



Universiteit
Leiden
The Netherlands

The trichodysplasia spinulosa-associated polyomavirus : infection, pathogenesis, evolution and adaptation

Kazem, S.

Citation

Kazem, S. (2015, June 17). *The trichodysplasia spinulosa-associated polyomavirus : infection, pathogenesis, evolution and adaptation*. Retrieved from <https://hdl.handle.net/1887/33437>

Version: Corrected Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/33437>

Note: To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/33437> holds various files of this Leiden University dissertation.

Author: Kazem, Siamaque

Title: The trichodysplasia spinulosa-associated polyomavirus : infection, pathogenesis, evolution and adaptation

Issue Date: 2015-06-17

VePyV1 CaPyV JCPyV *mPyV* SA12 RacPyV GHPyV BKPyV

BatSgV CPyV OraPyV1 FPyV MptV APPyV1 SqPyV LPyV CSLPyV

EPyV PRPyV1 APP_gV2 KIPyV **MiPyV** PtvPyV2c OtPyV1 **STLPyV**

MEPyV1 *JSgV09* SV40 TSPyV CoPyV1 PPPyV CPPyV HPyV12

MXPyV PtvPyV1a PDPyV EiPyV1 AtPPyV1 *HaPyV* (GGGPyV1

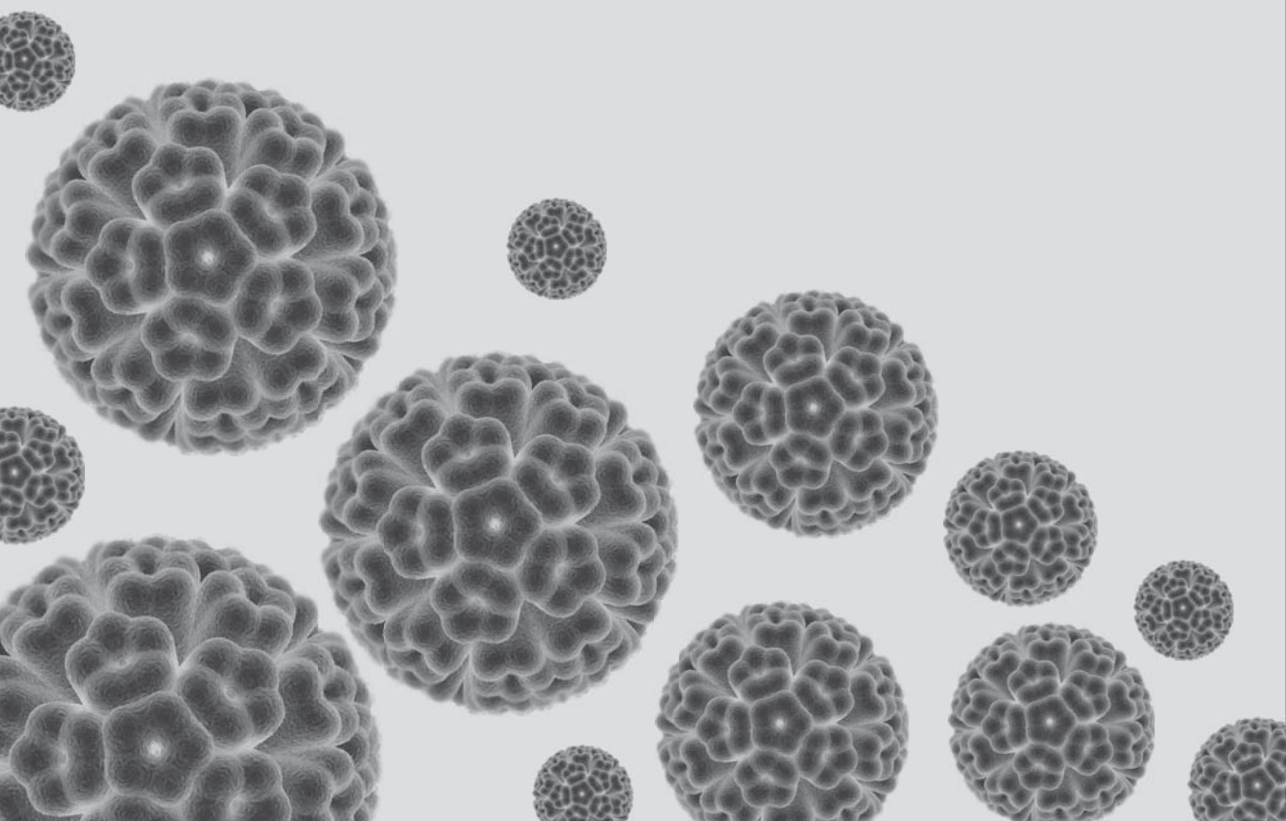
CdPyV DRPyV MWPyV APyV CAPyV1 HPyV7 ChPyV MasPyV

WUPyV *HEPyV6* BPyV MCPyV OraPyV2 MMPyV SLPyV HPyV10

VePyV1 CaPyV JCPyV mPyV SA12
BatPyV CPyV OraPyV1 FPyV MptV APPyV
EPyV PRPyV1 APPyV2 KIPyV MiPyV
MFPyV1 HPyV9 SV40 TSPyV CoPyV1
MXPyV PtvPyV1a PDPyV EiPyV1 At
CdPyV DRPyV MWPyV APyV CAPyV1
WUPyV HPyV6 BPyV MCPyV OraPyV2

Part III

TSPyV host adaptation and evolution



Chapter 5

Polyomavirus host adaptation

Interspecific adaptation by binary choice at *de novo* polyomavirus T antigen site through accelerated codon-constrained Val-Ala toggling within an intrinsically disordered region

Authors:

Chris Lauber^{1, 2*}

Siamaque Kazem^{1*}

Alexander Kravchenko^{3†}

Mariet Feltkamp¹

Alexander Gorbalenya^{1, 3, 4}

** These authors contributed equally to the paper as first authors*

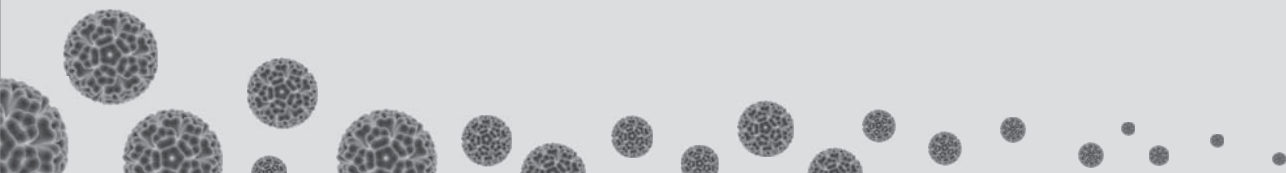
† Deceased; of blessed memory

Affiliations:

¹ Department of Medical Microbiology, Leiden University Medical Center, Leiden, The Netherlands, ² Institute for Medical Informatics and Biometry, Technische Universität Dresden, Dresden, Germany, ³ Belozersky Institute of Physico-Chemical Biology, Lomonosov Moscow State University, Moscow, Russia and ⁴ Faculty of Bioengineering and Bioinformatics, Lomonosov Moscow State University, Moscow, Russia

Published in:

Nucleic Acids Research, 2015 (*in-press*)



Abstract

It is common knowledge that conserved residues evolve slowly. We challenge generality of this central tenet of molecular biology by describing the fast evolution of a conserved nucleotide position that is located in the overlap of two open reading frames (ORFs) of polyomaviruses. The *de novo* ORF is expressed through either the ALTO protein or the Middle T antigen (MT/ALTO), while the ancestral ORF encodes the N-terminal domain of helicase-containing Large T (LT) antigen. In the latter domain the conserved Cys codon of the LXCXE pRB-binding motif constrains codon evolution in the overlapping MT/ALTO ORF to a binary choice between Val and Ala codons, termed here as codon-constrained Val-Ala (COCO-VA) toggling. We found the rate of COCO-VA toggling to approach the speciation rate and to be significantly accelerated compared to the baseline rate of chance substitution in a large monophyletic lineage including all viruses encoding MT/ALTO and three others. Importantly, the COCO-VA site is located in a short linear motif (SLiM) of an intrinsically disordered region, a typical characteristic of adaptive responders. These findings provide evidence that the COCO-VA toggling is under positive selection in many polyomaviruses, implying its critical role in interspecific adaptation, which is unprecedented for conserved residues.

Abstract

Introduction

Intrinsically disordered regions (IDRs), which are either not structured or may become structured upon interaction with diverse partners [1], have been identified in many proteins and implicated in various biological processes as adaptive responders [2 - 5]. They have a biased amino acid residue composition and evolved faster than structured proteins [6, 7], with exception of very small islands of relative conservation, known as short linear motifs (SLiM), that mediate protein-protein interactions [8].

IDRs are frequently encoded by overlapping open reading frames (ORFs) that evolved *de novo* by overprinting the ancestral ORFs and are common in viruses [9 - 13]. This overlapping of ORFs is accompanied by suppression of synonymous substitution rate in the ancestral ORFs (negative or purifying selection) compared to that of non-overlapping ORFs, indicative of codon constraints in the *de novo* ORFs due to their expression. The observed phenomenon has been extensively used for *in silico* identification of functional *de novo* ORFs [12, 14, 15], which often led to the elucidation of non-canonical expression mechanisms of these ORFs (e.g., [16 - 19]). Suppression of synonymous substitution rate is also reciprocally imposed on *de novo* ORFs by the overlapping ancestral gene. These observations led to analysis of relative rate change of substitutions in the *de novo* genes compared to ancestral or non-overlapping genes [12, 20 - 23].

This ORF-wide analysis has not been extended to individual codons of the *de novo* ORFs due to formidable technical challenges. A common approach to characterize site-specific evolution is to estimate deviation from the substitution rate under a model of neutral evolution for each codon of an ORF. Suppression and acceleration of the substitution rate is attributed to negative and positive selection, respectively, with positive selection being seen as the hallmark signature of adaptation during intra-species evolution [24]. One particular pattern of variation under positive selection is the frequent exchange of residues with pervasive return to the wild-type state, dubbed residue toggling [25]. Identification of codons under selection, either negative or positive, is part of the established evolutionary-based pipeline that informs functional characterization of proteins encoded in non-overlapping ORFs [2, 27]. However, the available techniques were not developed to untangle selection forces acting on the overlapping ORFs, which constrain evolution of each other. This may explain the lack of identification of *de novo* codon(s) under positive selection, despite broad recognition of a prominent role that the overlapping ORFs play in adaptation of viruses to host [12].

One of the largest and poorly characterized pairs of proteins encoded by overlapping ORFs is expressed by members of the fast growing *Polyomaviridae* family (**Supplementary Table S1**). These viruses cause latent infections in diverse mammals and birds, and in humans, some of these viruses have been responsible for different pathologies in immunocompromised individuals [28, 29]. Polyomaviruses employ multi-ORF double-stranded DNA (dsDNA) genomes of approximately 5 kb [30, 31]. Genomes of a large subset of polyoma-

viruses include two overlapping ORFs [15, 32, 33], designated here ORF2 and ORF5 (**Figure 1A**; for other designations see **Supplementary Text S1** and **Supplementary Table S2**). ORF2 encodes the second exon of the large T antigen (LT) that includes a helicase domain [30, 34]. ORF5 is expressed as a separate protein (ALTO) in Merkel cell polyomavirus (MCPyV) [15]; while it encodes the second exon of Middle T antigen (MT) in murine and hamster polyomaviruses (MPyV and HaPyV) [33, 35, 36]. The ORF5-encoded part of MT antigen is implicated in control of cell transformation [33, 35, 36], enriched with Pro residues [37, 38] and includes a C-terminal transmembrane domain [35] that is essential for the oncogenic function of MT [39]. This function and interaction of MT with different cellular proteins may be modulated by phosphorylation at several Ser, Thr and Tyr residues in rodent polyomaviruses [36]. We will use ORF5-plus and ORF5-less to refer to respective subsets of polyomaviruses; ORF5-plus viruses are also known as *Almipolyomaviruses* [15]. Likewise, and purely for the sake of uniformity, hereafter we have designated the ORF5-encoding product as MT/ALTO for all ORF5-plus polyomaviruses. Because ORF5 is conserved in only ORF5-plus polyomaviruses, while the overlapping part of ORF2 is found in all mammalian polyomaviruses [15], these ORFs are defined as *de novo* (ORF5) and ancestral (ORF2), according to Sabath *et al.*, [12]. ORF5-plus viruses form a large monophyletic cluster in one of the main branches of polyomavirus tree [15], dubbed *Orthopolyomaviruses I* [Ortho-I]; with three other branches being *Orthopolyomaviruses II*, *Malawipolyomaviruses* and *Wukipolyomaviruses* [40], although branch delineation and designation may vary in different studies [15, 41].

To understand the evolution of overlapping ORFs, we studied ORF2 and ORF5 at codon resolution. We found that one of the most conserved ORF5 codons, located in a SLiM of ORF5, experienced an accelerated evolutionary rate despite being strongly constrained to two amino acids by the overlapping ancestral ORF2. Using available and specially developed evolutionary-based approaches we revealed an unprecedented frequent toggling between these two residues during large-scale multi-species evolution in the Ortho-I clade of polyomaviruses. This analysis is, to our knowledge, the first to identify a conserved position of *de novo* protein under positive selection. Its results suggest a new IDR-mediated adaptation mechanism employed by many mammalian polyomaviruses with potential relevance to understanding adaptation of other viruses and organisms.

Materials and Methods

Datasets: viruses, sequences and alignments

Full-length genome sequences of 55 polyomaviruses available in the Genbank/RefSeq database on February 2013 (**Supplementary Table S1**) were downloaded into the Viralis platform [42]. When several genomes per species were available, the RefSeq sequence was chosen for presentation. The Muscle program [43] and ClustalW [44] were used to generate family-wide multiple amino-acid alignments for viral capsid protein (VP)1 encoded in ORF3, VP2 (ORF4) and LT (ORF2), followed by manual curation. For each of the three protein

alignments, strongly conserved blocks [45] were extracted using the Blocks Accepting Gaps Generator (BAGG) tool (www.genebee.msu.su/~antonov/bagg/cgi/bagg.cgi) to produce a concatenated multiple sequence alignment used for phylogenetic reconstruction and other analyses (see below). The ORF2-wide alignment was also mapped on the genome sequences, which were then translated in the alternative reading frame (RF -3) encoding ORF5 in twenty-two viruses of the ORF5-plus group to produce an ORF5 alignment. ORF5 size varies from 441 nucleotides (nts) to 846 nts, and ORF5 sequence conservation was detectable only in some subsets of polyomaviruses (**Supplementary Table S3**; data not shown; [15]).

For analysis of site-specific evolutionary selection by Datamonkey programs, we used 10 alignments of selected positions of the ORF5 and ORF2 (datasets, D1-D10). These 10 alignments represented different groups of viruses, including all mammalian polyomaviruses (D1 and D2), Ortho-I viruses (D3 and D4), ORF5-plus viruses (D5 and D6), ORF5-less viruses (D7), and three non-overlapping lineages of ORF5-plus viruses (D8-D10), each analyzed separately (see **Supplementary Table S4** for details). Using conservation considerations, some codons of ORF5 and ORF2 were selected, so all datasets included ORF5 codons while D2, D4, D6 and D7 included also ORF2 codons. For ORF5 of D1-D7, those codons were chosen whose overlapping codon in ORF2 (-1 frame) was aligned with no gaps across mammalian polyomaviruses. For ORF5 of D8-D10 and ORF2, most conserved codons in respective alignments were used after manual pruning of weakly aligned codons.

