- 1 Supplementary Materials and Methods
- 3 **1. Sample Acquisition and Pathology Review**
- 4 5

1.1 METHODS

6 **1.1.1 Sample Acquisition**

7 MPM tumors were collected and shipped to the Biospecimen Core Resource 8 (BCR) between September, 2012 and December, 2013. Qualifying tumor samples were 9 obtained from patients who had received no prior treatment for their disease 10 (chemotherapy or radiotherapy). Specimens were shipped overnight from 13 tissue 11 source sites (TSS) using a cryoport that maintained an average temperature of less than 12 -180°C. In addition to tumor samples, each frozen primary tumor specimen had a 13 companion normal tissue specimen (blood or blood components, including DNA 14 extracted at the tissue source site). Adjacent non-tumor tissue was also submitted for a 15 subset of cases.

16

17 Cases were staged according to the American Joint Committee on Cancer 18 (AJCC) staging system. Pathology quality control was performed on each tumor and 19 adjacent normal tissue (if available) specimen from a frozen section slide prepared either 20 by the BCR or by the TSS. Hematoxylin and eosin (H&E) stained sections from each 21 sample were subjected to independent pathology review to confirm that the tumor 22 specimen was histologically consistent with mesothelioma and the adjacent tissue 23 specimen contained no tumor cells. The percent tumor nuclei, percent necrosis, and 24 other pathology annotations were also assessed. Tumor samples with >60% tumor 25 nuclei and ≤20% necrosis were submitted for nucleic acid extraction.

26

27 **1.1.2 Sample Processing**

DNA and RNA were extracted and quality was assessed at the BCR. RNA and DNA were extracted from tumor and adjacent non-tumor tissue specimens using a modification of the DNA/RNA AllPrep kit (Qiagen). The flow-through from the Qiagen DNA column was processed using a *mir*Vana miRNA Isolation Kit (Ambion). This latter step generated RNA preparations that included RNA <200 nt suitable for miRNA analysis. DNA was extracted from blood using the QiaAmp DNA Blood Midi kit (Qiagen).

1 RNA samples were quantified by measuring Abs₂₆₀ with a UV spectrophotometer 2 and DNA guantified by PicoGreen assay. DNA specimens were resolved by 1% agarose 3 gel electrophoresis to confirm high molecular weight fragments. A custom Seguenom 4 SNP panel or the AmpFISTR Identifiler (Applied Biosystems) was utilized to verify that 5 tumor DNA and germline DNA representing a case were derived from the same patient. 6 Five hundred nanograms of each tumor and germline DNA were sent to Qiagen (Hilden, 7 Germany) for REPLI-g whole genome amplification using a 100 µg reaction scale. RNA 8 was analyzed via the RNA6000 Nano assay (Agilent) for determination of an RNA 9 Integrity Number (RIN), and only analytes with a RIN \geq 7.0 were included in this study. 10 Only cases yielding a minimum of 6.9 µg of tumor DNA, 5.15 µg RNA, and 4.9 µg of 11 germline DNA were included in this study.

12

13 **1.1.3 Sample Qualification**

The BCR received tumor samples with germline controls from a total of 187 cases, of which 87 cases qualified and were sent for further genomic analysis. Of the 100 that failed to qualify, 21 cases were disqualified prior to processing, 15 failed for pathology screening, and 64 cases failed due to molecular criteria.

18

Of the 15 that failed pathologic criteria, 14 failed for absence of tumor cells, and one 1 failed for necrosis. The majority of the 64 cases that failed molecular screening had low normal DNA yields (40 cases). The remaining cases had insufficient tumor DNA (5 cases) or low RNA integrity (19 cases). The difference of 8 samples represents those we removed at the very beginning of the AWG as 4 had neoadjuvant therapy, 3 were unpaired (i.e. tumor and normal appeared unmatched), and 1 did not have sufficient DNA.

26

27 1.1.4 Pathology Review

Aperio© scanned H&E stained slides provided by the tissue source sites from 79 tumors were reviewed according to the 2015 WHO classification as definite MPM subtyped as epithelioid, biphasic or sarcomatoid(1). Where subtyping could not be achieved, tumors were classified as not otherwise specified (NOS). H&E stained slides from frozen sections were reviewed in all 79 cases. In 73 cases, an additional H&E slide was provided from a representative formalin fixed paraffin embedded (FFPE) tissue block. Whenever possible, immunohistochemical staining results were obtained from primary pathology reports. Six pathologists participated in the pathology review and
 discrepant interpretations were resolved by consensus between two pathologists who
 reviewed all cases.

4

5 **1.2 PATHOLOGY REVIEW RESULTS**

6 Tumors were reclassified as summarized in Supplementary Figure 1A. A total of 74 tumors were accepted as MPM with 52 epithelioid (Supplementary Figure 1B&C), 13 8 biphasic (Supplementary Figure 2D&2E), 3 sarcomatoid (Supplementary Figure 2F&2G) 9 and 6 NOS. Five tumors were excluded either because the diagnosis of MPM could not 10 be confirmed (3 cases) or because the normal tissue submitted was not adequate (1 11 case), or because the frozen sample submitted as tumor contained only normal tissue (1 12 case).

13 The following immunohistochemical stains had been performed at the primary 14 institutions. Keratin showed positive staining with the following antibodies: AE1/AE3 15 (n=25), CK7 (n=18), CAM5.2 (n=7) and keratin (not further specified, n=6). The following 16 mesothelial stains were positive: calretinin (n=72), WT1 (n=55) and D2-40 (n=35). The 17 following carcinoma markers were negative: CEA-NOS (n=25), CEA-polyclonal (n=7), 18 CEA-monoclonal (n=7), CD15 (n=19, one focally positive tumor was excluded), BER-19 EP4 (n=25, 2 were focally positive), MOC-31 (n=17, 2 focally positive, 2 positive), and 20 TTF-1 (n=52, 1 positive case was excluded).

- 21
- 22

23 **2. Genome-wide LOH Validation Datasets**

24

25 **2.1 ICGC cohort description**

Tumor-normal matched DNA samples from Japanese MPM cases were analyzed by whole exome sequencing. Data was analyzed using an in-house variant caller, Karkinos (http://sourceforge.net/projects/karkinos/), to detect single nucleotide variations (SNV) and allelic somatic copy number alterations (sCNA). Among 80 Japanese cases (48 frozen tissue samples and 32 primary cell lines from cancer patients within 20 passages), two cases showed genome-wide LOH (Figure 2B in the main manuscript).

SNV detection and estimation of tumor cellularity were previously described(2, 3).
 We detected sCNA in an allelic manner and generated an allelic copy number plot from
 exome sequencing data using karkinos, which detects sCNA by calculating the ratio of

allele specific reads between matched tumor and normal samples at the positions with
 heterozygous SNPs, by normalization of raw data, adjustment of GC contents, wavelet
 de-noising, and multi-state HMM.

4 Both genome-wide LOH cases showed loss of one copy of most chromosomes, 5 whilst retaining two copies of chromosomes 5 and 7. Allelic copy number plots for 6 sample K2F2-A45 (primary cell line at passages 9), and sample K2F2-H60 (frozen tissue 7 sample) are shown in Figure 2B of the main manuscript. K2F2-A45 has a homozygous 8 deletion of SETDB1, and CDKN2A, but doesn't have a driver mutation in TP53. K2F2-9 H60 has a stopgain SNV of SETDB1, a splice site mutation of TP53, and a frameshift 10 deletion of NF2. No somatic mutations in BAP1, SETD2, and PBRM1 were detected in 11 either case.

12

13 **2.2 Brigham and Women's genome-wide LOH cohort description**

As cytogenetic data were not available on the TCGA MPM cases, we sought independent validation of this finding in a series of 916 MPM cases prospectively karyotyped at Brigham and Women's Hospital (BWH) in Boston, MA, between 1990 and 2013. Among these, 16 pleural MPM cases (1.7%) with a near-haploid karyotype (Supplementary Table S3) have been identified.

Whilst the mean age of the patients in the BWH MPM cohort was >65, the mean age of the near-haploid subset is 54. Moreover, this subgroup had a significantly higher percentage of female patients (10/16; 62.5%) compared to the overall BWH cohort (211/916; 23% female) (P = 0.005).

- 23
- 24

25

3. Whole Exome Sequencing and Multicenter Variant Calling

26

27 **3.1 Exome enrichment and sequencing**

Genomic libraries were prepared using the Illumina Paired End Sample Prep Kit following the manufacturer's instructions. Enrichment was performed as described previously(4) using the Agilent SureSelect Human All Exon 50Mb kit following the manufacturer's recommended protocol.

Each exome was sequenced using a 75bp paired-end protocol on an Illumina HiSeq DNA Analyzer, to produce approximately 10Gb of sequence per exome. Sequencing reads were aligned to the human genome (NCBI build 37) using the Burrows-Wheeler Aligner (BWA) algorithm with default settings(5). Reads which were unmapped, PCR-derived duplicates, or outside the targeted region of the genome, were excluded from the analysis. The remaining uniquely mapping reads provided 70–90% coverage over the targeted exons at a minimum depth of 30x.

5 6

3.2. Whole Exome Sequencing and Multicenter mutation calling

7 3.2.1 Wellcome Trust Sanger Institute

8 The CaVEMan (Cancer Variants through Expectation Maximization) algorithm 9 was used to call single-nucleotide substitutions(4). To call insertions and deletions, we 10 used split-read mapping implemented as a modification of the Pindel algorithm (6, 7). 11 Mutations were annotated to Ensembl version 58 using the VAGrENT algorithm(8). Post-12 processing of mutation calls was performed to remove recurrent artifacts from the set of 13 initial variant calls as described elsewhere(9). Briefly, regions of recurrent mis-mapping 14 or sequencing errors are removed through excluding regions near homopolymer tracts, 15 germline indels, or with highly recurrent errors in a panel of normal exomes. All indels 16 (n=548) and all putative driver substitutions (n=80) were reviewed by manual inspection 17 of the sequence data. For the indels, 303 variants were true positive, 173 false positive, 18 and 72 ambiguous. False positive indels were excluded from subsequent analyses. 19 From the subset of inspected substitutions 77/80 were true positive.

20

21 **3.2.2 Broad Institute**

The Firehose pipeline (http://www.broadinstitute.org/cancer/cga/Firehose) performed quality control (QC) on the BAM files, point mutation calling, small insertion and deletion detection, annotation of detected mutations, filtering for OxoG artifacts and filtering by "panel-of-normals". These steps are described in further detail below.

