

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

Tuna Toptan, ... , Yuan Chang, Patrick S. Moore

JCI Insight. 2020. <https://doi.org/10.1172/jci.insight.143003>.

Resource and Technical Advance In-Press Preview Virology

The challenge of discovering a completely new human tumor virus of unknown phylogeny or sequence depends on detecting viral molecules and differentiating them from host molecules in the virus-associated neoplasm. We developed differential peptide subtraction (DPS) using differential mass-spectrometry (dMS) followed by targeted analysis to facilitate this discovery. We validated this approach by analyzing Merkel cell carcinoma (MCC), an aggressive human neoplasm, in which ~80% of cases are caused by the human Merkel cell polyomavirus (MCV). Approximately 20% of MCC have a high mutational burden and are negative for MCV, but are microscopically indistinguishable from virus positive cases. Using 23 (12 MCV positive, 11 MCV negative) formalin-fixed MCC, DPS identified both viral and human biomarkers (MCV Large T antigen, CDKN2AIP, SERPINB5 and TRIM29) that discriminates MCV positive and negative MCC. Statistical analysis of 498,131 dMS features not matching the human proteome by DPS revealed 562 (0.11%) to be up-regulated in virus-infected samples. Remarkably, four (20%) of the top 20 candidate MS spectra originated from MCV T oncoprotein peptides and confirmed by reverse translation degenerate oligonucleotide sequencing. DPS is a robust proteomic approach to identify novel viral sequences in infectious tumors when nucleic acid-based methods are not feasible.

Find the latest version:

<https://jci.me/143003/pdf>



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

2

3 Tuna Toptan^{1,2}, Pamela S. Cantrell³, Xuemei Zeng³, Yang Liu³, Mai Sun³, Nathan A.

4 Yates^{3,4*}, Yuan Chang^{1,*}, Patrick S. Moore^{1,*}

5

6 ¹ Hillman Cancer Center, University of Pittsburgh, Pittsburgh, PA, 15232, USA

7 ² Institute of Medical Virology, University Hospital Frankfurt, Goethe University, Frankfurt

8 am Main, 60590, Germany

9 ³ Biomedical Mass Spectrometry Center, University of Pittsburgh, Pittsburgh, PA, 15261,

10 USA

11 ⁴ Department of Cell Biology, University of Pittsburgh, Pittsburgh, PA, 15261, USA

12 * Equal contribution

13

14 The authors have declared that no conflict of interest exists.

15 **ABSTRACT**

16

17 The challenge of discovering a completely new human tumor virus of unknown phylogeny

18 or sequence depends on detecting viral molecules and differentiating them from host

19 molecules in the virus-associated neoplasm. We developed differential peptide

20 subtraction (DPS) using differential mass-spectrometry (dMS) followed by targeted

21 analysis to facilitate this discovery. We validated this approach by analyzing Merkel cell

22 carcinoma (MCC), an aggressive human neoplasm, in which ~80% of cases are caused

23 by the human Merkel cell polyomavirus (MCV). Approximately 20% of MCC have a high

24 mutational burden and are negative for MCV, but are microscopically indistinguishable

25 from virus positive cases. Using 23 (12 MCV positive, 11 MCV negative) formalin-fixed

26 MCC, DPS identified both viral and human biomarkers (MCV Large T antigen,

27 CDKN2AIP, SERPINB5 and TRIM29) that discriminates MCV positive and negative MCC.

28 Statistical analysis of 498,131 dMS features not matching the human proteome by DPS

29 revealed 562 (0.11%) to be up-regulated in virus-infected samples. Remarkably, four

30 (20%) of the top 20 candidate MS spectra originated from MCV T oncoprotein peptides

31 and confirmed by reverse translation degenerate oligonucleotide sequencing. DPS is a

32 robust proteomic approach to identify novel viral sequences in infectious tumors when

33 nucleic acid-based methods are not feasible.

34

35 INTRODUCTION

36 Seven human viruses are responsible for approximately 15% of the tumor burden world-
37 wide. This phylogenetically heterogeneous group of viruses differ extensively in their
38 genome sizes, nucleic acid composition, and replication mechanisms (1). Likewise, the
39 discovery processes for each of these seven tumor viruses has varied and evolved closely
40 with technological advances, particularly in molecular biology. Epstein-Barr virus (EBV or
41 human herpesvirus (HHV4) (2)), a large double stranded DNA herpesvirus, was first
42 identified in 1964 based on classic microbiology detection practices in cell culture and
43 electron microscopy. Hepatitis B virus (HBV), unculturable at the time, was found by
44 serologic screening in 1965 (3). The discovery of human T-lymphotropic virus-1 (HTLV-1),
45 a retrovirus, was facilitated by reverse transcriptase assays in 1980 (4), and, although many
46 strains of human papillomaviruses (HPV) were already identified by 1983, cervical cancer-
47 associated HPV strains were only identified through strain-specific DNA Southern
48 hybridization studies (5). Hepatitis C virus (HCV), a flavivirus, was found by cDNA library
49 screening in 1989 (6).

50

51 Molecular subtractive techniques have been most recently used to determine the infectious
52 etiologies of Kaposi sarcoma (KS) and Merkel cell carcinoma (MCC). In 1993, fragments of
53 the Kaposi sarcoma herpesvirus (KSHV/HHV8) genome were cloned using
54 representational difference analysis (RDA), a DNA-based subtractive process which can
55 isolate foreign nucleic acids from the human genome (7, 8). In 2008, Merkel cell
56 polyomavirus (MCV) transcripts were found by digital transcriptome subtraction (DTS), an
57 *in silico* RNA subtractive process taking advantage of a timely expansion in sequencing
58 capabilities, databases and search engines (9, 10).

59

60 Virus-associated cancers are biological accidents, detrimental to both the host and the viral
61 pathogen (11). The cancer virus is generally not actively replicating (latency or pseudo-

62 latency) in cancerous cells, which would otherwise tend to kill the host cell. However, latent
63 viral transcript levels tend to be reduced relative to cellular or lytic viral transcript levels (12).
64 Latent viral proteins, on the other hand, can have exceptional stability (13), be expressed
65 by non-canonical translation (14, 15), and can circumvent cellular protein degradation
66 mechanisms (13). This is thought to be a viral strategy to reduce immunoproteasomal
67 peptide processing to escape host immune responses against latent viral proteins (16).
68 Based on these biologic features of tumor viruses, we pursued a protein-based detection
69 method that may be useful for tumors in which RNA is unavailable or in which viral transcript
70 levels are too low to be routinely detected. A protein-based virus discovery method using
71 cross-reactive antibodies to viral proteins has been described for polyomaviruses (pan-
72 polyomavirus immunohistochemistry test, P-PIT) (17, 18). Since P-PIT depends on
73 conserved epitope(s) within a class of known viruses, it cannot identify unique nucleic acid
74 or peptide sequences from a new agent.

75

76 We sought an unbiased approach for deep peptide sequencing to differentiate human from
77 foreign peptides belonging to novel viruses that can make use of archival pathology
78 samples. To achieve this, we developed a methodology called differential peptide
79 subtraction (DPS) using label-free differential mass-spectrometry (dMS) to quantify relative
80 peptide abundance between complex samples (19-21). The advantages of DPS are that it
81 is able to interrogate protein abundance, can identify novel peptides, makes use of formalin-
82 fixed, paraffin-embedded (FFPE) tissues, and, if no pathogen is found, can reveal unique
83 cellular protein biomarkers that may improve diagnosis and prognosis of a target disease.

84

85 RESULTS

86 Unbiased DPS was performed on polyomavirus-positive and -negative MCC. MCC is a
87 highly aggressive human skin cancer, 80% of which is etiologically associated with MCV
88 (9). Virus positive MCCs express viral small T (sT) and large T (LT) antigen oncoproteins
89 and have a low mutation burden (22). In this subset of MCCs, MCV genome clonally
90 integrates into the host chromosome and acquires mutations or deletions resulting in the
91 translation of C-terminally truncated LT proteins which vary in size from tumor to tumor. In
92 contrast, virus-negative MCC, although microscopically indistinguishable from virus-
93 positive MCC, carry high mutation burdens and driver somatic mutations that phenocopy
94 MCV infection (23). We used MCV negative (n=11) and positive (n=12) MCC FFPE tissues,
95 and processed them in a blinded fashion to determine if DPS can distinguish the presence
96 of a tumor virus *de novo* in human tissues without prior knowledge of the virus identity or
97 sequence. The dMS workflow that was used for sample processing and data analysis is
98 illustrated in the graphical abstract.

