

LATENCY, EXPRESSION AND SPLICING DURING INFECTION WITH HIV

Scott Sherrill-Mix

A DISSERTATION

in

Genomics and Computational Biology

Presented to the Faculties of the University of Pennsylvania

in

Partial Fulfillment of the Requirements for the

Degree of Doctor of Philosophy

2015

Supervisor of Dissertation:

Frederic D. Bushman, Ph.D., Professor of Microbiology

Graduate Group Chairperson:

Li-San Wang, Ph.D., Associate Professor of Pathology and Laboratory Medicine

Dissertation Committee:

Nancy Zhang, Ph.D. Associate Professor of Statistics

Yoseph Barash, Ph.D., Assistant Professor of Genetics

Kristen Lynch, Ph.D., Professor of Biochemistry and Biophysics

Michael Malim, Ph.D., Professor of Infectious Diseases, King's College London

LATENCY, EXPRESSION AND SPLICING DURING INFECTION WITH HIV

© COPYRIGHT

2015

Scott A. Sherrill-Mix

This work is licensed under the
Creative Commons Attribution
NonCommercial-ShareAlike 3.0
License

To view a copy of this license, visit

<http://creativecommons.org/licenses/by-nc-sa/3.0/>

Dedicated to William Maurer, Gayle Maurer & Michele Sherrill-Mix

ACKNOWLEDGEMENTS

I would like to thank

Rick

Bushman lab Chris wetlab collaborators

Friends in the Bushman lab and classmates in GCB.

My committee—Nancy Zhang, Yosephy Barash, Kristen Lynch and Michael Malim—have provided guidance and encouragement. Many faculty of GCB mentoring and teaching. Hannah Chervitz, Tiffany Barlow, Mali Skotheim, Caitlin Greig and Laurie Zimmerman for managing everything and helping manage the layers of bureaucracy. Funding from the HIV Immune Networks Team (HINT) consortium P01 AI090935 and NRSA computational genomics training grant T32 HG000046.

Ram Myers and Mike James

Xiaofen and Otto

...

ABSTRACT

LATENCY, EXPRESSION AND SPLICING DURING INFECTION WITH HIV

Scott Sherrill-Mix

Frederic D. Bushman, Ph.D.

Over 35 million people are living with human immunodeficiency virus (HIV-1). The mechanisms causing integrated provirus to become latent, the diversity of spliced viral transcripts and the cellular response to infection are not fully characterized and hinder the eradication of HIV-1. We applied high-throughput sequencing to investigate the effects of host chromatin on proviral latency and variation of expression and splicing in both the host and virus during infection.

To evaluate the link between host chromatin and proviral latency, we compared genomic and epigenetic features to HIV-1 integration site data for latent and active provirus from five cell culture models. Latency was associated with chromosomal position within individual models. However, no shared mechanisms of latency were observed between cell culture models. These differences suggest that cell culture models may not completely reflect latency in patients.

We carried out two studies to explore mRNA populations during HIV infection. Single-molecule amplification and sequencing revealed that the clinical isolate HIV_{89.6} produces at least 109 different spliced mRNAs. Viral message populations differed between cell types, between human donors and longitudinally during infection. We then sequenced mRNA from control and HIV_{89.6}-infected primary human T cells. Over 17 percent of cellular genes showed altered activity associated with infection. These gene expression patterns differed from HIV infection in cell lines but paralleled infections in primary cells. Infection with HIV_{89.6} increased intron retention in cellular genes and abundance of RNA from human endogenous retroviruses. We also quantified the

frequency and location of chimeric HIV-host RNAs. These two studies together provided a detailed accounting of both HIV_{89.6} and host expression and alternative splicing.

A more cost-effective method of detecting viral load would aid patients with poor access to healthcare. We developed improved methods for assaying HIV-1 RNA using loop-mediated isothermal amplification based on primers targeting regions of the HIV-1 genome conserved across subtypes. Combined with lab-on-a-chip technology, these techniques allow quantitative measurements of viral load in a point-of-care device targeted to resource-limited settings.

This work disclosed novel HIV-host interactions and developed techniques and knowledge that will aid in the study and management of HIV-1 infection.

TABLE OF CONTENTS

ABSTRACT	v
LIST OF TABLES.....	ix
LIST OF ILLUSTRATIONS	x
CHAPTER 1 : Introduction	1
1.1 Impact of HIV	1
1.2 The HIV virus.....	3
1.3 Integration and latency.....	7
1.4 HIV splicing	7
1.5 Host cell interactions	11
1.6 Repetitive elements	11
1.7 RNA detection	11
CHAPTER 2 : HIV latency and integration site placement in five cell-based models	12
2.1 Abstract.....	12
2.2 Background	13
2.3 Methods.....	14
2.4 Results	19
2.5 Conclusions	32
2.6 Availability of supporting data	34
2.7 Author's contributions	34
2.8 Acknowledgements	35
CHAPTER 3 : Dynamic regulation of HIV-1 mRNA populations analyzed by single-molecule enrichment and long-read sequencing	36
3.1 Abstract.....	36
3.2 Introduction.....	36
3.3 Materials and methods	39
3.4 Results	46
3.5 Discussion	52
3.6 Acknowledgements	56
CHAPTER 4 : Gene activity in primary T cells infected with HIV _{89.6} : intron retention and induction of distinctive genomic repeats	57
4.1 Abstract.....	57
4.2 Background	58
4.3 Methods.....	59
4.4 Results	63
4.5 Discussion	81
4.6 Conclusions	85

4.7 Availability of supporting data	85
4.8 Author's contributions	85
4.9 Acknowledgements	86
4.10 Additional Files	87
 CHAPTER 5 : A reverse transcription loop-mediated isothermal amplification assay optimized to detect multiple HIV subtypes.....	89
5.1 Abstract.....	89
5.2 Introduction.....	89
5.3 Methods.....	91
5.4 Results	93
5.5 Testing different primer designs	94
5.6 Discussion	101
 CHAPTER 6 : Conclusions and future directions	104
6.1 Latency and integration location	104
6.2 HIV-1 alternative splicing	106
6.3 Host expression during HIV infection	106
6.4 LAMP PCR and lab-on-a-chip	106
 APPENDICES	110
A.1 Generalized linear models of changes in use of mutually exclusive HIV-1 splice acceptors	110
A.2 Reproducible report of HIV integration sites and latency analysis	116
 BIBLIOGRAPHY.....	147

LIST OF TABLES

TABLE 2.1 : Integrations from <i>in vitro</i> models of latency.....	18
TABLE 2.2 : Genomic data available for comparison to integration sites.....	20
TABLE 4.1 : Samples and RNA-Seq sequencing coverage	64
TABLE 4.2 : Data used for meta-analysis of expression changes in HIV	65

LIST OF ILLUSTRATIONS

FIGURE 1.1 : The HIV lifecycle	4
FIGURE 1.2 : The HIV-1 genome.....	4
FIGURE 2.1 : Correlations of genomic features and latency.....	22
FIGURE 2.2 : Lasso regressions predicting latency.....	23
FIGURE 2.3 : Cellular expression and latency	25
FIGURE 2.4 : Strand orientation and latency	26
FIGURE 2.5 : Genes and latency	27
FIGURE 2.6 : Alphoid repeats and latency	28
FIGURE 2.7 : Acetylation and latency.....	30
FIGURE 2.8 : Shared expression status between near neighbors.....	31
FIGURE 3.1 : Mapping the splice donors and acceptors of HIV _{89.6}	38
FIGURE 3.2 : Spliced transcripts produced from HIV _{89.6}	45
FIGURE 3.3 : Novel transcripts utilizing acceptor A8c.....	50
FIGURE 3.4 : Temporal, cell type and donor variability in accumulation of HIV-1 messages.....	53
FIGURE 4.1 : Comparisons among studies quantifying cellular gene expres- sion after HIV infection.....	66
FIGURE 4.2 : Comparisons of the effect of HIV infection on gene expression to studies comparing subsets of immune cells	68
FIGURE 4.3 : Changes in the abundance of intronic regions with HIV infection	70
FIGURE 4.4 : Repeat categories enriched upon infection with HIV	72
FIGURE 4.5 : Characteristics of LTR12C sequences associated with induction upon infection with HIV _{89.6}	74
FIGURE 4.6 : Transcription and splicing of the HIV _{89.6} RNA	77
FIGURE 4.7 : Chimeric RNA sequences containing both human and HIV se- quences	80
FIGURE 5.1 : Amplification results for all RT-LAMP primer sets tested	95
FIGURE 5.2 : Subtype-agnostic RT-LAMP primers design	96
FIGURE 5.3 : Performance of the AceIN-26 primer set with different starting RNA concentrations	99
FIGURE 5.4 : Replicate tests of the ACeIN-26 primer set over six HIV subtypes	100
FIGURE 6.1 : Ebola RT-LAMP primers design	108

CHAPTER 1 : Introduction

1.1 Impact of HIV

In 1981, physicians began to notice a mysterious increase, often clustered in men who had sex with men or intravenous drug users, in the occurrences of Kaposi's sarcoma and pneumocystis pneumonia^{1–6}.

Kaposi's sarcoma was, until 1981, a rare cancer in the US found largely in elderly men with Jewish or Mediterranean ancestry⁷. Kaposi's sarcoma had also been seen in immunocompromised individuals^{8–10} and there were suggestions that it was a virus-associated cancer¹¹ although the causative human herpesvirus 8 would not be discovered for another decade^{12,13}.

Pneumocystis pneumonia was known to be caused by infection of the alveoli with the yeast-like fungus *Pneumocystis jirovecii*^{14,15}. Pneumocystis pneumonia was almost exclusively seen only in patients with suppressed immune systems or immune disorders and rarely if ever in immunocompetent individuals¹⁵.

The mechanism for this spike of opportunistic infections was clarified when researchers found severe T cells depletion and decreases in cellular immunity in these patients^{4–6,16,17}. The disease was eventually labeled acquired immunodeficiency syndrome (AIDS). However, the underlying cause remained unclear. Potential transmissions by transfusion^{18–20}, injection drug use^{4,17,21}, maternal transmission²² and both homosexual^{16,23} and heterosexual^{17,24} contact pointed towards an infectious agent. In 1983, a virus later named human immunodeficiency virus type 1 (HIV-1) was isolated from patient samples^{25–28} and soon detected in most immunodeficient patients^{28–31}.

Reports of AIDS and associated opportunistic infections in sub-Saharan Africa soon revealed widespread endemic infection^{32–35} and a great diversity of viruses^{36–41}. Retrospective studies suggested that the virus had been present in Europe and USA for

several decades^{42,43} and circulating for even longer in Africa^{33,44–48}. Isolates of HIV-1 from what is now Kinshasa, in the Democratic Republic of Congo, from as early as 1959 showed already existing viral diversification^{46–48}. Phylogenetic analyses of HIV-1 type M sequences estimate a most recent common ancestor in the early 1900s^{48–53}.

A virus similar to HIV-1 was observed in chimpanzees^{54,55} and survey of wild chimpanzees revealed a likely origin of HIV-1 was from a zoonotic transmission, perhaps from harvesting of chimpanzees for food^{56–61}. The zoonosis most likely occurred in southeast Cameroon^{62–64} from which the virus was transported down the Sangha River⁶⁵ to the city of Kinshasha, where the virus began its global spread^{38,48,53,66}. A combination of social upheaval, mobilization, urbanization and mass vaccination campaigns with unsterilized needles appear to have provided fuel for the growing epidemic^{53,67–69}.

In the early days of the epidemic before there were tests to detect the virus, the first sign of HIV was often the onset of AIDS. Opportunistic infections⁷⁰ and death often followed soon after. The median survival time after diagnosis with AIDS was about 1 year if untreated^{71,72}. Testing for HIV revealed that from the time of infection with HIV, patients had a median survival time of about 10 years without antiretroviral therapy^{73–75}. About 20 years life expectancy in HAART era^{74,76} although catching late and cofounding factors can increase mortality⁷⁶

The successful trial of the reverse transcriptase inhibitor azidothymidine provided the first hope for treatment HIV in 1987^{77–79} but it soon became apparent that the fast mutation rate of HIV^{80–86} and strong selection by drug therapy could quickly create drug-resistant forms of virus in patients receiving single drug therapy^{87–96}. Median survival time from AIDS diagnosis rose to only about 2 years with therapy^{72,78,97,98}. Sequential administration of multiple antiretroviral drugs^{99–102} did not greatly improve prognosis.

Synergistic combinations of antiretroviral drugs^{103–108} and the difficulty of evolving

multiple drug resistant mutations therapy^{109,110} [[more refs here]] meant that therapy using simultaneous combinations of drugs finally began to offer patients more hope of long term survival^{111–119}.

2/3 life expectancy in 2005^{120,121}. Life expectancy of HIV patients receiving consistent antiretrovirals now approaches control population^{122–124} but inflammation and drug toxicity, health care burden still a problem¹²⁵

HIV rebounds quickly after cessation of therapy¹²⁶ Groups and subtypes?

early establishment of latent infection¹²⁷ within 3 days¹²⁸ SAHA induces latent¹²⁹ new drug for latency disulfiram¹³⁰ modest induction of latent provirus¹²⁶ characterizing latent reservoir rare cells present in resting actice and macro¹³¹

1.2 The HIV virus

HIV is an enveloped single strand positive-sense retrovirus, an RNA virus which uses reverse transcription to create a DNA intermediate in host cells^{132,133}. Viral encoded integrase protein discovered¹³⁴ and mapped to 3' end of pol^{135–137} through mutation and loss of function. Mutations at ends also results in defective viruses¹³⁸.

The HIV genome encodes genes for at least two polyproteins and seven proteins:

gag Gag (group specific antigens) is cleaved by viral protease to produce matrix, capsid, nucleocapsid and p6 protein along with two small spacer peptides SP1 and SP2.

MA p17 MA (matrix)

CA p24 CA (capsid) is a myristoylated membrane protein.

NC p7 NC (nucleocapsid)

p6 p6 (protein 6 kda) is a small protein which appears to primarily recruit cellular

Cellular restriction factors

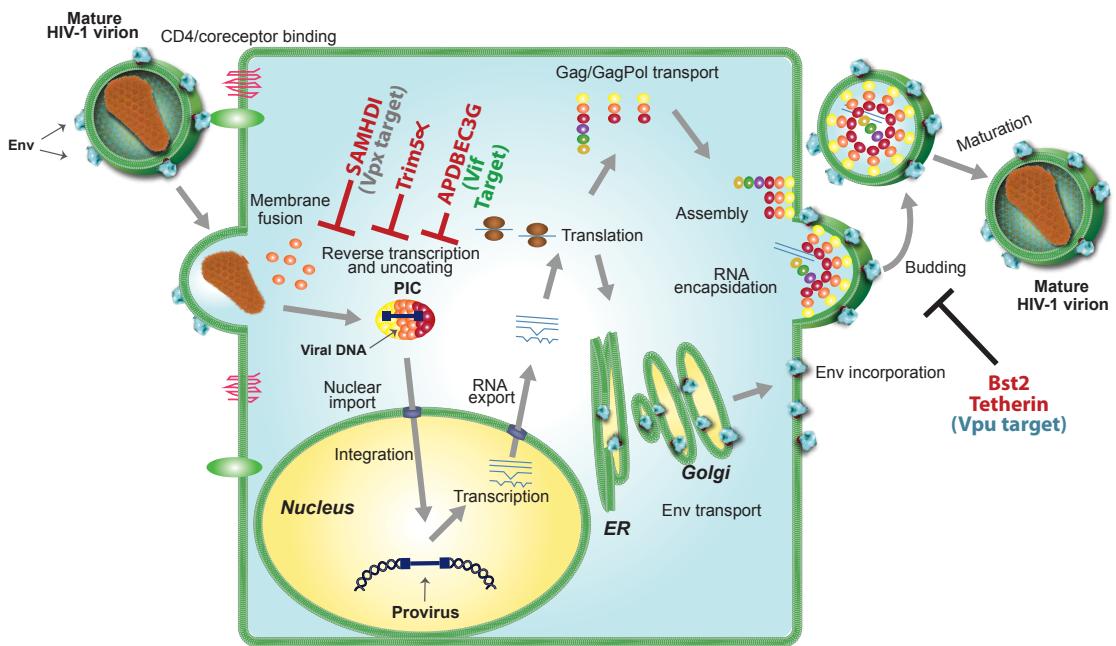


Figure 1.1: The HIV lifecycle/ [[CAPTION HERE]]

[[GENOME FIG]]

Figure 1.2: The HIV-1 genome [[CAPTION HERE]]

proteins to allow virion budding from the cell membrane^{139–141}.

pol Pol (polymerase) is cleaved by viral protease to produce reverse transcriptase, RNaseH, integrase and HIV protease.

RT p51 RT (reverse transcriptase)

RNase H p15 RNase H

IN p31 IN (integrase)

PR PR (protease)

env gp160 Env proteins traps CD4 in the endoplasmic reticulum^{142,143}

tat Tat (transactivator) protein is a transactivator of expression from the HIV-1 long terminal repeat^{144,145}. Virus does not replicate efficiently without this transactivation^{144,145}. and appears to regulate cellular expression such as downregulation of major histocompatibility complex type I expression¹⁴⁶.

rev Rev (regulator of expression of virion proteins) trans-activator protein shuttles between the nucleus and cytoplasm¹⁴⁷ and causes the export of partially spliced and unspliced viral transcripts^{148–152} from the nucleus through recognition of a structured RNA rev response element^{153,154}.

nef Nef (negative factor) is a myristoylated membrane associated protein¹⁵⁵ that is involved in multiple functions. Nef causes endocytosis of the HIV receptors CD4^{156–160} and CCR5¹⁶¹ and viral surveillance major histocompatibility complex molecules^{162–165}. Nef also induces T cell activation through interactions with signaling kinases and the T cell receptor^{166–170}. In contrast, Nef in most other primate lentiviruses inhibits activation and inflammation¹⁷¹ perhaps indicating that the gain of *vpu* in HIV-1 and its simian relatives allowed the loss of the immune inhibition traits of *nef* and contributes to the increased pathogenecity in

these viruses^{172,173}.

vpr Vpr (viral protein R)

vif Vif (virion infectivity factor) counteracts the cellular restriction factor APOBEC3G¹⁷⁴ by excluding APOBEC3G from incorporation into the virion¹⁷⁵ and causing APOBEC3G to be ubiquitinated and degraded^{176–178}. APOBEC3G is otherwise packaged into virions¹⁷⁹ and deaminates the HIV genome during reverse transcription causing G-to-A hypermutation^{179–182}.

vpu Vpu (viral protein U)^{183,184} is a small integral membrane protein which has two known functions; degradation of CD4 and downregulation of BST-2 from the cell membrane [[cites]]. Vpu causes cellular CD4 to be ubiquitinated and degraded^{185,186} which prevents interactions between progeny virus Env and host cell CD4^{159,160,187,188} and superinfection by other viruses¹⁵⁷ while also releasing Env proteins from CD4 interactions in the endoplasmic reticulum^{142,143}. Vpu also counteracts the cellular restriction factor BST-2, which would otherwise interfere with viral budding¹⁸⁹. Vpu does not appear to be found in the virion¹⁹⁰.

long terminal repeats at each end¹⁹¹. integrates in many locations¹⁹¹

The virus was sequenced in 1985^{192–195}.

Splice classes¹⁹³ and others?

retrovirus, two single stranded RNA genomes (recombination), RT, integration,

HIV diversification within a single patient in env loop¹⁹⁶ positive selection^{197,198} positive selection at same sites over time correlate with slower progression^{198,199} 1% distance increase diversity in env per year (decreases later in infection)²⁰⁰

half life of 2 days and 10^9 CD4 T cells per day^{201–203} macrophages? half life 1–4 weeks^{203 204}

1.3 Integration and latency

cells can produce virus without cell death *in vitro*²⁰⁵ first latent HIV²⁰⁶

1.4 HIV splicing

RNA splicing was first observed in adenovirus^{207,208}. Improved understanding of HIV and other viruses offers medical benefits. Although HAART treatments have greatly improved HIV prognosis, long-term survival of HIV patients remains reduced by at least a decade compared to the general population²⁰⁹. In addition HAART does not provide a permanent cure²¹⁰ and the long-term costs of treatment runs into hundreds of thousands of dollars over the lifetime of a patient^{211,212}. Induction or alteration of splicing has been suggested as a potential treatment^{213,214} and differential splicing appears to be one factor limiting cross species infection and the development of animal models of HIV²¹⁵.

Driven by a strong selective pressure for genome compactness^{216–218}, HIV and other lentiviruses subvert host cell alternative splicing pathways to allow tight packing of their genetic information. Through weak splice sites and overlapping reading frames (Figure ??), the virus manages to produce precise quantities of at least nine proteins and polyproteins from its single transcription start site and less than 10 kb genome²¹⁹.

As such an integral part of the virus life cycle^{220,221}, alteration of splicing poses a tempting therapeutic target. Inhibition of cellular splicing factors reduces viral reproduction in many genome-wide siRNA screens^{222–224} and several members of the spliceosome interact with viral proteins in affinity pulldowns²²⁵. Open reading frames in uncharacterized transcripts appear to produce epitopes useful for vaccine development²²⁶. Potential treatments altering viral splicing through small molecule inhibitors^{213,227} and gene therapy^{214,228} have restricted viral replication *in vitro*. However without methods to quantify viral splicing or a thorough quantification of splicing under varying conditions, the development of such treatments remains limited.

Viral proteins also interact with components of the cellular splicing complex^{225,229,230}. These interactions have been reported to change splicing in viral^{230–232} and cellular transcripts^{233,234} and raise the possibility that the virus has evolved to alter host splicing. A genome-wide study of changes in cellular splicing during HIV infection would greatly clarify this hypothesis but no such study has been performed.

Alternative splicing, the differential inclusion of exons and removal of introns from primary mRNA transcripts, allows rapid evolution of protein segments^{235–237} and drastic increases in the number of proteins generated by a single DNA sequence²³⁸. Many viruses subvert the splicing machinery of their eukaryotic hosts to modify their viral mRNA²³⁹.

In particular, it has previously been reported that HIV utilizes alternative splicing to generate more than 40 mRNA transcripts encoding at least 9 proteins and polyproteins from a genome smaller than 10kb²⁴⁰. A specific progression of viral transcripts appear in the cytoplasm of the host cell as infection progresses allowing a shift from regulatory protein production in early infection into virion production in late infection^{220,221,241}. Because HIV has only a single transcription start site, these transcriptional changes are driven by alternative splicing²¹⁹.

Although it plays such an essential role for the virus, only a single detailed census of viral splicing has been reported²⁴⁰. Due to limitations in technology, this study was limited to only the most abundant transcripts in lab-adapted HIV strains in cell culture²⁴⁰. Yet rare transcripts may play an important role in immune response²²⁶ and encode unknown proteins²⁴²; lab adapted HIV can differ markedly from viruses actually found in patients²⁴³; cell cultures often do not reflect *in vivo* conditions²⁴⁴; and splicing can vary between humans^{245,246} and cell types^{247,248}. Without a fuller characterization of transcripts under these relevant conditions, many aspects of viral splicing will remain poorly understood.

Alternative splicing may also play an unappreciated role in HIV-host interactions. Viral proteins interact with the splicing complex^{225,229,230} and alter splicing of some cellular transcripts^{233,234}. Yet, although infection has been shown to cause genome-wide changes in the expression of cellular genes^{249–253}, no genome-wide study of cellular alternative splicing during HIV infection has ever been reported. Such a genome-wide study of splicing changes might reveal a distortion of diverse cellular splicing which is adaptively advantageous to the virus.

Current sequencing advancements allow a much broader and deeper investigation of viral splicing. Targeted amplification with RainDance droplet PCR offer the potential to reduce size bias inherent in bulk PCR²⁵⁴. RNA-seq with Illumina sequencing allows extremely deep sequencing of cellular and viral transcripts with billions of bases of short read sequence^{255,256}. Single molecule sequencing with Pacific Biosciences provides reads approaching 20,000 bases^{257,258} that could characterize entire viral transcripts in one continuous read. By combining these technologies, viral and cellular transcripts could be interrogated to an unprecedented level.

A better understanding of viral splicing and viral effects on host splicing may bring therapeutic benefits. siRNA inhibition of splicing factors reduces HIV replication in many genome-wide screens^{222–224}. Alteration of viral splicing through small molecule inhibitor of SR protein kinases²¹³ and Splicing Factor 2²²⁷, shRNA against spliceosomal U7 snRNP²²⁸ and expression of modified spliceosomal U1 snRNP²¹⁴ show treatment potential *in vitro*. In addition, rare uncharacterized HIV transcripts and their encoded proteins appear to produce potent immune response in HIV patients²²⁶ thus offering potential targets for vaccine development. Yet without methods to characterize viral RNA and measure the effects of treatments on viral splicing, further development is inhibited.

Inclusion and exclusion of a particular stretch of RNA into an mRNA is determined by a balance of RNA secondary structure^{259–261}, chromatin structure²⁶², nucleosome

positioning²⁶³, histone marks²⁶⁴, previous splicings²⁶⁵, order of intron removal^{266,267} and enhancers²⁶⁸ and suppressors²⁶⁹ that bind specific motifs²⁷⁰. Together these factors create a precise controllable splicing code^{248,271,272}.

In HIV, splicing occurs between at least four splice donors and eight splice acceptors²¹⁹. Two splice donors, D1 and D4, are relatively strong while the remaining donors and all acceptors are fairly weak²⁷³. Several exonic splicing silencers^{274,275} and exon splicing enhancers^{276,277} and a single intronic splicing silencer²⁷⁸ in the viral genome interact with many human splicing factors, including hnRNPs A1^{275,278} H, F, 2H9, and A2²⁶⁰ and SR proteins SRp40^{276,279}, SRp75²⁷⁹, ASF/SF2²⁷⁶ and SC35²⁶⁰, to alter viral splicing^{219,280}.

Several viral proteins affect mRNA abundances. Rev causes export of unspliced viral mRNA that would otherwise be trapped in the nucleus²⁸¹ to be exported^{239,282} and may also interact with splicing factors to alter viral splicing²²⁹. The HIV protein Tat is best known for its transactivation of viral transcription^{144,283} and triggering apoptosis in uninfected cells^{284,285} but Tat also appears to independently affect alternative splicing of viral transcripts^{230–232,286}. Viral protein Vpr is known to cause cell cycle arrest²⁸⁷ and mediate nuclear import of the viral preintegration complex²⁸⁸. Vpr also appears to alter alternative splicing of some cellular transcripts^{233,234} and interact with the SMN complex²²⁵, which assembles spliceosomal snRNP²⁸⁹. Although all three of these proteins modify viral splicing, whether they also cause widespread alterations in cellular splicing is unknown.

Despite the critical role alternative splicing plays in viral replication, no genome-wide studies of lentiviral effects on cellular splicing or detailed censuses of viral alternative splicing in biologically relevant settings have been published.

RNA-seq offers a much broader view of alternative splicing than previously possible^{290,291} but Illumina sequencing has not yet been applied to the study of differential

splicing in host RNA of HIV-infected cells. There have been many studies of cellular expression using microarrays^{249–252,286} and Sage^{292,293} but only a single study using Illumina RNA-seq and alternative splicing changes were not reported²⁵³. Thus a potentially significant aspect of HIV-host interactions remains unknown.

The most extensive survey of HIV transcripts to date was published in 1993²⁴⁰. Technology at the time necessitated the use of Northern blots and RNA protection assays²⁴⁰ which can not distinguish multiple similarly sized transcripts or detect rare transcripts. This study also focused on a single lab adapted HIV_{NL4-3} strain in HeLa cell culture.

Many previous studies of viral splicing have used lab-adapted strains of HIV which often differ from patient isolates²⁴³ in cell cultures which often differ from primary cells²⁴⁴. Selection under cell culture conditions may quickly alter splicing patterns to down regulate proteins unneeded *in vitro*. Characterization of alternative splicing in biologically relevant cell types infected with clinical isolates of HIV are sorely needed.

1.5 Host cell interactions

1.6 Repetitive elements

1.7 RNA detection

First HIV antibody test^{30,31} First DNA test²⁹⁴

CHAPTER 2 : HIV latency and integration site placement in five cell-based models

2.1 Abstract

Background: HIV infection can be treated effectively with antiretroviral agents, but the persistence of a latent reservoir of integrated proviruses prevents eradication of HIV from infected individuals. The chromosomal environment of integrated proviruses has been proposed to influence HIV latency, but the determinants of transcriptional repression have not been fully clarified, and it is unclear whether the same molecular mechanisms drive latency in different cell culture models.

Results: Here we compare data from five different *in vitro* models of latency based on primary human T cells or a T cell line. Cells were infected *in vitro* and separated into fractions containing proviruses that were either expressed or silent/inducible, and integration site populations sequenced from each. We compared the locations of 6,252 expressed proviruses to those of 6,184 silent/inducible proviruses with respect to 140 forms of genomic annotation, many analyzed over chromosomal intervals of multiple lengths. A regularized logistic regression model linking proviral expression status to genomic features revealed no predictors of latency that performed better than chance, though several genomic features were significantly associated with proviral expression in individual models. Proviruses in the same chromosomal region did tend to share the same expressed or silent/inducible status if they were from the same cell culture model, but not if they were from different models.

Conclusions: The silent/inducible phenotype appears to be associated with chromosomal position, but the molecular basis is not fully clarified and may differ among *in vitro* models of latency.

2.2 Background

Highly active antiretroviral therapy (HAART) can suppress HIV-1 replication in infected patients, but the ability of HIV to persist as an inducible reservoir of latent proviruses^{131,295,296} obstructs eradication of the virus and functional cure²¹⁰. These latent proviruses are long lived^{297,298} and relatively invisible to the immune system^{131,204}. The potential for even a single virus to restart infection despite successful antiviral therapy means that it may be necessary to eliminate all latent proviruses to eradicate HIV from an infected person.

After integration, a positive feedback loop of Tat transactivation appears to partition proviral gene activity into either of two stable states^{299–301}—abundant Tat driving high proviral expression or little Tat leading to quiescent latency. Similar to the positional effect variegation observed in fruit fly chromosomal rearrangements^{302,303}, studies on cell clones with single integrations show that differing integration sites can have large differences in proviral expression^{304–306}. These data suggest that integration site location, along with the cellular environment^{306–309}, influences the balance between latency and proviral expression.

Associations between latency and genomic features have also been reported in collections of integration sites from cell culture models although the consistency of these effects across model systems and their relationships to latency in patients remains uncertain. Lewinski et al.³¹⁰ reported that proviruses integrated in gene deserts, alphoid repeats and highly expressed genes are more likely to have low expression. Shan et al.³¹¹ reported an association between latency and integration in the same transcriptional orientation as host genes. Pace et al.³¹² found that silent and expressed provirus integration sites differed in the abundance and expression levels of nearby genes, GC content, CpG islands and alphoid repeats. In model systems with defined integration sites, Lenasi et al.³¹³ reported decreased and Han et al.³¹⁴ reported increased viral transcription when the provirus is downstream of a highly expressed host gene.

Cell-based models of latency are important for many aspects of HIV research, including screening small molecules that can reverse latency and potentially allow eradication^{315,316}. Location-driven differences in expression are preserved even after demethylation and histone deacetylase treatment³⁰⁴, which suggests that integration location has the potential to confound “shock and kill” anti-latency treatments^{317,318}. A greater understanding of the effects of integration site location on latency could thus affect antiretroviral development.

To search for features of integration site associated with latency, we generated a set of inducible and expressed integration sites using a primary central memory CD4⁺ T cell model of latency^{319,320}, collected four previously reported integration site datasets and modeled the effects of genomic features near the integration site on the expression status of these proviruses. Although some genomic features associated with latency in individual models, no feature was consistently associated with proviral expression across all five cell culture models. However, closely neighboring proviruses within the same cellular model shared the same latency status much more often than expected by chance suggesting that chromosomal position of integration affects latency but that the mechanism remains unclear or differs between cell culture models. Thus these data help inform the design of experiments in HIV eradication research.

2.3 Methods

2.3.1 Integration sites

Naive CD4⁺ T cells were purified by negative selection from peripheral blood mononuclear cells. The cells were activated with anti-CD3 and anti-CD28 (+TGF-beta, anti-IL-12, and anti-IL-4) to generate “non-polarized” cells (the in vitro equivalent of central memory T cells). Five days after isolation, cells were infected with an NL4-3-based virus with GFP in place of Nef and the LAI envelope (X4) provided in trans at a concentration of 500 ng of p24 as measured by ELISA per million cells. Based on previous experience

with this model, this amount of p24 should produce an MOI of approximately 0.15. Cells were cultured in the presence of IL-2. Two days post-infection, cells were sorted for GFP+; this active population expresses GFP even when treated with flavopiridol, although for this study they were not treated. The inducible population was the set of GFP negative cells from the initial sort that, 9 days post-infection, were activated with anti-CD3 and anti-CD28 and sorted for GFP production.

Genomic DNA from the inducible and expressed populations was digested with MseI, ligated to an adapter, and amplified by ligation-mediated PCR essentially as in Wu et al.³²¹ and Mitchell et al.³²² except that the nested PCR primers included sequence for the Ion Torrent P1 adapter and adapter A sequence with a 5 base barcode sequence specific to the inducible or expressed conditions. Amplicons were sequenced using an Ion Torrent Personal Genome Machine (PGM) according to manufacturer's instructions using an Ion 316 chip and the Ion PGM 200 Sequencing kit (Life Technologies). The sequence reads were sorted into samples by barcode. All reads were required to match the expected 5' sequence with a Levenshtein edit distance less than 3 from the expected barcode, 5' primer and HIV long terminal repeat (LTR). The 5' primer and HIV sequence, along with the 3' primer if present, were trimmed from the read. Sequences with less than 24 bases remaining or containing any eight base window with an average quality less than 15 were discarded. Duplicate reads and reads forming an exact substring of a longer read were removed.

2.3.2 Analysis

All statistical analysis was performed in R 2.15.2³²³. The analyses are described in a reproducible report (Appendix A.2). The annotated integration site data necessary to perform the analyses and the compilable code to generate this reproducible report are provided as supplemental information³²⁴. The new Central Memory CD4⁺ data set was analyzed as in Berry et al.³²⁵. The integration patterns appeared similar to previously reported HIV integration site datasets³²⁶.

2.3.3 Previously published data

We collected integration sites from three previously reported studies (Table 2.1), for a total of four expressed versus silent/inducible pairs of samples. These studies used primary CD4⁺ T cells or Jurkat cells infected with HIV or HIV-derived constructs as cell culture models of latency. Flow cytometry allowed cells expressing viral encoded proteins to be sorted from non-expressing cells. In two of the studies, these non-expressing populations were stimulated to ensure that the provirus could be aroused from latency. Specific differences in protocol between the study sets are summarized below.

Jurkat. Lewinski et al.³¹⁰ infected Jurkat cells with a VSV-G pseudotyped, GFP-expressing pEV731 HIV construct (LTR-Tat-IRES-GFP)³⁰⁴ at an MOI of 0.1. The cells were sorted into GFP+ and GFP- two to four days after infection. GFP+ cells were sorted again two weeks after infection and cells that were again GFP+ were collected for integration site sequencing. GFP- cells were sorted for GFP negativity twice more than stimulated with TNFalpha. Cells that were GFP+ after stimulation were collected for integration site sequencing. DNA was digested with MseI or a combination of NheI, SpeI and XbaI, ligated to adapters for nested PCR, amplified and sequenced by Sanger capillary electrophoresis.

Bcl-2 transduced CD4⁺. Shan et al.³¹¹ transduced CD4⁺ T cells with Bcl-2, costimulated with bound anti-CD3 and soluble anti-CD28 antibodies, interleukin-2 and T cell growth factor and then infected with X4-pseudotyped GFP-expressing NL4-3- δ 6-drEGFP construct³²⁷ at an MOI of less than 0.1. DNA was extracted, digested with PstI and circularized³²⁸. HIV-human junctions were amplified by reverse PCR and sequenced using Sanger capillary electrophoresis.

Active CD4⁺ & Resting CD4⁺. Pace et al.³¹² spinoculated CD4⁺ T cells with HIV NL4-3 at an MOI of 0.1. After 96 hours, the cells were stained for intracellular Gag CD25,

CD69 and HLA-DR and sorted into four subpopulations based on activation state and Gag expression; activated Gag-, activated Gag+, resting Gag- and resting Gag+. The ability of the viruses to reactivate was not tested although previous studies have shown that the majority are likely inducible³²⁹. Genomic DNA was extracted and digested with restriction enzymes MseI and Tsp509 and ligated to adapters. Proviral LTR-host genome junctions were sequenced by 454 pyrosequencing after nested PCR.

All datasets were processed using the hiReadsProcessor R package³³⁰. Adaptor trimmed reads were aligned to UCSC freeze hg19 using BLAT³³¹. Genomic alignments were scored and required to start within the first three bases of a read with 98% identity. Alignments for a given read with a BLAT score less than the maximum score for that read were discarded. Reads giving rise to multiple best scoring genomic alignments were excluded, while reads with a single best hit were dereplicated and converged if within 5bp of each other. The Bcl-2 transduced CD4⁺ sample was sequenced from U3 in the 5' HIV LTR while the other samples were sequenced from U5 in the 3' LTR. To account for the 5 base duplication of host DNA caused by HIV integration, the chromosomal coordinates of the Bcl-2 transduced CD4⁺ sample were adjusted by ±4 bases.

To allow for alignment difficulties in the analysis of genomic repeats, reads with multiple best scoring alignments, along with the single best hit reads used above, were included in the repeat analyses. If any best scoring alignment for a read fell within a repeat, then that read was considered to map to that repeat.

