# The Transcriptome Architecture of Polyomaviruses

Jason Nomburg[1,2,3], Wei Zou[4], Thomas C. Frost[1,3], Chandreyee Datta[5,6,7,8], Shobha Vasudevan[5,6,7,8], Gabriel J. Starrett[9], Michael J. Imperiale[4,10], Matthew Meyerson[1,2,11]*, James A. DeCaprio[1,3,12]*

[1]Department of Medical Oncology, Dana-Farber Cancer Institute, Boston MA

[2]Broad Institute of MIT and Harvard, Cambridge, MA

[3]Harvard Program in Virology, Harvard University Graduate School of Arts and Sciences, Boston, MA

[4]Department of Microbiology and Immunology, University of Michigan, Ann Arbor, Michigan

[5]Massachusetts General Hospital Cancer Center, Harvard Medical School, 185 Cambridge St, CPZN4202, Boston, MA

[6]Department of Medicine, Massachusetts General Hospital and Harvard Medical School, Boston, MA

[7]Center for Regenerative Medicine, Massachusetts General Hospital, Harvard Medical School, Boston, MA

[8]Harvard Stem Cell Institute, Harvard University, Cambridge, MA

[9]Laboratory of Cellular Oncology, CCR, NCI, NIH, Bethesda, MD, USA

21    [10]Rogel Cancer Center, Ann Arbor, MI

22    [11]Department of Genetics, Harvard Medical School, Boston, MA

23    [12]Department of Medicine, Brigham and Women's Hospital, Harvard Medical School,

24    Boston, MA

25

26    *Correspondence to:

27    James A. DeCaprio

28    james_decaprio@dfci.harvard.edu

29

30    Matthew Meyerson

31    matthew_meyerson@dfci.harvard.edu

32

## Abstract

34   Polyomaviruses (PyV) are ubiquitous pathogens that can cause devastating human

35   diseases. Due to the small size of their genomes, PyV utilize complex patterns of RNA

36   splicing to maximize their coding capacity. Despite the importance of PyV to human

37   disease, their transcriptome architecture is poorly characterized. Here, we compare

38   short- and long-read RNA sequencing data from eight human and non-human PyV. We

39   provide a detailed transcriptome atlas for BK polyomavirus (BKPyV), an important

40   human pathogen, and the prototype PyV, simian virus 40 (SV40). We identify pervasive

41   wraparound transcription in PyV, wherein transcription runs through the polyA site and

42   circles the genome multiple times. Comparative analyses identify novel, conserved

43   transcripts that increase PyV coding capacity. One of these conserved transcripts

44   encodes superT, a T antigen containing two RB-binding LxCxE motifs. We find that

45   superT-encoding transcripts are abundant in PyV-associated human cancers. Together,

46   we show that comparative transcriptomic approaches can greatly expand known

47   transcript and coding capacity in one of the simplest and most well-studied viral families.

48

## Introduction

50   Polyomaviruses (PyV) are ubiquitous pathogens that can cause devastating human

51   diseases (Jiang et al., 2009a) including polyomavirus-associated nephropathy (PVAN),

52   hemorrhagic cystitis, and bladder cancer associated with BKPyV (Starrett et al., 2021),

53   progressive multifocal leukoencephalopathy caused by JCPyV, Merkel cell carcinoma

54   caused by Merkel cell polyomavirus (MCPyV), and dermatosis caused by human

55    polyomavirus 7 (HPyV7) (Jiang et al., 2009a; Nguyen et al., 2017). PyV have circular

56    double-stranded DNA genomes and express viral genes with distinct "early" and "late"

57    kinetics. Early and late transcripts are driven by a bi-directional central promoter, and

58    each terminate at their own polyA signal sequence located between the early and late

59    regions. The PyV early region encodes tumor or T antigens that promote cell cycle

60    progression and facilitate replication of the viral genome by host DNA polymerase. The

61    PyV late region originating from the common PyV promoter element on the opposite

62    genome strand encodes the structural proteins required for the generation of progeny

63    virions.

64

65    PyV transcripts undergo complex splicing to increase their coding capacity in the face of

66    their small ~5kb genomes. In addition to the major large T (LT) and small T (ST)

67    antigens, additional T antigen splice forms have been identified including transcripts that

68    generate truncated versions of LT (17kT, 57kT, truncT, and T' in SV40, MCPyV,

69    BKPyV, and JCPyV respectively), a "superT" antigen that contains a duplicated LxCxE

70    RB-binding motif in SV40, middle T (MT) in murine PyV (MPyV), and ALTO in MCPyV

71    (Abend et al., 2009; Carter et al., 2013; Freund et al., 1992; Kress et al., 1979; Shuda et

72    al., 2008; Smith et al., 1979; Trowbridge and Frisque, 1995; Zerrahn et al., 1993).

73    Although the diversity of late transcripts has been explored in SV40 (Good et al., 1988),

74    late transcript diversity in other PyV, including the major human pathogens, is poorly

75    characterized. To address this lack of knowledge of PyV transcription and to discover

76    unannotated biologically relevant PyV-encoded protein products, we used long- and

77    short-read RNA sequencing technologies to characterize the transcriptomes of eight

78    human and non-human PyVs.

79

# Results

**RNA sequencing expands PyV transcript diversity.**

82    To expand known PyV transcript diversity, we conducted a series of viral infections

83    followed by total or polyA short-read Illumina RNA sequencing (short-RNAseq (total)

84    and short-RNAseq (polyA) respectively) (**Figure 1A**). We integrated this newly

85    generated data with publicly available data from infected cell culture, human skin, and

86    other settings (**Table 1**). Viruses studied include SV40, BKPyV Dunlop variant and Dik

87    (wild type, or archetype), JCPyV, MPyV, MCPyV, HPyV7, and bark scorpion

88    polyomavirus 1 (BSPyV1).

89

90    For SV40 and the BKPyV Dunlop variant, which replicate robustly in cell culture, we

91    complemented short-read sequencing with Nanopore direct RNA sequencing

92    (dRNAseq) and PacBio Single-Molecule Real-Time sequencing (SMRTseq) (**Figure**

93    **1A**), two long-read sequencing approaches for polyA RNA with distinct library

94    preparations and sequencing strategies. Resultant RNAseq reads from long- and short-

95    read sequencing strategies were mapped against the viral reference genome and

96    grouped into transcript classes based on the presence of shared introns as detailed in

97    the Methods (**Figure 1B**). For SV40, viral transcripts represented 11.6% and 8.8% of

98    transcripts in dRNAseq and SMRTseq, respectively. For BKPyV Dunlop, viral transcripts

99    represented 28.6% and 27.8% of transcripts in dRNAseq and SMRTseq, respectively.

100   The total number of viral reads is detailed in **Figure S1A**. Transcripts within the same

101   class contain the same introns but may have distinct transcript start sites (TSSs) and

102   transcript end sites (TESs). For most transcriptomes, the majority of transcripts are

103   members of the first few transcript classes (**Figure S1B**). To filter out erroneous splice

104   sites, we required that all introns present in a dRNAseq or SMRTseq read must also be

105   supported by at least 5 splice junction-spanning reads in short-RNAseq (total) data.

106   Detailed information on this transcript class strategy is present in the Methods.

107

108   Comparison of read coverage from short-RNAseq (total), dRNAseq, and SMRTseq

109   revealed that dRNAseq and SMRTseq were relatively consistent with read coverage,

110   generally reflected expected patterns of exon usage (**Figure 1C**). In contrast, the read

111   coverage of short-RNAseq (total) was less representative of expected viral exon usage

112   and may reflect noise due to the amplification of smaller RNA fragments (**Figure 1C**).

113

114   For SV40 and BKPyV Dunlop, a transcript class (consisting of transcripts with shared

115   introns) was considered a bona-fide viral transcript if it was at least 0.1% of late or early

116   transcripts in dRNAseq or SMRTseq data as described in the Methods. For SV40,

117   which has detailed splice annotations (Good et al., 1988), we found that dRNAseq and

118   SMRTseq data are largely consistent with existing annotations. However, we identified

119   five previously unannotated SV40 transcripts that were supported by both long-read

6

120    sequencing approaches, plus one additional previously unannotated SV40 transcript

121    class supported by SMRTseq and short-RNAseq (total) (**Figure 1D, Figure 2).**

122

123    In contrast to SV40 and despite its clinical importance, BKPyV transcripts have been

124    poorly characterized. We identified a total of 23 transcripts, 21 of which are supported

125    by both dRNAseq and SMRTseq data and only six of which were previously identified

126    (Abend et al., 2009; Seif et al., 1979) (**Figure 1E, Figure 3**). While novel BKPyV late

127    transcripts are often analogous to the characterized wraparound and non-wraparound

128    transcripts previously identified in SV40, several additional and unexpected BKPyV

129    early transcripts were identified. For example, an atypically early TSS revealed a splice

130    donor that was used to generate transcript E3 (**Figure S4**). Early transcripts including

131    E6, E9, and E11 are conserved across numerous PyV and lead to formation of novel

132    ORFs - these are described in detail below.

133

134    We generated a comprehensive atlas of SV40 and BKPyV transcripts in **Figures S3 -**

135    **S8**. Watch plots display the structure of each identified transcript, and read pileups

136    show all transcripts identified in each transcript class. The relative abundance of each

137    transcript as well as exact splice coordinates and abundance information for each

138    identified transcript is provided in **Supplementary Tables 1 and 2**. Transcripts can also

139    be explored using an interactive Google Colab notebook

140    (https://colab.research.google.com/github/jnoms/SV40_transcriptome/blob/main/bin/cola

141    b/PyV_exploratory.ipynb). A comprehensive analysis of all splice sites detected in short-

7

142     read short-RNAseq (total) and short-RNAseq (polyA) in eight PyV studied is presented

143     in **Figure S9**.

144

145     To address the possibility that distinct transcript isoforms could be preferentially

146     translated, we performed polysome profiling of SV40-infected cells coupled with

147     dRNAseq of whole-cell and polysome-associated polyadenylated RNAs (**Figure 1F**).

148     The ribosome occupancy, determined as the ratio between a transcript's normalized

149     polysome abundance and its normalized whole-cell abundance, has a mean of slightly

150     above 1 for host transcripts (**Figure S2D**). We found 11.2% of reads in the whole-cell

151     fraction and 18.7% in the polysome fraction were viral, consistent with active translation

152     of viral transcripts. For late transcripts, the relative abundance in the whole-cell fraction

153     was tightly coupled to polysome relative abundance (**Figure 1G**), indicating limited

154     preferential translation of late transcripts. In contrast to late transcripts, we found that

155     the LT:ST ratio was 1.3:1 in the polysome fraction compared to a 3.4:1 ratio of LT:ST

156     transcripts in the whole-cell fraction, indicating preferential translation of ST during

157     infection.

158

159     **Wraparound transcription is conserved across diverse PyV.**

160     Long-read sequencing revealed the existence of many late transcripts that contain

161     multiple copies of a duplicated leader exon. Leader-leader splicing is due to

162     "wraparound transcription" of PyV transcripts that failed to terminate at the late

163     polyadenylation signal and continue to circle the genome repeatedly. PyV wrapround

164    transcription has been described previously although the structure and diversity of these

165    RNA species is unknown (Adami et al., 1989; Garren et al., 2015; Luo and Carmichael,

166    1991; Reddy et al., 1978). We investigated these transcripts in dRNAseq data from

167    SV40 and BKPyV. To supplement these data, we also performed dRNAseq on MPyV-

168    infected cells. Wraparound transcription, defined by the presence of repetitive copies of

169    a shared leader sequence, was found in long-read sequencing for all three PyVs

170    (**Figure 4A, B, C**: note the presence of the leader-leader or repeated exon near the "11

171    o'clock" position in watch plots). In addition to this leader sequence repetition, there are

172    diverse forms of wraparound transcripts that contain various combinations of

173    subsequent introns and encode for distinct viral proteins (**Figure 2, 3**). While only 3.6%

174    of SV40 transcripts originate from wraparound transcription, BKPyV and MPyV have

175    markedly higher rates at 25% and 41% respectively (**Figure 4D**).

176

177    Next, we inferred the presence of wraparound transcription in diverse PyV by identifying

178    short-RNAseq (total) reads that span the leader-leader junction (**Figure 4E**). Despite the

179    limited length of these short reads, leader-leader junctions can be accurately identified

180    within a single read through analysis of junction sites (**Figure 4F**). We found evidence of

181    wraparound transcription in all eight PyV investigated here. This includes HPyV7

182    RNAseq from infected human skin and RNAseq data from a scorpion containing the

183    highly divergent Bark scorpion polyomavirus 1 (BSPyV1), indicating that wraparound

184    transcription occurs *in vivo* and is widely conserved across PyV.

185

**186** **Pervasive premature polyadenylation of early transcripts in SV40, BKPyV, and**

**187** **MPyV.**

**188** We found that many early transcripts in SV40 and BKPyV underwent alternative

**189** polyadenylation (APA) earlier than the canonical polyA site as indicated by premature

**190** transcript end positions near 3 o'clock in the watch plots (**Figure S10A, B**). Early

**191** transcript APA had been previously identified in MPyV, where there is a canonical polyA

**192** signal sequence (AATAAA) within the LT ORF (Kamen et al., 1980b; Norbury and Fried,

**193** 1987). Indeed, dRNAseq identified APA of early transcripts in MPyV-infected cells

**194** (**Figure S10C, D**). In contrast to MPyV, APA in SV40 and BKPyV may be driven by

**195** alternative polyA signal sequences to the 5' of the APA site (ATTAAA in SV40,

**196** AAGAAA or TATAAA in BKPyV). Assessment of the cumulative incidence of early

**197** transcript termination shows abrupt increases in transcript termination ~1500nt

**198** upstream of the canonical polyA site in all three viruses (**Figure S10D**). This APA

**199** appears to be similarly abundant in LT and ST transcripts. We found that transcripts

**200** with APA still contain a full polyA tail that, while shorter than the polyA tails of transcripts

**201** that use the canonical polyA site, still tend to be longer than the polyA tails of host

**202** transcripts (**Figure S10E, S10F, S10G, S2C**). The polyA tail length of a spike-in control

**203** RNA with a known 30-adenine polyA tail was correctly estimated by dRNAseq (**Figure**

**204** **S2C**). We find that transcripts containing APA can associate with polysomes (**Figure**

**205** **S10H, I**), indicating that these transcripts are translated.

**206**

**207** **Comparative analysis of short-RNAseq (total) data reveals conserved,**

**208** **unannotated splice-forms that may generate variant ORFs.**

10

209    Next, we conducted a comparative analysis of PyV transcription from short-RNAseq

210    (total) data (**Figure S9**), with the hypothesis that data from diverse PyV could reveal

211    unannotated splice forms. This analysis led to the discovery of several unannotated but

212    conserved splicing events that have the potential to expand the coding capacity of PyV

213    (**Figure 5**).