Alignments of the conserved motifs in the N-terminal part of LT ORF2 and ORF5, partially described elsewhere [15], were produced and converted into logos. To produce alignments as input for the RNAz program, we converted codon ORF5-based alignments of four subsets of ORF5-plus and two subsets of ORF5-less polyomaviruses, into the respective nucleotide alignments (**Supplementary Table S3**).

Phylogeny reconstruction

Phylogenetic analyses were performed by using a Bayesian approach implemented in BEAST version 1.7.4 [46] and the Whelan and Goldman (WAG) amino acid substitution matrix [47]. Rate heterogeneity among sites was modeled using a gamma distribution with four categories, and a relaxed molecular-clock approach was tested against the strict molecular-clock approach [48] and was found to be superior. Markov chain Monte Carlo (MCMC) chains were run for 2 million steps and the first 10% were discarded as burn-in. Convergence of the runs was verified using the Tracer tool (<http://beast.bio.ed.ac.uk/tracer>).

Analysis of natural selection at codons

We have used Mixed Effects Model of Evolution (MEME) [49] and Fast, Unconstrained Bayesian AppRoximation (FUBAR) [27] at the Datamonkey website (<http://www.datamonkey.org>) [26] to test for natural selection at conserved ORF5 codons. In addition, we have screened for toggling at ORF5 residues using TOGGLE, an implementation of the residue toggling method developed for HIV-1 by Delport *et al.*, [25]. We have analyzed in total ten different datasets, D1-D10 (see above), capturing different positions and virus diversities

(see **Supplementary Table S4**). For each analyzed data set, selection of evolutionary model was performed automatically at the Datamonkey web site using default parameters prior to the analysis.

Analysis of COCO-VA toggling by BayesTraits

Evolution of non-synonymous replacements at the **Codon-Constrained Val-Ala** (COCO-VA) site of ORF5 was analyzed by BayesTraits package using the Multistate model (<http://www.evolution.rdg.ac.uk/BayesTraits.html>) [50]. This codon is constrained to encode either Ala or Val in all mammalian polyomaviruses due to the overlapping Cys codon of the LXCXE motif that is expressed in the LT ORF2 of these viruses. The analyzed polyomaviruses were divided into two groups based on whether or not they express ORF5: ORF5-plus and ORF5-less viruses, respectively. The COCO-VA site is expressed as part of ORF5 in ORF5-plus, but not in ORF5-less viruses.

To test whether Ala-Val trait transitions are statistically more frequent in the ORF5-plus lineage compared to ORF5-less viruses, we applied the BayesTraits multistate model using a single trait (Ala/Val). We ran the analysis for three virus datasets: the combined set of mammalian polyomaviruses as well as separately for ORF5-plus and ORF5-less viruses, with respective posterior tree samples obtained through independent BEAST analyses. We then compared the estimated Ala-to-Val and Val-to-Ala transition rates between the three datasets, including an average Ala-Val exchange rate (corresponding to the toggling rate) by plotting the distributions. Statistical significance of differences in Ala-Val exchange rates was assessed using log Bayes Factors that was calculated with the R package Bayes Factor (<http://bayesfactorppl.r-forge.r-project.org/>). As Ala-Val exchange is equivalent to T-C exchange at the second codon position of the COCO-VA codon (see “Results and Discussion”), we applied BayesTraits also to the third position of that codon as a control.

Statistical analyses of COCO-VA toggling using patristic distances

For each virus the smallest pair-wise patristic distance (SPAT) to a virus encoding the same amino acid (monomorphic pairs: Ala \leftrightarrow Ala and Val \leftrightarrow Val; monoSPAT) and to that encoding the different amino acid (polymorphic pair: Ala \leftrightarrow Val, polySPAT) was calculated. Patristic distances were extracted from the polyomavirus phylogeny using the package Analyses of Phylogenetics and Evolution (APE) in R language [51].

We estimated the rate of COCO-VA toggling as the ratio of monoSPAT to the sum of polySPAT and monoSPAT values; designated SPAT ratio hereafter. Due to limited virus sampling at the intra-species level, we applied a sliding window approach to compare SPAT ratios between ORF5-plus and ORF5-less viruses. A window size of 0.15 and a shift of 0.05 at the monoSPAT scale were used. A two-sample non-parametric Mann-Whitney U test was utilized to test for statistically significant differences between the two virus groups within a particular window. A deviation of distributions of SPAT ratios from the average toggling rate of 0.5 was assessed using the Wilcoxon rank-sum test.

To independently assess the partitioning of mammalian polyomaviruses into ORF5-plus and ORF5-less virus groups, we determined the ranking of a predefined two-set partitioning among all possible two-set partitioning of the same type for the 30 viruses with monoSPAT values smaller than the derived threshold of 0.35. These 30 viruses comprise 14 ORF5-plus and 16 ORF5-less viruses or 16 Ortho-I and 14 non-Ortho-I viruses. We calculated the difference of mean toggling rate values between the two groups in each of these partitionings and determined its ranking among the differences of mean toggling rates obtained for all other 14-16 or 16-14 partitioning of the 30 viruses, whose total was 145,422,675 possible partitionings (e.g. combinations).

General bioinformatics analyses

For selected phylogenetic lineages, alignments of ORF5 were converted HMM profiles and compared to each other using HHsearch [52] in both local and global alignment modes.

Sequence logos of selected alignments were produced using the WebLogo server [53, 54].

Secondary structure and disorder prediction of protein sequences were generated using the Disorder Prediction MetaServer, which reports consensus results of eight protein disorder predictor tools: DISEMBL [55], DISOPRED [56], DISpro [57], FoldIndex [58], GlobPlot2 [59], IUPred [60], RONN [61], and VSL2 [62], and two protein secondary structure predictor tools: PROFsec [63] and PSIPred [64] (<http://wwwnmr.cabm.rutgers.edu/bioinformatics/disorder/>). The prediction of disorder was considered significant if at least four predictors gave a hit.

Secondary RNA structures in ORF2/ORF5 overlapping region were predicted with the program RNAz in a region of about 300-900 bp flanking the region encoding LXCXE motif sequence [65]. The server uses an algorithm that detects thermodynamically stable and evolutionarily conserved RNA secondary structures in multiple RNA-sequence alignments on both RNA-strands, with number of sequences in alignments not exceeding six. If subsets were larger than six, they were reduced to a combination of six virus sequences. For structure prediction the default RNAz parameters of “Standard Analysis” were utilized, which scored in the overlapping windows of 120 alignment columns with step-size of 40 nucleotides (<http://rna.tbi.univie.ac.at/cgi-bin/RNAz.cgi?PAGE=1&TYPE=S>).

Proline enrichment in putative ORF5-encoded protein sequences was analyzed by use of a custom R script that counts Proline residues and visualizes the counts with respect to location in the protein sequence and premature stop codons in the case of ORF5-less viruses (www.R-project.org) [66].

Results and Discussion

Discovery of Codon-Constrained Val-Ala (COCO-VA) toggling in MT/ALTO

We were interested in understanding the evolution and function of the *de novo* ORF5. Only

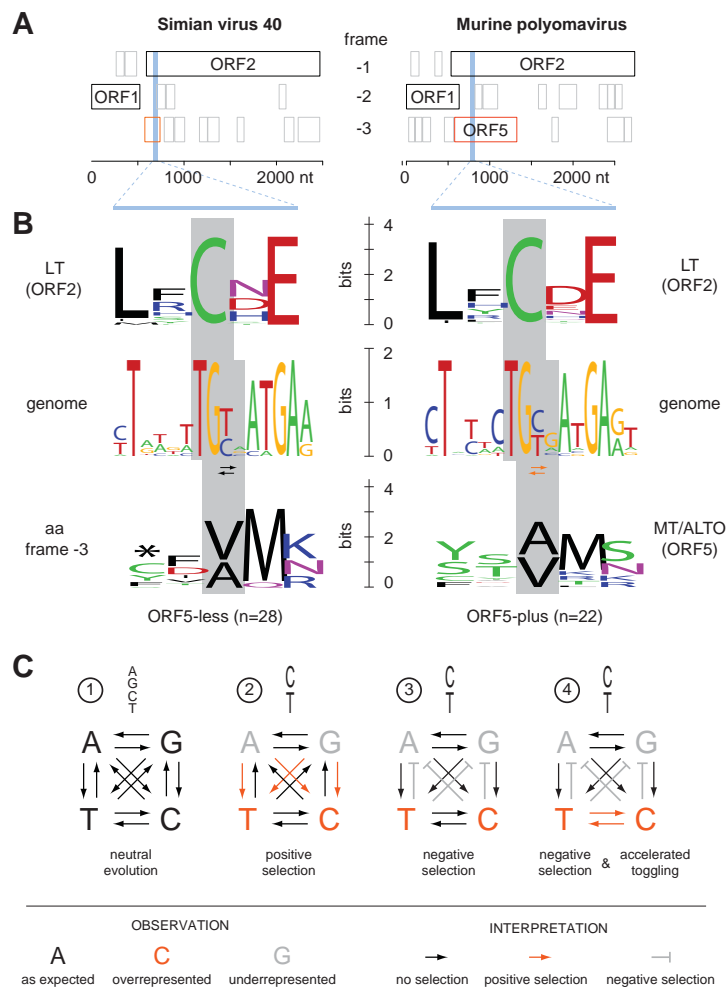


Figure 1. Toggling at the COCO-VA site in mammalian ORF5-plus and ORF5-less polyomaviruses. **(A)** ORF organization in three reading frames of the genomic region encoding the early genes is shown for Simian virus 40, SV-40 (left; NC_001669) and Murine polyomavirus, MPyV (right; Genbank accession NC_001515) representing ORF5-less and ORF5-plus polyomaviruses, respectively. The ORF2 frame was chosen as -1 frame for both viruses. ORF borders are defined here from stop to stop codon. Large expressed ORFs are boxed/outlined and named while other ORFs with a size of at least 75 nt are shown in grey. ORF5 of the ORF5-plus virus and one of its derivatives of the ORF5-less virus are highlighted in the -3 frame. The background highlighting indicates location of the LXCXE motif (an essential motif found in polyomaviruses and other viruses, and cellular proteins that mediates binding and inactivation of the cellular tumour-suppressor protein pRB [82 – 84]). **(B)** Shown are sequence logos of the LXCXE motif (top), the corresponding nucleotide sequence (middle) and the amino acid sequence translated from the ORF5 frame (bottom) for multiple alignment of the 28 ORF5-less (left) and 22 ORF5-plus viruses (right) analysed in this study using Viralis platform [42]. The asterisk indicates stop codons in the -3 frame of some ORF5-less viruses. See **M&M** section for other details. **(C)** Shown are four possible scenarios of evolution of a polynucleotide site under different selection regimes. In scenarios 2 to 4 different selection force(s) result in the same observed nucleotide diversity restricted to C or T. Scenarios 3 and 4 depict the COCO-VA toggling in ORF5-less and -plus viruses, respectively.

four short conserved motifs, designated ORF5m1 to ORF5m4, were evident in the ORF5-wide alignment (**Supplementary Figure S1**) due to an extremely high residue and two-fold size variation (see also below and Carter *et al.*, [15]). They are counterparts of four motifs of LT antigen in the overlapping part of ORF2. Remarkably, the most conserved 3rd aa residue of ORF5m2, identified in this study, has a restricted binary residue variation (Val/Ala) in both ORF5-plus and ORF5-less polyomaviruses (**Figure 1A** and **1B**). Val and Ala are encoded by eight **G(C/T)(A/G/C/T)** triplets which are the only codons compatible with the two **TG(C/T)** codons for conserved Cys of the LT LXCXE motif in the ancestral ORF2 (the two-nucleotide overlap between the Val/Ala and Cys codons is highlighted in bold). In other words, only variation at the second codon position of the COCO-VA codon (**C** or **T**) determines the encoded amino acid (Ala or Val) (**Figure 1B**). We named the observed phenomenon **Codon-Constrained Val-Ala** (COCO-VA) toggling.

The C/T variation represents only half of the full four-nucleotide variation possible at a polynucleotide position (**Figure 1C**). When each kind of nucleotide is equally frequent at a given position, it is likely to evolve at no selection (neutral evolution) (**Figure 1C1**), which may be found in the third codon positions of non-overlapping ORFs. In contrast, a restricted nucleotide variation, like C/T, may emerge as a result of selection, either positive (**Figure 1C2**) or negative (**Figure 1C3**), which is typically observed at the first and second positions of codons of non-overlapping ORFs. Evolutionary interpretation of the nucleotide variation is more complex in the overlapping ORFs, which may be subject to several evolutionary forces acting on each ORF. For instance, there is no doubt that the restricted C/T variation at the 2nd codon position of the COCO-VA site is due to negative selection in the alternative ORF2 to maintain the Cys residue. On the other hand, this restricted variation would be equally compatible with no selection or positive selection in ORF5, with the latter scenario leading to accelerated toggling between C and T (compare **Figure 1C3** and **Figure 1C4**). Therefore, we asked whether selection is involved in the COCO-VA toggling.

Phylogeny suggests accelerated COCO-VA toggling in ORF5-encoding polyomaviruses

First and in line with general reasoning [67], we note that conservation of the LXCXE Cys residue may not constrain the COCO-VA toggling. Second, if C and T nucleotides at the third position of the Cys codon are utilized unevenly, additional *non-ORF2* selection pressure(s), for instance on RNA, must be taken into account when analyzing the toggling. Third, the ORF2 LXCXE conservation in *both* ORF5-plus and -less polyomaviruses provided us with two contrasting virus groups that differ in relation to the COCO-VA site expression through ORF5. Consequently, the COCO-VA site is not expected to be under selection pressure in ORF5-less viruses (**Figure 1C3** scenario), while its evolution in ORF5-plus viruses may or may not be driven by selection depending on the functional importance of these residues (either **Figure 1C3** or **Figure 1C4** scenario). Fourth, the restricted binary choice of aa residues at the COCO-VA site compared to the full 20 amino acid (aa) residue variation simplifies the evolutionary analysis of its residue variation.

Taking all these considerations into account, we reasoned that the relative abundance of *either Ala or Val* in ORF5-plus compared to ORF5-less polyomaviruses would be indicative of selection on residue *type*. Since Ala and Val are similarly and evenly abundant at the COCO-VA site in the known ORF5-plus and ORF5-less mammalian polyomaviruses: 11 vs. 11 and 12 vs. 14 (**Figure 2**), respectively, no indication for selection is apparent. This observation indicates also that the COCO-VA site may not have experienced other, non-ORF2-related selection favoring one of the two nucleotides. Accordingly, we have not found conserved RNA secondary structure elements in this region (see **Supplementary Text S2** and **Supplementary Table S3, Supplementary Figure S2** and **S3**), which, potentially, could have been an alternative source of constraint on the non-synonymous substitution in ORF5.

Next, we investigated the *frequency* of COCO-VA toggling among polyomaviruses. In this and subsequent analyses, switching between Ala and Val residues was accounted with no regard to its direction: from Ala to Val or from Val to Ala. The analysis was limited to the interspecies comparisons. The rate of the COCO-VA toggling in the ORF5-less polyomaviruses provided a baseline rate of COCO-VA toggling that can be expected by chance mutation (neutral evolution). Comparison of this rate with that of the ORF5-plus viruses informed us about directional selection at the COCO-VA codon in the latter viruses.

