Genet.* 2007;81: 559–575. doi:10.1086/519795
- 21 55. 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang  
22 HM, et al. A global reference for human genetic variation. *Nature.* 2015;526: 68–74.  
23 doi:10.1038/nature15393

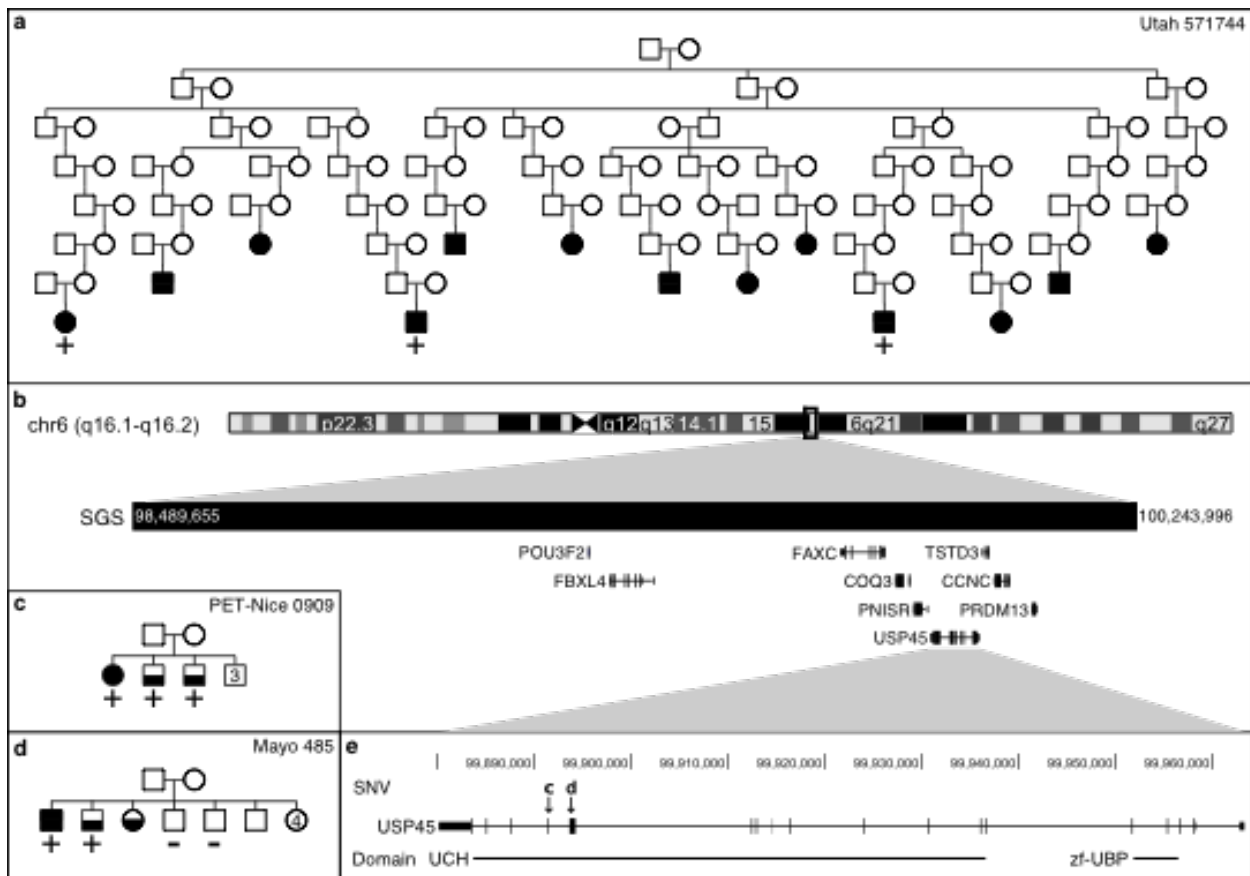
- 1 56. Abel HJ, Thomas A. Accuracy and Computational Efficiency of a Graphical Modeling  
2 Approach to Linkage Disequilibrium Estimation. *Stat Appl Genet Mol Biol*. 2011;10.  
3 doi:10.2202/1544-6115.1615
- 4 57. Matisse TC, Chen F, Chen W, De La Vega FM, Hansen M, He C, et al. A second-  
5 generation combined linkage physical map of the human genome. *Genome Res*.  
6 2007;17: 1783–6. doi:10.1101/gr.7156307
- 7 58. Lauritzen SL. *Graphical models*. Clarendon Press; 1996.
- 8 59. Linderman MD, Brandt T, Edelmann L, Jabado O, Kasai Y, Kornreich R, et al. Analytical  
9 validation of whole exome and whole genome sequencing for clinical applications. *BMC*  
10 *Med Genomics*. 2014;7: 20. doi:10.1186/1755-8794-7-20
- 11 60. Paila U, Chapman BA, Kirchner R, Quinlan AR. GEMINI: Integrative Exploration of  
12 Genetic Variation and Genome Annotations. Gardner PP, editor. *PLoS Comput Biol*.  
13 2013;9: e1003153. doi:10.1371/journal.pcbi.1003153
- 14 61. Neale BM, Rivas MA, Voight BF, Altshuler D, Devlin B, Orho-Melander M, et al. Testing  
15 for an unusual distribution of rare variants. *PLoS Genet*. 2011;7.  
16 doi:10.1371/journal.pgen.1001322

## 1 FIGURES

2 **Fig. 1. Adequacy of the gamma distribution.** The gamma distribution provides an adequate fit  
3 for multiple types of pedigrees. For example, HRP 549917 has gamma shape = 4.4 and rate =  
4 3.6 with good visual density (a) and CDF (b) fit, with lambda = 0.9. HRP 34955 has gamma  
5 shape = 2.8 and rate = 2.9 with good visual density (c) and CDF (d) fit, with lambda = 1.0.



1 **Fig. 2. Significant SGS, pedigrees, and segregating SNVs.** In pedigrees, MM cases are fully  
 2 shaded and MGUS cases are half shaded. Numbers indicate multiple individuals. a) Utah  
 3 pedigree sharing the genome-wide significant SGS. The three genotyped MM cases that are  
 4 SGS carriers are marked by +. Note. The genealogy extends beyond SEER cancer registry  
 5 data. MGUS are unknown in this pedigree. b) Genomic region of significant SGS. c) INSERM  
 6 pedigree carrying the stop gain SNV marked by “c” in box e. 1 MM and 2 MGUSs carry the  
 7 SNV. d) Mayo Clinic pedigree carrying the missense SNV marked by “d” in box e. 1 MM and 1  
 8 MGUS carry the SNV, but not 2 unaffected siblings. e) Risk candidate gene, *USP45*, has 2  
 9 segregating SNVs in the ubiquitin C-terminal hydrolase 2 (UCH) domain.



1 **Fig. 3. SGS with multiple lines of evidence.** a/b) Utah pedigrees carrying the overlapping  
2 SGSs on chr1p36.11-p35.1. + indicates the genotyped MM cases that are SGS carriers, -  
3 indicates genotyped and non-carriers. c) Weill Cornell pedigree with a segregating, missense  
4 SNV in *ARID1A* indicated by “c” in e. d) Genomic region of overlapping SGS. Dark black genes  
5 fall in both regions. e) 2 rare and segregating, missense SNVs were observed in whole-exome  
6 sequencing. SNV “b” is carried by the cases indicated with + in box b. SNV “c” is carried by the  
7 cases in box c.

1