99

100 Polyomavirus status was initially determined by MCV LT antigen immunohistochemical
101 staining (**Supplemental Table 1**). Proteins were extracted from FFPE tissues and digested
102 using filter-aided sample preparation (FASP)-based tryptic digestion and analyzed by nano-
103 flow liquid chromatography tandem mass spectrometry (nLC-MS/MS) (**Figure 1A**). High-
104 resolution full-scan (MS1) mass spectra and low-resolution tandem (MS2) mass spectra
105 were recorded on a hybrid ORBITrap Velos mass spectrometer (**Figure 1A**). Four types of
106 experimental samples were included in the experimental design: eleven MCV negative
107 samples, twelve MCV positive samples, nine sample processing replicates, and one
108 instrument control sample (**Figure 1A, Supplemental Figure 1**). A data set of 498,131
109 high-resolution MS1 features (**Supplemental Table 2**) was extracted from the raw mass
110 spectral data using the MaxQuant (v1.6.0.1) proteomic software package (24, 25).
111 Subsequently, all MS1 features that could be identified by searching the MS2 spectra

112 against a human Uniprot protein database (downloaded in February 2013 with 87,662
113 entries) were removed. The log₂ transformed intensities of the unidentified proteomic
114 features were analyzed with a Student's t test to select features that exhibit significant
115 differences in relative abundance between MCV-positive and MCV-negative tumor
116 samples. Filtering for spectral features with a p-value < 0.01 and at least a 10-fold higher
117 intensity in MCV-positive samples compared to MCV-negative samples returned 562
118 features. Targeted nLC-MS/MS analysis was used to collect MS₂ spectra for the 20 most
119 significant features ranked by ascending p-value. Manual *de novo* sequencing identified
120 amino acid sequence tags greater than five amino acids long for 11 of 20 selected features
121 (**Table 1**). A Blast search against UniProtKB revealed that four of these peptides matched
122 to the MCV T antigen protein sequence (**Table 1, Figure 1B**). MCV small T (sT) and Large
123 T (LT) antigens are derived from differentially-spliced transcripts and share a 78-amino acid
124 (aa) N-terminus, nucleotide 196-429, Frame 1). LT splices after the first exon into a C-
125 terminal exon (738 aa, nucleotide 861-3080, Frame 3), whereas sT-transcript reads through
126 the initial splice site (**Figure 1B**). The localization of four peptides identified by targeted
127 nLC-MS/MS analysis are shown in Figure 1B. The two peptides on the left (green-ID#14
128 and orange-ID#4) are common to sT and LT, whereas the other two (purple-ID#15 and
129 blue-ID#1) correspond to LT. There is a partial match to sT (DE) for the third peptide (purple-
130 ID#15) which spans the splice junction between exon 1 and 2 of LT. The relative abundance
131 of the identified MCV peptides in the infected vs. control samples (**Figure 1C**) show that a
132 human tumor virus in tumor tissues can generate sufficient protein to be identified *de novo*
133 from tumor tissue. Thus, the comparison of proteomic profiles from infected and control
134 tissues allows identification of new proteins without *a priori* knowledge of the protein
135 sequence.

136

137 In the case of a novel virus, dMS-identified peptides will not have a match in the databases,
138 nevertheless this information can facilitate the recovery of the viral genome sequence. To

139 this end, we sought to trace the non-human dMS identified peptides back to their genetic
140 origins by next generation sequencing (NGS) with cDNA libraries generated using
141 degenerate oligos based on the identified peptide sequences (**Supplemental Table 3**). In
142 designing degenerate primers, we aimed to avoid primer sequences with 6-fold codon
143 nucleotide variants (L, S and R) and to maximize the number of 2 fold codon variants (D,
144 E, Y, N, K) thereby maintaining moderate binding-specificity while reducing oligo
145 degeneracy. In line with this, peptide areas with “X” residues (Table 1) representing either
146 an L or an I (3-fold degeneracy sites) codons which are indistinguishable by MS were
147 excluded. In addition to non-human matches (ID#: 1,4, and 15, **Table 1**), we included
148 peptide ID# 3 which was only a partial match to human in a Blast search. Based on the
149 silico reverse translation, forward and reverse primers were designed for a total of four
150 peptides (**Supplemental Table 3**). First, we tested the binding efficiency of these
151 degenerate primers (DP) by a low-cycle RT-PCR (**Figure 2A**). Four different sets of
152 combinations of forward and reverse DP, and cDNA template from a MCV positive MCC
153 sample were used for the PCR reaction (**Figure 2A**). Combinations of Forward (F)4,
154 Reverse (R)1 and F4, R15 primers resulted in 400 and 200 bp PCR products respectively,
155 which were confirmed to be derived from MCV by sequencing (**Figure 2B**).

156

157 For the library generation to perform NGS analysis, degenerate oligos were fused to
158 SMART adaptor sequence (SMART-deg, 25 nt) (**Supplemental Table 4**) and used for
159 cDNA synthesis as described previously (26) (**Figure 2C**). To demonstrate the principle and
160 efficiency of this procedure, mixtures of SMART-degenerate oligos or a modified oligo(dT)
161 SMART primer were used to facilitate reverse transcription from viral or viral and host RNA,
162 respectively. Due to high degeneracy of these primers we sought to increase their specificity
163 and designed another set of primers by addition of a number of locked-nucleic acid (LNA)
164 modifications for indicated bases (**Supplemental Table 4**). Using these new LNA modified
165 primers, we generated MCC^{deg} MCC^{LNA-deg}, and (MCC^{polyA}) SMART-cDNAs which were then

166 processed into three Nextera Flex libraries and subsequently sequenced using NextSeq
167 500 (**Figure 2D**). 58-68 million reads per sample were obtained which were processed and
168 mapped to a combined reference index from GRCh38 and MCV (JF813003) annotations
169 (**Table 2**). Normalization procedures to account for different sequencing depths amongst
170 the three libraries include conversion of data to transcripts-per-million (TPM) read-outs and
171 trimmed mean of M values (TMM). We detected 7.3 and 2.6 times more MCV reads in
172 degenerate oligo primed RNAseq samples compared to polyA-based sequencing reads. In
173 addition, reads from MCC^{deg} and MCC^{LNA-deg} largely mapped upstream of the DP binding
174 sites within the T antigen region. Hence, this strategy can facilitate the identification of a
175 viral genome sequences even in cases where the dMS-peptides do not match to previously
176 identified pathogens.

177

178 Label-free dMS method not only identified differentially expressed viral peptides within a
179 complex mixture, but also proteins which can serve as prognostic biomarkers. A total of
180 17,921 unique human peptides from 2,832 corresponding protein groups were quantified
181 and the peptide intensity values were log₂ transformed (**Supplemental Table 5**). A
182 Student's t test was used for statistical comparison between MCV-positive and -negative
183 peptide intensity values. Significant proteins were selected if more than half of the identified
184 peptides from a protein were significant (p value < 0.05) and single peptide identifications
185 were excluded from the analysis. A total of 38 proteins showed significantly increased
186 abundance whereas eight proteins were decreased in abundance in MCV-positive samples.
187 The list of identified peptides for these proteins are included in **Supplemental Table 6**.

188

189 To validate differentially expressed human peptides as potential biomarkers, five MCV
190 positive and four MCV negative MCC tissue cores together with control tissues were used
191 to generate a tissue microarray and were analyzed for the expression of CDKN2AIP,
192 SERPINB5 and TRIM29 by immunohistochemistry (**Figure 3**). Consistent with dMS results,

193 we found loss of TRIM29 and SERPINB5 expression, and higher levels of CDKN2AIP
194 expression in all MCV positive MCC cases (**Figure 3A**). These results suggest a role of
195 MCV T antigens in the regulation of SERPINB5 and TRIM29 expression (**Figure 3B**).
196

197 **DISCUSSION**

198 In this study, we provide a nanoLC-MS/MS-based protocol to compare tissues and identify
199 differentially expressed peptides and potential prognostic markers. This is a
200 peptide/proteome subtraction process that is analogous to the mRNA/transcriptome
201 subtraction (DTS) originally used to discover MCV (9). Importantly, the high DPS de novo
202 identification rate for MCV peptides in the context of the entire human tumor tissue
203 proteome shows that this approach is promising. We anticipate that it can supplement RNA-
204 based analyses of suspected infectious cancers, especially for tumors in which it is difficult
205 to obtain sufficient RNA for sequencing.