- QC on BAM files: The sample cross-individual contamination levels were
 estimated using the ContEst program(10).
- 28 2. Somatic point mutation calling: The MuTect algorithm(11) was used to detect
 29 somatic single nucleotide variants (SNV).
- 30 3. Small insertion and deletion detection: The Indelocator algorithm
 31 (https://www.broadinstitute.org/cancer/cga/indelocator) was used to detect small
 32 insertions and deletions (InDel).
- 33 4. SNV and InDel annotations: variants detected by MuTect and Indelocator were
 34 annotated using Oncotator(12). Oncotator mapped somatic mutations to

 1
 respective genes, transcripts, and other relevant features. These annotations

 2
 correspond to the fields in the TCGA Mutation Annotation Format (MAF) files

 3
 version

 4
 (https://wiki.nci.nih.gov/display/TCGA/Mutation+Annotation+Format+(MAF)+Spec

ification).

5. Filtering for OxoG artifacts: 1741 G>T/C>A transversions that are a consequence
of heating, shearing, and oxidative damage to the DNA during genomic library
preparation(12) were filtered out of the Broad call set.

6. Filtering by "panel-of-normals": The sites of detected SNV and InDel were
examined against a panel of 8313 normal samples (PoN). This identified any
SNV or InDel that is a recurrent artifact and removed variants that were found
outside coding regions. As a result of the filtering, 865 SNVs and 15639 InDels
were removed, leaving a total 3361 mutations in the final MAF (3011 SNVs, 66
insertions and 284 deletions).

15

5

16 **3.2.3 University of British Columbia Cancer Research Centre**

17 Strelka(13) v1.0.6 was used to identify somatic SNVs, and short InDels from the 18 WES dataset. All parameters were set to defaults, with the exception of 19 "isSkipDepthFilters", which was set to 1 in order to skip depth filtration given the higher 20 coverage in the dataset. Blood or normal tissue was used as matched normal specimen. 21 The variants were subsequently annotated using SnpEff(14), and the COSMIC(15) v61 22 and dbSNP(16) v137 databases.

23

24 **3.2.4 Multicenter Mutation Calling**

25 SNV and InDels data from three centers, Broad Institute (BI: 4067 SNVs, 186 26 deletions, and 52 insertions), Sanger Institute (SI; 3687 SNVs, 434 deletions, and 114 27 insertions) and University of British Columbia (UBC; 8370 SNVs, 167 deletions, and 40 28 insertions) were used in the creation of the MCC MAF file. In order to create the set of 29 SNVs, any SNV that appeared in at least two of the three centers was included in the 30 MCC MAF file used for subsequent analysis. Also, any non-coding SNV or InDel was 31 removed from the final set. In cases where coding status from SI was specified, that 32 coding annotation was prioritized and used to determine the coding status of the SNV or 33 InDel; otherwise, BI's annotations were used. Consequently, 3011 SNVs were included 34 in the final MCC MAF file. For insertions and deletions, all SI's manually reviewed InDels

(278 deletions and 64 insertions) and 8 additional InDels (6 deletions and 2 insertions)
that weren't manually reviewed by SI, but called by both BI and UBC and manually
reviewed by BI, were included in the final set of InDel calls. Thus, 350 InDels (284
deletions and 66 insertions) were included in the final MCC MAF file.

5 Because BAP1 alterations are of specific interest in MPM, RNA-Seq InDels (18 6 InDels) and exome calls in BAP1 from BI, SI and UBC were closely analyzed by 7 Analysis Working Group members at the University of North Carolina. Combining the 8 calls across all centers yielded 26 BAP1 SNVs and InDels (9 SNVs and 17 deletions). 9 Three additional BAP1 mutations (two SNVs and one deletion) that were deemed to be 10 affecting cancer progression by the Analysis Working Group were included in the final 11 set of calls as well. The three BAP1 mutations were: intragenic SNV that affected BAP1 12 expression level in ZN-A9VS, low allelic fraction SNV in 3U-A98G, and 48 basepair long 13 intronic deletion in 3U-A98G. In all, 29 BAP1 mutations (11 SNVs and 19 deletions) were 14 in the MCC MAF file (See Supplementary Table S2A for details).

In order to determine the alternate and reference allele counts for mutations, BI and SI computed allele counts. For SNVs, alternate and reference allele counts were determined by the BI's Mutation Validator tool. For InDels, SI calculated allele counts for all but 8 InDels in the MCC MAF file. The allele counts for these remaining InDels, which were called by BI and UBC, were determined by BI's Mutation Validator tool.

20

21 **3.2.5. Mutational spectrum**

22 The observed mutational spectrum of somatic SNVs for each sample in the cohort 23 was extracted by considering each variant in its pyrimidine context (that is, C>A, C>G, 24 C>T, T>A, T>C and T>G). Each of the six types of base substitution were further split 25 into 16 subcategories by the reference base immediately 5' and 3' to the mutated base. 26 For each patient, counts of SNVs in each of these 96 channels then represent the 27 observed mutational spectrum. The dataset did not have sufficient information content 28 for a formal signature extraction for two reasons. First, the low number of variants in 29 each tumor (due to the low overall mutation burden combined with a footprint of only 30 30Mb in an exome) limits statistical power to detect signatures. Second, across patients 31 (allowing for the aforementioned low numbers), there was very little evidence for 32 variation in mutation spectrum, with the exception of the one hypermutated patient. 33 Given that formal signature extraction relies on having individual signatures distributed 34 unequally across patients, any such extraction would be unreliable. Although mutation burden was too low for signature extraction at the level of individual patients, we felt that an analysis comparing the mutational spectra combined across asbestos-exposed patients versus asbestos-unexposed was well powered. This analysis was performed using a chi-squared test.

5

6 **4. Copy number analysis**

7

8 **4.1 METHODS**

9 Affymetrix SNP 6.0 arrays were used to hybridize genomic DNA from each tumor and 10 normal sample using standard protocols at the Genome Analysis Platform of the Broad 11 Institute(17). Briefly, from raw CEL files, Birdseed was used to infer preliminary copy 12 number at each probe locus(18). For each tumor, tangent normalization was applied to 13 estimate genome-wide copy number. Tangent normalization is based on the observation 14 that the linear combination of all normal samples that are most similar to the tumor tends 15 to match the noise profile of the tumor better than any set of individual normal samples: 16 therefore used to divide this linear combination is the tumor signals 17 (19) (http://www.broadinstitute.org/cancer/cga/copynumber_pipeline). Individual copy-18 number estimates then underwent segmentation using Circular Binary Segmentation(20), 19 during which regions corresponding to germline copy number alterations were removed. 20 Ziggurat Deconstruction was then applied to assign a length and amplitude to each 21 identified copy number change, in a way that accounts for different copy number values 22 inferred across the locus from the heterogeneous cell population(21).

23

Allelic CN, whole genome doubling, subclonality, and purity and ploidy estimates were calculated using the ABSOLUTE algorithm(22) and CBS-derived segmented CN values were re-centered using the *In Silico* Admixture Removal (ISAR) procedure(23) and significant focal CN alterations were identified from using GISTIC 2.0.22(21).

28 CN clustering was based on total integer arm-level CN, normalized over 4 copies for 29 tumors estimated to have been whole genome doubled once, and 8 copies for tumors 30 estimated to have been whole genome doubled twice. Tumors were clustered based on 31 thresholded re-occurring alteration peaks GISTIC CN at from analysis 32 (all_lesions.conf_99.txt file). Clustering was done in R based on Manhattan distance 33 using Ward's method. Allelic CN derived from ABSOLUTE was used along with visual inspection of relative CN to determine regions of loss of heterozygosity (LOH) and
 homozygous deletions.

3 The FACETS algorithm(24) was used to perform allele-specific copy number and LOH 4 analysis from WES data. Briefly, read count information was extracted from paired 5 tumor-normal whole-exome sequencing BAM files. Total log-copy-ratio (logR) was 6 computed from the total read count across germline SNP sites in tumor versus normal. 7 GC-normalization was done using loess regression. Allelic imbalances were assessed 8 using the variant allele log-odds-ratio (logOR) at heterozygous SNP sites. A joint 9 segmentation analysis was applied to identify regions of the genome with copy number 10 alterations by extending the Circular Binary Segmentation (CBS) algorithm to a 11 bivariate change-point detection method based on the Hotelling T2 statistic. Integer 12 copy number was estimated using a genotype-mixture model correcting for tumor purity, 13 ploidy and clonal heterogeneity. Software is available at 14 https://github.com/mskcc/facets.

15

16 **4.2 RESULTS**

17 The SCNA landscape of MPM was characterized by frequent recurring focal and arm-18 level deletions, but no recurrent focal amplifications. Using ABSOLUTE-inferred allelic 19 CN information(22), we were able to determine if deletions at each locus were 20 homozygous, hemizygous, or heterozygous. Several tumor suppressor genes 21 previously found to be altered in MPM (25) were recurrently homozygously deleted, 22 including CDKN2A (36/74), BAP1 (12/74), NF2 (6/74), PBRM1 (3/74), SETD2 (3/74), 23 and PTEN (2/74). PTPRD and RBFOX1 were also recurrently deleted, but both have 24 previously been reported to be fragile sites in cancer genomes and their mRNA 25 expression did not correlate with CN(26, 27). Additionally, many of the focally deleted 26 tumor suppressors were also significantly recurrently mutated in the cohort (Fig. 1 in 27 main manuscript).

While CN and mutational status of *BAP1* and *SETD2* were not correlated with overall survival, *NF2* status (stratified by homozygous deletion, mutation, 1-copy LOH, 2-copy LOH, and wild-type; Cox model P=0.024) and *CDKN2A* status (stratified by homozygous deletion and other; Cox model P=7.3x10⁻⁶) conferred significant differences in overall survival. Biallelic inactivation of *NF2* was also significantly enriched in tumors of patients with no history of asbestos exposure (Chi-square P=0.0027). The 19/74 tumors harboring biallelic inactivation of *BAP1* were all
 epithelioid (Chi-square P=0.019).

3 Hierarchical clustering revealed 6 distinct CN clusters (Supplementary Figure 2A). One 4 group (cluster 1) consisted of tumors with focal deletions in chromosome 3p around the 5 BAP1/PBRM1/SETD2 locus as well as partial or whole deletion of chromosome 22, but 6 few other SCNAs. Cluster 4 consisted of tumors either lacking broad SCNAs entirely or 7 exhibiting high levels of genome-wide LOH. With one exception, no homozygous 8 deletions were found in this group. Cluster 2 consisted of 4 hyperdiploid/hypotriploid 9 samples with more arm-level gains than losses. Cluster 3 was enriched for CDKN2A 10 homozygous deletions but was otherwise generally chromosomally stable. Clusters 5 11 and 6 were both characterized by chromosomal instability, but tumors in cluster 5 had 12 partial telomeric loss of chromosome 1p, whereas those in cluster 6 tended to have 13 whole-arm or partial centromeric loss of 1p. The majority of 1p deletions affected only 14 the telomeric or the centromeric ends of 1p, suggesting positive selection for deletion of 15 either, or negative selection for deletion of the genomic region.