In the framework of this comparison, we have first mapped COCO-VA toggling on a Bayesian phylogenetic tree of polyomaviruses (**Figure 2**). Due to extreme sequence divergence of the ORF2/ORF5 overlap region in mammalian polyomaviruses (see above and Carter *et al.*, [15]), reliable alignment of this region is limited to four motifs of only ~30 residues in total (**Supplementary Figure S1**), which may not be sufficient for reliable phylogeny reconstruction. Therefore we choose to use a concatenated alignment of other conserved domains representing LT, VP1 and VP2 proteins and accounting for ~50% of genome for phylogeny inference. Large monophyletic groups on this tree were formed by viruses, which were recognized as similar in the ORF2/ORF5 overlapping region. Additionally, we have observed good agreement between topologies of separate branches of this tree, each representing closely related polyomaviruses, with trees of these same viruses using alignments of the ORF2/ORF5 overlap region (**Supplementary Figure S4**). These observations showed that the ORF2/ORF5 overlap region is likely to have coevolved with the LT, VP1 and VP2/3 proteins, whose tree was thus considered suitable for analysis of the COCO-VA toggling.

Subsequently, visual inspection of the tree revealed contrasting patterns of phylogenetic grouping for Ala- and Val-specific viruses in ORF5-plus and ORF5-less subsets of mammalian polyomaviruses, respectively (**Figure 2**). While Ala- and Val-specific viruses were largely intertwined in the first subset, they predominantly formed large residue-specific monophyletic groups in the second subset. This result was indicative of acceleration of the COCO-VA toggling in ORF5-plus viruses. To verify and extend this observation further, we have conducted additional evolutionary-based analyses using available and specially designed approaches.

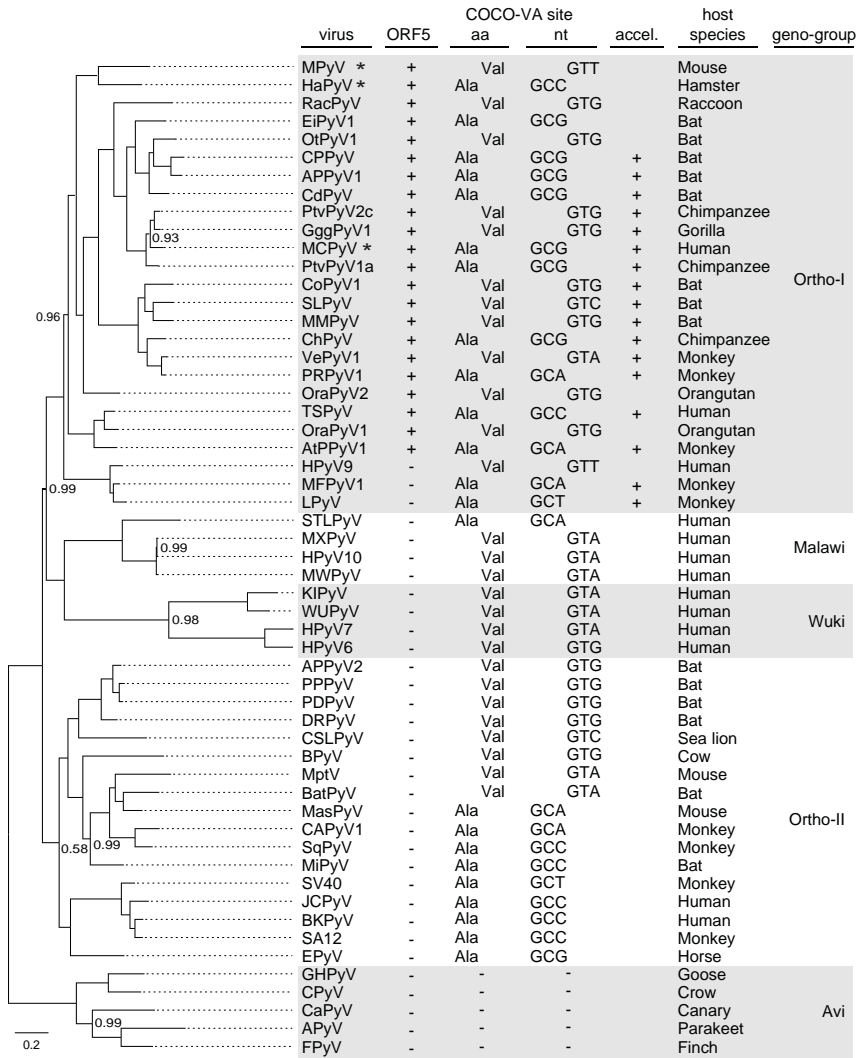


Figure 2. Polyomavirus phylogeny and ORF5 characteristics. Shown is a Bayesian phylogeny using BEAST version 1.7.4 [46] for 55 polyomaviruses (listed in **Supplementary Table S1**) based on conserved regions in the LT, VP1 and VP2 proteins (see **M&M** section for details). The numbers plotted in the tree show posterior probability support values for internal branching events <1. The scale bar is in average number of amino acid substitutions. Asterisks in the virus column indicate viruses for which the ORF5 expression has been demonstrated experimentally. The ORF5 column indicates the presence (+) or absence (-) of ORF5 in polyomaviruses genomes. The COCO-VA site column depicts the residue (Ala or Val) and corresponding codon at the COCO-VA site that is constrained by the Cys codon of the LT LXCXE motif in mammalian polyomaviruses (see **Figure 1** and **Supplementary Figure S1**). The acceleration column (accel.) labels viruses that experienced selection-driven acceleration at the COCO-VA site. The geno-group column depicts the phylogenetic distribution of polyomaviruses according to Feltkamp *et al.*, [40]. Please note that Carter *et al.*, 2013 [15] divided all mammalian polyomaviruses into two groups, monophyletic *Almipolyomaviruses* and paraphyletic *non-Almipolyomaviruses*, which correspond to ORF5-plus and ORF5-less polyomaviruses, respectively. The tree was pseudorooted at the branch connecting mammalian and avian (Avi) polyomaviruses (see **M&M** section for other details).

No evidence for positive selection at the COCO-VA site by conventional evolutionary analyses

We started with employing two most advanced and widely used programs, MEME [49] and FUBAR [27], developed for evolutionary analysis of residue variation. Also included was TOGGLE [25], which was specifically developed for analysis of residue toggling. These three programs are available through the Datamonkey website [26]. The employed programs differ in how they accommodate lineage- and site-specific variation in the analyzed dataset to infer patterns of evolution and deduce selection forces acting on individual codons. Since these tools were developed for the analysis of non-overlapping ORFs, only a single evolutionary force, if identified, is reported. Consequently, we did not expect that these programs could infer both purifying selection (due to Cys conservation of LXCXE) and positive selection (accelerated Val-Ala toggling) at the COCO-VA site of the ORF5-plus viruses as depicted in **Figure 1C4**. Rather, we asked whether the programs could provide evidence for either negative or positive selection at this site of ORF5-plus viruses and negative selection at this site of ORF5-less viruses. This type of inferences depends on the number and diversity of alignment positions under analysis. Due to the high sequence divergence of the ORF2/ORF5 overlapping region, we thus analyzed different subsets of mammalian polyomaviruses, in order to facilitate identification of selection forces. Specifically, the programs were applied to ten different alignments of ORF5, D1-D10, representing selected ORF5 codons, which may or may not be merged with ORF2 codons for different subsets (see **M&M** and **Supplementary Table S4**). In none of the thirty conducted analyses, the COCO-VA codon was identified to be under positive/diversifying selection, including toggling. Also the COCO-VA codon was not found to be negatively selected in analyses that included only ORF5-plus viruses, neither in the entire set nor its D5, D6, D8, D9 and D10 subsets. However, upon analysis of the other five virus sets by FUBAR, either including all ORF5-plus viruses along with other viruses (D1-D4) or including only ORF5-less viruses (D7), the COCO-VA codon was identified to be under purifying selection. In contrast, various other ORF5 codons were identified as being positively or negatively selected or be involved in toggling, in many of these analyses (**Supplementary Table S4**).

The lack of evidence for positive selection/toggling at the COCO-VA codon in ORF5-plus viruses according to these analyses could be either a true negative result (lack of the phenomenon) or a false negative result (failure to detect a signal due to systematic technical deficiency). As detailed below, we believe that the latter explanation is most likely. Indeed, the employed three programs operate under the assumption that the entire codon table of 61 varieties is available for evolution at every site in the analyzed alignments. Consequently, the eight different codons (four for Val and four for Ala) observed at the COCO-VA site were seen as severely restricted rather than representing the full spectrum allowed at this site (imposed by Cys conservation in the overlapping ORF2 codon). This misreading of the observed residue variation has profound implications for its evolutionary interpretation, since high diversity tends to be interpreted as a sign of positive selection, while restricted diversity is commonly associated with purifying selection during evolution of non-overlapping

ORFs. It is thus not surprising that the COCO-VA site evolution was qualified to be under negative selection in several tests by FUBAR. This result could be seen as evidence for the dominance of purifying selection at the COCO-VA site according to FUBAR.

Toggling at the COCO-VA site is significantly accelerated

Due to the above considerations, we decided to continue our testing of the toggling by applying an approach that could be free from the limitations of standard evolutionary based programs developed for non-overlapping ORFs. First, we sought to verify the elevated frequency of Ala-Val exchange in ORF5-plus compared to ORF5-less viruses that is apparent from polyomavirus phylogeny (**Figure 2**). To this end, we have applied the Multistate method of the BayesTraits package [50] to compare the COCO-VA site variation in mammalian ORF5-plus and ORF5-less polyomaviruses. The program employs continuous-time Markov models to estimate the transition rates between multiple states for a single trait (Ala/Val for COCO-VA site in this case) while it traverses a tree. The produced estimates take into account the uncertainty associated with tree reconstruction as it utilizes the full posterior tree sample. The estimated transition rate distribution was plotted for three virus datasets (**Figure 3** left). From this plot it is evident that the estimated Ala-Val exchange rate is more than 3 times higher (13.6 vs 4.3) for ORF5-plus viruses compared to ORF5-less viruses, with a 25-75% interquartile range of 9.5-19.1 and 2.5-7.7, respectively. This striking difference between the two datasets is strongly supported by a log Bayes Factor (logBF) of 3352.2, which is astronomically large and dwarfs the significance threshold of 2. As expected, the estimate for mammalian polyomaviruses was intermediate between those two with a 25-75% interquartile range of 5.6-13.3 (**Figure 3** left). As Ala-Val exchange at the COCO-VA position is equivalent to C-T exchange at the second codon position (see **Figure 2** and **M&M**) we have compared its exchange rate to that at the third codon position (**Figure 3** right). This position accepts all four nucleotides and its variation is primarily driven by selection in the overlapping ORF2 in which it occupies the first codon position of the subsequent residue. As may be expected, the exchange rate at this position (now averaged over four instead of two nucleotides) is comparable for ORF5-plus and ORF5-less viruses (median and 25-75% interquartile range: 10.2 and 8.5-12.3 vs. 9.8 and 8.0-11.9, respectively). Of notice, these numbers are still and consistently smaller than those of the Ala-Val exchange rate for ORF5-plus viruses.

Importantly, the observed difference at the second codon position (i.e., the Ala-Val exchange) may not be attributed to differences in virus diversity of the compared two datasets, whose distributions of smallest pair-wise patristic distance (SPAT) values (median value and 25-75% inter-quartile range: 0.23 and 0.13-0.32 vs. 0.28 and 0.12-0.46) were not different at a statistically significant level (Mann-Whitney U test; $p=0.42$). Consequently, we concluded that the COCO-VA toggling rate is significantly and genuinely accelerated in the ORF5-plus compared to the ORF5-less polyomaviruses. Since the COCO-VA site is expressed in ORF5-plus but not ORF5-less viruses (although see below), this result implies positive selection on the COCO-VA site in ORF5-plus viruses.

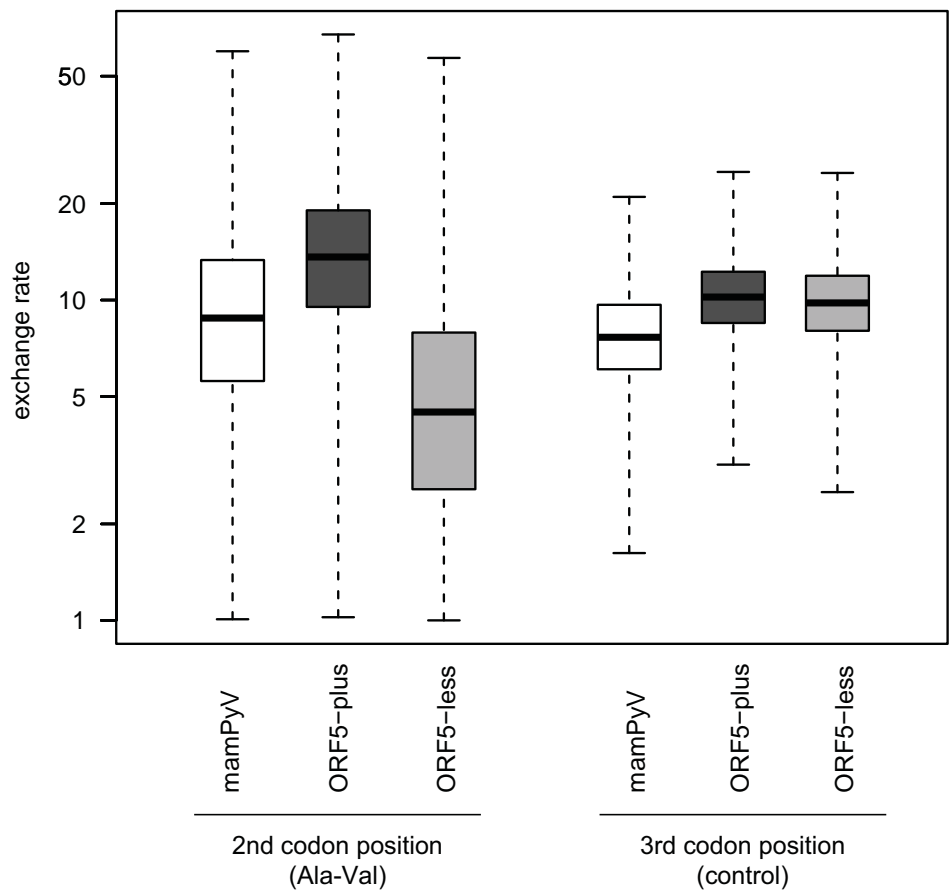


Figure 3. COCO-VA toggling is accelerated in ORF5-plus compared to ORF5-less viruses. Shown are results of BayesTraits multistate analysis of COCO-VA toggling rate in three groups of viruses. The distributions of estimated exchange rates at second (left side) and third (right side) codon position of the COCO-VA codon is shown. The exchange rate at the second codon position corresponds to the COCO-VA toggling while the rate at the third codon position serves as a control. The distributions are shown as Box-and-whisker graphs. The boxes span from the first to the third quartile and include the median (bold line), and the whiskers (dashed lines) extend to the extreme values.

Ratio approach to study accelerated COCO-VA toggling

To study the accelerated COCO-VA toggling further, we have developed a ratio approach remotely similar to that of comparing the ratio of non-synonymous to synonymous substitutions. We used the ratio of monoSPAT/(monoSPAT+polySPAT) values as a normalized measure of the COCO-VA toggling rate relative to Ala/Ala or Val/Val persistence, with polySPAT and monoSPAT resembling estimations of non-synonymous and synonymous substitutions, respectively (for group designations see **M&M**). Only the C/T variation at the 2nd codon position that controls Val-Ala exchange, rather than the entire codon for Ala/Val as it would have

been the case upon analysis of a non-overlapping ORF by a conventional technique, was analyzed in our test. We thus avoided complications to the analysis that would otherwise be caused by the unaccounted evolutionary pressure on the third position of Ala/Val codons by the ORF2 overlapping codon, where it occupies the first position of the subsequent residue (**Figure 1B**). An SPAT ratio of 0.5 indicates that the Ala/Val exchange rate matches that of Ala/Ala or Val/Val persistence during evolution of a particular lineage (hereafter, matching rate). Since amino acid residue persistence at the COCO-VA site in a pair of viruses may involve either no genetic change or synonymous substitution, the matching rate for toggling under the model of neutral evolution could be expected only at sufficiently large evolutionary distances when chance mutation, either synonymous or non-synonymous, is highly probable. Accordingly, persistence would dominate over toggling at smaller distances under this model, resulting in SPAT ratios smaller than 0.5. If positive selection is involved in toggling, increase of SPAT ratios compared to those expected under neutral evolution could be observed at sufficiently small evolutionary distances.