214

215    We found that PyVs including HPyV7, MPyV, BKPyV Dunlop, and MCPyV express a

216    transcript utilizing the LT first exon donor but an acceptor within the ST ORF leading to

217    the generation of the ST2 ORF (**Figure 5A**). This splice occurs in-frame in HPyV7 and

218    BKPyV resulting in an internal deletion within ST, while in MPyV and MCPyV this splice

219    lands out of frame and results in the addition of novel C-terminal amino acids. The ST2

220    splice is highly abundant in HPyV7 representing over 20% of spliced early transcripts

221    from HPyV7-infected human skin. ST2-encoding transcripts were detected in BKPyV

222    dRNAseq and SMRTseq data (transcript E6).

223

224    MPyV encodes MT in addition to the LT and ST antigens common with other PyV.

225    MPyV MT is generated from a splicing event that connects the ST ORF with an ORF in

226    the alternative frame of the LT second exon. To our surprise, we found that BKPyV

227    expresses low levels of a similar transcript containing a splice that connects the ST

228    ORF with an MT-like ORF likewise in an alternative frame of the LT second exon

229    (**Figure 5B**). This MT transcript was also detected in BKPyV dRNAseq and SMRTseq

230    data (transcript E9).

11

231

232    JCPyV encodes two VP1 variants, VP1Xs, that consist of the N-terminal region of VP1

233    with novel C-termini that make up as much as 30% of late spliced transcripts in JCPyV

234    (**Figure 5C**) and have been recently identified and validated by an independent group

235    (Saribas et al., 2018). We found that VP1X-encoding transcripts were also produced by

236    MCPyV, SV40, BKPyV, and MPyV, albeit at a lower abundance than in JCPyV. Except

237    for one JCPyV VP1X-encoding splice, these transcripts were generated from splicing of

238    wraparound transcripts that run through the late polyA signal sequence.

239

240    **SuperT, a T antigen containing two RB-binding motifs, is present in multiple PyV**

241    **and in PyV-associated human cancers.**

242    Studies in SV40-transformed cells previously identified a superT antigen with higher

243    molecular weight than LT, containing a duplicated region with two copies of the LxCxE

244    RB-binding motif (Eul and Patzel, 2013). We found that a superT-specific splice was

245    present in SV40, BKPyV (Dik and Dunlop variants), JCPyV, and MCPyV during viral

246    infection (**Figure 5D**). The superT-specific splice originates from a splice donor

247    canonically associated with a conserved truncated LT antigen (17kT in SV40, truncT in

248    BKPyV, 57kT in MCPyV, and T' in JCPyV), but uses the LT second exon acceptor

249    available due to wraparound transcription. We find evidence of superT in the dRNAseq

250    and SMRTseq data for SV40 and BKPyV Dunlop infections (transcripts E4 and E11

251    respectively). Western blot with an antibody reactive to LT in BKPyV Dik-infected cells

252    revealed a band with slightly higher molecular weight than LT that is consistent with

12

253     superT (**Figure 5F**). BKPyV Dik mutant M1, designed to remove ST by replacing the LT

254     intron with an intron from the plasmid pCI (**Figure 5E**), also generated a superT band of

255     expected size. BKPyV Dik mutant M2 was generated by removing the LT intron and

256     adding the pCI intron just 5' of the LT first exon. Should the truncT donor be used to

257     generate superT in this mutant, the only available acceptor is before the LT 1st exon,

258     which would result in the formation of an aberrantly larger superT due to the inclusion of

259     a second copy of the LT first exon (**Figure 5E**). short-RNAseq (polyA) analysis of cells

260     infected with BKPyV Dik WT, M1, or M2 show junctions consistent with this model

261     (**Figure S11**), and western blot revealed that the superT band in M2 is shifted to a

262     higher molecular weight (**Figure 5F**). Together, these data indicate that superT is

263     generated by BKPyV Dik during viral infection.

264

265     SuperT was initially identified as an unexplained higher-molecular weight T antigen

266     present in many SV40-transformed cell lines (Kress et al., 1979; Smith et al., 1979).

267     While superT can be generated during viral infection because of wraparound

268     transcription, in SV40-transformed cells it would be possible to yield pre-mRNAs that

269     can be spliced to form superT should the virus be integrated in tandem copies (**Figure

270     6A**). Indeed, we previously observed that MCPyV integration events in Merkel cell

271     carcinoma (MCC) often lead to partial duplications of the viral genome and result in the

272     tandem insertion of multiple copies of viral early genes (Starrett et al., 2020).

273     Furthermore, the duplicated region in superT includes the RB-binding LxCxE motif,

274     raising the possibility that superT can function as a potent oncogene. We therefore

275     asked if there is evidence of superT in PyV-associated human cancers.

13

276

277    To address this question, we first analyzed short-RNAseq (total) data from five BKPyV-

278    associated bladder cancers (Starrett et al., 2021). To our surprise, we found that short-

279    RNAseq (total) data from two replicates of one BKPyV-associated bladder cancer

280    contained a higher abundance of superT-specific splice than even the LT- or ST-specific

281    splices, suggesting that a large fraction of "LT" in this tumor is superT (**Figure 6B**). We

282    next analyzed short-RNAseq (polyA) data from a series of 30 MCPyV-positive MCCs

283    and found evidence of superT in six cases (**Figure 6B**). Notably, the total number of

284    viral reads in some MCPyV-positive but superT-null tumors was very low, leaving open

285    the possibility that sequencing depth was insufficient to identify the superT splice in

286    additional tumors. Using PCR and sanger sequencing, we confirmed the presence of

287    the superT splice in MCC tumor J45_440 (**Figure S12A**).

288

289    We hypothesized that superT may be generated by cis-splicing due to concatemeric

290    integration of multiple copies of the etiologic PyV in these tumors (**Figure 6A**). To

291    address this hypothesis, we investigated three MCCs (J45_440, J17_296, J11_285) for

292    which we possess short-read whole genome sequencing data. From J11_285, we were

293    able to assemble the entire integration site, showing that MCPyV is integrated in a

294    manner that could allow cis-splicing to generate superT (**Figure S12B**). For J45_440,

295    we assembled a single viral block integrated in chromosome 7 (**Figure S12C**). We

296    found that 1) there are likely 2 copies of the viral genome, and 2) the 5' viral integration

297    site appears to fall on chromosome 7 "after" the 3' viral integration site, observations

298    consistent with the existence of two copies of the viral genome in tandem separated by

14

299    a small segment of host DNA at this integration site. For J17_296, from the assembly,

300    we could infer three distinct segments of viral DNA with integration sites closely spaced

301    within chromosome 2 (**Figure S12D**), indicating a complex integration pattern. The

302    longer block contains two copies of the early region and can likely support superT

303    generation through cis-splicing. The LT ORF of MCPyV is often truncated by premature

304    stop codons or deletions in MCC. We found that a stop codon in J17_296 likely

305    prevents expression of superT, but no stop codons occur before the superT splice in

306    J45_440 or J11_285 (**Figure S12E**). Together, these data indicate that viral integration

307    sites often could support cis-splicing to generate superT.

308    Two recent studies have found evidence of circular RNAs (circRNAs) that may be

309    generated by MCPyV in MCC and may support the translation of ALTO (Abere et al.,

310    2020; Yang et al., 2021). Of note, the major circRNA splice is equivalent to our

311    proposed MCPyV superT splice - a short-RNAseq read spanning the proposed circRNA

312    junction cannot be differentiated from a read spanning the superT junction. However,

313    the MCC RNAseq samples in which we found superT are short-RNAseq (polyA), which

314    should select against the potential circRNA due to its lack of a polyA tail. Furthermore,

315    we detect the superT splice in short-RNAseq (polyA) of SV40 and BKPyV Dunlop

316    infections in cell culture (**Figure S9**), although at around ~2/3 of its relative abundance

317    in short-RNAseq (total). Finally, we identify full-length superT transcripts in dRNAseq

318    data, which is highly unlikely to sequence circRNA since it is not polyadenylated. This

319    leaves open the possibility that some superT-like splice in short-RNAseq (total) from

320    viral infection originates from circRNA but suggests that most are from linear transcripts

321    that contain a polyA tail.

15

322

## Discussion

324    Here, we show that leveraging multiple long- and short-read RNA sequencing

325    approaches across 8 polyomaviruses has allowed us to greatly expand known transcript

326    diversity of this viral family. Short read RNAseq has limited capacity to characterize

327    transcriptome diversity because only a small fraction of reads span splice junctions, and

328    these junctions often cannot be phased with other junctions or to the transcript start and

329    end sites. Integrating long-read sequencing has allowed sequencing of entire

330    transcripts, including phasing of splice sites and transcript start and end positions.

331    Recent studies have leveraged long read sequencing to shed light on exceptional

332    complexity in the transcriptomes of diverse RNA and DNA viruses (Balázs et al., 2017;

333    Depledge et al., 2019; Garalde et al., 2018; Keller et al., 2018; Kim et al., 2020;

334    Nomburg et al., 2020; Price et al., 2020). We have expanded these studies to show that

335    a comparative approach within a viral family can identify conserved transcripts that

336    extend viral coding capacity.

337

338    Historically, studies of PyV transcripts were limited by the sensitivity and resolution of

339    northern blots, or by the read length of short read sequencing. Despite these limitations,

340    studies in the 70's and 80's were able to cumulatively characterize several SV40 late

341    transcripts, including one containing leader-leader splicing (Ghosh et al., 1978; Good et

342    al., 1988; Reddy et al., 1978). In contrast to SV40, the architecture of BKPyV late

343    transcripts is poorly characterized - prior to this work, the two major classes of late

344     transcripts ("16S" and "19S", reflecting transcript size based on gradient sedimentation

345     properties) were the primary late transcript classifications (Seif et al., 1979). Only

346     recently did a study provide some evidence for leader-leader splicing in BKPyV (Zou et

347     al., 2020). While the late transcripts of most PyV are thought to encode the canonical

348     late viral proteins, a recent study in JCPyV identified two splice events that lead to the

349     generation of novel proteins containing the N-terminal region of VP1 - one of which was

350     validated through western blot (Saribas et al., 2018). We found that these transcripts

351     (deemed "VP1X") are highly expressed in JCPyV but are also expressed at lower level

352     in BKPyV, MPyV, MCPyV, and SV40.

353

354     Leader-leader splicing is known to be highly prevalent in MPyV, where as many as 12

355     leader exons have been observed on a single RNA (Kamen et al., 1980a; Legon et al.,

356     1979; Treisman, 1980) - in our data, we have identified over 15 leader exons in a single

357     transcript. Furthermore, leader-leader splicing is required for stable accumulation of

358     MPyV late transcripts, dependent on length but not nucleotide composition of the leader

359     (Adami et al., 1989). Despite these observations, the exact structure, diversity, and

360     conservation of wraparound transcripts was not understood. Here, we found that leader-

361     leader splicing and wraparound transcription occurs in all PyV studied, including in the

362     divergent Bark scorpion polyomavirus 1, and found that the prevalence of leader-leader

363     splicing varies significantly between PyV. It is possible that this variation reflects

364     differences in the strength of the late polyA signals of these PyV. We found a large

365     diversity of wraparound transcripts containing variable numbers of the leader sequence

366     and diverse patterns of subsequent exon usage.

17

367

368    While late and early transcripts are thought to primarily end at the canonical late or early

369    polyadenylation sites, studies previously observed APA of early transcripts in MPyV

370    (Kamen et al., 1980b; Norbury and Fried, 1987). Here, we likewise identify pervasive

371    APA of early SV40 and BKPyV Dunlop transcripts and find that SV40 early transcripts

372    with APA can associate with polysomes and are likely translated. In addition, polysome

373    profiling revealed that SV40 transcripts are higher abundance in polysome-associated

374    RNAs than in whole-cell RNA populations, indicating preferential translation of SV40

375    transcripts. The relative abundance of individual late viral transcripts in the polysome

376    closely reflected their whole-cell abundance - conversely, ST transcripts were

377    preferentially translated compared to LT transcripts. The mechanism driving this

378    difference needs further study, as these transcripts differ only by a minor difference in

379    splice donor usage.