2.3.4 Genomic features

A total of 140 whole genome features for CD4⁺ T-cells were gathered from data sources indicated in Table 2.2. For features encoded as peaks or hotspots, the log of the distance of each integration site to the nearest border was used for modeling. Integration sites from HIV 89.6 infection in primary CD4⁺ T cells³³² were used to count nearby

Title	Cell type	Virus	Time of harvest after infection	Sequencing	Generation of expressed vs. silent/inducible	Citation	Silent/inducible unique sites	Expressed unique sites
Jurkat	Jurkat cells	HIV vector pEV731 (LTR-Tat-IRES-GFP)	2 weeks	Sanger	TNF α , GFP expression	Lewinski et al. ³¹⁰	463 inducible	643
Bcl-2 transduced CD4 $^{+}$	Primary CD4 $^{+}$ T cells (Bcl-2 transduced)	HIV NL4-3- δ 6-drEGFP (inactivated <i>gag</i> , <i>vif</i> , <i>vpr</i> , <i>vpu</i> , <i>nef</i> and <i>env</i> replaced by GFP)	3 days + 3-4 weeks + 3 days	Sanger	anti-CD3, anti-CD28 antibodies, GFP expression	Shan et al. ³¹¹	446 inducible	273
Active CD4 $^{+}$	Primary active CD4 $^{+}$ T cells	HIV NL4-3	3 days	454	high vs. low Gag	Pace et al. ³¹²	1604 silent	1274
Resting CD4 $^{+}$	Primary resting CD4 $^{+}$ T cells	HIV NL4-3	3 days	454	high vs. low Gag	Pace et al. ³¹²	1942 silent	784
Central Memory CD4 $^{+}$	Primary central memory CD4 $^{+}$ T cells	HIV NL4-3 Δ Nef GFP	2 days/9 days	Ion-Torrent	anti-CD3, anti-CD28 antibodies, GFP expression	This paper	1729 inducible	3278

Table 2.1: HIV-1 integration datasets from *in vitro* models of latency where the proviruses were determined to be silent/inducible or expressed

integrations and determine a ± 20 bp position weight matrix for integration targets. Illumina RNA-Seq from active CD4 $^{+}$ cells (Chapter 4) was used to estimate raw cellular expression and fragments per kilobase of transcript per million mapped reads for genes as calculated by Cufflinks²⁹⁰. For sequence-based data like RNA-Seq and ChIP-Seq, the number of reads aligned within a $\pm 50, 500, 5,000$ 50,000 and 500,000 bp windows of each integration site were counted and log transformed. In addition, chromatin state classifications derived from a hidden Markov model based on histone marks and a few binding factors³³³ were included as binary variables. All data from previous genomic freezes were converted to hg19 using liftover³³⁴.

2.4 Results

The combination of integration site data newly reported here (set named “Central Memory CD4 $^{+}$ ”) with previously published data (sets named “Jurkat”, “Bcl-2 transduced CD4 $^{+}$ ”, “Active CD4 $^{+}$ ”, and “Resting CD4 $^{+}$ ”) provides a collection of 12,436 integration sites (Table 2.1) where the expression status of the provirus—silent/inducible or expressed—is known. In three of the datasets, Jurkat, Central Memory CD4 $^{+}$ and Bcl-2 transduced CD4 $^{+}$, the proviruses were sorted based on inducibility. In the Resting CD4 $^{+}$ and Active CD4 $^{+}$ datasets, cells were sorted only based on proviral expression. Previous studies have shown that most silent proviruses in this model system are inducible³²⁹.

2.4.1 Global model

If a genomic feature and latency are monotonically related then we should be able to detect this relationship using Spearman rank correlation. In addition if a feature has a consistent effect across models we should see a consistent pattern in the direction of correlation. A simple first look for correlation between genomic features (Table 2.2) and latency status yielded inconsistent results among the five samples with no variables having a significant Spearman rank correlation across all, or even four out of five, of

Group	Type	Source	Number	Types
T cell expression	RNA-Seq	Chapter 4	1	RNA
Jurkat expression	RNA-Seq	Encode ³³⁵	1	wgEncodeHudsonalphaRnaSeq
Integration sites	Locations	Berry et al. ³³²	1	sites
DNase sensitivity	DNA-Seq/peaks	Encode ³³⁵	1	wgEncodeOpenChromDnase
Methylation	DNA-Seq	336	1	Methyl
CpG	Locations	UCSC ³³⁷	1	cpgIslandExt
Sequence-based	Continuous	—	4	% GC, HIV PWM score, distance to centrosome, chromosomal position
Repeats	Locations	UCSC ³³⁷	16	DNA, LINE, Low_complexity, LTR, Other, RC, RNA, rRNA, Satellite, scRNA, Simple_repeat, SINE, snRNA, srpRNA, tRNA, alphoid
Histone features	ChIP-Seq/Peaks	Wang et al. ³³⁸	18	H2AK5ac, H2AK9ac, H2BK120ac, H2BK12ac, H2BK20ac, H2BK5ac, H3K14ac, H3K18ac, H3K23ac, H3K27ac, H3K36ac, H3K4ac, H3K9ac, H4K12ac, H4K16ac, H4K5ac, H4K8ac, H4K91ac
Histone features	ChIP-Seq/Peaks	Barski et al. ³³⁹	23	CTCF, H2AZ, H2BK5me1, H3K27me1, H3K27me2, H3K27me3, H3K36me1, H3K36me3, H3K4me1, H3K4me2, H3K4me3, H3K79me1, H3K79me2, H3K79me3, H3K9me1, H3K9me2, H3K9me3, H3R2me1, H3R2me2, H4K20me1, H4K20me3, H4R3me2, PolII
Chromatin state	Binary	Ernst and Kellis ³³³	51	state ₁ ,state ₂ ,...,state ₅₁
HATs and HDACs	ChIP-Seq	Wang et al. ³⁴⁰	11	Resting-HDAC1, Resting-HDAC2, Resting-HDAC3, Resting-HDAC6, Resting-p300, Resting-CBP, Resting-MOF, Resting-PCAF, Resting-Tip60, Active-HDAC6, Active-Tip60
Nucleosome	ChIP-Seq	Schones et al. ³⁴¹	2	Resting-Nucleosomes, Active Nucleosomes
UCSC genes	Locations	Hsu et al. ³⁴²	4	in gene, in gene (same strand), gene count, distance to nearest gene, in exon, in intron

Table 2.2: Genomic data available for comparison to HIV integration sites

the samples (Figure 2.1). This suggests that there is not a consistent simple monotonic relationship between the genomic variable and latency, or that any such correlations are modest and not detectable across all studies given the available statistical power. We return to some of the stronger trends below.

To investigate whether a combination of variables may affect latency, we fit a lasso-regularized logistic regression, as implemented in the R package `glmnet`³⁴³, to predict latency using the genomic variables. The relationship between silent/inducible status and each genomic variable was allowed to vary between models by including the interaction of genomic features with dummy variables indicating cellular model. The λ smoothing parameter of the lasso regression was optimized by finding the λ with lowest classification error in 480-fold cross validation and finding the simplest model with misclassification error within one standard error.

The proportion of silent/inducible sites varied between the samples. To avoid the model overfitting on this source of variation, an indicator variable for each sample was included in the base model. The base model with no genomic variables was selected as the best model by cross validation (Figure 2.2A). This suggest that there is not a consistent linear relationship between an additive combination of genomic variables and latency across all models.

When each dataset was fit individually with leave-one-out cross validation, improvements in cross-validated misclassification error were only observed in the Active CD4⁺ (5.8% decrease in misclassification error, standard error: 2.1) and Jurkat (6.7% decrease in misclassification error, standard error: 3.5) samples (Figure 2.2B-F). There was no overlap in variables selected for the Active CD4⁺ and Jurkat samples.

Finding little global association between latency and genomic features, we investigated whether predictors of latency reported previously by single studies were consistently associated with latency across studies.

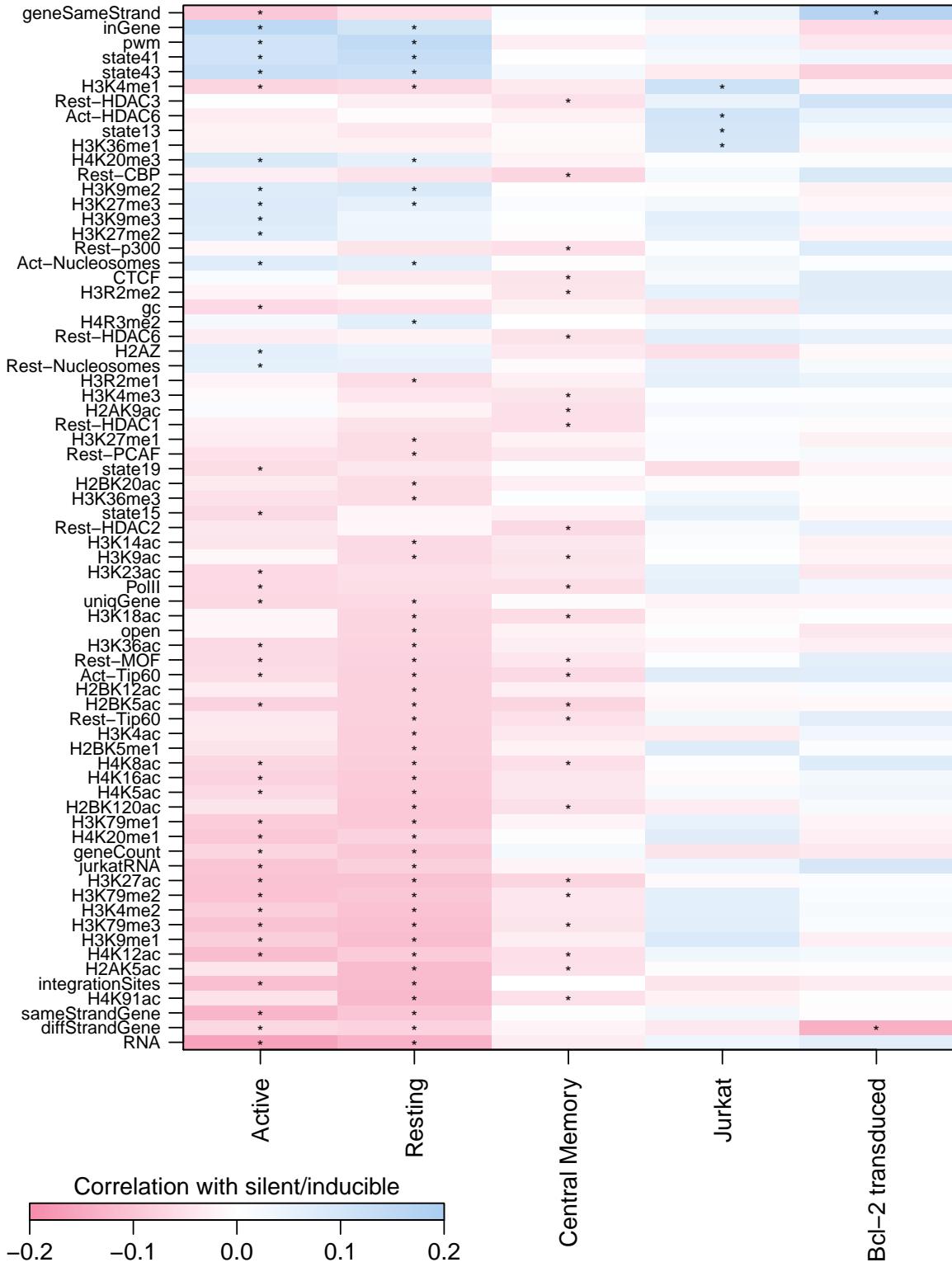


Figure 2.1: Spearman rank correlation between proviral expression status and genomic features. Only genomic features with at least one correlation with latency with a false discovery rate q -value < 0.01 (marked by asterisks) are shown.

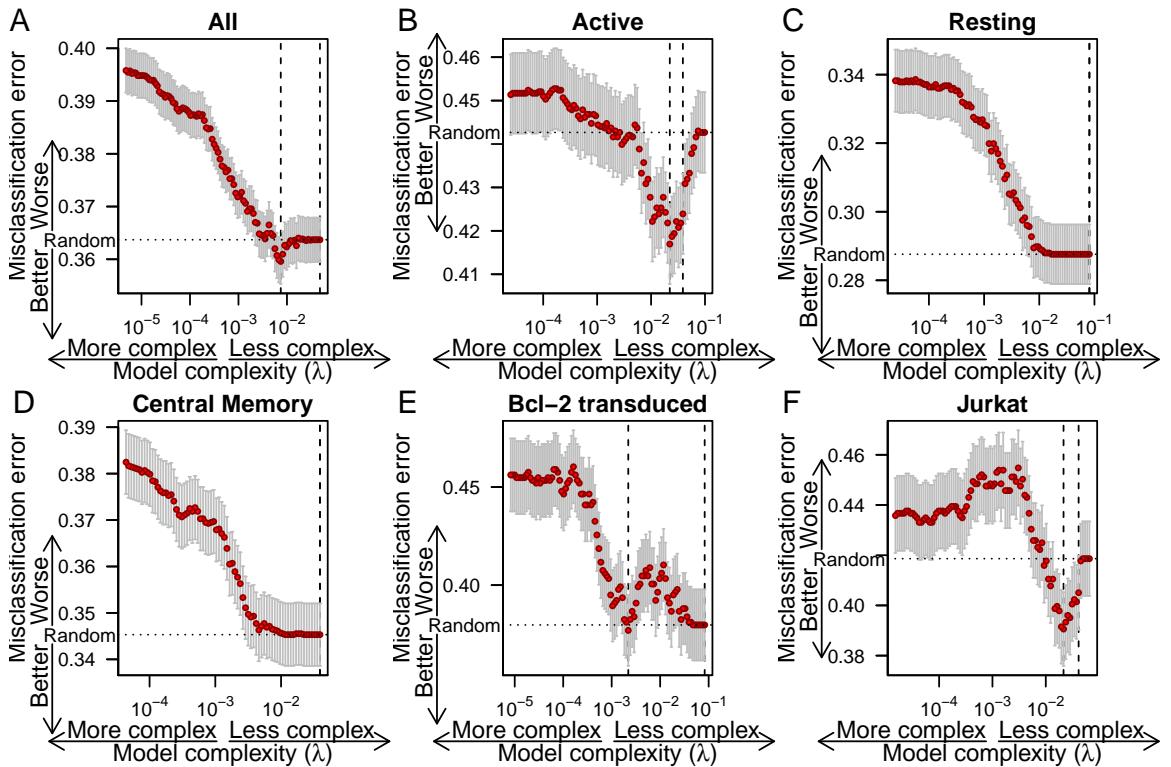


Figure 2.2: Misclassification error from cross validation for lasso regressions of silent/inducible status on genomic features as a function of λ , the regularization coefficient for the lasso regression, for all cell culture models combined and each individual cell culture model. The number of variables included and size of coefficients in the model increases to the left. Whiskers show the standard error of mean misclassification error. Dashed vertical lines indicate the minimum misclassification error and the simplest model within one standard error. Dotted horizontal line indicates the misclassification error expected from random guessing.

2.4.2 Cellular transcription

Model systems with defined integration sites show upstream transcription can interfere with viral transcription³⁴⁴ and that cellular transcription in the same orientation may interfere with viral transcription³¹³ or increase viral transcription³¹⁴ and in opposite orientations may decrease transcription³¹⁴. In integration site studies, integration outside genes appears to increase latency³¹⁰ but high transcription of nearby host cell genes may cause increased latency^{310,311}. In addition, Tat or other viral proteins may affect cellular transcription^{253,345}.

To look at transcription and latency, we ran a logistic regression of silent/inducible status on a quartic function of RNA expression, as determined by RNA-Seq reads within 5,000 bases in Jurkat cells for the Jurkat sample or CD4⁺ T cells for the remaining samples, interacted with indicator variables encoding cell culture model. There appears to be little agreement between samples (Figure 2.3). The Resting CD4⁺ and Active CD4⁺ datasets show an enrichment in silent proviruses in regions with low gene expression. The other three studies show the opposite or no relationship for low expression regions. The two samples showing increased silence in areas of low expression (Resting CD4⁺ and Active CD4⁺) are from a study that did not check whether inactive viruses could be activated. One possible explanation is that regions with low gene transcription may harbor proviruses that are not easily activated, though some other discrepancy between *in vitro* systems could also explain the difference. Both the Jurkat and Active CD4⁺ samples appear to increase in latency with increasing expression while the remaining three studies did not show a strong trend.

2.4.3 Orientation bias

Shan et al.³¹¹ reported that inducible proviruses were oriented in the same strand as the host cell genes into which they had integrated more often than chance. This orientation bias was still reproduced after our reprocessing of the Bcl-2 transduced

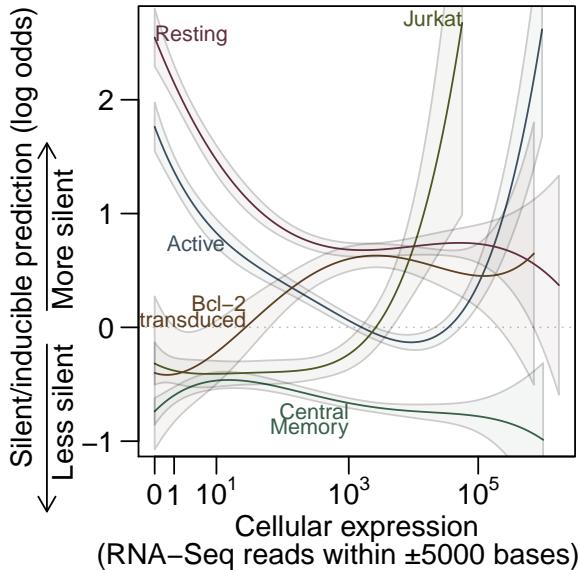


Figure 2.3: Predictions from a logistic regression of silent/inducible status on cellular RNA expression. High y-axis values are predicted to be silent/inducible. Dashed line shows where equal odds of silent/inducible and expressed are predicted. Solid lines show predictions from the regression for each sample and shaded regions indicate one standard error from the modeled predictions.

CD4⁺ sample from Shan et al.³¹¹. However, the proportion of provirus oriented in the same strand as host genes did not differ significantly from 50% in the other samples (Figure ??). Perhaps orientation bias and transcriptional interference are especially sensitive to parameters of the model system.

2.4.4 Gene deserts

Lewinski et al.³¹⁰ reported increased latency in gene deserts. In the collected data, integration outside known genes was associated with latency (Fisher's exact test, $p < 10^{-6}$). This seemed to largely be driven by the Active CD4⁺ and Resting CD4⁺ samples with significant association found individually in only those two samples (both $p < 10^{-8}$) and no significant association observed in the other three samples (Figure 2.5A). Looking only at integration sites outside genes, silent sites in the Resting CD4⁺ sample had a mean distance to the nearest gene 2.5 times greater than that of expressed sites (95% CI: 2.2–6.2×, $p < 10^{-6}$, Welch two sample t-test on log transformed distance) (Figure 2.5B). The Active CD4⁺ sample had a small difference that did not survive Bonferroni correction.

Lewinski et al.³¹⁰ also reported decreased latency near CpG islands and reasoned

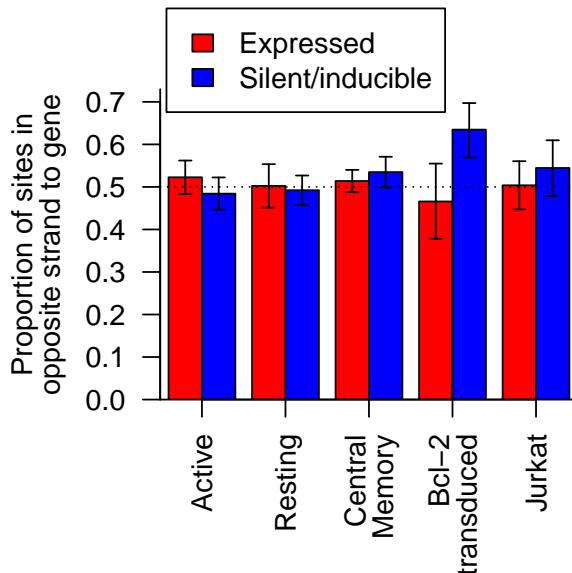


Figure 2.4: The proportion of provirus integrated in the opposite strand compared to cellular genes in silent/inducible (blue) and expressed (red) samples. Error bars show the 95% Clopper-Pearson binomial confidence interval.

this was tied to the increased latency in gene deserts. In the Resting CD4⁺ sample, silent sites were on average further from CpG islands than expressed sites (Bonferroni corrected Welch's two sample T test, $p = 0.006$), but there was no significant relationship between silent/inducible status and log distance to CpG island after Bonferroni correction if the integration site's location inside or outside of a gene was accounted for first (analysis of deviance).

2.4.5 Alphoid repeats

Alphoid repeats are repetitive DNA sequences found largely in the heterochromatin of centromeres³⁴⁶. Integration near heterochromatic alphoid repeats has been reported to associate with latency^{305,310,312}. Looking only at uniquely mapping sites, there was no statistically significant association between latency and location inside an alphoid repeat in pooled or individual samples (Fisher's exact test).

Since alphoid repeats are both problematic to assemble in genomes and difficult to map onto, we reasoned that some alphoid hits might be lost or miscounted in the filtering procedures of the standard workup. To counteract this, we treated each sequence read as an independent observation of a proviral integration and included sequence

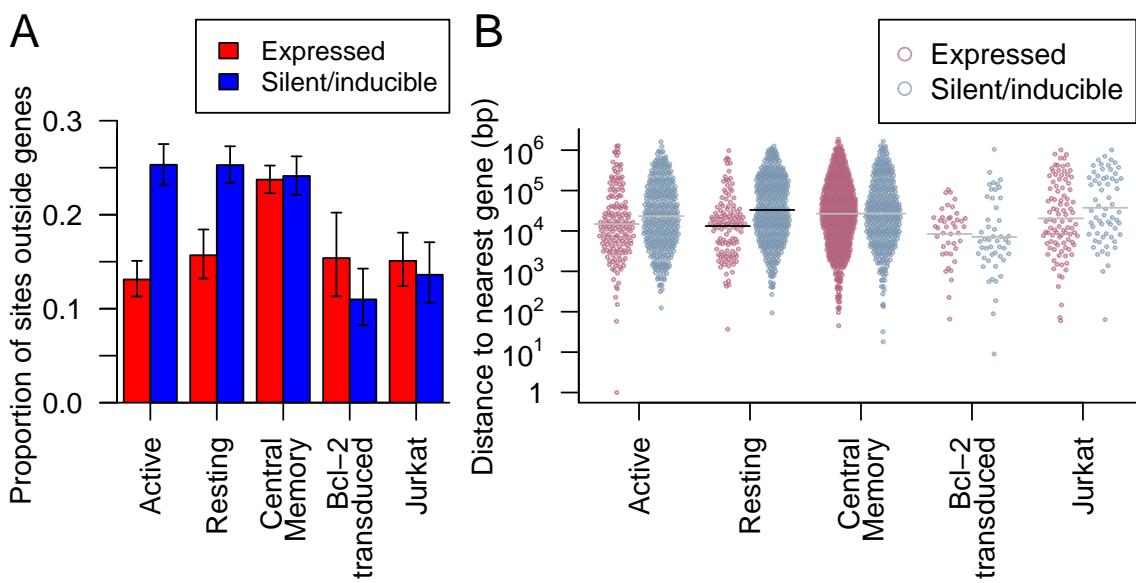


Figure 2.5: (A) The proportion of provirus integrated outside genes in silent/inducible (blue) and expressed (red) samples. Error bars show the 95% Clopper-Pearson binomial confidence interval. (B) The nearest distance to any gene for integration sites (points) outside genes in the five samples. Points are spread in proportion to kernel density estimates. Horizontal lines indicate sample means where there was a significant difference in means between silent/inducible and expressed provirus (black) or no significant difference (grey).

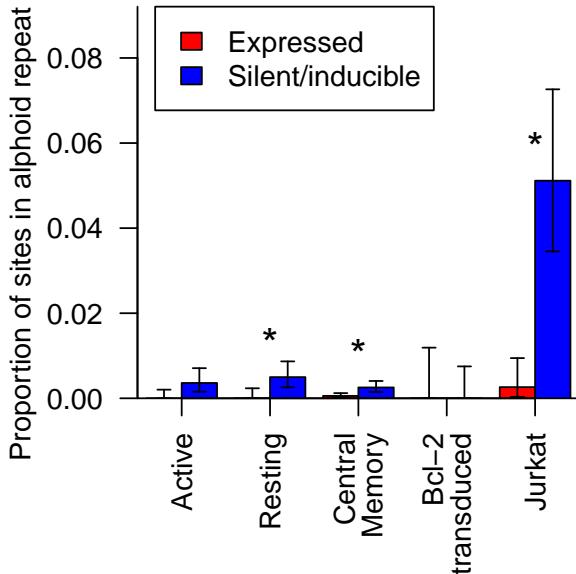


Figure 2.6: The proportion of integration sites with matches in alphoid repeats in silent/inducible (blue) and expressed (red) cells in five samples. Error bars show the 95% Clopper-Pearson binomial confidence interval. Asterisks indicate significant associations between integrations within an alphoid repeat and proviral expression status (Bonferroni corrected Fisher’s exact test $p < 0.05$).

reads with more than one best scoring alignment. For multiply aligned reads, we considered the read to have been inside an alphoid repeat if any of its best scoring alignments fell within a repeat. We found 74 reads with potential alphoid mappings. Integration inside alphoid repeats was significantly associated with the expression status of a provirus in the Resting CD4 $^{+}$, Jurkat and Central Memory CD4 $^{+}$ datasets (Bonferroni corrected Fisher’s exact test, all $p < 0.05$) and approached significance in the Active CD4 $^{+}$ dataset ($p = 0.053$) (Figure 2.6). The Bcl-2 transduced CD4 $^{+}$ data did not contain any integration sites in alphoid repeats, probably due to 1) the relatively low number of integration sites in the dataset and 2) to the requirement for cleavage at two PstI restriction sites, which are not found in the consensus sequence of alphoid repeats³⁴⁷. Of the 1340 repeat types in the RepeatMasker database³⁴⁷, only alphoid repeats achieved a significant association with proviral expression in more than two datasets.

2.4.6 Acetylation

Histone marks or chromatin remodeling, especially involving the key “Nuc-1” histone near the transcription start site in the viral LTR, appear to affect viral expression^{306,348,349}. Based on this effect, histone deacetylase inhibitors have been developed

as potential HIV treatments and show some promise in disrupting latency³¹⁸. In these genome-wide datasets, we do not have information on the state of individual LTR nucleosomes. However, repressive chromatin does seem to spread to nearby locations if not blocked by insulators^{302,303} and the state of neighboring chromatin could affect proviral transcription independently of provirus-associated histones.

We found that the number of ChIP-seq reads near an integration site from several histone acetylation marks (Figure 2.1) were associated with efficient expression in the Active CD4⁺, Resting CD4⁺ and Central Memory CD4⁺ samples. H4K12ac had the strongest association (Bonferroni corrected Fisher's method combination of Spearman's ρ , $p < 10^{-25}$) with silence/latency (Figure 2.7A).

Although the appearance of several significantly associated acetylation marks might suggest acetylation exerts a considerable effect on the expression of a provirus, there are strong correlations among these marks, so their effects may not be independent. To account for the correlations between these variables, we performed a principal component analysis (PCA) to convert the correlated acetylation marks into a series of uncorrelated principal components that capture much of the variance within a few components. Here, the first principal component explained 59% of the variance and the first ten components 84%. Several of these principal components again displayed significant associations with latency in the Active CD4⁺, Resting CD4⁺ and Central Memory CD4⁺ samples but no significant correlations in the Bcl-2 transduced CD4⁺ or Jurkat samples (Figure 2.7B). A logistic regression of expression status on the first ten principal components and sample did not reduce misclassification error from a base model including only sample in 480-fold cross validation (base model misclassification error: 36.4%, PCA model: 36.5%). This suggests that acetylation of neighboring chromatin does not exert strong effects on latency in all samples.

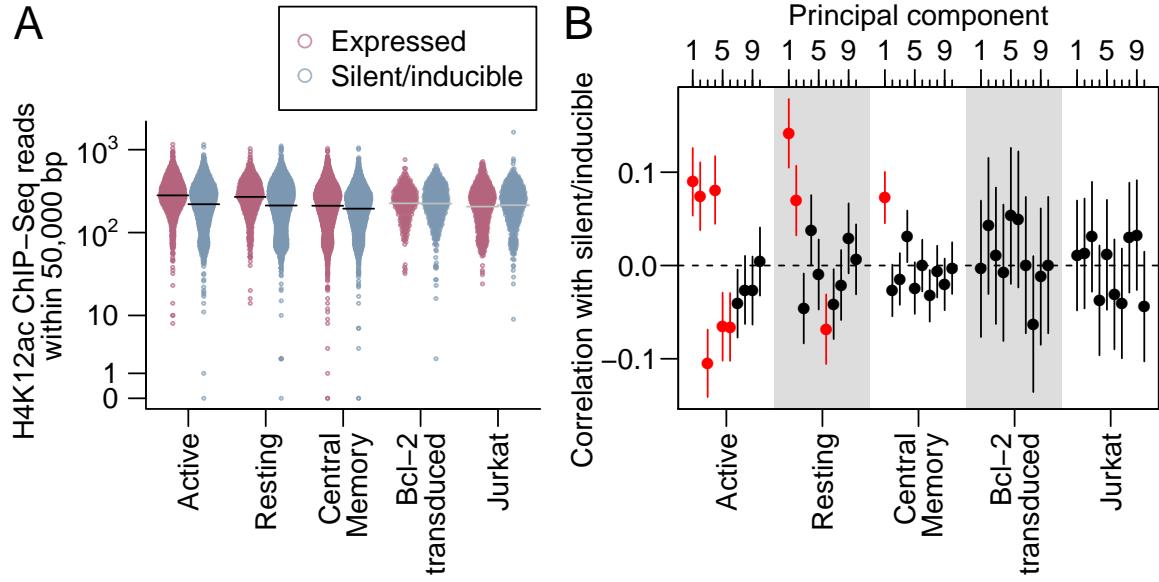


Figure 2.7: (A) The number of ChIP-seq reads for H4K12ac, the histone mark with the lowest Fisher's method p -value for correlation with latency, within 50,000 bases across the five samples. Integration sites (points) are spread in proportion to kernel density estimates. Horizontal lines indicate sample means where there was a significant difference (black) in means between silent/inducible and expressed provirus or no significant difference (grey). (B) The correlation (points) and its 95% confidence interval (vertical lines) between principal components of acetylation and silent/inducible status for each of the five samples. Red indicates correlations with a Bonferroni-corrected p -value < 0.05 .

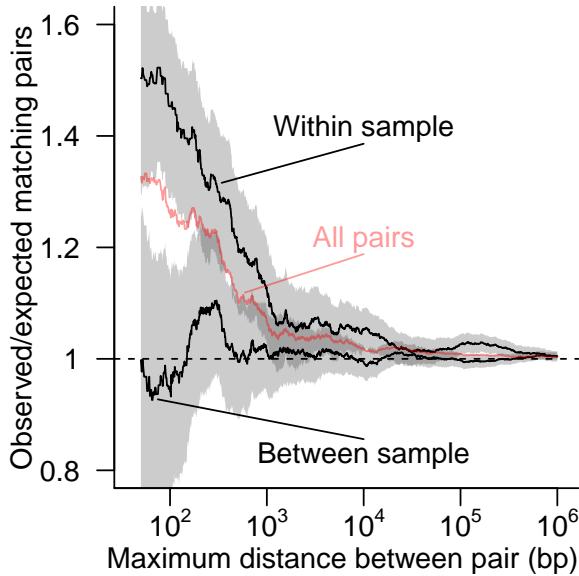


Figure 2.8: The ratio of the number of pairs of proviruses with matching expression status to the number of matches expected by random pairings given the frequency of silent/inducible proviruses. All possible pairs of proviruses integrated within a given distance of each other on the same chromosome (red line) were separated into two sets; one with both proviruses from within the same cell culture model and one with proviruses paired between two different cell culture models (black lines). The shaded region shows the 95% Clopper-Pearson binomial confidence interval for within and between sample pairings. The dashed horizontal line shows the ratio of 1 expected if there is no association between the expression status of neighboring proviruses.

2.4.7 Clustering

We reasoned that if there was a strong relationship between latency and chromosomal position, then integration sites that are near one another on the same chromosome should share the same expression status more often than expected by chance. To test this, we compared how often pairs of proviruses shared the same expression status in relation to the distance between the two sites (Figure 2.8). Pairs of sites with little distance between integration locations did share the same expression status more often than expected by chance (e.g. neighbors closer than 100bp, Fisher exact test $p = 0.0002$). Breaking out the data to separate between sample and within sample pairings showed that this matching was limited to neighbors within the same experimental model (Figure 2.8), emphasizing that chromosomal environment does appear to influence latency, but the factors involved differ among experimental models of latency.

2.5 Conclusions

Here we compared the latency status of HIV-1 proviruses in five model systems with the genomic features surrounding their integration sites. Surprisingly, no relationships between genomic features near the integration location and latency achieved significance in all models. Proviruses from the same cellular model integrated in nearby positions did share the same latency status much more often than predicted by chance, indicating the existence of local features influencing latency, but these were not consistent among models. This suggests that whatever features are affecting latency are highly local and model-specific, and that we may not have access to all relevant chromosomal features e.g. ^{350–353}.

In addition to differences in experimental conditions, methodological issues have the potential to obscure patterns. Examples include multiply infected cells, inactivated viruses and inaccurate assessment of HIV gene activity—each of these are discussed below.

A latent provirus integrated into the same cell as an expressed provirus will be erroneously sorted as expressed, potentially confounding analysis. A low multiplicity of infection (MOI) will help to avoid this problem, but there is still the potential for a significant proportion of the cells studied to contain multiple integrations. This problem arises because although cells with multiple integrations form a small proportion of total cells, most of the total are cells lacking an integrated provirus and thus are excluded by experimental design. For example, assuming integrations are Poisson distributed with an MOI of 0.1 (1 integration per 10 cells), 90.5% of cells will not contain a provirus, 9% of cells will contain one proviral integration and 0.5% of cells will contain multiple integrations. The cells without an integration are not amplified by HIV-targeted PCR leaving only 9.5% of the total cells. Of these cells actually under study, 4.9% will contain multiple integrations. Thus the signal from expressed proviruses may be muted by the presence of latent proviruses in the expressed population.

The replication cycle of HIV is error prone, and a significant proportion of virions contain mutated genomes⁸³. In studies that do not check for inducibility, mutant proviruses integrated in regions of the genome otherwise favorable to proviral expression can be sorted into the latent pool due to mutational inactivation. This problem of inactivated provirus is worse when latent provirus are rare and exacerbated further when looking at latency in the cells of HIV patients due to selective enrichment of inactivated proviruses incapable of spreading infection¹³¹. Here, the effects of mutation are minimized in the datasets that required inducible viral expression (Jurkat, Bcl-2 transduced CD4⁺, Central Memory CD4⁺) but may be a confounder in the two datasets that were sorted based on lack of viral expression only (Active CD4⁺, Resting CD4⁺).

Inaccurate staining or leaky markers may also result in misclassification of proviruses. False positives and false negatives will result in incorrectly sorted latent and expressed integrations. For example, if 5% of cells not containing Gag are labeled as Gag+ and there are an equal amount of latent and expressed integration sites, then 4.8% of integrations labeled expressed will actually be latent. If a category is rare, false staining has even greater potential to cause error. For example, if only 5% of sites are latent and a Gag stain has a false negative rate of 5%, then we would expect 48.7% of sites classified as latent to actually be mislabeled expressed integrations.

Attempts to induce latent proviruses in patients have so far focused on using histone deacetylase inhibitors, raising interest in associations with histone acetylation in these data. An important caveat in results from these genome-wide data is that histone modification near the integrated provirus may not be representative of modification within the provirus at the key “Nuc-1” nucleosome of the transcription start site³⁴⁹, though local correlations in chromatin states are well established from studies of position effect variegation^{302,303}. We found that some histone acetylation marks were significantly associated with viral expression in some but not all samples (Figures 2.1, 2.7). This lack of association may be due to a lack of power in these studies, but the

confidence intervals suggest that any correlations between acetylations and latency are unlikely to be strong. These weak correlations raise the possibility that there are populations of latent proviruses that are not associated with acetylation and may not be inducible by histone deacetylase inhibitors.

This study highlights that the choice of model system can have a large effect on measurements of latency. Further studies are needed to determine which *in vitro* models best reflect latency *in vivo*. Different cell models may report genuinely different mechanisms of latency. While we did see some relationship between histone acetylation and latency, paralleling a recent clinical trial of SAHA³¹⁸, associations with histone acetylation did not explain a large fraction of the difference between latent and expressed proviruses in any of the five models. One possible explanation is that there may be multiple mechanisms that maintain proviruses in a latent state. To be successful, shock-and-kill treatments must induce and destroy all latent proviruses to eliminate HIV from an infected individual, raising the question of whether multiple simultaneous inducing treatments will be necessary.

2.6 Availability of supporting data

Sequence reads from the Central Memory CD4⁺ sample reported here, the Resting CD4⁺ and Active CD4⁺ data reported by Pace et al.³¹², the Bcl-2 transduced CD4⁺ data reported by Shan et al.³¹¹ and reprocessed data originally reported by Lewinski et al.³¹⁰ are available at the Sequence Read Archive under accession number SRP028573.

2.7 Author's contributions

SS-M led the computational analysis, with assistance from CCB and NM. MKL, DL and JG analyzed integration sites using IonTorrent sequencing. MF, AB and VP prepared DNA from latent and activated T cells using the Central Memory CD4⁺ model. LS, RFS, MJP, LMA and UO'D contributed data and suggestions. SS-M, KEO and FDB planned the overall study, and SS-M and FDB wrote the paper. All authors read and

approved the final manuscript.

2.8 Acknowledgements

We would like to thank Werner Witke for assistance with IonTorrent sequencing. This work was supported in part by NIH grants R01 AI 052845-11 to FDB, R21AI 096993 and K02AI078766 to UO'D, 5T32HG000046 to SS-M, AI087508 to VP and R01AI038201 to JG, the Penn Genome Frontiers Institute, the University of Pennsylvania Center for AIDS Research (CFAR) P30 AI 045008 and the University of California, San Diego, CFAR P30 AI036214.

CHAPTER 3 : Dynamic regulation of HIV-1 mRNA populations analyzed by single-molecule enrichment and long-read sequencing

3.1 Abstract

Alternative RNA splicing greatly expands the repertoire of proteins encoded by genomes. Next-generation sequencing (NGS) is attractive for studying alternative splicing because of the efficiency and low cost per base, but short reads typical of NGS only report mRNA fragments containing one or few splice junctions. Here, we used single-molecule amplification and long-read sequencing to study the HIV-1 provirus, which is only 9700 bp in length, but encodes nine major proteins via alternative splicing. Our data showed that the clinical isolate HIV_{89.6} produces at least 109 different spliced RNAs, including a previously unappreciated ~1 kb class of messages, two of which encode new proteins. HIV-1 message populations differed between cell types, longitudinally during infection, and among T cells from different human donors. These findings open a new window on a little studied aspect of HIV-1 replication, suggest therapeutic opportunities and provide advanced tools for the study of alternative splicing.

3.2 Introduction

Alternative splicing greatly expands the information content of genomes by producing multiple mRNAs from individual transcription units. Approximately 95% of human genes with multiple exons encode RNA transcripts that are alternatively spliced, and mutations that affect alternative splicing are associated with diseases ranging from cystic fibrosis to chronic lymphoproliferative leukemia^{247,354–357}. Work to decipher an RNA ‘splicing code’ has revealed that multiple interactions between trans-acting factors and RNA elements determine splicing patterns, though regulation is little understood for most genes²⁴⁸.