The above considerations indicate that under the model of neutral evolution we could expect different SPAT ratios at small and large evolutionary distances. To verify this and define ranges for small and large evolutionary distances separating pairs of monomorphic viruses, we analyzed the difference between the matching rate and within-window distributions of SPAT ratios involving all mammalian polyomaviruses, which were plotted against monoSPAT values. Due to stochastic reasons, the estimated toggling rate may deviate from the actual rate for a virus. To address this limitation, we pooled SPAT ratios within a predefined window that was slid along the monoSPAT axis. Our analysis revealed that the Ala/Val exchange rate of polyomaviruses varies considerably, with very different median values being observed in two monoSPAT ranges (**Figure 4A**). In the monoSPAT range of 0-0.35, median SPAT ratios were consistently smaller compared to the matching rate, while in the monoSPAT range of 0.35-0.75, they were consistently larger than the matching rate. Accordingly, the entire monoSPAT range was split into two sub-ranges in our subsequent analyses. For evolutionary interpretations of the ratio test in subsequent analyses, we used the results obtained for ORF5-less viruses as a base-line, since Val-Ala toggling in these viruses is expected to experience no selection (**Figure 1C3** scenario) over the entire evolutionary distance range.

Accelerated COCO-VA toggling is associated with interspecies diversification of ORF5-plus polyomaviruses

Is the accelerated toggling a characteristic of the entire ORF5-plus viruses or its subsets? The difference between SPAT ratios in the distributions for ORF5-plus and ORF5-less viruses was statistically significant over the monoSPAT range of 0-0.35 (MWU test p -value=7e-06), but not over the 0.35-1.2 range (MWU test p -value=0.829) (**Figure 4B**). Importantly, this result may not be due to biases of the virus sampling which was comparable for ORF5-plus and -less viruses in the two distributions along the monoSPAT range (**Supplementary Figure S5**).

Consequently, the above observations indicate a selection-driven acceleration of COCO-VA toggling in the majority (fifteen out of twenty-two) of ORF5-plus polyomaviruses, each of which is separated from another monomorphic virus by a monoSPAT of 0.35 or smaller (**Figure 4**). For the remaining seven ORF5-plus viruses, each of which is separated from another monomorphic virus by a monoSPAT larger than 0.35, no accelerated COCO-VA toggling was observed. This could be either due to specifics of evolution or the unavailability of close monomorphic relatives of these viruses in the current sampling. If the former is true, the viruses with accelerated COCO-VA toggling may be expected to cluster in the tree, while a random phyletic distribution is likely otherwise. **Figure 2** shows that the fifteen viruses with accelerated COCO-VA toggling are scattered across the entire branch of ORF5-plus viruses. This observation implies that the accelerated COCO-VA toggling may involve *all* ORF5-plus viruses (all terminal nodes in the respective tree branch) thus presenting an extreme case of convergent evolution. An improved, much larger virus sampling, which includes closely related viruses for each analyzed virus species, will enable verification of this implication. Also, it may facilitate additional insights, including: a) refining the estimate of the monoSPAT threshold at which the COCO-VA toggling acceleration can be observed, and b) extending our analysis to poorly sampled intra-species diversity, in order to address the question whether COCO-VA toggling drives speciation or vice versa.

Could the observed difference between ORF5-plus and ORF5-less viruses in the monoSPAT range of 0-0.35 (**Figure 4B**) have emerged also under the evolutionary scenario that is alternative to that involving positive selection on ORF5-plus and no selection on ORF5-less viruses? If the COCO-VA site was under strong negative selection in ORF5-less viruses while being under either weak negative or no selection in ORF5-plus viruses, SPAT ratio of these viruses would differ. The following considerations make this scenario unlikely to be applicable to explain the data obtained in our study. First, this scenario implies that the COCO-VA site must be expressed in *all* ORF5-less viruses. These viruses include some of the most well characterized polyomaviruses, e.g. SV40, with no evidence for the expression of the COCO-VA site, although some of the poorly characterized ORF5-less viruses may indeed express this site (see below). We could also recall that ORF5-less viruses were defined as a group not having the property (ORF5) rather than having one, which would be required to link strong negative selection to the functional characteristic. Second, SPAT ratio of ORF5-plus viruses in the monoSPAT range of 0-0.35 is comparable to the matching rate ($p=0.390$ in Wilcoxon rank sum test; see **Supplementary Table S5**). This result is in the excellent agreement with positive selection acting on the COCO-VA site in ORF5-plus viruses, while it may not be reconciled with the weak negative selection hypothesis. On the other hand, it would in principle be compatible with neutral evolution of the COCO-VA site in ORF5-plus viruses under the condition that polyomaviruses have very high mutation rate. The estimates of this rate vary greatly and generally this aspect has not been fully resolved [68]. However, we note that the Val-Ala variation is already observed in several monophyletic subsets of ORF5-plus viruses which otherwise diverged little or modestly. This observation indicates that the

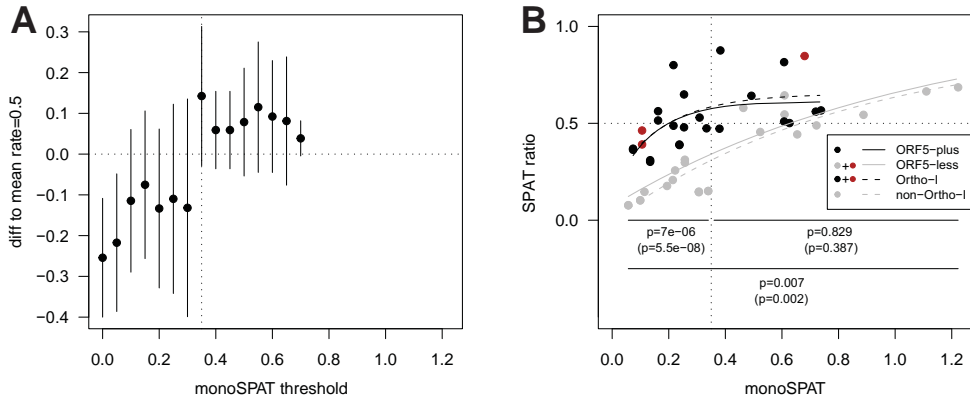


Figure 4. Accelerated toggling at the COCO-VA site in ORF5-plus and Ortho-I viruses. For the purpose of this analysis a representative set of 50 mammalian polyomaviruses (see **Supplementary Table S1**) was studied. Due to the lack of species demarcation criteria for polyomaviruses, we chose to consider viruses with different names as representing different species (dots in the plot). The only exception was made for MX polyomavirus, Human polyomavirus 10 and MW polyomavirus, which were represented only by the latter because of the very small distances that separate these three viruses. Two pairs of virus partitioning (subsets) of the mammalian polyomaviruses, based on the ORF5 presence and phylogeny, were considered. They and their colour codes are defined in the inset of panel **B**. (**A**) The partitioning of the monoSPAT scale at 0.35 was derived based on the drop of the mean difference of SPAT ratios to the matching rate. Here, a sliding window (size 0.15, shift 0.05) starting at monoSPAT of 0.0 was moved along the monoSPAT range to calculate within-window mean differences (dots) and associated standard deviations (vertical lines). See also **M&M** section and **Supplementary Figure S5** for other details. (**B**) The curves show the fit of a 3-parameter logistic function to each of four different subsets. The numbers below show P-values of Mann–Whitney U tests comparing the SPAT ratio distributions between ORF5-plus and ORF5-less viruses (*Orthopolyomavirus-I* and non-*Orthopolyomavirus-I*) for two monoSPAT ranges (0–0.35, 0.35–1.25). A horizontal dotted line is drawn at the matching rate, whose evolutionary interpretation is defined in the text.

Val-Ala variation may be among most frequent rather than average as would be expected under the neutral evolution scenario. This aspect could be studied most closely with the improved virus sampling. In conclusion, based on the available data the accelerated COCO-VA toggling due to positive selection is the most likely evolutionary scenario.

Accelerated COCO-VA toggling is most strongly associated with monophyletic Ortho-I viruses

The results described above provide evidence for accelerated COCO-VA toggling in the 0–0.35 monoSPAT range for ORF5-plus viruses. However, it is also evident that the distributions of ORF5-plus and –less SPAT ratios overlap with two ORF5-less viruses deviating considerably from their group-mates and instead fitting into the other group rather well (**Figure 4B**, red dots in the 0–0.35 monoSPAT range). This grouping with ORF5-plus viruses received strong statistical support when analyzing all of the 145,422,675 possible 16-by-14 combinations of the 30 viruses with monoSPAT values in the range of 0–0.35 (**Supplementary Text S3** and **Supplementary Figure S6**). Intriguingly, these two viruses along with another one for which no closely related monomorphic virus is available in the current virus sampling (**Figure 4B**,

red dot in the 0.35-1.20 monoSPAT range) form a sister lineage of ORF5-plus polyomaviruses at the root of the Ortho-I monophyletic group (**Figure 2**) [15, 40]. These results suggest that the accelerated toggling is most strongly associated with the Ortho-I group. Since the observed accelerated toggling is indicative of positive selection that may be realized only upon expression of the COCO-VA site in the two (three) poorly characterized basal Ortho-I viruses, such hypothesis must be considered. It could be achieved using a mechanism other than expression of the entire ORF5, for instance, through alternative splicing of mRNA(s) [69] that could fuse the COCO-VA site with other ORF(s). In the evolutionary framework, such expression of the COCO-VA site would be ancestral to those used by ORF5-plus viruses, implying that the COCO-VA site and the associated sequence motif could be a nucleation site for the subsequent ORF5 origin by ORF expansion [15].

COCO-VA toggling is located in a SLiM of an intrinsically disordered region

What could be the structural basis of COCO-VA toggling? Bioinformatics analyses indicate that MT/ALTO is a Pro residue rich IDR, whose motifs could form SLiMs (**Figures S1, S7, S8 and S9**) [15]. Thus, the interspecific toggling targets a SLiM, which is in line with the notion that IDRs evolving differently than structured protein regions [6, 7]. Since SLiMs promote protein folding in relatively flat energy landscapes [70] and mediate interactions with partners that are relatively weak [71, 72], difference between physico-chemical properties of the just two possible COCO-VA residues, Val and Ala, could be of significance. For instance, these residues have contrasting structural propensities, favoring the formation of either α -helix (Ala) or β -sheet (Val) [73], which might be used to promote alternative folding of MT/ALTO upon interaction with partner(s). Unfortunately, this hypothesis may not be tested using the available computational approaches.

Concluding Remarks

In complex protein networks, SLiMs are emerging as evolutionary adaptive transmitters of intracellular signals involving multiple interacting partners [2, 3, 13, 74]. Here we presented evidence for the evolutionary signature of adaptation in the otherwise uncharacterized SLiM of MT/ALTO. The effect of COCO-VA toggling on the SLiM may be similar to that of phosphorylation which could modulate SLiM activity considerably [75]. MT antigen of rodent polyomaviruses has been shown to interact through its ORF5-encoded part with numerous cellular targets involved in signal transduction [33, 35, 36, 39]. The function of ALTO, identified just recently, has not been resolved yet [15]. The described COCO-VA toggling is notable because of a unique combination of properties: it involves one of the just few conserved positions of the otherwise highly divergent MT/ALTO protein, and it may affect every species of Ortho-I polyomaviruses. These viruses are known to infect bats, rodents, monkeys, hominids and humans with apparently frequent host switching (**Figure 2**). Future studies should identify driving forces of the COCO-VA toggling to enable its comparison with intra-species residue toggling [25]. The latter is likely driven by the cellular immune

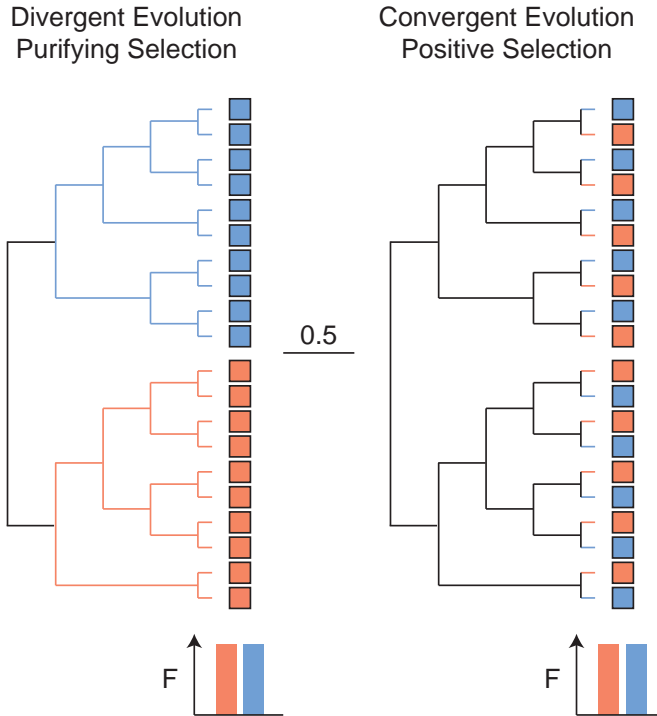


Figure 5. Contrasting modes of evolution at a conserved protein site accepting two residues. Shown are fictional examples of evolution of a conserved protein site with two-residue variation in families of structured (left panel) and unstructured (right) proteins, whose evolutionary scale of replacement at all sites was considerable (bar 0.5) and whose phylogeny is described by identical trees. In both cases, two residues are evenly distributed, each occupying 50% of terminal nodes. Residue type either clusters into two monophyletic groups (left) or is intertwined (right). The left panel depicts divergent evolution driven mostly by purifying selection as seen in many characterized structured proteins. The right panel depicts convergent evolution driven by positive selection as discovered at the COCO-VA site in the presented chapter and may be experienced at other sites in unstructured proteins.

response and occurs at much smaller time and divergence scales, and with the exchange of many residues. Practically, our study suggests that analysis of substitution rates can be applied to individual residues in overlapping ORFs. It extends the utility of the substitution rate analysis from mapping to dissecting functional elements in overlapping ORFs.

The described phenomenon also challenges common perception of conservation of proteins, which is believed to be inversely and universally correlated with the rate of evolution. Accordingly, sites accepting relatively few residues are classified conserved and evolving slowly under negative selection. Typically, such residues are critical for maintaining protein core and/or playing an essential role in the active site of structured proteins. Besides the sites that are strictly invariant, those that accept only two residues during large-scale evolution are among the most conserved. Exchange of these residues could happen due to either rare fixation of non-synonymous mutation that is driven by episodic positive selection

or residue drift. As a result, each of the two residues is likely to be associated with a large monophyletic clade in the tree (**Figure 5** left panel), the pattern that can be recognized by available programs (e.g. [76]). Examples of this type of evolution are plenty in many protein families. For instance, rare exchange of the catalytic nucleophile Cys and Ser residues in virus proteases with chymotrypsin-like fold [77] or phosphate-binding Ser and Thr residues in the Walker-box GKS/T motif of nucleotide-binding proteins [78], are notable. The above considerations indicate that in structured proteins, limited residue variation may largely be imposed by the molecular environment in which these proteins operate. In contrast, constraints on the genetic level is the chief factor determining residue variation in proteins encoded in overlapping ORFs. Consequently, this restricted residue variation in overlapping ORFs may not be linked to residue function in the manner described for structured proteins. Accordingly, overlapping ORFs predominantly encode unstructured proteins with their most conserved SLiMs mediating adaptation, a function that is commonly facilitated by the least conserved elements in structured proteins. Along the same line, we now provide evidence for the phylogenetic intertwining of viruses that employ, respectively, Val and Ala at the conserved COCO-VA site in the IDR of MT/ALTO. When depicted in a simplified form, this phylogenetic pattern can be contrasted with the clade-specific association of residues in a tree of structured proteins (compare right and left panels of **Figure 5**). This contrast is particularly striking since it is not evident in the cumulative frequency of residues at terminal nodes (bottom panels underneath of trees in **Figure 5**). Thus, this logos-style representation of residue conservation, which is very popular in functional studies, may not capture residue change and its role in adaptation. Only analysis in the context of phylogeny could do it, as demonstrated in this study.