380

381    In addition to the major early transcripts encoding LT and ST, other early transcripts

382    have been identified in some PyV. MPyV encodes MT, generated by a splice

383    connecting the ST ORF and an ORF overprinted with LT second exon. MT is a primary

384    oncogene in MPyV and was thought to be largely restricted to rodent PyVs (Gottlieb and

385    Villarreal, 2001). We found that BKPyV generates a MT-like ORF through splicing

386    connecting ST and an ORF similarly overprinted with the LT second exon, showing that

387    non-rodent PyVs may be capable of expressing MT-like ORFs. In addition, MPyV also

388    encodes a tinyT antigen consisting largely of the LT first exon, resulting from a splice

389    connecting the LT first exon donor and MT acceptor (Riley et al., 1997). We identified a

18

390   novel T antigen, ST2, that is generated from a splice from the LT first exon donor to a

391   splice acceptor within the ST ORF. This transcript is highly expressed in HPyV7 and

392   present at lower levels in BKPyV, MPyV, and MCPyV. Many PyV encode a truncated

393   variant of LT - this includes SV40 17kT, BKPyV truncT, MCPyV 57kT and JCPyV T'

394   proteins (Abend et al., 2009; Shuda et al., 2008; Trowbridge and Frisque, 1995; Zerrahn

395   et al., 1993). These transcripts contain a canonical LT splice and a subsequent splice

396   that removes a large portion of the LT ORF.

397

398   We found that the same secondary splice sites responsible for truncated LT variants

399   can be used to generate superT. superT was initially observed in many SV40-

400   transformed cell lines (Kress et al., 1979; Smith et al., 1979) - in a similar manner, we

401   find that concatemeric integration of BKPyV and MCPyV in human cancers can facilitate

402   the generation of superT. We also find that superT is generated in lytic infections of

403   SV40, BKPyV, MCPyV and JCPyV. Eul and colleagues have published several studies

404   proposing that SV40 superT can be generated by trans-splicing between two separate

405   pre-mRNAs in the context of artificial expression constructs encoding the SV40 early

406   region (20, 31, 32). However, we find that in MCC tumors that generate superT and for

407   which we can assemble the viral integration site, the viral genome is likely integrated in

408   tandem in a way that could facilitate the cis-splicing of pre-mRNA that spans multiple

409   genome copies. Thus, while we cannot rule out trans-splicing from these data, we

410   believe cis-splicing is more likely. Future studies are necessary to understand the

411   biology of superT including its oncogenic potential and ability to bind multiple RB

412   molecules. Finally, efforts should be taken to understand if superT is expressed by PyV

19

413     and contributes to disease in other contexts, such as by BKPyV in PVAN or JCPyV in

414     PML.

415

416     We show that complex, uncharacterized splicing events are used by PyV to expand

417     their protein coding capacity. Future work is necessary to understand the biological

418     function of these transcripts and proteins. It is possible that unannotated splicing we

419     identify here could be differentially abundant in other biological contexts, so it will be

420     important to investigate PyV splicing in other infection contexts and human diseases.

421     Future transcriptome analyses that integrate long and short reads from multiple viruses

422     may have utility to expand characterized transcript and coding capacity in other viral

423     families.

424

## Conclusions

426     We provide a comprehensive transcriptome atlas for the prototype PyV SV40, as well

427     as the critically important human pathogen BKPyV. Comparative analyses of PyV

428     transcriptomes reveals conserved splice events that may expand PyV coding capacity.

429     We find that superT, a transcript generated by SV40, BKPyV, JCPyV, and MCPyV that

430     encodes a T antigen containing two RB-binding LxCxE domains, is present in several

431     PyV-associated human cancers. Together, these data expand our understanding of PyV

432     transcriptomes and uncover unannotated PyV-encoded proteins of potential relevance

433     to human disease.

434

435

# Materials and Methods

**Data and code availability.**

All code used in this project can be found at the zenodo and github links below. The

zenodo repository also contains all processed data necessary to reproduce all analyses

and figures. The main processing steps used to process RNAseq data are present as

nextflow pipelines which call modular bash and python scripts.

Zenodo: https://doi.org/10.5281/zenodo.5593468

Github: https://github.com/jnoms/SV40_transcriptome

Furthermore, a series of interactive Google Colab notebooks can download all

processed data from Zenodo and completely reproduce all analyses and non-schematic

primary figures. The colab documents are stored on github at

https://github.com/jnoms/SV40_transcriptome/tree/main/bin/colab. Direct links to the

Google Colab documents are as follows:

Figure 1:

https://colab.research.google.com/github/jnoms/SV40_transcriptome/blob/main/bin/cola

b/Figure1.ipynb

453     Figure 4:

454     https://colab.research.google.com/github/jnoms/SV40_transcriptome/blob/main/bin/cola

455     b/Figure4.ipynb

456     Figure 6:

457     https://colab.research.google.com/github/jnoms/SV40_transcriptome/blob/main/bin/cola

458     b/Figure6.ipynb

459

460     A Google Colab notebook is available for interactive investigation of all SV40 and

461     BKPyV viral transcript classes, and does not require computational skills to use:

462     https://colab.research.google.com/github/jnoms/SV40_transcriptome/blob/main/bin/cola

463     b/PyV_exploratory.ipynb

464

465     All raw RNA sequencing data are available at the NCBI sequence read archive at

466     accession XXXXXX.

467

468     **Datasets**

469     Information on all samples and viruses (excluding tumors) can be found in **Table 1**.

| Table 1 | | | | |
|---------|----------------|----------------------|------------------|------|
| **Virus** | **Sequencing Type** | **Origin (Accession)** | **MOI / Timepoint** | **Host** |

| SV40 | dRNAseq (two replicates) | Generated here | MOI 1 / 48hpi | *C. Sabaeus* |
|---|---|---|---|---|
| SV40 | SMRTseq | Generated here | MOI 1 / 48hpi | *C. Sabaeus* |
| SV40 (polysome input/whole-cell) | dRNAseq | Generated here | MOI 1 / 44hpi | *C. Sabaeus* |
| SV40 (polysome) | dRNAseq | Generated here | MOI 1 / 44hpi | *C. Sabaeus* |
| SV40 | Short-RNAseq (total) | Generated here | MOI 1 / 48hpi | *C. Sabaeus* |
| SV40 | short-RNAseq (polyA) | Generated here | MOI 1 / 48hpi | *C. Sabaeus* |
| BKPyV (Dunlop) | dRNAseq | Generated here | MOI 0.5 / 3dpi | Human |
| BKPyV (Dunlop) | SMRTseq | Generated here | MOI 0.5 / 3dpi | Human |
| BKPyV (Dunlop) | Short-RNAseq (total) | Generated here | MOI 0.5 / 3dpi | Human |
| BKPyV (Dunlop) | short-RNAseq (polyA) | Generated here | MOI 0.5 / 3dpi | Human |
| BKPyV (Dik) WT | Short-RNAseq (total) | Generated here | MOI 1 / 5dpi | Human |
| BKPyV (Dik) WT | Short-RNAseq (polyA) | Generated here | MOI 1 / 5dpi | Human |
| BKPyV (Dik) M1 | Short-RNAseq (polyA) | Generated here | MOI 1 / 5dpi | Human |

| BKPyV (Dik) M2 | Short-RNAseq (polyA) | Generated here | MOI 1 / 5dpi | Human |
|---|---|---|---|---|
| MPyV | dRNAseq | Generated here | Unknown / 28hpi | Mouse |
| MPyV | Short-RNAseq (total) | Garren et al. (Garren et al., 2015) (SRR2043214) | MOI 50 / 36hpi | Mouse |
| JCPyV | Short-RNAseq (total) | Assetta et al. (Assetta et al., 2016) (SRR9967610) | Unknown / 9dpi | Human |
| MCPyV (Synthetic genome) | short-RNAseq (polyA) | Theiss et al. (Theiss et al., 2015) (EBI: ERS760222) | 200ng viral DNA / Unknown | Human |
| HPyV7 | Short-RNAseq (total) | Rosenstein et al. (Rosenstein et al., 2021) (SRR11488976, SRR11488977) | From infected human skin | Human |
| BSPyV1 | Short-RNAseq (total) | Identified by Schmidlin et al. (Schmidlin et al., 2021) (SRR5958578) | From whole scorpion | *C. sculpturatus* |

470

471 **Tumor samples**

472 The BKPyV-associated bladder cancer is sample TBC03 that has been described

473 (Starrett et al., 2021). This sample is stranded, short-RNAseq (total).

474

475 Merkel cell carcinoma samples: Sections of tissue were isolated from patient-derived

476 tumor biopsies and suspended in RNAlater (Thermo Fischer) until further processing.

24

477    RNA and DNA was extracted from each section via the AllPrep DNA/RNA kit (Qiagen).

478    Isolated RNA and DNA were each sequenced (PE150) on the NovaSeq 6000 platform

479    (Illumina) for a depth of 50 M reads or 60x genomic coverage per sample, respectively

480    (Novogene). RNAseq data are unstranded, short-RNAseq (polyA).

481

482    **SV40 infection and RNA extraction**

483    BSC40 cells (ATCC CRL-2761) were seeded on 150mm dishes at $5.37*10^6$ cells per

484    plate - about 70% confluence. After waiting 4 hours for the cells to adhere, cells were

485    infected with SV40 at MOI 1 as previously described (Tremblay et al., 2001) with slight

486    modification. In brief, maintenance media was removed, and each 150mm dish was

487    inoculated with 6mL of virus stock diluted in DMEM + 2% FBS. Infection was allowed to

488    proceed at 37°C, 5% CO2 for one hour, with the plates rocked every 15 minutes to

489    ensure adequate coverage of the solution over the cell monolayer. At the end of this

490    period, DMEM + 2% FBS was added to a final volume of 25mL per 150mm dish. Each

491    dish was then incubated at 37°C, 5% CO2 for 48 hours. RNA was extracted using the

492    QIAGEN RNeasy Mini Plus Kit (QIAGEN 74134). This total RNA was then subjected to

493    Nanopore direct RNA sequencing and Illumina total- and polyA-RNA sequencing as

494    described below.

495

496    **BKPyV infection and RNA extraction**

497    Archetype and rearranged BKPyV (Dik and Dunlop, respectively) were purified and

498    titrated as described (Jiang et al., 2009b). RPTE-hTERT cells (Zhao and Imperiale,

499  2019) were plated in 6-well plate and prechilled for 15 min at 4°C and infected with Dik

500  or Dunlop at a MOI of 1 and 0.5 fluorescence-forming unit (FFU)/cell, respectively. The

501  cells were incubated at 4°C for 1 h with gentle shaking every 15 min. The virus was

502  removed and fresh REGM medium was added to the cells. Dik and Dunlop infected

503  cells were collected at 120 hpi and 96 hpi, respectively. Total RNA was extracted using

504  the Direct-zol RNA MiniPrep kit (ZYMO Research, USA). This total RNA was then

505  subjected to Nanopore direct RNA sequencing and Illumina total- and polyA-RNA

506  sequencing as described below.

507

508  **MPyV infection and RNA extraction**

509  C57 mouse embryo fibroblasts (ATCC SCRC-1008) were plated on a 150mm dish at

510  40% confluence. After several hours of growth, the typical DMEM + 10% FBS media

511  was replaced with serum free DMEM. The next day, the crude viral stock was thawed at

512  37°C, incubated at 45°C for 20 minutes to facilitate the final liberation of virus into the

513  supernatant, and cell debris removed from the viral stock with centrifugation. The

514  prepared virus stock was then diluted 1:10 with an absorption buffer consisting of HBSS

515  with 10mM HEPES, 1% FBS, at pH 5.6. Media was removed from the target cells, and

516  6mL of diluted virus in absorption buffer was added. Infection was allowed to proceed at

517  37°C, 5% $CO_2$ for one hour, with the plates rocked every 15 minutes to ensure

518  adequate coverage of the solution over the cell monolayer. At the end of this period, the

519  absorption buffer was removed and DMEM + 2% FBS was added to a final volume of

520  25mL per 150mm dish. Cells were inoculated for 28 hours at 37°C, 5% $CO_2$, after which

521  RNA was extracted using TRIzol (ThermoFisher 15596026) according to the

26

522    manufacturer's instructions. This total RNA was then subjected to Nanopore direct RNA

523    sequencing as described below.

524

525    The virus stock used here was kindly provided by the lab of Robert Garcea. This virus

526    stock (viral strain NG59RA) was a crude supernatant from MPyV-infected cells originally

527    generated by the lab of Thomas Benjamin on 02/08/2011 and was of unknown titer.

528    This stock was subjected to a total of three freeze-thaw cycles before use.

529

530    **SV40 polysome profiling**

531    BSC40 cells were plated on 4 150mm dishes at 60% confluence. After waiting 4 hours

532    for the cells to adhere, cells were infected with SV40 at MOI 1 as reported above. At 44

533    hours post infection cell culture media was replaced with media containing 100ug/mL

534    cycloheximide and incubated for 5 minutes. Plates were placed on ice, media

535    discarded, and cells were scraped into PBS containing 100ug/mL cycloheximide. Cells

536    were spun down, the PBS discarded, and cells were lysed in a lysis buffer containing

537    10mM Tris (pH 8), 100mM KCl, 10mM MgCl2, 2mM DTT, 1% Triton X100, 100ug/mL

538    cycloheximide, and 1unit/uL SUPERase RNase inhibitor (Thermo AM2694). Lysates

539    were incubated on ice for 20 minutes with intermittent tapping, and then spun at

540    10,000g for 10 minutes at 4°C. The supernatant was loaded onto a 10-55% sucrose

541    gradient followed by ultracentrifugation (Beckman Coulter Optima XPN-100

542    ultracentrifuge) at 32,500 × rpm at 4 °C for 80 minutes in the SW41 rotor. Gradients

543    were prepared with a gradient mixer and pump. Samples were separated by density

544    gradient fractionation system (Biocomp Piston gradient fractionator IP). RNA was

545    extracted from reserved input ("whole-cell") lysate, as well as the polysome fraction

546    using TRIzol. Equal volumes of each fraction containing heavy polysomes (>2) was

547    pooled prior to extraction (Lee et al., 2020).

548

549    **Western blotting**

550    Infection of wildtype Dik and two Dik mutants in RPTE-hTERT cells was performed as

551    mentioned above. Protein samples were harvested in E1A buffer with protease and

552    phosphatase inhibitors, electrophoresed, transferred, and probed with large tumor

553    antigen antibody (pAb416) as previously described (Zhao and Imperiale, 2019).