16 In GISTIC 2.0 analysis of MPM with wild type or 2-hit BAP1 inactivation, we found no 17 significant differences in focal SCNAs, other than 3p focal deletions containing BAP1 18 itself. However, there were more overall SCNAs in tumors without BAP1 alterations. 19 Total numbers of amplifications and deletions in chromosomal arms were significantly 20 greater in tumors with wild type BAP1 (median 15.5 vs. 9.5, P<0.01). All tumors in CN 21 cluster 4, which contain few SCNAs apart from the 3p focal deletion containing BAP1 22 and chromosome 22/6q deletion, had at least one allele with a BAP1 alteration 23 (Supplementary Figure 2D).

FACETS and ABSOLUTE analyses identified three cases with extensive loss of heterozygosity (Supplementary Figure 2E, 2F and 2G) and these cases were further analyzed using complementary methodologies to characterize their molecular features.

- 27
- 28

29 **5.** Gene Expression (RNA sequencing)

30

31 **5.1 RNA library construction, sequencing, and analysis**

One µg of total RNA was converted to mRNA libraries using the Illumina mRNA
 TruSeq kit (RS-122-2001 or RS-122-2002) following the manufacturer's directions.
 Libraries were sequenced 48x7x48bp on the Illumina HiSeq 2000 as previously

1 described(28), with FASTQ files generated by CASAVA. RNA reads were aligned to the 2 hg19 genome assembly using MapSplice(29) 0.7.4. Gene expression was quantified for 3 transcript models corresponding to the TCGA **GAF2.1** (http://tcgathe 4 data.nci.nih.gov/docs/GAF/GAF.hg19.June2011.bundle/outputs/TCGA.hg19.June2011.g 5 af), using RSEM(30) and normalized within-sample to a fixed upper quartile. Further 6 details on this processing are available in the Description file at the DCC, under the 7 V2_MapSpliceRSEM workflow (https://tcga-8 data.nci.nih.gov/tcgafiles/ftp_auth/distro_ftpusers/anonymous/tumor/meso/cgcc/unc.edu/ 9 illuminahiseq_rnaseqv2/rnaseqv2/unc.edu_MESO.IlluminaHiSeq_RNASeqV2.mage-10 tab.1.0.0/DESCRIPTION.txt) CGHUB at or 11 (https://cghub.ucsc.edu/docs/tcga/UNC_mRNAseq_summary.pdf).

12

13 **5.2 RESULTS**

14 Unsupervised gene expression using variably and highly expressed genes was 15 used to identify 4 mRNA-based clusters. Using ClaNC, we identified a set of 1,800 16 genes that strongly classified each of the four clusters and associated them with clinical 17 and genomic features (Supplementary Figure 3A). Cluster 1 was predominantly 18 epithelioid, and significantly less likely to have loss of CDKN2A/B (Fisher's exact test, 19 p<0.001, Supplementary Figure 3B), contained 5/8 SETD2 mutations (Fisher's Exact 20 Test, P=0.03), had high expression of innate immune cells, MSLN (concordant with 21 serum levels), and had the best overall survival (Supplementary Figure 3C). In contrast, 22 cluster 3 predominantly contained non-epithelioid histologies and had significantly poorer 23 outcomes compared to cluster 1 (HR 4.3, P=0.00009). Samples in this cluster 24 expressed genes involved in NOTCH and WNT signaling and had high levels of cell 25 cycle gene expression. This group of patients was also significantly older than the other 26 three (t test, P=0.0008, Supplementary Figure 3D). Cluster 4, while epithelioid, was 27 similar to 3 - with high expression of cell cycle genes and a poor overall survival 28 compared to 1 (HR 6.8, P=0.00003). Cluster 2 was defined by high expression of BAP1 29 and CLDN genes. This group was significantly less likely to have BAP1 alterations (t 30 test, P=0.001, Supplementary Figure 3A, 3B, 3E), and had similar survival to cluster 1, 31 with a younger age at diagnosis. Using the gene signatures from Bueno(25) and de 32 Reynies(31) for validation, samples clustered similarly to our TCGA analysis, with the 33 highest concordance observed for non-epithelioid samples.

34

1

3

2 6.

DNA Methylation Profiling

4 6.1 METHODS

5 The Illumina Infinium HM450 array(32) was used following standard protocols. 6 Briefly, genomic DNA (1 µg) was treated with sodium bisulfite and recovered using the 7 Zymo EZ DNA methylation kit (Zymo Research, Irvine, CA). Bisulfite-converted DNA 8 samples were amplified, fragmented and hybridized to BeadChips, followed by a locus-9 specific base extension with labeled nucleotides (cy3 and cy5). BeadArrays were 10 scanned and the raw data imported into custom programs in the R computing language 11 for pre-processing and calculation of DNA methylation β values for each probe and 12 sample. Quality control and probe exclusions were performed using standard protocols 13 as described(33).

14 Unsupervised consensus clustering was performed in the Bioconductor package 15 ConsensusClusterPlus(34) v1.22.0, with Euclidean distance and partitioning around 16 medoids (PAM) and was applied to the DNA methylation data using the most variable 17 1% of CpG probes. Fisher's exact test was used to test for associations of DNA 18 methylation clusters with other platform clusters and significantly mutated genes.

19 Leukocyte fraction was estimated as described previously(22). As a source of 20 leukocyte DNA methylation level, we used DNA methylation data of peripheral blood 21 mononuclear cells from six healthy donors(35) (GSE35069).

22 To identify CpG probes associated with BAP1 and SETD2 status, we used 23 empirical Bayes-modified t-tests as implemented in the limma package(36). The 24 correlations between SETD2 and BAP1 status and DNA methylation was strong enough 25 to support two criteria for probe selection - signature genes with FDR<0.01, and the 26 mean difference of β values between altered and wild type samples of more than 0.3.

27

28 6.2 RESULTS

29 Three robust DNA methylation clusters were identified based on the most 30 variable CpG loci on the Illumina array (Supplementary Figure 4A). Cluster 3 had a 31 higher leukocyte fraction and lower purity than the other two and was associated with 32 miRNA cluster 5, IncRNA cluster 4, PARADIGM cluster 4, and iCluster 4 (Supplementary 33 Figure 4A and 4B). Cluster 2 was enriched with SETD2 mutations, which have been 34 associated with altered DNA methylation in cancer(33, 37). To explore further the

1 correlation of SETD2 mutations and DNA methylation, we looked for CpG sites that 2 associated with SETD2 status and found more than 200 differentially methylated CpG 3 sites (Supplementary Table S6A). The association between BAP1 status and DNA 4 methylation was strong enough to support two criteria for probe selection, signature 5 genes were selected to have FDR<0.01 and the mean difference of beta values between 6 altered and wild type samples was more than 0.3. We found 84 CpG probes that had 7 significantly higher methylation in BAP1 inactivated samples (Supplementary Figure 4C). 8 For most of the probes, higher methylation was associated with low expression for the 9 corresponding gene (Supplementary Table S6B).

- 10
- 11

12 **7.** Noncoding RNA (IncRNAs and miRNAs) expression analyses

13

14 **7.1 METHODS**

15 **7.1.1 RNA-seq read mapping**

16 RNAseq FASTQ files for TCGA data (n=74) were downloaded from CGHub(38), 17 while those for from the validation cohort (Bueno(25), n=211) were downloaded from the 18 European Genome-phenome Archive (EGAS00001001563). For both datasets, FASTQ 19 files were processed using STAR(39) v2.3.0e with the following parameters: 20 minimum/maximum intron size set to 30 and 500,000, respectively, noncanonical, 21 unannotated junctions were removed, the maximum of tolerated mismatches set to 10, 22 and the outSAMstrandField intron motif option enabled. The Cuffdiff command included 23 with Cufflinks(40) v2.0.2 was used to calculate the fragments per kilobase of exon per 24 million fragments mapped (FPKM) with upper quartile normalization, fragment bias 25 correction, and multiread correction enabled. All other options were set to default. 26 Ensembl v82 gene annotations were used.

We generated microRNA sequence (miRNA-seq) data using methods described previously(41, 42) with miRBase v16 annotations, and assigned 5p and 3p mature strand names using miRBase v20.

30

31 7.1.2 Unsupervised consensus clustering

For the TCGA cohort, we extracted 347 expressed (mean FPKM ≥1) and highly
 variable (95th percentile, FPKM variance) IncRNAs from a normalized abundance matrix
 of 8167 IncRNAs (7671 lincRNA and 496 processed transcripts). We identified groups

of samples with similar abundance profiles by unsupervised consensus clustering with
 ConsensusClusterPlus v1.36.0. Calculations were performed using Pearson correlations,
 hierarchical clustering, 20,000 iterations, and a random gene fraction of 0.975 in each
 iteration.

5 For the Bueno validation cohort, we clustered 420 expressed, high-variance 6 IncRNA profiles using Pearson correlations, partitioning around mediods, 5000 iterations, 7 and a random 0.95 gene fraction in each iteration. We selected a four-cluster solution for 8 both cohorts.

9 For miRNA mature strand data, we selected an input file containing reads-per-10 million (RPM) data for the 303 (25% of 1212) most-variant 5p or 3p mature strands. We 11 used unsupervised consensus clustering with ConsensusClusterPlus v1.36.0, using a 12 Pearson distance, hierarchical clustering, 20000 iterations and a gene fraction of 0.95.

To generate a heatmap for subtypes, we first used a SAM multiclass analysis(43) to identify differentially abundant lncRNAs or miRNAs (FDR<0.05), which we filtered to retain lncRNAs or miRs that were expressed at least a mean of 5 FPKM or 25 RPM respectively(44). We transformed each row of the matrix by log10(RPM+1), then used the pheatmap R package (v1.0.2) to scale and cluster only the rows, and to generate the heatmap.

An FDR threshold of 0.05 was set for both miRNA and IncRNA.

19

20

21 **7.1.3 Maximal-search associations with survival**

For both the TCGA and the validation cohorts, we identified miRNAs and IncRNAs that were associated with overall survival. We used R functions from CutoffFinder(45) v2.1 to determine the RPM or FPKM value that stratified samples into two groups(46), then adjusted the corrected p values for multiple testing with the Benjamini-Hochberg method.

27

28 **7.1.4 Correlations with EMT scores**

We used MatrixEQTL(47) to calculate Spearman correlations (FDR<0.05) between RNAseq-based EMT scores (Supplementary Section 10) and miR RPMs or IncRNA FPKMs.

32

33 7.2 RESULTS

We assessed two types of noncoding RNAs that may offer molecular insights into
 MPM(48): microRNAs (miRNAs) and long noncoding RNAs (IncRNAs).

miRNA mature strands (miRs) are associated with MPM (48-51). Here, we used unsupervised consensus clustering to identify five miR subtypes (Figure 6H). These were associated with purity (P= 1.0×10^{-7}), leukocyte fraction (P= 1.1×10^{-6}), EMT scores (P= 1.7×10^{-7}), and 5-year survival (P= 6.6×10^{-4}). They were concordant with subtypes for IncRNAs (see below, P= 1.2×10^{-12}), mRNA (P= 1.2×10^{-14}), iCluster (P= 1.3×10^{-12}), PARADIGM (P= 3.0×10^{-13}), and were associated with subtypes for SCNA (P= 2.6×10^{-4}), and DNA methylation (P= 1.6×10^{-4}).

