

1 **Proteomic approach to discover human cancer viruses from formalin-fixed tissues**

2

3 Tuna Toptan<sup>1,2</sup>, Pamela S. Cantrell<sup>3</sup>, Xuemei Zeng<sup>3</sup>, Yang Liu<sup>3</sup>, Mai Sun<sup>3</sup>, Nathan A. Yates<sup>3,4\*</sup>,

4 Yuan Chang<sup>1,\*</sup>, Patrick S. Moore<sup>1,\*</sup>

5

6 **Supplemental Tables and Figures**

7 **Supplemental Table 1.** Immunohistochemistry staining results of MCC cases used for this

8 study.

9

MCC Cases	CM2B4
R10-115	+++
R10-121	+++
R10-164	-
R11-43	-
R11-65	+++
R12-31	+++
R14-02	+++
R14-05	-
R14-07	-
R14-33	-
R14-34	-
R15-03	-
R15-11	+++
R15-12	-
R15-13	+++
R15-22	+++
R15-23	+++
R16-03	-
R16-41	+++
R16-42	-
R16-67	+++
R16-68	-
R16-69	+++

10

11 **Supplemental Table 2:** List of high-resolution MS1 features without a corresponding human  
 12 peptide sequence identification.

13

14 **Supplemental Table 3.** Degenerate and LNA-modified degenerate primers used for PCR  
 15 analysis in **Figure 2**

#ID	Peptide	BlastP	Reverse translation	Forward DP	Reverse DP	LNA-Reverse DP
1	AYEYGP <del>NPH</del> (..)NSR	MCV	GCNTAYGARTAYGGNCCNAA <del>YCC</del> NCAY	gaRtaYggNccNaaYcc	ggRttNggNccRtaYtc	g+gRt+iN+ggNccRtaYtc
3	LQPVKc <del>TGAR</del>	Human/ Chimpanzee	YTNCARCCNGTNAARTGYACNGGNGCNMGN	ccNgtNaaRtgYacNgg	ccNgtRcaYttNacNgg	c+cN+gt+RcaYttNacNgg
4	XXEXAPNCYGN <del>X</del> PXMK	MCV	GCNCCNAA <del>Y</del> TGYTAYGGNAAY	gcNccNaaYtgYtaYgg	ccRtaRcaRttNggNgc	cc+Rt+aR+caRttNggNgc
15	(..) <del>DEVDEAP</del> ( <del>IL</del> )YGT <del>TK</del>	MCV	GAYGARGTNGAYGARGC <del>NCCN</del>	gaYgaRgtNgaYgaRgc	gcYtcRtcNacYtcRtc	g+cY+tc+RtcNacYtcRtc

16  
 17 DP: degenerate primer. Peptide identification numbers (#ID) from **Table 1** are given in the first  
 18 column.

19

20 **Supplemental Table 4.** SMART-Degenerate and LNA-modified degenerate primers used for  
 21 NGS library generation.

#ID	Reverse DP	SMART Reverse DP	LNA-Reverse DP	SMART-LNA-Reverse DP
1	ggRttNggNccRtaYtc	aagcagtggtatcaacgcagagtlac	ggRttNggNccRtaYtc	aagcagtggtatcaacgcagagtlac
3	ccNgtRcaYttNacNgg	aagcagtggtatcaacgcagagtlac	ccNgtRcaYttNacNgg	aagcagtggtatcaacgcagagtlac
4	ccRtaRcaRttNggNgc	aagcagtggtatcaacgcagagtlac	ccRtaRcaRttNggNgc	aagcagtggtatcaacgcagagtlac
15	gcYtcRtcNacYtcRtc	aagcagtggtatcaacgcagagtlac	gcYtcRtcNacYtcRtc	aagcagtggtatcaacgcagagtlac

22  
 23 DP: degenerate primer. Peptide identification numbers (#ID) from Table 1 are given in the first  
 24 column. Positions for the LNA base modifications are indicated with (+). LNA : Locked nucleic  
 25 acid (LNA). SMART primer sequence is highlighted in blue.