206

207 The top 20 unsupervised candidate MS feature sequences (after differential and database
208 subtraction that were present in MCV-positive but not MCV-negative samples) were
209 manually determined, a laborious process. These 20 peptide sequences were then aligned
210 to the human proteome by BLAST-P, which revealed 4 of these 20 peptides to be of MCV
211 origin. These four peptides map to the N-terminus of the MCV T antigen oncoprotein
212 complex, including peptides common to sT and LT and to the beginning of the second exon
213 in LT, that are common to the coding regions of the truncated LT proteins found in all the
214 MCV-positive MCC tumors (**Figure 1B**).

215

216 Modern virus discovery only requires a discovery of a single unique nucleotide sequence to
217 recover the entire viral sequence by gene walking. We show that starting from three unique
218 peptides, NGS of degenerate cDNA from MCV positive MCC tumor library recovers unique
219 viral nucleic acid sequences that can allow full viral characterization. Although this approach
220 proved to be more efficient than poly-A NGS, LNA modifications to the oligos used for cDNA
221 generation did not seem to improve the outcome. We anticipate that sequentially performing
222 these steps (first DPS on formalin-fixed tumor tissues followed by degenerate NGS of

223 candidate peptide coding sequences using well-accessioned tumor tissue RNA libraries) is
224 a viable strategy to find and characterize human tumor viruses, particularly in rare tumors.
225

226 DPS relies on comparison of a viral cancer proteome to a matched control non-viral tumor
227 proteome. Other known paired viral/non-viral tumors that could be similarly tested include
228 head-and-neck carcinoma, nasopharyngeal carcinoma, Burkitt lymphoma, and
229 hepatocellular carcinoma (1). In our study spectral features with a p-value < 0.01 and at
230 least a 10-fold higher intensity in MCV-positive tumors returned 562 features. Targeted nLC-
231 MS/MS analysis was used to collect MS2 spectra for the 20 most significant features ranked
232 by ascending p-value which enabled the identification of candidate viral peptides. Effective
233 ranking and prioritization is important because *de novo* sequencing remains a largely
234 tedious and slow manual process. The entire protocol from the processing of blinded
235 specimens to the unbiased identification of the viral protein consumes less than three weeks
236 of laboratory time.

237

238 For some of the cancer types, however, a well-defined control group might not be available.
239 In such cases, statistical analysis for hierarchical clustering of the samples might be useful.
240 To specifically address this potential problem, we used un-supervised hierarchical
241 clustering to investigate the possibility of using proteomic profiles to accurately classify MCV
242 samples into two groups, viral positive and viral negative. The best classification result was
243 obtained using the proteomic profiles for proteins associated with virus-related biological
244 processes. A total of 19 out of 22 samples were correctly classified (86% accuracy,
245 **Supplemental Figure 2**). MCV features remained significantly different between the two
246 cluster groups despite the drop in their significance ranking, supporting the potential of
247 applying un-supervised clustering for classification of samples with unknown viral status.

248

249 An alternative approach in the absence of matched negative control tissues is the
250 generation of a reference database comprising human MS/MS peptide features. This
251 reference database could then be used for DPS in silico subtraction of universal “human”
252 peptides from tumor MS/MS profiles, leaving candidate “non-human” peptide sequences.
253 Such a MS/MS database (the proteome equivalent of the nucleotide RefSeq database)
254 does not currently exist. Such a database would also be highly dependent on machine and
255 sample characteristics, as well as biological characteristics (e.g., single nucleotide
256 polymorphisms, post-transcriptional and post-translational modifications) that would make
257 universal comparisons difficult. Following subtraction, degenerate NGS using the candidate
258 non-human tandem MS spectra to design oligonucleotides could be used to search for viral
259 sequences. This strategy would not only circumvent the need for well-matched histological
260 tissue controls but also reduce the cost, time and manual labor needed in evaluation steps.
261 As with nucleotide DTS (10), an in silico DPS analysis may miss a viral pathogen if
262 commensal or endogenous virus peptide features are mistakenly assigned as “human” in
263 the comparison database.

264

265 Even when no new virus is found, DPS has utility for identifying human protein biomarkers.
266 We identified 38 human proteins significantly increased and eight proteins decreased in
267 MCV-positive vs. MCV-negative MCC samples, including SERPINB5, a reported tumor
268 suppressor also known as mammary serine protease inhibitor (MASPIN) (27), and TRIM29,
269 a ubiquitin E3 ligase that may act as a scaffold protein in the DNA damage response (28).
270 Loss of TRIM29 expression promotes invasion of skin squamous cell carcinoma cells by
271 altering distribution of keratins (29). Loss of expression of these two proteins might
272 contribute to a more aggressive disease course for MCV-negative compared to MCV-
273 positive MCC; however, larger cohort studies are needed to confirm these initial findings.
274 These and other differentially expressed proteins can be readily examined as potential

275 prognostic biomarkers for MCV-positive and MCV-negative MCC tumors or as biomarkers
276 to differentiate MCV-negative MCC from other small round cell neuroectodermal cancers.
277
278 DPS also offers advantages over RNAseq-only searches for cases where latency-
279 associated viral transcripts are significantly less abundant than cellular transcripts. At
280 present, DPS is more time-consuming than NGS and requires tissue-matched negative
281 control samples. Thus, it should be seen as an extension rather than a replacement for
282 RNAseq analysis in virus discovery. DPS, however, has a critical advantage in making use
283 of archival tumor FFPE tissues in which RNA is degraded. Development of a platform-
284 independent human MS/MS reference database may markedly expand the potential for
285 uncovering new human pathogens using DPS.

286 **METHODS**

287 **Cell line, tissues, tissue microarray generation and immunohistochemistry**

288 HEK293 cells (ATCC, Manassas, VA, USA) were maintained DMEM (10-013, Cellgro,
289 Manassas, VA, USA) supplemented with 10% FBS (Sigma-Aldrich, St. Louis, MO, USA).

290 MCV positive, MCV negative MCC tumors were obtained from cooperative human tissue
291 network (CHTN). Based on the MCV LT expression levels determined by CM2B4 staining,
292 11 MCV negative and 12 MCV positive tumors were selected for the differential Mass
293 Spectrometry (dMS) study. Among those cores from 5 MCV positive, 4 MCV negative tumor
294 FFPE blocks and a series of normal tissues (spleen, colon, brain, prostate, skin, adrenal
295 gland, kidney, lung, uterus and tonsil) were used to generate a tissue microarray at UPMC
296 Hillman Cancer Center Tissue and Research Pathology services.

297 Slides were deparaffinized in xylene and rehydrated in a series of ethanol solutions.
298 Endogenous peroxidase activity was blocked by incubation of the slides with 3% hydrogen
299 peroxide for 15 min. Epitope retrieval was performed using 1 mM EDTA buffer pH 8.0 at
300 125°C for 3 min, and 90°C for 15 sec in an antigen retrieval chamber (Decloaking chamber,
301 Biocare medical, Pacheco, CA, USA). After blocking (Protein block, serum free, Dako, Ely,
302 UK), monoclonal antibody CM2B4 generated by standard methods of immunizing mice with
303 KLH-derivatized SRSRKPSSNASRGA peptide from the MCV T antigen (22) (0.6 µg/ml
304 mAb, 1:1500,) and commercial antibodies CDKN2AIP (1:400, sc-81841, Santa Cruz, CA,
305 USA), MASPIN/SERPINB5 (1:400 sc-271694), ATDC/TRIM29 (1:400 sc-376125) were
306 diluted in (1 % BSA, 0.1 % gelatin, 0.5 % Triton-X, 0.05 % sodium azide in PBS pH 7.4)
307 were applied to each section overnight at 4°C in a humidified chamber. Following extensive
308 rinsing steps in TBS, sections were incubated with mouse Envision Polymer (Dako) for 30
309 min at room temperature, reacted with deaminobenzidine (DAB, Dako) and counterstained
310 with hematoxylin (Dako). Images were acquired using Olympus microscope AX70
311 (Olympus Co., Ltd., Tokyo, Japan). All other chemicals were purchased from Sigma-Aldrich.

312 **Sample selection, preparation for dMS**

313 A total of 23 formalin-fixed, paraffin-embedded (FFPE) Merkel Cell carcinoma (MCC) tissue
314 samples were selected on the basis of immunohistochemical staining that determined the
315 presence/absence of MCV. The samples were anonymized to blind analysts from the MCV
316 status until the proteomic sample preparation and mass spectrometric analysis were
317 complete. Samples were sectioned to a 10 μ m thickness using a microtome and stored on
318 standard microscope slides.