LncRNA expression can be more specific for cell type than coding gene expression(52-54). To our knowledge, only one expression lncRNA profiling report, based on NCode long noncoding microarrays, is available for MPM(55). Recently, multiplatform data that included transcriptome sequencing, and analyses focused on protein-coding genes, were reported for a large MPM cohort(25); this 'Bueno' transcriptome data served as an independent lncRNA data source.

16 For the TCGA cohort, we used unsupervised consensus clustering with 17 transcriptome sequence data to identify 4 IncRNA subtypes (Figure 6A). These were 18 associated with leucocyte fraction (P=3.2x10⁻⁴), EMT scores (P=1.8x10⁻⁴), and 5-year 19 survival (log-rank $P=1.4x10^{-4}$) (Figures 6A-D). They were strongly concordant with 20 subtypes from mRNA ($P=3.5x10^{-17}$), iCluster ($P=7.7x10^{-14}$), PARADIGM ($P=2.1x10^{-21}$), 21 and miRNA (P=1.4x10⁻¹²), and were associated with subtypes for SCNA (P=0.011), and 22 DNA methylation (P=0.025). Samples in IncRNA cluster 1 (n=20, 27%) had better 23 survival, while cluster 4 had a relatively high leukocyte fraction and low purity, many of 24 the non-epithelioid samples, and relatively shorter survival.

25 We noted that the iCluster, PARADIGM, miRNA and IncRNA subtypes that had 26 the best and worst survival were strongly concordant.

For the Bueno validation dataset(25), we identified 4 lncRNA subtypes that were concordant with the four RNA-seq-based subtypes reported in that work (P= $3.4x10^{-28}$) and were associated with 5-year survival (P= $1.1x10^{-3}$) (Figure 6E). For both the TCGA and the Bueno cohorts, lncRNAs that were differentially abundant between the bettersurvival subtype and other samples included those known to be associated with cancers in general (e.g. H19, LINC00152, MEG3) or with MPM in particular: NEAT1 and SNHG8(55), and GAS5(56).

1 EMT is a factor in MPM(31, 57). Many miRNAs(58-60) and IncRNAs(61) have 2 been associated with EMT. We noted that the miR-based and IncRNA-based 3 unsupervised clusters that had better survival had the lowest median EMT scores. In the 4 TCGA cohort, we identified miRs and lncRNAs that were statistically correlated (FDR < 5 0.05) with RNA-seq-based EMT scores (Table S5A). As expected, miR-200 family 6 members, members of an Xq27.3 genomic miRNA cluster(62), and miR-29 family 7 members(63) were correlated with the scores, as was the oncogenic LINC00152(64) 8 (ρ=0.41).

9 Next, we assessed miRs and IncRNAs that were differentially abundant in the 10 subtypes that had better survival (Figures 6F,G,L; Table S5B). For the TCGA cohort, 11 miR-126, 143-3p and 145-5p(65, 66) were less abundant in miR cluster 1, while miR-12 193a-3p(67) was more abundant. Among the differential IncRNAs in both the TCGA 13 (IncRNA cluster 1) and the Bueno (IncRNA cluster 4) cohorts (Figures 6F,G), the most 14 highly differential IncRNA RP11-263K4.5 showed a positive fold-change of 63.6 in 15 TCGA and 9.1 in the validation datasets, suggesting that this IncRNA may be a tumor 16 suppressor in MPM.

17 Finally, for all samples and then for only epithelioid samples, we identified 18 miRNAs that were statistically associated with survival in the TCGA cohort, and IncRNAs 19 that were that were associated with survival in both TCGA and Bueno cohorts (Table 20 S5C). For the full TCGA cohort, miRs with significant adjusted p-values (e.g. miR-514a-21 3p, 508-3p, 29b-2-5p, 101-3p, Hazard Ratio<1) included those that were negatively 22 correlated to EMT scores, and those that were relatively abundant in the miR cluster with 23 the best survival. For the epithelioid cases in the TCGA cohort, miR-148b-3p (P=5.7x10⁻ 24 ⁴) and miR-148a-5p (P=9.9x10⁻⁴) were statistically significant, though miR-148b-3p was 25 not statistically significant for the full cohort. miR-148b-3p is well known in lung and other 26 cancers(68). A number of IncRNAs were statistically associated with survival in the full 27 TCGA and validation cohorts. For example, GS1-600G8.5 (P=1.1x10⁻³ and 0.036, 28 respectively, Hazard Ratio>1) was less abundant in the good-survival IncRNA subtypes 29 in both cohorts, and was positively associated with TCGA EMT scores. Taken together, 30 the above results suggest that both miRNAs and IncRNAs may be important regulators 31 of EMT and of survival in MPM.

- 32
- 33

34 8. Microbial sequences

1

2 8.1 METHODS

3 We screened RNA and DNA sequence reads with a microbial detection pipeline 4 based on BioBloomTools (BBT), v1.2.4.b1, a Bloom filter-based method for rapidly 5 classifying RNA-seq or DNA-seq read sequences(69). We ran BBT in paired-end mode 6 to screen FASTQ files from 74 tumor RNA-seg libraries, 74 tumor and their respective 7 normal whole exome libraries within the TCGA MPM cohort. In a single-pass scan for 8 each library, BBT categorized each read pair as matching the human, a unique microbial, 9 more than one (multi-match), or no-match filters. We then calculated a reads-per-million 10 (RPM) abundance for each filter.

11 To detect integration of human herpes virus (HHV) and human papilloma virus 12 (HPV), we performed de novo assembly with ABySS v1.3.4. We assembled only the 13 reads classified by BBT, then merged the sets with Trans-ABySS(70) 1.4.8 to generate a 14 contig working set. We identified breakpoint locations by using BLAT(71) v34 to align to the GRCh37/hg19 and to either 109 HHV or 268 HHV reference sequences. We 15 16 retained contig alignments in which: a) the aligned human and viral sequences summed 17 to at least 90% of the contig length, and b) the human and viral aligned overlapped by 18 less than 50%. Human breakpoint coordinates were annotated against RefSeq and 19 UCSC gene annotations(72). Breakpoints that had at least 3 spanning mate-pair reads 20 or 5 flanking mate-pair reads were considered potential integration sites.

21

22 8.2 RESULTS

We assessed microbes using RNAseq and WES data for 74 MPM samples (Supplementary Table S7). The RNAseq libraries returned signals for common microbial contaminants. We found no evidence for genomic integration in 9 libraries that returned positive RPMs: one RNAseq library that showed weak signals for HPV and 8 tumor WES libraries that showed weak signals for HHV4 (EBV). Notably, there were no hits for polyomavirus (including SV40) sequences in any of the libraries.

- 29
- 30

31 9. Reverse-Phase Protein Array analysis

- 32
- 33 9.1 METHODS

Protein was extracted using RPPA lysis buffer from tumors and RPPA was performed as described previously(73-77). SDS-reduced lysates were adjusted to 1 µg/µL and manually serially diluted. An Aushon Biosystems 2470 arrayer (Burlington, MA) printed 1,056 samples on nitrocellulose-coated slides (Grace Bio-Labs), which were probed with 219 validated primary antibodies followed by corresponding secondary antibodies.

7 Signal was captured using a colorimetric, DAB-based DakoCytomation-catalyzed 8 system. Slides were scanned in a CanoScan 9000F. Spot intensities were analyzed and 9 quantified using Array-Pro Analyzer (Media Cybernetics Washington DC) to generate 10 spot signal intensities (Level 1 data). SuperCurveGUI(75, 77) was used to estimate EC_{50} 11 values of the proteins in each dilution series(73). A QC metric(77) was returned to help 12 determine the quality of each slide: if the score was less than 0.8 on a 0-1 scale, the 13 slide was omitted. In most cases, the staining was repeated and the highest QC scoring 14 slide was used for analysis (Level 2 data).

Protein measurements were corrected for loading(75, 77, 78) using median centering across antibodies (level 3 data). Final selection of antibodies was also driven by the availability of high quality antibodies, as assessed by specificity, sensitivity and dynamic range for quantification(79). In total, 219 antibodies and 52 MPM samples were used for the analysis. RPPA arrays were quantitated and processed as described previously(73, 75). Raw data (level 1), SuperCurve nonparameteric model fitting (level 2), and loading corrected data (level 3) were deposited at the DCC.

22

23 9.2 RESULTS

MPM samples were consensus clustered using 1-Pearson correlation as the distance metric and Wards' linkage algorithm. We identified five robust clusters with differential pathway expression (Supplementary Figure 5A and 5B).

The role of cell signaling networks in MPM was illustrated by computing 12 pathway scores described previously(80). The analysis showed that there were differences among the RPPA clusters for many of these pathways (Supplementary Figure 5C).

RPPA cluster 1 (n=11) showed higher EMT, Hormone receptor, RAS/MAPK,
 Breast reactive and Core reactive pathway expressions. Significantly higher
 FIBRONECTIN and COLLAGEN VI, and lower E-CADHERIN and BETA CATENIN
 expression levels contributed to higher EMT in 1, which also showed high apoptosis

1 activity; lower expression in PCNA associated with Cell cycle pathway; CHK1PS345, 2 CHK2PT68, MRE11 and RCC1 associated with high DNA damage response pathway 3 activity: and high expression of MYH11 from the Breast reactive pathway activity. Cluster 4 2 showed relatively high Hormone signaling (breast) pathway activity, which was due to 5 a higher expression of BCL2. Cluster 3 showed relatively low apoptosis, EMT, Hormone 6 receptor and Core reactive activity, and high PI3K/AKT and TSC/mTOR pathway activity. 7 whilst cluster 4 showed high Cell cycle, DNA damage response and EMT activity, 8 coupled with low RAS/MAPK and Breast reactive activity, and had the worst prognosis. 9 Cluster 5 (n=11) had the best prognosis. It showed high activity of DNA damage 10 response pathway, and low activity of EMT, RAS/MAPK, TSC/mTOR, Breast reactive 11 and Core reactive pathways. The most differentially expressed proteins between clusters 12 5 (best prognosis) and 4 (worst prognosis) were PAI1, CYCLINB1, CAVEOLIN1, EPPK1 13 and EEF2K. The worst prognosis cluster had significantly higher expression of PAI1, has 14 been associated with poor prognosis in previous studies, and could be a potential 15 therapeutic target in MPM(81, 82). CYCLINB1(83), CAVEOLIN1, EPPK1 and EEF2K(84) 16 could also be prognostically and/or therapeutically relevant.

17

18

19 **10. Epithelial-Mesenchymal Transition Analysis**

20

21 **10.1 METHODS**

74 MPM samples were scored based on expression of epithelial-mesenchymal
transition (EMT) signature genes using a method previously developed by our group(85).
Briefly, the EMT score for each sample is calculated as the mean expression of epithelial
markers subtracted from the mean expression of mesenchymal markers. Higher EMT
scores correlate with a more mesenchymal profile.

