

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

3

4 Rosalie G. Waller<sup>1¶</sup>, Todd M. Darlington<sup>1¶</sup>, Xiaomu Wei<sup>2</sup>, Michael J. Madsen<sup>1</sup>, Alun Thomas<sup>1</sup>,  
5 Karen Curtin<sup>1</sup>, Hilary Coon<sup>1</sup>, Venkatesh Rajamanickam<sup>1</sup>, Justin Musinsky<sup>3</sup>, David Jayabalan<sup>2</sup>,  
6 Djordje 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  
7 Galia<sup>7</sup>, 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>,  
8 Robert 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>

9

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

18

19 \*Corresponding author

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

21

22 <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 apply this new Shared Genomic Segment  
7 (SGS) method in 11 extended, Utah, multiple myeloma (MM) HRPs, and subsequent exome  
8 sequencing in SGS regions of interest in 1063 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

## 19 **AUTHOR SUMMARY**

20           Although family-based studies demonstrate inherited variants play a role in many  
21 common and complex diseases, finding the genes responsible remains a challenge. High-risk  
22 pedigrees, or families with more disease than expected by chance, have been helpful in the  
23 discovery of variants responsible for less complex diseases, but have not reached their potential  
24 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 apply  
4 our method in Multiple Myeloma, a heritable and complex cancer of plasma cells. We identified  
5 two genes *USP45* and *ARID1A* that fall within shared segments with compelling statistical  
6 evidence. Exome sequencing of these genes revealed likely-damaging variants inherited in  
7 Myeloma high-risk families, suggesting these genes likely play a role in development of  
8 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

## 14 **INTRODUCTION**

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

23 Here we develop a HRP strategy based on our previous Shared Genomic Segment  
24 (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 55 multiplex MM or MM/MGUS pedigrees to perform subsequent  
8 targeted searches at the variant level. MM is a complex cancer of the plasma cells with 30,330  
9 new cases annually (incidence 6.5/100,000 per year) [17]. Despite survival dramatically  
10 increasing from 25.8% in 1980 to 48.5% in 2012, MM remains a cancer with one of the lowest 5-  
11 year survival rates in adult hematological malignancies [17]. MM is preceded by a condition  
12 referred to as monoclonal gammopathy of undetermined significance (MGUS). Evidence for the  
13 familial clustering of MM is consistently replicated [18–21], as is its clustering with MGUS [22–  
14 25]. Genetic pedigree studies in MM are scarce as it remains a challenge to acquire samples in  
15 pedigrees due to rarity and low survival rates. The Utah MM HRPs are one of only a few  
16 pedigree resources worldwide and contains unparalleled multi-generational high-risk pedigrees.  
17 Thus far, no segregating risk variants have been identified for MM.

18

## 19 **RESULTS**

### 20 **Pedigree analysis strategy**

21 We developed a gene mapping strategy, based on the SGS method [16,26], that  
22 accounts for intra-familial heterogeneity and multiple testing. The basic SGS method identifies  
23 all genomic segments shared identical-by-state (sharing without regard to inheritance) between  
24 a 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 at every position throughout the genome are identified and nominal p-values  
15 assigned. Across subsets, an optimization procedure is performed at every marker across the  
16 genome to identify the segment with the most significant sharing evidence. All shared segments  
17 selected by the optimization procedure, and their respective p-values, comprise the final  
18 optimized SGS results.

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

1 represent the distribution (Fig. 1). Based on the fitted distribution,  $Y \sim \Gamma(k, \sigma)$ , where  $k$  and  $\sigma$  are  
2 the shape and rate parameters, we apply the Theory of Large Deviations; previously applied to  
3 successfully model genome-wide fluctuations in linkage analysis [27]. The significance  
4 threshold,  $T$ , accounts for multiple testing of optimized segments across the genome, and is  
5 found by solving Eq. 1:

$$6 \quad \mu(X) = [C + 2GX]\alpha(X) \quad (1)$$

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

12 In general, we found that the fitted distributions produced stable significance thresholds  
13 after 100,000-300,000 simulations (Table 1). Typically, threshold determination requires 1,000-  
14 3,000 CPU hours per pedigree, increasing with the number of subsets and separating meioses  
15 between pedigree cases. For example, in pedigree UT-571744, 300k simulations genome-wide  
16 (2,513,408 segments) took 1,275 CPU hours on tangent nodes featuring Intel Xeon E5-2650  
17 processors. Once significance thresholds are established, subset/segment combinations of  
18 potential interest are identified and additional simulations are restricted to those combinations to  
19 gain the required p-value resolution. For these subsequent targeted simulations, we use a  
20 marginalized LD map specific for the segment of interest, dramatically reducing the analysis  
21 time. For example, in pedigree UT-571744, 600M simulations on one segment took 325 CPU  
22 hours on tangent nodes featuring Intel Xeon E5-2650 processors. See S1 Fig. for an overview  
23 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

## 2 **Application to Utah, MM HRP**s

3 We applied our new pedigree analysis strategy to 11 Utah MM HRP

4 s using high-density OMNI Express SNP array genotype data. Each pedigree was selected to contain excess MM (4-

5 37 MM total per pedigree), had 2-4 sampled MM cases with genotype data, and 8-23 meioses

6 per pedigree between the sampled cases. After quality control, a consistent set of 678,447

7 SNPs were used for all SGS analyses. The total number of shared segments for each pedigree

8 across all subsets ranged from 638,525 to 6,765,500 (larger pedigrees with more subsets

9 producing larger numbers of segments). After optimization,  $Y = -\log_{10}(p)$  for 6,697 to 10,369

10 segments were fit to gamma distributions for each pedigree, and used to determine genome-

11 wide significant and suggestive thresholds (Eq. 1). The genome-wide significant thresholds

12 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

13 Table).

14 A genome-wide significant, 1.8 Mb shared segment ( $p = 3.3 \times 10^{-6}$ ) was observed in

15 pedigree UT-571744. All three genotyped MM cases, separated by 20 meioses, share the

16 segment (Fig. 2a and Table 2). The segment is located at chromosome 6q16 (98.49-100.24 Mb;

17 hg19) and includes 9 genes: *POU3F2*, *FBXL4*, *FAXC*, *COQ3*, *PNISR*, *USP45*, *TSTD3*, *CCNC*,

18 and *PRDM13* (Figure 2b).