The integrated HIV-1 provirus is ~9700 bp in length and has a single transcription start site, but according to the published literature yields at least 47 different mRNAs encoding 9 proteins or polyproteins, making HIV an attractive model for studies of alternative splicing²⁴⁰. HIV mRNAs fall into three classes: the unspliced RNA genome, which encodes Gag/Gag-Pol; partially spliced transcripts, ~4 kb in length, encoding Vif, Vpr, a one-exon version of Tat, and Env/Vpu; and completely spliced mRNAs of roughly 2 kb encoding Tat, Rev and Nef (Figure 3.1A). Additional rare ‘cryptic’ splice donors (5' splice sites) and acceptors (3' splice sites) contribute even more mRNAs^{242,358–362}. A complex array of positive and negative cis-acting elements surrounding each splice site regulates the relative abundance of the HIV-1 mRNAs, and disrupting the balance of message ratios impairs viral replication in several models^{219,222,223,227,363–366}. Studies have suggested strain-specific splicing patterns may exist^{240,367,368}. However, detailed studies of complete message populations have not been reported for clinical isolates of HIV-1.

Several groups have demonstrated tissue- and differentiation-specific splicing of cellular genes^{247,369,370}. Importantly for HIV, these include changes during T-cell activation^{371,372}, raising the question of how cell-specific splicing affects HIV replication. While most studies of HIV-1 splicing have been conducted in cell lines using lab-adapted viral strains, limited works in PBMCs from infected patients, monocytes and macrophages have suggested that differences may indeed exist in relevant cell types^{358,367,373,374}. Moreover, human splicing patterns differ between individuals, but such polymorphisms have not been investigated in the context of HIV infection^{245,246}.

Here, we use deep sequencing to comprehensively characterize the transcriptome of an early passage clinical isolate, HIV_{89.6}³⁷⁵, in primary CD4⁺ T cells from seven human donors and in the human osteosarcoma (HOS) cell line. Many deep sequencing techniques provide short reads, which rarely query more than a single exon-exon junction. To distinguish the full structure of HIV-1 mRNAs, which can contain several splice

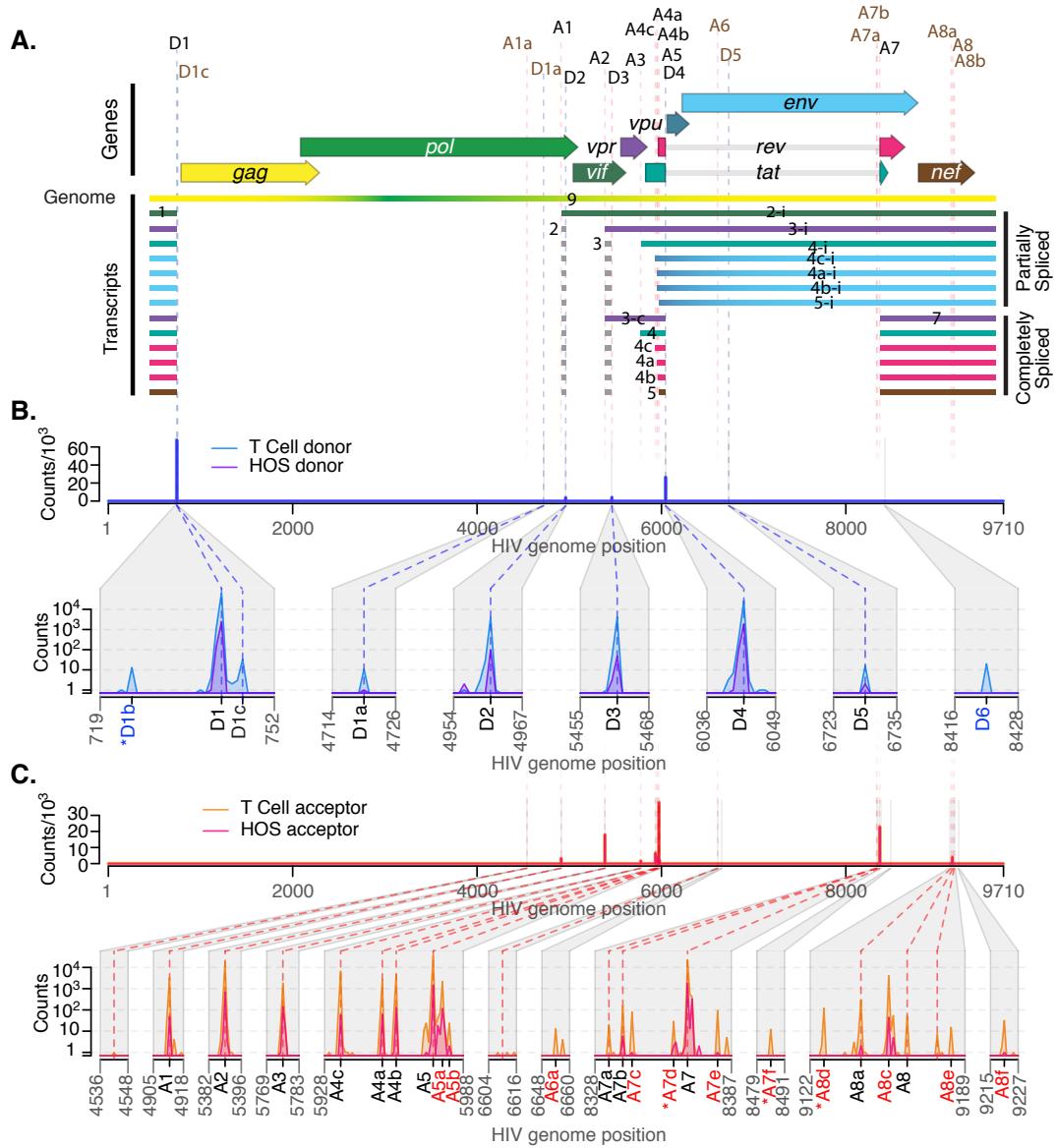


Figure 3.1: Mapping the splice donors and acceptors of HIV_{89.6}. PacBio sequence reads of HIV_{89.6} cDNA from infected HOS-CD4-CCR5 (HOS) and CD4⁺ T cells were aligned to the HIV_{89.6} genome shown in (A). Exons of the conserved HIV-1 transcripts are colored according to the encoded gene. Conserved (black) and published cryptic (brown) splice donors ('D') and acceptors ('A') are shown. Gaps in HIV-1 sequence alignments with at least one end located at a published or verified splice donor or acceptor were defined as introns. For each base of the HIV_{89.6} genome, the number of sequence reads in which that base occurred at the 5'-end (B) or 3'-end (C) of an intron is plotted for each cell type. Putative splice donors and acceptors were defined as loci that were found in at least 10 reads at the 5'- and 3'-ends of introns in sequence alignments from T-cell infections. Regions containing splice sites are enlarged for clarity. Asterisks indicate putative splice sites that are adjacent to dinucleotides other than the consensus GT and AG.

junctions, we used Pacific Biosciences (PacBio) sequencing technology, which yields read lengths up to 10 kb²⁵⁷. We used RainDance Technologies single-molecule PCR enrichment to preserve ratios of RNAs during preparation of sequencing templates. We identified previously published and novel HIV-1 transcripts and determined that HIV_{89.6} encodes a minimum of 109 different splice forms. These included a new size class of transcripts, some of which contain novel open reading frames (ORFs) that encode new proteins. We also found significant variation between cell types, over time during infection of HOS cells and among individuals. These data reveal unanticipated complexity and dynamics in HIV-1 message populations, begin to clarify a little studied dimension of HIV-1 replication and suggest possible targets for therapeutic interventions.

3.3 Materials and methods

3.3.1 Cell culture and viral infections

HIV_{89.6} was generated by transfection and subsequent expansion in SupT1 cells. Primary T cells were isolated by the University of Pennsylvania Center for AIDS Research Immunology core and confirmed to be homozygous for the wild-type CCR5 allele as shown in Supplementary Table S1 and described in Supplementary Methods. HOS-CD4-CCR5 cells^{376,377} were obtained through the AIDS Research and Reference Reagent Program, Division of AIDS, NIAID, NIH from Dr Nathaniel Landau. Single round infections in T cells and HOS-CD4-CCR5 cells were performed using standard methods (see Supplementary Methods).

3.3.2 RNA and reverse transcription

Total cellular RNA was purified using the Illustra RNA kit (GE Life Sciences, Fairfield, CT, USA) from 5×10^6 cells per infection. Viral cDNA was made using a reverse transcription primer complementary to a sequence in U3 (RTprime, Supplementary Table S2). We used Superscript III reverse transcriptase (Invitrogen) in the presence of

RNaseOUT (Invitrogen) to conduct first-strand cDNA synthesis from equal amounts of total cellular RNA from each HOS-CD4-CCR5 time point (15.2 µg) and from each T-cell infection (3 µg) according to the manufacturer's instructions for gene-specific priming of long cDNAs, and then treated with RNaseH (Invitrogen). We checked for full reverse transcription of the longest (unspliced) viral cDNAs by PCR using primers that bind in the first major intron of HIV_{89.6} (keo003, keo004, Supplementary Table S2, data not shown).

3.3.3 Bulk RT-PCR and cloning

Transcripts were amplified from cellular RNA using the Onestep RT-PCR kit (Qiagen) with primer pairs keo056/keo057 and keo058/keo059 (Supplementary Table S2) with the following amplification: 5 cycles of 30 s at 94°C, 12 s at 56°C, 40 s at 72°C; then 30 cycles of 30 s at 94°C, 14 s at 56°C, 40 s at 72°C; and finally 10 min at 72°C. For verification of dynamic changes, primers F1.2 and R1.2 were used with 35 cycles of 30 s at 94°C, 30 s at 56°C and 45 s at 72°C followed by 10 min at 72°C. Products were resolved on agarose gels (Nusieve 3:1, Lonza for verification of dynamic changes, Invitrogen for cloning) stained with ethidium-bromide (Sigma) for visualization, or SYBR Safe DNA gel stain (Invitrogen) for cloning (keo056/keo057 amplified material). DNA was purified using Qiaquick gel extraction kit (Qiagen) and cloned using the TOPO TA cloning kit (Invitrogen). Plasmid DNA was prepared using Qiaprep Spin Miniprep kit (Qiagen). Inserts were identified and verified using Sanger sequencing. The cDNAs for *tat*^{8c}, *tat* (1 and 2 exon), *ref*, *rev* and *nef*, and the transcript with exon structure 1-5-8c were cloned into the expression vector pIRES2-AcGFP1 (Clonetech) as described in Supplementary Methods.

3.3.4 Assays of protein activity and HIV replication

Activity and HIV replication assays were performed as described in Supplementary Methods. Tat activity expressed from each cDNA was measured in TZM-bl cells³⁷⁸ (gift

of Dr Robert W. Doms). Rev activity was assayed in HEK-293T cells co-transfected with pCMVGagPol-RRE-R, a reporter plasmid from which Gag and Pol are expressed in a Rev-dependent manner (gift of David Rekosh)³⁷⁹. Intracellular and released supernatant p24 was measured from cells transfected with expression constructs and infected with HIV_{89.6}.

3.3.5 Western blotting

HEK-293T cells were transfected with expression constructs and treated with MG132 (EMD Chemicals) to inhibit the proteasome or DMSO (Supplementary Methods). Proteins were detected by immunoblotting using a mouse antibody that recognizes the carboxy terminus of HIV-1 Nef diluted 1:1000 in 5% milk (gift of Dr James Hoxie)³⁸⁰. Horseradish peroxidase (HRP)-conjugated secondary rabbit-anti-mouse antibody (p0260, DAKO) was used for detection with SuperSignal West Pico Chemiluminescent Substrate (Thermo Scientific). Beta-tubulin was used as a loading control, detected by the HRP-conjugated antibody (ab21058, Abcam).

3.3.6 Single-molecule amplification

Amplification was performed by RainDance Technologies using a protocol similar to that previously reported (detailed description in Supplementary Methods)²⁵⁴. Amplification was carried out in droplets to suppress competition between amplicons. PCR droplets were generated on the RDT 1000 (RainDance Technologies) using the manufacturer's recommended protocol. The custom primer libraries for this study contained 18 (HOS-CD4-CCR5 cells) or 20 (primary T cells) PCR primer pairs designed to amplify different HIV RNA isoforms (Supplementary Table S2).

3.3.7 Single-molecule sequencing

DNA amplification products from the RainDance PCR droplets were converted to SM-RTbell templates using the PacBio RS DNA Template Preparation Kit. Sequencing

was performed by Pacific Biosciences using the PacBio SMRT sequencing technology as described²⁵⁷. Sequence information was acquired during real time as the immobilized DNA polymerase translocated along the template molecule. Prior to sequence acquisition, hairpin adapters were ligated to each DNA template end so that DNA polymerase could traverse DNA molecules multiple times during rolling circle replication (SMRTbell template sequencing³⁸¹), allowing error control by calculating the consensus ('circular consensus sequence' or CCS). For raw reads, the average length was 2860 nt, and 10% were > 5000 nt. After condensing into consensus reads, the mean read length was 249.5 nt, due to the use of a shorter Pacific Biosciences sequencing protocol to accommodate the small size of many amplicons. Consensus reads of 1% were > 1100 nt. Sequencing data were collected in 45-min movies.

3.3.8 Data analysis

Raw reads were processed to produce CCSs. Raw reads were also retained to help in primer identification and to avoid biasing against long reads. Reads were aligned against the human genome using Blat³³¹. Misprimed reads matching the RT primer, reads with a CCS length shorter than 40 nt or raw length shorter than 100 nt and reads matching the human genome were discarded. Filtered reads were aligned against the HIV_{89.6} reference genome. Potential novel donors and acceptors were found by filtering putative splice junctions in the Blat hits for a perfect sequence match 20 bases up- and downstream of the junction, ignoring homopolymer errors, and requiring that one end of the junction be a known splice site. Local maximums within a 5-nt span with > 9 such junctions were called as novel splice sites.

Filter-passed reads were aligned against all expected fragments based on primers and known and novel junctions. Primers were identified in CCS reads by an edit distance ≤ 1 from the primer in the start or end of the read, in raw reads by an edit distance ≤ 5 from a concatenation of the primer, hairpin adapter and the reverse complement of the primer, and in both types of reads by a Blat hit spanning an entire expected fragment.

Gaps in Blat hits were ignored if \leq 10 bases long or in regions of likely poor read quality \leq 20 bases long where an inferred insertion of unmatched bases in the read occurred at the same location as skipped bases in the reference. Any Blat hits with a gap $>$ 10 nt remaining in the query read were discarded. If HIV sequence was repeated in a given read (likely due to PacBio circular sequencing), the alignments were collapsed into the union of the coverage. Gaps in the HIV sequence found in uninterrupted query sequence were called as tentative introns. Splice junctions were assigned to conserved or previously identified (published or in this work) splice sites and reads appearing to contain donors or acceptors further than 5 nt away from these sites were discarded. Reads with Blat hits outside the expected primer range were discarded from that primer grouping. The assigned primer pair, observed junctions and exonic sequence were used to assign each read to a given spliceform (specific transcript structure) or set of possible spliceforms. Partial sequences that did not extend through both primers were assigned to specific transcripts if the read contained enough information to rule out all other spliceforms or if all other possible spliceforms contained rare ($< 1\%$ usage) donors or acceptors (Supplementary Table S3). Otherwise, the read was called indeterminate.

To calculate the ratios of transcripts within the partially spliced class, we counted the number of reads for each assigned spliceform amplified by primer pair 1.3 and divided by the total number of assigned partially spliced reads amplified with these primers (Supplementary Figure S1 and Supplementary Table S2). Assigned sequences amplified with primer pairs 1.4 and 4.1 (full-length cDNAs, T cells only) were used to calculate ratios of transcripts within each of the two completely splice classes (~ 2 and ~ 1 kb). To compare ratios of ~ 2 kb transcripts calculated within reads from primer pairs 1.4 and 4.1, we normalized ratios from pair 4.1 to the *nef* 2 transcript (containing exons 1, 5 and 7). Due to size biases inherent in the approach, we did not compare across size classes, and unspliced transcripts were not included in ratio analysis. For all ratio analysis, transcripts including cryptic or novel junctions were counted only if they appeared in at least five reads, otherwise they were excluded from the analysis and from the count of

total assigned reads.

To estimate the minimum total number of transcripts present, partial sequence reads were included. Each exon-exon junction occurring in at least five reads and not previously assigned to a particular transcript (Figure 3.2) was counted as evidence of an additional transcript (47 additional junctions were detected, see Supplementary Table S4). If two such junctions could conceivably occur in a single mRNA, we counted only one unless we could verify from sequence reads that they were amplified from separate cDNAs, resulting in 31 additional transcripts. The minimum transcript number calculated by a greedy algorithm treating introns as events in a scheduling problem agreed with the above calculation.

Several groups have demonstrated tissue- and differentiation-specific splicing of cellular genes^{247,369,370}. Importantly for HIV, these include changes during T-cell activation^{371,372}, raising the question of how cell-specific splicing affects HIV replication. While most studies of HIV-1 splicing have been conducted in cell lines using lab-adapted viral strains, limited works in PBMCs from infected patients, monocytes and macrophages have suggested that differences may indeed exist in relevant cell types^{358,367,373,374}. Moreover, human splicing patterns differ between individuals, but such polymorphisms have not been investigated in the context of HIV infection^{245,246}.

For studies of transcript dynamics, reads from primer pairs 1.2, 1.3 and 1.4 containing junctions between D1 or any donor and each of five mutually exclusive acceptors, A3, A4c, A4a, A4b, A5 and A5a, were collected and their ratios calculated.

3.3.9 Statistical analysis

Statistical modeling was performed using generalized linear modeling as described in Supplementary Report S2. All analyses were performed in R 2.14.0 (R Development Core)³²³.

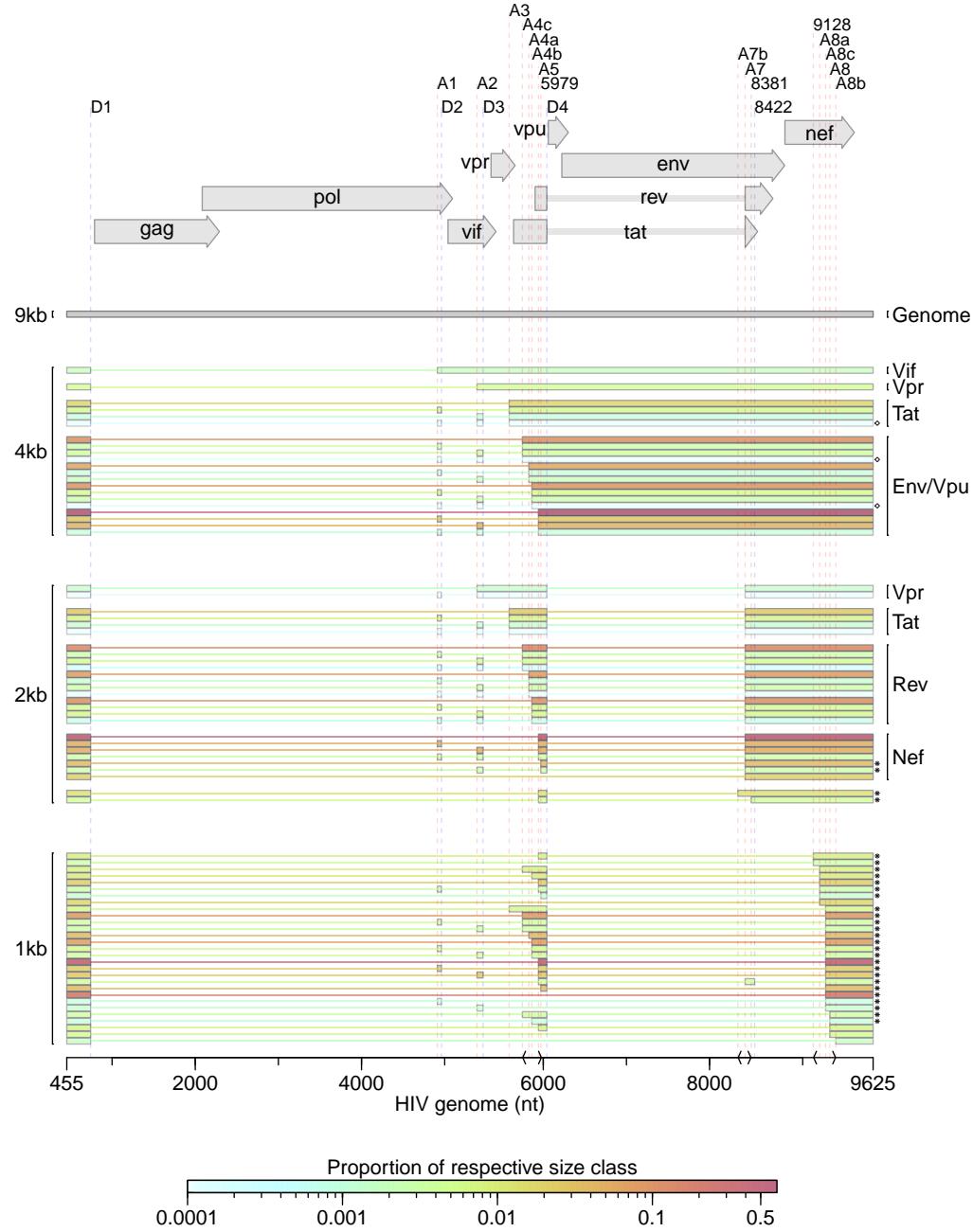


Figure 3.2: HIV_{89.6} transcripts in T cells for which the full message structure was determined are shown arranged by size class. Thick bars correspond to exons and thin lines to excised introns. For the well-conserved transcripts, encoded proteins are indicated. The relative abundance of each transcript within its size class is indicated by color. Asterisks denote transcripts that have not been reported previously to our knowledge. Of the 47 conserved HIV-1 transcripts, three were detected in fewer than five reads (indicated with ◊) and two messages were not detected and are not shown (one encoding Vpr and one encoding Env/Vpu). Depicted non-conserved transcripts (using novel or cryptic splice sites) were each detected in at least five independent sequence reads across samples from at least two different human T-cell donors.

3.3.10 Data access

Sequence data is available in the SRA database with the following accession numbers: SRP014319.

3.4 Results

3.4.1 Sequencing HIV-1 transcripts produced in primary T cells and HOS cells

In order to characterize HIV-1 transcript populations, we prepared viral cDNA from primary CD4⁺ T cells of seven different healthy human donors infected in vitro with HIV_{89.6}, an early passage dual-tropic clade-B clinical isolate (Supplementary Figure S1, human donor data in Supplementary Table S1)³⁷⁵. We also studied HIV messages produced in infected HOS cells engineered to express CD4 and CCR5 (HOS-CD4-CCR5) because these cells support efficient HIV replication and engineered variants are widely used in HIV research. HOS cells were harvested at 18, 24 and 48 hours post infection (hpi) to investigate longitudinal changes during infection, and for comparison to 48 h infected T cells.

To preserve the relative proportions of template molecules while amplifying the cDNA, we used RainDance Technologies' single-molecule micro-droplet based PCR²⁵⁴. Droplet libraries containing multiple overlapping primer pairs were designed to query all message forms and allow later calculation of relative abundance (Supplementary Table S2 and Supplementary Figure S1). Each primer was unique so that sequences could be assigned to a specific primer pair, which helped reconstruct the origin of sequence reads and deduce message structures. Amplified DNA products were sequenced using Single Molecule Real-Time (SMRT) technology from Pacific Biosciences^{257,381}. We obtained 847 492 filtered reads of amplified HIV-1 transcripts in primary CD4⁺ T cells and 89 350 in HOS cells. The longest sequenced continuous stretch of HIV-1 cDNA was 2629 bp.

3.4.2 Splice donors and acceptors

We aligned PacBio reads containing HIV sequences to the HIV_{89.6} genome and identified candidate introns as recurring gaps in our sequences. Using this approach, we observed splicing at each of the widely conserved major splice donors and acceptors and several published cryptic sites (Figure 3.1A, hereafter referred to by their identifications shown in this figure, ‘D’ for donors, ‘A’ for acceptors).

In addition, we identified 13 putative novel splice sites: 2 donors and 11 acceptors (Figure 3.1 and Supplementary Table S3). In order to be selected as a bona fide splice site and remove artifacts possibly created by recombination during sample preparation, we required that the new acceptor or donor was observed spliced to previously reported splice donors or acceptors in > 10 sequence reads in CD4⁺ T cells. The most frequently used novel splice site was an acceptor that we have termed A8c because it lies near A8, A8a and A8b (discussed in detail below). Additional novel sites are further discussed in Supplementary Report S1.

Most of the new splice sites adhered to consensus sequences for the standard spliceosome (Supplementary Table S3). However, there appeared to be one splice donor upstream of D1 with a cytidine in place of the usual uracil 2 nt downstream of the splice site. Similar ‘GC donors’ appear in 1% of known splice junctions in humans³⁸². Of the novel splice acceptors, three were preceded by dinucleotides other than the consensus AG. Alternative dinucleotides are used infrequently as splice acceptors^{383–386}; however, it is possible that our deep sequencing method allowed us to observe rare events.

3.4.3 Structures of spliced HIV_{89.6} RNAs

To quantify the populations of HIV-1 transcripts, we aligned all reads to the collection of 47 well-established spliced HIV-1 transcripts and detected 45 of them (Figure 3.2). We additionally aligned reads to the HIV_{89.6} genome allowing all possible combinations of splice junctions—canonical, cryptic or novel—determined from the sequencing data

(Figure 3.1), yielding an additional 32 complete transcripts, 19 of which were novel. The data also provide evidence for more novel splice junctions but in incomplete sequences, implying the existence of additional new transcripts (Supplementary Table S4 and Supplementary Report S1). The full data set taken together provides evidence for least 109 different HIV_{89.6} transcripts in primary T cells.

Amplification primers that isolated the two main classes of spliced messages allowed us to determine the ratios of mRNAs in each (Figure 3.2 and Supplementary Table S5). Within the partially spliced class of transcripts, *env/vpu*, *tat* (1-exon), *vpr* and *vif* messages existed in an average ratio of 96:4:< 1:< 1 in CD4⁺ T cells. The ratio of *nef:rev:tat:vpr* within the ~2 kb transcript class was 64:33:3:< 1. Consistent with previous reports, the most abundant transcript in each class contained the splice junction from D1 to A5 (D1^A5)—an *env/vpu* transcript contributing 64% of the partially spliced class, and a completely spliced *nef* transcript contributing 47% of ~2 kb messages (Figure 3.2)^{240,387}. The relatively low abundance of transcripts encoding Tat suggests that Tat sufficiently stimulates HIV transcription elongation at low concentrations, or that the *tat* transcripts must be efficiently translated. Due to biases inherent in the reverse transcription step, we could only compare transcripts within each size class, and we note that our methods have not been validated for empirical quantification. However, the ratios were roughly confirmed using overlapping sequence reads obtained with alternate primer pairs and by end point RT-PCR analysis of HIV-1 RNAs (data not shown).

Exons 2 and 3 are non-coding exons whose inclusion in transcripts other than *vif* and *vpr* has no known function. We found that they were included in other messages infrequently, each in ~7–8% of transcripts in the ~2 kb completely spliced class of transcripts and 5% of partially spliced transcripts accumulating in T cells. This is consistent with previous measurements in the partially spliced class but much lower than has been estimated for completely spliced transcripts in HeLa cells, suggesting

cell-type-specific splicing patterns may influence inclusion of these exons²⁴⁰.

3.4.4 A novel ~1 kb class of completely spliced transcripts

Primers placed near the 5'- and 3'-ends of the HIV_{89.6} genome amplified a second class of completely spliced transcripts ~1 kb in length. In place of A7, these transcripts use a set of little studied splice acceptors located ~800 bp downstream within the 3'-TR. Two groups have previously observed splicing from D1 to acceptors A8, A8a and A8b in this region, yielding messages of this size class in patient samples; however, none of these could be translated to a protein of significant length^{358,362}. We determined the complete structure of 29 members of the 1-kb class (Figure 3.2 and Supplementary Table S5). The most abundant messages observed in this class use the novel acceptor A8c to define their terminal exon. For HIV89.6, acceptor A8c was used nearly as frequently as A7, which gives us the 2-kb class of transcripts (Supplementary Table S3), and this was supported by end point RT-PCR analysis (data not shown).

Acceptor A8c is not well conserved in HIV-1/SIVcpz (14%), although it is conserved in clade G viruses (> 95%) and most HIV-2/SIVsmm genomes (86%)³⁸⁸. This is due to the poor conservation of an adenosine at the wobble base position of the 123rd codon (proline) of the Nef reading frame, which creates the AG dinucleotide generally required at splice acceptors. Since any base at this position would code for proline, there does not seem to be strong selection for a splice acceptor here. However, A8c is displaced from nearby well-conserved (> 90%) cryptic acceptors A8a and A8b by multiples of 3 bp (12 and 21 bp, respectively), so splicing to any of these three acceptors would create similar ORFs. All HIVs and SIVs maintain at least one of these three acceptors, suggesting possible function³⁸⁸. We confirmed that the 1 kb transcripts using A8a, A8b and A8c were present in infected HOS and T cells by end point RT-PCR using additional primer pairs and by Sanger sequencing of cloned transcripts (Figure 3.3A and B; data not shown).

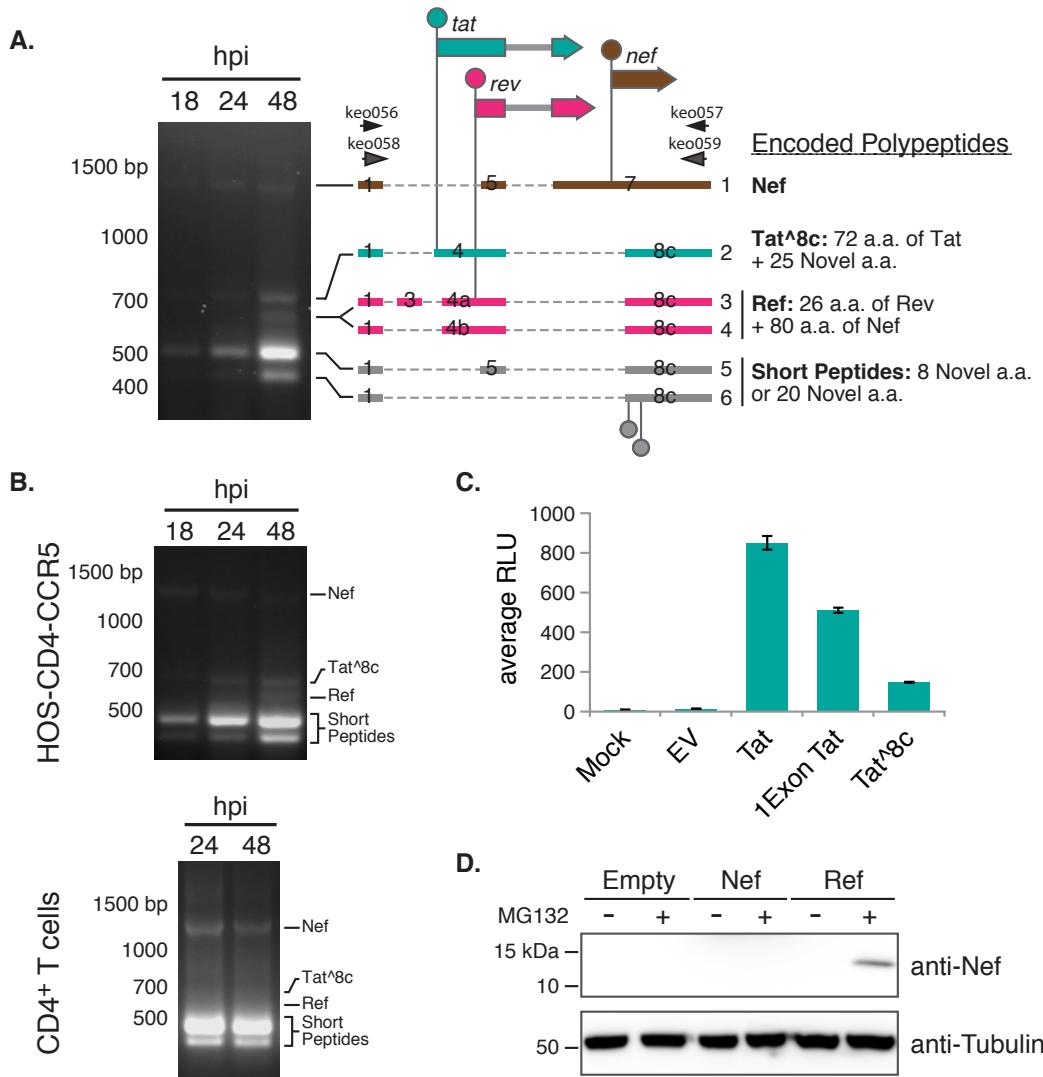


Figure 3.3: HIV_{89.6} transcripts were amplified by RT-PCR using RNA from infected HOS-CD4-CCR5 cells with primers keo056 and keo057. Major bands detected after gel electrophoresis were cloned from the 48 hpi sample and message structures determined by Sanger sequencing. Thick bars represent exons and dashed lines excised introns. Genes are shown above (not to scale) with start codons indicated by circles. Messages 1, 2, 4 and 5 were cloned into expression plasmids for activity assays. (B) Confirmation of presence of the ~1 kb message RNAs in HOS-CD4-CCR5 and primary CD4⁺ T cells (human donor 1, harvested 24 and 48 hpi). An independent primer pair (keo058 and keo059) was used to amplify transcripts by RT-PCR. (C) Tat activity was measured in Tzm-bl cells as Tat-dependent luciferase production after transient transfection with expression plasmids. (D) Western blot showing expression of protein of the predicted size for Ref (12.5 kb) in cells transfected with the Ref expression construct and treated with proteasome inhibitor MG132, detected by an antibody recognizing the carboxy-terminus of Nef. Expression plasmid encoding Nef was included to control for possible expression of partial Nef peptides or breakdown products from the Nef ORF.

The 1-kb transcript containing exons 1, 4 and 8c (1-4-8c, where exon 8c begins at A8c and extends to the poly-adenylation site) encodes the first exon of Tat followed by 25 novel amino acids (termed Tat^{8c}). Tat^{8c} showed activity when overexpressed in cells containing a Tat reporter construct (Figure 3.3C, nucleotide and amino acid sequences in Supplementary Table S6). Transcripts with exon structures 1-4a/b/c-8c encode a novel fusion of the amino-terminal 26 amino acids of Rev and the carboxy-terminal 80 amino acids of Nef, hereafter referred to as Ref. We did not detect Rev activity on overexpression of the *ref* transcript, and Ref did not appear to interfere with the normal function of Rev or with HIV replication (Supplementary Figure S2). Ref was detectable by western blot using antibodies targeting the C terminus of Nef after inhibition of the proteasome, suggesting that the fusion is expressed but not stable (Figure 3.3D). Thus, Ref has the potential to encode a new epitope potentially relevant in immune detection of HIV. The transcripts with exon structures 1-5-8c and 1-8c encode at most a short peptide, and so are candidates for acting as regulatory RNAs.

3.4.5 Temporal dynamics of transcript populations

To assess longitudinal variation, we investigated HIV_{89.6} transcript populations during the course of a single round of infection in HOS-CD4-CCR5 cells. A sensitive method for comparison among conditions involves quantifying utilization of six mutually exclusive splice acceptors A3, A4c, A4a, A4b, A5 and a novel acceptor just downstream of A5 termed A5a. Splicing at these acceptors determines the relative levels of messages encoding Tat and Env/Vpu in the partially spliced class and messages encoding Tat, Rev and Nef in the completely spliced class.

We observed longitudinal changes in the levels of these messages in HOS cells over 12–48 h that were statistically significant ($p < 10^{-10}$; generalized linear model described in Supplementary Report S2). This pattern was especially evident in junctions involving donor 1 spliced to each of these acceptors (Figure 3.4A). Most dramatically, transcripts with splicing junctions between D1 and A3 (tat messages) increased with time ($p <$

10^{-10}), while D1^A4b junctions (used in *env/vpu* or *rev* messages) were used reciprocally less ($p < 10^{-10}$). Such kinetic changes affecting specific transcripts both with and without the Rev-response element cannot be explained by the accumulation of Rev, and they may reflect differential transcript stability or HIV-induced alterations to the host splicing machinery. Temporal changes in HOS cells were confirmed using end point RT-PCR and analysis after electrophoresis on ethidium-stained gels (Figure 3.4B).

3.4.6 Cell-type-specific splicing patterns

We also compared splicing between T cells and HOS cells and found significant cell type differences ($p < 10^{-10}$). For example, while transcripts with D1^A5 junctions were dominant in both cell types, messages using the D1^A4c splice junction (encoding Env/Vpu or Rev) made up the bulk of the remaining transcripts in T cells but were a minor species in HOS-CD4-CCR5 cells. Likewise, Tat messages (using A3), which were quite abundant in HOS cells at all time points, contributed relatively little to populations of transcripts in primary T cells harvested at 48 hpi (Figure 3.4A). We also used end point PCR and analysis on ethidium-bromide-stained gels to confirm that the relative ratios of transcripts containing junctions to A3, A4a, A4b and A4c were different in HOS and T cells (Figure 3.4B).

3.4.7 Human variation in HIV-1 splicing

Quantitative comparisons also revealed modest differences in splicing between primary CD4⁺ T cells isolated from different human donors that were statistically significant ($p < 10^{-10}$) under a generalized linear model (Figure 3.4A). The magnitudes of predicted differences were small, all $< 33\%$ and most $< 10\%$.