27

28 **10.2 RESULTS**

EMT is known to be common in MPM(31, 57). Across multiple tumor types, MPM possessed the second highest overall EMT score after sarcoma (Fig. 7B). Although the vast majority of MPM tumors had EMT scores greater than 0, indicating higher expression of mesenchymal versus epithelial genes, the individual scores varied significantly. Histology correlated with EMT score, with epithelioid MPM possessing the lowest EMT score (Figure 7A). However, a few epithelioid tumors possessed high EMT

1 scores), atypical for this histology and suggesting distinct epithelial and mesenchymal 2 biologies exist within otherwise similar histologies that may predict distinct responses to 3 targeted therapies. Across all histologies, increasing EMT scores were significantly 4 correlated with higher leukocyte fraction (r=0.30; P=0.008). EMT is known to contribute 5 to immune escape, and we have previously observed an association between EMT 6 score and increased expression of immune checkpoint and other potentially targetable 7 immune genes. In light of these data and recent clinical data suggesting a subset of 8 MPM patients respond to immune checkpoint blockade, we investigated whether EMT 9 was associated with expression of immune genes in MPM. As in other solid tumors(86), 10 EMT score was significantly associated with the expression of many potentially 11 targetable immune checkpoint genes, including OX40 ligand (OX40L), Transforming 12 Growth Factor Beta 1 (TGFB1), B7-H3 (aka Cluster of Differentiation 276, CD276), 13 OX40 receptor (OX40) and Programmed Death-Ligand 2 (PD-L2) (P<0.001), while the 14 V-domain Ig suppressor of T cell activation (VISTA), a negative regulator of T cell 15 proliferation and T-cell cytokine production(87) was strongly associated with low EMT 16 score (r=-0.476, P=1.56x10⁻⁵). However, VISTA expression was higher in MPM 17 compared to all other tumor types available in the TCGA due to high VISTA expression 18 in most epithelioid MPM samples. Thus, EMT score could be used to predict enrichment 19 of specific immune checkpoint targets for selection from already available 20 immunotherapy agents. The correlation between EMT and immune checkpoint gene 21 expression also supports combinatorial strategies to target both immune checkpoint and 22 EMT simultaneously. Additionally, these data suggest that epithelioid histology, with its 23 relatively low EMT score, may predict response to targeting of VISTA, a negative 24 immune checkpoint molecule that, when blocked or removed, appears to slow tumor 25 growth in a mechanism distinct from PD-1 signaling in preclinical models(88, 89). As 26 previously observed(90), we found that EMT scores were also significantly associated 27 with subgroup classification via unsupervised analysis across multiple platforms 28 including integrative platforms such as iCluster and PARADIGM, as well as individual 29 platforms like mRNA, miRNA, lncRNA, methylation and RPPA (Fig. 7A).

- 30
- 31
- 32

11. iCluster Supplementary Information

33

34 **11.1 METHODS**

We utilized iCluster(91), which formulates the problem of subgroup discovery as
 a joint multivariate regression of multiple data types with reference to a set of common
 latent variables that represent the underlying tumor subtypes 1-4.

.

4 Five molecular assay platforms - SCNA, DNA methylation, mRNA expression, 5 IncRNA and miRNA expression were provided as input. Data were pre-processed using 6 the following procedures. CBS-segmented SCNA data was further reduced to a set of 7 non-redundant regions as described(92). For methylation data, the median absolute 8 deviation was employed to select the top 4000 most variable CpG sites after β -mixture 9 quantile normalization(93). Methylation probes with >20% or more missing data and 10 those corresponding to SNP and autosomal chromosomes were removed. For mRNA, 11 IncRNA, and mature strand miRNA sequence data, poorly expressed genes were 12 excluded based on median-normalized counts, and variance filtering led to a list of 13 reduced features for clustering. mRNA, IncRNA and miRNA expression features were 14 log2 transformed, normalized and scaled before using as an input to iCluster.

15

16 **11.2. VALIDATION COHORTS**

17 **11.2.1 Bueno Cohort**

18 Normalized gene expression profiles for 211 MPM samples previously 19 described(25) (Bueno cohort) were log-transformed and row-centered. For each sample, 20 Pearson correlation was computed between the profile and the centroids of the four 21 iClusters across the detected 2,606 of the 2,807 signature genes; the sample was 22 assigned to the iCluster with the highest correlation. Kaplan-Meier plots were generated 23 and multivariable Cox regression analysis was performed on the assigned tumor 24 samples, where iCluster assignment, histology, and age were considered as covariates. 25 We performed a similar analysis for the 141 epithelioid samples, employing 2,123 of the 26 2,292 signature genes for this histological subtype. Kaplan-Meier survival curves were 27 generated and Cox regression analysis was performed on the assigned epithelioid tumor 28 samples, where iCluster assignment and age were considered as covariates.

29

No significant predictive value of iCluster assignments was confirmed for the whole 211-case cohort (Supplementary Figure 6A and 6B), which could be attributed to the vastly different histological type distributions between the TCGA and the Bueno cohorts. However, when the analysis was restricted to the epithelioid cases, cases assigned to epithelioid iCluster 1 had a significant survival advantage, even when
 adjusted for age (Fig. 4D, Supplementary Figure 6C)

3

4 **11.2.2 Lopez-Rios Cohort:**

5 We log-transformed and row-centered the gene expression data (Affymetrix 6 U133A arrays) from 52 previously published MPM cases(94) (Lopez-Rios) for 1,995 7 genes from the 2,807 signature genes present on the Affymetrix chip. The cohort size 8 was too small to allow histology-based stratification, so we did not perform an 9 epithelioid-only analysis. Pearson correlation coefficients were computed between each 10 sample's profile and the centroids of the four iCluster groups. Each sample was 11 assigned to the iCluster with the highest correlation (Supplementary Figures 6D and 6E). 12 Log rank test was performed on the 49 samples with survival data, and demonstrated a 13 significant survival difference between clusters (Supplementary Figure 6F). Multivariable 14 Cox regression analysis was performed, with histology and age as covariates 15 (Supplementary Figure 6G).

- 16
- 17

18

12. PARADIGM supplementary Information

19

20 **12.1 METHODS**

21 We used median centered, log scaled mRNA expression and SCNA 22 GISTIC2(21) data to calculate inferred pathway activity levels using PARADIGM(95). 23 Integrative PARADIGM analysis of the 74 MPM cases clustered using 24 ConsensusClusterPlus(34) identified 4 distinct clusters. To compare Cluster 4 (worst 25 prognosis) to Cluster 1 (best prognosis), we ran PATHMARK(37) on the statistically 26 significant differential activities obtained from SAM to extract connected components of 27 the global PARADIGM regulatory network. Activities that fall outside 2 standard 28 deviations outside of the empirical distribution of the statistically significant differentials 29 are included the final result. A network connection is extracted if both vertices in that 30 connection pass the filter. Networks are then visualized using Cytoscape(96) and 31 CircleGraph(97).

32

33 **12.2 RESULTS**

1 Cluster 4 patients have significantly worse survival and show upregulation in 2 AURKA, E2F targets, G2M checkpoints, as well as PI3K and mTOR pathways. This 3 cluster faithfully recapitulates platform-specific clustering. Cluster 1 patients have the 4 best prognosis, and have upregulation in EGFR signaling, whilst kinase subnetworks are 5 downregulated. There is an enrichment of epithelioid patients in this cluster, however, 6 even upon correction, the group remains distinct. Patients are also more likely to have 7 undergone pneumonectomy. Patients in clusters 2 and 3 have similar prognosis, but are 8 genomically distinct.

- 9
- 10

11 **13. Regulome Explorer**

13 **13.1. METHODS**

14 13.1.1 Integrated Analysis and Interactive Exploration

15 We have integrated all of the data types produced by TCGA and described in this 16 paper into a single "feature matrix" for MPM. From this comprehensive dataset, 17 significant pairwise associations have been inferred and can be visually explored using 18 interactive web Regulome Explorer, an application 19 (http://explorer.cancerregulome.org). This application allows interactive exploration of 20 significant associations between molecular features, between molecular features and 21 derived numeric features, and between molecular and categorical features, such as 22 clinical parameters or cluster assignments. In addition to associations inferred directly 23 from the TCGA data, additional sources of information are integrated into the 24 visualization (e.g., NCBI Gene, miRBase, UCSC Genome Browser, etc).

25

26 **13.1.2 Feature Matrix Construction and Pairwise Statistical Significance**

A feature matrix was constructed using all available clinical, sample, and molecular data for 74 unique patient/tumor samples. The molecular data includes all analytical platforms described here. For mRNA and miRNA expression, quantification files were log2 transformed, and filtered to remove low-variability targets (bottom 25%). For methylation data, probes were filtered to remove the bottom 25%. For somatic mutations, several binary features indicating the presence or absence of a mutation in each sample were generated. Statistical association among data types was evaluated by pairwise comparisons within the feature matrix. P values for the associations between and among clinical and molecular data were computed using appropriate statistical tests for each pair. To account for multiple-testing bias, the p value was adjusted using the Bonferroni correction.

- 6
- 7
- 8

1 14. VISTA Immunohistochemistry

2

3 14.1. Methods

4 Two epithelioid MPM cases contributed to the TCGA cohort by MSKCC (TCGA-SC-5 A6LQ-01 and TCGA-SC-A6LM-01) were selected for VISTA immunohistochemical 6 studies based on availability of additional FFPE tumor tissue. Blocks were sectioned, de-7 paraffinized, and stained with a fully automated system (Benchmark ULTRA; Ventana 8 Medical Systems, Tucson, AZ), using the rabbit monoclonal anti-VISTA antibody, clone 9 D1L2G, 0.1 µg/mL (Cell Signaling Technology, Danvers, MA, USA). Additionally, VISTA 10 expression was assessed in normal mesothelial lining from pleura and benign pleuritis 11 with reactive mesothelial proliferation (derived from the Department of Pathology at 12 MSKCC). Spleen and colon FFPE tissue was used as positive control for antibody 13 specificity. 14 15 16 17 REFERENCES 18 19 1. Travis WD, Brambilla, E., Burke, A.P., Marx A., Nicholson A.G. WHO 20 Classification of Tumours of the Lung, Pleura, Thymus and Heart. 4 ed. Lyon: 21 International Agency for Research on Cancer; 2015. 22 2. Kakiuchi M, Nishizawa T, Ueda H, Gotoh K, Tanaka A, Hayashi A, et al. 23 Recurrent gain-of-function mutations of RHOA in diffuse-type gastric carcinoma. Nature 24 genetics. 2014;46:583-7. 25 Totoki Y, Tatsuno K, Covington KR, Ueda H, Creighton CJ, Kato M, et al. Trans-3. 26 ancestry mutational landscape of hepatocellular carcinoma genomes. Nature genetics. 27 2014;46:1267-73. 28 4. Varela I, Tarpey P, Raine K, Huang D, Ong CK, Stephens P, et al. Exome 29 sequencing identifies frequent mutation of the SWI/SNF complex gene PBRM1 in renal 30 carcinoma. Nature. 2011;469:539-42. 31 Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler 5. 32 transform. Bioinformatics. 2010;26:589-95. 33

33 6. Ye K, Schulz MH, Long Q, Apweiler R, Ning Z. Pindel: a pattern growth approach
34 to detect break points of large deletions and medium sized insertions from paired-end
35 short reads. Bioinformatics. 2009;25:2865-71.