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

Family	Cases	Me	Position	Len	p	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 <sup>^</sup>			<i>ARID1A</i>	p.Ser090Gly	MS	0.0002
Cornell MM12	2	4	1:27,089,712 <sup>`</sup>			<i>ARID1A</i>	p.Met890Val	MS	0.0001

**Legend:** Cases – total MM and MGUS cases (number of MGUS); Me – meioses; Position – build HG19, <sup>^</sup>rs752026201, <sup>`</sup>rs140664170; Len – length in mega-bases; p – SGS p-value, (significant and suggestive genome-wide thresholds were 3.8x10<sup>-6</sup> and 8.5x10<sup>-5</sup> for UT 571744, 3.5x10<sup>-6</sup> and 4.6x10<sup>-5</sup> for UT-576834, and 6.2x10<sup>-6</sup> and 1.2x10<sup>-4</sup> for UT 260), <sup>‡</sup>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 gnomAD individuals.

1  
2 We also identified two HRPs, UT-576834 and UT-260, with overlapping shared  
3 segments at 1p36.11 (Fig. 3). A 8.9 Mb (24.39-33.30 Mb, p = 3.0x10<sup>-4</sup>) segment was observed  
4 in 3 of the 4 genotyped MM cases in UT-576834, shared across 12 meioses (Fig. 3b and Table  
5 2). A nested 1.2 Mb shared segment (26.22-27.38 Mb; p = 2.1x10<sup>-4</sup>) segregated to 3 MM cases  
6 separated by 16 meioses in UT-260 (Fig. 3a and Table 2). The overlapping segment contains  
7 30 genes (Fig. 3d).

8  
9 **Exome follow-up of shared segments in HRPs**

10 Whole-exome sequencing (WES) data was interrogated, targeted to the identified SGS  
11 region, to identify potential risk variants in the pedigree sharers in the HRP and in a broader set  
12 of 44 pedigrees. WES data was available for: 28 cases from the 11 extended Utah HRPs; and



1 126 exomes from 44 densely clustered MM/MGUS families from Mayo Clinic Rochester, Weill  
2 Cornell, Memorial Sloan Kettering Cancer Center, International Agency for Research on  
3 Cancer, and INSERM France (S2 Table). Prioritization was used to identify variants that were:  
4 in the target segment; rare (alternate allele frequency, AAF<0.001 in the non-Finnish, European,  
5 gnomAD individuals), potentially deleterious (variant impact predicted to be high or moderate);  
6 and observed recurrently in the appropriate segment sharers (if observed in the segment  
7 discovery pedigree).

8 At 6q16, no rare, potentially deleterious coding risk variants were shared by the 3 UT-  
9 571744 MM cases in the 1.8 Mb genome-wide significant segment, indicating non-coding  
10 regulatory variants may be responsible for MM risk in this pedigree. However, two, rare coding  
11 and potentially deleterious single nucleotide variants (SNVs) were identified in two MM/MGUS  
12 families (Fig. 2c-e and Table 2). Both SNVs are in the hydrolase domain of *USP45*: a stop gain  
13 (p.Gln691\*) shared by 3 sibling cases (1 MM and 2 MGUS) in an INSERM family (PET-Nice  
14 0909) and a missense SNV (p.Gln621Glu) shared by 2 siblings (1 MM and 1 MGUS) but not  
15 their 2 screened unaffected siblings in Mayo family 485. Coverage of these positions in ExAC  
16 sequence data is high (> 99% of the 60,706 ExAC samples had at least 10x read coverage) and  
17 neither variant was observed. Collating the SGS evidence in UT 571744 (genome-wide rate of  
18  $\mu=0.0423$ ) with the sequence findings, correcting for 11 SGS pedigrees, the 45 pedigrees  
19 interrogated for sequence variants, and the 9 genes in the SGS region, we estimate the rate of  
20 observing all these findings at the 6q16 region by chance is low ( $\pi=0.01$ , see Methods) and  
21 study-wide significant.

22 Pedigree exomes in the 1.2 Mb segment at 1p36.11 revealed two, rare and potentially  
23 deleterious SNVs. The first in discovery pedigree UT-576834: a missense SNV (rs752026201,  
24 p.Ser90Gly, AAF = 0.0002 in gnomAD) in *ARID1A* (Fig. 3e) shared by 3 of the 4 Utah MM  
25 cases, concordant with the segment sharing pattern. A second rare, missense SNV in *ARID1A*

1 (rs140664170, p.Met890Val, AAF = 0.0001 in gnomAD) was found to be carried by a pair of MM  
2 cousins in Weill-Cornell family 12 (Fig. 3c and e, and Table 2). Based on the ExAC data,  
3 *ARID1A* is extremely intolerant to missense variants and loss of function (LoF) SNVs [28].  
4

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

6 Our SGS findings and pedigree WES identify *USP45* and *ARID1A* as candidate genes  
7 for inherited MM risk. We further investigated shared segments and WES for evidence  
8 supporting the complexes *USP45* and *ARID1A* are involved in. Here we further expanded our  
9 WES to: 186 MM/MGUS cases (early onset MM/MGUS or familial MGUS) from our collaborative  
10 group, 733 sporadic MM cases from dbGaP [29], and 964 controls [30].

11 *USP45* is an essential DNA repair regulator, de-ubiquitylating *ERCC1* to allow for DNA  
12 translocation of the *ERCC1-ERCC4* endonuclease [31,32]. This endonuclease is a part of the  
13 global genome nucleotide-excision repair (GG-NER) incision complex, a 22 protein complex  
14 essential to removing lesions from DNA and cancer prevention [33–36] (S3 Table). We  
15 reviewed SGS results in the Utah HRP at the location of these 22 genes and identified a  
16 genome-wide suggestive segment in pedigree UT-34955 (S2 Fig.). This HRP identified a 0.8 Mb  
17 segment at 19q13 (45.71-46.51 Mb; hg19), containing 31 genes including *ERCC1* and *ERCC2*  
18 (S2 Fig. and S4 Table). The segment is shared by 3 MM cases separated by 16 meioses ( $p =$   
19  $6.6 \times 10^{-5}$ ). No rare, coding variants were identified from the WES in the 3 MM cases in UT-  
20 34955, nor in the remaining 44 pedigrees/families. We interrogated the 23 GG-NER genes in  
21 our 919 MM/MGUS exomes. This identified a ClinVar-annotated pathogenic, missense SNV in  
22 *ERCC4* (p.Arg799Trp) in one early-onset MM case and one sporadic MM case, and a stop-gain  
23 SNV in *ERCC3* (p.Arg574Ter), in the same domain as a ClinVar-annotated pathogenic variant,  
24 in a second early-onset MM case (S4 Table). Further, burden testing in all MM cases vs controls

1 was significant in 2 of the 23 GG-NER genes: *GTF2H1* and *DDB1* after correcting for multiple  
2 testing (S3 Table). The occurrence of two significantly burdened genes (at  $\alpha=0.0022$ ) from 23  
3 genes is unexpected ( $p=0.0011$ , Binomial(23,0.0022)).