3.5 Discussion

Use of single-molecule enrichment and long-read single-molecule sequencing has made possible the most complete study to date of the composition of HIV-1 message popula-

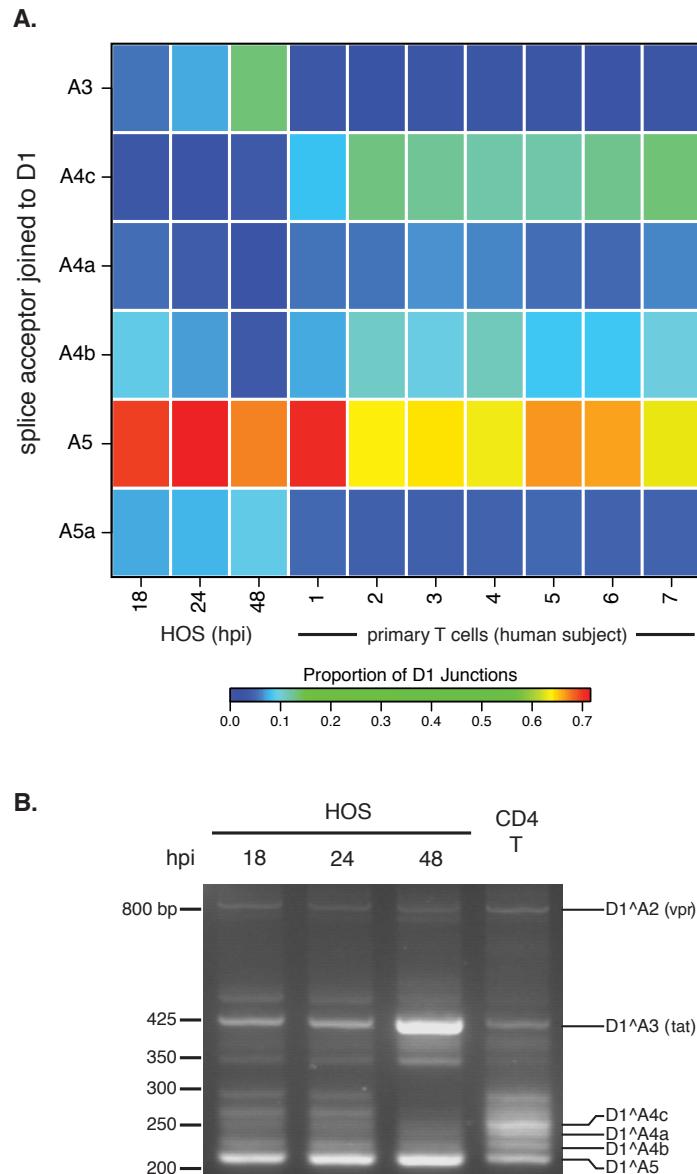


Figure 3.4: Temporal, cell type and donor variability in accumulation of HIV-1 messages. (A) In order to highlight changes in ratios of HIV-1 transcripts accumulating over time during infection and between HOS-CD4-CCR5 cells and primary T cells, we used PacBio read counts to calculate proportions of transcripts with splicing from the first major splice donor, D1, to each of the mutually exclusive acceptors: A3, A4c, A4a, A4c, A5 and the novel putative acceptor A5a. The heat map shows average data for T cell and HOS cell samples in columns with the color tiles indicating the proportion of D1 splicing to each of the mutually exclusive acceptors (rows), according to the color scale shown. (B) Reverse transcription and bulk PCR amplification of HIV_{89.6} transcripts from HOS cells and primary T cells from one human subject (subject 3) resolved by agarose gel electrophoresis and stained with ethidium bromide verified temporal and cell type changes shown in (A).

tions, revealing several new layers of regulation. Studies of the low-passage HIV_{89.6} isolate in a relevant cell type showed numerous differences from studies of lab-adapted HIV strains in transformed cell lines, highlighting the importance of studying the most relevant models. These data also illustrate the limitations of gel-based assays for studying HIV-1 message population. Multiple different combinations of HIV-1 exons yield mRNAs of similar sizes that are easily confused in typical assays using gel electrophoresis. Thus, in many settings the more detailed information provided by single-molecule amplification and single-molecule DNA sequencing is more useful.

Using these methods, we have detected significant variations between HIV message populations generated in T cells from different human donors. The differences were modest compared to those observed between cell types or time points, perhaps not surprisingly since any human polymorphisms strongly affecting mRNA processing might interfere with normal gene expression. However, because tight calibration of message levels is important to HIV-1, the observed differences in message ratios might affect HIV-1 acquisition or disease progression. The variation in observed transcripts could also be affected by different kinetics of infection in T cells from the different donors. In either case, these data suggest that human polymorphisms may exist that affect HIV-1 message populations in infected individuals, providing a new candidate mechanism connecting human genetic variation with measures of HIV disease.

Sequences from the 89.6 viral strain revealed a class of small (~1 kb) completely spliced transcripts, most contributed by splicing to a new poorly conserved acceptor A8c. These encoded two new proteins, one of which had Tat activity, and we showed that another, a Rev-Nef fusion termed Ref, could be detected in cells. HIV_{89.6} is a particularly cytotoxic virus isolated from the CSF of a patient, and it forms unusually large syncitia in macrophages³⁷⁵. The abundance of 1-kb transcripts produced by this virus provides a possible explanation for its unique properties. In addition to the novel acceptor A8c, we have also identified 3 putative novel splice donors and 11 putative novel acceptors,

which require further studied to clarify possible functions.

The wealth of new messages found here in HIV_{89.6} and in other HIV-1 isolates suggests there may be ongoing evolution of novel splice sites and new ORFs. Because splice acceptors in HIV-1 are weak²¹⁹, mutations creating sequences that even slightly resemble the 3' splice site consensus may be occasionally recruited as novel acceptors, creating new mRNAs. In fact, new splice signals may evolve with relative ease—it has been estimated that reasonable matches to the consensus for splice donors, acceptors and branch-point sites occur within random sequence every 290, 490 and 24 bp, respectively³⁸⁹, though sequence substitutions in HIV are usually also constrained by overlapping viral coding regions. We and others have observed appearance of novel exons within the major HIV-1 introns^{242,359,360}. Such long stretches of RNA relatively devoid of competing splice sites may be particularly poised to evolve new signals. On the other hand, most of the putative novel splice acceptors we observed clustered near previously identified acceptors in HIV-1, suggesting that conserved cis-acting splicing signals may recruit factors that act promiscuously on new nearby sequences. Clusters of splice sites might also provide redundancies that protect vital messages, as suggested previously^{390,391}. Frequent evolution of new splice sites may allow viruses to test out new combinations of exons, potentially yielding new RNAs and proteins, like those reported here. However, such novelty must compete with immune constraints—unstable novel polypeptides like Ref can be targeted to the proteasome and presented on MHC molecules as new epitopes for immune recognition.

HIV has likely evolved to produce calibrated message populations in T cells which seem to be altered with relative ease, as in infection in HOS cells, suggesting that therapeutic disruption of correct splicing may be feasible. A few studies have begun to explore small molecule therapy to disrupt HIV-1 splicing^{227,364}. Several factors could be responsible for the differences we observed between HOS and T cells, including hnRNP A/B and H, SC35, SF2/ASF and SRp40^{276,392}. Inhibition of SF2/ASF has already been shown to

abrogate HIV-1 replication in vitro²²⁷. Thus the lability seen here for function of these factors suggests they may be attractive antiretroviral targets.

3.6 Acknowledgements

We would like to thank the University of Pennsylvania Center for AIDS Research (CFAR) for preparation of viral stocks and isolation of primary CD4⁺ T cells; James A. Hoxie, Ronald G. Collman, Jianxin You, Robert W. Doms, Paul Bates, David Rekosh and members of the Bushman laboratory for reagents, helpful discussion and technical expertise. F.D.B., K.T., D.L., E.S., K.E.O. and R.M. conceived and designed the experiment. K.E.O. and R.C.A. carried out sample preparation and experimental validation. P.D. and J.O. performed single-molecule amplification. K.T. and S.W. performed sequencing. S.S.-M., K.E.O. and M.B. analyzed the data. K.E.O., F.D.B. and S.S.-M. wrote the manuscript.

CHAPTER 4 : Gene activity in primary T cells infected with HIV_{89.6}: intron retention and induction of distinctive genomic repeats

4.1 Abstract

Background: HIV infection has been reported to alter cellular gene activity, but published studies have commonly assayed transformed cell lines and lab-adapted HIV strains, yielding inconsistent results. Here we carried out a deep RNA-Seq analysis of primary human T cells infected with the low passage HIV isolate HIV_{89.6}.

Results: Seventeen percent of cellular genes showed altered activity 48 hours after infection. In a meta-analysis including four other studies, our data differed from studies of transcription after HIV infection of cell lines but showed more parallels with infections of primary cells. We found a global trend toward retention of introns after infection, suggestive of a novel cellular response to infection. HIV_{89.6} infection was also associated with activation of human endogenous retroviruses (HERVs) and several retrotransposons, of interest as possible novel antigens that could serve as vaccine targets. The most highly activated group of HERVs was a subset of the ERV-9, a group not reported previously to be induced by HIV. Analysis showed that activation was associated with a particular variant of an ERV-9 long terminal repeat that contains an indel near the U3-R border. These data also allowed quantification of > 70 splice forms of the HIV_{89.6} RNA and specified the main types of chimeric HIV_{89.6}-host RNAs. Comparison to 147,281 integration site sequences from the same infected cells allowed quantification of authentic versus artifactual chimeric reads (0.1% of the total), showing that 5' read-in, splicing out of HIV_{89.6} from the D4 donor and 3' read-through were the most common HIV_{89.6}-host cell chimeric RNA forms.

Conclusions: Analysis of RNA abundance after infection of primary T cells with the low passage HIV_{89.6} isolate disclosed multiple novel features of HIV-host interactions,

notably intron retention and induction of transcription of distinctive retrotransposons and endogenous retroviruses.

4.2 Background

HIV replication requires integration of a cDNA copy of the viral RNA genome into cellular chromosomes, followed by transcription and splicing to yield viral mRNA. Alternative splicing allows the small 9.1 kb HIV genome to generate at least 108 mRNA transcripts encoding at least 9 proteins and polyproteins^{192,219,240,361,393,394}. During replication, HIV also reprograms cellular transcription and splicing. For example, the virus-encoded Vpr protein arrests the cell cycle^{287,395–397} and the viral Tat protein binds to P-TEFb and alters transcript at the HIV promoter and some cellular promoters^{398–403}.

Multiple studies suggest that cells detect HIV infection and respond by inducing interferon-regulated, apoptotic and stress response pathways^{253,404–411}. Several studies have also suggested that HIV infection disrupts normal cellular splicing pathways^{374,411}. However, results have varied with many experimental parameters, including target cell type, HIV isolate and the duration of infection. Many of the published studies focused on infections with lab-adapted HIV strains in transformed cell lines^{251,253,293,404,411,412}, and so results may not be fully reflective of infections in patients.

In this study, we sought to generate data more resembling HIV replication in patients by analyzing transcriptional responses after infection of primary T cells with HIV_{89.6}, a low passage patient isolate³⁷⁵. This represents a continuation of a long term effort to understand HIV-host cell interactions at the transcriptional level that began with analysis of transcription by HIV_{89.6} in primary T cells using Pacific Biosciences long read single molecule sequencing³⁹⁴. Our strategy here was to analyze a single time after infection in depth, analyzing over 1 billion sequence reads from HIV_{89.6} infected and uninfected host cells. These data were then combined with 147,281 unique integration site sequences from the same infections and the Pacific Biosciences data on HIV_{89.6}

transcription to 1) elucidate effects of HIV infection on host cell mRNA abundances and splicing, 2) characterize viral message structure in detail and 3) probe the nature of the chimeras formed between host cell and viral RNAs.

4.3 Methods

4.3.1 Cell culture and viral infections

HIV_{89.6} stocks were generated by the University of Pennsylvania Center for Aids Research. 293T cells were transfected with a plasmid encoding an HIV_{89.6} provirus, and harvested virus was passaged in SupT1 cells once. Viral stocks were quantified by measuring p24 antigen content. Primary CD4⁺ T cells were isolated by the University of Pennsylvania Center for AIDS research Immunology Core from apheresis product from a single healthy male donor (ND365) using the RosetteSep Human CD4⁺ T Cell Enrichment Cocktail (StemCell Technologies).

T cells were stimulated for 3 days at 0.5×10^6 cells per milliliter in R10 media (RPMI 1640 with GlutaMAX (Invitrogen) supplemented with 10% FBS (Sigma-Aldrich) with 100 units U/mL recombinant IL2 (Novartis) + 5 μ g/mL PHA-L (Sigma-Aldrich)). Cells were infected in triplicate and mock infections were performed in duplicate. For each infection, 6.6×10^6 cells were mixed with 1.32 μ g HIV_{89.6} in a total volume of 2.25 mL. Infection mixtures was split into three wells of a 6 well plate for spinoculation at 1200 g for 2 hr at 37°C. Cells were incubated an additional 2 hr at 37°C. Cells were then pooled into flasks and volume was increased to a total of 12 mL. Spreading infection was allowed to proceed 48 hr at 37°C, after which cells were harvested. 1×10^6 cells were harvested for flow cytometry, and 6×10^6 cells were pelleted following two washes in PBS for nucleic acid extraction. Genomic DNA and total RNA were isolated from 6×10^6 T cells per infection using the AllPrep DNA/RNA Mini Kit (Qiagen) with Qiashredder columns (Qiagen) for homogenization according to the manufacturer's instructions. DNA was eluted in 140 μ L elution buffer. RNA samples were treated with DNase prior

to elution in 40 μ L water.

4.3.2 Analysis of HIV_{89.6} integration sites in primary T cells

Integration site sequences were determined for DNA fractions from the above infections after ligation mediated PCR³³². A total of 147,281 unique integration site sequences were determined. An analysis of integration site distributions for these samples was reported in Berry et al.³³².

4.3.3 mRNA sequencing

Messenger RNA was isolated and amplified from purified total cellular RNA (3 μ L or approximately 9 μ g from each uninfected sample, 25 μ L or approximately 3 μ g from each infected sample) using the Illumina TruSeq RNA sample preparation kit according to manufacturer's protocol. SuperScript III (Invitrogen) was used for reverse transcription. Each sample was tagged with a separate barcode and sequenced on an Illumina HiSeq 2000 using 100-bp paired-end chemistry.

4.3.4 Flow cytometry

To assess percent infected cells, 1×10^6 cells per infection were stained for flow cytometry. All staining incubations were at room temperature. Cells were first washed in PBS and then twice in FACS wash buffer (PBS, 2.5% FBS, 2 mM EDTA). Cells were fixed and permeabilized with CytoFix/CytoPerm (BD) for 20 minutes and washed with Perm-Wash Buffer (BD) before staining with anti-HIV-Gag-PE (Beckman Coulter) for 60 min. Finally cells were washed in FACS wash buffer and resuspended in 3% PFA. Samples were run on a LSRII (BD) and analyzed with FlowJo 8.8.6 (Treestar). Cells were gated as follows: lymphocytes (SSC-A by FSC-A), then singlets (FSC-A by FSC-H), then by Gag expression (FSC-A by Gag).

4.3.5 Analysis

Reads were aligned to the human genome using a combination of BLAT³³¹ and Bowtie⁴¹³ through the Rum pipeline⁴¹⁴. Estimates of fragments per kilobase of transcript per million mapped reads and changes in expression for cellular genes were calculated by Cufflinks²⁹⁰. Reads found to contain sequence similar to the HIV genome using a suffix tree algorithm were aligned against the HIV_{89.6} genome using BLAT³³¹. All statistical analyses were performed in R 3.1.2³²³. RNA-Seq reads from Chang et al.²⁵³ were downloaded from the Sequence Read Archive (SRP013224) and aligned using the Rum pipeline.

Gene lists were obtained from the supplementary materials of four other studies of differential gene expression during HIV infection^{253,293,409,415}. We called genes differentially expressed in Li et al.⁴¹⁵ if they had a reported $p < 0.01$ or in Lefebvre et al.²⁹³, Chang et al.²⁵³ and Imbeault et al.⁴⁰⁹ if they had an adjusted $p < 0.05$. We called genes as differentially expressed in our own study if the adjusted $p < 0.01$. For the comparison of differentially expressed genes regardless of direction in figure 4.1 (below the diagonal), it was unclear exactly how many genes were studied in each study so we assumed a background of the 14,192 genes (the number of genes which could be tested for significance in our data).

We obtained transcriptional profiles comparing immune cell subsets from the Molecular Signatures Database⁴¹⁶. MSigDB set names from the MSigDB used in Figure 4.2A were GSE10325 LUPUS CD4 TCELL VS LUPUS BCELL, GSE10325 CD4 TCELL VS MYELOID, GSE10325 CD4 TCELL VS BCELL, GSE10325 LUPUS CD4 TCELL VS LUPUS MYELOID, GSE3982 MEMORY CD4 TCELL VS TH1, GSE22886 CD4 TCELL VS BCELL NAIIVE, GSE11057 CD4 CENT MEM VS PBMC, GSE11057 CD4 EFF MEM VS PBMC, GSE3982 MEMORY CD4 TCELL VS TH2 and GSE11057 PBMC VS MEM CD4 TCELL and in Figure 4.2B were GSE36476 CTRL VS TSST ACT 72H MEMORY CD4 TCELL OLD, GSE10325 CD4 TCELL VS LUPUS CD4 TCELL, GSE22886 NAIIVE

CD4 TCELL VS 12H ACT TH1, GSE3982 CENT MEMORY CD4 TCELL VS TH1, GSE17974 CTRL VS ACT IL4 AND ANTI IL12 48H CD4 TCELL, GSE24634 IL4 VS CTRL TREATED NAIVE CD4 TCELL DAY5, GSE24634 NAIVE CD4 TCELL VS DAY10 IL4 CONV TREG, GSE1460 CD4 THYMOCYTE VS THYMIC STROMAL CELL and GSE1460 INTRATHYMIC T PROGENITOR VS NAIVE CD4 TCELL ADULT BLOOD.

We downloaded the RepeatMasker track from the UCSC genome browser⁴¹⁷ and used the SAMtools library⁴¹⁸ to assign reads to the repeat regions. HERV-K age estimates were obtained from the supplementary materials of Subramanian et al.⁴¹⁹.

We used a Bayesian estimate of the ratio of expression in uninfected and HIV infected samples to account for sampling effort and differing expression in genomic regions. We modeled the observed counts as a binomial distribution with a flat beta prior ($\alpha = 1, \beta = 1$) separately for uninfected and infected samples. We then Monte Carlo sampled the two posterior distribution to estimate the posterior distribution of the ratio. For introns, the number of binomial successes was set to the number of reads mapped to the intron and the number of trials was the total number of reads observed in the genes overlapping that intron. For repeat regions, the number of binomial successes was set to the number of reads mapped to that region and the number of trials was the total number of reads mapped to the human genome.

To estimate determinants of LTR12C expression, we fit a logistic regression for which LTR12C increased in expression with HIV_{89.6} infection (95% Bayesian credible interval > 1) on to characteristics of the LTR12C regions. We extracted all the LTR12C regions from the human genome and determined the U3-R boundary using a ends free alignment of the previously reported U3-R border^{420–424} against the sequences. Regions less than 1,000 bases long were discarded. Previous studies disagreed about the location of the LTR12C transcription start site and it appears that transcription may start in several places^{421,422}. We took the 5' most site that had agreement between

studies (transcription starting with TGGCAACCC). We split the sequences into short, medium and long length classes based on an indel about 70 bases upstream from the transcription start site. For each length class, we generated a consensus sequence and counted the Levenshtein edit distance between the consensuses and each corresponding sequence. We also counted the number of NFY motifs (CCAAT or ATTGG), MZF1 motifs (GTGGGGA) and GATA2 motifs (GATA or TATC) in the entire U3 region or checked if any of the three motifs was present in the 150 bases upstream of the TSS. A final regression model was selected using stepwise regression with an AIC cutoff of 5. For display, the LTR12C sequences were aligned with MUSCLE⁴²⁵.

The abundance of the HIV RNA size classes was estimated as described in Additional File 5. These estimates were then multiplied by the within size class proportions estimated by Ociejewa et al.³⁹⁴ using PacBio sequencing of HIV_{89.6} to yield proportions over 78 measured HIV_{89.6} RNAs.

4.4 Results

4.4.1 Infections studied

HIV_{89.6}, a clade B primary clinical isolate³⁷⁵, was used to infect primary CD4⁺ T cells from a single human donor in three replicate infections. For comparison, two additional replicates from the same donor were mock infected. Samples were harvested after 48 hours of infection, which allowed for widespread infection in the primary T cell cultures, though some cells may be infected secondarily by viruses produced in the first round. Thus cultures probably were not tightly synchronized but did have extensive representation of infected primary T cells. From these samples, we obtained 1,161,705,678 101-bp reads from primary CD4⁺ T cells from a single donor; 1,021,207,853 were mapped to the human genome and 24,783,844 to the HIV_{89.6} provirus (Table 4.1). Below we first discuss the influence of infection on cellular gene activity and RNA splicing, then analyze HIV RNAs and lastly analyze chimeras formed between HIV and cellular

Sample	Infection rate (%)	Reads	Human reads	HIV reads	% HIV	% HIV in infected
Uninfected-1	—	232,450,106	212,391,460	—	—	—
Uninfected-2	—	235,048,212	203,760,783	—	—	—
Infected-1	37.5	234,378,088	199,871,662	10,219,315	4.86	13.0
Infected-2	26	226,078,422	198,436,507	7,322,556	3.56	13.7
Infected-3	21	233,750,850	205,747,441	7,241,973	3.40	16.2

Table 4.1: Samples used in this study, their infection rates and sequencing depth.

RNAs.

4.4.2 Changes in gene activity in primary T cells upon infection with HIV^{89,6}

Changes in host cell gene expression have been reported during HIV infection^{251,252,293,404,409,410} and differences in expression have been observed associated with the stage⁴¹⁵ and progression⁴²⁶ of disease. Here we observed significant changes in gene expression (false discovery rate corrected $q < 0.01$) in 3,142 genes, 17.1% of expressed cellular genes (Additional file 1). The genes with most extreme increases, all $> 6 \times$ fold higher, during HIV infection included IFI44L, RSAD2, HMOX1, MX1, USP18, IGJ, OAS1, CMPK2, DDX60, IFI44, IFI6, IFNG and CCL3. All of these have been reported to be involved in innate immunity⁴²⁷ or are interferon inducible⁴²⁸, highlighting a strong innate immune response in the cells studied. Genes with the largest decreases, all $> 3 \times$ fold lower, were GNG4, GPA33, IL6R, CCR8, RORC, AFF2 and CCR2.

Many gene ontology categories were significantly enriched for differentially expressed genes (Additional file 2). Notably upregulated with infection were genes involved in apoptosis, immune responses and cytokine production (all $q < 10^{-4}$) and down-regulated were genes involved in viral gene expression, nonsense-mediated decay and translation elongation and termination (all $q < 10^{-19}$). These changes suggest that the cells responded to HIV infection with the induction of inflammatory, interferon regulated and apoptotic responses, patterns posited from several previous studies^{253,293,404–410,412,429}. Several genes were activated that were characteristic of other hematopoietic lineages, e.g. hemoglobin β , CD8, CD20 and CD117, while several CD4 $^+$ T cell specific genes, e.g.

Cell type	HIV type	Differentially expressed genes (Up/Down)	Study
Primary CD4 ⁺ T	HIV _{89.6}	3393 (1756/1637)	This study
Primary CD4 ⁺ T	NL4-3 BAL-IRES-HSA	228 (182/46)	Imbeault et al. ⁴⁰⁹
Lymph node biopsies	Acute infection	448 (383/65)	Li et al. ⁴¹⁵
SupT1	HIV _{LAI}	4997 (2666/2331)	Chang et al. ²⁵³
SupT1	NL4-3Δenv-eGFP/VSV-G	579 (212/367)	Lefebvre et al. ²⁹³

Table 4.2: Data from this study and four others used for meta-analysis of human gene expression changes during HIV infection

CD4 and CD3, were downregulated, potentially consistent with de-differentiation of infected and bystander cells. We return to this point in the discussion.

4.4.3 Comparison of transcriptional profiles from HIV_{89.6} infection of primary T cells to data on HIV infection in other cell types

We sought to identify the transcriptional responses that were most conserved upon HIV infection and so collected and analyzed data from four other studies of transcription in HIV-infected cells (Table 4.2). These included two studies of infection of the SupT1 cell line^{253,293}, a study of primary CD4⁺ T cells⁴⁰⁹ and a study of lymphatic tissue in acutely viremic patients⁴¹⁵. Genes were scored as increased or decreased in activity after infection, and the amount of agreement was compared among the different studies.

No gene was called as differentially expressed in all five studies. Eight genes were differentially expressed in the same direction in 4 out of 5 studies; AQP3 and EPHX2 were down-regulated with HIV infection and CD70, EGR1, FOS, ISG20, RGS16 and SAMD9L were up-regulated. A full listing is provided in Additional file 4. Several of the up-regulated genes are known to be interferon inducible, again emphasizing the role of innate immune pathways.

For each pair of studies, we compared whether they agreed on the identities of differentially expressed genes and whether they agreed on the direction of change (Figure 4.1). The estimated alterations in gene activity showed notable differences in the responses to infection in primary cells versus the SupT1 cell line. The two SupT1

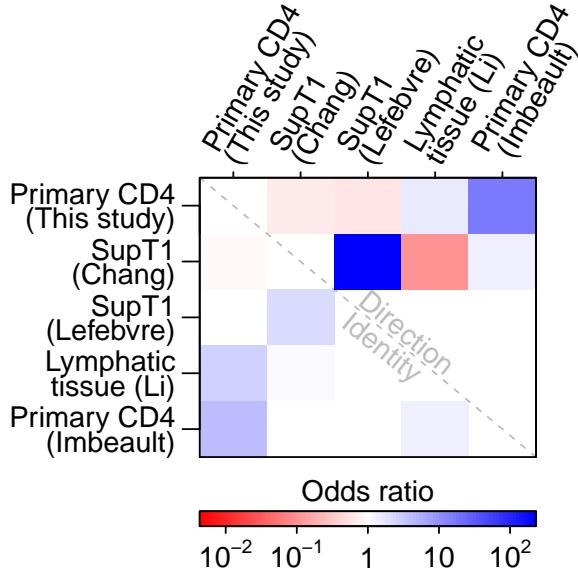


Figure 4.1: Comparisons among studies quantifying cellular gene expression after HIV infection. For each pair of studies, the association between up- and down-regulation calls was measured for genes identified by both studies as differentially expressed (above the diagonal). As another comparison, we also measured the agreement between studies for which genes were called differentially expressed regardless of direction (below the diagonal). The color scale shows the conservative (i.e. closest to 1) boundary of the confidence interval of the odds ratio with blue indicating a positive association and red a negative association between studies. For confidence intervals overlapping 1, the value was set to 1. Therefore all colored squares indicate significant associations.

studies were significantly similar ($p < 10^{-15}$) to each other but were not significantly associated (Lefebvre et al.²⁹³, $p = 0.2$) or were negatively associated (Chang et al.²⁵³, $p = 10^{-7}$) with data from lymphatic tissue in acute HIV patients. The primary T cell study reported here was significantly associated with the second study in primary cells ($p < 10^{-15}$) and with a study of lymphatic tissue from patients acutely infected with HIV ($p = 0.003$). Our primary T cell data was negatively associated with the SupT1 studies (both $p < 10^{-3}$). This documents significant differences in responses to HIV infection between infected primary cells and SupT1 cells and suggests that results of infections in primary cells more closely align with actual acute HIV infections in patients. SupT1 cells might be expected to respond to infection differently than primary cells since they have several nonsynonymous mutations in innate immunity genes⁴³⁰, have blocks in immune signaling pathways⁴³¹ and fail to activate many interferon stimulated genes during HIV infection⁴¹⁰.

4.4.4 Comparison of the HIV infected cell transcriptional profiles to additional experimental T cell profiles

To investigate the transcriptional changes in more depth, we compared the results of the five studies of HIV infection to transcriptional profiles comparing immune cell subsets available at the Molecular Signatures Database (MSigDB)⁴¹⁶. The MSigDB reports genes that are increased or decreased in relative expression for each of 185 pairs of transcriptional profiles involving CD4⁺ T cells. We compared the lists of affected genes in each pair to genes altered in activity by HIV infection. Those pairs of studies with the most significant associations with HIV_{89.6} data are shown in Figure 4.2A. For comparison, the associations with the four other HIV transcriptional profiling studies mentioned above are shown as well.

The most significant associations for our data showed gene expression in HIV_{89.6}-infected cells moving away from typical T cell expression patterns and towards patterns more similar to B cells, myeloid cells and bulk peripheral blood mononuclear cells (all Fisher's $p < 10^{-15}$) (Figure 4.2A). These changes were also seen, although to a lesser extent, in the Imbeault et al.⁴³² study which also used primary CD4⁺ T cells.

For comparison, we also extracted those profiles most strongly associated with the transcriptional data on lymphatic tissue of HIV patients⁴¹⁵. The profiles showed patterns similar to strongly stimulated T cells, autoimmune disease and to the Th1 T cell subset (all $p < 0.01$) (Figure 4.2B). Our data in primary CD4⁺ T cells paralleled the changes seen in lymphatic tissue. These transcriptional changes again highlight the strong immune response generated by HIV infection in primary cells.

4.4.5 Intron retention

Cells respond to infection by shutting down macromolecular synthesis at multiple levels^{433–437}, so we investigated whether cells also showed perturbations in splicing efficiency after infection. As a probe, we created a database of cellular genomic re-

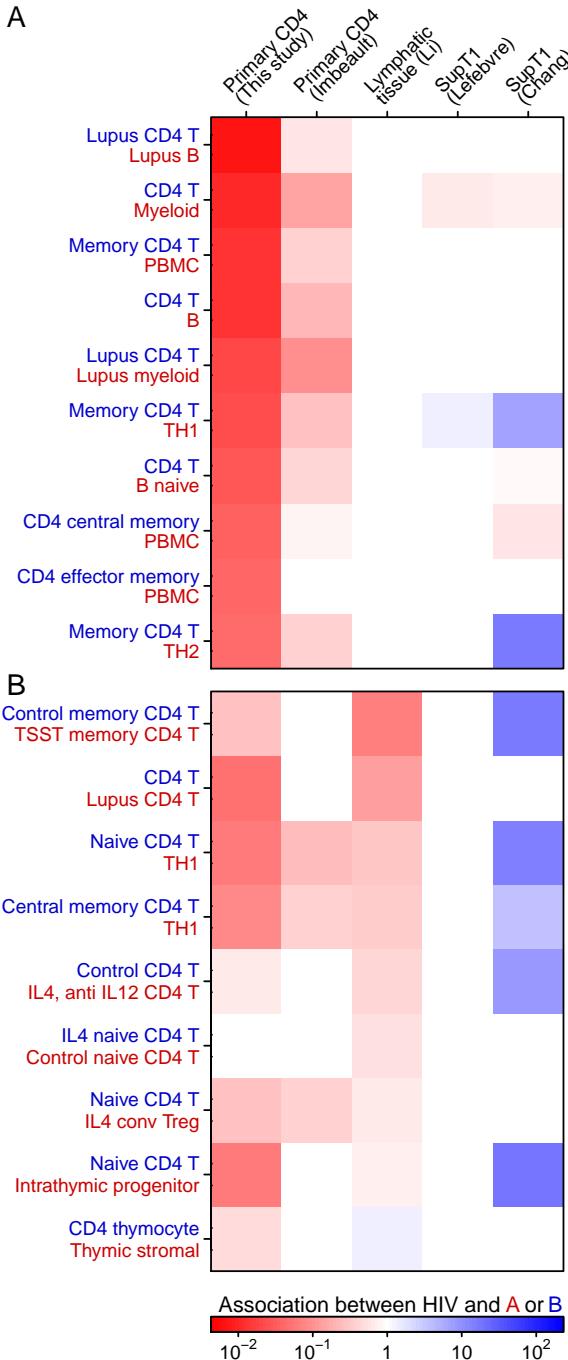


Figure 4.2: Comparisons of the effect of HIV infection on gene expression to studies comparing subsets of immune cells. The MSigDB database was used to extract 185 sets of differentially expressed genes from pairs of transcriptional profiling studies of immune cell subsets involving CD4⁺ T cells. For each pair of studies, we used Fisher's exact test to measure the association between up- and down-regulation calls for genes identified as differentially expressed in both our HIV study and the comparator immune subsets. A) The transcriptional profiles with strongest associations with changes observed in our study of HIV_{89.6} infection of primary T cells. Blue indicates a positive association between changes seen in HIV infected cells and the first immune subset (text colored blue) while red indicates a positive association with the second immune subset (text colored red). The color scale shows the conservative (i.e. closest to 1) boundary of the confidence interval of the odds ratio. For confidence intervals overlapping 1, the value was set to 1. Therefore all colored squares indicate significant associations. B) As in A, but showing the transcriptional profiles most strongly associated with changes observed in lymph node biopsies from acutely infected patients⁴¹⁵.

gions annotated exclusively as exons or introns in all spliceforms in the UCSC gene database³⁴² and quantified expression in these regions in infected and uninfected cells. We found a significant increase in intronic sequences relative to exonic sequence (Wilcoxon $p < 10^{-15}$) (Figure 4.3A). This increase in intronic sequence was reproducible between replicates in our study (Kendall's $\tau=0.42$, $p < 10^{-15}$) (Figure 4.3B). We reanalyzed RNA-Seq data from Chang et al.²⁵³ and also documented intron retention which correlated with the changes seen in our data (Kendall's $\tau=0.12$, $p < 10^{-15}$) (Figure 4.3C).

A possible artifactual explanation for enrichment of intronic sequences could involve greater DNA contamination in the infected cells samples. That is, if the relative amount of DNA differed between treatments, the amount of apparent intronic sequences could also differ due to sequencing of contaminating DNA. To examine whether DNA contamination was abundant in our samples, we compiled a collection of 27 large gene desert regions, defined here as 1) regions outside the centrosome and first and last cytoband, 2) containing less than 1% unknown sequence, 3) containing no genes annotated in UCSC genes³⁴², 4) containing no repeats annotated in the repeatMasker database³⁴⁷ and 5) spanning more than 100 kb. No reads were mapped to these 41 Mb of gene deserts in any sample, arguing against explanations based on DNA contamination. Thus these data indicate that intron retention was increased in these cell populations upon HIV infection, revealing a previously undisclosed aspect of the host cell transcriptional response to infection.

Previous studies have reported changes in the expression and localization of splicing factors with HIV infection^{374,438,439}. In our data, HIV_{89.6} infection significantly altered the expression of genes involved in RNA splicing ($p = 2 \times 10^{-7}$) and nonsense-mediated decay ($p < 10^{-15}$). Genes related to nonsense-mediated decay genes showed a strong pattern of lowered RNA abundance, with 71 out of 118 annotated genes significantly lower in expression after infection. These patterns suggest potential mechanisms for the intron retention observed here.

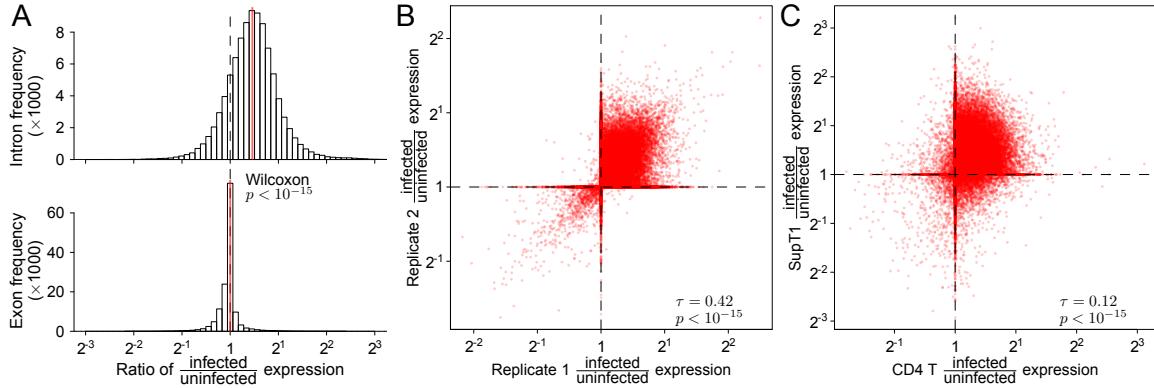


Figure 4.3: Changes in the abundance of intronic regions with HIV infection. Expression of intronic and exonic regions was quantified as the proportion of reads mapping within the intron/exon out of the total reads mapping to the transcription units overlapping that intron/exon. A) Comparison of the ratios of expression between infected and uninfected replicates in exclusively intronic or exonic regions of transcription units. B) Reproducibility of intron retention between replicates. Each point quantifies the change in expression with HIV infection for a specific intronic region. The x-axis shows changes in gene activity accompanying infection for one set of replicates (Infected-1 and Infected-2 vs. Uninfected-1) and the y-axis shows the same data for different replicates (Infected-3 vs. Uninfected-2). C) Reproducibility of intron retention between studies. The plot is arranged as in B but with all data from our study combined on the x-axis and corresponding data from Chang et al.²⁵³ on the y-axis.

4.4.6 Induction of transcription from HERVs and LINEs by HIV_{89.6} infection

HIV infection has been reported to induce expression of certain HERVs, particularly HERV-K^{440–442}, and LINE and Alu transposable elements⁴⁴³, providing candidate markers of infection and possible vaccine targets. Thus we analyzed our data in primary T cells infected with HIV_{89.6} to investigate the expression of HERVs, LINEs and other repeated sequences. Figure 4.4A shows a comparison of the association between changes in expression with HIV_{89.6} infection and the various genomic repeat types over varying levels of differential expression. At high levels of expression, ERV-9 (odds ratio at 4× expression: 152, 95% CI:82.5–259) and its long terminal repeat LTR12C (odds ratio at 4× expression: 144, 95% CI: 98.2–207) are the only repeats highly associated with upregulation during HIV infection. Looking at genomic repeats with any significant increase, the expression of many recently acquired genomic repeats, including L1HS, LTR5_Hs (a human specific LTR of HERV-K), AluYa5, AluYg6 and

SVA_D and SVA_F, were associated with HIV_{89.6} infection (Figure 4.4B).

We saw a relationship between the age of genomic repeats and its likelihood of being induced by HIV_{89.6} infection. The most highly enriched repeats were associated with relatively recent hominid-specific repeat classes as annotated by the RepeatMasker database (repeat classes with $p < 10^{-50}$ odds ratio: 31.6, 95% CI: 8.88–112). In HERV-K (HML-2), the most recently active endogenous retrovirus in the human genome^{419,444,445}, we saw that integrations unique to the human genome⁴¹⁹ were more likely to be differentially expressed than older HERV-Ks (odds ratio: 5.38, 95% CI: 1.93–16.0).

Previous RNA-Seq studies of cellular expression during HIV infection in transformed cell lines did not report increases in HERV mRNA^{253,293}. To investigate this difference, we downloaded and analyzed the RNA-Seq data from Chang et al.²⁵³, which quantified gene activity in transformed SupT1 cells infected with a lab-adapted strain of HIV. We found a much higher level of HERV expression in their data in both HIV infected cells and uninfected controls than in primary cells (Figure 4.4C). We suspect that in SupT1 cells, as with many cancerous cells^{446–450}, the baseline expression of transposons and endogenous retroviruses is higher than in primary cells, masking further induction by HIV infection.