554    **RNA sequencing**

555    The concentration of total RNA was determined using the Qubit Fluorometer with the

556    Qubit RNA HS Assay Kit (ThermoFisher Q32852). RNA quality was then assessed on

557    an Agilent Bioanalyzer and the RNA 6000 Pico Kit (Agilent 5067-1513). PolyA RNA was

558    isolated using the NEBNext Poly(A) mRNA Magnetic Isolation Module (NEB E7490S)

559    with an input of 5ug of total RNA - for SV40 and BKPyV Dunlop, up to 8 total reactions

560    were used to yield sufficient polyA RNA (500ng) for subsequent protocols. In the case of

561    MPyV, due to limited amounts of total RNA, three reactions were used to yield roughly

562    100ng of polyA RNA. PolyA RNA concentration was determined again using the Qubit

563    RNA HS Assay Kit (ThermoFisher Q32852). PolyA RNA was then concentrated to 9uL

564    using a centriVap.

565

28

566    500ng of polyA RNA (or, in the case of MPyV, 100ng) in 9uL was then processed using

567    the Nanopore Direct RNA sequencing kit (SQK-RNA002). Resultant libraries were

568    sequenced for up to 24 hours on a MinION using an R9.4.1 flow cell.

569

570    In the case of polysome profiling: extracted RNA from the input and polysomes were

571    separately subjected to 5 reactions each of the NEBNext Poly(A) mRNA Magnetic

572    Isolation Module using 5ug RNA input per reaction. All resultant polyA RNA was then

573    processed using the Nanopore Direct RNA sequencing kit (SQK-RNA002). Resultant

574    libraries were sequenced for up to 24 hours on a MinION-Mk1C using an R9.4.1 flow

575    cell.

576

577    Illumina total RNA sequencing and polyA RNA sequencing of SV40-, BKPyV Dunlop-,

578    and BKPyV Dik-infected cells was conducted by Novogene Corporation Inc. The QC for

579    the RNA samples was performed using Qubit and Bioanalyzer instruments. Libraries

580    were then prepared using NEBNext Ultra II with RiboZero Plus kit (for short-RNAseq

581    (total)) and NEBNext Ultra II with PolyA Selection kit (for short-RNAseq (polyA)). Both

582    library approaches are strand-specific. Library quality and concentration was assessed

583    with Labchip and qPCR. Libraries were sequenced on NovaSeq6000 using PE150

584    sequencing.

585

586    PacBio SMRT sequencing of SV40 and BKPyV Dunlop was conducted by the Georgia

587    Genomics and Bioinformatics Core. Each sample was subjected to IsoSeq library

29

588    preparation and sequenced on an individual 8M SMRT cell for 26 hours on a Sequel-II

589    machine.

590

**Initial Sequence Processing**

592    Raw Nanopore dRNAseq reads from standard SV40, BKPyV Dunlop, and MPyV

593    infections were basecalled with Guppy version 4.2.2 with the following command:

594    guppy_basecaller -i fast5 -s basecalled --flowcell FLO-MIN106 --kit SQK-RNA002 -r --

595    trim_strategy rna --reverse_sequence true --u_substitution true --

596    cpu_threads_per_caller 10

597

598    Raw Nanopore dRNAseq reads from polysome profiling of SV40 transcripts were

599    basecalled on a MinION-Mk1C using MinKNOW version 21.02.2.

600

601    PacBio SMRTseq subreads were processed using ccs (version 6.0.0). Full-length,

602    nonchimeric reads were then generated using the lima (version 2.0.0) and Isoseq3

603    (version 3.4.0) packages provided by PacBio.

604

605    Stranded Illumina short-RNAseq (total) and short-RNAseq (polyA) reads were

606    processed in the following way: Files containing read 1 (R1) and read 2 (R2) were

607    trimmed and adapters removed using Trim Galore! (Krueger, 2016). Next, reads in R1

608    files were reverse complemented to orient the reads correctly relative to the transcript of

609    origin, and all read headers in the R1 and R2 files were labeled with "_1" or "_2"

610    respectively. The R1 and R2 files were then concatenated. This Illumina processing

611    pipeline is available in process_illumina.nf.

612

613    The MCC tumor RNAseq assessed in this manuscript were short-RNAseq (polyA) that

614    were NOT stranded. This means that the strand of origin of each read is unknown. To

615    address this uncertainty, the complement AND reverse complement of both R1 and R2

616    were concatenated into the final FASTQ file. As described below in the section

617    "Processing of short-read short-RNAseq (total) and short-RNAseq (polyA) span files",

618    future processing kept the most-likely alignment strand for each read.

619

620    **Sequence Alignment and Processing**

621    Most long-read sequencing data and Illumina sequencing data were aligned to the

622    appropriate viral genome using Minimap2 (Li, 2018). The exceptions are the short-

623    RNAseq (total) JCPyV data from Assetta et al. (Assetta et al., 2016) and the HPyV7

624    data from Rosenstein et al. (Rosenstein et al., 2021) - these samples contained

625    sequencing reads of 101bp or shorter and were instead mapped with STAR (Dobin et

626    al., 2013). All non-primary alignments were discarded. Sequence alignments in BAM

627    format were then converted to BED using bedtools (Quinlan and Hall, 2010). Here,

628    bedtools considers any Minimap2- or STAR-called intron ("N" cigar flag) as an intron to

629    split alignment segments. Parameters for alignment and bed conversion can be found in

630    minimap2.sh and star.sh.

31

631

632  To capture transcripts that originate from a pre-mRNA that circled the viral genome

633  more than once, and therefore contain repetitive sequences, all alignments were

634  conducted against concatenated copies of the viral genome. In the case of short-read

635  short-RNAseq (total) and short-RNAseq (polyA), the reference consisted of two

636  concatenated copies of the viral genome. For long-read dRNAseq and SMRTseq, the

637  reference consisted of twenty concatenated copies of the viral genome.

638

639  Because the references consisted of multiple copies of the same viral genome, mapped

640  reads were assigned to a random copy of the genome. Therefore, all reads in resultant

641  BED files were "slid" such that they started in the first genome copy of the reference

642  using bed_slide_wraparound_reads.py.

643

644  All reference genomes can be found in resources/ref directory of the associated github

645  repository. All references used contain the PyV late region at the start/5' end of the

646  reference on the "+" or sense strand, with the early region on the antisense or "-" strand.

647  The concatenated references are based on the following reference genomes collected

648  from NCBI, with any modifications listed:

649  -  SV40: NC_001669.1. The first 100 nucleotides were moved to the end of the

650     sequence.

651  -  BKPyV: KP412983.1

652  -  JCPyV: NC_001699.1

653    -   MPyV: NC_001515.2. The sequence was reverse-complemented to orient the

654        late region towards the start of the reference.

655    -   MCPyV: NC_010277.2

656    -   HPyV7: NC_014407.1

657    -   BSPyV1: LN846618.1

658

659    Next, a span file was generated from each slid BED file using bed_to_span.py. This

660    script splits each read into "spans", where each span is an exon or an intron with all

661    positions relative to the viral genome. The introns are defined by the Minimap2- or

662    STAR-called introns ("N" cigar flag) as mentioned above. All regions between the start

663    and end of the reads that are not introns were called as distinct exons. Transcripts were

664    clustered into transcript classes based on introns as discussed below. A "tidy" output

665    span file was then generated that contains the name, strand, and transcript class of a

666    given read, with separate lines for the start and end of each span (e.g., exon or intron)

667    within the sequencing read.

668

669    **Alignment of repetitive regions**

670    Reads that originate from a transcript that circles the genome more than once can be

671    detected because there is one or more repetitive regions within the read. Alignment

672    against multi-copy reference genomes (20 copies in the case of dRNAseq and

673    SMRTseq) as described above sufficiently captured most of these transcripts, with

674    some exceptions. First, BKPyV SMRTseq data had a poor alignment rate of the leader

33

675  exon in late WA transcripts - this means that WA transcripts are underrepresented in the

676  BKPyV SMRTseq data. Second, alignment of superT and superT* transcripts from

677  SMRTseq and dRNAseq data was generally poor, with the repetitive region often failing

678  to map via Minimap2. Potential superT and superT* reads in dRNAseq and SMRTseq

679  data were identified through assessment of BAM files following mapping. Early reads

680  that contain a CIGAR flag showing an insertion of 100 bases or more were flagged, and

681  up to 50 of these transcripts were manually investigated through online BLASTN

682  (Johnson et al., 2008) against the viral reference genome. Reads supporting superT in

683  SV40 dRNAseq data, superT* in SV40 SMRTseq data, and superT in BKPyV SMRTseq

684  data were initially missing from Minimap2 alignments but were identified via this

685  approach. One transcript of each type was then repaired upon data import to R such

686  that these transcripts are represented in downstream visualizations - these actions are

687  clearly marked in UTILS_import_data.R. Thus, superT and superT* in SMRTseq and

688  dRNAseq data are underrepresented in abundance plots (**Figure S3D, S4D**) and read

689  pileups (**Figure S6, S8**) compared to their actual abundance in the cell due to these

690  alignment challenges.

691

692  **Generation of transcript classes**

693  Transcript classes were generated during processing of BED files using

694  bed_to_span.py. Each transcript class consists of sequencing reads that contain the

695  same combination of introns. The transcript class number is based on the abundance of

696  transcripts within a transcript class - e.g., transcript class 1 contains more transcripts

697  than transcript class 2, and so on. Transcript class generation is similar for both long-

34

698    and short-read sequencing data, although short reads usually (but not always) tend to

699    contain a maximum of one intron. Notably, transcript class assignment is independent of

700    the transcript start and end positions, meaning that there can be heterogeneity of

701    transcript start and end positions within a transcript class. For all SMRTseq and

702    dRNAseq data, for a transcript class to be generated all introns contained within the

703    transcript class were required to be supported by at least 5 junction-spanning reads

704    within a short-RNAseq (total) dataset. For SV40 and BKPyV Dunlop SMRTseq and

705    dRNAseq data, the short-RNAseq (total) data was generated from RNA from the same

706    extraction. SV40 dRNAseq replicate 2 was corrected with the short-RNAseq (total) data

707    from the first SV40 replicate. For the MPyV dRNAseq data, short-RNAseq (total) data

708    from Garren et al. was used. If a transcript contained an intron that was not supported

709    by at least 5 junction-spanning reads in the Illumina dataset, it was discarded. We opted

710    to use this filtering strategy rather than implementing long-read correction because

711    correction algorithms were unable to cope with wraparound transcripts.

712

713    There were limited circumstances where dRNAseq or SMRTseq transcript classes were

714    removed manually during processing - this occurred to four transcript classes that made

715    it through filtering. In these circumstances, alignments were deemed to be artifactual

716    due to Minimap2 alignment errors. These instances are clearly programmatically

717    marked in UTILS_import_data.R with specific rationale for each action.

718

719    **Splice coordinate system**

35

720   All splice or intron positions marked in any figure or table of this manuscript are **0-**

721   **indexed positions of the intron**. To convert these coordinates to the 1-

722   indexed/absolute position of the intron on the viral genome, add 1 to the intron start

723   position. For example, for the intron 276-1600, viral genome nucleotide # 277 is the first

724   nucleotide within the intron, and viral genome nucleotide # 1600 is the last nucleotide

725   within the intron.

726

727   **Processing of short-RNAseq (total) and short-RNAseq (polyA) span files**

728   The majority of the short-read RNAseq data investigated here used a strand-specific

729   sequencing strategy (except for the MCC tumor RNAseq). With this strategy, the strand

730   of origin for the transcript yielding each read is known, and a read can be correctly

731   assigned to the sense ("+" / late) or antisense ("-" / early) strand. However, a fraction of

732   transcripts can be inaccurately stranded due to artifacts during library preparation.

733   When there were many more late reads than early reads in a short-read dataset, a

734   prohibitive fraction of "early" reads would be reads from late transcripts that were

735   incorrectly stranded due to this artifact. To address this issue, short reads that aligned

736   to the + strand were required to either start or end within the late region (defined as the

737   first ½ of the genome), and short-reads that aligned to the - strand were required to

738   either start or end within the early region (defined as the second ½ of the genome).

739

740   **Transcript identification**

36

741 For **Figure 1** and all supplementary figures, a SV40 or BKPyV transcript was identified

742 and assigned a transcript ID if it was at least 0.1% of early or late strands in dRNAseq

743 or SMRTseq data with one exception - SV40 transcript L8 had been previously

744 identified and was kept despite being at only 0.06% abundance. Existing transcript

745 names, where available, were taken from relevant studies (Abend et al., 2009; Good et

746 al., 1988; Seif et al., 1979; Zerrahn et al., 1993). This assignment occurred from the

747 span files, meaning that all sequencing reads in question were previously required to

748 contain introns that were supported by at least 5 short-RNAseq (total) junction-spanning

749 reads. For SV40, for which there was two dRNAseq replicates, identification of a

750 sequencing read at 0.1% or greater in just one replicate was sufficient.

751

752 Transcript IDs (e.g., E1, E2, E3…, L1, L2, L3,…) consist of the kinetic class (E: Early, or

753 L: Late) of the identified transcript followed by an integer value in ascending order of

754 abundance. This abundance value was calculated by ordering the transcripts in order of

755 the maximum observed relative abundance in dRNAseq or SMRTseq data.

756

757 Of note, the relative abundance of transcripts between dRNAseq and SMRTseq data is

758 skewed by distinct read-length biases between the two approaches. The dRNAseq

759 approach has a 3' bias and a bias towards shorter transcripts, while SMRTseq library

760 preparation resulted in preferential sequencing of transcripts closer to ~2500bp in

761 length. Resultant differences in the length of aligned reads can be seen in **Figure S1C**.