4 ARID1A is a member of the SWI/SNF chromatin remodeling complex, a 15 gene  
5 complex involved in DNA transcription regulation [37] (see S5 Table). Members of this complex  
6 are mutated in >20% of malignancies [38–40], but are extremely intolerant to LoF and missense  
7 variation [41] (S5 Table). We reviewed SGS results in the Utah HRP at the location of these 15  
8 genes and identified a marginal, genome-wide suggestive segment in pedigree UT-549917  
9 shared by 4 MM cases across 21 meioses ( $p = 2.17 \times 10^{-5}$ , S3 Fig. and S6 Table). This 1.5 Mb  
10 segment at chr3p21.1-p21.2 (52.01-53.56 Mb; hg19) contains 32 genes including *PBRM1* from  
11 the SWI/SNF complex. No coding variants were identified in this gene in UT-549917, nor in the  
12 remaining 44 pedigrees/families. Burden testing was significant for 7 of the 15 genes in the  
13 complex after correcting for multiple testing: *ARID1A*, *ARID1B*, *SMARCA4*, *ACTL6A*,  
14 *SMARCD3*, *SMARCC2*, and *SMARCE1* (S5 Table). The occurrence of seven significantly  
15 burdened genes (at  $\alpha=0.0033$ ) from 15 genes is unexpected by chance ( $p=2.7 \times 10^{-14}$ ,  
16 Binomial(15,0.0033)).

17

## 18 **DISCUSSION**

19 We developed a novel strategy to identify segregating chromosomal segments shared  
20 by subsets of cases in HRP. It focuses on extended HRP that are singularly powerful to  
21 identify significant genetic segregation. Our strategy allows for genetic heterogeneity within such  
22 pedigrees and provides formal significance thresholds for valid interpretation. Previously,  
23 extended HRP have not delivered on their potential in complex traits because in common,  
24 complex traits, HRP are likely enriched for multiple susceptibility variants and may capture both

1 familial and sporadic cases in their branches. Our optimization strategy over subsets is  
2 attractive because it allows for heterogeneity without prior knowledge of genetic similarities or  
3 deep phenotyping. This new statistic also identifies the sharers and clearly delimits the shared  
4 region, making follow-up interrogation straight-forward. This is a distinct advantage over  
5 standard linkage analysis and previous pairwise SGS methods where neither sharers or the  
6 region are defined [42].

7         Application of the method to extended MM pedigrees demonstrated the utility of this new  
8 method and illustrated that the segments identified were used successfully to narrow the search  
9 for risk variants in smaller pedigrees, allowing for an overall strategy that can utilize both large  
10 pedigrees and smaller families together for discovery (Table 2, Fig. 2 and Fig. 3). Post-hoc,  
11 additional value can be gained from demographic and/or clinical data on the sharing subsets  
12 shedding light on other shared characteristics that may aid future mapping. Also, we note that in  
13 the absence of any significant findings, genome-wide SGS results can be used as genomic  
14 annotations of segregation evidence for more heuristic approaches.

15         While we identified several rare, potentially deleterious coding variants of interest,  
16 several of the SGS discovery pedigrees had no coding variants that satisfied prioritization  
17 criteria. We believe this will be characteristic of complex traits and that regulatory variants will  
18 also play a substantial role. Mutations with strong causal likelihood found in other disease  
19 cohorts may focus the search for regulatory variation to particular genes within a shared  
20 segment, as with *USP45* in MM. In the absence of such compelling evidence, a return to  
21 pedigree segregation methods will provide identification of statistically compelling regions which  
22 can concentrate efforts to identify and characterize regulatory risk variants. Future work will  
23 include targeted sequencing of the promising MM SGS identified to investigate non-coding  
24 variants that may play a role in MM risk in these families. Our proposed method is a new

1 analytic tool with the potential to reinvigorate the use of extended HRPs in the identification of  
2 risk variants that contribute to common, complex disease.

3 Multiple myeloma is a malignancy of the plasma cells that has been shown to be familial  
4 [43]. Consistent with a role for genetics, case-control studies have been successful in identifying  
5 association signals for 17 low-risk variants [44–48]. However, despite consistent evidence for  
6 familial clustering, our study is the first to explore high-risk MM pedigrees. Using the unique  
7 genealogical database available in Utah, we identified and studied extended MM HRPs. We  
8 identified a genome-wide significant segment containing *USP45*, an important regulator of DNA  
9 repair (Fig. 2 and Table 2), and a genome-wide suggestive segment harboring other genes in  
10 the GG-NER incision complex (*ERCC1* and *ERCC2*). Exome sequencing in a collaborative  
11 resource of high-risk families and early-onset cases revealed four rare, potentially deleterious  
12 coding variants; two novel variants in *USP45* segregating in two pedigrees and two variants in  
13 early-onset cases in *ERCC3* and *ERCC4*, the latter annotated as pathogenic in ClinVar. Burden  
14 testing including sporadic MM, and comparing to controls, identified significant enrichment for  
15 variants in MM cases in 2 of the 23 GG-NER genes in the protein endonuclease regulation  
16 complex.

17 In particular, the functional literature supports *USP45* as a candidate cancer risk gene.  
18 *USP45* has been shown to deubiquitylate *ERCC1*, a catalytic subunit of the *ERCC1-ERCC4*  
19 DNA repair endonuclease (*ERCC4* also known as XPF) [31]. This endonuclease is a critical  
20 regulator of DNA repair processes [34]. The complex repairs recombination, double strand  
21 break, and inter-strand crosslink by cutting DNA overhangs around a lesion, degrades 3' G-rich  
22 overhangs in telomere maintenance, and plays a role in cancer prevention and in tumor  
23 resistance to chemotherapy [31,34]. Mouse models have shown *USP45* knockout cells have  
24 higher levels of ubiquitylated *ERCC1* and that cells are hypersensitive to UV radiation and DNA  
25 inter-strand cross-links, repair of UV-induced DNA damage, and *ERCC1* translocation to DNA

1 damage is impaired [31]. Hence, the deubiquitylase activity of USP45 is important for  
2 maintaining the DNA repair ability of ERCC1-ERCC4. In total, these observations implicate the  
3 GG-NER incision complex and specifically the interaction of USP45 and the disruption of the  
4 ERCC1-ERCC4 role in DNA repair as a mechanism of potential importance in MM risk.