Raine KM, Hinton J, Butler AP, Teague JW, Davies H, Tarpey P, et al. cgpPindel:
 Identifying Somatically Acquired Insertion and Deletion Events from Paired End

Sequencing. Current protocols in bioinformatics / editoral board, Andreas D Baxevanis
 [et al]. 2015;52:15 7 1-2.

Menzies A, Teague JW, Butler AP, Davies H, Tarpey P, Nik-Zainal S, et al.
 VAGrENT: Variation Annotation Generator. Current protocols in bioinformatics / editoral
 board, Andreas D Baxevanis [et al]. 2015;52:15 8 1-1.

9. Papaemmanuil E, Cazzola M, Boultwood J, Malcovati L, Vyas P, Bowen D, et al.
7 Somatic SF3B1 mutation in myelodysplasia with ring sideroblasts. The New England
8 journal of medicine. 2011;365:1384-95.

9 10. Cibulskis K, McKenna A, Fennell T, Banks E, DePristo M, Getz G. ContEst:
10 estimating cross-contamination of human samples in next-generation sequencing data.
11 Bioinformatics. 2011;27:2601-2.

11. Cibulskis K, Lawrence MS, Carter SL, Sivachenko A, Jaffe D, Sougnez C, et al.
Sensitive detection of somatic point mutations in impure and heterogeneous cancer
samples. Nature biotechnology. 2013;31:213-9.

12. Ramos AH, Lichtenstein L, Gupta M, Lawrence MS, Pugh TJ, Saksena G, et al.
 Oncotator: cancer variant annotation tool. Human mutation. 2015;36:E2423-9.

13. Saunders CT, Wong WS, Swamy S, Becq J, Murray LJ, Cheetham RK. Strelka:
accurate somatic small-variant calling from sequenced tumor-normal sample pairs.
Bioinformatics. 2012;28:1811-7.

14. Cingolani P, Platts A, Wang le L, Coon M, Nguyen T, Wang L, et al. A program
for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff:
SNPs in the genome of Drosophila melanogaster strain w1118; iso-2; iso-3. Fly.
2012;6:80-92.

Forbes SA, Tang G, Bindal N, Bamford S, Dawson E, Cole C, et al. COSMIC (the
Catalogue of Somatic Mutations in Cancer): a resource to investigate acquired mutations
in human cancer. Nucleic acids research. 2010;38:D652-7.

Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, et al.
dbSNP: the NCBI database of genetic variation. Nucleic acids research. 2001;29:308-11.

17. McCarroll SA, Kuruvilla FG, Korn JM, Cawley S, Nemesh J, Wysoker A, et al.
Integrated detection and population-genetic analysis of SNPs and copy number variation.
Nature genetics. 2008;40:1166-74.

18. Korn JM, Kuruvilla FG, McCarroll SA, Wysoker A, Nemesh J, Cawley S, et al.
Integrated genotype calling and association analysis of SNPs, common copy number
polymorphisms and rare CNVs. Nature genetics. 2008;40:1253-60.

35 19. Cancer Genome Atlas Research N. Integrated genomic analyses of ovarian36 carcinoma. Nature. 2011;474:609-15.

37 20. Olshen AB, Venkatraman ES, Lucito R, Wigler M. Circular binary segmentation
38 for the analysis of array-based DNA copy number data. Biostatistics. 2004;5:557-72.

- 1 21. Mermel CH, Schumacher SE, Hill B, Meyerson ML, Beroukhim R, Getz G.
- 2 GISTIC2.0 facilitates sensitive and confident localization of the targets of focal somatic 3 copy-number alteration in human cancers. Genome Biol. 2011;12:R41.
- 4 22. Carter SL, Cibulskis K, Helman E, McKenna A, Shen H, Zack T, et al. Absolute
 5 quantification of somatic DNA alterations in human cancer. Nat Biotechnol. 2012;30:4136 21.
- Zack TI, Schumacher SE, Carter SL, Cherniack AD, Saksena G, Tabak B, et al.
 Pan-cancer patterns of somatic copy number alteration. Nature genetics. 2013;45:113440.
- Shen R, Seshan VE. FACETS: allele-specific copy number and clonal
 heterogeneity analysis tool for high-throughput DNA sequencing. Nucleic Acids Res.
 2016;44:e131.
- Bueno R, Stawiski EW, Goldstein LD, Durinck S, De Rienzo A, Modrusan Z, et al.
 Comprehensive genomic analysis of malignant pleural mesothelioma identifies recurrent
 mutations, gene fusions and splicing alterations. Nat Genet. 2016;48:407-16.
- Bignell GR, Greenman CD, Davies H, Butler AP, Edkins S, Andrews JM, et al.
 Signatures of mutation and selection in the cancer genome. Nature. 2010;463:893-8.
- Rajaram M, Zhang J, Wang T, Li J, Kuscu C, Qi H, et al. Two Distinct Categories
 of Focal Deletions in Cancer Genomes. PloS one. 2013;8:e66264.
- 20 28. Cancer Genome Atlas Research N. Comprehensive genomic characterization of21 squamous cell lung cancers. Nature. 2012;489:519-25.
- 22 29. Wang K, Singh D, Zeng Z, Coleman SJ, Huang Y, Savich GL, et al. MapSplice:
 accurate mapping of RNA-seq reads for splice junction discovery. Nucleic acids
 research. 2010;38:e178.
- 25 30. Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data
 26 with or without a reference genome. BMC bioinformatics. 2011;12:323.
- de Reynies A, Jaurand MC, Renier A, Couchy G, Hysi I, Elarouci N, et al.
 Molecular classification of malignant pleural mesothelioma: identification of a poor
 prognosis subgroup linked to the epithelial-to-mesenchymal transition. Clin Cancer Res.
 2014;20:1323-34.
- 31 32. Bibikova M, Barnes B, Tsan C, Ho V, Klotzle B, Le JM, et al. High density DNA
 32 methylation array with single CpG site resolution. Genomics. 2011;98:288-95.
- 33. Cancer Genome Atlas Research N. Comprehensive molecular profiling of lung
 adenocarcinoma. Nature. 2014;511:543-50.
- 35 34. Wilkerson MD, Hayes DN. ConsensusClusterPlus: a class discovery tool with
 36 confidence assessments and item tracking. Bioinformatics. 2010;26:1572-3.

1 35. Reinius LE, Acevedo N, Joerink M, Pershagen G, Dahlen SE, Greco D, et al.

2 Differential DNA methylation in purified human blood cells: implications for cell lineage 3 and studies on disease susceptibility. PloS one. 2012;7:e41361.

36. Smyth GK. Limma: Linear Models for Microarray Data. In: Gentleman R, Carey,
V., Dudoit, S., Irizarry, R., Huber, W., editor. Bioinformatics and Computational Biology
Solutions using R and Bioconductor. New York: Springer; 2005. p. 397-420.

7 37. Cancer Genome Atlas Research N. Comprehensive molecular characterization
 8 of clear cell renal cell carcinoma. Nature. 2013;499:43-9.

9 38. Wilks C, Cline MS, Weiler E, Diehkans M, Craft B, Martin C, et al. The Cancer
10 Genomics Hub (CGHub): overcoming cancer through the power of torrential data.
11 Database : the journal of biological databases and curation. 2014;2014.

39. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR:
ultrafast universal RNA-seq aligner. Bioinformatics. 2013;29:15-21.

40. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, et al.
Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and
isoform switching during cell differentiation. Nature biotechnology. 2010;28:511-5.

41. Chu A, Robertson G, Brooks D, Mungall AJ, Birol I, Coope R, et al. Large-scale
profiling of microRNAs for The Cancer Genome Atlas. Nucleic acids research.
2016;44:e3.

42. Cancer Genome Atlas N. Comprehensive molecular portraits of human breast
 tumours. Nature. 2012;490:61-70.

43. Li J, Tibshirani R. Finding consistent patterns: a nonparametric approach for
 identifying differential expression in RNA-Seq data. Statistical methods in medical
 research. 2013;22:519-36.

44. Mullokandov G, Baccarini A, Ruzo A, Jayaprakash AD, Tung N, Israelow B, et al.
High-throughput assessment of microRNA activity and function using microRNA sensor
and decoy libraries. Nature methods. 2012;9:840-6.

45. Budczies J, Klauschen F, Sinn BV, Gyorffy B, Schmitt WD, Darb-Esfahani S, et
al. Cutoff Finder: a comprehensive and straightforward Web application enabling rapid
biomarker cutoff optimization. PloS one. 2012;7:e51862.

46. Altman DG, Lausen B, Sauerbrei W, Schumacher M. Dangers of using "optimal"
cutpoints in the evaluation of prognostic factors. J Natl Cancer Inst. 1994;86:829-35.

47. Shabalin AA. Matrix eQTL: ultra fast eQTL analysis via large matrix operations.
Bioinformatics. 2012;28:1353-8.

48. Quinn L, Finn SP, Cuffe S, Gray SG. Non-coding RNA repertoires in malignant
 pleural mesothelioma. Lung Cancer. 2015;90:417-26.

1 49. Ramirez-Salazar EG, Salinas-Silva LC, Vazquez-Manriquez ME, Gayosso-

2 Gomez LV, Negrete-Garcia MC, Ramirez-Rodriguez SL, et al. Analysis of microRNA

- 3 expression signatures in malignant pleural mesothelioma, pleural inflammation, and
- 4 atypical mesothelial hyperplasia reveals common predictive tumorigenesis-related

5 targets. Experimental and molecular pathology. 2014;97:375-85.

6 50. Reid G. MicroRNAs in mesothelioma: from tumour suppressors and biomarkers
 7 to therapeutic targets. Journal of thoracic disease. 2015;7:1031-40.

- 8 51. Truini A, Coco S, Alama A, Genova C, Sini C, Dal Bello MG, et al. Role of
- 9 microRNAs in malignant mesothelioma. Cellular and molecular life sciences : CMLS.
 2014;71:2865-78.

52. Hon CC, Ramilowski JA, Harshbarger J, Bertin N, Rackham OJ, Gough J, et al.
An atlas of human long non-coding RNAs with accurate 5' ends. Nature. 2017;543:199204.

14 53. Mele M, Mattioli K, Mallard W, Shechner DM, Gerhardinger C, Rinn JL.