Since Cys is one of the least frequent amino acid residues and none of the other residues can constrain evolution in the -3 RF (or +1 RF) to only two residues, the described codon-constrained accelerated toggling might be viewed as an extremely exotic phenomenon limited to polyomaviruses. We believe that this perception is biased for several reasons. First of all, the (unknown) diversity of the Virus Universe is expected to be many orders of magnitude larger than the number of currently recognized few thousand virus species [79, 80]. This implies a good chance of discovering COCO-VA toggling in other viruses in the future. Furthermore, accelerated toggling might involve more than two residues at a site that could still be considered conservative relative to many other sites. Such constraint could be imposed by conserved amino acids other than Cys in the overlapping ORF or, if non-overlapping ORF is involved, by a different genetic mechanism, e.g. RNA structure, or even by a partner or partners interacting with an IDR site. Thus, the described COCO-VA toggling may represent an extreme case of common evolution of individual residues in IDRs of proteins, making it potentially relevant to understanding biology and pathology of adaptation of many organisms.

Acknowledgement

The authors thank Jelle J. Goeman, Sander Kooijman, Andrey M. Leontovich, and Els van der Meijden for useful discussions, and Igor Sidorov and Dmitry Samborskiy for help with Viralis. A.E.G. is member of the Netherlands Bioinformatics Center (NBIC) Faculty.

Funding

This study was funded by the Leiden University Medical Center, The Netherlands, partially through the Collaborative Agreement in Bioinformatics between Leiden University Medical Center and Moscow State University (MoBiLe Program), and Leiden University Fund.

References

1. Wright PE, Dyson HJ. Intrinsically unstructured proteins: re-assessing the protein structure-function paradigm (1999) *J Mol. Biol.* 293: 321-331.
2. Dyson HJ, Wright PE. Intrinsically unstructured proteins and their functions (2005) *Nat Rev Mol Cell Biol.* 6: 197-208.
3. Dunker AK, Silman I, Uversky VN, Sussman JL. Function and structure of inherently disordered proteins (2008) *Curr Opin Struct Biol.* 18: 756-764.
4. Tompa P. The interplay between structure and function in intrinsically unstructured proteins (2005) *FEBS Lett.* 579: 3346-3354.
5. Xue B, Dunker AK, Uversky VN. Orderly order in protein intrinsic disorder distribution: disorder in 3500 proteomes from viruses and the three domains of life (2012) *J Biomol. Struct. Dyn.* 30: 137-149.
6. Brown CJ, Johnson AK, Daughdrill GW. Comparing models of evolution for ordered and disordered proteins (2010) *Mol. Biol. Evol.* 27: 609-621.
7. Brown CJ, Johnson AK, Dunker AK, Daughdrill GW. Evolution and disorder (2011) *Curr. Opin. Struct. Biol.* 21: 441-446.
8. Davey NE, Van Roey K, Weatheritt RJ, Toedt G, Uyar B, Altenberg B, Budd A, Diella F, Dinkel H, Gibson TJ. Attributes of short linear motifs (2012) *Mol. Biosyst.* 8: 268-281.
9. Keese PK, Gibbs A. Origins of genes: "big bang" or continuous creation? (1992) *Proc Natl Acad Sci U S A* 89: 9489-9493.
10. Firth AE, Brown CM. Detecting overlapping coding sequences in virus genomes (2006) *BMC Bioinformatics.* 7: 75.
11. Belshaw R, Pybus OG, Rambaut A. The evolution of genome compression and genomic novelty in RNA viruses (2007) *Genome Res.* 17: 1496-1504.
12. Sabath N, Wagner A, Karlin D. Evolution of viral proteins originated de novo by overprinting (2012) *Mol. Biol. Evol.* 29: 3767-3780.
13. Pushker R, Mooney C, Davey NE, Jacque JM, Shields DC. Marked variability in the extent of protein disorder within and between viral families (2013) *PLoS. One.* 8: e60724.
14. Ling R, Pate AE, Carr JP, Firth AE. An essential fifth coding ORF in the sobemoviruses (2013) *Virology* 446: 397-408.
15. Carter JJ, Daugherty MD, Qi X, Bheda-Malge A, Wipf GC, Robinson K, Roman A, Malik HS, Galloway DA. Identification of an overprinting gene in Merkel cell polyomavirus provides evolutionary insight into the birth of viral genes (2013) *Proc Natl Acad Sci U. S. A.* 110: 12744-12749.
16. Firth AE, Atkins JF. Evidence for a novel coding sequence overlapping the 5'-terminal approximately 90 codons of the gill-associated and yellow head okavirus envelope glycoprotein gene (2009) *Virol J* 6: 222.

17. Firth AE, Atkins JF. A case for a CUG-initiated coding sequence overlapping torovirus ORF1a and encoding a novel 30 kDa product (2009) *Virol J* 6: 136.
18. Loughran G, Firth AE, Atkins JF. Ribosomal frameshifting into an overlapping gene in the 2B-encoding region of the cardiovirus genome (2011) *Proc Natl Acad Sci U S A* 108: E1111-E1119.
19. Fang Y, Treffers EE, Li Y, Tas A, Sun Z, van der Meer Y, de Ru AH, van Veelen PA, Atkins JF, Snijder EJ, Firth AE. Efficient -2 frameshifting by mammalian ribosomes to synthesize an additional arterivirus protein (2012) *Proc Natl Acad Sci U S A* 109: E2920-E2928.
20. Mizokami M, Orito E, Ohba K, Ikeo K, Lau JY, Gojobori T. Constrained evolution with respect to gene overlap of hepatitis B virus (1997) *J Mol. Evol.* 44 Suppl 1: S83-S90.
21. Hughes AL, Hughes MA. Patterns of nucleotide difference in overlapping and non-overlapping reading frames of papillomavirus genomes (2005) *Virus Res.* 113: 81-88.
22. de Groot S, Mailund T, Hein J. Comparative annotation of viral genomes with non-conserved gene structure (2007) *Bioinformatics.* 23: 1080-1089.
23. Sabath N, Landan G, Graur D. A method for the simultaneous estimation of selection intensities in overlapping genes (2008) *PLoS ONE* 3: e3996.
24. Kryazhimskiy S, Plotkin JB. The population genetics of dN/dS (2008) *PLoS Genet.* 4: e1000304.
25. Delport W, Scheffler K, Seoighe C. Frequent toggling between alternative amino acids is driven by selection in HIV-1 (2008) *PLoS Pathog.* 4: e1000242.
26. Kosakovsky Pond SL, Frost SD. Datamonkey: rapid detection of selective pressure on individual sites of codon alignments (2005) *Bioinformatics* 21: 2531-2533.
27. Murrell B, Moola S, Mabona A, Weighill T, Sheward D, Kosakovsky Pond SL, Scheffler K. FUBAR: a fast, unconstrained bayesian approximation for inferring selection (2013) *Mol. Biol. Evol.* 30: 1196-1205.
28. Dalianis T, Hirsch HH. Human polyomaviruses in disease and cancer (2013) *Virology* 437: 63-72.
29. Kazem S, van der Meijden E, Feltkamp MC. The trichodysplasia spinulosa-associated polyomavirus: virological background and clinical implications (2013) *APMIS* 121: 770-782.
30. Decaprio JA, Garcea RL. A cornucopia of human polyomaviruses (2013) *Nat Rev Microbiol.* 11: 264-276.
31. DeCaprio JA, Imperiale MJ, Major EO. Polyomaviruses. In: Fields *VIROLOGY* (2013). Knipe DM, Howley PM (editors). Philadelphia: Wolters Kluwer / Lippincott Williams & Wilkins: 1633-1661.
32. Kazem S, Lauber C, van der Meijden E, Kooijman S, Bialasiewicz S, Wang RC, Gorbalenya AE, Feltkamp MC. Global circulation of slowly evolving trichodysplasia spinulosa-associated polyomavirus and its adaptation to the human population through alternative T antigens (2013) *J Neurovirol.* 19: 298-299.

33. Hutchinson MA, Hunter T, Eckhart W. Characterization of T antigens in polyoma-infected and transformed cells (1978) *Cell* 15: 65-77.
34. Topalis D, Andrei G, Snoeck R. The large tumor antigen: a “Swiss Army knife” protein possessing the functions required for the polyomavirus life cycle (2013) *Antiviral Res.* 97: 122-136.
35. Courtneidge SA, Goutebroze L, Cartwright A, Heber A, Scherneck S, Feunteun J. Identification and characterization of the hamster polyomavirus middle T antigen (1991) *J. Virol.* 65: 3301-3308.
36. Fluck MM, Schaffhausen BS. Lessons in signaling and tumorigenesis from polyomavirus middle T antigen (2009) *Microbiol. Mol. Biol. Rev.* 73: 542-563.
37. Magnusson G, Nilsson MG, Dilworth SM, Smolar N. Characterization of polyoma mutants with altered middle and large T-antigens (1981) *J. Virol.* 39: 673-683.
38. Yi X, Freund R. Deletion of proline-rich domain in polyomavirus T antigens results in virus partially defective in transformation and tumorigenesis (1998) *Virology* 248: 420-431.
39. Cheng J, DeCaprio JA, Fluck MM, Schaffhausen BS. Cellular transformation by Simian Virus 40 and Murine Polyoma Virus T antigens (2009) *Semin Cancer Biol.* 19: 218-228.
40. Feltkamp MC, Kazem S, van der Meijden E, Lauber C, Gorbalenya AE. From Stockholm to Malawi: recent developments in studying human polyomaviruses (2013) *J Gen Virol* 94: 482-496.
41. Johne R, Buck CB, Allander T, Atwood WJ, Garcea RL, Imperiale MJ, Major EO, Ramqvist T, Norkin LC. Taxonomical developments in the family Polyomaviridae (2011) *Arch. Virol.* 156: 1627-1634.
42. Gorbalenya AE, Lieutaud P, Harris MR, Coutard B, Canard B, Kleywegt GJ, Kravchenko AA, Samborskiy DV, Sidorov IA, Leontovich AM, Jones TA. Practical application of bioinformatics by the multidisciplinary VIZIER consortium (2010) *Antiviral Res.* 87: 95-110.
43. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput (2004) *Nucleic Acids Res.* 32: 1792-1797.
44. Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H, Valentin F, Wallace IM, Wilm A, Lopez R, Thompson JD, Gibson TJ, Higgins DG. Clustal W and Clustal X version 2.0 (2007) *Bioinformatics.* 23: 2947-2948.
45. Castresana J. Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis (2000) *Mol. Biol. Evol.* 17: 540-552.
46. Drummond AJ, Suchard MA, Xie D, Rambaut A. Bayesian phylogenetics with BEAUti and the BEAST 1.7 (2012) *Mol Biol Evol* 29: 1969-1973.
47. Whelan S, Goldman N. A general empirical model of protein evolution derived from multiple protein families using a maximum-likelihood approach (2001) *Mol. Biol. Evol.* 18: 691-699.

48. Drummond AJ, Ho SY, Phillips MJ, Rambaut A. Relaxed phylogenetics and dating with confidence (2006) *PLoS Biol* 4: e88.
49. Murrell B, Wertheim JO, Moola S, Weighill T, Scheffler K, Kosakovsky Pond SL. Detecting individual sites subject to episodic diversifying selection (2012) *PLoS Genet.* 8: e1002764.
50. Pagel M, Meade A, Barker D. Bayesian estimation of ancestral character states on phylogenies (2004) *Syst. Biol* 53: 673-684.
51. Paradis E, Claude J, Strimmer K. APE: Analyses of Phylogenetics and Evolution in R language (2004) *Bioinformatics.* 20: 289-290.
52. Soding J. Protein homology detection by HMM-HMM comparison (2005) *Bioinformatics.* 21: 951-960.
53. Schneider TD, Stephens RM. Sequence logos: a new way to display consensus sequences (1990) *Nucleic Acids Res.* 18: 6097-6100.
54. Crooks GE, Hon G, Chandonia JM, Brenner SE. WebLogo: a sequence logo generator (2004) *Genome Res.* 14: 1188-1190.
55. Linding R, Jensen LJ, Diella F, Bork P, Gibson TJ, Russell RB. Protein disorder prediction: implications for structural proteomics (2003) *Structure.* 11: 1453-1459.
56. Ward JJ, Sodhi JS, McGuffin LJ, Buxton BF, Jones DT. Prediction and functional analysis of native disorder in proteins from the three kingdoms of life (2004) *J Mol. Biol.* 337: 635-645.
57. Deng X, Eickholt J, Cheng J. PreDisorder: ab initio sequence-based prediction of protein disordered regions (2009) *BMC. Bioinformatics* 10: 436.
58. Prilusky J, Felder CE, Zeev-Ben-Mordehai T, Rydberg EH, Man O, Beckmann JS, Silman I, Sussman JL. FoldIndex: a simple tool to predict whether a given protein sequence is intrinsically unfolded (2005) *Bioinformatics* 21: 3435-3438.
59. Linding R, Russell RB, Neduva V, Gibson TJ. GlobPlot: Exploring protein sequences for globularity and disorder (2003) *Nucleic Acids Res.* 31: 3701-3708.
60. Dosztanyi Z, Csizmok V, Tompa P, Simon I. IUPred: web server for the prediction of intrinsically unstructured regions of proteins based on estimated energy content (2005) *Bioinformatics* 21: 3433-3434.
61. Yang ZR, Thomson R, McNeil P, Esnouf RM. RONN: the bio-basis function neural network technique applied to the detection of natively disordered regions in proteins (2005) *Bioinformatics* 21: 3369-3376.
62. Vucetic S, Brown CJ, Dunker AK, Obradovic Z. Flavors of protein disorder (2003) *Proteins* 52: 573-584.
63. Rost B, Yachdav G, Liu J. The PredictProtein server (2004) *Nucleic Acids Res.* 32: W321-W326.
64. Buchan DW, Ward SM, Lobley AE, Nugent TC, Bryson K, Jones DT. Protein annotation and modelling servers at University College London (2010) *Nucleic Acids Res.* 38: W563-W568.