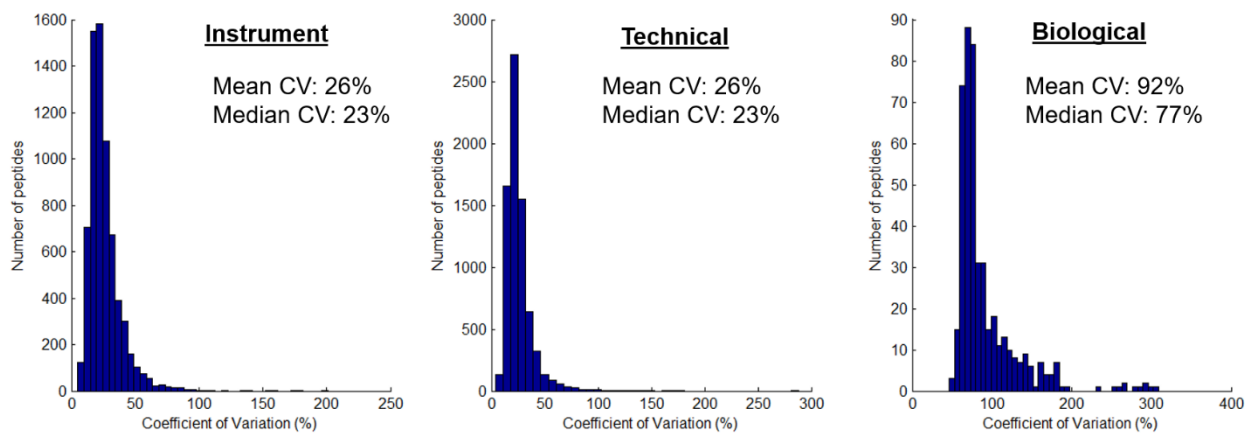
26

27 **Supplemental Table 5.** List of identified human peptides and associated quantification  
 28 information.

29 **Supplemental Table 6.** List of peptides from proteins with significant difference between MCV  
 30 (+) and control samples.

31 **Supplemental Figure 1**

32

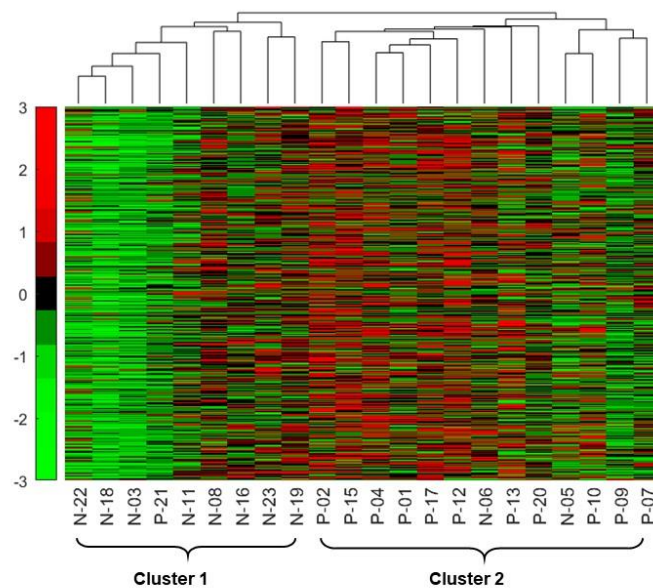


33

34 Histogram distribution of the coefficient of variation (CV%) for instrument replicates, sample  
35 preparation replicates (technical) and biological samples. For each sample type, only peptides  
36 with no missing values were included in the calculation.