15 Chromatin environment, transcriptional regulation, and splicing distinguish lincRNAs and 16 mRNAs. Genome Res. 2017;27:27-37.

17 54. Nguyen Q, Carninci P. Expression Specificity of Disease-Associated IncRNAs:
 18 Toward Personalized Medicine. Curr Top Microbiol Immunol. 2016;394:237-58.

19 55. Wright CM, Kirschner MB, Cheng YY, O'Byrne KJ, Gray SG, Schelch K, et al.
20 Long non coding RNAs (IncRNAs) are dysregulated in Malignant Pleural Mesothelioma
21 (MPM). PloS one. 2013;8:e70940.

56. Renganathan A, Kresoja-Rakic J, Echeverry N, Ziltener G, Vrugt B, Opitz I, et al.
GAS5 long non-coding RNA in malignant pleural mesothelioma. Mol Cancer.
2014;13:119.

57. Fassina A, Cappellesso R, Guzzardo V, Dalla Via L, Piccolo S, Ventura L, et al.
Epithelial-mesenchymal transition in malignant mesothelioma. Mod Pathol. 2012;25:8699.

28 58. Abba ML, Patil N, Leupold JH, Allgayer H. MicroRNA Regulation of Epithelial to
 29 Mesenchymal Transition. Journal of clinical medicine. 2016;5.

59. Diaz-Martin J, Diaz-Lopez A, Moreno-Bueno G, Castilla MA, Rosa-Rosa JM,
Cano A, et al. A core microRNA signature associated with inducers of the epithelial-tomesenchymal transition. The Journal of pathology. 2014;232:319-29.

Khanbabaei H, Teimoori A, Mohammadi M. The interplay between microRNAs
 and Twist1 transcription factor: a systematic review. Tumour biology : the journal of the
 International Society for Oncodevelopmental Biology and Medicine. 2016;37:7007-19.

36 61. Dhamija S, Diederichs S. From junk to master regulators of invasion: IncRNA

37 functions in migration, EMT and metastasis. International journal of cancer.

38 2016;139:269-80.

1 62. Ma N, Zhang W, Qiao C, Luo H, Zhang X, Liu D, et al. The Tumor Suppressive 2 Role of MiRNA-509-5p by Targeting FOXM1 in Non-Small Cell Lung Cancer. Cellular 3 physiology and biochemistry : international journal of experimental cellular physiology,

4 biochemistry, and pharmacology. 2016;38:1435-46.

63. Cicchini C, de Nonno V, Battistelli C, Cozzolino AM, De Santis Puzzonia M,
Ciafre SA, et al. Epigenetic control of EMT/MET dynamics: HNF4alpha impacts DNMT3s
through miRs-29. Biochimica et biophysica acta. 2015;1849:919-29.

8 64. Zhao J, Liu Y, Zhang W, Zhou Z, Wu J, Cui P, et al. Long non-coding RNA
9 Linc00152 is involved in cell cycle arrest, apoptosis, epithelial to mesenchymal transition,
10 cell migration and invasion in gastric cancer. Cell cycle. 2015;14:3112-23.

Andersen M, Grauslund M, Ravn J, Sorensen JB, Andersen CB, Santoni-Rugiu E.
Diagnostic potential of miR-126, miR-143, miR-145, and miR-652 in malignant pleural
mesothelioma. The Journal of molecular diagnostics : JMD. 2014;16:418-30.

Andersen M, Trapani D, Ravn J, Sorensen JB, Andersen CB, Grauslund M, et al.
Methylation-associated Silencing of microRNA-126 and its Host Gene EGFL7 in
Malignant Pleural Mesothelioma. Anticancer research. 2015;35:6223-9.

17 67. Williams M, Kirschner MB, Cheng YY, Hanh J, Weiss J, Mugridge N, et al. miR193a-3p is a potential tumor suppressor in malignant pleural mesothelioma. Oncotarget.
2015;6:23480-95.

Sui C, Meng F, Li Y, Jiang Y. miR-148b reverses cisplatin-resistance in non small cell cancer cells via negatively regulating DNA (cytosine-5)-methyltransferase
 1(DNMT1) expression. J Transl Med. 2015;13:132.

69. Chu J, Sadeghi S, Raymond A, Jackman SD, Nip KM, Mar R, et al. BioBloom
tools: fast, accurate and memory-efficient host species sequence screening using bloom
filters. Bioinformatics. 2014;30:3402-4.

70. Robertson G, Schein J, Chiu R, Corbett R, Field M, Jackman SD, et al. De novo
assembly and analysis of RNA-seq data. Nature methods. 2010;7:909-12.

28 71. Kent WJ. BLAT--the BLAST-like alignment tool. Genome research. 2002;12:656-29 64.

30 72. Kuhn RM, Haussler D, Kent WJ. The UCSC genome browser and associated
 31 tools. Briefings in bioinformatics. 2013;14:144-61.

Tibes R, Qiu Y, Lu Y, Hennessy B, Andreeff M, Mills GB, et al. Reverse phase
protein array: validation of a novel proteomic technology and utility for analysis of
primary leukemia specimens and hematopoietic stem cells. Molecular cancer
therapeutics. 2006;5:2512-21.

T4. Liang J, Shao SH, Xu ZX, Hennessy B, Ding Z, Larrea M, et al. The energy
sensing LKB1-AMPK pathway regulates p27(kip1) phosphorylation mediating the
decision to enter autophagy or apoptosis. Nature cell biology. 2007;9:218-24.

- 1 75. Hu J, He X, Baggerly KA, Coombes KR, Hennessy BT, Mills GB. Non-parametric 2 quantification of protein lysate arrays. Bioinformatics. 2007;23:1986-94.
- 3 76. Hennessy BT, Lu Y, Poradosu E, Yu Q, Yu S, Hall H, et al. Pharmacodynamic
 4 markers of perifosine efficacy. Clinical cancer research : an official journal of the
 5 American Association for Cancer Research. 2007;13:7421-31.
- 6 77. Coombes K, Neely, S., Joy, C. et al. SuperCurve Package. R Package Version
 7 1.4.1. 1.4.1 ed. Houston, Texas: MD Anderson Cancer Center; 2011.
- 8 78. Gonzalez-Angulo AM, Hennessy BT, Meric-Bernstam F, Sahin A, Liu W, Ju Z, et
 9 al. Functional proteomics can define prognosis and predict pathologic complete
 10 response in patients with breast cancer. Clinical proteomics. 2011;8:11.
- 79. Hennessy BT, Lu Y, Gonzalez-Angulo AM, Carey MS, Myhre S, Ju Z, et al. A
 Technical Assessment of the Utility of Reverse Phase Protein Arrays for the Study of the
 Functional Proteome in Non-microdissected Human Breast Cancers. Clinical proteomics.
 2010;6:129-51.
- 15 80. Akbani R, Ng PK, Werner HM, Shahmoradgoli M, Zhang F, Ju Z, et al. A pancancer proteomic perspective on The Cancer Genome Atlas. Nature communications.
 2014;5:3887.
- 18 81. Bajou K, Noel A, Gerard RD, Masson V, Brunner N, Holst-Hansen C, et al.
 19 Absence of host plasminogen activator inhibitor 1 prevents cancer invasion and
 20 vascularization. Nat Med. 1998;4:923-8.
- 82. Takayama Y, Hattori N, Hamada H, Masuda T, Omori K, Akita S, et al. Inhibition
 of PAI-1 Limits Tumor Angiogenesis Regardless of Angiogenic Stimuli in Malignant
 Pleural Mesothelioma. Cancer Res. 2016;76:3285-94.
- 83. Hayashi M, Madokoro H, Yamada K, Nishida H, Morimoto C, Sakamoto M, et al.
 A humanized anti-CD26 monoclonal antibody inhibits cell growth of malignant
 mesothelioma via retarded G2/M cell cycle transition. Cancer Cell Int. 2016;16:35.
- 27 84. Liu R, Proud CG. Eukaryotic elongation factor 2 kinase as a drug target in cancer,
 28 and in cardiovascular and neurodegenerative diseases. Acta Pharmacol Sin.
 29 2016;37:285-94.
- 85. Mak MP, Tong P, Diao L, Cardnell RJ, Gibbons DL, William WN, et al. A PatientDerived, Pan-Cancer EMT Signature Identifies Global Molecular Alterations and Immune
 Target Enrichment Following Epithelial-to-Mesenchymal Transition. Clin Cancer Res.
 2016;22:609-20.
- 86. Chen L, Gibbons DL, Goswami S, Cortez MA, Ahn YH, Byers LA, et al.
 Metastasis is regulated via microRNA-200/ZEB1 axis control of tumour cell PD-L1
 expression and intratumoral immunosuppression. Nature communications. 2014;5:5241.
- 37 87. Lines JL, Pantazi E, Mak J, Sempere LF, Wang L, O'Connell S, et al. VISTA is an
 38 immune checkpoint molecule for human T cells. Cancer research. 2014;74:1924-32.

88. Le Mercier I, Chen W, Lines JL, Day M, Li J, Sergent P, et al. VISTA Regulates
 the Development of Protective Antitumor Immunity. Cancer research. 2014;74:1933-44.

Build Strategy
Liu J, Yuan Y, Chen W, Putra J, Suriawinata AA, Schenk AD, et al. Immunecheckpoint proteins VISTA and PD-1 nonredundantly regulate murine T-cell responses.
Proc Natl Acad Sci U S A. 2015;112:6682-7.

6 90. Cancer Genome Atlas N. Comprehensive genomic characterization of head and 7 neck squamous cell carcinomas. Nature. 2015;517:576-82.

8 91. Shen R, Olshen AB, Ladanyi M. Integrative clustering of multiple genomic data
9 types using a joint latent variable model with application to breast and lung cancer
10 subtype analysis. Bioinformatics. 2009;25:2906-12.

Mo Q, Wang S, Seshan VE, Olshen AB, Schultz N, Sander C, et al. Pattern
discovery and cancer gene identification in integrated cancer genomic data. Proceedings
of the National Academy of Sciences of the United States of America. 2013;110:4245-50.

93. Pidsley R, CC YW, Volta M, Lunnon K, Mill J, Schalkwyk LC. A data-driven
approach to preprocessing Illumina 450K methylation array data. BMC Genomics.
2013;14:293.

17 94. Lopez-Rios F, Chuai S, Flores R, Shimizu S, Ohno T, Wakahara K, et al. Global
18 gene expression profiling of pleural mesotheliomas: overexpression of aurora kinases
19 and P16/CDKN2A deletion as prognostic factors and critical evaluation of microarray20 based prognostic prediction. Cancer Res. 2006;66:2970-9.

95. Vaske CJ, Benz SC, Sanborn JZ, Earl D, Szeto C, Zhu J, et al. Inference of
patient-specific pathway activities from multi-dimensional cancer genomics data using
PARADIGM. Bioinformatics. 2010;26:i237-45.

96. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al.
Cytoscape: a software environment for integrated models of biomolecular interaction networks. Genome Res. 2003;13:2498-504.