65. Gruber AR, Neubock R, Hofacker IL, Washietl S. The RNAz web server: prediction of thermodynamically stable and evolutionarily conserved RNA structures (2007) *Nucleic Acids Res.* 35: W335-W338.
66. Goodman SN. Toward evidence-based medical statistics. 2: The Bayes factor (1999) *Ann Intern Med.* 130: 1005-1013.
67. Rogozin IB, Spiridonov AN, Sorokin AV, Wolf YI, Jordan IK, Tatusov RL, Koonin EV. Purifying and directional selection in overlapping prokaryotic genes (2002) *Trends Genet.* 18: 228-232.
68. Firth C, Kitchen A, Shapiro B, Suchard MA, Holmes EC, Rambaut A. Using time-structured data to estimate evolutionary rates of double-stranded DNA viruses (2010) *Molecular Biology and Evolution* 27: 2038-2051.
69. Zheng ZM. Viral oncogenes, noncoding RNAs, and RNA splicing in human tumor viruses (2010) *Int. J. Biol. Sci.* 6: 730-755.
70. Fisher CK, Stultz CM. Constructing ensembles for intrinsically disordered proteins (2011) *Curr Opin Struct Biol* 21: 426-431.
71. Van Roey K, Gibson TJ, Davey NE. Motif switches: decision-making in cell regulation (2012) *Curr Opin. Struct Biol* 22: 378-385.
72. Perkins JR, Diboun I, Dessailly BH, Lees JG, Orengo C. Transient protein-protein interactions: structural, functional, and network properties (2010) *Structure.* 18: 1233-1243.
73. Kabsch W, Sander C. Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features (1983) *Biopolymers* 22: 2577-2637.
74. Kovacs E, Tompa P, Liliom K, Kalmar L. Dual coding in alternative reading frames correlates with intrinsic protein disorder (2010) *Proc Natl Acad Sci U. S. A.* 107: 5429-5434.
75. Deribe YL, Pawson T, Dikic I. Post-translational modifications in signal integration (2010) *Nat. Struct Mol. Biol* 17: 666-672.
76. Gu X, Zou Y, Su Z, Huang W, Zhou Z, Arendsee Z, Zeng Y. An update of DIVERGE software for functional divergence analysis of protein family (2013) *Mol. Biol Evol.* 30: 1713-1719.
77. Gorbalenya AE, Snijder EJ. Viral cysteine proteinases (1996) *Persp. Drug Discov. Design* 6: 64-86.
78. Koonin EV, Gorbalenya AE. Tale of two serines (1989) *Nature* 338: 467-468.
79. King AMQ, Adams MJ, Carstens EB, Lefkowitz EJ. *Virus Taxonomy: Ninth Report of the International Committee on Taxonomy of Viruses.* London: Academic Press; 2012.
80. Breitbart M, Rohwer F. Here a virus, there a virus, everywhere the same virus? (2005) *Trends Microbiol.* 13: 278-284.
81. Komarova AV, Haenni AL, Ramirez BC. Virus versus host cell translation love and hate stories (2009) *Adv. Virus Res.* 73: 99-170.

82. Kazem S, van der Meijden E, Wang RC, Rosenberg AS, Pope E, Benoit T, Fleckman P, Feltkamp MC. Polyomavirus-associated trichodysplasia spinulosa involves hyperproliferation, pRB phosphorylation and upregulation of p16 and p21 (2014) PLoS ONE 9: e108947.
83. Decaprio JA. How the Rb tumor suppressor structure and function was revealed by the study of Adenovirus and SV40 (2009) Virology. 384: 274-284.
84. de Souza RF, Iyer LM, Aravind L. Diversity and evolution of chromatin proteins encoded by DNA viruses (2010) Biochim. Biophys. Acta 1799: 302-318.

Supplementary Data

Supplementary Text

Supplementary Text S1. Nomenclature of open reading frames of mammalian polyomaviruses

Polyomaviruses use genomes with two non-overlapping protein-coding regions that are expressed either early or late in infection, respectively [1]. Each region includes open reading frames (ORFs) that produce – for some through alternative splicing – early proteins Small, Middle and Large T antigen (ST, MT and LT) and the late capsid proteins (VP2/VP3 and VP1), respectively (**Supplementary Table S2**) [2]. ST and LT are ubiquitous in all mammalian polyomaviruses. MT, however, was known only for the mouse (MPyV) and hamster (HaPyV) polyomaviruses. Recently, a third early protein called ALTO was described for MCPyV that is homologous to the C-terminal domain of MT [3].

Polyomaviruses use alternative pre-messengerRNA (primary transcript) splicing to produce mRNAs that direct synthesis of proteins. The currently used nomenclature focuses on annotation of proteins and the respective transcripts. When ORFs are named, which happens only occasionally in literature, protein-based designation is used. The produced designation could be confusing for ORFs that encode more than one protein. For instance, the first ORF in the early region is commonly called ST ORF, while it also encodes the first exon of MT and LT protein. To address this complexity, we have used in this study a rational nomenclature of ORFs designations that is independent from names of proteins/transcripts. Its rational is similar to that used to design ORFs nomenclature in other virus families with similarly complex relations between ORFs, transcripts and proteins [4]. We defined regions flanked by two stop codons as ORFs in three different reading frames. Two pairs of ubiquitous early and late ORFs were designated ORF1 and ORF2, and ORF3 and ORF4, respectively, while an optional early ORF was designated ORF5. In this regard, the early primary transcript can be alternatively spliced to merge (parts of) ORF1 and ORF2 for directing the synthesis of LT, and (parts of) ORF1 and ORF5 for directing the synthesis of MT. Alternatively, ORF1 and ORF5 can be translated directly through internal start codons, respectively, resulting into synthesis of ST and ALTO (**Supplementary Table S2**). The developed ORF nomenclature can accommodate the identification of new alternatively spliced transcripts as well as proteins expressed by canonical and non-canonical mechanisms.

Supplementary Text S2. No conserved RNA secondary structure elements are evident around COCO-VA genomic site

Conserved RNA secondary structures constrain evolution of the respective genomic regions of viruses involved. To clarify whether toggling at the COCO-VA site was affected by the presence of conserved RNA secondary structures in polyomaviruses, we used the RNAz web-based program to predict functional RNA structures based on two criteria: i) evolutionary

conservation and ii) thermodynamic stability [5]. Due to overall poor conservation of ORF5 in mammalian polyomaviruses, the analysis was conducted in monophyletic subsets of polyomaviruses for which reliable alignments were possible to produce (see also [3]). In total, nucleotide alignments for six subsets of the ORF5-plus and ORF5-less polyomaviruses were generated (**Supplementary Table S3** and **Figure 3**) and analyzed in a genomic region around the LXCXE overlapping COCO-VA site (see **M&M** for details).

Several RNA secondary structures in the tested region with a probability higher than the cutoff value of $p=0.5$ were predicted by the RNAz program in the subsets 1, 2 and 3 of the ORF5-plus viruses (**Supplementary Figure S2**). Majority of these secondary structures were predicted on the negative-strand of the input RNA-alignments. No conserved RNA secondary structures were identified on the positive or negative-stands of other subsets. Importantly, one of the predicted RNA structures closely corresponds to the experimentally validated pre-microRNA located on the negative-strand of MCPyV RNA [6, 7] (**Supplementary Figure S3**), lending further and independent support to the results of RNA structure analysis by RNAz. Since COCO-VA site is NOT base-paired in either of the predicted RNA structures, we concluded that the evolution of the COCO-VA site is not likely constrained by RNA secondary structure elements in mammalian polyomaviruses.

Supplementary Text S3. Accelerated COCO-VA toggling is most associated with monophyletic Ortho-I viruses

In addition to assessing statistical significance of difference between SPAT ratio distributions of ORF5-plus and ORF5-less viruses, we sought to assess the combinatory scale of 14/16 partitioning that would give another, top-down perspective on the probability value obtained for the difference. To this end, we ranked the difference among differences for all possible two-side partitioning of this type (14 vs. 16) for 30 viruses in the 0-0.35 mono-SPAT range. In total, 145,422,675 values were obtained and plotted as a histogram revealing highly symmetrical bell-like distribution of values (**Supplementary Figure S6**), whose difference of group means of toggling rate ratios for all possible combinations was virtually zero ($2.7e-17$). The ORF5-plus/ORF5-less partitioning was ranked #284 (counted from the left side) with a difference of means of -0.2726 between the compared datasets. This high rank corresponds to a p-value of $1.95e-06$ that was close to the $7e-06$ value obtained in MWU test. The high ranking of the ORF5-plus/ORF5-less partitioning strongly supported the prior conclusions.

In the view of such high ranking, why this partitioning did not outrank all other partitionings? Inspection of partitioning with higher ranks showed that, like ORF5-plus viruses, many of them and including the number 1 (difference of means of -0.3051), involved a large subset of phylogenetically compact Ortho-I viruses, which was contrasted against mainly non-Ortho-I viruses. This observation prompted us to compare Ortho-I vs. non-Ortho-I viruses. Since this partitioning compares 16 to 14 viruses, its ranking could be derived from the density distribution already used to rank ORF5-plus/ORF5-less partitioning

(14 vs. 16 viruses), now from its right tail (**Supplementary Figure S6**). Its analysis showed that the Ortho-I/non-Ortho-I partitioning had a difference of mean toggling rate values of 0.2993 that ranked number three. This ranking is a two-order improvement over the ranking of ORF5-plus/ORF5-less partitioning and corresponded to a p-value of $2.1\text{e-}08$, which was close to $5.5\text{e-}08$ in the MWU test (**Figure 4B**). The difference in values of mean toggling rates for the Ortho-I/non-Ortho-I partitioning was only marginally worse than those for two top-ranking 16/14 partitioning, which had 0.3007 and 0.2996 values, respectively. This result showed that the accelerated toggling is most strongly associated with the tree-based Ortho-I/non-Ortho-I partitioning.

Sequence alignments and the polyomavirus phylogeny used in this study can be found at <https://github.com/chrtin/COCOVAatoggling>.

Supplementary Tables

Supplementary Table S1. Polyomavirus representatives used in this study*

Accession	Abbreviation	Virus name	Host
NC_001515	MPyV	Murine polyomavirus	Mouse
NC_001663	HaPyV	Hamster polyomavirus	Hamster
JQ178241	RacPyV	Raccoon polyomavirus	Raccoon
NC_020068	EiPyV1	Eidolon polyomavirus-1	Bat
NC_020071	OtPyV1	Otomops polyomavirus-1	Bat
JQ958889	CPPyV	Carollia perspicillata polyomavirus	Bat
JQ958887	APPyV1	Artibeus planirostris polyomavirus	Bat
NC_020067	CdPyV	Cardioderma polyomavirus	Bat
HQ385749	PtPyV2c	Pan troglodytes verus polyomavirus-2c	Chimpanzee
HQ385752	GggPyV1	Gorilla gorilla gorilla polyomavirus-1	Gorilla
NC_010277	MCPyV	Merkel cell polyomavirus	Human
HQ385746	PtPyV1a	Pan troglodytes verus polyomavirus1a	Chimpanzee
NC_020065	CoPyV1	Chaerephon polyomavirus-1	Bat
JQ958888	SLPyV	Sturnira lilium polyomavirus	Bat
JQ958893	MMPyV	Molossus molossus polyomavirus	Bat
NC_014743	ChPyV	Chimpanzee polyomavirus	Chimpanzee
NC_019844	VePyV1	Vervet monkey polyomavirus-1	Monkey
NC_019850	PRPyV1	Ptilocolobus rufomitratus polyomavirus-1	Monkey
FN356901	OraPyV2	Orangutan polyomavirus	Orangutan
NC_014361	TSPyV	Trichodysplasia spinulosa-associated polyomavirus	Human
FN356900	OraPyV1	Orangutan polyomavirus	Orangutan
NC_019853	AtPPyV1	Ateles paniscus polyomavirus-1	Monkey
NC_015150	HPyV9	Human polyomavirus-9	Human
NC_019851	MFPyV1	Macaca fascicularis polyomavirus-1	Monkey
NC_004763	LPyV	African green monkey polyomavirus	Monkey
NC_020106	STLPyV	STL polyomavirus	Human
JX259273	MXPyV	MX polyomavirus	Human
JX262162	HPyV10	Human polyomavirus-10	Human
JQ898291	MWPyV	MW polyomavirus	Human
NC_009238	KIPyV	KI polyomavirus	Human
NC_009539	WUPyV	WU Polyomavirus	Human
NC_014407	HPyV7	Human polyomavirus-7	Human
NC_014406	HPyV6	Human polyomavirus-6	Human
JQ958890	APPyV2	Artibeus planirostris polyomavirus	Bat
JQ958891	PPPyV	Pteronotus parnellii polyomavirus	Bat
NC_020070	PDPyV	Pteronotus polyomavirus	Bat
JQ958892	DRPyV	Desmodus rotundus polyomavirus	Bat
NC_013796	CSLPyV	California sea lion polyomavirus	Sea lion
NC_001442	BPyV	Bovine polyomavirus	Cow
NC_001505	MptV	Murine pneumotropic virus	Mouse
NC_011310	BatPyV	Myotis polyomavirus VM-2008	Bat
AB588640	MasPyV	Mastomys polyomavirus	Mouse
NC_019854	CAPyV1	Cebus albifrons polyomavirus-1	Monkey
NC_009951	SqPyV	Squirrel monkey polyomavirus	Monkey
NC_020069	MiPyV	Miniopterus polyomavirus	Bat
NC_001669	SV40	Simian virus-40	Monkey
NC_001699	JCPyV	JC polyomavirus	Human
NC_001538	BKPyV	BK polyomavirus	Human
NC_007611	SA12	Simian virus-12	Monkey
JQ412134	EPyV	Equine polyomavirus	Horse
NC_004800	GHPyV	Goose hemorrhagic polyomavirus	Goose
NC_007922	CPyV	Crow polyomavirus	Crow
GU345044	CaPyV	Canary polyomavirus	Canary
NC_004764	APyV	Budgerigar fledgling disease virus-1	Parakeet
NC_007923	FPyV	Finch polyomavirus	Finch

*Viruses are listed according to phylogenetic tree in **Figure 2**.

Supplementary Table S2. ORF designations used in this study and by others

Open reading frame nomenclature		Encoding	Reference
This study	By others		
ORF1	ST-antigen	Small T antigen	[8]
	MT-antigen	Middle T antigen 1 st exon	[9]
	LT-antigen	Large T antigen 1 st exon	[10]
ORF2	LT-antigen	Large T antigen 2 nd exon	[10]
ORF3	VP2	VP2	[11]
	VP3	VP3	[12]
ORF4	VP1	VP1	[12]
ORF5	MT-antigen	Middle T antigen 2 nd exon	[9]
	ALTO protein	ALTO protein	[3]

Supplementary Table S3. Virus subset compositions whose alignments were used for RNA secondary structure prediction

Subset	Polyomavirus	ORF5
1	GggPyV1, PtvPyV2c, PtvPyV1a, MCPyV APPyV1, CPPyV, CdPyV, OtPyV1, EiPyV1, RacPyV	plus
2	VePyV1, ChPyV, PRPyV1, MMPyV, SLPyV, CoPyV1	plus
3	TSPyV, AtPPyV1, OraPyV1	plus
4	MPyV, HaPyV	plus
5	HPyV9, MFPyV1, LPyV	plus
6	KIPyV, WUPyV, HPyV6, HPyV7, STLPyV, HPyV10, MXPpyV, MWPyV, JCPyV, SA12, SV40, BKPyV, APPyV2, PPPyV, DRPyV, PDPyV, CAPyV1, SqPyV, MiPyV, MptV, MasPyV, BatPyV	less

Supplementary Table S4. W Results of natural selection analysis of conserved ORF5 positions using Datamonkey

Dataset	Genome region ¹	Viruses ²	Codons ³		MEME ⁴		FUBAR ⁴		TOGGLE ⁴		
			ORF5	ORF2	#pos	COCOVA	#pos	#neg	COCOVA	#toggling	COCOVA
D1	ORF5	mamPyV	34	0	12	none	17	1	neg	15	none
D2	ORF5+ORF2	mamPyV	34	474	15	none	14	2	neg	14	none
D3	ORF5	Ortho-I	34	0	8	none	10	1	neg	7	none
D4	ORF5+ORF2	Ortho-I	34	473	6	none	5	3	neg	7	none
D5	ORF5	ORF5+	34	0	5	none	8	0	none	6	none
D6	ORF5+ORF2	ORF5+	34	473	7	none	7	2	none	5	none
D7	ORF5+ORF2	ORF5-	34	474	15	none	11	4	neg	11	none
D8	ORF5	lineageA	121	0	30	none	15	20	none	2	none
D9	ORF5	lineageB	183	0	14	none	4	15	none	1	none
D10	ORF5	lineageC	314	0	2	none	1	2	none	n.a. ⁵	n.a. ⁵

¹ ORF2: C-terminal part of ORF2 not overlapping with ORF5

² mamPyV: mammalian polyomaviruses; ORF5+: ORF5-plus viruses; ORF5-: ORF5-less viruses; lineageA: MCPyV, GgPyV1, PtvPyV1a, EIPyV, OtPyV1, CPPyV1, APPyV1, CdpPyV, RacPyV; lineageB: VePyV1, PRPyV1, ChPyV, MMPyV1, SLPyV, CoPyV1; lineageC: TSPyV (3 isolates), OraPyV1, AtPyV1

³ The number and origin of codons used to infer phylogeny. For ORF5 of D1-D7, only codons whose overlapping codon in ORF2 (-1 frame) was aligned with no gaps across all mammalian polyomaviruses were chosen. For ORF2 of D2, D4, D6 and D7, and ORF5 of D8-D10, weakly aligned columns were removed manually from the alignment.