319 **Preparation of FFPE tissue for mass spectrometry**

320 Deparaffinization was achieved with two xylene washes (3 min each), rehydrated with serial
321 ethanol washes (100, 100, 95, and 70% for 1 minute each), and LC-MS grade water twice
322 for 3 minutes each. After deparaffinization, 100 μ l lysis buffer (300 mM Tris pH 8.0, 100 mM
323 DTT, 2% SDS) was added to each tissue sample, followed by 30 min sonication, 1 h
324 incubation at 95°C and 2 h incubation at 65°C. After centrifugation at 17,000 g for 10 min,
325 the supernatants containing the extracted proteins were transferred to new eppendorf tubes
326 and Pierce 660 nm Protein Assay kit with the IDCR packet (ThermoFisher Scientific,
327 Waltham, MA, USA) was used to determine the total protein content.

328 Sample aliquots containing 30 μ g of total protein were digested with trypsin using the Filter
329 Aided Sample Preparation (FASP) protocol (30). In brief, the protein samples were added
330 to YM30 Microcon microcentrifuge filters (Millipore, Darmstadt, Germany) and washed three
331 times with 200 μ l of urea buffer (100 mM Tris-HCl pH 8.0, 8 M urea), each with 15 min
332 centrifugation at 14,000 g. Alkylation was performed by incubating at room temperatures
333 for 20 minutes in 100 μ l of urea buffer with 20 mM iodoacetamide. Samples were then
334 washed 3 times with 100 μ l urea buffer and then 3 times with 100 μ l 50 mM ammonium
335 bicarbonate, each with 10 min centrifugation at 14,000 g. 1.2 μ g Sequencing Grade TPCK-
336 treated trypsin (Promega, Fitchburg, WI, USA) were then added to each sample for
337 overnight digestion in a humidified 37°C incubator. The resultant peptides were desalted
338 using C18 Supelco cartridges (Supelco, Bellefonte, PA, USA), SpeedVac dried and then

339 reconstituted in 30 μ l 0.1% formic acid for analysis. All other chemicals were purchased
340 from Sigma-Aldrich.

341 Quality control samples were used to evaluate variability introduced by proteomic sample
342 processing and mass spectrometric analysis. A set of 9 sample processing controls were
343 created by combining equal amounts of undigested protein from the 23 extracted FFPE
344 samples and processed alongside the experimental samples. A pool FFPE protein extract
345 was divided into nine aliquots and processed together with the experimental samples to
346 assess sample preparation performance. A pooled instrument control sample was
347 generated by combining equal volumes of all the digested samples and analyzed multiple
348 times to monitor the stability of the mass spectrometer system over time. All sample
349 identities were blinded to eliminate analyst bias and processed using a balanced block
350 design to reduce variability introduced during sample processing and nLC-MS/MS analysis.
351 The mean coefficient of variation (CV) for all quantified human peptides was used to
352 characterize the biological (CV~90%) and technical (CV~30%) variation in the individual
353 and replicate samples, respectively (**Supplemental Figure 1**).

354 **Mass spectrometry and data processing**

355 Complex mixtures of proteolytic peptides (0.2 μ g for each injection) were analyzed by nano-
356 flow liquid chromatography tandem mass spectrometry (nLC-MS/MS) with a nano Acquity
357 UHPLC (Waters Corporation, Milford, MA) interfaced to a hybrid Orbitrap Velos Pro mass
358 spectrometer (ThermoFisher Scientific). Peptide separation was carried out on a C18
359 PicoChip™ 25 cm column (New Objective, Woburn, MA) with a 66 min linear gradient of 2-
360 35 % solvent B (acetonitrile/0.1% formic acid) at a 300 nL/min flowrate. The mass
361 spectrometer was operated in positive ionization mode with an electrospray voltage of 1.9
362 kV and capillary temperature of 275°C. Ion sampling and accumulation was controlled with
363 automatic gain control (AGC) and maximum injection time settings of 1,000,000 and 500 ms
364 for full-scan high-resolution (MS1) mass spectra, and 5,000 and 100 ms for the low-
365 resolution ion trap tandem (MS2) mass spectra, respectively. Data-dependent acquisition

366 recorded a full-scan MS1 spectrum at a resolution setting of 60,000 followed by 13 MS2
367 spectra at normalized collision energy setting of 35 with dynamic exclusion enabled.
368 Separate nLC-MS/MS analyses that collect MS2 spectra on pre-defined precursor ions
369 were performed using an isolation width of 2 m/z units and a relative collision energy setting
370 of 35.

371 The raw mass spectrometry data were analyzed with MaxQuant software version 1.6.0.1
372 (24) that incorporates the Andromeda (25) protein identification search engine and label-
373 free quantification tools. MS2 spectra were searched against the UniProt human proteome
374 database (February 2013 release) using standard ORBltrap parameters and a reversed
375 decoy database strategy that limits false peptide identifications rates to 1% or less. Briefly,
376 a precursor mass tolerance setting of 20 and 4.5 ppm were used for the first and main
377 database search, respectively. A mass tolerance setting of 0.5 Da was used for the MS2
378 fragment ions. Search enzyme specificity was defined as trypsin with maximum of two
379 missed cleavages, fixed Cysteine carbamidomethylation, and variable methionine oxidation
380 and protein N-terminal acetylation modifications. A minimum peptide length setting of 7 was
381 used and the maximum number of modification per peptide was limited to 5. The “match
382 between runs” and “matched unidentified” settings were enabled to prompt quantification of
383 high-resolution MS1 features, regardless of the peptide sequence identification status.

384 Raw MS data files, together with MaxQuant quantification results, have been deposited to
385 the ProteomXchange consortium via the MassIVE partner repository with the data set
386 identifier PXD021520 and URL <http://proteomecentral.proteomexchange.org/cgi/GetDataset?ID= PXD021520>.

388 **RNA extraction and SMART-library generation**

389 Total RNA was isolated from MCV positive MCC tumor (R1165) using TRIzol (Ambion Inc,
390 Austin, TX, USA) and was treated with TURBO DNase (ThermoFisher Scientific). RNA
391 quality was examined by 2100 Bioanalyzer (Agilent Technologies, Santa Clara, CA, USA)
392 before (RIN value 5.3) and after ribosomal RNA depletion using RiboMinus Eukaryote kit

393 (ThermoFisher) according to the manufacturer's recommendations. Ribosome depleted
394 samples were subsequently used for MCC-SMART library preparation. Libraries were
395 prepared using the SMARTer PCR cDNA synthesis kit (Clontech, Mountain View, CA, USA)
396 according to the manufacturer's recommendations with the following modifications:
397 Switching Mechanism at 5' End of RNA Template (SMART) fusion primers were designed,
398 which have the SMART sequence (5'-AAGCAGTGGTATCAACGCAGAGTAC-3') added to
399 the 5' end of each dMS-identified MCV- or human-specific degenerate reverse primer listed
400 in Supplementary Table 3. dMS-SMART-degenerate primer mix or a modified oligo(dT)
401 primer (3' SMART CDS primer IIA) were used to prime first strand cDNA synthesis.
402 Reaction mixtures consisting of 3.5 µl of RNA (~300 ng), 1 µl of 24 µM SMART primer mix
403 (1.2 µM final concentration for each), or 1 µl of 12 µM 3'SMART CDS primer IIA (5'-
404 AAGCAGTGGTATCAACGCAGAGTACT₍₃₀₎N₋₁N-3', where N=A,C,G, or T and N₋₁=A,G, or
405 C) were heated at 72°C for 3 min, and then the temperature was lowered to 47°C (0.1°C/min
406 slope) for 2 min before the addition of 5.5 µl of master mix (2 µl of 5x first-strand buffer,
407 0.25 µl of 100 mM DTT, 1 µl of 10 mM dNTP mixture, 1 µl of 12 µM SMARTer IIA
408 oligonucleotide, 0.25 µl of 40 U/µl RNase inhibitor, and 1 µl of 100 U/µl SMARTScribe RT).
409 cDNA synthesis reaction mixtures of clinical specimens were incubated at 47°C for a total
410 of 90 min, terminated at 70°C for 10 min, and brought to 4°C before the addition 0.1 µl
411 RNase H (5 U/µl, New England Biolabs, Ipswich, MA, USA). Reaction mixtures were
412 incubated at 37°C for 20 min, subsequently kept at 4°C, and adjusted to 50 µl with water.
413 SMART cDNA was amplified by long-distance PCR on a thermocycler as follows using
414 Advantage II reagents (Clontech): 7.5 µl SMART cDNA, 7.5 µl 10x Advantage 2 PCR buffer,
415 1.5 µl 50xdNTP mix (10 mM), 1.5 µl 5' PCR primer IIA (12 µM), 1.5 µl 50xAdvantage 2
416 polymerase mix, and 55.5 µl water (total of 75 µl). Reaction mixtures were cycled as follows:
417 step 1, 95°C for 1 min; 35 cycles of step 2, 95°C for 15 s, 65°C for 30 s, and 68°C for 3 min
418 followed by hold at 4°C.