97. Wong CK, Vaske CJ, Ng S, Sanborn JZ, Benz SC, Haussler D, et al. The UCSC
Interaction Browser: multidimensional data views in pathway context. Nucleic Acids Res.
2013;41:W218-24.

- 31
- 32
- 33
- 34
- 35 36
- 37
- 38
- 39
- 40
- 41

1 Supplementary Figures and Figure Legends

2

3 Supplementary Figure S1: Pathology Review. A: Reclassification based on 4 pathology review. **B-G**: Histologic patterns of MPM: Epithelioid with tubules of 5 uniform tumor cells showing abundant eosinophilic cytoplasm (B); epithelioid 6 consisting of solid sheets of epithelioid cells (C); biphasic with bundles of spindle-7 shaped tumor cells surrounded by nests of epithelioid tumor cells (D); biphasic 8 with malignant spindled tumor cells (left) and epithelioid cells (right) (E); 9 sarcomatoid with cytologically malignant spindle cells showing hyperchromatic 10 nuclei (F); sarcomatoid with malignant spindle cells infiltrating among fat cells (G).

11

12 Supplementary Figure S2: Copy Number Analyses. A: Unsupervised 13 hierarchical clustering of DNA copy number profiles of Mesotheliomas into six 14 SCNA cluster groups. In the heatmap, SCNAs in tumors (x axis) are plotted by 15 genomic location (y axis). Red bars designate regions of amplification, while 16 blue bars designate deleted regions. **B:** The heatmap shows SCNAs in tumors 17 grouped by either those with biallelic inactivation of BAP1 or those that have two 18 active wild type BAP1 alleles. C: CDKN2A vs. MTPA mRNA expression. Co-19 deleted samples tend to exhibit lower levels of expression in both genes. D: 20 Number of inactivated BAP1 alleles in copy number clusters in A. For a X^2 test across all cluster groups, p = 0.028. E: Allele-specific copy number analysis of 21 22 TCGA-MQ-A6BR whole-exome sequencing data. Top panel shows total copy 23 log-ratio, middle panel shows allelic log-odds-ratio revealing allelic imbalances. 24 Red lines are joint segmentation using the FACETS algorithm. Integer copy 25 number estimates (total and minor) along with cellular fraction (cf-em) estimates 26 are shown in the bottom panel. F: Allele-specific copy number analysis of TCGA-27 SC-A6LP whole-exome sequencing data. G: Allele-specific copy number 28 analysis of TCGA-UD-AAC1 whole-exome sequencing data.

29

30 **Supplementary Figure S3: Gene Expression Analyses.** Distinct expression 31 subtypes of MPM associated with histology, clinical, and genomic features. **A**:

1 Unsupervised gene expression clustering defined four distinct mRNA-based 2 subtypes of MPM. A predictive gene list of 1,800 genes are displayed in the 3 heatmap with selected genes/pathways that are expressed in each group listed. 4 Annotation tracks are displayed for the four expression clusters (K1-K4) and 5 histology classification. K3 is enriched in non-epithelioid histologies. B: Selected 6 genomic features significantly associated with expression subtypes. BAP1 7 alterations annotation track includes single nucleotide variants (SNV), small and 8 large insertions and deletions (INDEL), and structural variants (SV), loss of 9 heterozygosity (LOH) events, and copy number (CN) loss events (blue). CN, 10 GISTIC scores – pink, gain; light blue, loss; dark blue, deep loss. C: Kaplan-11 Meier survival curves of probability of overall survival with overall significance 12 testing by log rank test and cox proportional hazards analysis using the good 13 outcome group K4 as the reference. **D**: Patient age of diagnosis by expression 14 subtype. Median and mean ages per subtype displayed below plot. Asterisk 15 denotes a pvalue of 0.0008 for a t-test comparing ages of K3 vs all other 16 subtypes. E: Normalized BAP1 expression levels across subtypes. BAP1 17 alterations are indicated by a red box and wildtype BAP1 by black circles. F and 18 **G**: Using gene sets from two previously published datasets: Bueno et al. (**F**) and 19 de Reynies et al. (G), we observe samples clustering similar to our classification. 20 **H**: Gene expression significantly associated with *BAP1* inactivation.

21

22 Supplementary Figure S4: DNA Methylation. A: Heatmap showing beta values 23 of 74 samples ordered by unsupervised clusters of DNA methylation data. 24 Samples are presented in columns and the CpG loci are presented in rows. An 25 annotation panel on the right of the heatmap indicates location of CpG locus 26 relative to CpG island status and gene. Annotation bars at the top of the heatmap 27 show genomic variables associated with the methylation clusters. Features 28 marked with (*) are significantly associated with DNA methylation clusters 29 (Fisher's Exact test p<0.01). **B**: Box plots of the leukocyte fraction and tumor 30 purity by DNA methylation clusters. C: Heatmap showing differentially methylated 31 probes between BAP1 altered and WT samples. An annotation panel on the right of the heatmap indicates location of CpG locus relative to CpG island status and
gene. An annotation bar at the top corresponds to the number of *BAP1*alterations in a sample.

4

5 Supplementary Figure S5: RPPA analysis and mRNA immune signatures. 6 A: Consensus matrix (k=5) for clustering 52 MPMs, based on 1-Pearson 7 correlation as distance metric and Wards' linkage algorithm. 80% samples used 8 in each of the 1000 iterations. **B**: Heatmap of consensus clusters from RPPA 9 Data. The columns and rows represent the 52 samples and 219 cancer-10 associated antibodies respectively. Five clusters identified with significant 11 differential protein expression and pathway activity. The annotation bars shown 12 above the heatmap were not used for clustering. C: Box plots of pathway scores 13 for each of the five RPPA clusters. P-values based on the ANOVA test. 14 Pathways differentially expressed between the clusters include Apoptosis, Cell 15 cycle, EMT, Hormone receptor, RAS/MAPK, TSC/mTOR, Breast reactive and 16 Core reactive pathways. D: Kaplan-Meier survival curves for the five RPPA 17 sample clusters. P-values based on the G-rho family of Harrington and Fleming 18 [13] tests to evaluate the difference between two or more survival curves. The 19 proteins differentially expressed between good prognosis cluster C5 and the poor 20 prognosis cluster C4 were PAI1, CYCLINB1, CAVEOLIN1, EPPK1 and EEF2K 21 as shown in panel B. E: GSEA-based mRNA immune signatures across iCluster 22 clusters.

23

24 Supplementary Figure S6: iCluster validation analyses. A: Correlation of 25 Bueno cohort (n=211, all histologies) to iCluster centroids. B: Kaplan-Meier plot 26 of survival analysis by Bueno iClusters across all histologies. C: Multivariate Cox 27 regression of OS by Bueno iClusters adjusted by histology and age. D: Silhouette 28 plot of Lopez-Rios cohort (n=52) iClusters. E: Breakdown of histology in Lopez-29 Rios iClusters. F: Kaplan-Meier plot of survival analysis by Lopez-Rios iClusters 30 (n=49). G: Multivaraiate Cox regression of OS by Lopez-Rios iClusters adjusted 31 by histology and age.

1

2 Supplementary Figure S7: Integrated mRNA and SCNA PARADIGM analysis.

3 A: Top heatmap shows molecular subtypes derived through unsupervised 4 analysis of PARADIGM results. Major pathway groups characterizing these 5 subtypes are annotated to the right of the heatmap. Middle and bottom heatmaps 6 show mRNA expression and SCNA within PARADIGM-derived subtypes. B: 7 Survival analysis of the PARADIGM-derived molecular subtypes. Kaplan–Meier 8 plot shows significant difference in survival of the best-surviving and the worst-9 surviving clusters. C: Differential activity pathway for BAP1 wild-type vs. 10 inactivated samples based on PARADIGM results. Red nodes indicate genes 11 that are active in BAP1 wild-type MPM. Blue nodes indicate genes with low 12 activity in the same samples. The size of the nodes is inversely proportional to 13 the differential p-value (larger nodes indicate higher significance). D: Top 14 differentially active pathways between samples with BAP1 inactivation vs. wild-15 type. PARADIGM IPLs for the genes in each pathway are aggregated per-patient 16 and ANOVA test is used to compute significance in distributions of aggregated 17 IPLs between the two groups (BAP1 wild-type and inactivated) of samples.

18

19

20

Initial Diagnosis



Final Diagnosis by our pathology review subgroup













Cox Regression with iCluster, Histology and Age (n=211)

Characteristic	HR (95% CI)	Р
iCluster group (ref iCluster 1)		
iCluster2	1.06 (0.70, 1.60)	
iCluster3	0.94 (0.62, 1.43)	0.93
iCluster4	0.94 (0.63, 1.38)	
Histology (ref Epithelioid)		
Non-Epithelioid	1.45 (1.06, 2.00)	0.023
Age (continuous)	1.03 (1.01, 1.04)	2.92e-4

Cox Regression with iCluster and Age restricted to epithelioid histology (n=141)

Characteristic	HR (95% CI)	Р
iCluster group (ref iCluster 1)		
iCluster2	2.36 (1.40, 4.00)	
iCluster3	2.75 (1.64, 4.61)	1.24e-4
iCluster4	2.43 (1.47, 4.04)	
Age (continuous)	1.03 (1.02, 1.05)	8.98e-5



Average silhouette width · 0.09

Cox Regression with iCluster, Age and Histology (n=49)

Characteristic	HR (95% CI)	Р
iCluster group (ref		
iCluster 1)		
iCluster2	0.74 (0.28, 1.95)	
iCluster3	2.38 (0.99, 5.69)	0.10
iCluster4	1.53 (0.67, 3.49)	
Age (continuous)	1.01 (0.98, 1.05)	0.39
Histology (ref Epithelioid)		
Non-Epithelioid	2.52 (1.14, 5.59)	0.03
	(· · · · · · · · · · · · · · · · · · ·	

Cox Regression with iCluster, Histology and Age (n=74)

Characteristic	HR (95% CI)	Р
iCluster group (ref iCluster 1)		
iCluster 4	5.71 (2.49-13.10)	
iCluster 3	1.49 (0.63-3.52)	4.6E-04
iCluster 2	2.50 (1.14-5.51)	
Histology (ref Epithelioid)		
Non-Epithelioid	2.15 (1.17-3.94)	0.02
Age (continuous)	0.98 (0.95-1.02)	0.37

Cox Regression with Paradigm Clusters, Histology and Age (n=74)

Characteristic	HR (95% CI)	Р
Paradigm group (ref cluster 1)		
Paradigm Cluster 4	23.63 (7.57-73.80)	
Paradigm Cluster 3	3.58 (1.48-8.68)	6.4E-07
Paradigm Cluster 2	2.76 (1.23-6.17)	
Histology (ref Epithelioid)		
Non-Epithelioid	1.60 (0.78-03.26)	0.20
Age (continuous)	0.99 (0.96-1.03)	0.70



Α