5 Our strategy also identified shared segments overlapping at chr1p36.11 in two Utah  
6 pedigrees containing *ARID1A* (Fig. 3 and Table 2) and a genome-wide suggestive segment in a  
7 third pedigree harboring another gene in the SWI/SNF complex (*PBRM1*). For the SWI/SNF  
8 complex, exome sequencing revealed two rare, potentially deleterious variants in *ARID1A*  
9 segregating in two pedigrees. Burden testing provided further evidence for enrichment of  
10 variants in *ARID1A* specifically, and in 7 of the 15 genes in the complex. As a component of the  
11 SWI/SNF chromatin remodeling complex, ARID1A facilitates gene activation by assisting  
12 transcription machinery gain access to gene targets [49]. Based on the patterns of mutations in  
13 tumor cells, *ARID1A* likely functions as a tumor-suppressor [50]. Members of the SWI/SNF  
14 chromatin remodeling complexes are mutated in 20% of malignancies [38], but are extremely  
15 intolerant to LoF and missense variation [41] (S5 Table). Blockage of chromatin remodeling may  
16 sustain cancer development [39]. Aberrant chromatin remodeling contributes to the  
17 pathogenesis of ovarian clear-cell carcinoma [50]. It has previously been shown that *ARID1A* is  
18 intolerant to variation (LoF and missense mutations) [28], consistent with its prominent somatic  
19 role in multiple tumors [38,50,51], including hematological malignancies [52–54]. These  
20 observations implicate the SWI/SNF chromatin remodeling complex, and specifically *ARID1A* in  
21 MM risk.

22 This study has limitations. First, the method is applicable only to extended HRPs that are  
23 singularly effective for identifying segregating segments (15 meioses between cases is optimal  
24 [16]). The method is not directly applicable to the many smaller family-based resources that  
25 have been gathered in the complex trait field and may therefore result in findings from single

1 large pedigrees that are private and difficult to replicate. However, as illustrated in our example,  
2 in a collaborative setting containing both extended HRPs and smaller families, the approach can  
3 be mutually beneficial. Second, our observation of two borderline genome-wide suggestive  
4 overlapping segments at 1p36 led to our identification of *ARID1A* as a potential candidate risk  
5 gene and illustrates the potential for discoveries using overlapping subthreshold evidence.  
6 However, it raises analytical questions of how to systematically identify such segments. This  
7 segment would have been ignored based on strict individual-pedigree thresholds and highlights  
8 an important area for further methodological development. Third, as in all family-based genetic  
9 studies our method is susceptible to inaccurate pedigree structures and poorly matched control  
10 populations. However, relationship and ethnicity checks are standard protocol and mitigate the  
11 possibility of error. Finally, this study is observational and cannot describe causation. We have  
12 identified two complexes, several genes and specific variants as compelling candidates involved  
13 in MM risk, but further functional studies will be required to determine and characterize the  
14 mechanisms involved in risk.

15 In conclusion, we have developed a strategy for gene mapping in complex traits that  
16 accounts for heterogeneity within HRPs and formally corrects for multiple testing to allow for  
17 statistically rigorous discovery. We applied this strategy to MM, a complex cancer of plasma  
18 cells, and identified multiple shared segments containing genes in nucleotide excision repair  
19 and SWI/SNF chromatin remodeling. Exome follow-up supported these segments in both the  
20 Utah large HRPs and smaller families from other sites. Our study offers a novel technique for  
21 HRP gene mapping and demonstrates its utility to narrow the search for risk-variants in complex  
22 traits.

## 1    **METHODS**

### 2    **SGS Analysis in Utah, Myeloma HRP**

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

12  
13    **Quality control.** Only bi-allelic SNPs were considered. Genotypes and individual call-rates  
14    were used to ensure high quality data. PLINK was used to remove SNPs with  $< 95\%$  call rate  
15    across individuals [55]. The final SNP set contained 678,447 single nucleotide variants. After  
16    SNP removal for low call rates, individuals were removed based on  $< 90\%$  call rate across the  
17    genome, or if they failed the PLINK sex check. One MM case was removed. The QC'ed SNP  
18    data were transformed to match strand orientation of the 1000Genomes. PLINK relationship  
19    estimates were assessed against pedigree structure from the UPDB to identify any potential  
20    issues with pedigree structure. None were found.

21  
22    **Probability of sharing a segment.** SGS analysis identifies contiguous SNPs that are  
23    shared identical-by-state (IBS) by cases in a HRP and assigns an empirical probability of  
24    chance ancestral sharing [26]. First, a set of cases in a HRP are defined and all segments of



1 contiguous SNPs shared IBS are identified. All shared segments > 20 SNPs are considered.  
2 Lengths shorter than 20 are commonly shared between unrelated individuals. Second,  
3 population-based data (here we used CEU and GBR data from the 1000Genomes Project [56])  
4 are used to estimate a graphical model for linkage disequilibrium (LD) [57], providing a  
5 probability distribution of chromosome-wide haplotypes in the population. Third, pairs of  
6 haplotypes are randomly assigned to pedigree founders according to the haplotype distribution.  
7 Founders are individuals whose parents are not specified in the pedigree. For chromosome-  
8 wide haplotype simulations the full chromosome LD model is used. Fourth, Mendelian  
9 segregation and recombination are simulated to generate genotypes for all pedigree members.  
10 The Rutgers genetic map [58] is used for a genetic map for recombination, with interpolation  
11 based on physical base pair position for SNPs not represented. Steps two through four create  
12 one simulated data set, a random sample from the null hypothesis. This process is repeated  
13 hundreds of thousands to millions of times for each subset.

14 Each shared segment in the real data (step one) is compared to the simulated segments  
15 at the precise genomic location. The number of times the null segment equals or encompasses  
16 the observed segment is counted and divided by the total number of simulations to generate the  
17 empirical nominal p-value for the observed shared segment. The simulations continue until a p-  
18 value has been estimated to a required resolution, or until it surpasses a defined significance  
19 threshold. To facilitate this in an efficient manner, we follow-up specific segments using  
20 marginal distributions from the LD model, established using standard graphical modeling  
21 methods [59]. The marginalized LD model encompassing only the region of interest, but  
22 capturing relevant LD to accurately simulate genotypes from this region alone. This reduction in  
23 markers vastly increases the speed in which simulations are generated. The graphical model  
24 estimation, marginalization, and simulation processes are computationally efficient requiring  
25 time and storage that is linear with the number of SNPs being considered.

1 **Heterogeneity optimization.** We systematically perform SGS analysis on each subset of  
2 cases in a HRP. If required, the number of subsets can be limited by meioses or subset size.  
3 This may be necessary for common traits with large full sets. A lower limit of 10 meioses is a  
4 good rule of thumb for reducing the computational burden of subset assessment. At each  
5 marker position across the genome, the optimized segment is the one minimizing the p-value  
6 across all subsets considered. All segments selected by the optimization procedure, and their  
7 respective p-values, comprise the final optimized SGS results.

8  
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,  $\sigma$  rate 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:

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

15  
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:

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

20  
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,  
23 our estimated  $k$  and  $\sigma$  gamma parameters, where  $T = 10^{-X\sigma/2}$ ,  $X = 2Y/\sigma \sim \chi_{2k}^2$ , and the genetic  
24 length of the genome (matched to that used in the gene-drop) to utilize Eq. 3 and derive  $\mu(X)$

1 thresholds. Solving for  $\mu(X) = 0.05$  and  $\mu(X) = 1$  produced significance and suggestive  
2 thresholds, respectively. These thresholds are remarkably stable after a few hundred thousand  
3 simulations. For pedigrees with very large numbers of meioses (>50) between the full case-set  
4 a larger number of simulations may be required.