419 Amplified SMART cDNA was purified with AMP-Pure magnetic beads (Beckman Coulter
420 Genomics, Brea, CA, USA) using a ratio of 1.8x beads to sample according to the
421 manufacturer's recommendations. Libraries were eluted in 30 µl of 10 mM Tris-Cl (pH 7.5)
422 and then quantified on Agilent 2100 Bioanalyzer (Agilent) reagents.

423

424 **NGS library generation, sequencing and analysis.**

425 Nextera DNA Flex kit was used to generate libraries from SMART-cDNA templates following
426 the manufacturer's instructions and sequencing was carried out on a NextSeq500 platform
427 (Illumina Inc., San Diego, CA, USA) for 2x75 paired-end reads. Fastq files were imported
428 into CLC Genomics Workbench 20.0 software (Qiagen, Hilden, Germany), paired-end
429 reads 1 and 2 were merged, and duplicate reads were removed. Reads were filtered for Q-
430 scores above 30, trimmed for quality (limit = 0.05) and ambiguity (2-nt maximum), and the
431 illumina and SMART adaptor sequence were removed. Reads below 20 nt were discarded
432 and paired-end reads were aligned to combined reference index from GRCh38 (Hg38) and
433 MCV (JF813003) or to individual reference genomes. The following alignment settings were
434 applied: mismatch = 2, insertion = 3, deletion = 3, length fraction = 0.8, and similarity fraction
435 = 0.8. Sequencing data is deposited at NCBI GEO Platform accession number GSE157610.

436 **Statistics**

437 Feature selection was based on a combination of statistical significance and fold change
438 difference. A two-sided equal variance Student's t-test on the log₂ transformed intensities
439 was used to determine the significance of difference between MCV positive and negative
440 samples for all high-resolution MS1 features that consist of an "isotope group" without a
441 corresponding human peptide sequence identification. Zero peptide intensity values were
442 imputed with one-tenth of the global minimum of non-zero values to enable log₂
443 transformation and fold change calculation. Unidentified MS1 features with at least 10 fold
444 increase in MCV positive samples were ranked in an ascending order according to the
445 Student's t-test p-values. 20 unidentified MS1 features with the highest significance were

446 subject to targeted nLC-MS/MS analysis. The targeted MS2 spectra were interpreted by
447 manual de novo sequence analysis (31) (**Supplemental Figure 3-6**) and the identified
448 sequence was confirmed with synthesized peptide standards. A representative select ion
449 chromatogram depicts the relative abundance of Feature 1 peptide AYEYGPMPH(158)NSR
450 in individual MCV positive and negative patient samples (**Supplemental Figure 7**).

451 **Study approval**

452 Tissues were obtained from cooperative human tissue network (CHTN) and examined
453 under the University of Pittsburgh Institutional Review Board IRB 86-22: UPCI Tissue
454 Banking Protocol.

455

456

457 **AUTHOR CONTRIBUTIONS**

458 T.T., Y.C., P.S.M., and N.A.Y. designed the experiments. T.T., P.S.C., Y.L., and X.Z.
459 performed the experiments. T.T. performed RNAseq, IHC, PCR and data analysis. P.S.C.,
460 M.S., and X.Z. performed differential mass spectrometry and related data analysis. N.A.Y.,
461 Y.C. and P.S.M supervised the project. T.T., P.S.C., X.Z., N.A.Y., Y.C., and P.S.M. wrote
462 the paper. Y.C. and P.S.M. contributed equally to this work.

463

464 **ACKNOWLEDGEMENTS**

465 This project was supported by the National Institutes of Health (grant numbers R35
466 CA197463 to P. S. M. and CA170354 to Y. C.). P. S. M. and Y. C. are additionally supported
467 as American Cancer Society Research Professors and receive support from the Pittsburgh
468 Foundation (to P. S. M.) and the University of Pittsburgh Medical Center Foundation (to Y.
469 C.). T.T. was supported in part by University of Pittsburgh Skin SPORE Career
470 Enhancement Program Funding (NIH P50 CA121973-09) and Hillman Cancer Center Pilot
471 Project Grant for Cancer Proteomics. Proteomics analysis was performed by the Hillman
472 Cancer Center Proteomics Facility supported in part by award P30CA047904. This project
473 used the University of Pittsburgh Health Sciences Sequencing Core at UPMC Children's
474 Hospital of Pittsburgh for library generation and Illumina sequencing.

475

476 **REFERENCES**

- 477 1. Chang Y, Moore PS, and Weiss RA. Human oncogenic viruses: nature and
478 discovery. *Philos Trans R Soc Lond B Biol Sci.* 2017;372(1732).
- 479 2. Epstein MA, Achong BG, and Barr YM. VIRUS PARTICLES IN CULTURED
480 LYMPHOBLASTS FROM BURKITT'S LYMPHOMA. *Lancet (London,*
481 *England).* 1964;1(7335):702-3.
- 482 3. Blumberg BS, Alter HJ, and Visnich S. A "NEW" ANTIGEN IN LEUKEMIA
483 SERA. *Jama.* 1965;191:541-6.
- 484 4. Poiesz BJ, Ruscetti FW, Gazdar AF, Bunn PA, Minna JD, and Gallo RC.
485 Detection and isolation of type C retrovirus particles from fresh and cultured
486 lymphocytes of a patient with cutaneous T-cell lymphoma. *Proc Natl Acad*
487 *Sci U S A.* 1980;77(12):7415-9.
- 488 5. Dürst M, Gissmann L, Ikenberg H, and zur Hausen H. A papillomavirus
489 DNA from a cervical carcinoma and its prevalence in cancer biopsy
490 samples from different geographic regions. *Proc Natl Acad Sci U S A.*
491 1983;80(12):3812-5.
- 492 6. Choo QL, Kuo G, Weiner AJ, Overby LR, Bradley DW, and Houghton M.
493 Isolation of a cDNA clone derived from a blood-borne non-A, non-B viral
494 hepatitis genome. *Science.* 1989;244(4902):359-62.
- 495 7. Chang Y, Cesarman E, Pessin MS, Lee F, Culpepper J, Knowles DM, et al.
496 Identification of herpesvirus-like DNA sequences in AIDS-associated
497 Kaposi's sarcoma. *Science.* 1994;266(5192):1865-9.
- 498 8. Lisitsyn N, Lisitsyn N, and Wigler M. Cloning the differences between two
499 complex genomes. *Science.* 1993;259(5097):946-51.

- 500 9. Feng H, Shuda M, Chang Y, and Moore PS. Clonal integration of a
501 polyomavirus in human Merkel cell carcinoma. *Science*.
502 2008;319(5866):1096-100.
- 503 10. Feng H, Taylor JL, Benos PV, Newton R, Waddell K, Lucas SB, et al.
504 Human transcriptome subtraction by using short sequence tags to search
505 for tumor viruses in conjunctival carcinoma. *J Virol*. 2007;81(20):11332-40.
- 506 11. Moore PS, and Chang Y. Why do viruses cause cancer? Highlights of the
507 first century of human tumour virology. *Nature Reviews Cancer*.
508 2010;10(12):878-89.
- 509 12. Sarid R, Flore O, Bohenzky RA, Chang Y, and Moore PS. Transcription
510 mapping of the Kaposi's sarcoma-associated herpesvirus (human
511 herpesvirus 8) genome in a body cavity-based lymphoma cell line (BC-1). *J*
512 *Virol*. 1998;72(2):1005-12.
- 513 13. Kwun HJ, da Silva SR, Shah IM, Blake N, Moore PS, and Chang Y.
514 Kaposi's sarcoma-associated herpesvirus latency-associated nuclear
515 antigen 1 mimics Epstein-Barr virus EBNA1 immune evasion through
516 central repeat domain effects on protein processing. *J Virol*.
517 2007;81(15):8225-35.
- 518 14. Kwun HJ, Toptan T, Ramos da Silva S, Atkins JF, Moore PS, and Chang Y.
519 Human DNA tumor viruses generate alternative reading frame proteins
520 through repeat sequence recoding. *Proc Natl Acad Sci U S A*.
521 2014;111(41):E4342-9.
- 522 15. Toptan T, Fonseca L, Kwun HJ, Chang Y, and Moore PS. Complex
523 alternative cytoplasmic protein isoforms of the Kaposi's sarcoma-associated

524 herpesvirus latency-associated nuclear antigen 1 generated through
525 noncanonical translation initiation. *J Virol.* 2013;87(5):2744-55.