⁴ #: number of ORF5 codons; pos: positive/diversifying selection; neg: negative/purifying selection; none: no selection/toggling detected; COCOVA: COCO-VA codon in ORF5. When using TOGGLE, only ORF5 codons have been screened for toggling

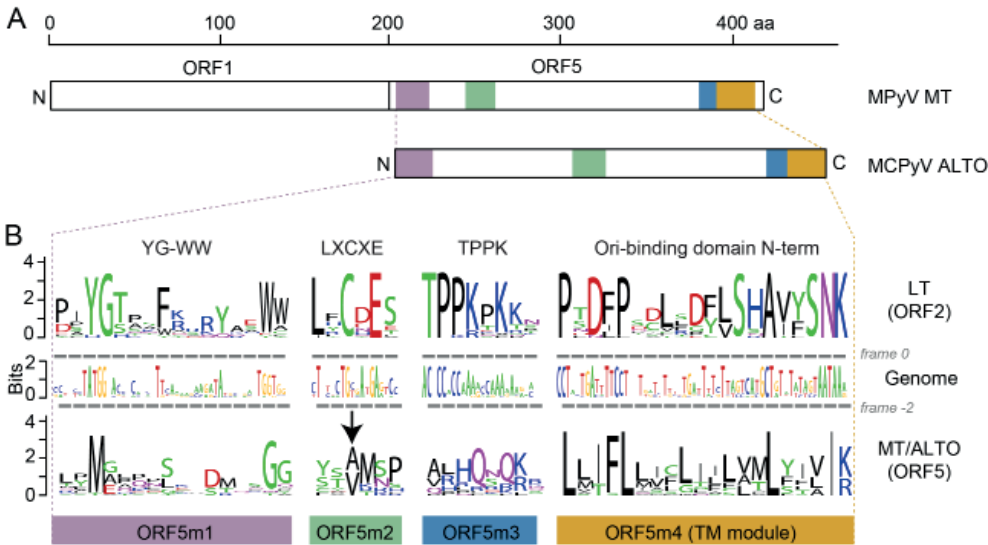
⁵ n.a., non-available: no results were obtained due to small (<10) number of tree branches

Supplementary Table S5. Shown are p-values of a Wilcoxon rank sum test for deviation of the mean of a distribution of SPAT ratios from the matching rate

Virus group	MonoSPAT range		
	0 - 0.35	0.35 - 1.25	0 - 1.25
ORF5-plus	0.390	0.039	0.656
ORF5-less	3.1e-5	0.064	0.003
Ortho-I	0.211	0.020	0.672
non-Ortho-I	1.2e-4	0.129	0.001

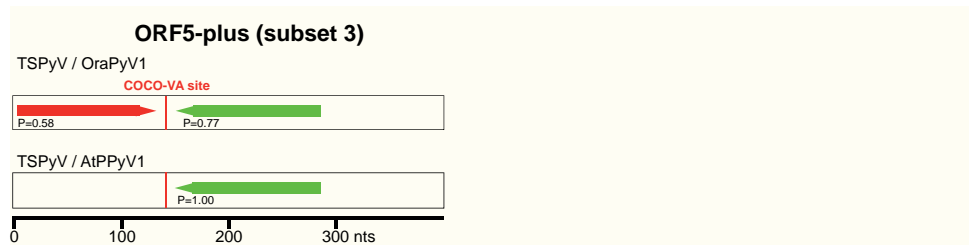
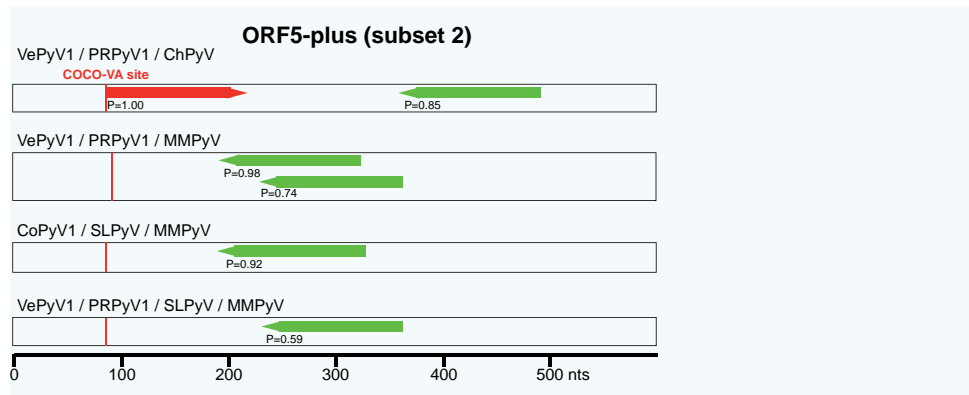
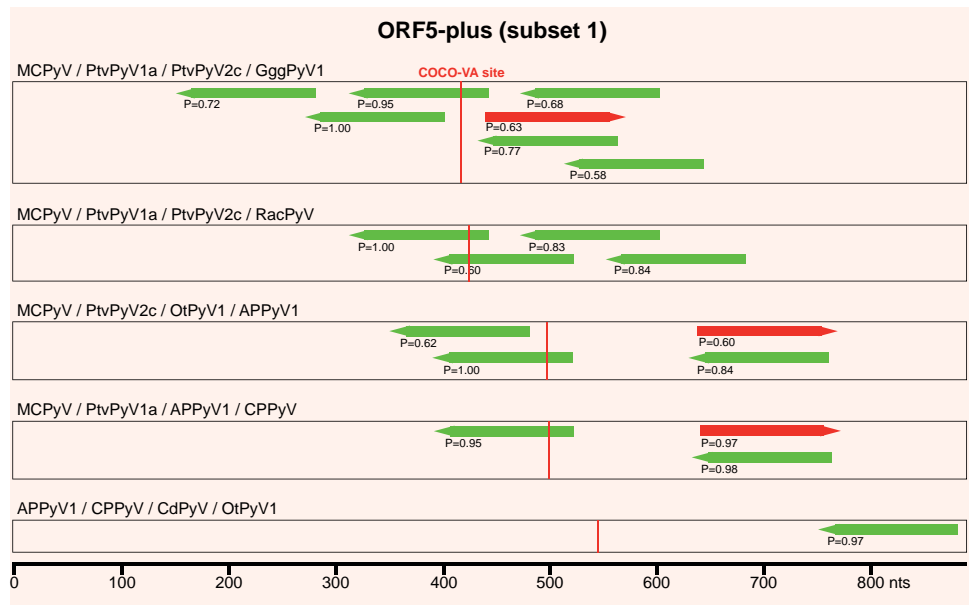
Bold, statistically significant p-values ($\alpha=0.05$)

Supplementary Figures

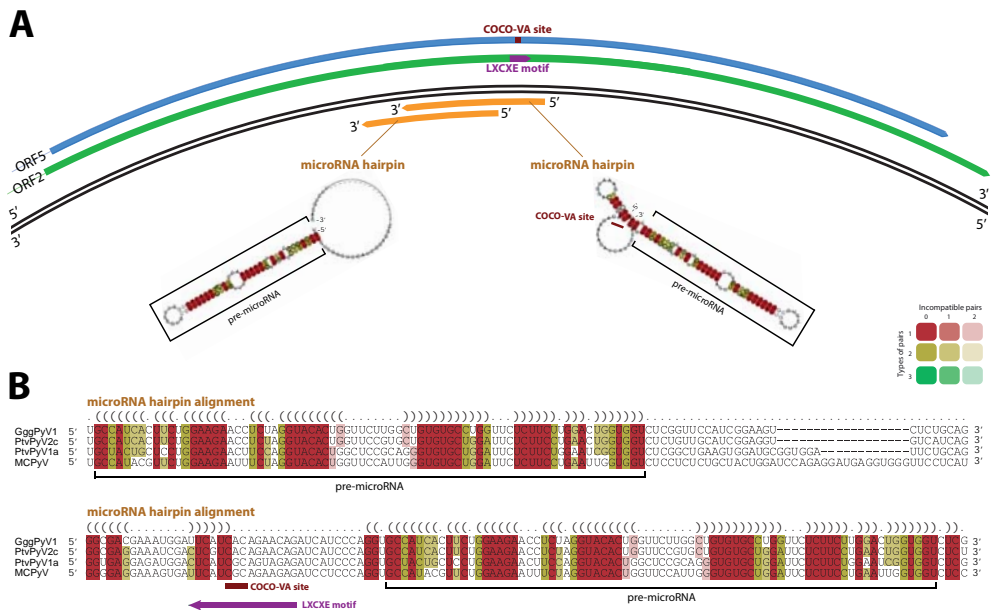


Supplementary Figure S1. Conserved features of ORF5 of the Ortho-I group. **(A)** domain organization of two different proteins, partially or fully encoded in ORF5, that were experimentally characterized in MPyV and MCPyV. The ORF1-encoded MT domain is not further specified. The ORF5-encoded domain includes four conserved motifs, ORF5m1 - ORF5m4, colored differently and detailed in panel **B**. The ORF5 remains open upstream of the start codon for MCPyV ALTO. **(B)** Conserved motifs in the ORF5/ORF2 overlap. Presented are sequence logos of four conserved motifs in LT ORF2 (bottom), their counterparts in MT/ALTO ORF5 (top), and the corresponding genome region (middle). The logos are based on alignments of the 22 viruses of the Ortho-I group (see **M&M**). Gray bars indicate the codon structure of the LT ORF2 and MT/ALTO ORF5 reading frames. The arrow indicates the COCO-VA toggling position.

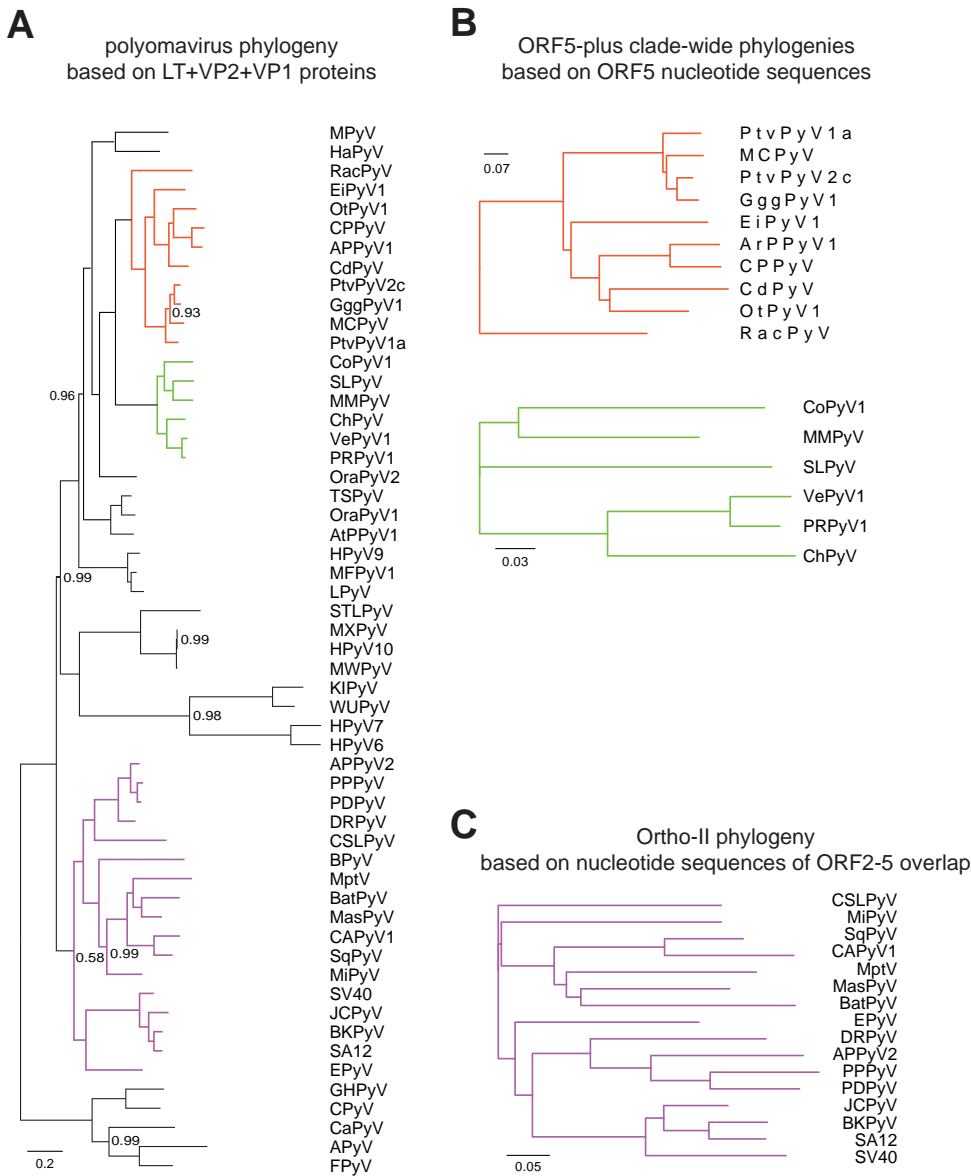
5



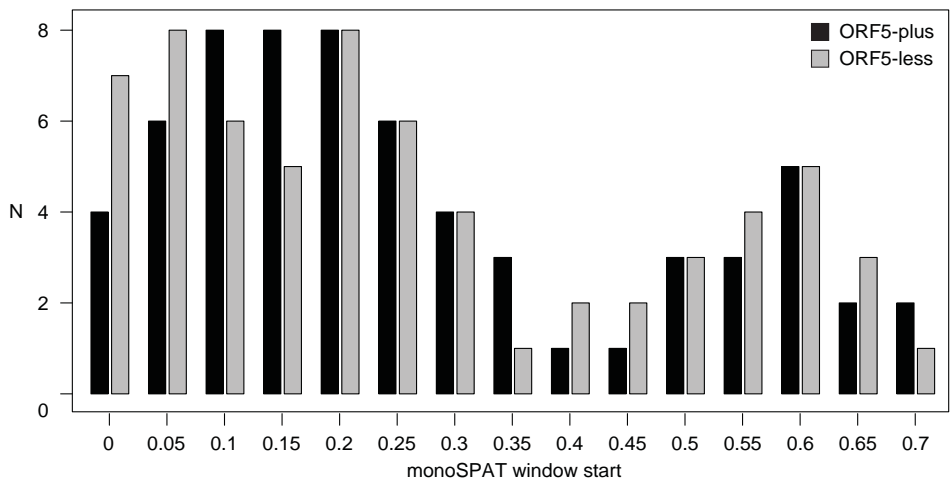
Supplementary Figure S2. Prediction of RNA secondary structure consensus using RNAz. ORF5-plus viruses from subset 1 (top panel), subset 2 (middle panel) and subset 3 (bottom panel) alignment demonstrated regions of RNA secondary structures in several combination of viruses. The predicted secondary structures from the negative-strand (green bars) and positive-strand (red bars), within the given alignment of indicated viruses with a probability value higher than 0.5, are shown. Position of the COCO-VA in the given alignment is illustrated by horizontal red line.