5

6 **Software availability.** The SGS program is available for download at  
7 <https://gitlab.com/camplab/sgs> and <https://gitlab.com/camplab/jps>. The main architecture is  
8 written in Java. Probability assessments can be multi-threaded, but the largest parallelization  
9 gains are achieved by running independent analyses across chromosomes.

10

## 11 **Targeted sequencing**

12 **Participants.** WES data were interrogated in the regions defined by the shared segments of  
13 interest. WES data was available on 964 controls [30] and 1,063 MM or MGUS cases including:  
14 28 MM from the 11 Utah HRP; 70 MM and 46 MGUS from 44 densely clustered families (each  
15 containing at least 2 MM or at least 1 MM and 1 MGUS); 186 genetically-enriched MM/MGUS  
16 (148 MM and 38 MGUS) including early-onset and MGUS clustering in families; and 733  
17 sporadic MM cases from dbGaP [29]. Of the 44 densely MM/MGUS high-risk families, 25 were  
18 ascertained by INSERM, France (36 MM, 38 MGUS), 9 by Mayo Clinic, Minnesota (10 MM, 8  
19 MGUS, 10 unaffected family members), 6 by Memorial Sloan Kettering Cancer Center, New  
20 York (14 MM), 3 by International Agency for Research on Cancer, France (8 MM), and 1 by  
21 Weill Cornell, New York (2 MM). Most of the families had both MM and MGUS cases (32  
22 families total) and 12 families only had MM cases sequenced. Six families had at least one  
23 unaffected relative sequenced. (See S2 Table.) All individuals in the Utah HRP and the all but  
24 three of the densely clustered families were of non-Finish European descent.

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 [60]. 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 (GEnome MINIng) database for ease of querying [61]. Additional functional annotations  
9 available in the GEMINI suite include CADD, ANNOVAR, conservation, location, and if the  
10 variant was listed in OMIM.

11  
12 **Variant prioritization.** A GEMINI query was developed to identify variants which were: high  
13 or medium impact; AAF < 0.001 in the non-Finnish, European, gnomAD individuals; and within  
14 the shared segments of interest. Genes harboring segregating variants in at least two high-risk  
15 pedigrees (the discovery pedigree and/or the 44 high-risk pedigrees from collaborating sites)  
16 were considered candidate susceptibility genes. These criteria were selected to maintain  
17 findings that were unlikely by chance after accounting for both the SGS and sequencing stages  
18 of the study. From ExAC exomes, the number of medium/high impact variants with AAF<0.001  
19 per person per gene is 0.0016 [28]. The probability of identifying segregating variants in at least  
20 two pedigrees in the same gene can be approximated with a probability from a Binomial(45,  
21 0.0016), which equals 0.0024. To account for the multiple genes in the SGS region, a second  
22 probability from Binomial(G, 0.0024) can be used to estimate the probability of observing two  
23 segregating variants by chance in G genes. With a threshold of AAF<0.001, the probability of

1 observing at least one gene that harbors 2 variants that segregate in high-risk pedigrees ( $\phi$ )  
2 remains unexpected by chance ( $<0.05$ ) for up to a reasonable large number of genes ( $G=20$ ).

3

4 **Joint Assessment of SGS and Sequencing.** We assessed the overall rate of expectation for  
5 the joint experiments of the SGS and pedigree sequencing findings as  $\pi = 11 \times \mu \times \phi$ , where  $\mu$   
6 is the fully corrected genome-wide rate for the SGS region identified, and  $\phi$  is the fully corrected  
7 probability of the sequencing findings based on the number of genes in the SGS region, as  
8 described above.

9

10 **Burden testing.** Burden testing was performed on jointly called and processed WES from  
11 1,063 MM/MGUS cases and 964 unaffected controls for the 23 genes in the GG-NER incision  
12 complex (including *USP45*) and 15 genes in the SWI/SNF chromatin remodeling complex. The  
13 GEMINI software [61] was used to perform a c-alpha test [62] with 1000 permutations. Only  
14 variants with AAF  $< 0.05$  and high or moderate predicted impact were included in the analysis.

15

## 16 **ACKNOWLEDGMENTS**

17 We thank the DNA Sequencing Core Facility and Genomics Core Facility at the University of  
18 Utah, and the computational resources and staff expertise provided by Scientific Computing at  
19 the Icahn School of Medicine at Mount Sinai. Data collection was made possible, in part, by the  
20 Utah Population Database and the Utah Cancer Registry. We thank the participants and their  
21 families who make this research possible.

## 1 REFERENCES

- 2 1. Maher B. Personal genomes: The case of the missing heritability. *Nature*. 2008;456: 18–  
3 21. doi:10.1038/456018a
- 4 2. Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorf LA, Hunter DJ, et al. Finding the  
5 missing heritability of complex diseases. *Nature*. 2009;461: 747–53.  
6 doi:10.1038/nature08494
- 7 3. Eichler EE, Flint J, Gibson G, Kong A, Leal SM, Moore JH, et al. Missing heritability and  
8 strategies for finding the underlying causes of complex disease. *Nat Rev Genet*.  
9 2010;11: 446–50. doi:10.1038/nrg2809
- 10 4. Miki Y, Swensen J, Shattuck-Eidens D, Futreal PA, Harshman K, Tavtigian S, et al. A  
11 strong candidate for the breast and ovarian cancer susceptibility gene BRCA1. *Science*.  
12 1994;266: 66–71. Available: <http://www.ncbi.nlm.nih.gov/pubmed/7545954>
- 13 5. Wooster R, Neuhausen SL, Mangion J, Quirk Y, Ford D, Collins N, et al. Localization of  
14 a breast cancer susceptibility gene, BRCA2, to chromosome 13q12-13. *Science*.  
15 1994;265: 2088–90. Available: <http://www.ncbi.nlm.nih.gov/pubmed/8091231>
- 16 6. Vance JM, Pericak-Vance MA, Yamaoka LH, Speer MC, Rosenwasser GO, Small K, et  
17 al. Genetic linkage mapping of chromosome 17 markers and neurofibromatosis type I.  
18 *Am J Hum Genet*. 1989;44: 25–9. Available:  
19 <http://www.ncbi.nlm.nih.gov/pubmed/2491777>
- 20 7. Cannon-Albright LA, Goldgar DE, Meyer LJ, Lewis CM, Anderson DE, Fountain JW, et  
21 al. Assignment of a locus for familial melanoma, MLM, to chromosome 9p13-p22.  
22 *Science*. 1992;258: 1148–52. Available: <http://www.ncbi.nlm.nih.gov/pubmed/1439824>