37

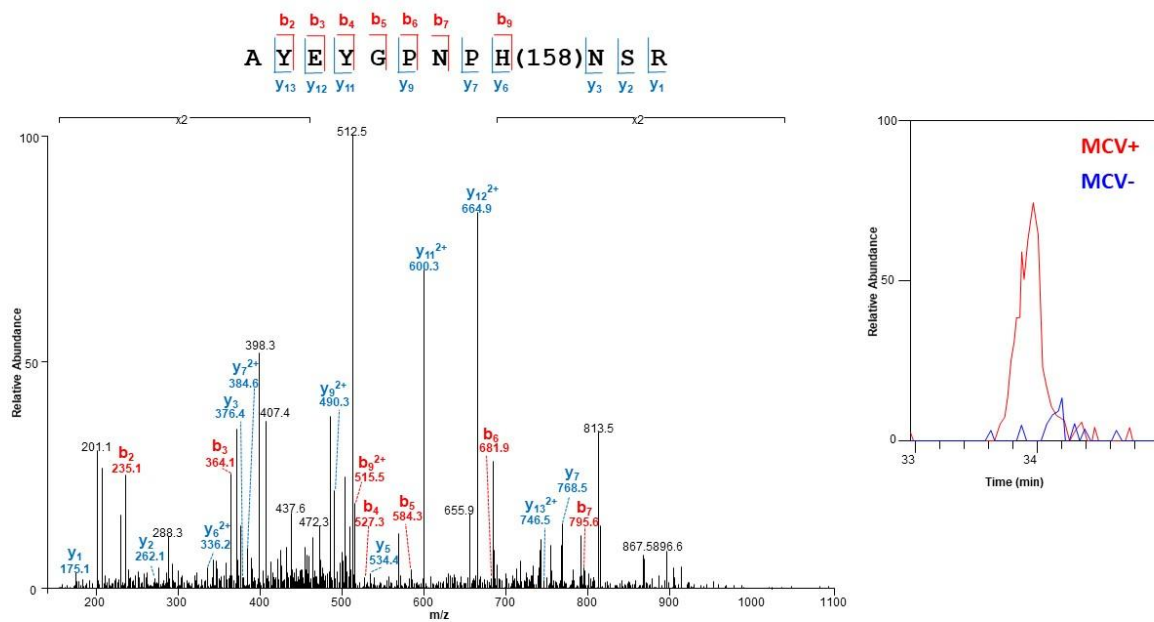
38 **Supplemental Figure 2**  
39



40  
41 Hierarchical clustering of samples based on the intensity values of peptides from viral related  
42 proteins. A total of 235 identified proteins are involved in viral related biological processes  
43 including viral transcription, viral process, viral translation termination and re-initiation, viral  
44 mRNA export from host cell nucleus, viral release from host cell, viral entry into host cells and  
45 virion assembly according to David Bioinformatics Resources 6.8, NIAID/NIH. Sample distance  
46 was based on Spearman's correlation, and the linkage was based on average linkage. The  
47 intensity values were standardized for each peptide with Z-transformation. Green and red colors  
48 in the clustering represent under- and over-expression, respectively. Sample #14 was excluded  
49 to its low overall intensities for all identified peptides. Peptides with more than 50% missing  
50 values were also excluded from clustering analysis. N: MCV negative, P: MCV positive samples.  
51

52 **Supplemental Figure 3**

53

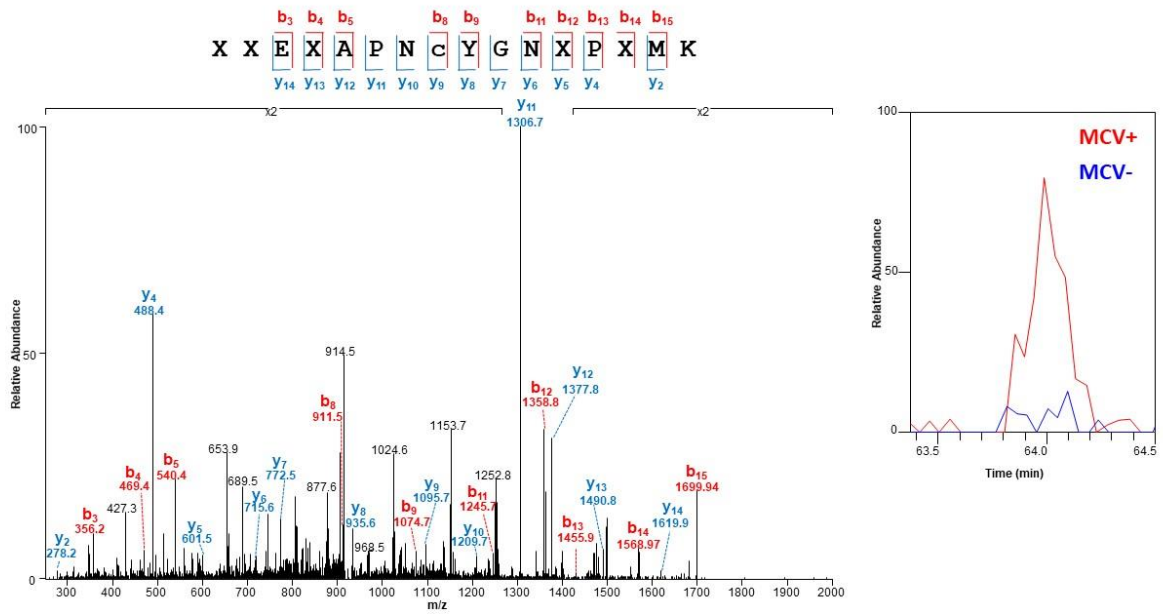


54

55 Tandem mass spectrum of proteomic feature #1 from Table 1. Spectrum was acquired by  
 56 targeted collision activated dissociation (CAD). The tryptic peptide ion m/z 521.6 (p-value 2.26E-  
 57 19 and fold change 4316.9) has a manual de novo amino acid sequence tag  
 58 AYEYGNPH(158)NSR. The amino acid sequence gap at 158 could be GT, TG, AS, or SA.  
 59 Representative MCV positive and negative shown as MS selected ion chromatogram.