We observed heterogeneous expression among ERV-9/LTR12C sequences and so investigated the primary sequence determinants. We observed that ERV-9/LTR12C has three variants of differing length in the U3 region just upstream of the transcription start site (Figure 4.5A), an important region for transcription initiation⁴²¹. The U3 region of LTR12C also contains multiple motifs for transcription factors NFY, GATA2 and MZF1⁴²⁴. To clarify factors affecting expression levels, we counted the number of motifs matching these transcription factors, assigned each LTR12C to one of the length classes, counted the number of mutations away from the consensus for that length class and checked for integration in a transcription unit. We then carried out a regression analysis to test the effects of these variables on LTR12C differential expression. We

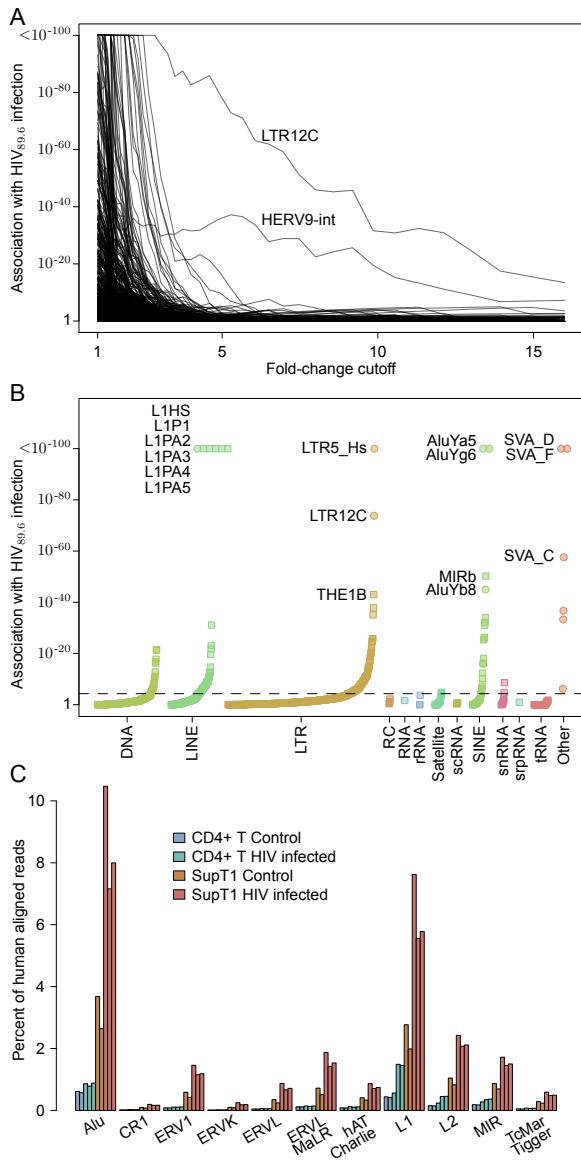


Figure 4.4: Repeat categories enriched upon infection with HIV. A) The association of repeat regions differentially expressed after HIV_{89.6} infection of primary T cells observed for varying thresholds of differential expression. The threshold used to call a gene differentially expressed based on the Bayesian posterior median was varied and Fisher's exact test was used to assess whether any genomic repeats had a significant association with this differential expression. Note that only ERV-9 (annotated as HERV9-int in the RepeatMasker database) and its corresponding long terminal repeat LTR12C were significantly associated with large changes in expression. B) Enrichment of repeat categories in regions differentially expressed (Bayesian 95% credible interval > 1) between HIV-infected and control CD4⁺ T cells. The repeated sequences are ordered on the x-axis by the extent of induction within each class, the y-axis shows the p-value for upregulation after infection. The dashed line indicates a Bonferroni corrected *p* value of 0.05. (C) The proportion of human mapped reads that align within classes of genomic repeats for data from primary CD4⁺ T cells from this study and SupT1 cells from Chang et al.²⁵³. A single read mapping multiple times to a given category was only counted once.

found that HIV_{89.6} induced transcription was more likely with the fewer mutations away from consensus, the number of locations matching the NFY transcription factor binding motif (CCAAT) and LTRs containing the short length variant of the 3' U3 region. The presence of a MZF1 motif near the transcription start site decreased transcription (Figure 4.5B).

4.4.7 HIV mRNA synthesis and splicing

Over 24 million Illumina reads mapped to HIV_{89.6}, yielding an average coverage of over 240,000-fold. Reads mapping to HIV_{89.6} comprised between 3.4–4.8% of mapped reads in the infected samples (Table 4.1). Assuming HIV-infected cells contain the same amount of mRNA as uninfected cells and adjusting for rates of infection ranging between 21–37.5% (Table 4.1), we estimate that HIV transcripts comprise between 13.0–16.2% of the total polyadenylated mRNA nucleotides in infected cells 48 hours after initial infection. This parallels previous estimates of around 10%⁴⁵¹ at 48 hours postinfection, 38% at 24 hours²⁵³ or 30% after 72 hours⁴⁰⁴.

Over 47,257 single reads spanned previously reported HIV splice junctions, allowing a quantitative assessment of donor and acceptor utilization (Figure 4.6A). As expected from previous studies^{240,394}, the most abundant junctions were D1-A5 and D4-A7. We confirmed the use of unusual splice acceptors A8c and A5a, previously reported in HIV_{89.6}³⁹⁴. In our data, we also see a higher abundance of D1-A1 and D1-A2 splice junctions than might be expected^{240,394}, although previous studies reported proportional abundance within size classes, making comparisons between size classes uncertain.

A 3' bias is apparent in our sequencing data (Additional file 5). This could be due to the poly-A capture step of the protocol where any break in the RNA would result in distal 5' sequences being lost⁴⁵². We used sequence reads from the large unspliced HIV intron 1 to measure this bias using a regression of the log of the number of fragments with a 5'-most end starting at a given position against the distance of that position from the

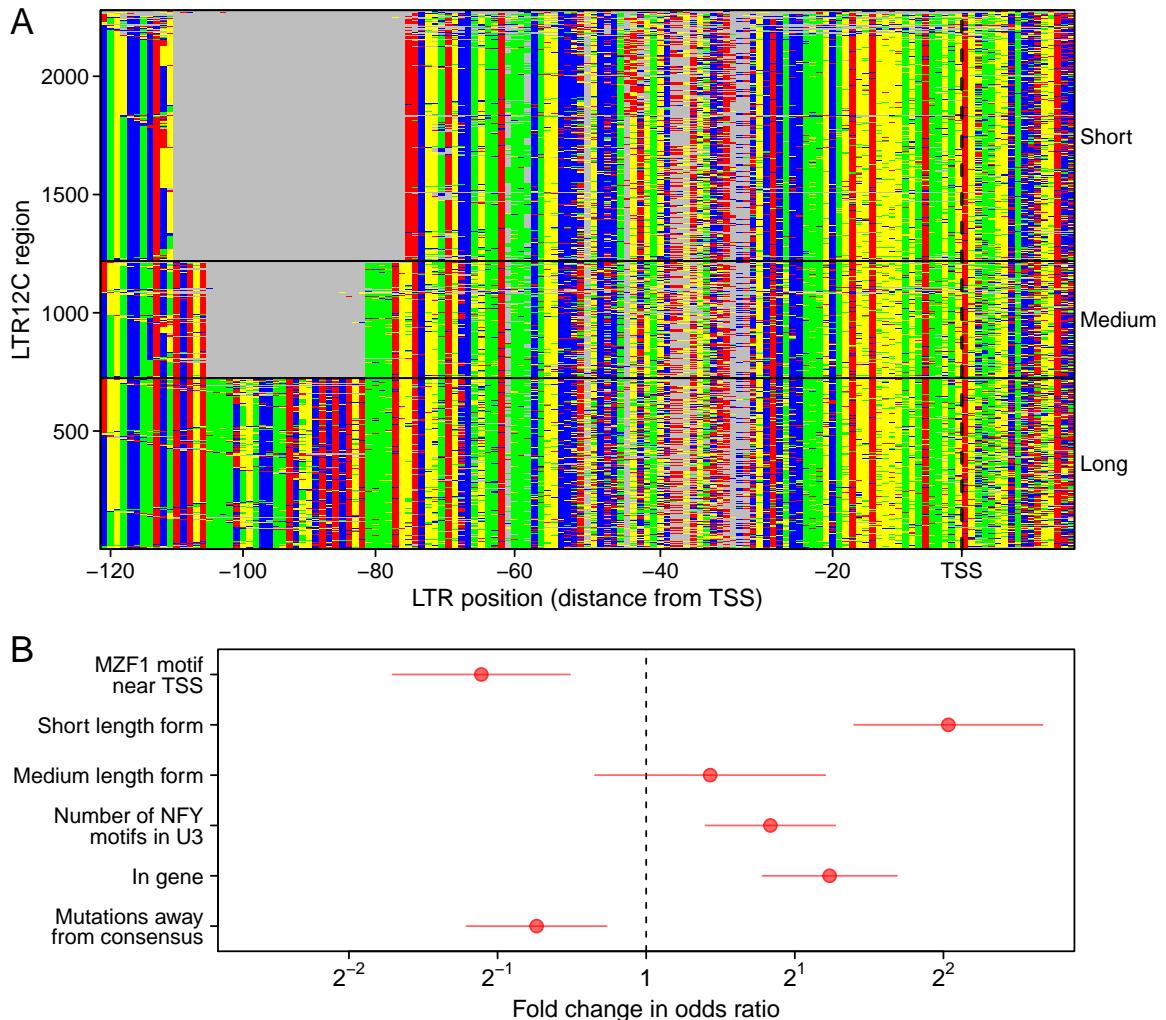


Figure 4.5: Characteristics of LTR12C sequences associated with induction upon infection of primary T cells with HIV_{89.6}. A) An alignment of the 3' end of the U3 region of repeats annotated as ERV-9 LTR12C. Each row is a LTR sequence and each column a base in that sequence colored by nucleotide identity. Three distinct classes are visible with a short, medium and long form. Mutations away from the consensus can also be seen. B) The coefficients (points) and ± 1.96 standard errors (horizontal lines) of a logistic regression comparing differential expression of LTR12C to the presence of MZF1 and NFY motifs, short/medium/long length alternate forms of the U3-R region, mutations away from the consensus for each length form and integration inside a transcription unit. The coefficient shown for mutations away from consensus is for a 10 mutation difference and the coefficient shown for NFY motifs is for a change of 5 additional motifs. All other coefficients are for binary values.

viral polyadenylation site, yielding an estimated probability of breakage of 0.021% per base (Additional file 5). Given this rate of termination, there is only a 14% chance of reaching the 5' end of the 9171 nt unspliced HIV genome ($(1 - 0.00021)^{9171}$).

Ocwieja et al.³⁹⁴ determined the relative abundance of HIV_{89.6} of similarly sized transcripts using PacBio single molecule sequencing, but were not able to estimate the relative abundance of all transcripts due to a sequencing bias favoring shorter transcripts. For this reason, relative abundances could only be specified within message size classes (i.e. the 4 kb, 2 kb and unexpectedly a 1 kb size class as well) and the overall quantitative abundances were unknown. The RNA-Seq data reported here are unable to determine complete transcript abundance because the short read length does not allow reconstruction of multiply spliced messages but do permit estimation of size class abundances after correcting for 3' bias (Additional file 5). Thus the PacBio data reported by Ocwieja et al.³⁹⁴ and the Illumina data reported here can be combined together to determine complete relative abundance of all HIV_{89.6} transcripts (Figure 4.6B).

The most abundant HIV mRNAs were the unspliced HIV genome (37.6%), a transcript encoding Nef (D1-A5-D4-A7: 15.5%), two 1 kb size class transcripts (D1-A5-D4-A8c: 10.6%, D1-A8c: 4.9%) and two Rev-encoding transcripts (D1-A4c-D4-A7: 4.2%, D1-A4b-D4-A7: 3.1%). The function of this large amount of 1 kb transcript is unknown. These two 1 kb transcripts do not appear to encode significant open reading frames although other 1 kb transcripts can encode a Rev-Nef fusion³⁹⁴.

Using these abundances, we can estimate the number of HIV_{89.6} genomes in these primary T cells 48 hours after infection. To determine the proportion of the mRNA nucleotides from viral transcripts, we multiplied the estimated abundances by their transcript lengths. Unspliced genome transcripts appear to form 79% of the mRNA nucleotides from HIV_{89.6} transcripts. Assuming T cells contain at least 0.1 pg of mRNA then an infected cell should contain at least 0.011 pg of unspliced HIV transcript

$(0.1\text{pg} \times 0.14 \frac{\text{HIV mRNA nt}}{\text{cell mRNA nt}} \times 0.79 \frac{\text{unspliced mRNA nt}}{\text{HIV mRNA nt}})$ or, assuming 9171 bases of RNA weigh about 5×10^{-6} pg, at least 2200 HIV genomes at 48 hour post infection. This estimate roughly agrees with previous estimates of HIV production per cell^{451,453,454}.

4.4.8 Human-HIV chimeric reads

The suggestion that HIV integration may disrupt cellular cancer-associated genes and thereby promote cell proliferation^{455–458} has focused attention on the range of novel message types formed when HIV integrates within transcription units^{324,356,459–461}. Chimeric reads containing HIV and cellular sequence are also of clinical interest due to the potential of lentiviral vectors to trigger oncogenesis in gene therapy patients through insertional mutagenesis^{462–465}.

In our data, 80,045 reads contained sequences matching to both HIV and human genomic DNA, but a considerable complication arises because chimeras can be formed artifactually during the preparation of libraries for sequence analysis^{466–473}. Many of the chimeric sequences in our data contained junctions between the HIV and human sequence where the ends of the human and HIV sequence were similar and potentially complementary (Figure 4.7A). This raises the concern that some of these chimeras could be products of in vitro recombinations during the reverse transcription, amplification and sequencing processes. Template switching between sequences with shared similarity is a well established property of retroviral reverse transcriptase enzymes used in RNA-Seq library preparation^{474–476}. Priming off incomplete transcripts during DNA synthesis is another potential source of chimeric transcripts^{466,467,477,478}. Failing to account for chimeras can hinder interpretation of deep sequencing data^{468–473}.

Also consistent with artifactual chimera formation, 7,354 reads (9.2% of chimeric messages) contained HIV sequences joined to human mitochondrial sequences, yet HIV proviruses have not previously been found integrated in mitochondrial DNA³⁵⁶. To probe this further, we used ligation-mediated PCR to recover integration site junctions

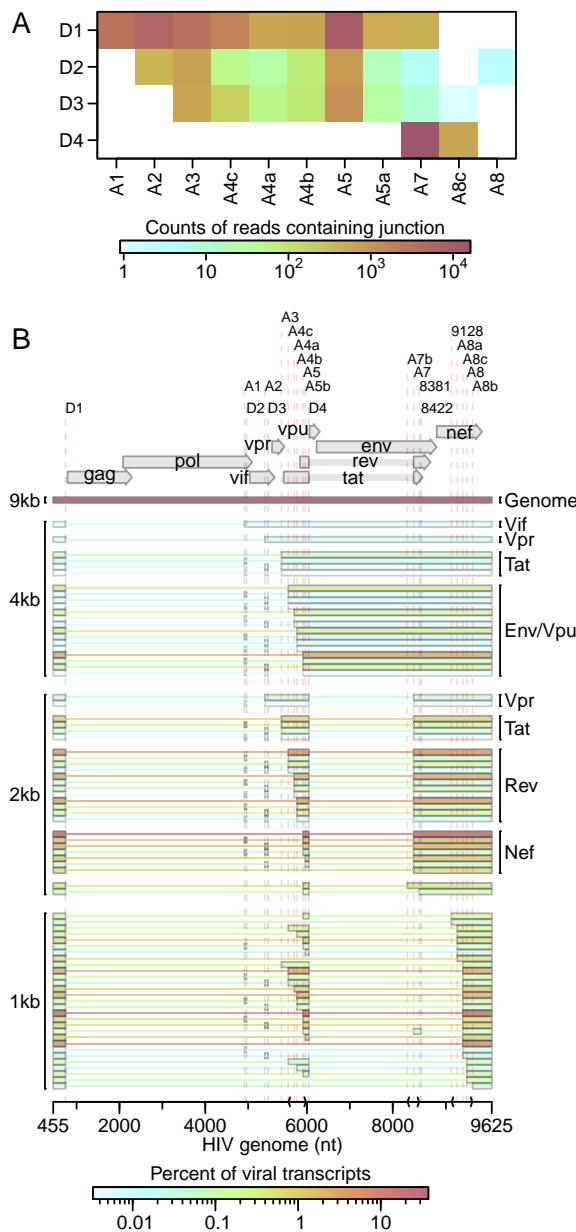


Figure 4.6: Transcription and splicing of the HIV_{89.6} RNA. A) Junctions between HIV splice donors and acceptors observed in the RNA-Seq data. Acceptors are shown as the columns and donors as the rows with the coloring indicating the frequency of each pairing. B) The relative abundance of all HIV_{89.6} transcripts as determined by a combination of PacBio sequencing³⁹⁴ and Illumina sequencing. Message structures were generated by targeted long read single molecule sequencing, which allowed association of multiple splice junctions in single sequence reads. The Illumina short read sequencing allowed normalization of message abundances between size classes. The inferred HIV message population is shown colored by relative abundance.

from the same infected cell populations analyzed by RNA seq, yielding 147,281 unique integration sites (Figure 4.7B)³³². No integrations in mitochondrial DNA were detected. We conclude that chimeric HIV-mitochondrial sequence reads in the RNA-seq data represent artifacts of library construction and so used these chimeras as an assay to evaluate subsequent data filtering steps. We reasoned that reads without sequence similarity at junctions between human and HIV mapping were less likely to be artifacts caused by template switching. Filtering to only reads where no overlap and no unknown intervening sequence was present between human and HIV portions left 2181 junctions and reduced the proportion of reads containing mitochondrial DNA to 2.4%. Of the remaining HIV-human chimeric reads, the HIV portion of 605 sequences bordered the 3' or 5' end of HIV or an HIV splice donor or acceptor. Filtering to these more likely authentic junctions left only 2 (0.3%) chimeric reads containing mitochondrial sequence. This decrease in likely mitochondrial artifacts suggests that the filtering was effective. The high rate of mitochondrial chimeras in the unfiltered sequences raises the concern that artifacts may easily distort results in studies using similar amplification and sequencing techniques.

Chimeric messages composed of HIV and cellular RNA sequences can be formed by cellular gene transcription reading into the integrated provirus, by HIV transcription reading out through the viral polyadenylation site or by splicing between human and viral splice sites. In our filtered data, the predominant forms appear to be derived from reading through the HIV polyadenylation signal into the surrounding DNA (78%), splicing out of the viral D4 splice donor to join to human slice acceptors (17%) and reading into the HIV 5' LTR from human sequence (4.0%) (Figure 4.7C). No splice site other than D4 had more than two chimeric reads observed.

The filtered chimeric reads had many traits consistent with biological chimera formation. The reads containing HIV D4 joined to human sequences had the characteristics expected of splicing—72.1% of the chimeric junctions mapped to known human acceptors

and 96.1% mapped to a location immediately preceded by the AG consensus of human mRNA acceptors. The reads containing the 5' or 3' LTR border were almost exclusively (93%) found in transcription units, with odds of being in a gene 2.8-fold (95% CI: 1.6–3.2×) higher than integration sites from the same sample. The 5' or 3' chimeras were also more likely to be located in an exon than integration sites even after excluding any integration or chimera not located in a transcription unit (odds ratio: 2.1×, 95% CI: 1.6–2.6×).

We next compared whether the human and viral segments of chimeric reads agreed or disagreed in orientation (i.e. strand transcribed) for reads with the human portion mapped within annotated transcription units. The sequencing technique used here does not preserve strand information, but we can check whether the strand of a sequence read agrees or disagrees with the annotated gene strand and compare this to the observed strand of the HIV portion of the read. We found a strong association between the orientation of the human and HIV portions of chimeric reads within 3' and 5' chimeras (odds ratio: 6.2×, 95% CI: 3.9–10.2×). This highly significant enrichment of HIV and human genes in the same orientation (Fisher's exact test $p < 10^{-15}$) might indicate that antisense HIV RNA is rapidly degraded by a response to double-stranded RNA or that polymerases oriented in opposing directions interfere with one another during elongation. Chimeras involving HIV splice donor D4 were even more highly enriched for matching orientations (odds ratio: 52.5×, 95% CI: 12.1–307×) suggesting that pairing with human splice acceptors may add an additional constraint on the orientation of D4 chimeric reads.

Based on these data, we can propose a lower bound on the relative abundance of chimeras. If we assume that our filtering removed nearly all artifacts so that we have few false positives, then our estimate should be lower than the true proportion of chimeras. In our data, only $\frac{604}{12,689,879} = 0.0048\%$ of reads containing sequence mapping to HIV also contained identifiable chimeric junctions. However, this is an underestimate

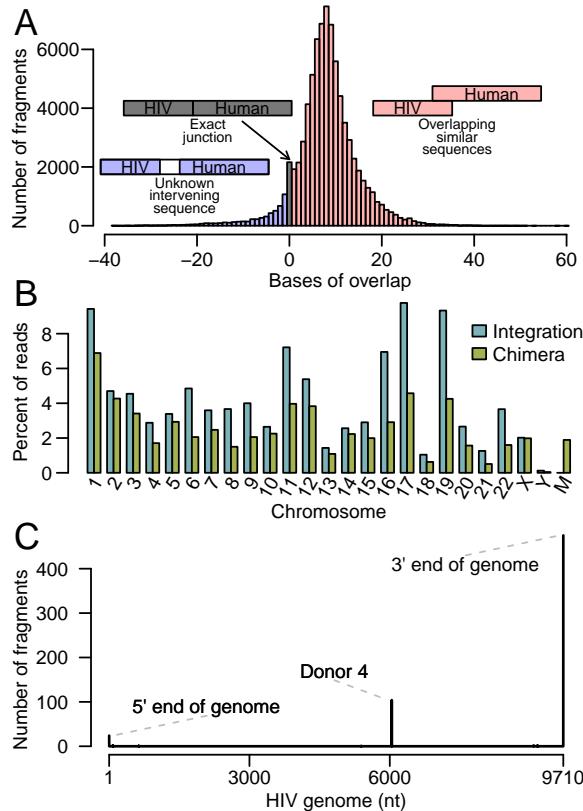


Figure 4.7: Analysis of chimeric RNA sequences containing both human and HIV sequences. A) The length of overlapping sequence (regions of complementarity potentially favoring chimera formation) matching both human and HIV at inferred chimeric junctions. The x-axis shows the length of the overlap and the y-axis shows the frequency of chimeric junctions with the indicated extent of overlap. B) Chromosomal distribution of uniquely mapping HIV integration sites from the same infections of primary T cells and comparison to uniquely mapping human sequences in chimeric reads observed in RNA-Seq. Note that the mitochondrial genome, denoted as M, has no authentic integration sites but does have extensive matches to chimeric junctions found in the RNA-Seq data. C) Counts of the location in the HIV genome of the HIV-human junctions in filtered chimeric reads.

because in an HIV-derived mRNA, any fragment of the sequence will be mappable to HIV, while for a chimeric sequence only a read spanning the HIV-human junction will allow identification of a chimera. If we assume that 25 bases of sequence are necessary to map to human or HIV sequence, then, with the 100-bp reads used here, only read fragments starting between 75- and 25-bp downstream of the chimeric junction will be identifiable. If we assume the average chimeric mRNA sequences is at least 2 kb long, then a read from a chimeric sequence has at most a $\frac{50}{2000} = 2.5\%$ chance of containing a mappable junction. Thus, a lower bound for the proportion of HIV mRNA that also contain human-derived sequences is $0.2\% (\frac{0.0048\%}{2.5\%})$. Looking only at splicing from HIV donor D4, we saw 16,843 reads containing a junction from D4 to an HIV acceptor and 104 reads from D4 to human sequence. Thus, in our data, 0.6% of D4 splice products form junctions with human acceptors instead of HIV acceptors.

4.5 Discussion

Here we used RNA-Seq to analyze mRNA accumulation and splicing in primary T cells infected with the low passage isolate HIV_{89.6}. We did not carry out dense time series analysis, compare different human cell donors or compare different perturbations of the infections—instead, we focused on generating a dense data set at a single time point. We analyzed replicate infected cell and control samples to allow discrimination of within-condition versus between-condition variation and assessed differences using a series of bioinformatic approaches. Many previous studies have used microarray technology or RNA-Seq to study gene activity in HIV-infected cells^{251,253,293,404–412}, usually analyzing infections of transformed cell lines or laboratory adapted strains of HIV-1. Here we present what is to our knowledge the deepest RNA-Seq data set reported for infection in primary T cells using a low passage HIV isolate (HIV_{89.6}). This data set was paired with a set of 147,281 unique integration site sequences extracted from the same infections, which were critical to our ability to quality control chimeric reads. An advantage of studies using cell lines and laboratory adapted strains is that often a high percent of cell infection can be achieved, whereas in this study we achieved only around ~30% infection. However, we report distinctive features of the transcriptional response not seen in studies of HIV infections in cell lines. Novel in this study are 1) identification of intron retention as a consequence of HIV infection, 2) the finding of activation of ERV-9/LTR12C after HIV infection, 3) generation of a quantitative account of the structures and abundances of over 70 HIV_{89.6} messages and 4) clarification of the predominant types of HIV-host transcriptional chimeras. These findings are discussed below.

Broad changes in host cell mRNA abundances were evident after infection, with over 17% of expressed genes changing significantly in activity. Changes included expected response to viral infection, apoptosis and T cell activation. Although it is not possible here to separate the response of infected and bystander cells, this study highlights the

drastic changes in cellular expression caused by HIV-1 infection. In a meta-analysis including four previously published studies, no gene was detected as differentially expressed in all five studies and only a handful of genes appeared in four out of five studies. Further analysis showed that expression changes appear to be cell type specific, raising concerns that studies using cell lines may not fully reflect host cell responses in *in vivo* infections.

Unexpectedly, intronic sequences were more common in the RNA-Seq data from cells after HIV_{89.6} infection than in mock infected cells. The mechanism is unclear. It is possible that the splicing machinery is reduced in activity after 48 hours of infection, perhaps as a part of the antiviral response of infected and bystander cells. HIV infection does appear to alter expression and localization of some splicing factors^{374,439}. In addition, we saw a large reduction in the abundance of mRNA from nonsense-mediated decay related genes, perhaps indicating that RNA surveillance is loosened thus allowing more unspliced or aberrantly spliced transcripts. Alternatively, fully spliced mRNAs might be more rapidly degraded after infection, possibly by interferon-mediated induction of RNaseL⁴⁷⁹. A speculative possibility is that HIV_{89.6} encodes a factor that alters cellular splicing or promotes mRNA degradation to optimize splicing and translation of viral messages.

Infection resulted in increased expression of specific cellular repeated sequences. HERVs, in particular HERV-K, have previously been observed to show increased RNA accumulation with HIV infection^{440–442,480} and possibly represent vaccine targets because of their production of distinctive proteins^{446,480–484}. Here, though we saw modest increases in HERV-K expression, ERV-9 had the greatest change in expression (33 LTR12C and 14 ERV-9 annotated regions with greater than 4× change in expression). Previous RNA-Seq studies of HIV infection in cell lines did not report increases in HERV expression^{253,293} but this difference is likely due to a much higher baseline expression of HERVs in transformed cell lines. We also observed increases in LINE

and Alu element transcription, as has been reported previously⁴⁴³, and expression changes in ERV-9/LTR12C expression associated with transcription factor motifs and U3 variants.

Many of the repeated sequence elements that were induced by HIV_{89.6} infection are relatively recently integrated in the human genome. The reason for this pattern is unclear. It may be that older elements have accumulated more mutations, resulting in an inactivation of transcriptional signals. Alternatively, perhaps the elements that are induced have been recruited for transcriptional control of cellular functions, so that their transcriptional activity is preserved evolutionarily^{423,485,486}.

Comparison of results of sequencing HIV_{89.6} messages using long-read single molecule sequencing (Pacific Biosciences) and dense short read sequencing (Illumina data reported here) allowed a full quantitative accounting of more than 70 HIV_{89.6} splice forms. The full length unspliced HIV RNA comprised 37.6% of all messages, corresponding to about 2000 genomes per cell. Notably abundant messages included those encoding Nef (D1-A5-D4-A7: 15.5%) and two Rev-encoding transcripts (D1-A4c-D4-A7: 4.2%, D1-A4b-D4-A7: 3.1%). The full set of messages is summarized in Figure 4.6B. Our previous analysis revealed an unusually prominent 1 kb size class. HIV_{89.6} encodes a rare splice acceptor (A8c) within Nef responsible for formation of the short messages. Our data indicated that two members of the 1-kb size class, D1-A5-D4-A8c and D1-A8c, accounted for 10.6% and 4.9% of all messages. The 1 kb size class as a whole accounted for fully 20% of messages. Most HIV/SIV variants appear to encode an acceptor near this position, suggesting a potential unknown function for these short spliced forms^{358,362,394}.

After filtering, we detected a sizeable number of apparently authentic chimeras containing both HIV and cellular sequences, allowing comparison to examples of host-cell modification by integration. Mechanisms of insertional activation have been studied intensively in animal models of transformation and in adverse events in human gene therapy. One of the most common mechanisms involves insertion of a retroviral en-

hancer near a cellular promoter, so that the rate of initiation is increased and normal cellular messages are increased in abundance. However, another common mechanism involves formation of chimeric messages involving both cellular and viral/vector sequences. In HIV infection, examples of insertion in the Bach2 and MKL2 genes have been associated with long term persistence of particular cell clones^{455–458}. In these cells, proviruses were integrated within the cellular transcription unit, and the transcriptional direction of the integrated provirus was the same as that of Bach2 or MKL2. This would allow formation of a fusion of the 5' HIV sequences with 3' Bach2 sequences, potentially involving the most common events seen here (either 3' read out or splicing from HIV D4 to a cellular exon). However, a closely studied example of clonal expansion in a successful lentiviral vector gene therapy for beta-thalassemia was associated with expansion of a cell clone harboring an integrated vector within the transcription unit of HMGA2. In this case the message spliced into the vector and terminated, removing a negative regulatory sequence normally present in the 3' end HMGA2 message⁴⁶². A targeted study in vitro of chimeric message formation by lentiviral vectors showed examples of multiple types of read-in and -out and splice-in and -out⁴⁶⁴, which may have been more frequent and more varied than for HIV_{89.6} proviruses studied here. The lack of splicing or reading into HIV in this study may be a reflection of the high rate HIV transcription in these infected cells—because HIV was so highly expressed, there would be more opportunities for polymerase to splice out of or read through the HIV genome than to read or splice in. The vast majority of HIV proviruses in expanded clones in well-suppressed patients now appear to be defective⁴⁵⁸—going forward, it will be of interest to investigate whether these HIV proviruses are damaged in ways that promote formation of chimeric transcripts.

Lastly, we note that several features of the transcriptional response to HIV_{89.6} infection were suggestive of de-differentiation away from T cell specific expression patterns. The increase in expression of cellular HERVs and LINEs is characteristic of cells in early development. Specific HERVs and transposons, including ERV-9/LTR12C and HERV-

K, have been implicated in regulating gene activity early in development^{423,485,487–490}. Several genes related to other hematopoietic cell types showed elevated RNA abundance after HIV_{89.6} infection. These data are of interest given the finding that patients undergoing long term ART can contain long lived T cell clones that may contribute to the latent reservoir^{458,491–494}. Possibly the transcriptional responses seen in infected primary T cells here are reflective of processes leading to formation of the long-lived latently-infected cells with stem-like properties.

4.6 Conclusions

Infections of primary T cells with a low passage HIV isolate show several distinctive features compared with previously published data using T cell lines and/or lab-adapted HIV strains. We found strong changes in expression in genes related to immune response and apoptosis similar to studies of HIV infection in patient samples and primary cells but different from studies performed in SupT1 cell lines. Notable changes after infection included intron retention and activation of recently integrated retrotransposons and endogenous retroviruses, in particular LTR12C/ERV-9. We also present complete absolute estimation of over 70 messages from HIV_{89.6} and specify the major virus-host chimeras as read out from the 3' end of the provirus and splicing from viral splice donor 4 to cellular acceptors.

4.7 Availability of supporting data

RNA-Seq reads from this study are available at the Sequence Read Archive under accession number SRP055981. The integration site data is available at the Sequence Read Archive under accession number SRP057555.

4.8 Author's contributions

KEO performed the infections and sequencing. SS-M analyzed the data. SS-M, KEO and FDB planned the overall study, and SS-M and FDB wrote the paper. All authors

read and approved the final manuscript.

4.9 Acknowledgements

We would like to thank the University of Pennsylvania Center for AIDS Research (P30 AI045008) for preparation of viral stocks and isolation of primary CD4⁺ T cells; Ronald G. Collman and members of the Bushman laboratory for reagents, helpful discussion and technical expertise. This work was funded by NIH grant R01 AI052845, the HIV Immune Networks Team (HINT) consortium P01 AI090935 and NRSA computational genomics training grant T32 HG000046.

4.10 Additional Files

4.10.1 Additional file 1 — Analysis of genes differentially expressed during HIV_{89.6} infection of primary CD4⁺ T cells

Output from CuffLinks analysis of the RNA-Seq data organized in a csv file with columns UCSC gene ID, gene symbol, status of test, FPKM in uninfected and infected samples, the log₂ fold change, test statistic and false discovery rate adjusted *p*-value.

4.10.2 Additional file 2 — Analysis of Gene Ontology categories associated with differential expression during HIV_{89.6} infection of primary CD4⁺ T cells

Counts of differentially expressed genes for each Gene Ontology category. Columns are the name of the category, the numbers of genes differentially up- and downregulated or not significantly changed and odds ratios and *p*-values from Fisher's exact tests.

4.10.3 Additional file 4 — Genes called as up- or downregulated by studies of expression during HIV infection

Genes called as differentially expressed in the five studies analyzed in the meta-analysis of differential expression with HIV infection. Columns are the study, the gene name(s) and whether the differential expression was up or down.

4.10.4 Additional file 5 — Estimating relative abundance of HIV_{89.6} message size classes using RNA-Seq data

A) RNA-Seq coverage of the HIV_{89.6} genome for the replicates in this study. Each replicate is indicated by a different color. The HIV genome is shown on the x-axis and the number of reads that aligned to each position is shown on the y-axis. Black line indicates the 0.021% coverage decrease per base distance from the 3' end of the mRNA estimated from a least squares fit on the read counts in the first intron. B)

Diagram of the segments of the HIV_{89.6} RNA present in each of 9 kb, 4 kb, 2 kb and 1 kb size class. C) The proportion of reads mapped to each of the segments of the HIV_{89.6} genome shown in B adjusted by the length of the segment. Each replicate is shown by a different color. D) Corrected representation of RNA segments from the different size classes. Because cDNA synthesis was primed from the polyA tail, more 3' sequences are recovered preferentially. Using the bias estimate from A, we adjusted each genome segment by the inverse of the bias predicted based on its distance from the 3' end of the mRNA. Corrected proportions for the indicated RNA segments are shown colored by replicate. E) The proportion of each size class was inferred using the estimates in D by calculating the difference between segments. Replicates are indicated by color.

CHAPTER 5 : A reverse transcription loop-mediated isothermal amplification assay optimized to detect multiple HIV subtypes

5.1 Abstract

Diagnostic methods for detecting and quantifying HIV RNA have been improving, but efficient methods for point-of-care analysis are still needed, particularly for applications in resource-limited settings. Detection based on reverse-transcription loop-mediated isothermal amplification (RT-LAMP) is particularly useful for this, because when combined with fluorescence-based DNA detection, RT-LAMP can be implemented with minimal equipment and expense. Assays have been developed to detect HIV RNA with RT-LAMP, but existing methods detect only a limited subset of HIV subtypes. Here we report a bioinformatic study to develop optimized primers, followed by empirical testing of 44 new primer designs. One primer set (ACeIN-26), targeting the HIV integrase coding region, consistently detected subtypes A, B, C, D, and G. The assay was sensitive to at least 5000 copies per reaction for subtypes A, B, C, D, and G, with Z-factors of above 0.69 (detection of the minor subtype F was found to be unreliable). There are already rapid and efficient assays available for detecting HIV infection in a binary yes/no format, but the rapid RT-LAMP assay described here has additional uses, including 1) tracking response to medication by comparing longitudinal values for a subject, 2) detecting of infection in neonates unimpeded by the presence of maternal antibody, and 3) detecting infection prior to seroconversion.

5.2 Introduction

Despite the introduction of efficient antiretroviral therapy, HIV infection and AIDS continue to cause a worldwide health crisis⁴⁹⁵. Methods for detecting HIV infection have improved greatly with time⁴⁹⁶—today rapid assays are available that can detect HIV infection in a yes-no format using a home test kit that detects antibodies in saliva.

Viral load assays that quantify viral RNA with quick turn-around time are widely available in the developed world. However, quantitative viral load assays are not commonly available with actionable time scales in much of the developing world. This motivates the development of new rapid and quantitative assays that can be used at the point of care with minimal infrastructure^{497,498}.

One simple and quantitative detection method involves reverse transcription-based loop mediated isothermal amplification (RT-LAMP)⁴⁹⁹. In this method, a DNA copy of the viral RNA is generated by reverse transcriptase, and then isothermal amplification is carried out to increase the amount of total DNA. Primer binding sites are chosen so that a series of strand displacement steps allow continuous synthesis of DNA without requiring thermocycling. Reaction products can be detected by adding an intercalating dye to reaction mixtures that fluoresces only when bound to DNA, allowing quantification of product formation by measurement of fluorescence intensity. Such assays can potentially be packaged in simple self-contained devices and read out with no technology beyond a cell phone.

RT-LAMP assays for HIV-1 have been developed previously and reported to show high sensitivity and specificity for subtype B, the most common HIV strain in the developed world^{498,500,501}. Another recent study reported RT-LAMP primer set optimized for the detection of HIV variants circulating in China⁵⁰², and another on confirmatory RT-LAMP for group M viruses⁵⁰³. Assays have also been developed for HIV-2⁵⁰⁴. A complication arises in using available RT-LAMP assays due to the variation of HIV genomic sequences among the HIV subtypes^{505,506}, so that an RT-LAMP assay optimized for one viral subtype may not detect viral RNA of another subtype⁵⁰⁷. Tests presented below show that many RT-LAMP assays are efficient for detecting subtype B, for which they were designed, but often performed poorly on other subtypes. Subtype C infects the greatest number of people worldwide, including in Sub-Saharan Africa, where such RT-LAMP assays would be most valuable, motivating optimization for subtype

C. Several additional non-B subtypes are also responsible for significant burdens of disease world-wide⁵⁰⁸.

Here we present the development of an RT-LAMP assay capable of detecting HIV-1 subtypes A, B, C, D, and G. We first carried out a bioinformatic analysis to identify regions conserved in all the HIV subtypes. We then tested 44 different combinations of RT-LAMP primers targeting this region in over 700 individual assays, allowing identification of a primer set (ACeIN-26) that was suitable for detecting these subtypes. We propose that the optimized RT-LAMP assay may be useful for quantifying HIV RNA copy numbers in point-of-care applications in the developing world, where multiple different subtypes may be encountered.