- 1 8. Leppert M, Dobbs M, Scambler P, O'Connell P, Nakamura Y, Stauffer D, et al. The gene  
2 for familial polyposis coli maps to the long arm of chromosome 5. *Science*. 1987;238:  
3 1411–3. Available: <http://www.ncbi.nlm.nih.gov/pubmed/3479843>
- 4 9. Nishisho I, Nakamura Y, Miyoshi Y, Miki Y, Ando H, Horii A, et al. Mutations of  
5 chromosome 5q21 genes in FAP and colorectal cancer patients. *Science*. 1991;253:  
6 665–9. Available: <http://www.ncbi.nlm.nih.gov/pubmed/1651563>
- 7 10. Ng SB, Buckingham KJ, Lee C, Bigham AW, Tabor HK, Dent KM, et al. Exome  
8 sequencing identifies the cause of a mendelian disorder. *Nat Genet*. 2010;42: 30–5.  
9 doi:10.1038/ng.499
- 10 11. Ng SB, Bigham AW, Buckingham KJ, Hannibal MC, McMillin MJ, Gildersleeve HI, et al.  
11 Exome sequencing identifies MLL2 mutations as a cause of Kabuki syndrome. *Nat*  
12 *Genet*. 2010;42: 790–3. doi:10.1038/ng.646
- 13 12. McClellan J, King M-C. Genetic Heterogeneity in Human Disease. *Cell*. 2010;141: 210–  
14 217. doi:10.1016/j.cell.2010.03.032
- 15 13. Mitchell KJ. What is complex about complex disorders? *Genome Biol*. 2012;13: 237.  
16 doi:10.1186/gb-2012-13-1-237
- 17 14. Li X, Montgomery SB. Detection and Impact of Rare Regulatory Variants in Human  
18 Disease. *Front Genet*. 2013;4. doi:10.3389/fgene.2013.00067
- 19 15. Albert FW, Kruglyak L. The role of regulatory variation in complex traits and disease. *Nat*  
20 *Rev Genet*. 2015;16: 197–212. doi:10.1038/nrg3891
- 21 16. Knight S, Abo RP, Abel HJ, Neklason DW, Tuohy TM, Burt RW, et al. Shared Genomic  
22 Segment Analysis: The Power to Find Rare Disease Variants. *Ann Hum Genet*. 2012;76:  
23 500–509. doi:10.1111/j.1469-1809.2012.00728.x
- 24 17. Myeloma - SEER Stat Fact Sheets [Internet]. Available:  
25 <https://seer.cancer.gov/statfacts/html/mulmy.html>

- 1 18. Cannon-Albright LA, Thomas A, Goldgar DE. Familiality of cancer in Utah. *Cancer Res.*  
2 1994;54: 2378–2385.
- 3 19. Landgren O, Linet MS, McMaster ML, Gridley G, Hemminki K, Goldin LR. Familial  
4 characteristics of autoimmune and hematologic disorders in 8,406 multiple myeloma  
5 patients: A population-based case-control study. *Int J Cancer.* 2006;118: 3095–3098.  
6 doi:10.1002/ijc.21745
- 7 20. Albright F, Teerlink C, Werner TL, Cannon-Albright LA. Significant evidence for a  
8 heritable contribution to cancer predisposition: a review of cancer familiality by site. *BMC*  
9 *Cancer.* BioMed Central Ltd; 2012;12: 138. doi:10.1186/1471-2407-12-138
- 10 21. Schinasi LH, Brown EE, Camp NJ, Wang SS, Hofmann JN, Chiu BC, et al. Multiple  
11 myeloma and family history of lymphohaematopoietic cancers: Results from the  
12 International Multiple Myeloma Consortium. *Br J Haematol.* England; 2016;175: 87–101.  
13 doi:10.1111/bjh.14199
- 14 22. Landgren O, Kristinsson SY, Goldin LR, Caporaso NE, Blimark C, Mellqvist U-H, et al.  
15 Risk of plasma cell and lymphoproliferative disorders among 14621 first-degree relatives  
16 of 4458 patients with monoclonal gammopathy of undetermined significance in Sweden.  
17 *Blood.* 2009;114: 791–5. doi:10.1182/blood-2008-12-191676
- 18 23. Greenberg AJ, Rajkumar SV, Vachon CM. Familial monoclonal gammopathy of  
19 undetermined significance and multiple myeloma: epidemiology, risk factors, and  
20 biological characteristics. *Blood.* 2012;119: 5359–66. doi:10.1182/blood-2011-11-  
21 387324
- 22 24. Greenberg AJ, Rajkumar SV, Larson DR, Dispenzieri A, Therneau TM, Colby CL, et al.  
23 Increased prevalence of light chain monoclonal gammopathy of undetermined  
24 significance (LC-MGUS) in first-degree relatives of individuals with multiple myeloma. *Br*  
25 *J Haematol.* 2012;157: 472–5. Available: <http://www.ncbi.nlm.nih.gov/pubmed/22629552>



- 1 25. Vachon CM, Kyle RA, Therneau TM, Foreman BJ, Larson DR, Colby CL, et al.  
2 Increased risk of monoclonal gammopathy in first-degree relatives of patients with  
3 multiple myeloma or monoclonal gammopathy of undetermined significance. *Blood*.  
4 2009;114: 785–90. doi:10.1182/blood-2008-12-192575
- 5 26. Thomas A, Camp NJ, Farnham JM, Allen-Brady K, Cannon-Albright LA. Shared  
6 Genomic Segment Analysis. Mapping Disease Predisposition Genes in Extended  
7 Pedigrees Using SNP Genotype Assays. *Ann Hum Genet*. 2008;72: 279–287.  
8 doi:10.1111/j.1469-1809.2007.00406.x
- 9 27. Lander E, Kruglyak L. Genetic dissection of complex traits: guidelines for interpreting  
10 and reporting linkage results. *Nat Genet*. 1995;11: 141–147.
- 11 28. Lek M, Karczewski KJ, Minikel E V., Samocha KE, Banks E, Fennell T, et al. Analysis of  
12 protein-coding genetic variation in 60,706 humans. *Nature*. *Nature Research*; 2016;536:  
13 285–291. doi:10.1038/nature19057
- 14 29. Myeloma data downloaded from the dbGaP web site under accessions:  
15 phs000348.v2.p1 and phs000748.v4.p3. [Internet].
- 16 30. Control data downloaded from the dbGaP web site under accessions:  
17 phs000209.v13.p3, phs000276.v2.p1, phs000179.v5.p2, phs000298.v3.p2,  
18 phs000424.v6.p1, phs000653.v2.p1, phs000687.v1.p1, phs000814.v1.p1, and  
19 phs000806.v1.p1.
- 20 31. Perez-Oliva AB, Lachaud C, Szyniarowski P, Muñoz I, Macartney T, Hickson I, et al.  
21 USP45 deubiquitylase controls ERCC1-XPF endonuclease-mediated DNA damage  
22 responses. *EMBO J*. 2015;34: 326–43. doi:10.15252/emj.201489184
- 23 32. USP45 in the GG-NER Incision Complex [Internet]. Available:  
24 <http://www.reactome.org/PathwayBrowser/#/R-HSA-5696398&SEL=R-HSA-5696465&PATH=R-HSA-73894>  
25