762 The TSS distribution of SV40 late transcripts varies between transcript classes, while

763    the late TSS distribution tends to be similar across transcript classes in BKPyV (**Figure**

764    **S2A**).

765

766    **Calculation of sequencing coverage**

767    To determine the sequencing coverage for each sample (as in **Figure 1C**), BAM files

768    from alignment were "slid" such that all transcripts must start in the first genome copy of

769    the reference using bam_slide_wraparound_reads.py, in a similar manner as the beds

770    were slid as described above. Forward and reverse strand reads were split, and the

771    depth was calculated using the command `samtools depth -aa -d0` separately for

772    forward and reverse reads. These processing steps are present in bam_coverage.nf.

773    During plotting, the coverage for each strand was normalized to the maximum coverage

774    at any position (e.g., the maximum coverage of the late and early strands was set to 1).

775

776    **Watch plots**

777    Each panel of a watch plot represents information for a single transcript class. The

778    center "arms" of these plots are histograms detailing the distribution of start (blue) and

779    end (red) positions for the transcripts within the transcript class. These histograms are

780    normalized to the highest abundance position. The outer ring of each watch plot shows

781    the viral ORF map. Each inner grey ring indicates the number of genomes spanned - all

782    transcripts are displayed moving outwards from the center. Red segments indicate the

783    **exons** of each transcript class. The first exon starts on the most-inner grey ring at the

784    most common transcript start site for the transcript class, and the last exon ends on the

785    most-outer grey ring at the most common transcript end site for the transcript class. The

786    3' end of the transcript is indicated by the red arrow at the end of the last exon. Thus,

787    the transcripts spiral outwards from the center in the direction of the red arrow. **Figure 4**

788    contains a schematic key describing watch plots.

789

790    **Read pileup plots**

791    Each square/rectangular panel of a read pileup plot shows the reads present in a single

792    transcript class. The arrows at the top of each panel indicate the viral ORF map, with

793    dashed lines indicating the end of each genome copy. Next, the lines indicate

794    histograms of the transcript start (blue) and end (red) sites for the transcripts within the

795    transcript class. Below the x-axis, each row indicates a single sequencing read. The

796    spans in red indicate the exons inferred from a sequencing read, while the spans in pink

797    indicate the introns/splice junctions. Sometimes the distribution of transcript end

798    positions for a transcript class can be obscured by the thickness of the transcript lines -

799    the histograms should always be consulted to assess abundance.

800

801    For SV40 dRNAseq watch and pileups: There were two SV40 dRNAseq replicates.

802    Watch plots and read pileups are based on replicate 1, although missing transcripts that

803    were identified in replicate 2 but not 1 were also plotted.

804

805    **Short-read intron plots** (**Figure S9, S11**)

806    In these plots, lines indicate specific introns. The upper and lower horizontal arrows

807    indicate the viral ORF map - often, these ORF maps will indicate two concatenated viral

808    reference genomes. The circles above or below each ORF map indicate the percentage

809    of early or late introns that fall at each genome position. Early introns and percentages

810    are colored red, while late introns and percentages are colored blue.

811

812    **polyA tail length**

813    polyA tail length was determined from dRNAseq data using the `polya` command of

814    Nanopolish (Loman et al., 2015). To determine the polyA distribution of host transcripts,

815    sequencing reads were aligned to the human GRCh38 (for BKPyV samples), *C.*

816    *Sabaeus* (for SV40 samples), or mouse (for MPyV) cDNA transcriptomes downloaded

817    from ensembl. Only reads with a Nanopolish QC tag of "PASS" were considered for

818    downstream polyA tail length analyses.

819

820    The dRNAseq library preparation included the addition of the "RNA Control Standard"

821    (RCS), which is a synthetic RNA based on yeast ENO2 containing a 30-adenine polyA

822    tail. dRNAseq samples were mapped against ENO2 to assess the polyA tail length

823    distribution of this control.

824

825    The cumulative incidence of transcript termination (**Figure S10D**) was calculated by

826    determining, for each early read, how far the read's transcript end site is from the

827    canonical polyA site position for each virus.

828

829    **Polysome profiling analysis**

830    To determine the ribosome occupancy of host genes, dRNAseq reads were aligned to

831    the *C. Sabaeus* cDNA transcriptome downloaded from ensembl. The number of reads

832    mapped to each transcript was extracted with `samtools idxstats`. Transcripts were

833    filtered to include only those with at least 10 reads in both polysome and input fractions.

834    The normalized abundance of each transcript in each fraction was defined as (# of

835    mapped reads)/(total number of virus and host mapped reads). Ribosome occupancy of

836    each transcript was determined as (normalized abundance in polysome)/(normalized

837    abundance in whole-cell), where a value of >1 indicates preferential translation.

838

839    Ribosome occupancy of individual viral transcripts could not be calculated because of

840    increased rates of transcript truncation in the polysome fraction compared to the whole-

841    cell fraction. This was indicated by a nearly doubled proportion of unspliced reads with

842    premature 5' ends in the polysome fraction compared to the whole-cell fraction, and

843    likely indicates transcript degradation during sucrose centrifugation or fraction collection.

844    Because viral transcripts are mostly identical and vary largely at a 5' splice site,

845    elevated transcript truncation decreased the observed abundance of individual viral

41

846 transcripts in the polysome fraction and make ribosome occupancy calculations for

847 individual viral transcripts unreliable.

848

849 **MCC440 superT PCR and sanger sequencing**

850 Anchored poly-dT primers (Life Technologies) were used for specific reverse-

851 transcription of full-length mRNA into cDNA. Primers were designed to uniquely amplify

852 the super-LT junction through exploitation of repetitive sequences. Primer sequences

853 were as follows (5' -> 3'); Forward: CTGGACTGGGAGTCTGAAGC, Reverse:

854 ACCCCTCCTCCATTCTCAAGA. Q5 polymerase (NEB) with standard reaction

855 conditions was used for amplification.

856

857 **Generation of integrated PyV structures and viral variant calling**

858 Tumor WGS was aligned against a fusion reference genome containing hg38 and

859 Merkel cell polyomavirus (NC_010277) using bowtie2 with default parameters.

860 Integrated virus assembly graphs and annotations were generated using Oncovirus

861 tools (https://github.com/gstarrett/oncovirus_tools). Assembly graphs were then

862 manually interpreted to create linear integration structures for PyV-associated MCC.

863

864 Point mutations were called in the PyV genomes using lofreq with default parameters

865 (https://csb5.github.io/lofreq/) (PMID: 23066108). Lofreq output was functionally

866 annotated with SnpEff (http://pcingola.github.io/SnpEff/) (PMID: 22728672) using the

867    relevant GenBank gene annotations for the above genomes. Variants were plotted out

868    in R with the ggplot2 package.

869

870

# List of abbreviations

872    APA - Alternative polyadenylation

873    BKPyV - BK Polyomavirus

874    BSPyV1 - Bark scorpion polyomavirus 1

875    dRNAseq - Nanopore direct RNA sequencing

876    HPyV7 - Human polyomavirus 7

877    JCPyV - JC Polyomavirus

878    LT - Large T antigen

879    MCC - Merkel cell carcinoma

880    MCPyV - Merkel cell polyomavirus

881    MPyV - Murine polyomavirus

882    MT - Middle T antigen

883    ORF - Open reading frame

884    SMRTseq - PacBio SMRT sequencing

885    ST - Small T antigen

886    SV40 - Simian virus 40

887    TES - Transcript end site

888    TSS - Transcript start site

889    PVAN - Polyomavirus-associated nephropathy

890    PyV - Polyomavirus

891

892

# Declarations

893

**Competing interests**

894

895    M.M. receives research support from Bayer, Janssen, Ono; consults for Bayer, Interline,

896    Isabl; and receives patent royalties from Labcorp and Bayer. J.A.D. has received

897    research support from Rain Therapeutics, Inc. and is a consultant for Rain

898    Therapeutics, Inc. and Takeda, Inc.

899

**Funding**

900

908

**Authors' Contributions**

910    Conceptualization - J.N., M.M., J.A.D.

911    Data curation - J.N.

912    Formal analysis - J.N.

913    Funding acquisition - M.M, J.A.D., M.J.I., S.V., G.J.S.

914    Investigation - J.N., W.Z., T.C.F., C.D.

915    Methodology - J.N., M.M., J.A.D.

916    Project administration - J.N.

917    Resources - J.N., M.M., J.A.D., T.C.F., W.Z., M.J.I, S.V., C.D.

918    Software - J.N.

919    Supervision - M.M., J.A.D.

920    Validation - J.N.

921    Visualization - J.N.

922    Writing, Original draft - J.N., M.M., J.A.D.

923    Writing, review & editing - J.N., W.Z., T.C.F., C.D., S.V., G.J.S., M.J.I., M.M., J.A.D.

924

934

# Figure legends

936    **Figure 1** - RNA sequencing expands known SV40 and BKPyV transcript diversity.

937    **A.** Overview of experimental procedures. Cells were infected with a polyomavirus,

938        and RNAs extracted. RNA was sequenced using long-read (Nanopore dRNAseq

939        and PacBio SMRTseq) and short-read (Illumina short-RNAseq (total) and short-

940        RNAseq (polyA)). Transcripts were analyzed, and the impact of observed splice

941        events on viral open reading frames was assessed.

942    **B.** Mechanism of transcript clustering in this study. Transcripts were aligned to the

943        viral genome and grouped into transcript classes based on the presence of

944      shared introns. Thus, within a transcript class there may be variation in the exact

945      transcript start and end positions. This clustering strategy was used for both long-

946      and short-RNAseq data.

947   **C.** Viral RNA sequence coverage for SV40 and BKPyV as determined from

948      dRNAseq, SMRTseq, and short-RNAseq (total) data. The Y axis indicates the

949      scaled coverage, with X axis indicating the position on the viral genome.

950      Coverage for late transcripts (mapping to the + strand) is above the x axis, while

951      coverage for early transcripts (mapping to the - strand) is below the x axis.

952      Coverage is scaled separately for each strand such that the maximum observed

953      coverage for each strand is 1. Arrows at the top of the plot indicate the positions

954      of viral genes.

955   **D-E**. UpSet plot indicating the overlap between existing transcript annotations,

956   dRNAseq data, and SMRTseq data for SV40 (**D**) and BKPyV Dunlop (**E**). Bars

957   indicating overlap with existing transcript annotations are black, while those

958   indicating no overlap with existing annotations are blue. These blue bars indicate the

959   number of novel, unannotated transcripts identified.

960   **F**. Overview of polysome profiling of SV40-infected cells. BSC40 cells were

961      infected with SV40. Cells were lysed, and a portion of the lysate was subjected to

962      dRNAseq (representative of the RNA content of the whole cell). The remaining

963      lysates was centrifuged through a sucrose gradient, after which fractions

964      containing RNA associated with two or more ribosomes were pooled and

965      subjected to dRNAseq.

47

966     **G**. Relative abundance of SV40 early and late transcripts in the whole-cell and

967     polysome fractions of SV40-infected cells. Y-axis indicates the percentage of

968     early or late transcripts and is log scale. X axis indicates each transcript, with

969     black dots indicating each transcript's whole-cell relative abundance and red dots

970     indicating each transcript's polysome relative abundance.

971

972     **Figure 2** - Annotated and novel SV40 transcripts.

973     **A.** Transcripts are shown relative to the viral genome. Each line is a viral transcript,

974     with red lines indicating exons and dashed blue lines indicating introns. Spokes

975     indicate the positions of common splice donors and splice acceptors. Transcripts

976     that were annotated prior to this study are on a yellow background, and novel

977     transcripts are on a while background. Wraparound transcription that results in

978     multiple copies of a region is annotated with double lines, and the number of

979     copies is indicated in parentheses. The line labeled "pA" indicates the

980     approximate position of the polyA signal sequence.

981     **Figure 3** - Annotated and novel BKPyV transcripts.

982     **A.** Transcripts are shown relative to the viral genome. Each line is a viral transcript,

983     with red lines indicating exons and dashed blue lines indicating introns. Spokes

984     indicate the positions of common splice donors and splice acceptors. Transcripts

985     that were annotated prior to this study are on a yellow background, and novel

986     transcripts are on a while background. Wraparound transcription that results in

987     multiple copies of a region is annotated with double lines, and the number of

48

988       copies is indicated in parentheses. The line labeled "pA" indicates the

989       approximate position of the polyA signal sequence.

990

991       **Figure 4** - Pervasive wraparound transcription across PyV

992       **A-C**. Watch plots indicating the top 4 highest abundance late wraparound

993       transcript classes in dRNAseq data from SV40 (A), BKPyV Dunlop (B), and

994       MPyV (C). The outer ring of each watch plot indicates the position of the viral

995       ORFs. The inner arms are histograms detailing the distribution of transcript starts

996       (in blue) and ends (in red) for transcripts within each transcript class. The red

997       segments indicate exons. Transcripts start in the innermost ring - a second or

998       third ring indicates that the pre-mRNA that generated the transcript must have

999       circled the viral genome multiple times. The 3' end of the transcript and the

1000       direction in which these plots are oriented is indicated by the red arrow at the end

1001       of the last exon segment. The red exon segments start at the most common

1002       transcript start site within the transcript class, and end at the most common

1003       transcript end site within the class. The watch plot key shows an example of the

1004       path of the pre-mRNA for SV40 transcript class L6_I.

1005    **D.** Bar plots indicating the percentage of late transcripts that span a given number of

1006       genome lengths in SV40, BKPyV Dunlop, and MPyV dRNAseq data.

1007    **E.** The leader-leader junction, that connects the pre-mRNA from one genome to the

1008       subsequent wraparound, was identified in Illumina short-RNAseq (total) data.