5.3 Methods

5.3.1 Viral strains used in this study

Viral strains tested included HIV-1 92/UG/029 (Uganda) (subtype A, NIH AIDS Reagent program reagent number 1650), HIV-1 THRO (subtype B, plasmid derived, University of Pennsylvania CFAR)⁵⁰⁹, CH269 (subtype C, plasmid derived, University of Pennsylvania CFAR)⁵⁰⁹, UG0242 (subtype D, University of Pennsylvania CFAR), 93BRO20 (subtype F, University of Pennsylvania CFAR), HIV-1 G3 (subtype G, NIH AIDS Reagent program reagent number 3187)⁵¹⁰.

Viral stocks were prepared by transfection and infection. Culture supernatants were cleared of cellular debris by centrifugation at 1500g for 10 min. The supernatant containing virus was then treated with 100 U DNase (Roche) per 450 uL virus for 15 min at 30°C. RNA was isolated using the QiaAmp Viral RNA mini kit (Qiagen GmbH, Hilden, Germany). RNA was eluted in 80 uL of the provided elution buffer and stored at -80°C.

Concentration of viral RNA copies was calculated from p24 capsid antigen capture

assay results provided by the University of Pennsylvania CFAR or the NIH AIDS-reagent program. In calculating viral RNA copy numbers, we assumed that all p24 was incorporated in virions, all RNA was recovered completely from stocks, 2 genomes were present per virion, 2000 p24 molecules per viral particle, and the molecular weight of HIV-1 p24 was 25.6 kDa.

5.3.2 Assays

RT-LAMP reaction mixtures (15 μ L) contained 0.2 μ M each of primers F3 and B3 (if a primer set used multiple B3 primers, mixture contained 0.2 μ M of each); 1.6 μ M each of FIP and BIP primers (if a primer set had multiple FIP primes, reaction mixture contained 0.8 μ M of each FIP primer); and 0.8 μ M each of LoopF and LoopB primers; 7.5 μ L OptiGene Isothermal Mastermix ISO-100nd (Optigene, UK), ROX reference dye (0.15 μ L from a 50X stock), EvaGreen dye (0.4 μ L from a 20X stock; Biotium, Hayward, CA); HIV RNA in 4.7 μ L; AMV reverse transcriptase (10U/ μ L) 0.1 μ L and water to 15 μ L In most cases where two primer sets were combined, the total primer concentration within the reaction was doubled such that the above individual primer molarities were maintained. For the mixture ACeIN-26+F-IN (S2 Table, line 46), the total primer concentration was not doubled—the F-IN primer set comprised 25% of the total primer concentration, and the ACeIN-26 primer set comprised 75% of the total primer concentration with the ratios of primers listed above preserved. This mixture was combined 1:1 with the ACe-PR primer set (S2 Table, line 47) such that total primer concentration in the final mixture was doubled.

Amplification was measured using the 7500-Fast Real Time PCR system from Applied Biosystems with the following settings: 1 minute at 62°C; 60 cycles of 30 seconds at 62°C and 30 seconds at 63°C. Data was collected every minute. Product structure was assessed using dissociation curves which showed denaturation at 83°C. Products from selected amplification reactions were analyzed by agarose gel electrophoresis and showed a ladder of low molecular weight products (data not shown).

Product synthesis was quantified as the cycle of threshold for 10% amplification. Z-factors⁵¹¹ were calculated from tests of 24 replicates using the ACeIN26 primer set in assays with viral RNA of each subtype. No detection after 60 min was given a value of 61 min in the Z-factor calculation.

5.4 Results

5.4.1 Testing published RT-LAMP primer sets against multiple HIV subtypes

We first assessed the performance of existing RT-LAMP assays on RNA samples from multiple HIV subtypes. We obtained viral stocks from HIV subtypes A, B, C, D, F, and G, estimated the numbers of virions per ml, and extracted RNA. RNAs were mixed with RT-LAMP reagents which included the six RT-LAMP primers, designated F3, B3, FIP, BIP, LF and LB⁴⁹⁹. Reactions also contained reverse transcriptase, DNA polymerase, nucleotides and the intercalating fluorescent EvaGreen dye, which yields a fluorescent signal upon DNA binding. DNA synthesis was quantified as the increase in fluorescence intensity over time, which yielded a typical curve describing exponential growth with saturation (examples are shown below). Results are expressed as threshold times (T_t) for achieving 10% of maximum fluorescence intensity at the HIV RNA template copy number tested.

In initial tests, published primer sets targeting the HIV-1 subtype B coding regions for capsid (CA), protease (PR), and reverse transcriptase (RT) (named B-CA, B-PR and B-RT) were assayed in reactions with RNAs from four of the subtypes. Results with each primer set tested are shown in Figure 5.1 in heat map format, where each tile summarizes the results of tests of 5000 RNA copies. Primers and their groupings into sets are summarized in S1 and S2 Tables, average assay results are in S3 Table, and raw assay data is in S4 Table. Assays (Figure 5.1, top) with the B-CA, B-PR and B-RT primer sets detected subtypes B and D at 5000 RNA copies with threshold times less than 20 min. However, assays with B-CA and B-RT detected subtypes C and F with

threshold times > 50 min, indicating inefficient amplification and the potential for poor separation between signal and noise. B-PR did not detect subtype C at all. In an effort to improve the breadth of detection, we first tried mixing the B-PR primers, which detected clade F (albeit with limited efficiency) with the B-CA and B-RT primers (Figure 5.1 and S3 and S4 Tables). In neither case did this provide coverage of all four clades tested. We thus did not test these primer sets on RNAs from the remaining subtypes and instead sought to develop primer sets targeting different regions of the HIV genome.

5.4.2 Primer design strategy

To design primers that detected multiple HIV subtypes efficiently, we analyzed alignments of HIV genomes (downloaded from the Los Alamos National Laboratory site⁵⁰⁵) for regions with similarity across most viruses, revealing that a segment of the pol gene encoding IN was particularly conserved (Figure 5.2A). A total of six primers are required for each RT-LAMP assay⁴⁹⁹. We used the EIKEN primer design tool to identify an initial primer set targeting this region. In further analysis, positions in the alignments were identified within primer landing sites that commonly contained multiple different bases. Primer positions were manually adjusted to avoid these bases when possible, and when necessary mixtures were formulated containing each of these commonly occurring bases (S1 and S2 Tables). An extensive series of variants targeting the IN coding region was tested empirically in assays containing RNAs from multiple subtypes (5000 RNA copies per reaction, over 700 total assays; S3 and S4 Tables). Based on initial results, primers were further modified by adjusting the primer position or addition of locked nucleic acids as described below.

5.5 Testing different primer designs

Our first design, ACeIN-1 (“ACe” for “All Clade” and “IN” for “integrase”), targeted the HIV IN coding region and contained multiple bases at selected sites to broaden

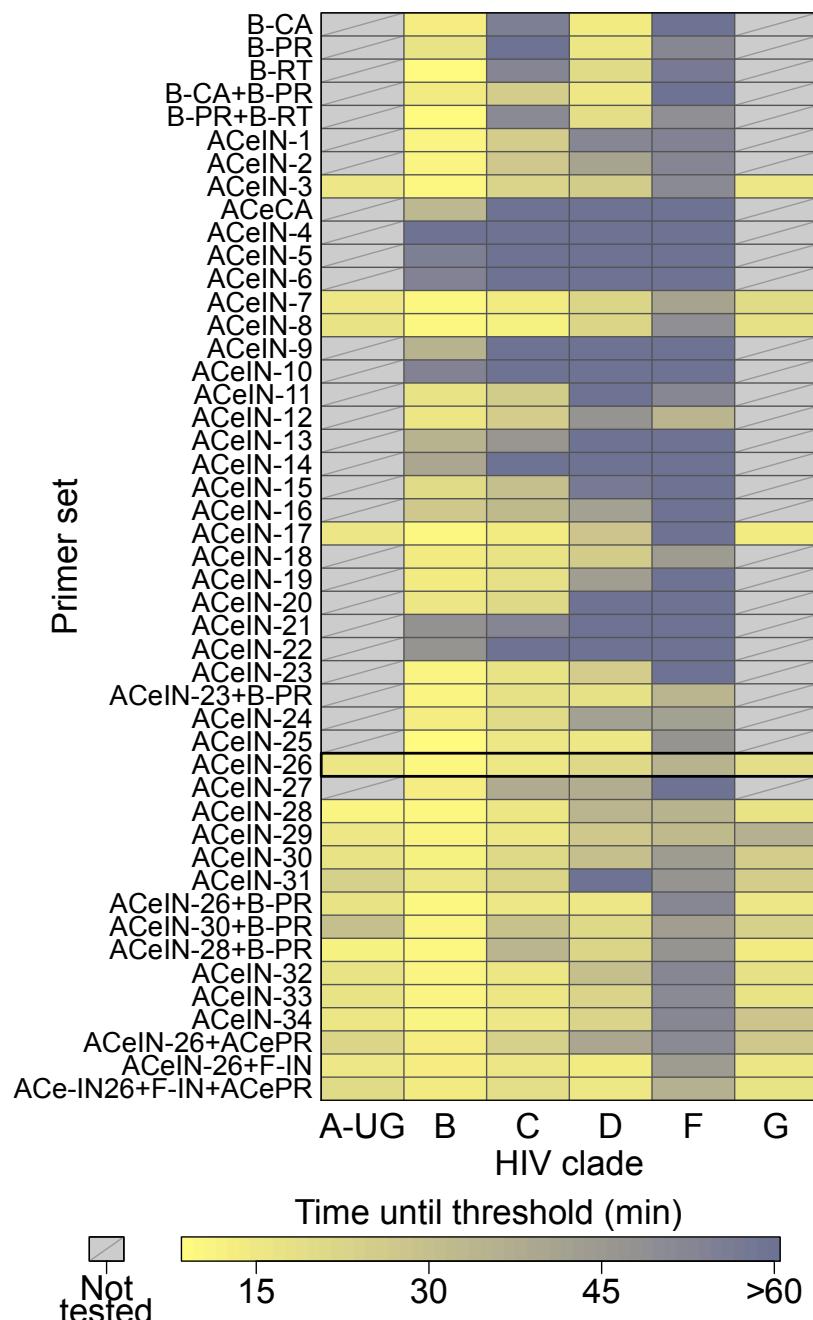


Figure 5.1: Summary of amplification results for all the RT-LAMP primer sets tested in this study. The data is shown as a heat map, with more intense yellow coloring indicating shorter amplification times (key at bottom). Primer sets tested are named along the left of the figure. Primer sequences, and their organization into LAMP primer sets, are cataloged in S1 and S2 Tables. The raw data and averaged data are collected in S3 and S4 Tables. ACeIN-26 primer set (highlighted) had one of the best performances across the subtypes and a relatively simple primer design.

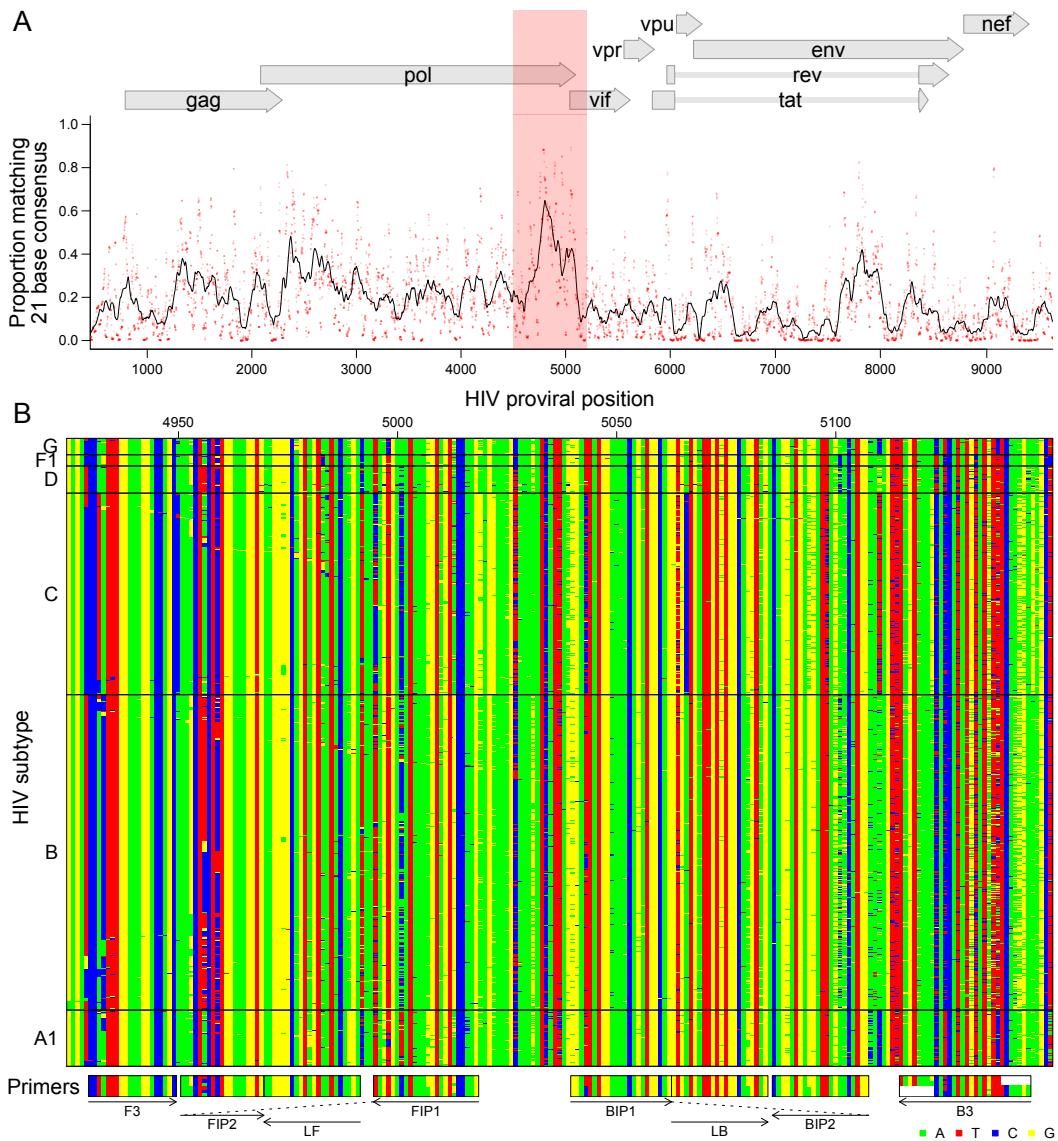


Figure 5.2: Bioinformatic analysis to design subtype-agnostic RT-LAMP primers. A) Conservation of sequence in HIV. HIV genomes ($n = 1340$) from the Los Alamos National Laboratory collection were aligned and conservation calculated. The x-axis shows the coordinate on the HIV genome, the y-axis shows the proportion of sequences matching the consensus for each 21 base segment of the genome (red points). The black line shows a 101 base sliding average over these proportions. The vertical red shading shows the region targeted for LAMP primer design that was used as input into the EIKEN primer design tool. Numbering is relative to the HIV_{89.6} sequence. B) Aligned genomes, showing the locations of the ACeIN-26 primers. Sequences are shown with DNA bases color-coded as shown at the lower right. Each row indicates an HIV sequence and each column a base in that sequence. Horizontal lines separate the HIV subtypes (labeled at left). Arrows indicate the strand targeted by each primer. Primers targeting the negative strand of the virus are shown as reverse compliments for ease of viewing.

detection (Figure 5.1). ACeIN-2 and -3 have primers (B3) with slightly different landing sites. Tests showed that the mixture of primers allowed amplification with a shorter threshold time than did either alone (Figure 5.1).

We also tried to design a new primer set to the CA coding region (Figure 5.1, ACeCA) but found that the set only amplified clade B, and not efficiently. Thus this design was abandoned.

ACeIN-3 through-6 were altered by inserting a polyT sequence between the two different sections of FIP and BIP in various combinations, a modification introduced with the goal of improving primer folding, but these designs performed quite poorly (Figure 5.1).

Because the FIP primer appeared to bind the region with most variability among clades, we tried variations that bound to several nearby regions. These were tried with and without the polyT containing BIP and FIP primers in various combinations (Figure 5.1, ACeIN-7 through-22). We also tried mixing all of the variations of FIP together (ACeIN-23; S2 Table). The ACeIN-23 primer set was tried as a mixture with the B-PR set to try to capture clade F, yielding a relatively effective primer set (Figure 5.1, ACeIN-23+B-PR).

In an effort to increase affinity, an additional G/C pair was added to F3 and tested with various other IN primers (Figure 5.1, ACeIN-24 through-31). Testing showed improvement, with ACeIN-26 showing particularly robust amplification.

In a second effort to increase primer affinities, we substituted locked nucleic acids (LNAs) for selected bases that were particularly highly conserved among subtypes (Figure 5.1, ACeIN-30, -31, -32, -33, and -34). Some improvement was shown over the non-LNA containing bases. However, the ACeIN-26 primer set was as effective as or better than any LNA containing primer sets.

In further tests, the ACeIN-26, -28 and -30 primers were tested combined with the

ACePR primer set (a slightly modified version of the B-PR primer set, S2 Table, row 2, designed to accommodate a wider selection of HIV-1 subtypes) but no improvement was seen and efficiency may even have fallen for some subtypes. We also designed a primer set that matched exactly to the targeted sequences found in the problematic subtype F, and mixed this set with the ACeIN-26 primers. However, no improvement was seen (Figure 5.1, mixtures with F-IN set). Mixing the ACeIN-26 primers with both the ACePR and F-specific primers did yield effective primer sets (Figure 5.1, ACeIN26+F-IN and ACeIN26+F-IN+ACePR). However, amplification efficiency was not greatly improved over the ACeIN-26 primer set, so we proceeded with the simpler ACeIN-26 primer set (Figure 5.2B) in further studies.

5.5.1 Performance of the optimized RT-LAMP assay

The ACeIN-26 RT-LAMP primer set was next tested to determine the minimum concentration of RNA detectable under the reaction conditions studied (Figure 5.3). RNA template amounts were titrated and time to detection quantified. Tests showed detection after less than 20 min of incubation for 50 copies of subtypes A or B, detection after less than 30 min for 5000 copies for C, D, and G, and detection after less than 20 min for 50,000 copies for F.

For clinical implementation the reliability of an assay is critical. This is commonly summarized as a Z-factor⁵¹¹, which takes into account both the separation in means between positive and negative samples and the variance in measurement of each. An assay with a Z-factor above 0.5 is judged to be an excellent assay. Z-factors for detection of each of the subtypes at 5000 RNA copies per reaction were > 0.50 for subtypes A, B, C, D, and G, respectively (Figure 5.4, n = 24 replicates per test). Detection of subtype F at 5000 copies per reaction was sporadic, showing a much lower Z-factor. Therefore our ACeIN-26 RT-LAMP primer set appears well suited to detect 5000 copies of subtypes A, B, C, D and G.

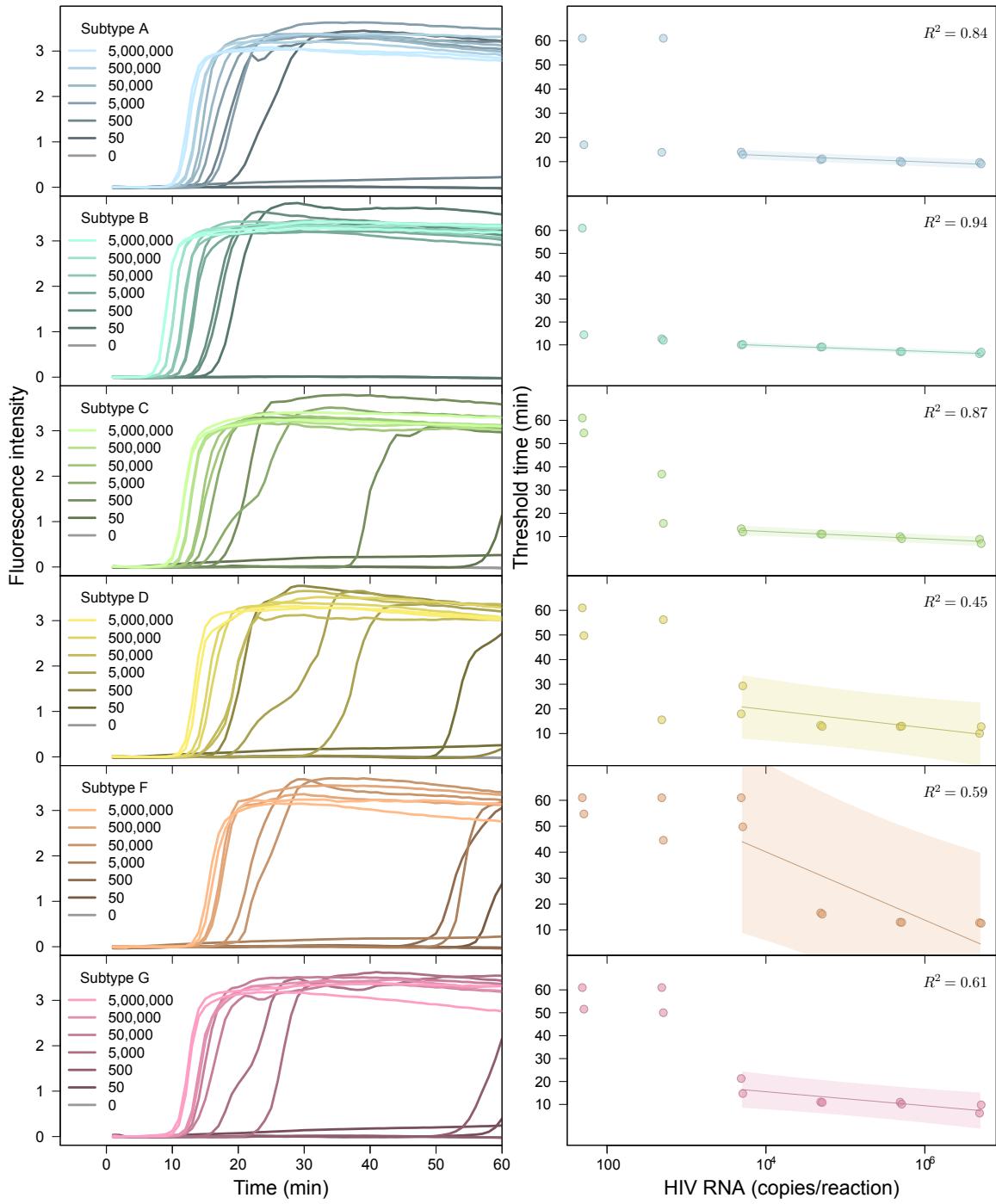


Figure 5.3: Performance of the AceIN-26 primer set with different starting RNA concentrations. Tests of each subtype are shown as rows. In each lettered panel, the left shows the raw accumulation of fluorescence signal (y-axis) as a function of time (x-axis); the right panel shows the threshold time (y-axis) as a function of log RNA copy number (x-axis) added to the reaction. In the right hand panels, values were dithered where two points overlapped to allow visualization of both.

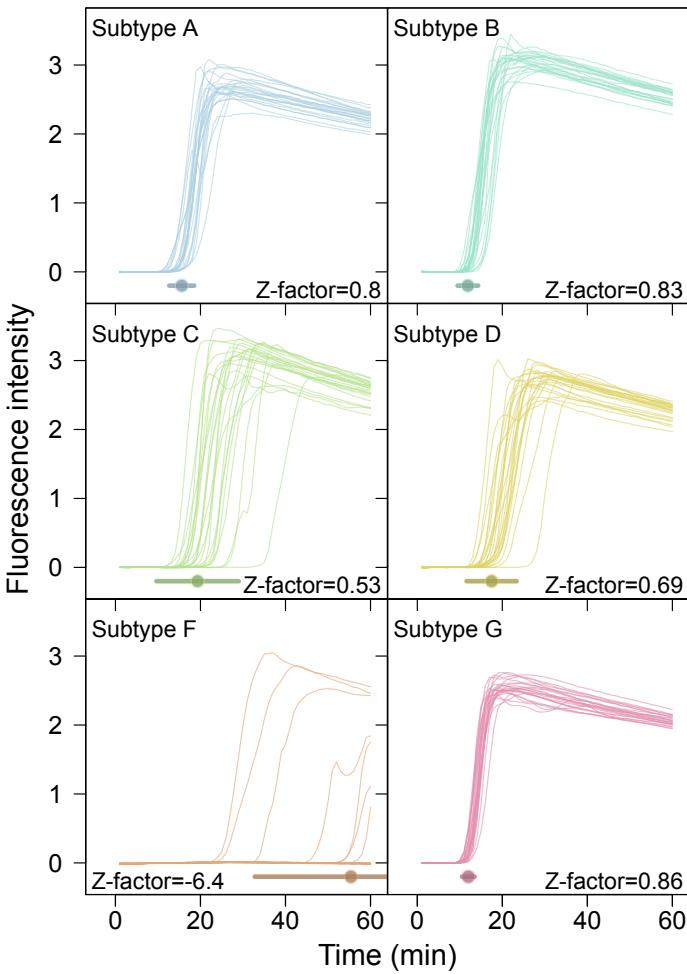


Figure 5.4: Examples of time course assays, displaying replicate tests of RT-LAMP primer set ACeIN-26 tested over six HIV subtypes, used in Z-factor calculations. A total of 5000 RNA copies were tested in each $15 \mu\text{L}$ reaction. Time is shown on the x-axis, Fluorescence intensity on the y-axis. Replicates are distinguished using an arbitrary color code. Z-factor values and standard deviations are shown on each panel.

5.6 Discussion

Here we present an RT-LAMP assay optimized to identify multiple HIV subtypes. Infections with subtype B predominate in most parts of the developed world, but elsewhere other subtypes are more common⁵⁰⁸. Thus nucleic acid-based assays for use in the developing world need to query HIV subtypes more broadly. Previously reported RT-LAMP assays, while effective at detecting subtype B, commonly showed poor ability to detect at least some of the HIV subtypes, including C, which is common in the developing world (Figure 5.1). Here we first carried out an initial bioinformatic survey to identify regions conserved across all HIV subtypes that could serve as binding sites for RT-LAMP primers. We then tested primer sets targeting these regions empirically for efficiency. Testing 44 different primer sets revealed that assays containing ACeIN-26 were effective in detecting 5000 copies of RNA from subtypes A, B, C, D, and G within 30 minutes of incubation. For these five subtypes, the times of incubation to reach the threshold times were not too different, which simplifies interpretation when the subtype in the sample is unknown. Regardless of the efficiency, these assays can be applied to longitudinal studies of changes in viral load within an individual. We propose that RT-LAMP assays based on the ACeIN-26 primer set can be useful world-wide for assaying HIV-1 viral loads in infected patients.

There are several limitations to our study. Subtypes A, B, C, D, and G were detected efficiently and showed Z-factors above 0.5, but subtype F was detected reliably only with higher template amounts, probably due to more extensive mismatches with the ACeIN-26 primer set. Subtype F is estimated, however, to comprise only 0.59% of all infections globally⁵⁰⁸, though it is common in some regions. For many of the common circulating recombinant forms, such as AE and BC, the target site for ACeIN-26 is from a subtype known to be efficiently detected, though in some cases the efficiency of detection is not easy to predict and will need to be tested. We did not test subtypes beyond A, B, C, D, F and G, and we did not attempt to assess multiple different variants

within each subtype. Thus, while we do know that our RT-LAMP assays are more widely applicable than many of those reported previously, we do not know whether they are able to detect all strains efficiently. In addition, although we carried out more than 700 assays in this study, there remain multiple parameters that could be optimized further, such as primer concentrations, salt type and concentration, temperature, and divalent metal concentrations, so there are likely further opportunities for improvement. Also, possible effects of RNA quality on assay performance were not tested rigorously.

A particularly important parameter for further optimization is primer sequence. Several groups have recently published primer sets optimized for broad detection of different HIV lineages^{502,503}, offering opportunities for creating sophisticated primer blends with increased breadth of detection. However, in developing such mixtures, it will be important to monitor for possible complicating interactions of primers with each other. As an example of ongoing development of mixtures, we found that addition of another primer to the ACeIN-26 set that was matched to a common subtype C lineage allowed improved detection of subtype C variants (S1 Report). In order to improve detection of subtype F, which was suboptimal with ACeIN-26, additional primer sets could be mixed to specifically target subtype F, though the ones we tried so far did not work well. It will be useful to explore the performance of broader primer mixtures in future work.

Today rapid assays are available that can report infection efficiently, for example by detecting anti-HIV antibodies in oral samples—however, the nucleic acid-based method presented here has additional potential uses. We envision combining the RT-LAMP assay with simple point-of-care devices for purifying blood plasma⁴⁹⁷ and quantitative analysis of accumulation of fluorescent signals⁵¹². In one implementation of the technology, cell phones could be used to capture and analyze results, thereby minimizing equipment costs. Point-of-care devices are available facilitating the concentration of viral RNA from blood plasma or saliva⁵¹² to allow the detection of the 1000 RNA copy threshold that the WHO defines as virological treatment failure (World Health

Organization, Consolidated ARV guidelines, June 2013). Together, these methods will allow assessment of parameters beyond just the presence/absence of infection. Quantitative RT-LAMP assays should allow tracking of responses to medication, detection in neonates (where immunological tests are confounded by presence of maternal antibody), and early detection before seroconversion.

CHAPTER 6 : Conclusions and future directions

In this dissertation, we described studies to characterize the nature of HIV-1 latency, characterize expression and alternative splicing and host cell response to infection and develop alternative methods for detection of infection and quantification of viral loads.

A common theme was that cell lines and *in vitro* models of these replication steps often disagree with each other and with primary cell data.

6.1 Latency and integration location

In Chapter 2, we showed that the chromosomal location of integration affects proviral latency but the mechanisms appear to differ between cell culture models. Similarly a recent study of nine cell culture models found that no single model reliably predicted the performance of activating compounds in *ex vivo* tests of latently infected cells from HIV patients⁵¹³. This suggests that either some cell culture models do not accurately reflect latency in patients or that there are diverse subsets of cells with differing mechanisms of latency within patients.

Cell culture models are currently used to screen potentially therapeutic compounds^{130,513}. If some cell culture models are not representative of *in vivo* conditions then potential treatments may be discarded or marked for development erroneously. Further comparisons between additional cell culture models and additional replicates of existing models might allow discrimination between batch/lab effects and reveal patterns in model behavior. Comparison with cells extracted from patients or infected lab animals might offer a gold standard comparison although it is difficult to obtain large amounts of cells and difficult to distinguish defective provirus from latent provirus in such populations.

Various treatments are now being considered for the reactivation of latent provirus⁵¹³. To further understand the mechanisms of these treatments, it would be informative to

compare the features of latent provirus induced by a given treatment to latent viable provirus remaining uninduced. Repeated cell sorting and integration site sequencing might provide insight on mechanism. For example, one could first sort out cells with active provirus, then treat with the potential latency modulator and sort out cells with newly active provirus and then treat with a strong inducer or alternative stimuli and sort out cell with newly activated provirus. This would give subsets of cells where latent proviruses had been activated by treatment and cells with provirus which were not activated by treatment but still inducible. Synergies between treatments could be assessed and the location of integration sites could be determined and used to locate patterns of genomic features correlated with induction for each treatment.

Current efforts at “shock and kill” therapy focus on histone deacetylase inhibitors. If there are diverse mechanisms of latency within patients then much of the latent reservoir may remain unactivated by single target therapies. Clinical trials with histone deacetylase inhibitors have shown some small increases in viral RNA but little decrease in the latent reservoir of HIV^{318,514? ,515}. It appears that the majority of viable latent provirus from patient cells are not reactivated by current therapies¹²⁶. These results are particularly worrisome since 10,000-fold or more reductions of the latent reservoir are likely to be necessary to functionally cure HIV⁵¹⁶.

In Chapter 2, we used publicly available genomic data. Perhaps there is some chromosomal feature with a strong association with latency but the data is not currently available or varies greatly between cell populations. More varieties of annotations are rapidly becoming available⁵¹⁷⁻⁵²¹. Decreasing sequencing costs⁵²²⁻⁵²⁴ may also make it feasible to measure more epigenetic features in the exact cell population of interest. Repeating analysis similar to Chapter 2, perhaps by simply rerunning the reproducible report in Appendix A.2, with new data would allow any new features to be monitored for correlations with latency.

6.2 HIV-1 alternative splicing

In Chapter 3, we

Further clarification using more detailed sequencing in more time points, cell types and strains of HIV-1 and other lentiviruses rema

PacBio was bad. Figure? Do better

In addition an important subset of HIV are the founder viruses transmitted between hosts^{525,526}. These viruses are not well studied and perhaps their splicing and gene expression differ from the rest of the viral swarm of late-term patients.

6.3 Host expression during HIV infection

Non polyadenylated RNA. Strand specific sequencing. Longer reads and longer fragments.

Localization nucleus vs cytoplasm

Cell types, macrophages

Infection, sorting

Cell lines bad

Endogenous retrovirus

6.4 LAMP PCR and lab-on-a-chip

In Chapter 5, we report a loop mediated isothermal amplification system using primers optimized to detect most subtypes of HIV-1. An alternative to a single broadly targeted primer set would be to design separate primer sets targeted specifically to each subtype so that a positive amplification would then be able to discriminate viral subtype. Different viral subtypes can have different rates of disease progression^{527–530},

transmission dynamics^{531–533} and response to treatment^{534–536}. Simple low-cost devices with multiple reactions chambers could be used to both identify viral subtype and estimate viral load^{537,538} and allow modified treatment decisions.

A LAMP chip with subtype-specific primers would also allow the detection of some superinfections. Superinfection of a single individual with multiple distinct strains of HIV is common in high risk individuals^{493,539–542} and the general population⁵⁴³. Superinfection can lead to disease progression^{544–549} or drug resistance⁵⁵⁰. Superinfection also allows recombination between divergent strains^{539,545,546,548,551} and this rapid exchange of genetic information can lead to more fit recombinant strains and worsen the global epidemic^{57,62,546,552,553}. LAMP detection of superinfection could allow early intervention and suppression in superinfected individuals.

The techniques described in Chapter 5 also allow for rapid development of detection assays for novel pathogens. For example, in a recent outbreak in West Africa, Zaire ebolavirus has infected over 26,000 confirmed, probable and suspected cases and caused over 11,000 reported deaths^{554–556}. Early detection and quarantine are essential to the control of this epidemic⁵⁵⁷. Amplification of Ebola virus nucleic acid through polymerase chain reaction is the best diagnostic test currently available but the necessary resources are often not available in these resource-poor regions^{558,559}. Antigen-based tests are quicker and available at the point-of-care but are not as accurate or sensitive as polymerase chain reaction tests and are still in limited supply⁵⁵⁹. Loop-mediated isothermal amplification offers the potential for rapid, sensitive and efficient detection of Ebola RNA but currently available LAMP primers⁵⁶⁰ do not match the outbreak strain. Using sequences from the recent outbreak^{554,561} and the methods described in Chapter 5, we designed primers to match all known Zaire ebolavirus 6.1. These primer combined with simple lab-on-a-chip devices for purifying blood plasma⁴⁹⁷ and imaging fluorescent signals^{512,537} could allow rapid point-of-care detection of Ebolavirus.

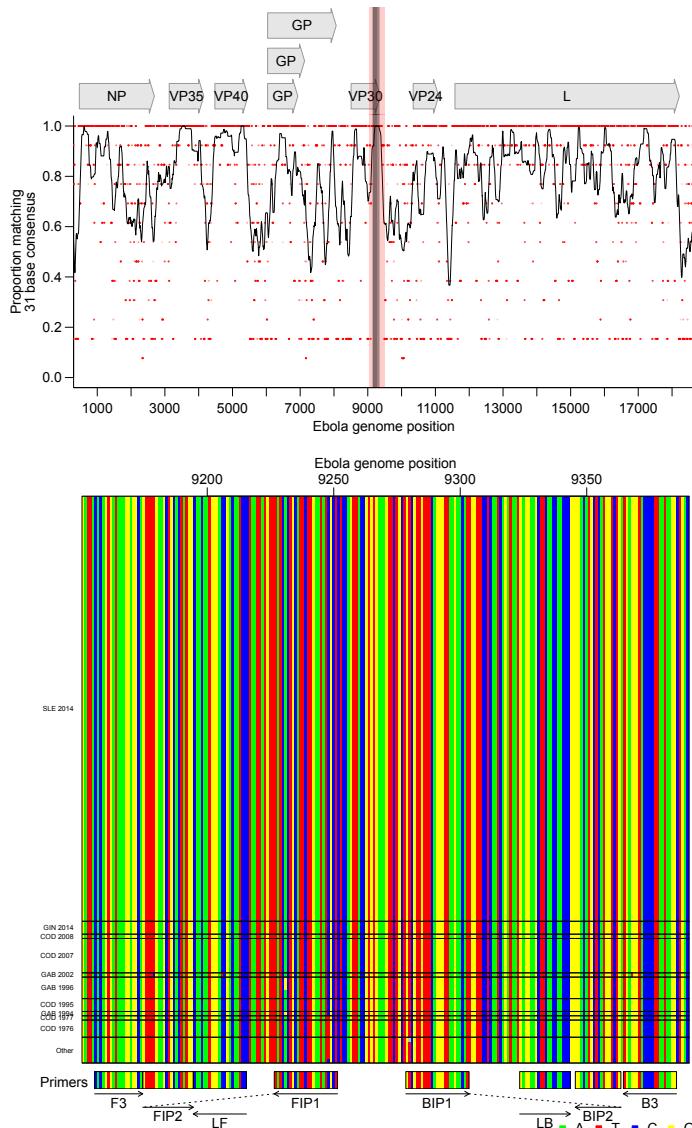


Figure 6.1: Bioinformatic analysis to design Ebola RT-LAMP primers. A) Conservation of sequence in Ebola. Ebola genomes ($n = 131$) from Genbank and sequences from the recent Zaire Ebolavirus outbreak⁵⁵⁴ were aligned and conservation calculated. The x-axis shows the coordinate on the Ebola genome, the y-axis shows the proportion of sequences matching the consensus for each 21 base segment of the genome (red points). The black line shows a 101 base sliding average over these proportions. The vertical red shading shows the region targeted for LAMP primer design that was used as input into the EIKEN primer design tool. Numbering is relative to the Ebola Mayinga sequence. B) Aligned genomes, showing the locations of the preliminary primers. Sequences in the red shaded region in A are shown, with DNA bases color-coded as shown at the lower right. Each row indicates an HIV sequence and each column a base in that sequence. Horizontal lines separate Ebolavirus outbreaks (labeled at left). Arrows indicate the strand targeted by each primer. Primers targeting the negative strand of the virus are shown as reverse compliments for ease of viewing.