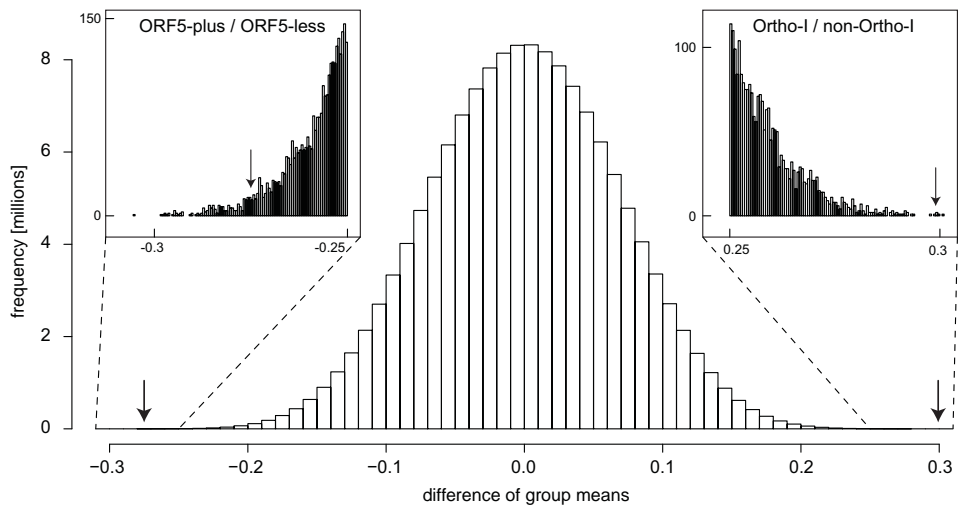
Supplementary Figure S3. Pre-microRNA hairpin prediction by RNAz. **(A)** Two overlapping predicted RNA secondary structures (orange arrows) in the input alignments of subset 1 covering the ORF2-ORF5 overlapping region of ~800 nucleotides, and the corresponding region relative to the LXCXE motif/COCO-VA site, are shown. Blue arrow indicates ORF5 encoded amino acids, and green arrow depicts ORF2 encoded amino acids. The significance levels of RNAz predicted secondary structure hits are color-coded. Colors indicate the number of different types of base pairs that support stabilizing selection on the structures (inset: color table-combination). Circles indicate variable positions in the stems, also known as loops. **(B)** Structure and color annotated output alignment of the two RNA secondary structures shown in **A**, which were predicted in MCPyV, GgPyV1, PtvPyV2c and PtvPyV1a input alignment. This region corresponds to the pre-microRNA shown experimentally for MCPyV [6, 7]. The dot-bracket notation on top of each alignment indicates loop-stem, respectively, of the RNA secondary structure consensus. Region relative to the LXCXE motif and COCO-VA site sequence is shown below the bottom alignment.



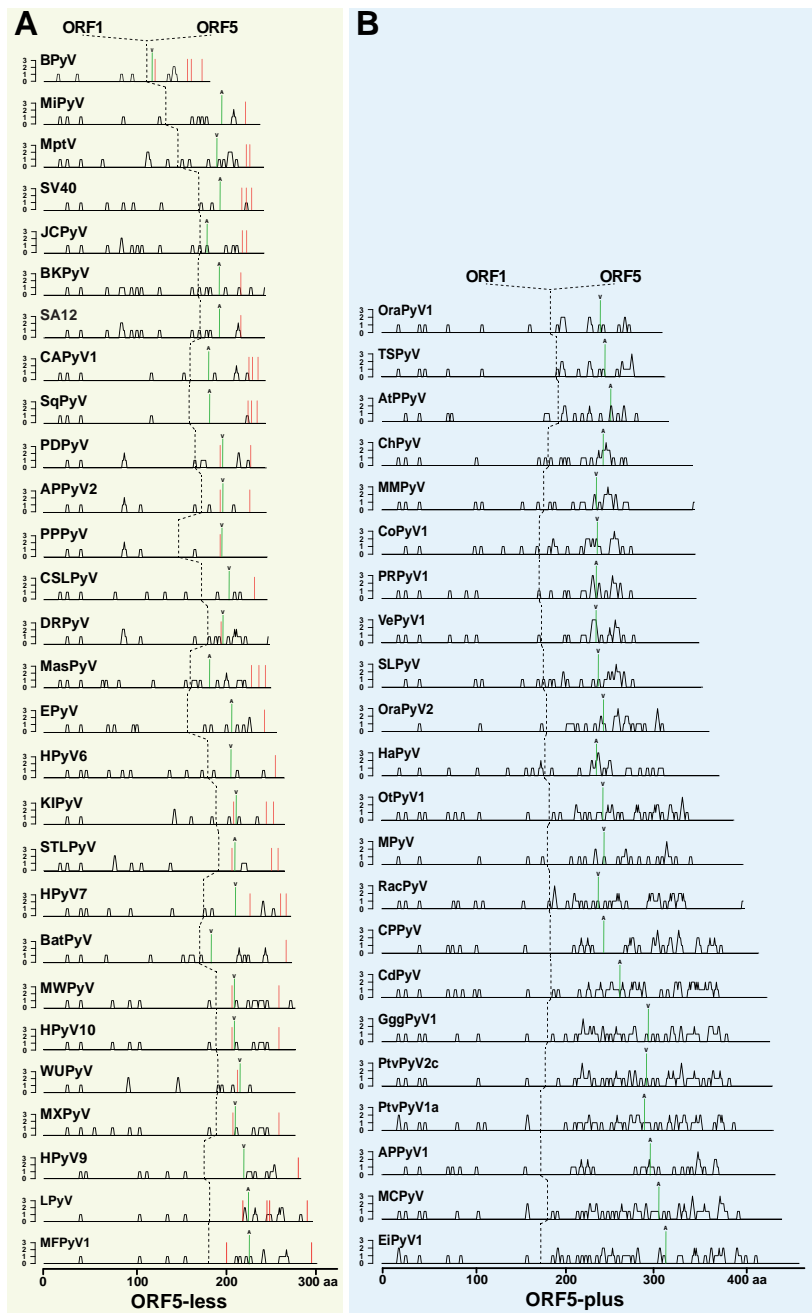
Supplementary Figure S4. Congruence of lineage-specific phylogenies with the polyomavirus tree. Compared is the topology of the polyomavirus tree estimated using LT, VP1, and VP2 proteins (**A**) with that of two ORF5-plus monophyletic lineages independently estimated using ORF5 codon alignments (**B**) and of Ortho-II viruses estimated using nucleotide alignments of the ORF2-ORF5 overlapping genomic region (**C**). The lineages are indicated by color. The trees in **B** and **C** are neighbor joining trees used in the Datamonkey analyses.



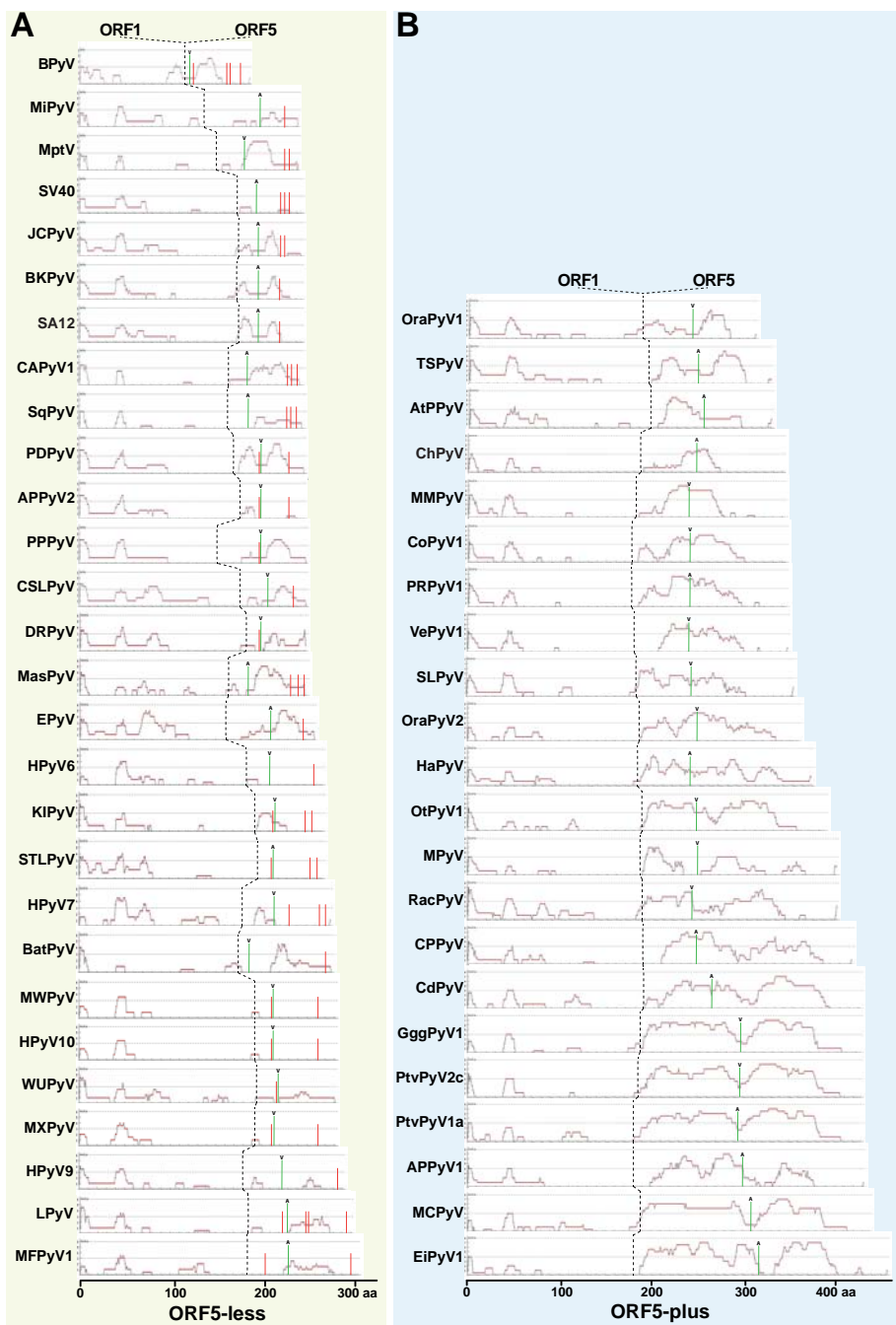
Supplementary Figure S5. Virus sampling in relation to monoSPAT. The number of ORF5-plus and ORF5-less viruses per sliding window (size: 0.15, shift: 0.05) along the monoSPAT scale is shown.



Supplementary Figure S6. Ranking test. Shown are differences of mean SPAT ratios between groups for all possible 14-16 partitionings of the 30 viruses with monoSPAT values in the range of 0-0.35. These are 145,422,675 combinations in total. Values for the Ortho-I/non-Ortho-I and for the ORF5-plus/ORF5-less partitionings are indicated by arrows and insets.



Supplementary Figure S7. Enrichment of Proline residues near the COCO-VA site in ORF5-less (A) and ORF5-plus (B) viruses. Shown is the number and location of Proline residues along the translated putative MT polypeptide sequence of the mammalian polyomaviruses analyzed in this study. The green and red bars indicate the COCO-VA site and termination codons, respectively. The dashed line indicates the ORF1-ORF5 splice junction. ORF5-plus viruses show an enrichment of Proline residues at the ORF5-encoded part of MT, which is a characteristic of intrinsically disordered protein regions.



Supplementary Figure S8. The COCO-VA site is embedded in an intrinsically disordered protein region in ORF5-less (A) and ORF5-plus (B) viruses. Shown is the predication of protein disorder along the translated putative MT polypeptide sequence of the mammalian polyomaviruses analyzed in this study. The green and red bars indicate the COCO-VA site and termination codons, respectively. The dashed line indicates the ORF1-ORF5 splice junction. Large parts of the ORF5-encoded part of MT are predicted to be disordered for ORF5-plus viruses.



Supplementary Figure S9. Protein secondary structure and COCO-VA site in ORF5-less (A) and ORF5-plus (B) viruses. Consensus of two protein secondary structure predictors is shown along the translated putative MT polypeptide sequence of the mammalian polyomaviruses. Red and green bars indicate termination codons and COCO-VA site, respectively, and dashed lines suggest the putative ORF1-ORF5 splice junction. On top, a color-coded scoring table for Helices and Sheets is shown.

Supplementary References

1. Broekema NM, Imperiale MJ. miRNA regulation of BK polyomavirus replication during early infection (2013) *Proc. Natl. Acad. Sci. U. S. A* 110: 8200-8205.
2. Decaprio JA, Garcea RL. A cornucopia of human polyomaviruses (2013) *Nat Rev Microbiol.* 11: 264-276.
3. Carter JJ, Daugherty MD, Qi X, Bheda-Malge A, Wipf GC, Robinson K, Roman A, Malik HS, Galloway DA. Identification of an overprinting gene in Merkel cell polyomavirus provides evolutionary insight into the birth of viral genes (2013) *Proc Natl Acad Sci U. S. A.* 110: 12744-12749.
4. Komarova AV, Haenni AL, Ramirez BC. Virus versus host cell translation love and hate stories (2009) *Adv. Virus Res.* 73: 99-170.
5. Gruber AR, Neubock R, Hofacker IL, Washietl S. The RNAz web server: prediction of thermodynamically stable and evolutionarily conserved RNA structures (2007) *Nucleic Acids Res.* 35: W335-W338.
6. Seo GJ, Chen CJ, Sullivan CS. Merkel cell polyomavirus encodes a microRNA with the ability to autoregulate viral gene expression (2009) *Virology* 383: 183-187.
7. Lee S, Paulson KG, Murchison EP, Afanasiev OK, Alkan C, Leonard JH, Byrd DR, Hannon GJ, Nghiem P. Identification and validation of a novel mature microRNA encoded by the Merkel cell polyomavirus in human Merkel cell carcinomas (2011) *J. Clin. Virol.* 52: 272-275.
8. Cho US, Morrone S, Sablina AA, Arroyo JD, Hahn WC, Xu W. Structural basis of PP2A inhibition by small t antigen (2007) *PLoS Biol.* 5: e202.
9. Fluck MM, Schaffhausen BS. Lessons in signaling and tumorigenesis from polyomavirus middle T antigen (2009) *Microbiol. Mol. Biol. Rev.* 73: 542-563.
10. Topalis D, Andrei G, Snoeck R. The large tumor antigen: a “Swiss Army knife” protein possessing the functions required for the polyomavirus life cycle (2013) *Antiviral Res.* 97: 122-136.
11. Stehle T, Yan Y, Benjamin TL, Harrison SC. Structure of murine polyomavirus complexed with an oligosaccharide receptor fragment (1994) *Nature* 369: 160-163.
12. Griffith JP, Griffith DL, Rayment I, Murakami WT, Caspar DL. Inside polyomavirus at 25-A resolution (1992) *Nature* 355: 652-654.