60

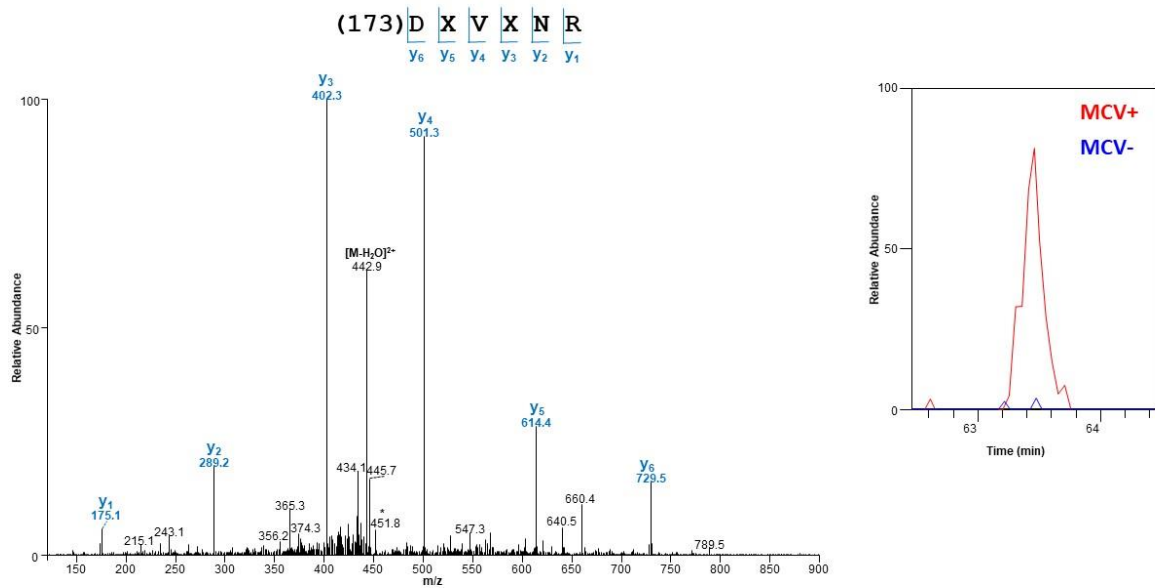
61 **Supplemental Figure 4**



62  
 63 Tandem mass spectrum of proteomic feature #4 from Table 1. Spectrum was acquired by  
 64 targeted collision activated dissociation (CAD). The tryptic peptide ion *m/z* 923.5 (p-value 5.11E-  
 65 07 and fold change 50.6) has a manual *de novo* amino acid sequence tag  
 66 **XXEXAPNcYGNXPXMK**. The amino acids “X” refer to isoleucine or leucine and “c” is a cysteine  
 67 residue with fixed modification carbamidomethylation (+57.02). The first two amino acids could  
 68 have been XX, EP, or EP; however, EP and PE amino acids were excluded on the basis of  
 69 mass accuracy. Representative MCV positive and negative shown as MS selected ion  
 70 chromatogram.  
 71