6.4.1 Conclusions

These studies contribute to the study and treatment of HIV-1 by revealing aspects of latency, expression and host response. They highlight the importance of primary cell models and the effects that host cell can have on viral processes. With rapidly increasing sequencing throughput, studies like those presented in this thesis offer the opportunity for a deeper and broader understanding of HIV-1 biology and host response and further development of diagnostics and therapeutics.

APPENDIX A.1 : Generalized linear models of changes in use of mutually exclusive HIV-1 splice acceptors

Reads splicing from D1 to one of five mutually exclusive acceptors, D3, D4c, D4a, D4b, D5, and D5a, in three primers, 1.2, 1.3 and 1.4, were collected. Since these data are based on counts, we modeled them as Poisson distributed with an extra variance term allowing for additional variance using a quasi-Poisson generalized linear model with log link. We accounted for differences in sequencing effort by including the total number of D1 to mutually exclusive acceptors reads in each primer-sample as an offset. Differences in the read counts a) over time,b) between human donor and c) cell type were analyzed separately. A term was included for each acceptor and its interaction with the variable of interest. The models included primer and replicate terms and their individual interactions with acceptor to account for any confounding factors.

A.1.1 HOS vs T Cells

R command:

```
glm(count~offset(log(total)) + acceptor:primer + acceptor:isHos  
+ acceptor, data = mutEx[mutEx$time == 48, ],  
family = 'quasipoisson')
```

Difference between HOS and T cells may be confounded by run differences between early sequencing and later sequencing. Verification by agarose gel (Figure 3.4b) suggest that these differences are likely biological.

Variable	Df	Deviance	Resid. Df	Resid. Dev	F	Pr(>F)
NULL	395	138 330				
acceptor	5	133 985	390	4345	9004	$<2.2 \times 10^{-16}$
acceptor:primer	12	751	378	3594	21.03	$<2.2 \times 10^{-16}$
acceptor:isHos	6	2466	372	1127	138.1	$<2.2 \times 10^{-16}$

So after accounting for primer-acceptor bias, the difference between HOS and T cells is significant.

The interesting terms in the model are:

Variable	Estimate	Std. Error	t value	Pr(> t)
acceptorA3:isHosTRUE	1.4717	0.065 86	22.35	$<2.2 \times 10^{-16}$
acceptorA4a:isHosTRUE	-0.9449	0.1246	-7.583	2.73×10^{-13}
acceptorA4b:isHosTRUE	-0.9285	0.1059	-8.767	$<2.2 \times 10^{-16}$
acceptorA4c:isHosTRUE	-1.228	0.1066	-11.51	$<2.2 \times 10^{-16}$
acceptorA5:isHosTRUE	0.090 82	0.026 08	3.483	0.000 555
acceptorA5a:isHosTRUE	0.6308	0.079 40	7.945	2.33×10^{-14}

So it appears A3 is up; A4c, A4a and A4b are down; A5 is up a little and A5a up in HOS.

A.1.2 HOS Over Time

R command:

```
glm(value~offset(log(total)) + acceptor + acceptor:primer
+ acceptor:time, data=mutEx[mutEx$isHos, ],
family ='quasipoisson')
```

Looking only within HOS, we see a significant linear effect of time:

Variable	Df	Deviance	Resid. Df	Resid. Dev	F	Pr(>F)
NULL	53	17962				
acceptor	5	17710	48	252.2	6698	$<2.2 \times 10^{-16}$
acceptor:primer	12	18.0	36	234.2	2.834	0.01018
acceptor:time	6	217.8	30	16.4	68.65	3.57×10^{-16}

We are assuming that a particular acceptor will have the same change in all three primers here.

The interesting terms are:

Variable	Estimate	Std. Error	t value	Pr(> t)
acceptorA3:time	0.02477	0.001778	13.93	1.22×10^{-14}
acceptorA4a:time	-0.01621	0.002812	-5.765	2.69×10^{-6}
acceptorA4b:time	-0.02526	0.002271	-11.12	3.62×10^{-12}
acceptorA4c:time	0.015867	0.003050	5.202	1.32×10^{-5}
acceptorA5:time	-0.001918	0.0006313	-3.038	0.0049
acceptorA5a:time	0.004919	0.001969	2.499	0.0182

So A3, A4c and A5a increase over time and A4a, A4b and A5 decrease over time. All of these coefficients are with a log link and linear and so multiplicative. That means that for example A3 will increase 2.5%/hour ($\exp(0.0247)$) or equivalently 81% (1.025^{24}) over 24hours.

A.1.3 Between Human Comparison

R command:

```
glm(value~offset(log(total)) + acceptor + acceptor:run
+ acceptor:primer + acceptor:subject,
data=mutEx[!mutEx$isHos,], family = 'quasipoisson')
```

In humans, we added a term to account for any potential run bias between the three replicates. Subject refers to the seven human blood donors from which T cells were collected:

Variable	Df	Deviance	Resid. Df	Resid. Dev	F	Pr(>F)
NULL	377	128 430				
acceptor	5	126 446	372	1985	19 598	$<2.2 \times 10^{-16}$
acceptor:run	12	136	360	1849	8.792	1.77×10^{-14}
acceptor:primer	12	850	348	998	54.91	$<2.2 \times 10^{-16}$
acceptor:subject	36	597	312	401	12.86	$<2.2 \times 10^{-16}$

So after accounting for any run and primer bias, subject ID has a statistically significant effect on our observed counts. If we compare everything to subject 7, the interesting terms are:

Variable	Estimate	Std. Error	t value	Pr(> t)
acceptorA3:subject6	-0.001 399	0.072 86	-0.019	0.9847
acceptorA4a:subject6	-0.112 90	0.049 44	-2.284	0.023 07
acceptorA4b:subject6	-0.054 33	0.040 38	-1.345	0.1795
acceptorA4c:subject6	0.028 29	0.033 60	0.842	0.4005
acceptorA5:subject6	0.016 83	0.016 00	1.051	0.2939
acceptorA5a:subject6	-0.030 85	0.060 92	-0.506	0.6129
acceptorA3:subject5	-0.077 67	0.074 23	-1.046	0.2962
acceptorA4a:subject5	-0.1144	0.049 82	-2.296	0.0223
acceptorA4b:subject5	-0.0684	0.040 90	-1.672	0.0956
acceptorA4c:subject5	-0.085 85	0.034 75	-2.471	0.0140
acceptorA5:subject5	0.038 88	0.016 16	2.406	0.0167
acceptorA5a:subject5	0.078 77	0.060 38	1.304	0.1930
acceptorA3:subject4	-0.1849	0.095 78	-1.931	0.0544
acceptorA4a:subject4	0.071 86	0.057 91	1.241	0.2156
acceptorA4b:subject4	0.126 20	0.047 14	2.677	0.0078
acceptorA4c:subject4	-0.100 21	0.043 03	-2.329	0.0205
acceptorA5:subject4	-0.001 16	0.019 69	-0.059	0.9531
acceptorA5a:subject4	0.023 46	0.073 53	0.319	0.7499
acceptorA3:subject3	-0.003 51	0.086 65	-0.041	0.9677
acceptorA4a:subject3	0.071 07	0.055 64	1.277	0.2024
acceptorA4b:subject3	0.006 46	0.046 99	0.138	0.8907
acceptorA4c:subject3	-0.063 34	0.040 76	-1.554	0.1212
acceptorA5:subject3	0.010 52	0.018 87	0.557	0.5776
acceptorA5a:subject3	-0.070 95	0.072 85	-0.974	0.3309
acceptorA3:subject2	-0.2329	0.091 76	-2.539	0.0116
acceptorA4a:subject2	0.024 05	0.056 43	0.426	0.6702
acceptorA4b:subject2	0.1107	0.045 35	2.441	0.0152
acceptorA4c:subject2	0.021 76	0.039 52	0.551	0.5823
acceptorA5:subject2	-0.003 760	0.018 69	-0.201	0.8407
acceptorA5a:subject2	-0.1608	0.073 51	-2.187	0.0295
acceptorA3:subject1	0.095 36	0.065 56	1.454	0.1468
acceptorA4a:subject1	0.029 32	0.044 31	0.662	0.5087
acceptorA4b:subject1	-0.2144	0.038 43	-5.578	5.28×10^{-8}
acceptorA4c:subject1	-0.3974	0.033 85	-11.74	$<2.2 \times 10^{-16}$
acceptorA5:subject1	0.091 44	0.014 70	6.221	1.58×10^{-9}
acceptorA5a:subject1	0.027 47	0.055 94	0.491	0.6238

So there were small but significant effects between subjects especially between subject 1 and subjects 2–7. Interestingly T cells were collected from apheresis product in subject

1 and from whole blood in subjects 2–7 although why this would affect later assays is unknown.

APPENDIX A.2 : Reproducible report of HIV integration sites and latency analysis

A.2.1 Supplementary data

Additional File 2 is a gzipped csv file that includes a row for each uniquely mapped provirus and its surrounding genomic annotations. The csv file should have 12436 rows (excluding header) with 6252 expressed and 6184 latent proviruses.

```
integrationData <- read.csv("AdditionalFile2.csv.gz",
  stringsAsFactors = FALSE)

nrow(integrationData)

## [1] 12436

table(integrationData$isLatent)

##
## FALSE    TRUE
##   6252   6184
```

A.2.2 Lasso regression

The lasso regressions take a while to run so I've turned down the number of cross validations here (set eval=FALSE below to completely skip this step). Leave one out and 480-fold cross validation were used in the paper but processing may take a few days without parallel processing. Lasso regression requires the R glmnet package.

```

notFitColumns <- c("id", "chr", "pos", "strand", "sample", "isLatent")

samples <- unique(as.character(integrationData$sample))

sampleMatrix <- do.call(cbind, lapply(samples, function(x)
  integrationData$sample ==
  x) )

colnames(sampleMatrix) <- gsub(" ", "_", samples)

interact <- function(predMatrix, columns, addNames = NULL) {
  out <- do.call(cbind, lapply(1:ncol(columns), function(x)
    predMatrix *
    columns[, x]))
  if (!is.null(addNames)) {
    if (length(addNames) != ncol(columns)) {
      stop(simpleError("Names not same length as columns"))
    }
    colnames(out) <- sprintf("%s_%s", rep(addNames, each =
      ncol(predMatrix)),
      rep(colnames(predMatrix), length(addNames)))
  }
  return(out)
}

fitData <- as.matrix(integrationData[, !colnames(integrationData)
  %in%
  notFitColumns])

```

```
library(glmnet)

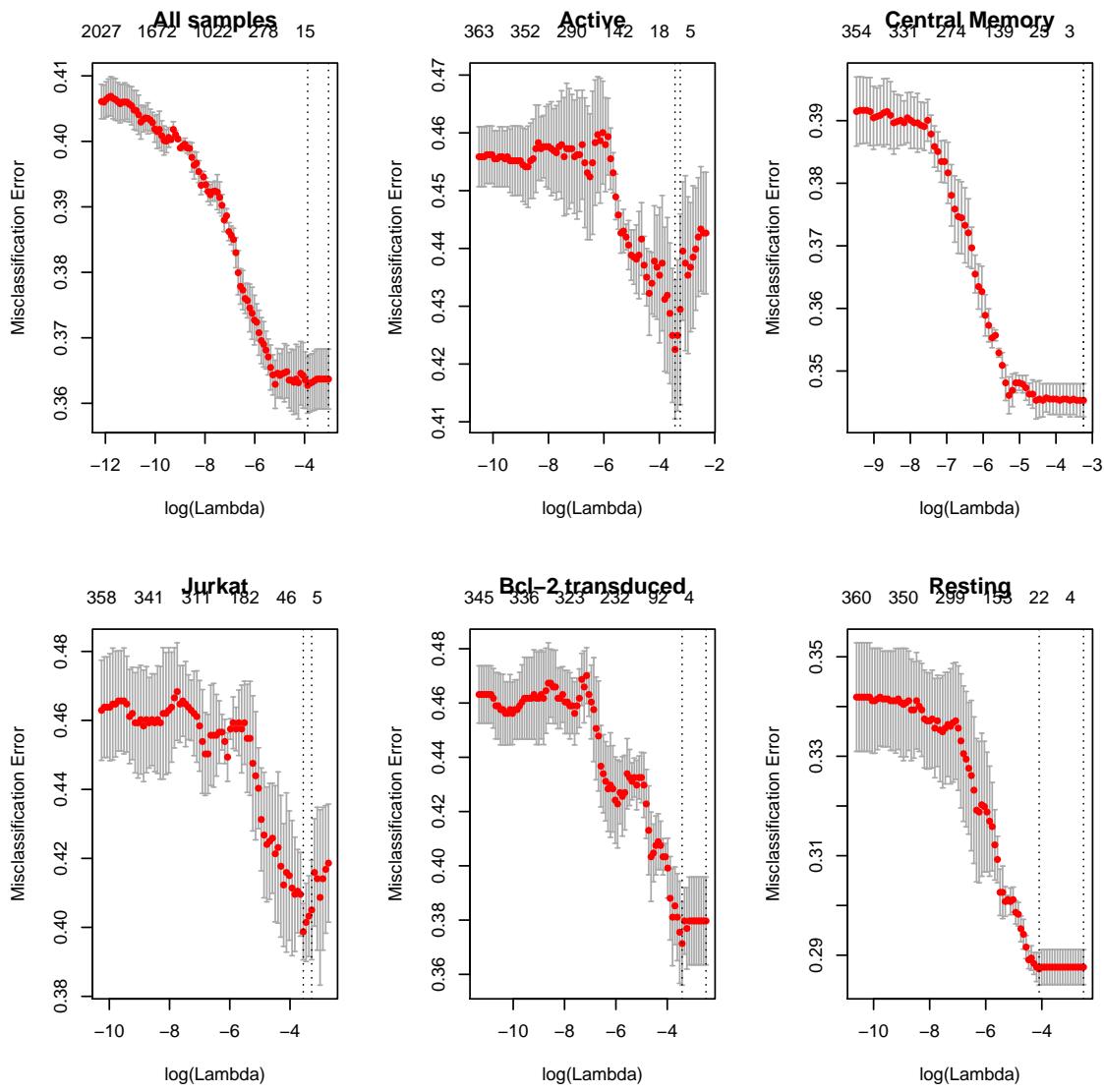
penalties <- rep(1, ncol(fitData2))

penalties[ncol(fitData2) - (ncol(sampleMatrix):1) + 1] <- 0

lassoFit <- cv.glmnet(fitData2, integrationData$isLatent, family
= "binomial",
type.measure = "class", nfolds = 3, penalty.factor =
penalties)

seperateFits <- lapply(samples, function(x) cv.glmnet(fitData[
integrationData$sample ==
x, ], integrationData$isLatent[integrationData$sample ==
x], family = "binomial", type.measure = "class", nfolds = 3))

names(seperateFits) <- samples
```



A.2.3 Correlation

We looked for correlation between the genomic variables and expression status of the proviruses.

```
corMat <- apply(fitData, 2, function(x) sapply(samples, function(
  y) {
    selector <- integrationData$sample == y
```

```

if (sd(x[selector]) == 0)
  return(0)

isLatent <- integrationData[selector, "isLatent"]
cor(as.numeric(isLatent), x[selector], method = "spearman")
} )

quantile(corrMat, seq(0, 1, 0.1))

##          0%        10%        20%        30%
## -0.185223020 -0.081555830 -0.048938130 -0.030895834
##          40%        50%        60%        70%
## -0.018053321 -0.005613895  0.003580982  0.017822483
##          80%        90%       100%
##  0.036694554  0.062003356  0.170642314

```

If we looked for genomic variables consistently correlated or anti-correlated with proviral expression status with an FDR q-value less than 0.01, no variable was significantly correlated in more than 3 samples.

```

pMat <- apply(fitData, 2, function(x) sapply(samples, function(y)
{
  selector <- integrationData$sample == y
  if (sd(x[selector]) == 0)
    return(NA)
  isLatent <- integrationData[selector, "isLatent"]
  cor.test(as.numeric(isLatent), x[selector], method =
  "spearman",
  exact = FALSE)$p.value
})

```

```

} )

adjustPMat <- pMat

adjustPMat[, ] <- p.adjust(pMat, "fdr")

downPMat <- upPMat <- adjustPMat

downPMat[corMat > 0] <- 1

upPMat[corMat < 0] <- 1

table(apply(upPMat < 0.01 & !is.na(upPMat), 2, sum))

##
##    0    1    2    3
## 298  27  38  10

table(apply(downPMat < 0.01 & !is.na(downPMat), 2, sum))

##
##    0    1    2    3
## 216  36  63  58

```

A.2.4 RNA expression

We fit a logistic regression to a polynomial of log RNA-Seq reads within 5000 bases from Jurkat cells for the Jurkat sample and T cells for the rest.

```

rna <- ifelse(integrationData$sample == "Jurkat",
               integrationData$log_jurkatRNA,
               integrationData$rna_5000)

```

```

rna2 <- rna^2

rna3 <- rna^3  #

rna4 <- rna^4

glmData <- data.frame(isLatent = integrationData$isLatent, sample
= integrationData$sample,
rna, rna2, rna3, rna4)

glmMod <- glm(isLatent ~ sample * rna + sample * rna2 + sample *
rna3 + sample * rna4, data = glmData, family = "binomial")

summary(glmMod)

##
## Call:
## glm(formula = isLatent ~ sample * rna + sample * rna2 + sample
##
##       * rna3 + sample * rna4, family = "binomial", data = glmData)
##
## Deviance Residuals:
##      Min        1Q    Median        3Q        Max
## -2.2899 -0.9864 -0.8676  1.0960  1.6007
##
## Coefficients:
##                               Estimate Std. Error z value
## (Intercept)               1.7623655  0.2138859  8.240
## sampleBcl-2 transduced   -2.1625912  0.7061524 -3.062
## sampleCentral Memory     -2.5010063  0.2437685 -10.260
## sampleJurkat              -2.0800202  0.2836871 -7.332

```

```

## sampleResting          0.7840481  0.3312247  2.367
## rna                   -0.6567268  0.2344422 -2.801
## rna2                  0.1387703  0.0770589  1.801
## rna3                  -0.0167219  0.0094076 -1.777
## rna4                  0.0007572  0.0003845  1.969
## sampleBcl-2 transduced:rna 0.5750186  0.6366537  0.903
## sampleCentral Memory:rna  0.9067758  0.2750955  3.296
## sampleJurkat:rna       0.5294036  0.3867163  1.369
## sampleResting:rna      0.0366276  0.3436248  0.107
## sampleBcl-2 transduced:rna2 -0.0369353  0.1878816 -0.197
## sampleCentral Memory:rna2 -0.2106715  0.0915492 -2.301
## sampleJurkat:rna2      -0.0766215  0.1641153 -0.467
## sampleResting:rna2     -0.0760450  0.1086998 -0.700
## sampleBcl-2 transduced:rna3 0.0032503  0.0213743  0.152
## sampleCentral Memory:rna3  0.0237064  0.0112661  2.104
## sampleJurkat:rna3      0.0042183  0.0263910  0.160
## sampleResting:rna3     0.0153132  0.0128711  1.190
## sampleBcl-2 transduced:rna4 -0.0002532  0.0008267 -0.306
## sampleCentral Memory:rna4 -0.0009877  0.0004627 -2.135
## sampleJurkat:rna4      0.0001725  0.0014215  0.121
## sampleResting:rna4     -0.0008049  0.0005119 -1.572
##                                     Pr(>|z|)
## (Intercept) < 2e-16 ***
## sampleBcl-2 transduced 0.00219 **
## sampleCentral Memory   < 2e-16 ***
## sampleJurkat           2.27e-13 ***
## sampleResting          0.01793 *

```

```

## rna          0.00509 ** 
## rna2         0.07173 . 
## rna3         0.07549 . 
## rna4         0.04891 * 
## sampleBcl-2 transduced:rna   0.36643 
## sampleCentral Memory:rna    0.00098 *** 
## sampleJurkat:rna            0.17101 
## sampleResting:rna           0.91511 
## sampleBcl-2 transduced:rna2 0.84415 
## sampleCentral Memory:rna2   0.02138 * 
## sampleJurkat:rna2           0.64059 
## sampleResting:rna2          0.48419 
## sampleBcl-2 transduced:rna3 0.87913 
## sampleCentral Memory:rna3   0.03536 * 
## sampleJurkat:rna3           0.87301 
## sampleResting:rna3          0.23415 
## sampleBcl-2 transduced:rna4 0.75939 
## sampleCentral Memory:rna4   0.03280 * 
## sampleJurkat:rna4           0.90339 
## sampleResting:rna4          0.11585 
## --- 
## Signif. codes: 
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
## 
## (Dispersion parameter for binomial family taken to be 1) 
## 
## Null deviance: 17240  on 12435  degrees of freedom

```

```
## Residual deviance: 15874 on 12411 degrees of freedom
## AIC: 15924
##
## Number of Fisher Scoring iterations: 4
```

A.2.5 Strand orientation

We used a Fisher's exact test to check if silent/inducible proviruses were enriched when integrated in the same strand orientation as cellular genes.

```
selector <- integrationData$inGene == 1

strandTable <- with(integrationData[selector, ], table(ifelse(
  isLatent,
  "Silent/Inducible", "Active"), ifelse(inGeneSameStrand ==
  1, "Same", "Diff"), sample))

apply(strandTable, 3, fisher.test)

## $Active
##
##      Fisher's Exact Test for Count Data
##
## data: array(newX[, i], d.call, dn.call)
## p-value = 0.06061
## alternative hypothesis: true odds ratio is not equal to 1
## 95 percent confidence interval:
##  0.7219466 1.0081995
## sample estimates:
## odds ratio
```

```
##  0.8532127
##
##
## $`Bcl-2 transduced`
##
##      Fisher's Exact Test for Count Data
##
## data: array(newX[, i], d.call, dn.call)
## p-value = 2.177e-05
## alternative hypothesis: true odds ratio is not equal to 1
## 95 percent confidence interval:
##  1.446896 2.872562
## sample estimates:
## odds ratio
##  2.036148
##
##
## $`Central Memory`
##
##      Fisher's Exact Test for Count Data
##
## data: array(newX[, i], d.call, dn.call)
## p-value = 0.2907
## alternative hypothesis: true odds ratio is not equal to 1
## 95 percent confidence interval:
##  0.9386167 1.2320238
## sample estimates:
```

```
## odds ratio
##      1.07529
##
##
## $Jurkat
##
##      Fisher's Exact Test for Count Data
##
## data: array(newX[, i], d.call, dn.call)
## p-value = 0.1674
## alternative hypothesis: true odds ratio is not equal to 1
## 95 percent confidence interval:
##  0.9207548 1.5699893
## sample estimates:
## odds ratio
##      1.202007
##
##
## $Resting
##
##      Fisher's Exact Test for Count Data
##
## data: array(newX[, i], d.call, dn.call)
## p-value = 0.5732
## alternative hypothesis: true odds ratio is not equal to 1
## 95 percent confidence interval:
##  0.7825231 1.1405158
```

```
## sample estimates:  
## odds ratio  
## 0.9447415
```

A.2.6 Acetylation

To reduce correlation between acetylation marks, we generated the first ten principal components of the acetylation data and ran a logistic regression against them. We compared the cross validated performance of this regression with a base model only including which dataset the integration site came from. The cross-validation here has been reduced for efficiency but 480-fold cross-validation was used in the paper.

```
acetyl <- integrationData[, !grepl("logDist", colnames(  
integrationData)) &  
grepl("ac", colnames(integrationData))]  
  
acetylPCA <- princomp(acetyl)  
  
cumsum(acetylPCA$sdev[1:10]^2/sum(acetylPCA$sdev^2))  
  
## Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6  
## 0.5947268 0.6786611 0.7267433 0.7610502 0.7833616 0.7964470  
## Comp.7 Comp.8 Comp.9 Comp.10  
## 0.8093295 0.8215027 0.8299358 0.8372584  
  
cv.glm <- function(model, K = nrow(thisData), subsets = NULL) {  
  modelCall <- model$call  
  thisData <- eval(modelCall$data)  
  n <- nrow(thisData)  
  if (is.null(subsets))
```

```

    subsets <- split(1:n, sample(rep(1:K, length.out = n)))

preds <- lapply(subsets, function(outGroup) {
  subsetData <- thisData[-outGroup, , drop = FALSE]
  predData <- thisData[outGroup, , drop = FALSE]
  thisModel <- modelCall
  thisModel$data <- subsetData
  return(predict(eval(thisModel), predData))
})

pred <- unlist(preds)[order(unlist(subsets))]

subsetId <- rep(1:K, sapply(subsets, length))[order(unlist(
  subsets))]

return(data.frame(pred, subsetId))
}

inData <- data.frame(isLatent = integrationData$isLatent, sample
= as.factor(integrationData$sample),
acetylPCA$score[, 1:10])

modelPreds <- cv.glm(glm(isLatent ~ sample + Comp.1 + Comp.2 +
Comp.3 + Comp.4 + Comp.5 + Comp.6 + Comp.7 + Comp.8 + Comp.9 +
Comp.10, family = "binomial", data = inData), K = 5)

basePreds <- cv.glm(glm(isLatent ~ sample, family = "binomial",
data = inData), subsets = split(1:nrow(inData),
modelPreds$subsetId),
K = 5)

modelCorrect <- sum((modelPreds$pred > 0) ==
integrationData$isLatent)

```

```

baseCorrect <- sum((basePreds$pred > 0) ==
  integrationData$isLatent)

prop.test(c(baseCorrect, modelCorrect), rep(nrow(integrationData),
  ,
  2))

## 

##      2-sample test for equality of proportions with
##      continuity correction

## 

## data: c(baseCorrect, modelCorrect) out of rep(nrow(
## integrationData), 2)
## X-squared = 0.00017372, df = 1, p-value = 0.9895
## alternative hypothesis: two.sided
## 95 percent confidence interval:
## -0.01187726 0.01219890
## sample estimates:
## prop 1    prop 2
## 0.6362978 0.6361370

```

A.2.7 Gene deserts

We used Fisher's exact test to look for an association between integration outside a gene and proviral expression status.

```

geneTable <- table(integrationData$isLatent,
  integrationData$inGene,
  integrationData$sample)

```

```
apply(geneTable, 3, fisher.test)

## $Active

##
##      Fisher's Exact Test for Count Data
##
## data: array(newX[, i], d.call, dn.call)
## p-value < 2.2e-16
## alternative hypothesis: true odds ratio is not equal to 1
## 95 percent confidence interval:
## 0.3629548 0.5446204
## sample estimates:
## odds ratio
## 0.4452621
##
## 
## 
## $`Bcl-2 transduced`

##
##      Fisher's Exact Test for Count Data
##
## data: array(newX[, i], d.call, dn.call)
## p-value = 0.1052
## alternative hypothesis: true odds ratio is not equal to 1
## 95 percent confidence interval:
## 0.9203418 2.3478599
## sample estimates:
## odds ratio
## 1.472224
```

```
##  
##  
## $`Central Memory`  
##  
## Fisher's Exact Test for Count Data  
##  
## data: array(newX[, i], d.call, dn.call)  
## p-value = 0.7803  
## alternative hypothesis: true odds ratio is not equal to 1  
## 95 percent confidence interval:  
## 0.8525329 1.1253952  
## sample estimates:  
## odds ratio  
## 0.9791165  
##  
##  
## $Jurkat  
##  
## Fisher's Exact Test for Count Data  
##  
## data: array(newX[, i], d.call, dn.call)  
## p-value = 0.5443  
## alternative hypothesis: true odds ratio is not equal to 1  
## 95 percent confidence interval:  
## 0.7909269 1.6167285  
## sample estimates:  
## odds ratio
```

```

##      1.127836
##
##
## $Resting
##
##      Fisher's Exact Test for Count Data
##
## data:  array(newX[, i], d.call, dn.call)
## p-value = 3.071e-08
## alternative hypothesis: true odds ratio is not equal to 1
## 95 percent confidence interval:
##  0.4384828 0.6864112
## sample estimates:
## odds ratio
## 0.5500205

```

We used a two-sample t-test to investigate whether there was a significant difference in distance to the nearest gene between expressed and silent/inducible proviruses integrated outside genes.

```

geneDistData <- integrationData[!integrationData$inGene, c(
  "isLatent",
  "logDist_nearest", "sample")]

by(geneDistData, geneDistData$sample, function(x) t.test(
  logDist_nearest ~
  isLatent, data = x))

## geneDistData$sample: Active

```

```

## Welch Two Sample t-test
##
## data: logDist_nearest by isLatent
## t = -2.4539, df = 287.73, p-value = 0.01472
## alternative hypothesis: true difference in means is not equal
## to 0
## 95 percent confidence interval:
## -0.80738340 -0.08867607
## sample estimates:
## mean in group FALSE mean in group TRUE
## 9.608737 10.056767
##
## -----
## geneDistData$sample: Bcl-2 transduced
##
## Welch Two Sample t-test
##
## data: logDist_nearest by isLatent
## t = 0.40978, df = 86.2, p-value = 0.683
## alternative hypothesis: true difference in means is not equal
## to 0
## 95 percent confidence interval:
## -0.6309351 0.9586004
## sample estimates:
## mean in group FALSE mean in group TRUE
## 9.036872 8.873039

```

```

## 
## -----
## geneDistData$sample: Central Memory
##
##      Welch Two Sample t-test
##
## data:  logDist_nearest by isLatent
## t = -0.07188, df = 861.61, p-value = 0.9427
## alternative hypothesis: true difference in means is not equal
## to 0
## 95 percent confidence interval:
## -0.2371374 0.2203819
## sample estimates:
## mean in group FALSE  mean in group TRUE
## 10.19225          10.20063
##
## -----
## geneDistData$sample: Jurkat
##
##      Welch Two Sample t-test
##
## data:  logDist_nearest by isLatent
## t = -1.8217, df = 139.56, p-value = 0.07064
## alternative hypothesis: true difference in means is not equal
## to 0
## 95 percent confidence interval:
## -1.26342086 0.05167979

```

```

## sample estimates:

## mean in group FALSE mean in group TRUE

##          9.925782      10.531652

## -----
## geneDistData$sample: Resting

## Welch Two Sample t-test

## data: logDist_nearest by isLatent
## t = -5.1275, df = 193.49, p-value = 7.096e-07
## alternative hypothesis: true difference in means is not equal
## to 0
## 95 percent confidence interval:
## -1.2687917 -0.5638568
## sample estimates:
## mean in group FALSE mean in group TRUE
##          9.489931      10.406255

```

To check for a relationship between silent/inducible status and distance to CpG islands, we used a two sample t-test on the logged distance and saw a significant difference between silent/inducible and expressed proviruses (before accounting for a correlation between being near CpG islands and in genes)

```

t.test(integrationData$logDist_cpg ~ integrationData$isLatent)

## Welch Two Sample t-test

```

```

## 

## data: integrationData$logDist_cpg by integrationData$isLatent
## t = -2.0233, df = 12381, p-value = 0.04306
## alternative hypothesis: true difference in means is not equal
## to 0

## 95 percent confidence interval:
## -0.105657514 -0.001675563

## sample estimates:

## mean in group FALSE mean in group TRUE
## 10.16362 10.21728

sapply(unique(integrationData$sample), function(x) with(
  integrationData[integrationData$sample ==
    x, ], p.adjust(t.test(logDist_cpg ~ isLatent)$p.value, method
    = "bonferroni",
    n = 5)))

##          Active   Central Memory        Jurkat
## 0.512040457 1.000000000 1.000000000
## Bcl-2 transduced      Resting
## 1.000000000 0.005866539

```

Many CpG islands are found near genes. To account for this relationship, we used an ANOVA test including whether the integration site was inside a gene prior to including CpG islands. After including integration inside genes, CpG islands were not significantly associated with silent/inducible status of the proviruses with all samples grouped or individually after Bonferonni correction for multiple comparisons.

```

anova(with(integrationData, glm(isLatent ~ I(logDist_nearest ==
0) + logDist_cpg, family = "binomial")), test = "Chisq")

## Analysis of Deviance Table

## 

## Model: binomial, link: logit

## 

## Response: isLatent

## 

## Terms added sequentially (first to last)

## 

##                               Df Deviance Resid. Df Resid. Dev
## NULL                           12435      17240
## I(logDist_nearest == 0)     1    26.2682    12434      17213
## logDist_cpg                  1     1.1328    12433      17212
## 
## Pr(>Chi)
## NULL
## I(logDist_nearest == 0) 2.971e-07 ***
## logDist_cpg                 0.2872
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

sapply(unique(integrationData$sample), function(x) {
  p.adjust(anova(with(integrationData[integrationData$sample ==
x, ], glm(isLatent ~ I(logDist_nearest == 0) +
logDist_cpg,
family = "binomial")), test = "Chisq")["logDist_cpg",

```

```

    "Pr(>Chi)"], method = "bonferroni", n = 5)
}

##          Active   Central Memory        Jurkat
## 1.0000000 1.0000000 1.0000000
## Bcl-2 transduced      Resting
## 1.0000000 0.2007788

```

A.2.8 Alphoid repeats

When analyzing repetitive elements, we treated each read as an independent observation and included reads with multiple alignments to the genome. Additional File 3 is a gzipped csv file containing a row for each read with multiple alignments and one row for each dereplicated integration site with a single alignment with the count variable indicating the number of reads dereplicated to that integration site. There should be 26,190 rows (excluding header) with 14,494 rows of expressed provirus and 11,696 rows of silent/inducible provirus.

```

repeats <- read.csv("AdditionalFile3.csv.gz", check.names = FALSE
,
stringsAsFactors = FALSE)

nrow(repeats)

## [1] 26190

summary(repeats$isLatent)

##      Mode   FALSE     TRUE     NA 's
## logical 14494    11696      0

```

```
notRepeatColumns <- c("id", "isLatent", "sample", "count")
```

To analyze whether there was an association between proviral expression status and integration within alphoid repeats, we used Fisher's exact test with a Bonferroni correction for five samples. For comparison, we looked at the association between proviral expression and the other repeats in the RepeatMasker database. We did not Bonferroni correct for the multiple repeat types so that the repeats could be compared with the analysis of alphoid repeats (for which we had an a priori hypothesis for an association with latency).

```
dummyX <- rep(c(TRUE, FALSE), 2)

dummyY <- rep(c(TRUE, FALSE), each = 2)

repeatData <- repeats[, !colnames(repeats) %in% notRepeatColumns]

repeatData <- repeatData[, apply(repeatData, 2, sum) > 0]

testRepeats <- function(x, repeats) {
  sapply(samples, function(thisSample, repeats) {
    selector <- repeats$sample == thisSample
    repLatent <- rep(repeats$isLatent[selector],
                     repeats$count[selector])
    repRepeat <- rep(x[selector], repeats$count[selector])
    fisher.test(table(c(dummyX, repLatent), c(dummyY,
                                                repRepeat)) -
                1)$p.value
  }, repeats)
}
```

```

repeatPs <- apply(repeatData, 2, testRepeats, repeats[,
  notRepeatColumns])

table(apply(repeatPs * 5 < 0.05, 2, sum))

## 
##    0     1     2     3
## 611   76   15    1

which(apply(repeatPs * 5 < 0.05, 2, sum) >= 3)

## ALR/Alpha
##          178

p.adjust(repeatPs[, "ALR/Alpha"], "bonferroni")

##          Active   Central Memory        Jurkat
## 5.026890e-02   3.940207e-03   1.027189e-08
## Bcl-2 transduced           Resting
## 1.000000e+00   2.424896e-02

```

A.2.9 Neighbors

We looked at all pairs of viruses on the same chromosome separated by no more than a given distance, e.g. 100 bases, either with all samples pooled or split between within sample pairs or between sample pairs.

```

allNeighbors <- data.frame(id1 = 0, id2 = 0)[0, ]

ids <- 1:nrow(integrationData)

for (chr in unique(integrationData$chr)) {

```

```

chrSelector <- integrationData$chr == chr

neighborPairs <- data.frame(id1 = rep(ids[chrSelector], sum(
  chrSelector)),
  id2 = rep(ids[chrSelector], each = sum(chrSelector)))

neighborPairs <- neighborPairs[neighborPairs$id1 <
  neighborPairs$id2,
]

allNeighbors <- rbind(allNeighbors, neighborPairs)

}

allNeighbors$dist <- abs(integrationData$pos[allNeighbors$id1] -
  integrationData$pos[allNeighbors$id2])

allNeighbors$latent1 <- integrationData$isLatent[allNeighbors$id1
]

allNeighbors$latent2 <- integrationData$isLatent[allNeighbors$id2
]

allNeighbors$sample1 <- integrationData$sample[allNeighbors$id1]

allNeighbors$sample2 <- integrationData$sample[allNeighbors$id2]

allNeighbors <- allNeighbors[allNeighbors$dist <= 1e+06, ]

```

The expected number of matching pairs was calculated as $\sum_{j \in \text{samples}} n_{j,d}(\theta_{j,d}\theta_{\neg j,d} + (1 - \theta_{j,d})(1 - \theta_{\neg j,d}))$ for between sample, $\sum_{j \in \text{samples}} n_{j,d}(\theta_{j,d}^2 + (1 - \theta_{j,d})^2)$ for within sample and $n_d(\theta_d^2 + (1 - \theta_d)^2)$ for all pairs, where $n_{j,d}$ is the number of pairs of proviruses separated by no more than d base pairs where the first provirus is from sample j , $\theta_{j,d}$ is the proportion of silent/inducible proviruses in sample j appearing in at least one

pair of proviruses separated by less than d base pairs and $\neg j$ means all samples except sample j .

```

dists <- unique(round(10^seq(1, 6, 1)))

pairings <- do.call(rbind, lapply(dists, function(x, allNeighbors
) {
  inSelector <- allNeighbors$dist <= x & allNeighbors$sample1
  ==
  allNeighbors$sample2
  outSelector <- allNeighbors$dist <= x & allNeighbors$sample1
  !=
  allNeighbors$sample2
  allSelector <- allNeighbors$dist <= x
  out <- data.frame(dist = x, observedIn = sum(allNeighbors[
    inSelector,
    "latent1"] == allNeighbors[inSelector, "latent2"]),
    observedOut = sum(allNeighbors[outSelector,
    "latent1"] == allNeighbors[outSelector, "latent2"]),
    observedAll = sum(allNeighbors[allSelector, "latent1"] ==
      allNeighbors[allSelector, "latent2"])), totalIn = sum(
      inSelector),
    totalOut = sum(outSelector), totalAll = sum(allSelector))
  out$expectedIn <- sum(with(allNeighbors[inSelector, ], sapply
    (samples,
    function(x) {
      inLatent <- c(latent1[sample1 == x], latent2[sample2
      ==
      x]) [!duplicated(c(id1[sample1 == x], id2[sample2
      ]))]
```

```

      ==
      x] ))]

if (length(inLatent) == 0) return(0)

return(sum(sample1 == x) * (mean(inLatent)^2 + mean(!
inLatent)^2))

})))

out$expectedOut <- sum(with(allNeighbors[outSelector, ],
sapply(samples, function(x) {

inLatent <- c(latent1[sample1 == x], latent2[sample2
==

x]) [ !duplicated(c(id1[sample1 == x], id2[sample2
==

x]))]

outLatent <- c(latent1[sample1 != x], latent2[sample2
!=

x]) [ !duplicated(c(id1[sample1 != x], id2[sample2
!=

x]))]

if (length(inLatent) == 0) return(0)

return(sum(sample1 == x) * (mean(inLatent) * mean(
outLatent) +
mean(!inLatent) * mean(!outLatent)))
})))

out$expectedAll <- sum(with(allNeighbors[allSelector, ],
{
allLatent <- c(latent1, latent2) [ !duplicated(c(id1,
id2))]
```

```

        return(length(latent1) * (mean(allLatent)^2 + mean(!
            allLatent)^2))
    } )
}

return(out)
}, allNeighbors))

rownames(pairings) <- pairings$dist

```

To look for more matches than expected by random pairing between neighboring proviruses, we used a one sample Z-test of proportion to compare the observed number of matching pairs with the expected proportion of pairs.

```

combinations <- c(All = "All", `Between sample` = "Out", `Within
sample` = "In")

lapply(combinations, function(x, pairing) {
    vars <- sprintf(c("observed%s", "expected%s", "total%s"),
                    x)
    expectedProb <- pairing[, vars[2]]/pairing[, vars[3]]
    prop.test(pairing[, vars[1]], pairing[, vars[3]], p =
        expectedProb)
}, pairings["100", ])

## $All
##
##      1-sample proportions test with continuity correction
##
## data:  pairing[, vars[1]] out of pairing[, vars[3]], null
## probability expectedProb

```

```
## X-squared = 13.002, df = 1, p-value = 0.0003111
## alternative hypothesis: true p is not equal to 0.5000141
## 95 percent confidence interval:
## 0.5586837 0.6962353
## sample estimates:
## p
## 0.63
##
##
## $`Between sample`
##
## 1-sample proportions test with continuity correction
##
## data: pairing[, vars[1]] out of pairing[, vars[3]], null
## probability expectedProb
## X-squared = 0.21919, df = 1, p-value = 0.6397
## alternative hypothesis: true p is not equal to 0.4836763
## 95 percent confidence interval:
## 0.3570532 0.5572662
## sample estimates:
## p
## 0.4554455
##
##
## $`Within sample`
##
## 1-sample proportions test with continuity correction
```

```
##  
## data: pairing[, vars[1]] out of pairing[, vars[3]], null  
## probability expectedProb  
## X-squared = 24.446, df = 1, p-value = 7.644e-07  
## alternative hypothesis: true p is not equal to 0.5561437  
## 95 percent confidence interval:  
## 0.7140170 0.8776751  
## sample estimates:  
## p  
## 0.8080808
```

A.2.10 Compiling this document

This document was generated using R's Sweave function (<http://en.wikipedia.org/wiki/Sweave>). If you would like to regenerate this document, download Additional Files 2, 3 and 4 from Sherrill-Mix et al.³²⁴ and make sure the files are all in the same directory and named AdditionalFile2.csv.gz, AdditionalFile3.csv.gz and AdditionalFile4.Rnw. Then compile by going to that directory and using the commands:

```
R CMD Sweave AdditionalFile4.Rnw  
pdflatex AdditionalFile4.tex
```

Note that you will need R and L^AT_EX (and the R package glmnet if you would like to rerun the lasso regressions) installed.