526 16. Tellam J, Sherritt M, Thomson S, Tellam R, Moss DJ, Burrows SR, et al.
527 Targeting of EBNA1 for rapid intracellular degradation overrides the
528 inhibitory effects of the Gly-Ala repeat domain and restores CD8+ T cell
529 recognition. *J Biol Chem.* 2001;276(36):33353-60.

530 17. Rigatti LH, Toptan T, Newsome JT, Moore PS, and Chang Y. Identification
531 and Characterization of Novel Rat Polyomavirus 2 in a Colony of X-SCID
532 Rats by P-PIT assay. *mSphere.* 2016;1(6).

533 18. Toptan T, Yousem SA, Ho J, Matsushima Y, Stabile LP, Fernandez-
534 Figueras MT, et al. Survey for human polyomaviruses in cancer. *JCI*
535 *Insight.* 2016;1(2).

536 19. Meng F, Wiener MC, Sachs JR, Burns C, Verma P, Paweletz CP, et al.
537 Quantitative analysis of complex peptide mixtures using FTMS and
538 differential mass spectrometry. *J Am Soc Mass Spectrom.* 2007;18(2):226-
539 33.

540 20. Wiener MC, Sachs JR, Deyanova EG, and Yates NA. Differential mass
541 spectrometry: a label-free LC-MS method for finding significant differences
542 in complex peptide and protein mixtures. *Anal Chem.* 2004;76(20):6085-96.

543 21. Zhao X, Deyanova EG, Lubbers LS, Zafian P, Li JJ, Liaw A, et al.
544 Differential mass spectrometry of rat plasma reveals proteins that are
545 responsive to 17beta-estradiol and a selective estrogen receptor modulator
546 PPT. *J Proteome Res.* 2008;7(10):4373-83.

547 22. Shuda M, Arora R, Kwun HJ, Feng H, Sarid R, Fernandez-Figueras MT, et
548 al. Human Merkel cell polyomavirus infection I. MCV T antigen expression

549 in Merkel cell carcinoma, lymphoid tissues and lymphoid tumors. *Int J*
550 *Cancer*. 2009;125(6):1243-9.

551 23. Harms PW, Harms KL, Moore PS, DeCaprio JA, Nghiem P, Wong MKK, et
552 al. The biology and treatment of Merkel cell carcinoma: current
553 understanding and research priorities. *Nat Rev Clin Oncol*.
554 2018;15(12):763-76.

555 24. Cox J, and Mann M. MaxQuant enables high peptide identification rates,
556 individualized p.p.b.-range mass accuracies and proteome-wide protein
557 quantification. *Nat Biotechnol*. 2008;26(12):1367-72.

558 25. Tyanova S, Temu T, and Cox J. The MaxQuant computational platform for
559 mass spectrometry-based shotgun proteomics. *Nat Protoc*.
560 2016;11(12):2301-19.

561 26. Berg MG, Yamaguchi J, Alessandri-Gradt E, Tell RW, Plantier JC, and
562 Brennan CA. A Pan-HIV Strategy for Complete Genome Sequencing. *J Clin*
563 *Microbiol*. 2016;54(4):868-82.

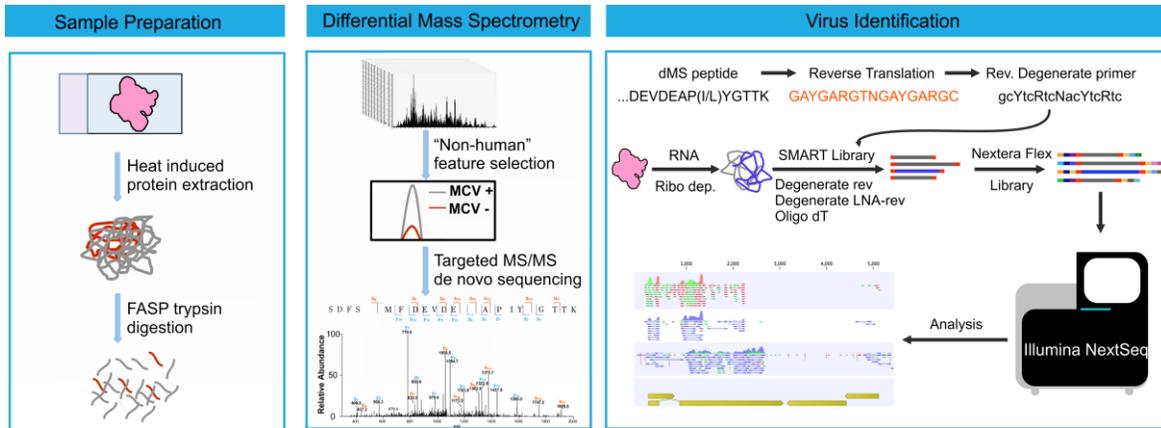
564 27. Berardi R, Morgese F, Onofri A, Mazzanti P, Pistelli M, Ballatore Z, et al.
565 Role of maspin in cancer. *Clin Transl Med*. 2013;2(1):8.

566 28. Masuda Y, Takahashi H, Sato S, Tomomori-Sato C, Saraf A, Washburn
567 MP, et al. TRIM29 regulates the assembly of DNA repair proteins into
568 damaged chromatin. *Nat Commun*. 2015;6:7299.

569 29. Yanagi T, Watanabe M, Hata H, Kitamura S, Imafuku K, Yanagi H, et al.
570 Loss of TRIM29 Alters Keratin Distribution to Promote Cell Invasion in
571 Squamous Cell Carcinoma. *Cancer Res*. 2018;78(24):6795-806.

- 572 30. Wiśniewski JR. Proteomic sample preparation from formalin fixed and
573 paraffin embedded tissue. *Journal of visualized experiments : JoVE*.
574 2013(79).
- 575 31. Kinter M, and Sherman N. *Protein Sequencing and Identification Using*
576 *Tandem Mass Spectrometry*. 2000:64-116.
- 577
- 578

579 **FIGURES and FIGURE LEGENDS**



580

581 **Graphical Abstract**

582 Schematic representation of sample preparation workflow, mass spectrometry analysis

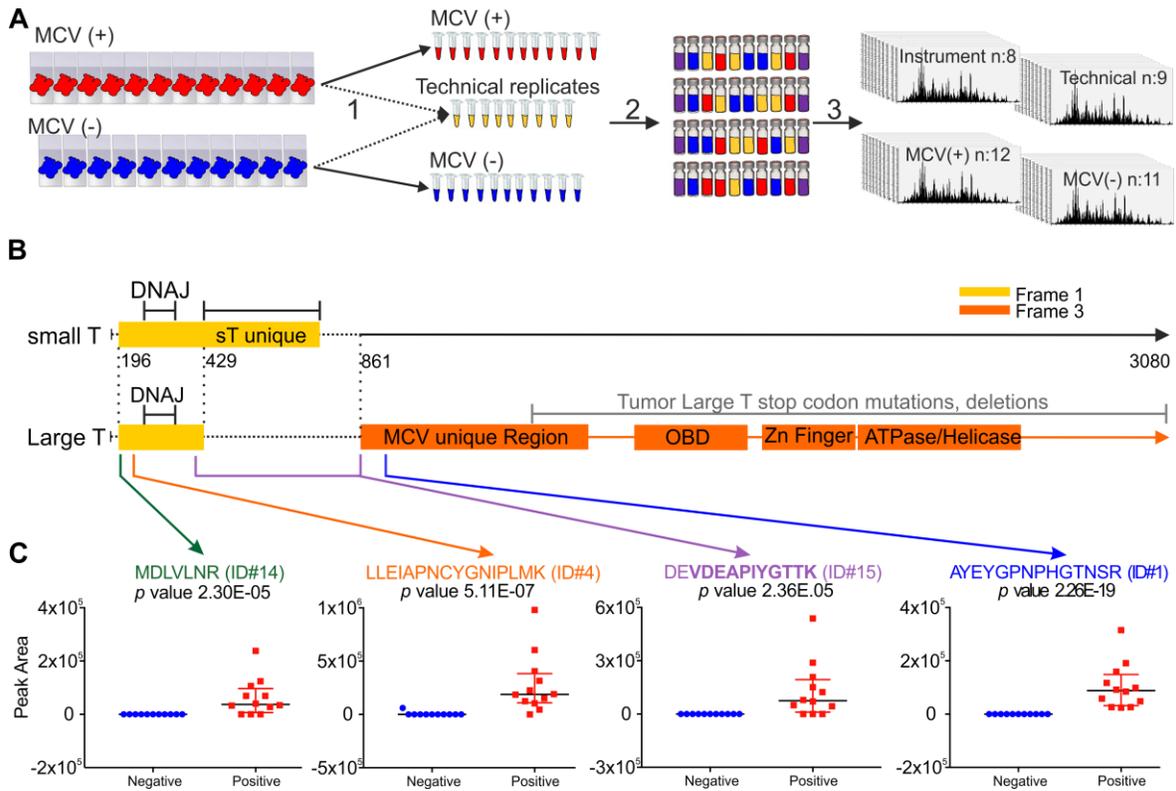
583 and viral DNA sequence identification from proteomic sequence tags.