- 1 33. Marteijn JA, Lans H, Vermeulen W, Hoeijmakers JHJ. Understanding nucleotide excision  
2 repair and its roles in cancer and ageing. *Nat Rev Mol Cell Biol.* 2014;15: 465–81.  
3 doi:10.1038/nrm3822
- 4 34. Kirschner K, Melton DW. Multiple roles of the ERCC1-XPF endonuclease in DNA repair  
5 and resistance to anticancer drugs. *Anticancer Res.* 2010;30: 3223–3232. doi:30/9/3223  
6 [pii]
- 7 35. Friedberg EC. How nucleotide excision repair protects against cancer. *Nat Rev Cancer.*  
8 2001;1: 22–33. doi:10.1038/35094000
- 9 36. Christmann M, Tomicic MT, Roos WP, Kaina B. Mechanisms of human DNA repair: an  
10 update. *Toxicology.* 2003;193: 3–34. Available:  
11 <http://www.ncbi.nlm.nih.gov/pubmed/14599765>
- 12 37. SWI/SNF Chromatin Remodeling Complex [Internet]. Available:  
13 <http://www.reactome.org/PathwayBrowser/#/R-HSA-5696398&PATH=R-HSA-73894>
- 14 38. Biegel JA, Busse TM, Weissman BE. SWI/SNF chromatin remodeling complexes and  
15 cancer. *Am J Med Genet C Semin Med Genet.* 2014;166C: 350–66.  
16 doi:10.1002/ajmg.c.31410
- 17 39. Romero O a, Sanchez-Cespedes M. The SWI/SNF genetic blockade: effects in cell  
18 differentiation, cancer and developmental diseases. *Oncogene.* Nature Publishing  
19 Group; 2014;33: 2681–9. doi:10.1038/onc.2013.227
- 20 40. Roberts CWM, Orkin SH. The SWI/SNF complex - chromatin and cancer. *Nat Rev*  
21 *Cancer.* Nature Publishing Group; 2004;4: 133–142. Available:  
22 <http://dx.doi.org/10.1038/nrc1273>
- 23 41. Lek M. Analysis of protein-coding genetic variation in 60,706 humans. 2015;  
24 doi:10.1101/030338

- 1 42. Cai Z, Camp NJ, Cannon-Albright L, Thomas A. Identification of regions of positive  
2 selection using Shared Genomic Segment analysis. *Eur J Hum Genet. Nature*  
3 *Publishing Group*; 2011;19: 667–671. doi:10.1038/ejhg.2010.257
- 4 43. Morgan GJ, Johnson DC, Weinhold N, Goldschmidt H, Landgren O, Lynch HT, et al.  
5 Inherited genetic susceptibility to multiple myeloma. *Leukemia*. 2014;28: 518–24.  
6 doi:10.1038/leu.2013.344
- 7 44. Broderick P, Chubb D, Johnson DC, Weinhold N, Försti A, Lloyd A, et al. Common  
8 variation at 3p22.1 and 7p15.3 influences multiple myeloma risk. *Nat Genet. Nature*  
9 *Publishing Group*; 2011;44: 58–61. doi:10.1038/ng.993.Common
- 10 45. Chubb D, Weinhold N, Broderick P, Chen B, Johnson DC, Försti A, et al. Common  
11 variation at 3q26.2, 6p21.33, 17p11.2 and 22q13.1 influences multiple myeloma risk. *Nat*  
12 *Genet. Nature Publishing Group*; 2013;45: 1221–1225. doi:10.1038/ng.2733
- 13 46. Weinhold N, Johnson DC, Chubb D, Chen B, Försti A, Hosking FJ, et al. The CCND1  
14 c.870G>A polymorphism is a risk factor for t(11;14)(q13;q32) multiple myeloma. *Nat*  
15 *Genet*. 2013;45: 522–5. doi:10.1038/ng.2583
- 16 47. Swaminathan B, Thorleifsson G, Jöud M, Ali M, Johnsson E, Ajore R, et al. Variants in  
17 ELL2 influencing immunoglobulin levels associate with multiple myeloma. *Nat Commun*.  
18 2015;6: 7213. doi:10.1038/ncomms8213
- 19 48. Mitchell JS, Li N, Weinhold N, Försti A, Ali M, Duin M Van, et al. Genome-wide  
20 association study identifies multiple susceptibility loci for multiple myeloma. *Nat*  
21 *Commun*. 2016;7: 12050. doi:10.1038/ncomms12050
- 22 49. Nie Z, Xue Y, Yang D, Zhou S, Deroo BJ, Archer TK, et al. A specificity and targeting  
23 subunit of a human SWI/SNF family-related chromatin-remodeling complex. *Mol Cell*  
24 *Biol*. 2000;20: 8879–88. Available: <http://www.ncbi.nlm.nih.gov/pubmed/11073988>

- 1 50. Jones S, Wang T-L, Shih I-M, Mao T-L, Nakayama K, Roden R, et al. Frequent  
2 mutations of chromatin remodeling gene ARID1A in ovarian clear cell carcinoma.  
3 Science. 2010;330: 228–31. doi:10.1126/science.1196333
- 4 51. Hodges C, Kirkland JG, Crabtree GR. The Many Roles of BAF (mSWI/SNF) and PBAF  
5 Complexes in Cancer. Cold Spring Harb Perspect Med. 2016;6.  
6 doi:10.1101/cshperspect.a026930
- 7 52. Puente XS, Beà S, Valdés-Mas R, Villamor N, Gutiérrez-Abril J, Martín-Subero JI, et al.  
8 Non-coding recurrent mutations in chronic lymphocytic leukaemia. Nature. 2015;526:  
9 519–524. doi:10.1038/nature14666
- 10 53. Lunning MA, Green MR. Mutation of chromatin modifiers; an emerging hallmark of  
11 germinal center B-cell lymphomas. Blood Cancer J. 2015;5: e361.  
12 doi:10.1038/bcj.2015.89
- 13 54. Choi J, Goh G, Walradt T, Hong BS, Bunick CG, Chen K, et al. Genomic landscape of  
14 cutaneous T cell lymphoma. Nat Genet. 2015;47: 1–11. doi:10.1038/ng.3356
- 15 55. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: A  
16 Tool Set for Whole-Genome Association and Population-Based Linkage Analyses. Am J  
17 Hum Genet. 2007;81: 559–575. doi:10.1086/519795
- 18 56. 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang  
19 HM, et al. A global reference for human genetic variation. Nature. 2015;526: 68–74.  
20 doi:10.1038/nature15393
- 21 57. Abel HJ, Thomas A. Accuracy and Computational Efficiency of a Graphical Modeling  
22 Approach to Linkage Disequilibrium Estimation. Stat Appl Genet Mol Biol. 2011;10.  
23 doi:10.2202/1544-6115.1615