72 **Supplemental Figure 5**

73



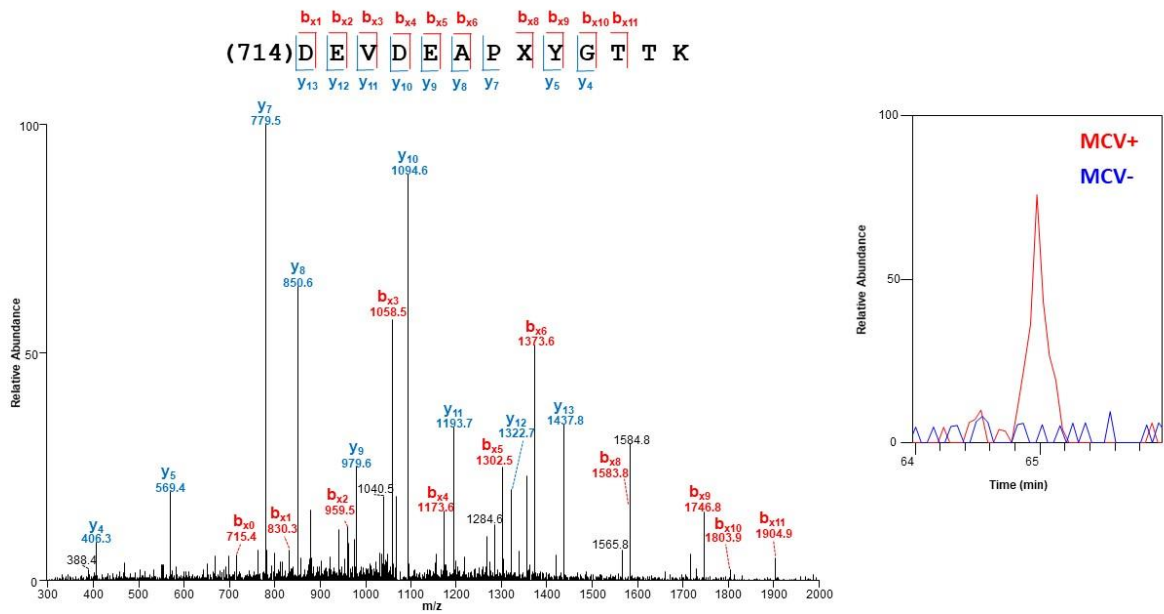
74

75

76 Tandem mass spectrum of proteomic feature #14 from Table 1. Spectrum was acquired by  
 77 targeted collision activated dissociation (CAD). The tryptic peptide ion  $m/z$  451.7 ( $p$ -value  $2.30E-$   
 78  $05$  and fold change  $2583.7$ ) has a manual *de novo* amino acid sequence tag (173)**DXVXNR**.  
 79 The amino acids “X” refer to isoleucine or leucine. Not enough information from the spectrum  
 80 could determine the N-terminal amino acids. Representative MCV positive and negative shown  
 81 as MS selected ion chromatogram.

82

83 **Supplemental Figure 6**



84

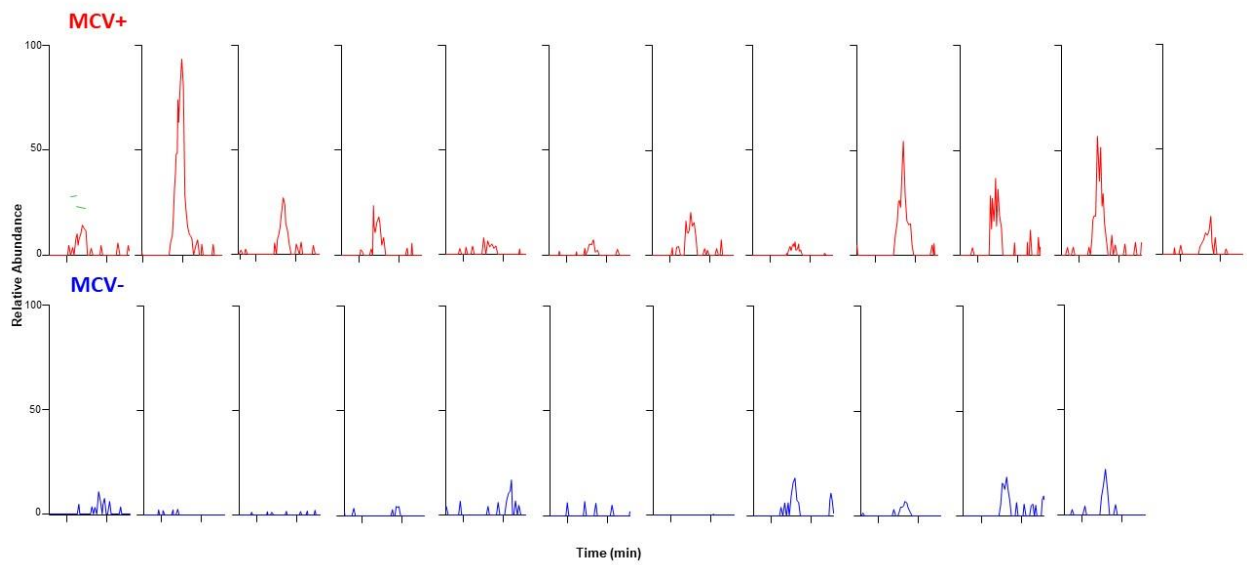
85 Tandem mass spectrum of proteomic feature #15 from Table 1. Spectrum was acquired by  
 86 targeted collision activated dissociation (CAD). The tryptic peptide ion  $m/z$  1076.5 ( $p$ -value  
 87  $2.36E-05$  and fold change 5405.8) has a manual *de novo* amino acid sequence tag  
 88 (714)**DEVDEAPXYGTTK**. The amino acid “X” refer to isoleucine or leucine. Not enough  
 89 information from the spectrum could determine the N-terminal amino acids. Representative  
 90 MCV positive and negative shown as MS selected ion chromatogram.

91



92 **Supplemental Figure 7**

93



94

95 The abundance of peptide **AYEYGNPH(158)NSR** (p-value 2.26E-19 and fold change 4316.9)  
96 in individual patient samples, shown as MS selected ion chromatogram from MCV positive and  
97 MCV negative groups. Each panel represents one individual FFPE sample.

98