584

585
586
587
588

Figure 1

DPS can detect *de novo* the presence of a tumor virus

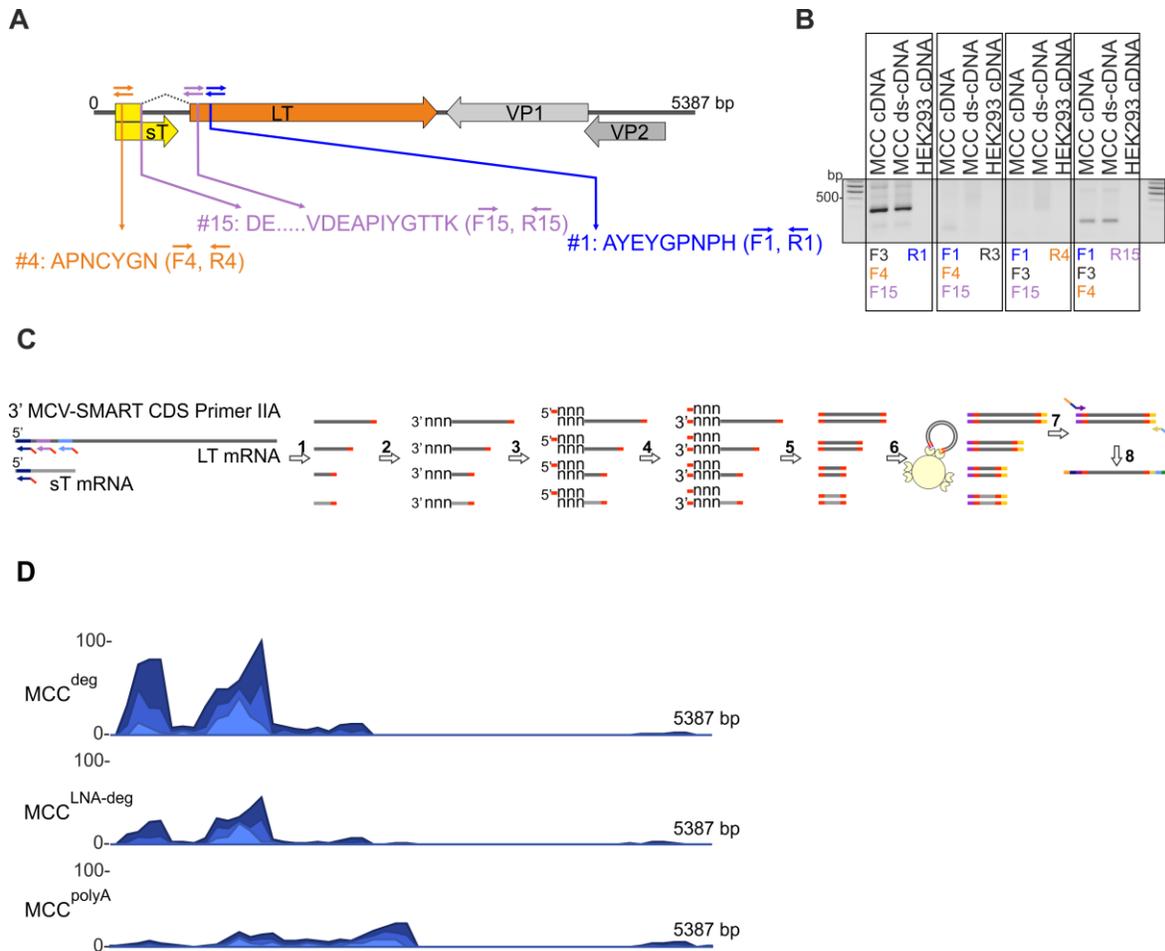


589

590 **A.** Workflow for dMS sample processing and instrumental analysis. 1. Deparaffinization,
591 antigen retrieval and lysis: 10 µL from each sample (n:23) was combined and aliquoted
592 into nine technical replicates. 2. FASP digestion step: Each sample was normalized to 30
593 µg. A total of 750 fmol of ovalbumin was added as an internal standard. A pooled
594 instrument control was made by combining 5 µL from sample (n:32). Samples (n:33) were
595 re-ordered. 3. nLC-MS/MS Analysis Step: Injection of ~0.2 µg on to C18 Picochip column
596 Orbitrap Velos Pro and analysis. **B.** Schematic illustration of MCV T antigen transcripts.
597 Small T (yellow-Frame 1) and Large T (yellow-Frame 1, orange-Frame 3) transcripts from
598 the early region including start, splice and termination sites are shown. Both small t and
599 Large T encode DNAJ domain. Small T and MCV unique domains, origin binding (OBD),
600 zinc finger, ATPase and helicase domains are depicted. The location of mutations and
601 deletions found in MCC tumor Large T are highlighted with a gray line. Positions of the
602 four MCV peptides identified by dMS analysis are indicated with green, orange, purple,
603 and blue arrows. **C.** Dot plots for the relative abundance of identified viral peptides in MCV
604 positive (red, n:12) vs. negative (blue, n:11) MCC samples. Peptides and their rankings
605 (Table 1) are shown in green (ID#14), orange (ID#4), purple (ID#15) and blue (ID#1).

606 Middle long lines indicates the mean and error bars represent standard deviation. P-
607 values were based on two-sided equal variance Student's t test.
608

609 **Figure 2**
 610
 611 **dMS-identified peptides facilitate identification of viral sequences by NGS with**
 612 **cDNA libraries generated using degenerate oligos**



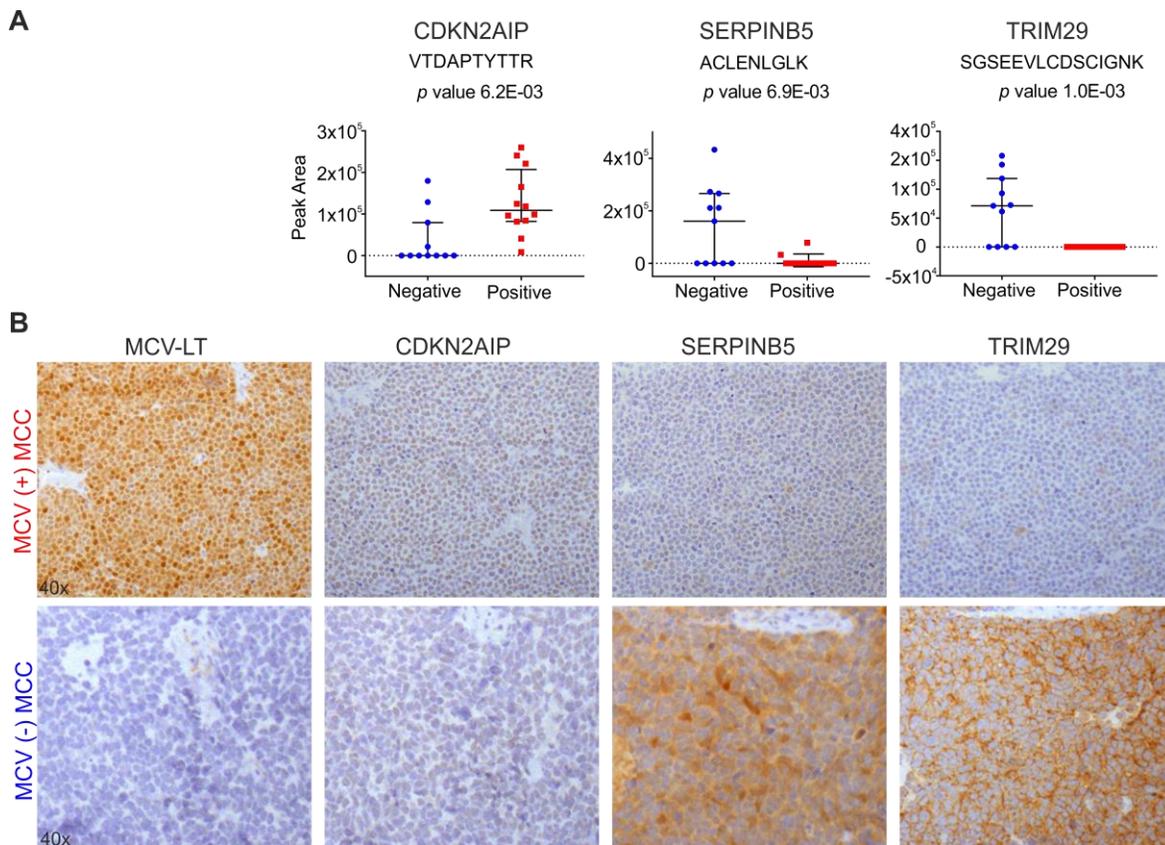
613
 614 **A.** Schematic illustration of the MCV genome. Early (Large T (LT) yellow-orange, small T (sT) (yellow)) and late (VP1 (light gray), VP2 (dark gray)) region open reading frames are
 615 shown. The corresponding positions of the three MCV peptides identified by dMS
 616 (Features #4, #15 and #1) and degenerate primer binding sites are shown in orange,
 617 purple and blue arrows respectively. **B.** RNA extracted from MCC tissue (R16-67) or
 618 HEK293 cells were subjected to cDNA synthesis with random hexamers and additionally
 619 second strand synthesis for the MCC sample (ds-cDNA). cDNAs were amplified using the
 620 indicated combinations of degenerate primers (**Supplemental Table 3**) corresponding to
 621 the peptide sites highlighted in light blue (F1, R1), violet (F15, R15) and orange (F4, R4).
 622 F: forward, R: reverse. F3 and R3 (black) are non-MCV primers. **C.** Library generation
 623 using SMART oligos and Nextera-flex: **1.** 3' SMART CDS Primer IIA (Supplementary
 624 Table 4) mediated 1st strand synthesis. **2.** Tailing by RT. In the cDNA reaction, non-
 625 templated bases (nnn) are added to the ends of nascent cDNA by the terminal transferase
 626