- 1 58. Matise TC, Chen F, Chen W, De La Vega FM, Hansen M, He C, et al. A second-  
2 generation combined linkage physical map of the human genome. *Genome Res.*  
3 2007;17: 1783–6. doi:10.1101/gr.7156307
- 4 59. Lauritzen SL. *Graphical models*. Clarendon Press; 1996.
- 5 60. Linderman MD, Brandt T, Edelmann L, Jabado O, Kasai Y, Kornreich R, et al. Analytical  
6 validation of whole exome and whole genome sequencing for clinical applications. *BMC*  
7 *Med Genomics*. 2014;7: 20. doi:10.1186/1755-8794-7-20
- 8 61. Paila U, Chapman BA, Kirchner R, Quinlan AR. GEMINI: Integrative Exploration of  
9 Genetic Variation and Genome Annotations. Gardner PP, editor. *PLoS Comput Biol.*  
10 2013;9: e1003153. doi:10.1371/journal.pcbi.1003153
- 11 62. Neale BM, Rivas MA, Voight BF, Altshuler D, Devlin B, Orho-Melander M, et al. Testing  
12 for an unusual distribution of rare variants. *PLoS Genet*. 2011;7.  
13 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  $k = 4.4$  and  $\sigma = 3.6$  with good  
4 visual density (a) and CDF (b) fit, with  $\lambda = 0.9$ . (Goodness of fit was estimated with  $\lambda$ , the  
5 median of empirical chi-squared distribution divided by the median of the expected chi-squared  
6 distribution.) HRP 34955 has  $k = 2.8$  and  $\sigma = 2.9$  with good visual density (c) and CDF (d) fit,  
7 with  $\lambda = 1.0$ .

8  
9 **Fig. 2. Significant SGS, pedigrees, and segregating SNVs.** In pedigrees, MM cases are fully  
10 shaded and MGUS cases are half shaded. Numbers indicate multiple individuals. a) Utah  
11 pedigree, 571744, sharing the genome-wide significant SGS. The pedigree is trimmed to allow  
12 for viewing (37 MM confirmed cases are known in this pedigree, 3 were ascertained and  
13 genotyped). + indicates the genotyped MM cases that are SGS carriers, - indicates genotyped  
14 and non-carriers, no carrier status indicates not genotyped. Note – the genealogy extends  
15 beyond SEER cancer registry data. MGUS are unknown in this pedigree. b) Genomic region of  
16 significant SGS. c) INSERM pedigree carrying the stop gain SNV marked by “c” in box e. 1 MM  
17 and 2 MGUSs carry the SNV. d) Mayo Clinic pedigree carrying the missense SNV marked by  
18 “d” in box e. 1 MM and 1 MGUS carry the SNV, but not 2 unaffected siblings. e) Risk candidate  
19 gene, *USP45*, has 2 segregating SNVs in the ubiquitin C-terminal hydrolase 2 (UCH) domain.

20  
21 **Fig. 3. SGS with multiple lines of evidence.** a/b) Utah pedigrees carrying the overlapping  
22 SGSs on chr1p36.11-p35.1. + indicates the genotyped MM cases that are SGS carriers, -  
23 indicates genotyped and non-carriers, no carrier status indicates not genotyped. c) Weill Cornell  
24 pedigree with a segregating, missense SNV in *ARID1A* indicated by “c” in e. d) Genomic region

1 of overlapping SGS. Dark black genes fall in both regions. e) 2 rare and segregating, missense  
2 SNVs were observed in whole-exome sequencing. SNV “b” is carried by the cases indicated  
3 with + in box b. SNV “c” in carried by the cases in box c.

4

## 5 **SUPPORTING INFORMATION**

6 **S1 Fig. SGS analysis workflow.** Overview of the strategy pipeline. Genotypes can be  
7 generated from a high-density SNP array, or by extracting SNVs from whole-genome  
8 sequencing. CEU and GBR genotypes (unrelated individuals only) from the 1000Genomes  
9 Project are generally used as population controls. Dotted boxes represent steps done per-  
10 pedigree. Dash-dot boxes represent steps done on all subsets of cases within a pedigree.  
11 Dashed box contains step repeated for each simulation. Abbreviations: SNP – single nucleotide  
12 polymorphism; SGS – shared genomic segment; LD – linkage disequilibrium; PED – pedigree  
13 file (contains relationships and genotypes).

14

15 **S2 Fig. Genome-wide suggestive segment contains *ERCC1*.** a) Utah pedigree carrying  
16 the genome-wide suggestive SGS at chr19q13.32. + indicates the genotyped MM cases that  
17 are SGS carriers, - indicates genotyped and non-carriers, no carrier status indicates not  
18 genotyped. b) Genomic region captured by the SGS. *ERCC1* and *ERCC2* are contained.

19

20 **S3 Fig. Shared segment containing *PBRM1*.** a) Pedigree Utah 549917 carries a genome-  
21 wide suggestive SGS at chr3p21.2-p21.1. + indicates the genotyped MM cases that are SGS  
22 carriers, - indicates genotyped and non-carriers, no carrier status indicates not genotyped. b)  
23 Genome region captured by the SGS including *PBRM1*, a component of the SWI/SNF  
24 chromatin remodeling complex.

1 **S1 Table. Genome-wide thresholds and segments.**

2

3 **S2 Table. Whole-exome sequenced families.** Total MM, MGUS, and controls in each  
4 pedigree and from each site.

5

6 **S3 Table. GG-NER Incision Complex genes.** Burden testing results (based on 1063  
7 MM/MGUS cases and 964 unaffected controls), SGS and prioritized SNV results, and tolerance  
8 to missense and loss of function variants (based on ExAC population data).

9

10 **S4 Table. Evidence for endonuclease regulation of DNA repair.**

11

12 **S5 Table. SWI/SNF Complex genes.** Burden testing results (based on 1063 MM/MGUS  
13 cases and 964 unaffected controls), SGS and prioritized SNV results, and tolerance to  
14 missense and loss of function variants (based on ExAC population data).

15

16 **S6 Table. Evidence for SWI/SNF chromatin remodeling.**