1009       The intron in question is plotted as a black line in this plot, with the x axis

49

1010      indicating the genomic position of the intron. The top late wraparound transcript

1011      for each virus was plotted. The gene map indicates the approximate gene

1012      position and is accurate for SV40 - the exact position of the viral genes varies

1013      between viruses. Percentages indicate the percentage of late junction-spanning

1014      transcripts that support the plotted wraparound leader-leader junction.

1015    **F.** Schematic illustrating how leader-leader wraparound transcription can be

1016      detected from short read short-RNAseq (total). Leader-leader splicing can be

1017      seen as a repetitive exon in watch plots from long-read RNAseq data. Ultimately,

1018      there was an original processed mRNA in the cell that contained two tandem

1019      leader sequences. When this transcript of origin is sequenced via short read

1020      sequencing, reads will be generated across its length. A minority of these reads

1021      will span the leader-leader junction, and mapping against the viral reference

1022      genome can be used to uncover leader-leader splicing.

1023  **Figure 5** - Detection of novel, conserved splicing events that expand PyV coding

1024  capacity.

1025  **A-D**. Schematics illustrating identified ORFs. Each row is a reading frame (except for

1026  ST and the LT 1st exon, which are in the same frame), and unannotated amino acids are

1027  represented by grey boxes. The measured intron is indicated by the red arrow. Colored

1028  ORFs are annotated, while grey ORFs are unannotated. Percentages on the right side

1029  of the figure are the percentage of spliced viral transcripts on the same strand as

1030  determined from short-read short-RNAseq (total) data. Numbers after each virus name

1031  indicate the transcript class within each short-RNAseq (total) dataset. The measured

1032  intron is indicated by the red arrow.

1033    A) ST2: This ORF is generating from a splicing event that uses the LT first exon

1034    donor and an acceptor within the ST ORF. In HPyV7 and BKPyV Dunlop, the splice

1035    lands in frame and results in an internal deletion within ST. In MPyV and MCPyV the

1036    splice lands out of frame, resulting in an ORF that contains the N-terminal region of

1037    ST and novel amino acids at the C terminus.

1038    B) MT: MPyV encodes a MT following splicing connecting the end of the ST ORF

1039    with an ORF in an alternate frame of the LT second exon. In BKPyV, a similar splice

1040    occurs connecting ST with an MT-like ORF in an alternative frame of the LT second

1041    exon.

1042    C) VP1X: JCPyV encodes two VP1X ORFs generated by splicing within VP1 and

1043    landing in an alternative frame of VP1, or earlier in the late region due to wraparound

1044    transcription. While predominant in JCPyV, VP1X is likewise present in many other

1045    PyV.

1046    D) superT: The superT-specific splice utilizes the splice donor canonically

1047    associated with truncated T antigens such as 17kT in SV40 and truncT in BKPyV.

1048    Due to wraparound transcription, a LT second exon acceptor is available to the 3' of

1049    this donor and acts as the acceptor. For the superT ORF to form, an initial LT splice

1050    is required. Ultimately, superT contains a duplication in part of the LT second exon

1051    that includes the RB-binding LxCxE motif.

1052  **E.** Schematics detailing BKPyV Dik isolates used for querying the existence of superT.

1053  BKPyV WT is wild type virus. M1 contains a LT intron that has been replaced with an

1054  intron from the plasmid pCI. Both WT and M1 are expected to generate LT and superT

1055    of expected sizes. M2 has a completely removed LT intron, and the pCI intron is located

1056    directly 5' of the LT ORF. M2 is expected to encode LT of expected size, but a larger

1057    superT variant due to incorporation of a second copy of the LT first exon.

1058    **F.** Western blot of cells infected with BKPyV Dik WT, M1, or M2 and probed with an

1059    antibody reactive against LT. The lower molecular weight band is LT, and the higher

1060    molecular weight bands are consistent with superT.

1061

1062    **Figure 6** - Detection superT-encoding transcripts in PyV-associated cancers

1063        **A.** Schematic detailing the generation of superT during lytic infection as compared

1064            to from integrated virus in cancer. During viral infection, the RNA polymerase can

1065            circle the viral genome multiple times, resulting in a pre-mRNA that can be

1066            spliced to generate superT. In the case of host integration, a polyomavirus can

1067            be integrated in tandem copies such that a pre-mRNA is generated with more

1068            than one copy of the viral early region. This pre-mRNA can be similarly spliced to

1069            generate a superT transcript.

1070        **B.** Heatmap indicating the abundance of the superT, ST, and LT introns from

1071            RNAseq data from two replicates of a BKPyV-positive bladder cancer and six

1072            MCPyV-associated MCCs. Percentages indicate the percentage of spliced early

1073            viral reads for each sample. The splice measured in each row is indicated by the

1074            red arrow in the schematics on the right side of the figure.

1075

1076    **Figure S1** - Sequencing statistics

52

1077     **A.** The number of reads for all datasets studied here. For long-read dRNAseq and

1078     SMRTseq, this number includes spliced and unspliced reads. Because short

1079     reads are only useful for transcript characterization when they span a splice

1080     junction, the counts for short-reads represent the number of splice-junction-

1081     spanning reads.

1082     **B.** The cumulative percentage of transcripts in each number of transcript classes, by

1083     strand. The X-axis indicated the total number of transcript classes. The Y axis

1084     indicates the cumulative percentage of transcripts within those transcript classes.

1085     These plots indicate that most transcripts in most samples are contained within

1086     the first few transcript classes.

1087     **C-E**. The alignment length distribution of early, late, spliced, and unspliced

1088     transcripts for dRNAseq and SMRTseq data from SV40 (**C**), BKPyV Dunlop (**D**), and

1089     MPyV (**E**). The X axis indicates the length of a read's alignment, while the Y axis

1090     indicates the density/percentage of transcripts with a given alignment length. This

1091     plot shows that dRNAseq and SMRTseq data sample from RNA populations of

1092     different length.

1093

1094   **Figure S2** - Transcript start sites and polyA tail lengths.

1095     **A, B**. The distribution of transcript start sites for late (**A**) and early (**B**) transcripts for

1096     SV40 (left column), BKPyV Dunlop (middle column), and MPyV (right column). The

1097     arrows indicate the viral ORF positions.

1098      **C.** The distribution of polyA tail lengths for the 30-adenine ENO2 control (black),

1099           host (red), and viral (yellow) transcripts for SV40, BKPyV Dunlop, and MPyV.

1100           The X axis indicates the length of the polyA tail, while the Y axis indicates the

1101           density/percentage of transcripts with each length.

1102      **D.** Ribosome occupancy of host transcripts in SV40-infected cells. Each grey dot is

1103           a host transcript. The red, blue, and black dots are specifically noted host

1104           transcripts. Ribosome occupancy is on the Y axis, while the X axis does not hold

1105           value. Lines on the violin plot indicate $1^{st}$, $2^{nd}$, and $3^{rd}$ quartiles.

1106

1107      **Figure S3** - SV40 transcriptome atlas, watch plots

1108           **A-C**. Watch plots indicating all identified transcripts in SV40. (**A**) and (**B**) show

1109           transcripts that were identified in both dRNAseq and SMRTseq data, while (**C**)

1110           shows transcripts identified in SMRTseq only.

1111      **E.** Barplots that show the abundance of each transcript type in the dRNAseq and

1112           SMRTseq data. Here, there are two dRNAseq bars (one per replicate). The Y

1113           axis indicates the percentage of transcripts of the same strand. As discussed in

1114           the methods, alignment of superT and superT* was challenging, so the actual

1115           abundance of these transcripts is higher than reported here.

1116

1117      **Figure S4** - BKPyV Dunlop transcriptome atlas, watch plots

1118    **A-C**. Watch plots indicating all identified transcripts in BKPyV Dunlop. (**A**) and

1119    (**B**) show transcripts that were identified in both dRNAseq and SMRTseq data,

1120    while (**C**) shows transcripts identified in dRNAseq only.

1121    **F.**  Barplots that show the abundance of each transcript type in the dRNAseq and

1122    SMRTseq data. The Y axis indicates the percentage of transcripts of the same

1123    strand. As discussed in the methods, alignment of superT and superT* was

1124    challenging, so the actual abundance of these transcripts is higher than reported

1125    here.

1126

1127    **Figure S5** - SV40 transcriptome atlas, late transcript read pileups

1128    **A, B**. Read pileups showing the late transcripts identified in SV40 dRNAseq (A)

1129    and SMRTseq (B). The arrows at the top of the plot indicate the viral ORF

1130    positions. Below the X axis, each row is an individual transcript, with exons

1131    indicated in red and splice junctions/introns indicated in pink. Above the X axis

1132    are histograms indicating the transcript start (blue) and transcript end (red) sites.

1133    (U: unspliced).

1134

1135    **Figure S6** - SV40 transcriptome atlas, early transcript read pileups

1136    **A, B**. Read pileups showing the early transcripts identified in SV40 dRNAseq (A)

1137    and SMRTseq (B). The arrows at the top of the plot indicate the viral ORF

1138    positions. Below the X axis, each row is an individual transcript, with exons

1139      indicated in red and splice junctions/introns indicated in pink. Above the X axis

1140      are histograms indicating the transcript start (blue) and transcript end (red) sites.

1141      (U: unspliced).

1142

1143   **Figure S7** - BKPyV Dunlop transcriptome atlas, late transcript read pileups

1144      **A, B**. Read pileups showing the late transcripts identified in BKPyV Dunlop

1145      dRNAseq (A) and SMRTseq (B). The arrows at the top of the plot indicate the

1146      viral ORF positions. Below the X axis, each row is an individual transcript, with

1147      exons indicated in red and splice junctions/introns indicated in pink. Above the X

1148      axis are histograms indicating the transcript start (blue) and transcript end (red)

1149      sites. (U: unspliced).

1150

1151   **Figure S8** - BKPyV Dunlop transcriptome atlas, early transcript read pileups

1152      **A, B**. Read pileups showing the early transcripts identified in BKPyV Dunlop

1153      dRNAseq (A) and SMRTseq (B). The arrows at the top of the plot indicate the

1154      viral ORF positions. Below the X axis, each row is an individual transcript, with

1155      exons indicated in red and splice junctions/introns indicated in pink. Above the X

1156      axis are histograms indicating the transcript start (blue) and transcript end (red)

1157      sites. (U: unspliced).

1158

1159   **Figure S9** - Intron plots for all datasets studied

56

1160    A. Intron plots generated from short-read RNAseq. The arrows at the top and

1161       bottom of each panel indicate the position of viral ORFs. The lines indicate

1162       specific introns identified in the RNAseq data, with the 5' end on the top and the

1163       3' end on the bottom. The blue color indicates late transcripts, with red indicating

1164       early transcripts. The size of the circles above and below the viral ORF maps

1165       indicate the percentage of junction-spanning reads with a 5' end (on top) or 3'

1166       end (on bottom) at that position. Junctions are plotted if they are at least 1% of

1167       early or late transcripts, except for the SV40 pA superT junction (transcript class

1168       3) which is just below threshold but is of interest.

1169    B. Another representation of intron plots for each virus. The top arrows indicate the

1170       position of viral ORFs. The X axis indicates the genomic position for each splice.

1171       The Y axis indicates a single transcript class, with that class' intron plotted as a

1172       line. The percentage of early or late transcripts is indicated with the numeric

1173       percentage. Junctions are plotted if they are at least 1% of early or late

1174       transcripts, except for the SV40 pA superT junction (transcript class 3) which is

1175       just below threshold but is of interest.

1176

1177    **Figure S10** - Alternative polyadenylation of early transcripts in SV40, BKPyV, and

1178    MPyV.

1179       **A-C**. Watch plots indicating the LT and ST transcripts for SV40 (**A**), BKPyV

1180       Dunlop (**B**), and MPyV (**C**). The focus of these plots is the distribution of

1181       transcript end positions, which are the inner red arms. The region of APA of

1182       highlighted in blue, with the canonical transcript end sites highlighted in red.

57

1183      **D**. A cumulative incidence plot of transcript termination in SV40 (blue), BKPyV

1184      Dunlop (red), and MPyV (green). The X axis indicates the distance to the

1185      canonical polyA site, while the Y axis indicates the percentage of transcripts that

1186      have terminated by that position.

1187      **E-G**. Density plots showing the distribution of polyA tail lengths for LT and ST

1188      transcripts that end at the canonical site (solid) or undergo APA (dashed) for

1189      SV40 (**E**), BKPyV Dunlop (**F**), and MPyV (**G**). The x axis indicates the length of

1190      the polyA tail, while the Y axis indicates the density/proportion of transcripts with

1191      the given length.

1192

1193    **Figure S11** - short-RNAseq (polyA) analysis of BKPyV Dik WT, M1, and M2

1194    **A-C**. Intron plots generated from short-read (polyA) RNAseq of cells infected with

1195    BKPyV WT, or the M1 or M2 mutants. The arrows at the top and bottom of each panel

1196    indicate the position of viral ORFs relative to the standard BKPyV genome - note that

1197    the genomes of mutants M1 and M2 are altered as indicated in Figure 5E. The lines

1198    indicate specific introns identified in the RNAseq data, with the 5' end on the top and the

1199    3' end on the bottom. The size of the circles above and below the viral ORF maps

1200    indicate the percentage of junction-spanning reads with a 5' end (on top) or 3' end (on

1201    bottom) at that position. Only early junctions that are at least 1% of early early

1202    transcripts are plotted. The superT junction is colored in gold. (**A**) Intron plot for BKPyV

1203    Dik WT. (**B**) Intron plot for BKPyV Dik M1. (**C**) Intron plot for BKPyV Dik M2.

1204

58

1205    **Figure S12** - superT in MCPyV-associated MCC

1206    A. Sanger sequencing of an RT-PCR product from MCC J45_440, showing the

1207       superT-specific junction.

1208    B. A schematic detailing the MCC 285 MCPyV integration site, showing how it is

1209       possible that superT is generated via cis-splicing.