627 activity of reverse transcriptase (RT). **3.** SMARTer IIA oligo anneals to non-templated
628 bases at cDNA ends (nnn). **4.** Template switch and extension at 3' end. The RT
629 polymerase switches strands to transcribe the complement of the oligonucleotide, leaving
630 the SMART adaptor at both ends of cDNA. **5.** Long-distance PCR with single 5' PCR
631 Primer IIA amplifies libraries. **6.** Bead-linked transposomes mediate the simultaneous
632 fragmentation of ds-cDNA and the addition of Illumina sequencing primers using Nextera-
633 flex. **7.** Reduced-cycle PCR amplification amplifies sequencing ready DNA fragments and
634 adds indexes and adapters. **8.** Sequencing-ready fragments are washed and pooled. **D.**
635 **NGS coverage maps of MCC-RNAseq libraries.** RNAseq reads were obtained from
636 three different samples to compare the efficiency of MCV read recovery using various
637 primer pool sets for cDNA and library generation (**Supplementary Table 3**). Ribo-
638 depletetted MCC RNA (R11-65) was subjected to cDNA synthesis with SMART-degenerate
639 oligo pool (MCC^{deg}), LNA modified SMART-degenerate oligo pool-SMART (MCC^{LNA-deg}),
640 and modified oligo dT-SMART (MCC^{polyA}) and then subjected to library generation using
641 Nextera-flex application. Standardized coverage depths (reads) for comparison purposes
642 are indicated on the Y-axis for each alignment.
643

644
645
646
647

Figure 3

DPS can identify differentially expressed human peptides as potential biomarkers



648

649 **A.** Dot plots for the relative abundance of identified human peptides in MCV positive (red,
650 n: 12) vs. negative (blue, n:11) MCC tumor samples. Middle long lines indicates the mean
651 and error bars represent standard deviation. P values were based on two-sided equal
652 variance Student's t test. **B.** Immunohistochemistry staining of MCC TMA. R10-115 and
653 R15-03 are representative MCV positive (upper panel) and negative (lower panel) MCC
654 cases respectively. According to the IHC staining results we detected SERPINB5 and
655 TRIM29 in MCV negative cases and in none of the MCV positive cases as predicted by
656 dMS analysis. MCV LT expression was detected using CM2B4 is a monoclonal antibody.
657

658 TABLES

659 **Table 1.** List of the top 20 significant proteomic features.
660

#ID	m/z	Charge	Mass	Calibrated retention time (min)	p value	Fold change	Peptide	Gene Name	Organism
1	521.571	3	1561.691	34.2	2.26E-19	4316.9	AYEYGNP(TG, GT, AS, SA)SR	T antigen	MCV
2	390.542	3	1168.603	33.6	5.84E-09	3552.4	Unable to obtain amino acid sequence		N/A
3	565.309	1	564.302	35.2	1.99E-07	15.7	LQPVKcTGAR	PTTG1P	Human/ Chimpanzee
4	923.484	2	1844.953	64.1	5.11E-07	50.6	XXEXA(PN, NP)cYGNXPXMK	T antigen	MCV
5	521.829	2	1041.643	57.4	6.65E-07	2770.5	DLIVATIIVK	ATIC	Human
6	579.612	3	1735.814	47.4	6.70E-07	2873.2	Unable to obtain amino acid sequence		N/A
7	454.926	3	1361.756	55.1	4.78E-06	10.4	Unable to obtain amino acid sequence		N/A
8	801.390	3	2401.147	64.7	2.02E-05	3048.5	(TX, XT)QFVDWY(SW, WS)EK		N/A
9	458.572	3	1372.694	31.0	2.06E-05	2028.4	Unable to obtain amino acid sequence		N/A
10	696.990	3	2087.950	48.6	2.09E-05	7199.7	NPSTVEAFDLAQSNSEHSR	PFAS	Human
11	406.229	3	1215.665	50.2	2.12E-05	4496.2	mKFNKK	U65	Human
12	449.743	2	897.471	43.2	2.23E-05	21636.5	AVLYNYR	C3	Human
13	738.883	2	1475.751	62.8	2.28E-05	41.5	DIINEEEVQFLK	AARS1	Human
14	451.742	2	901.469	63.5	2.30E-05	2583.7	(173)DXVXNR	T antigen	MCV
15	1076.475	2	2150.935	65.0	2.36E-05	5405.8	(714)DEVDEAPXYGTTK	T antigen	MCV
16	750.355	2	1498.695	51.9	2.46E-05	5616.0	STTSTIESFAAQEK	LUC7L3	Human
17	674.691	3	2021.051	68.9	2.48E-05	1711.9	VLFPGNSTQYNILEGLEK	MAP1B	Human
18	517.227	2	1032.440	25.2	2.70E-05	24.1	Unable to obtain amino acid sequence		N/A
19	565.326	3	1692.956	65.6	5.11E-05	13.2	QSAEXXDXK		N/A
20	519.513	4	2074.024	31.7	2.11E-04	4433.1	Unable to obtain amino acid sequence		N/A

661

662 #ID, feature ID/Rank; m/z, monoisotopic mass to charge from MaxQuant (v1.6.0.1) output;

663 Charge, the charge-state of the precursor ion; Mass, the predicted monoisotopic mass of

664 the identified peptide sequence; Calibrated retention time (min), the recalibrated retention

665 time in minutes in the elution profile of the precursor ion from MaxQuant (v1.6.0.1) output;

666 p value, Student's t test p value after log2 transformation of peak area; Peptide, amino

667 acid sequence associated with selected feature. X=isoleucine or leucine; c= Cysteine

668 carbamidomethylated (+57.02); m= Methionine oxidation (+15.99).

669

670 **Table 2.** Summary of normalized read counts obtained from NGS analysis. Counts per
 671 million (CPM) reads following trimmed mean of M values (TMM) adjustment are calculated
 672 for each sample and gene. Transcripts per million (TPM) normalized values are indicated.
 673

Sample	Gene	CPM-TMM adjusted	TPM
MCC-deg	sT	2,37	214,38
	LT	13,51	237,62
	VP1	0	0
	VP2	0,16	11,04
MCC-LNA.deg	sT	0,14	15,51
	LT	5,38	113,08
	VP1	0	0
	VP2	0,14	11,98
MCC-polyA	sT	0	0
	LT	2,1	95,96
	VP1	0	0
	VP2	0,11	19,31

674
 675
 676
 677
 678