1210    C. The assembled viral block in MCC tumor J45_440. This integration site is based

1211       on de-novo assembly using short whole genome sequencing reads. Despite only

1212       assembling one viral block, we found that 1) there are likely 2 copies of the viral

1213       genome, and 2) the 5' viral integration site appears to fall on chromosome 7

1214       "after" the 3' viral integration site, observations consistent with the existence of

1215       two copies of the viral genome in tandem separated by a small segment of host

1216       DNA at this integration site.

1217    D. The assembled viral blocks in MCC tumor J17_296. The longest block contains

1218       two copies of the early region.

1219    E. Lollipop plots showing identified SNPs in the MCPyV genomes of J45_440,

1220       J17_296, and J11_285. The gene-map below the figure indicates the position of

1221       viral ORFs. Each lollipop is colored according to the nucleotide substitution

1222       identified.

1223

1224

1225  Abend, J.R., Joseph, A.E., Das, D., Campbell-Cecen, D.B., and Imperiale, M.J. (2009).

1226  A truncated T antigen expressed from an alternatively spliced BK virus early mRNA.

1227  The Journal of general virology *90*, 1238.

1228  Abere, B., Zhou, H., Li, J., Cao, S., Toptan, T., Grundhoff, A., Fischer, N., Moore, P.S.,

1229  and Chang, Y. (2020). Merkel Cell Polyomavirus Encodes Circular RNAs (circRNAs)

1230  Enabling a Dynamic circRNA/microRNA/mRNA Regulatory Network. Mbio *11*, e03059-

1231  03020.

1232  Adami, G., Marlor, C., Barrett, N., and Carmichael, G.G. (1989). Leader-to-leader

1233  splicing is required for efficient production and accumulation of polyomavirus late

1234  mRNAs. Journal of virology *63*, 85-93.

1235  Assetta, B., De Cecco, M., O'Hara, B., and Atwood, W.J. (2016). JC polyomavirus

1236  infection of primary human renal epithelial cells is controlled by a type I IFN-induced

1237  response. MBio *7*, e00903-00916.

1238  Balázs, Z., Tombácz, D., Szűcs, A., Snyder, M., and Boldogkői, Z. (2017). Long-read

1239  sequencing of the human cytomegalovirus transcriptome with the Pacific Biosciences

1240  RSII platform. Scientific data *4*, 1-7.

1241  Carter, J.J., Daugherty, M.D., Qi, X., Bheda-Malge, A., Wipf, G.C., Robinson, K.,

1242  Roman, A., Malik, H.S., and Galloway, D.A. (2013). Identification of an overprinting

1243  gene in Merkel cell polyomavirus provides evolutionary insight into the birth of viral

1244  genes. Proceedings of the National Academy of Sciences *110*, 12744-12749.

1245   Depledge, D.P., Srinivas, K.P., Sadaoka, T., Bready, D., Mori, Y., Placantonakis, D.G.,

1246   Mohr, I., and Wilson, A.C. (2019). Direct RNA sequencing on nanopore arrays redefines

1247   the transcriptional complexity of a viral pathogen. Nature communications *10*, 1-13.

1248   Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P.,

1249   Chaisson, M., and Gingeras, T.R. (2013). STAR: ultrafast universal RNA-seq aligner.

1250   Bioinformatics *29*, 15-21.

1251   Eul, J., and Patzel, V. (2013). Homologous SV40 RNA trans-splicing: a new mechanism

1252   for diversification of viral sequences and phenotypes. RNA biology *10*, 1689-1699.

1253   Freund, R., Sotnikov, A., Bronson, R.T., and Benjamin, T.L. (1992). Polyoma virus

1254   middle T is essential for virus replication and persistence as well as for tumor induction

1255   in mice. Virology *191*, 716-723.

1256   Garalde, D.R., Snell, E.A., Jachimowicz, D., Sipos, B., Lloyd, J.H., Bruce, M., Pantic, N.,

1257   Admassu, T., James, P., and Warland, A. (2018). Highly parallel direct RNA sequencing

1258   on an array of nanopores. Nature methods *15*, 201-206.

1259   Garren, S.B., Kondaveeti, Y., Duff, M.O., and Carmichael, G.G. (2015). Global analysis

1260   of mouse polyomavirus infection reveals dynamic regulation of viral and host gene

1261   expression and promiscuous viral RNA editing. PLoS pathogens *11*, e1005166.

1262   Ghosh, P., Reddy, V., Swinscoe, J., Lebowitz, P., and Weissman, S. (1978).

1263   Heterogeneity and 5′-terminal structures of the late RNAs of simian virus 40. Journal of

1264   molecular biology *126*, 813-846.

1265    Good, P.J., Welch, R.C., Ryu, W.-S., and Mertz, J.E. (1988). The late spliced 19S and

1266    16S RNAs of simian virus 40 can be synthesized from a common pool of transcripts.

1267    Journal of virology *62*, 563-571.

1268    Gottlieb, K.A., and Villarreal, L.P. (2001). Natural biology of polyomavirus middle T

1269    antigen. Microbiology and molecular biology reviews *65*, 288-318.

1270    Jiang, M., Abend, J.R., Johnson, S.F., and Imperiale, M.J. (2009a). The role of

1271    polyomaviruses in human disease. Virology *384*, 266-273.

1272    Jiang, M., Abend, J.R., Tsai, B., and Imperiale, M.J. (2009b). Early events during BK

1273    virus entry and disassembly. Journal of virology *83*, 1350-1358.

1274    Johnson, M., Zaretskaya, I., Raytselis, Y., Merezhuk, Y., McGinnis, S., and Madden,

1275    T.L. (2008). NCBI BLAST: a better web interface. Nucleic acids research *36*, W5-W9.

1276    Kamen, R., Favaloro, J., and Parker, J. (1980a). Topography of the three late mRNA's

1277    of polyoma virus which encode the virion proteins. Journal of virology *33*, 637-651.

1278    Kamen, R., Favaloro, J., Parker, J., Treisman, R., Lania, L., Fried, M., and Mellor, A.

1279    (1980b). Comparison of polyoma virus transcription in productively infected mouse cells

1280    and transformed rodent cell lines. Paper presented at: Cold Spring Harbor symposia on

1281    quantitative biology (Cold Spring Harbor Laboratory Press).

1282    Keller, M.W., Rambo-Martin, B.L., Wilson, M.M., Ridenour, C.A., Shepard, S.S., Stark,

1283    T.J., Neuhaus, E.B., Dugan, V.G., Wentworth, D.E., and Barnes, J.R. (2018). Direct

1284    RNA sequencing of the coding complete influenza A virus genome. Scientific reports *8*,

1285    1-8.

1286 Kim, D., Lee, J.-Y., Yang, J.-S., Kim, J.W., Kim, V.N., and Chang, H. (2020). The

1287 architecture of SARS-CoV-2 transcriptome. Cell *181*, 914-921. e910.

1288 Kress, M., May, E., Cassingena, R., and May, P. (1979). Simian virus 40-transformed

1289 cells express new species of proteins precipitable by anti-simian virus 40 tumor serum.

1290 Journal of virology *31*, 472-483.

1291 Krueger, F. (2016). Babraham Bioinformatics-Trim Galore.

1292 Lee, S., Micalizzi, D., Truesdell, S.S., Bukhari, S.I., Boukhali, M., Lombardi-Story, J.,

1293 Kato, Y., Choo, M.-K., Dey-Guha, I., and Ji, F. (2020). A post-transcriptional program of

1294 chemoresistance by AU-rich elements and TTP in quiescent leukemic cells. Genome

1295 biology *21*, 1-23.

1296 Legon, S., Flavell, A.J., Cowie, A., and Kamen, R. (1979). Amplification in the leader

1297 sequence of late polyoma virus mRNAs. Cell *16*, 373-388.

1298 Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics

1299 *34*, 3094-3100.

1300 Loman, N.J., Quick, J., and Simpson, J.T. (2015). A complete bacterial genome

1301 assembled de novo using only nanopore sequencing data. Nature methods *12*, 733-

1302 735.

1303 Luo, Y., and Carmichael, G.G. (1991). Splice site skipping in polyomavirus late pre-

1304 mRNA processing. Journal of virology *65*, 6637-6644.

1305 Nguyen, K.D., Lee, E.E., Yue, Y., Stork, J., Pock, L., North, J.P., Vandergriff, T.,

1306 Cockerell, C., Hosler, G.A., and Pastrana, D.V. (2017). Human polyomavirus 6 and 7

1307    are associated with pruritic and dyskeratotic dermatoses. Journal of the American

1308    Academy of Dermatology *76*, 932-940. e933.

1309    Nomburg, J., Meyerson, M., and DeCaprio, J.A. (2020). Pervasive generation of non-

1310    canonical subgenomic RNAs by SARS-CoV-2. Genome medicine *12*, 1-14.

1311    Norbury, C.J., and Fried, M. (1987). Polyomavirus early region alternative poly (A) site:

1312    3'-end heterogeneity and altered splicing pattern. Journal of virology *61*, 3754-3758.

1313    Price, A.M., Hayer, K.E., McIntyre, A.B., Gokhale, N.S., Abebe, J.S., Della Fera, A.N.,

1314    Mason, C.E., Horner, S.M., Wilson, A.C., and Depledge, D.P. (2020). Direct RNA

1315    sequencing reveals m 6 A modifications on adenovirus RNA are necessary for efficient

1316    splicing. Nature communications *11*, 1-17.

1317    Quinlan, A.R., and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for comparing

1318    genomic features. Bioinformatics *26*, 841-842.

1319    Reddy, V.B., Ghosh, P.K., Lebowitz, P., and Sherman M, W. (1978). Gaps and

1320    duplicated sequences in the leaders of SV40 16S RNA. Nucleic acids research *5*, 4195-

1321    4214.

1322    Riley, M.I., Yoo, W., Mda, N.Y., and Folk, W.R. (1997). Tiny T antigen: an autonomous

1323    polyomavirus T antigen amino-terminal domain. Journal of virology *71*, 6068-6074.

1324    Rosenstein, R.K., Pastrana, D.V., Starrett, G.J., Sapio, M.R., Hill, N.T., Jo, J.-H., Lee,

1325    C.-C.R., Iadarola, M.J., Buck, C.B., and Kong, H.H. (2021). Host-Pathogen Interactions

1326    in Human Polyomavirus 7–Associated Pruritic Skin Eruption. The Journal of

1327    investigative dermatology *141*, 1344-1348. e1348.

1328  Saribas, A.S., DeVoto, J., Golla, A., Wollebo, H.S., White, M.K., and Safak, M. (2018).

1329  Discovery and characterization of novel trans-spliced products of human polyoma JC

1330  virus late transcripts from PML patients. Journal of cellular physiology *233*, 4137-4155.

1331  Schmidlin, K., Kraberger, S., Cook, C., DeNardo, D.F., Fontenele, R.S., Van Doorslaer,

1332  K., Martin, D.P., Buck, C.B., and Varsani, A. (2021). A novel lineage of polyomaviruses

1333  identified in bark scorpions. bioRxiv.

1334  Seif, I., Khoury, G., and Dhar, R. (1979). The genome of human papovavirus BKV. Cell

1335  *18*, 963-977.

1336  Shuda, M., Feng, H., Kwun, H.J., Rosen, S.T., Gjoerup, O., Moore, P.S., and Chang, Y.

1337  (2008). T antigen mutations are a human tumor-specific signature for Merkel cell

1338  polyomavirus. Proceedings of the National Academy of Sciences *105*, 16272-16277.

1339  Smith, A.E., Smith, R., and Paucha, E. (1979). Characterization of different tumor

1340  antigens present in cells transformed by simian virus 40. Cell *18*, 335-346.

1341  Starrett, G.J., Thakuria, M., Chen, T., Marcelus, C., Cheng, J., Nomburg, J., Thorner,

1342  A.R., Slevin, M.K., Powers, W., and Burns, R.T. (2020). Clinical and molecular

1343  characterization of virus-positive and virus-negative Merkel cell carcinoma. Genome

1344  medicine *12*, 1-22.

1345  Starrett, G.J., Yu, K., Golubeva, Y., Lenz, P., Piaskowski, M.L., Peterson, D., Dean, M.,

1346  Israni, A., Hernandez, B.Y., Tucker, T.C.*, et al.* (2021). Common Mechanisms of Virus-

1347  Mediated Oncogenesis in Bladder Cancers Arising In Solid Organ Transplant

1348  Recipients. medRxiv, 2021.2011.2011.21266080.

1349    Theiss, J.M., Günther, T., Alawi, M., Neumann, F., Tessmer, U., Fischer, N., and

1350    Grundhoff, A. (2015). A comprehensive analysis of replicating Merkel cell polyomavirus

1351    genomes delineates the viral transcription program and suggests a role for mcv-miR-M1

1352    in episomal persistence. PLoS pathogens *11*, e1004974.

1353    Treisman, R. (1980). Characterisation of polyoma late mRNA leader sequences by

1354    molecular cloning and DNA sequence analysis. Nucleic acids research *8*, 4867-4888.

1355    Tremblay, J.D., Sachsenmeier, K.F., and Pipas, J.M. (2001). Propagation of wild-type

1356    and mutant SV40. In SV40 Protocols (Springer), pp. 1-7.

1357    Trowbridge, P.W., and Frisque, R.J. (1995). Identification of three new JC virus proteins

1358    generated by alternative splicing of the early viral mRNA. Journal of neurovirology *1*,

1359    195-206.

1360    Yang, R., Lee, E.E., Kim, J., Choi, J.H., Kolitz, E., Chen, Y., Crewe, C., Salisbury, N.J.,

1361    Scherer, P.E., and Cockerell, C. (2021). Characterization of ALTO-encoding circular

1362    RNAs expressed by Merkel cell polyomavirus and trichodysplasia spinulosa

1363    polyomavirus. PLoS pathogens *17*, e1009582.

1364    Zerrahn, J., Knippschild, U., Winkler, T., and Deppert, W. (1993). Independent

1365    expression of the transforming amino-terminal domain of SV40 large I antigen from an

1366    alternatively spliced third SV40 early mRNA. The EMBO journal *12*, 4739-4746.

1367    Zhao, L., and Imperiale, M.J. (2019). Establishing Renal Proximal Tubule Epithelial-

1368    Derived Cell Lines Expressing Human Telomerase Reverse Transcriptase for Studying

1369    BK Polyomavirus. Microbiology resource announcements *8*, e01129-01119.

1370    Zou, W., Vue, G.S., Assetta, B., Manza, H., Atwood, W.J., and Imperiale, M.J. (2020).

1371    Control of archetype BK polyomavirus microRNA expression. Journal of Virology *95*,

1372    e01589-01520.

1373

1374

Figure 1



Figure 1 - RNA sequencing expands known SV40 and BKPyV transcript diversity.

A. Overview of experimental procedures. Cells were infected with a polyomavirus, and RNAs extracted. RNA was sequenced using long-read (Nanopore dRNAseq and PacBio SMRTseq) and short-read (Illumina short-RNAseq (total) and short-RNAseq (polyA)). Transcripts were analyzed, and the impact of observed splice events on viral open reading frames was assessed.

B. Mechanism of transcript clustering in this study. Transcripts were aligned to the viral genome and grouped into transcript classes based on the presence of shared introns. Thus, within a transcript class there may be variation in the exact transcript start and end positions. This clustering strategy was used for both long- and short-RNAseq data.

C. Viral RNA sequence coverage for SV40 and BKPyV as determined from dRNAseq, SMRTseq, and short-RNAseq (total) data. The Y axis indicates the scaled coverage, with X axis indicating the position on the viral genome. Coverage for late transcripts (mapping to the + strand) is above the x axis, while coverage for early transcripts (mapping to the - strand) is below the x axis. Coverage is scaled separately for each strand such that the maximum observed coverage for each strand is 1. Arrows at the top of the plot indicate the positions of viral genes.

D-E. UpSet plot indicating the overlap between existing transcript annotations, dRNAseq data, and SMRTseq data for SV40 (D) and BKPyV Dunlop (E). Bars indicating overlap with existing transcript annotations are black, while those indicating no overlap with existing annotations are blue. These blue bars indicate the number of novel, unannotated transcripts identified.

F. Overview of polysome profiling of SV40-infected cells. BSC40 cells were infected with SV40. Cells were lysed, and a portion of the lysate was subjected to dRNAseq (representative of the RNA content of the whole cell). The remaining lysates was centrifuged through a sucrose gradient, after which fractions containing RNA associated with two or more ribosomes were pooled and subjected to dRNAseq.

G. Relative abundance of SV40 early and late transcripts in the whole-cell and polysome fractions of SV40-infected cells. Y-axis indicates the percentage of early or late transcripts and is log scale. X axis indicates each transcript, with black dots indicating each transcript's whole-cell relative abundance and red dots indicating each transcript's polysome relative abundance.

## Figure 2



Figure 2 - Annotated and novel SV40 transcripts.
A. Transcripts are shown relative to the viral genome. Each line is a viral transcript, with red lines indicating exons and dashed blue lines indicating introns. Spokes indicate the positions of common splice donors and splice acceptors. Transcripts that were annotated prior to this study are on a yellow background, and novel transcripts are on a while background. Wraparound transcription that results in multiple copies of a region is annotated with double lines, and the number of copies is indicated in parentheses. The line labeled "pA" indicates the approximate position of the polyA signal sequence.

# Figure 3



**Figure 3 - Annotated and novel BKPyV transcripts.**

A. Transcripts are shown relative to the viral genome. Each line is a viral transcript, with red lines indicating exons and dashed blue lines indicating introns. Spokes indicate the positions of common splice donors and splice acceptors. Transcripts that were annotated prior to this study are on a yellow background, and novel transcripts are on a while background. Wraparound transcription that results in multiple copies of a region is annotated with double lines, and the number of copies is indicated in parentheses. The line labeled "pA" indicates the approximate position of the polyA signal sequence.

**Figure 4 - Pervasive wraparound transcription across PyV**

A-C. Watch plots indicating the top 4 highest abundance late wraparound transcript classes in dRNAseq data from SV40 (A), BKPyV Dunlop (B), and MPyV (C). The outer ring of each watch plot indicates the position of the viral ORFs. The inner arms are histograms detailing the distribution of transcript starts (in blue) and ends (in red) for transcripts within each transcript class. The red segments indicate exons. Transcripts start in the innermost ring - a second or third ring indicates that the pre-mRNA that generated the transcript must have circled the viral genome multiple times. The 3' end of the transcript and the direction in which these plots are oriented is indicated by the red arrow at the end of the last exon segment. The red exon segments start at the most common transcript start site within the transcript class, and end at the most common transcript end site within the class. The watch plot key shows an example of the path of the pre-mRNA for SV40 transcript class L6_I.

D. Bar plots indicating the percentage of late transcripts that span a given number of genome lengths in SV40, BKPyV Dunlop, and MPyV dRNAseq data.

E. The leader-leader junction, that connects the pre-mRNA from one genome to the subsequent wraparound, was identified in Illumina short-RNAseq (total) data. The intron in question is plotted as a black line in this plot, with the x axis indicating the genomic position of the intron. The top late wraparound transcript for each virus was plotted. The gene map indicates the approximate gene position and is accurate for SV40 - the exact position of the viral genes varies between viruses. Percentages indicate the percentage of late junction-spanning transcripts that support the plotted wraparound leader-leader junction.

F. Schematic illustrating how leader-leader wraparound transcription can be detected from short read short-RNAseq (total). Leader-leader splicing can be seen as a repetitive exon in watch plots from long-read RNAseq data. Ultimately, there was an original processed mRNA in the cell that contained two tandem leader sequences. When this transcript of origin is sequenced via short read sequencing, reads will be generated across its length. A minority of these reads will span the leader-leader junction, and mapping against the viral reference genome can be used to uncover leader-leader splicing.

**Figure 5 - Detection of novel, conserved splicing events that expand PyV coding capacity.**

A-D. Schematics illustrating identified ORFs. Each row is a reading frame (except for ST and the LT 1st exon, which are in the same frame), and unannotated amino acids are represented by grey boxes. The measured intron is indicated by the red arrow. Colored ORFs are annotated, while grey ORFs are unannotated. Percentages on the right side of the figure are the percentage of spliced viral transcripts on the same strand as determined from short-read short-RNAseq (total) data. Numbers after each virus name indicate the transcript class within each short-RNAseq (total) dataset. The measured intron is indicated by the red arrow.

A) ST2: This ORF is generating from a splicing event that uses the LT first exon donor and an acceptor within the ST ORF. In HPyV7 and BKPyV Dunlop, the splice lands in frame and results in an internal deletion within ST. In MPyV and MCPyV the splice lands out of frame, resulting in an ORF that contains the N-terminal region of ST and novel amino acids at the C terminus.

B) MT: MPyV encodes a MT following splicing connecting the end of the ST ORF with an ORF in an alternate frame of the LT second exon. In BKPyV, a similar splice occurs connecting ST with an MT-like ORF in an alternative frame of the LT second exon.

C) VP1X: JCPyV encodes two VP1X ORFs generated by splicing within VP1 and landing in an alternative frame of VP1, or earlier in the late region due to wraparound transcription. While predominant in JCPyV, VP1X is likewise present in many other PyV.

D) superT: The superT-specific splice utilizes the splice donor canonically associated with truncated T antigens such as 17kT in SV40 and truncT in BKPyV. Due to wraparound transcription, a LT second exon acceptor is available to the 3' of this donor and acts as the acceptor. For the superT ORF to form, an initial LT splice is required. Ultimately, superT contains a duplication in part of the LT second exon that includes the RB-binding LxCxE motif.

E. Schematics detailing BKPyV Dik isolates used for querying the existence of superT. BKPyV WT is wild type virus. M1 contains a LT intron that has been replaced with an intron from the plasmid pCI. Both WT and M1 are expected to generate LT and superT of expected sizes. M2 has a completely removed LT intron, and the pCI intron is located directly 5' of the LT ORF. M2 is expected to encode LT of expected size, but a larger superT variant due to incorporation of a second copy of the LT first exon.

F. Western blot of cells infected with BKPyV Dik WT, M1, or M2 and probed with an antibody reactive against LT. The lower molecular weight band is LT, and the higher molecular weight bands are consistent with superT.

**Figure 6**



Figure 6 - Detection superT-encoding transcripts in PyV-associated cancers

A. Schematic detailing the generation of superT during lytic infection as compared to from integrated virus in cancer. During viral infection, the RNA polymerase can circle the viral genome multiple times, resulting in a pre-mRNA that can be spliced to generate superT. In the case of host integration, a polyomavirus can be integrated in tandem copies such that a pre-mRNA is generated with more than one copy of the viral early region. This pre-mRNA can be similarly spliced to generate a superT transcript. B. Heatmap indicating the abundance of the superT, ST, and LT introns from RNAseq data from two replicates of a BKPyV-positive bladder cancer and six MCPyV-associated MCCs. Percentages indicate the percentage of spliced early viral reads for each sample. The splice measured in each row is indicated by the red arrow in the schematics on the right side of the figure.

Fig S1

Figure S2

FigS3: SV40 Transcriptome

FigS4: BKPyV (Dunlop) Transcriptome

Fig S5: SV40 Late Transcript Pileups

## Fig S6: SV40 Early Transcript Pileups

Fig S7: BKPyV (Dunlop) Late Transcript Pileups

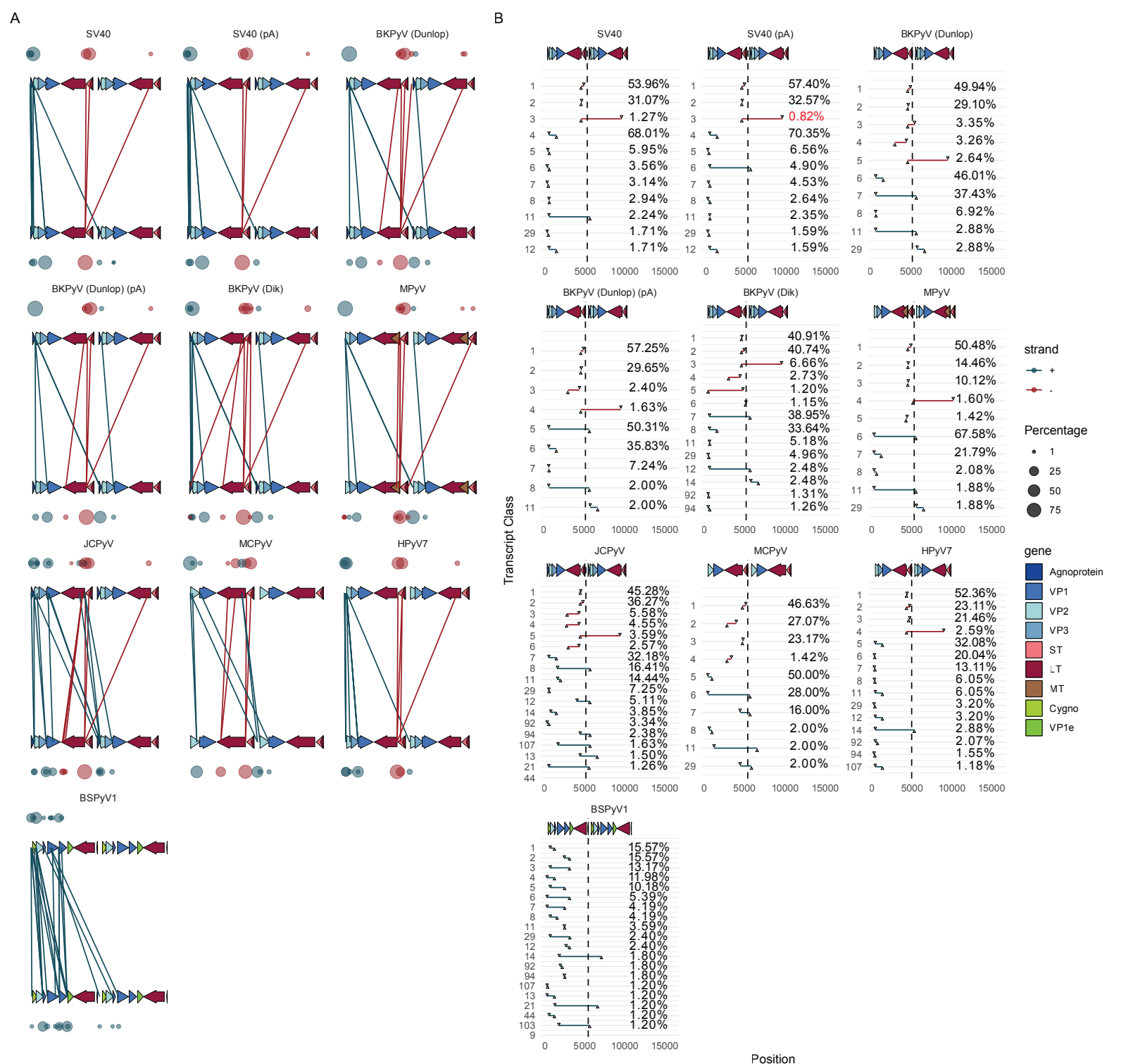Fig S8: BKPyV (Dunlop) Early Transcript Pileups

Fig S9

Illumina Introns >1% of Junction-spanning reads (by strand)

Figure S10

Fig S11

BKPyV Dik short-RNAseq (polyA)

Fig S12