BIBLIOGRAPHY

- [1] M. S. Gottlieb, H. M. Schanker, P. T. Fan, A. Saxon, J. D. Weisman, and I. Pozalski. 1981. *Pneumocystis pneumonia—Los Angeles.* *MMWR Morb Mortal Wkly Rep*, 30:250–252
- [2] A. Friedman-Kien, L. Laubenstein, M. Marmor, K. Hymes, J. Green, A. Ragaz, J. Gottlieb, F. Muggia, R. Demopoulos, and M. Weintraub. 1981. *Kaposi's sarcoma and Pneumocystis pneumonia among homosexual men—New York City and California.* *MMWR Morb Mortal Wkly Rep*, 30:305–308
- [3] K. B. Hymes, T. Cheung, J. B. Greene, N. S. Prose, A. Marcus, H. Ballard, D. C. William, and L. J. Laubenstein. 1981. *Kaposi's sarcoma in homosexual men—a report of eight cases.* *Lancet*, 2:598–600. doi: 10.1016/S0140-6736(81)92740-9
- [4] H. Masur, M. A. Michelis, J. B. Greene, I. Onorato, R. A. Stouwe, R. S. Holzman, G. Wormser, L. Brettman, M. Lange, H. W. Murray, and S. Cunningham-Rundles. 1981. *An outbreak of community-acquired Pneumocystis carinii pneumonia: initial manifestation of cellular immune dysfunction.* *N Engl J Med*, 305:1431–1438. doi: 10.1056/NEJM198112103052402
- [5] F. P. Siegal, C. Lopez, G. S. Hammer, A. E. Brown, S. J. Kornfeld, J. Gold, J. Hassett, S. Z. Hirschman, C. Cunningham-Rundles, and B. R. Adelsberg. 1981. *Severe acquired immunodeficiency in male homosexuals, manifested by chronic perianal ulcerative herpes simplex lesion.* *N Engl J Med*, 305:1439–1444. doi: 10.1056/NEJM198112103052403
- [6] M. S. Gottlieb, R. Schroff, H. M. Schanker, J. D. Weisman, P. T. Fan, R. A. Wolf, and A. Saxon. 1981. *Pneumocystis carinii pneumonia and mucosal candidiasis in previously healthy homosexual men: evidence of a new acquired cellular immunodeficiency.* *N Engl J Med*, 305:1425–1431. doi: 10.1056/NEJM198112103052401
- [7] Y. Laor and R. A. Schwartz. 1979. *Epidemiologic aspects of American Kaposi's sarcoma.* *J Surg Oncol*, 12:299–303. doi: 10.1002/jso.2930120403
- [8] M. B. Klein, F. A. Pereira, and I. Kantor. 1974. *Kaposi Sarcoma complicating systemic lupus erythematosus treated with immunosuppression.* *Arch Dermatol*, 110:602–604. doi: 10.1001/archderm.1974.01630100058014
- [9] B. D. Myers, E. Kessler, J. Levi, A. Pick, J. B. Rosenfeld, and P. Tikvah. 1974. *Kaposi sarcoma in kidney transplant recipients.* *Arch Intern Med*, 133:307–311. doi: 10.1001/archinte.1974.00320140145017
- [10] S. B. Kapadia and J. R. Krause. 1977. *Kaposi's sarcoma after long-term alkylating agent therapy for multiple myeloma.* *South Med J*, 70:1011–1013

- [11] B. Safai and R. A. Good. 1981. Kaposi's sarcoma: a review and recent developments. *CA Cancer J Clin*, 31:2–12. doi: 10.3322/canjclin.31.1.2
- [12] Y. Chang, E. Cesarman, M. S. Pessin, F. Lee, J. Culpepper, D. M. Knowles, and P. S. Moore. 1994. Identification of herpesvirus-like DNA sequences in AIDS-associated Kaposi's sarcoma. *Science*, 266:1865–1869. doi: 10.1126/science.7997879
- [13] F. Sitas, H. Carrara, V. Beral, R. Newton, G. Reeves, D. Bull, U. Jentsch, R. Pacella-Norman, D. Bourboulia, D. Whitby, C. Boshoff, and R. Weiss. 1999. Antibodies against human herpesvirus 8 in black South African patients with cancer. *N Engl J Med*, 340:1863–1871. doi: 10.1056/NEJM199906173402403
- [14] B. A. Burke and R. A. Good. 1973. Pneumocystis carinii infection. *Medicine (Baltimore)*, 52:23–51
- [15] W. T. Hughes. 1977. Pneumocystis carinii pneumonia. *N Engl J Med*, 297: 1381–1383. doi: 10.1056/NEJM197712222972505
- [16] J. Gerstoft, A. Malchow-Møller, I. Bygbjerg, E. Dickmeiss, C. Enk, P. Halberg, S. Haahr, M. Jacobsen, K. Jensen, J. Mejer, J. O. Nielsen, H. K. Thomsen, J. Søndergaard, and I. Lorenzen. 1982. Severe acquired immunodeficiency in European homosexual men. *Br Med J (Clin Res Ed)*, 285:17–19
- [17] H. Masur, M. A. Michelis, G. P. Wormser, S. Lewin, J. Gold, M. L. Tapper, J. Giron, C. W. Lerner, D. Armstrong, U. Setia, J. A. Sender, R. S. Siebken, P. Nicholas, Z. Arlen, S. Maayan, J. A. Ernst, F. P. Siegal, and S. Cunningham-Rundles. 1982. Opportunistic infection in previously healthy women. Initial manifestations of a community-acquired cellular immunodeficiency. *Ann Intern Med*, 97:533–539
- [18] A. Ammann, M. Cowan, D. Wara, H. Goldman, H. Perkins, R. Lanzerotti, J. Gullett, A. Duff, S. Dritz, and J. Chin. 1982. Possible transfusion-associated acquired immune deficiency syndrome (AIDS) — California. *MMWR Morb Mortal Wkly Rep*, 31:652–654
- [19] N. Ehrenkranz, J. Rubini, R. Gunn, C. Horsburgh, T. Collins, U. Hasiba, W. Hathaway, W. Doig, R. Hopkins, and J. Elliott. 1982. Pneumocystis carinii pneumonia among persons with hemophilia A. *MMWR Morb Mortal Wkly Rep*, 31:365–367
- [20] M.-C. Poon, A. Landay, J. Alexander, W. Birch, M. Eyster, H. Al-Mondhiry, J. Ballard, E. Witte, C. Hayes et al. 1982. Update on acquired immune deficiency syndrome (AIDS) among patients with hemophilia A. *MMWR Morb Mortal Wkly Rep*, 31:644–6, 652
- [21] J. B. Greene, G. S. Sidhu, S. Lewin, J. F. Levine, H. Masur, M. S. Simberkoff, P. Nicholas, R. C. Good, S. B. Zolla-Pazner, A. A. Pollock, M. L. Tapper, and R. S. Holzman. 1982. *Mycobacterium avium-intracellulare*: a cause of disseminated

- life-threatening infection in homosexuals and drug abusers. *Ann Intern Med*, 97: 539–546
- [22] R. O'Reilly, D. Kirkpatrick, C. B. Small, R. Klein, H. Keltz, G. Friedland, K. Bromberg, S. Fikrig, H. Mendez et al. 1982. Unexplained immunodeficiency and opportunistic infections in infants—New York, New Jersey, California. *MMWR Morb Mortal Wkly Rep*, 31:665–667
- [23] S. Fannin, M. Gottlieb, J. Weisman, E. Rogolsky, T. Prendergast, J. Chin, A. Friedman-Kien, L. Laubenstein, S. Friedman, and R. Rothenberg. 1982. A cluster of Kaposi's sarcoma and Pneumocystis carinii pneumonia among homosexual male residents of Los Angeles and Orange Counties, California. *MMWR Morb Mortal Wkly Rep*, 31:305–307
- [24] C. Harris, C. B. Small, G. Friedland, R. Klein, B. Moll, E. Emeson, I. Spigland, N. Steigbigel, R. Reiss, S. Friedman, and R. Rothenberg. 1983. Immunodeficiency among female sexual partners of males with acquired immune deficiency syndrome (AIDS) — New York. *MMWR Morb Mortal Wkly Rep*, 31:697–698
- [25] F. Barré-Sinoussi, J. C. Chermann, F. Rey, M. T. Nugeyre, S. Chamaret, J. Gruest, C. Dauguet, C. Axler-Blin, F. Vézinet-Brun, C. Rouzioux, W. Rozenbaum, and L. Montagnier. 1983. Isolation of a T-lymphotropic retrovirus from a patient at risk for acquired immune deficiency syndrome (AIDS). *Science*, 220:868–871
- [26] R. C. Gallo, P. S. Sarin, E. P. Gelmann, M. Robert-Guroff, E. Richardson, V. S. Kalyanaraman, D. Mann, G. D. Sidhu, R. E. Stahl, S. Zolla-Pazner, J. Leibowitch, and M. Popovic. 1983. Isolation of human T-cell leukemia virus in acquired immune deficiency syndrome (AIDS). *Science*, 220:865–867. doi: 10.1126/science.6601823
- [27] M. Popovic, M. G. Sarngadharan, E. Read, and R. C. Gallo. 1984. Detection, isolation, and continuous production of cytopathic retroviruses (HTLV-III) from patients with AIDS and pre-AIDS. *Science*, 224:497–500. doi: 10.1126/science.6200935
- [28] J. A. Levy, A. D. Hoffman, S. M. Kramer, J. A. Landis, J. M. Shimabukuro, and L. S. Oshiro. 1984. Isolation of lymphocytopathic retroviruses from San Francisco patients with AIDS. *Science*, 225:840–842. doi: 10.1126/science.6206563
- [29] R. C. Gallo, S. Z. Salahuddin, M. Popovic, G. M. Shearer, M. Kaplan, B. F. Haynes, T. J. Palker, R. Redfield, J. Oleske, and B. Safai. 1984. Frequent detection and isolation of cytopathic retroviruses (HTLV-III) from patients with AIDS and at risk for AIDS. *Science*, 224:500–503. doi: 10.1126/science.6200936
- [30] M. G. Sarngadharan, M. Popovic, L. Bruch, J. Schüpbach, and R. C. Gallo. 1984. Antibodies reactive with human T-lymphotropic retroviruses (HTLV-III) in the serum of patients with AIDS. *Science*, 224:506–508. doi: 10.1126/science.6324345

- [31] B. Safai, M. G. Sarngadharan, J. E. Groopman, K. Arnett, M. Popovic, A. Sliski, J. Schüpbach, and R. C. Gallo. 1984. Seroepidemiological studies of human T-lymphotropic retrovirus type III in acquired immunodeficiency syndrome. *Lancet*, 1:1438–1440. doi: 10.1016/S0140-6736(84)91933-0
- [32] N. Clumeck, F. Mascart-Lemone, J. de Maubeuge, D. Brenez, and L. Marcelis. 1983. Acquired immune deficiency syndrome in Black Africans. *Lancet*, 1:642. doi: 10.1016/S0140-6736(83)91808-1
- [33] N. Clumeck, J. Sonnet, H. Taelman, S. Cran, and P. Henrivaux. 1984. Acquired immune deficiency syndrome in Belgium and its relation to Central Africa. *Ann NY Acad Sci*, 437:264–269. doi: 10.1111/j.1749-6632.1984.tb37144.x
- [34] P. Van de Perre, D. Rouvroy, P. Lepage, J. Bogaerts, P. Kestelyn, J. Kayihigi, A. C. Hekker, J. P. Butzler, and N. Clumeck. 1984. Acquired immunodeficiency syndrome in Rwanda. *Lancet*, 2:62–65. doi: 10.1016/S0140-6736(84)90240-X
- [35] P. Piot, T. C. Quinn, H. Taelman, F. M. Feinsod, K. B. Minlangu, O. Wobin, N. Mbendi, P. Mazebo, K. Ndangi, and W. Stevens. 1984. Acquired immunodeficiency syndrome in a heterosexual population in Zaire. *Lancet*, 2:65–69. doi: 10.1016/S0140-6736(84)90241-1
- [36] J. N. Nkengasong, W. Janssens, L. Heyndrickx, K. Fransen, P. M. Ndumbe, J. Motte, A. Leonaers, M. Ngolle, J. Ayuk, and P. Piot. 1994. Genotypic subtypes of HIV-1 in Cameroon. *AIDS*, 8:1405–1412
- [37] J. Louwagie, W. Janssens, J. Mascola, L. Heyndrickx, P. Hegerich, G. van der Groen, F. E. McCutchan, and D. S. Burke. 1995. Genetic diversity of the envelope glycoprotein from human immunodeficiency virus type 1 isolates of African origin. *J Virol*, 69:263–271
- [38] N. Vidal, M. Peeters, C. Mulanga-Kabeya, N. Nzilambi, D. Robertson, W. Ilunga, H. Sema, K. Tshimanga, B. Bongo, and E. Delaporte. 2000. Unprecedented degree of human immunodeficiency virus type 1 (HIV-1) group M genetic diversity in the Democratic Republic of Congo suggests that the HIV-1 pandemic originated in Central Africa. *J Virol*, 74:10498–10507. doi: 10.1128/JVI.74.22.10498-10507. 2000
- [39] A. Rambaut, D. L. Robertson, O. G. Pybus, M. Peeters, and E. C. Holmes. 2001. Human immunodeficiency virus. Phylogeny and the origin of HIV-1. *Nature*, 410: 1047–1048. doi: 10.1038/35074179
- [40] C. Yang, B. Dash, S. L. Hanna, H. S. Frances, N. Nzilambi, R. C. Colebunders, M. St Louis, T. C. Quinn, T. M. Folks, and R. B. Lal. 2001. Predominance of HIV type 1 subtype G among commercial sex workers from Kinshasa, Democratic Republic of Congo. *AIDS Res Hum Retroviruses*, 17:361–365. doi: 10.1089/0889220150503726

- [41] M. L. Kalish, K. E. Robbins, D. Pieniazek, A. Schaefer, N. Nzilambi, T. C. Quinn, M. E. St Louis, A. S. Youngpairoj, J. Phillips, H. W. Jaffe, and T. M. Folks. 2004. Recombinant viruses and early global HIV-1 epidemic. *Emerg Infect Dis*, 10: 1227–1234. doi: 10.3201/eid1007.030904
- [42] S. S. Frøland, P. Jenum, C. F. Lindboe, K. W. Wefring, P. J. Linnestad, and T. Böhmer. 1988. HIV-1 infection in Norwegian family before 1970. *Lancet*, 1: 1344–1345. doi: 10.1016/S0140-6736(88)92164-2
- [43] R. F. Garry, M. H. Witte, A. A. Gottlieb, M. Elvin-Lewis, M. S. Gottlieb, C. L. Witte, S. S. Alexander, W. R. Cole, and W. Drake, Jr. 1988. Documentation of an AIDS virus infection in the United States in 1968. *JAMA*, 260:2085–2087. doi: 10.1001/jama.1988.03410140097031
- [44] I. C. Bygbjerg. 1983. AIDS in a Danish surgeon (Zaire, 1976). *Lancet*, 1:925. doi: 10.1016/S0140-6736(83)91348-X
- [45] J. Vandepitte, R. Verwilghen, and P. Zachee. 1983. AIDS and cryptococcosis (Zaire, 1977). *Lancet*, 1:925–926. doi: 10.1016/S0140-6736(83)91349-1
- [46] A. J. Nahmias, J. Weiss, X. Yao, F. Lee, R. Kodsi, M. Schanfield, T. Matthews, D. Bolognesi, D. Durack, and A. Motulsky. 1986. Evidence for human infection with an HTLV III/LAV-like virus in Central Africa, 1959. *Lancet*, 1:1279–1280. doi: 10.1016/S0140-6736(86)91422-4
- [47] T. Zhu, B. T. Korber, A. J. Nahmias, E. Hooper, P. M. Sharp, and D. D. Ho. 1998. An African HIV-1 sequence from 1959 and implications for the origin of the epidemic. *Nature*, 391:594–597. doi: 10.1038/35400
- [48] M. Worobey, M. Gemmel, D. E. Teuwen, T. Haselkorn, K. Kunstman, M. Bunce, J.-J. Muyembe, J.-M. M. Kabongo, R. M. Kalengayi, E. Van Marck, M. T. P. Gilbert, and S. M. Wolinsky. 2008. Direct evidence of extensive diversity of HIV-1 in Kinshasa by 1960. *Nature*, 455:661–664. doi: 10.1038/nature07390
- [49] B. Korber, M. Muldoon, J. Theiler, F. Gao, R. Gupta, A. Lapedes, B. H. Hahn, S. Wolinsky, and T. Bhattacharya. 2000. Timing the ancestor of the HIV-1 pandemic strains. *Science*, 288:1789–1796. doi: 10.1126/science.288.5472.1789
- [50] M. Salemi, K. Strimmer, W. W. Hall, M. Duffy, E. Delaporte, S. Mboup, M. Peeters, and A. M. Vandamme. 2001. Dating the common ancestor of SIVcpz and HIV-1 group M and the origin of HIV-1 subtypes using a new method to uncover clock-like molecular evolution. *FASEB J*, 15:276–278. doi: 10.1096/fj.00-0449fje
- [51] P. M. Sharp, E. Bailes, R. R. Chaudhuri, C. M. Rodenburg, M. O. Santiago, and B. H. Hahn. 2001. The origins of acquired immune deficiency syndrome viruses: where and when? *Philos Trans R Soc Lond B Biol Sci*, 356:867–876. doi: 10.1098/rstb.2001.0863

- [52] K. Yusim, M. Peeters, O. G. Pybus, T. Bhattacharya, E. Delaporte, C. Mulanga, M. Muldoon, J. Theiler, and B. Korber. 2001. Using human immunodeficiency virus type 1 sequences to infer historical features of the acquired immune deficiency syndrome epidemic and human immunodeficiency virus evolution. *Philos Trans R Soc Lond B Biol Sci*, 356:855–866. doi: 10.1098/rstb.2001.0859
- [53] N. R. Faria, A. Rambaut, M. A. Suchard, G. Baele, T. Bedford, M. J. Ward, A. J. Tatem, J. D. Sousa, N. Arinaminpathy, J. Pépin, D. Posada, M. Peeters, O. G. Pybus, and P. Lemey. 2014. The early spread and epidemic ignition of HIV-1 in human populations. *Science*, 346:56–61. doi: 10.1126/science.1256739
- [54] M. Peeters, C. Honoré, T. Huet, L. Bedjabaga, S. Ossari, P. Bussi, R. W. Cooper, and E. Delaporte. 1989. Isolation and partial characterization of an HIV-related virus occurring naturally in chimpanzees in Gabon. *AIDS*, 3:625–630
- [55] T. Huet, R. Cheynier, A. Meyerhans, G. Roelants, and S. Wain-Hobson. 1990. Genetic organization of a chimpanzee lentivirus related to HIV-1. *Nature*, 345: 356–359. doi: 10.1038/345356a0
- [56] E. Bowen-Jones and S. Pendry. 1999. The threat to primates and other mammals from the bushmeat trade in Africa, and how this threat could be diminished. *Oryx*, 33:233–246. doi: 10.1046/j.1365-3008.1999.00066.x
- [57] B. H. Hahn, G. M. Shaw, K. M. De Cock, and P. M. Sharp. 2000. AIDS as a zoonosis: scientific and public health implications. *Science*, 287:607–614. doi: 10.1126/science.287.5453.607
- [58] M. Peeters, V. Courgnaud, B. Abela, P. Auzel, X. Pourrut, F. Bibollet-Ruche, S. Loul, F. Liegeois, C. Butel, D. Koulagna, E. Mpoudi-Ngole, G. M. Shaw, B. H. Hahn, and E. Delaporte. 2002. Risk to human health from a plethora of simian immunodeficiency viruses in primate bushmeat. *Emerg Infect Dis*, 8:451–457. doi: 10.3201/eid0805.010522
- [59] N. D. Wolfe, T. A. Prosser, J. K. Carr, U. Tamoufe, E. Mpoudi-Ngole, J. N. Torimiro, M. LeBreton, F. E. McCutchan, D. L. Birx, and D. S. Burke. 2004. Exposure to nonhuman primates in rural Cameroon. *Emerg Infect Dis*, 10:2094–2099. doi: 10.3201/eid1012.040062
- [60] N. D. Wolfe, W. Heneine, J. K. Carr, A. D. Garcia, V. Shanmugam, U. Tamoufe, J. N. Torimiro, A. T. Prosser, M. Lebreton, E. Mpoudi-Ngole, F. E. McCutchan, D. L. Birx, T. M. Folks, D. S. Burke, and W. M. Switzer. 2005. Emergence of unique primate T-lymphotropic viruses among central African bushmeat hunters. *Proc Natl Acad Sci U S A*, 102:7994–7999. doi: 10.1073/pnas.0501734102
- [61] M. L. Kalish, N. D. Wolfe, C. B. Ndongmo, J. McNicholl, K. E. Robbins, M. Aidoo, P. N. Fonjungo, G. Alemnji, C. Zeh, C. F. Djoko, E. Mpoudi-Ngole, D. S. Burke, and

- T. M. Folks. 2005. Central African hunters exposed to simian immunodeficiency virus. *Emerg Infect Dis*, 11:1928–1930. doi: 10.3201/eid1112.050394
- [62] F. Gao, E. Bailes, D. L. Robertson, Y. Chen, C. M. Rodenburg, S. F. Michael, L. B. Cummins, L. O. Arthur, M. Peeters, G. M. Shaw, P. M. Sharp, and B. H. Hahn. 1999. Origin of HIV-1 in the chimpanzee Pan troglodytes troglodytes. *Nature*, 397:436–441. doi: 10.1038/17130
- [63] B. F. Keele, F. Van Heuverswyn, Y. Li, E. Bailes, J. Takehisa, M. L. Santiago, F. Bibollet-Ruche, Y. Chen, L. V. Wain, F. Liegeois, S. Loul, E. M. Ngole, Y. Bienvéne, E. Delaporte, J. F. Y. Brookfield, P. M. Sharp, G. M. Shaw, M. Peeters, and B. H. Hahn. 2006. Chimpanzee reservoirs of pandemic and nonpandemic HIV-1. *Science*, 313:523–526. doi: 10.1126/science.1126531
- [64] F. Van Heuverswyn, Y. Li, E. Bailes, C. Neel, B. Lafay, B. F. Keele, K. S. Shaw, J. Takehisa, M. H. Kraus, S. Loul, C. Butel, F. Liegeois, B. Yangda, P. M. Sharp, E. Mpoudi-Ngole, E. Delaporte, B. H. Hahn, and M. Peeters. 2007. Genetic diversity and phylogeographic clustering of SIVcpzPtt in wild chimpanzees in Cameroon. *Virology*, 368:155–171. doi: 10.1016/j.virol.2007.06.018
- [65] P. M. Sharp and B. H. Hahn. 2008. AIDS: prehistory of HIV-1. *Nature*, 455: 605–606. doi: 10.1038/455605a
- [66] D. Vangroenweghe. 2001. The earliest cases of human immunodeficiency virus type 1 group M in Congo-Kinshasa, Rwanda and Burundi and the origin of acquired immune deficiency syndrome. *Philos Trans R Soc Lond B Biol Sci*, 356: 923–925. doi: 10.1098/rstb.2001.0876
- [67] A. Chitnis, D. Rawls, and J. Moore. 2000. Origin of HIV type 1 in colonial French Equatorial Africa? *AIDS Res Hum Retroviruses*, 16:5–8. doi: 10.1089/088922200309548
- [68] J. D. de Sousa, V. Müller, P. Lemey, and A.-M. Vandamme. 2010. High GUD incidence in the early 20 century created a particularly permissive time window for the origin and initial spread of epidemic HIV strains. *PLoS One*, 5:e9936. doi: 10.1371/journal.pone.0009936
- [69] J. D. de Sousa, C. Alvarez, A.-M. Vandamme, and V. Müller. 2012. Enhanced heterosexual transmission hypothesis for the origin of pandemic HIV-1. *Viruses*, 4:1950–1983. doi: 10.3390/v4101950
- [70] R. D. Moore and R. E. Chaisson. 1996. Natural history of opportunistic disease in an HIV-infected urban clinical cohort. *Ann Intern Med*, 124:633–642. doi: 10.7326/0003-4819-124-7-199604010-00003
- [71] R. Rothenberg, M. Woelfel, R. Stoneburner, J. Milberg, R. Parker, and B. Truman. 1987. Survival with the acquired immunodeficiency syndrome. Experience with

5833 cases in New York City. *N Engl J Med*, 317:1297–1302. doi: 10.1056/NEJM198711193172101

- [72] S. Vella, M. Giuliano, P. Pezzotti, M. G. Agresti, C. Tomino, M. Floridia, D. Greco, M. Moroni, G. Visco, and F. Milazzo. 1992. Survival of zidovudine-treated patients with AIDS compared with that of contemporary untreated patients. *JAMA*, 267: 1232–1236. doi: 10.1001/jama.1992.03480090080031
- [73] M. M. Deschamps, D. W. Fitzgerald, J. W. Pape, and W. Johnson, Jr. 2000. HIV infection in Haiti: natural history and disease progression. *AIDS*, 14:2515–2521
- [74] K. M. Harrison, R. Song, and X. Zhang. 2010. Life expectancy after HIV diagnosis based on national HIV surveillance data from 25 states, United States. *J Acquir Immune Defic Syndr*, 53:124–130. doi: 10.1097/QAI.0b013e3181b563e7
- [75] Collaborative Group on AIDS Incubation and HIV Survival. 2000. Time from HIV-1 seroconversion to AIDS and death before widespread use of highly-active antiretroviral therapy: a collaborative re-analysis. *Lancet*, 355:1131–1137. doi: 10.1016/S0140-6736(00)02061-4
- [76] Antiretroviral Therapy Cohort Collaboration, M. Zwahlen, R. Harris, M. May, R. Hogg, D. Costagliola, F. de Wolf, J. Gill, G. Fätkenheuer, C. Lewden, M. Saag, S. Staszewski, A. d'Arminio Monforte, J. Casabona, F. Lampe, A. Justice, V. von Wyl, and M. Egger. 2009. Mortality of HIV-infected patients starting potent antiretroviral therapy: comparison with the general population in nine industrialized countries. *Int J Epidemiol*, 38:1624–1633. doi: 10.1093/ije/dyp306
- [77] M. A. Fischl, D. D. Richman, M. H. Grieco, M. S. Gottlieb, P. A. Volberding, O. L. Laskin, J. M. Leedom, J. E. Groopman, D. Mildvan, and R. T. Schooley. 1987. The efficacy of azidothymidine (AZT) in the treatment of patients with AIDS and AIDS-related complex. A double-blind, placebo-controlled trial. *N Engl J Med*, 317:185–191. doi: 10.1056/NEJM198707233170401
- [78] M. A. Fischl, D. D. Richman, D. M. Causey, M. H. Grieco, Y. Bryson, D. Mildvan, O. L. Laskin, J. E. Groopman, P. A. Volberding, and R. T. Schooley. 1989. Prolonged zidovudine therapy in patients with AIDS and advanced AIDS-related complex. *JAMA*, 262:2405–2410. doi: 10.1001/jama.1989.03430170067030
- [79] P. A. Volberding, S. W. Lagakos, M. A. Koch, C. Pettinelli, M. W. Myers, D. K. Booth, H. Balfour, Jr, R. C. Reichman, J. A. Bartlett, M. S. Hirsch, and The AIDS Clinical Trial Group. 1990. Zidovudine in asymptomatic human immunodeficiency virus infection. A controlled trial in persons with fewer than 500 CD4-positive cells per cubic millimeter. *N Engl J Med*, 322:941–949. doi: 10.1056/NEJM199004053221401
- [80] B. H. Hahn, G. M. Shaw, M. E. Taylor, R. R. Redfield, P. D. Markham, S. Z. Salahuddin, F. Wong-Staal, R. C. Gallo, E. S. Parks, and W. P. Parks. 1986.

Genetic variation in HTLV-III/LAV over time in patients with AIDS or at risk for AIDS. *Science*, 232:1548–1553. doi: 10.1126/science.3012778

- [81] B. D. Preston, B. J. Poiesz, and L. A. Loeb. 1988. Fidelity of HIV-1 reverse transcriptase. *Science*, 242:1168–1171. doi: 10.1126/science.2460924
- [82] J. D. Roberts, K. Bebenek, and T. A. Kunkel. 1988. The accuracy of reverse transcriptase from HIV-1. *Science*, 242:1171–1173. doi: 10.1126/science.2460925
- [83] L. M. Mansky and H. M. Temin. 1995. Lower in vivo mutation rate of human immunodeficiency virus type 1 than that predicted from the fidelity of purified reverse transcriptase. *J Virol*, 69:5087–5094
- [84] L. M. Mansky. 1996. The mutation rate of human immunodeficiency virus type 1 is influenced by the vpr gene. *Virology*, 222:391–400. doi: 10.1006/viro.1996.0436
- [85] M. E. Abram, A. L. Ferris, W. Shao, W. G. Alvord, and S. H. Hughes. 2010. Nature, position, and frequency of mutations made in a single cycle of HIV-1 replication. *J Virol*, 84:9864–9878. doi: 10.1128/JVI.00915-10
- [86] V. Achuthan, B. J. Keith, B. A. Connolly, and J. J. DeStefano. 2014. Human immunodeficiency virus reverse transcriptase displays dramatically higher fidelity under physiological magnesium conditions in vitro. *J Virol*, 88:8514–8527. doi: 10.1128/JVI.00752-14
- [87] B. A. Larder, G. Darby, and D. D. Richman. 1989. HIV with reduced sensitivity to zidovudine (AZT) isolated during prolonged therapy. *Science*, 243:1731–1734. doi: 10.1126/science.2467383
- [88] B. A. Larder and S. D. Kemp. 1989. Multiple mutations in HIV-1 reverse transcriptase confer high-level resistance to zidovudine (AZT). *Science*, 246: 1155–1158. doi: 10.1126/science.2479983
- [89] S. Land, G. Terloar, D. McPhee, C. Birch, R. Doherty, D. Cooper, and I. Gust. 1990. Decreased in vitro susceptibility to zidovudine of HIV isolates obtained from patients with AIDS. *J Infect Dis*, 161:326–329. doi: 10.1093/infdis/161.2.326
- [90] C. A. Boucher, M. Tersmette, J. M. Lange, P. Kellam, R. E. de Goede, J. W. Mulder, G. Darby, J. Goudsmit, and B. A. Larder. 1990. Zidovudine sensitivity of human immunodeficiency viruses from high-risk, symptom-free individuals during therapy. *Lancet*, 336:585–590. doi: 10.1016/0140-6736(90)93391-2
- [91] D. D. Richman, J. M. Grimes, and S. W. Lagakos. 1990. Effect of stage of disease and drug dose on zidovudine susceptibilities of isolates of human immunodeficiency virus. *J Acquir Immune Defic Syndr*, 3:743–746
- [92] D. D. Richman, J. C. Guatelli, J. Grimes, A. Tsiantis, and T. Gingras. 1991. Detection of mutations associated with zidovudine resistance in human immun-

- odeficiency virus by use of the polymerase chain reaction. *J Infect Dis*, 164: 1075–1081. doi: 10.1093/infdis/164.6.1075
- [93] J. E. Fitzgibbon, R. M. Howell, C. A. Haberzettl, S. J. Sperber, D. J. Gocke, and D. T. Dubin. 1992. Human immunodeficiency virus type 1 pol gene mutations which cause decreased susceptibility to 2',3'-dideoxycytidine. *Antimicrob Agents Chemother*, 36:153–157. doi: 10.1128/AAC.36.1.153
- [94] D. D. Richman, D. Havlir, J. Corbeil, D. Looney, C. Ignacio, S. A. Spector, J. Sullivan, S. Cheeseman, K. Barringer, and D. Pauletti. 1994. Nevirapine resistance mutations of human immunodeficiency virus type 1 selected during therapy. *J Virol*, 68:1660–1666
- [95] R. Schuurman, M. Nijhuis, R. van Leeuwen, P. Schipper, D. de Jong, P. Collis, S. A. Danner, J. Mulder, C. Loveday, and C. Christopherson. 1995. Rapid changes in human immunodeficiency virus type 1 RNA load and appearance of drug-resistant virus populations in persons treated with lamivudine (3TC). *J Infect Dis*, 171:1411–1419. doi: 10.1093/infdis/171.6.1411
- [96] J. C. Schmit, L. Ruiz, B. Clotet, A. Raventos, J. Tor, J. Leonard, J. Desmyter, E. De Clercq, and A. M. Vandamme. 1996. Resistance-related mutations in the HIV-1 protease gene of patients treated for 1 year with the protease inhibitor ritonavir (ABT-538). *AIDS*, 10:995–999
- [97] T. Creagh-Kirk, P. Doi, E. Andrews, S. Nusinoff-Lehrman, H. Tilson, D. Hoth, and D. W. Barry. 1988. Survival experience among patients with AIDS receiving zidovudine. Follow-up of patients in a compassionate plea program. *JAMA*, 260: 3009–3015. doi: 10.1001/jama.1988.03410200065027
- [98] R. D. Moore, J. Keruly, D. D. Richman, T. Creagh-Kirk, and R. E. Chaisson. 1992. Natural history of advanced HIV disease in patients treated with zidovudine. *AIDS*, 6:671–677
- [99] J. O. Kahn, S. W. Lagakos, D. D. Richman, A. Cross, C. Pettinelli, S. H. Liou, M. Brown, P. A. Volberding, C. S. Crumpacker, and G. Beall. 1992. A controlled trial comparing continued zidovudine with didanosine in human immunodeficiency virus infection. *N Engl J Med*, 327:581–587. doi: 10.1056/NEJM199208273270901
- [100] D. I. Abrams, A. I. Goldman, C. Launer, J. A. Korvick, J. D. Neaton, L. R. Crane, M. Grodesky, S. Wakefield, K. Muth, and S. Kornegay. 1994. A comparative trial of didanosine or zalcitabine after treatment with zidovudine in patients with human immunodeficiency virus infection. *N Engl J Med*, 330:657–662. doi: 10.1056/NEJM199403103301001
- [101] M. D. de Jong, M. Loewenthal, C. A. Boucher, I. van der Ende, D. Hall, P. Schipper, A. Imrie, H. M. Weigel, R. H. Kauffmann, and R. Koster. 1994. Alternating

- nevirapine and zidovudine treatment of human immunodeficiency virus type 1-infected persons does not prolong nevirapine activity. *J Infect Dis*, 169:1346–1350. doi: 10.1093/infdis/169.6.1346
- [102] J. C. Schmit, J. Cogniaux, P. Hermans, C. Van Vaeck, S. Sprecher, B. Van Remoortel, M. Witvrouw, J. Balzarini, J. Desmyter, E. De Clercq, and A. M. Vandamme. 1996. Multiple drug resistance to nucleoside analogues and nonnucleoside reverse transcriptase inhibitors in an efficiently replicating human immunodeficiency virus type 1 patient strain. *J Infect Dis*, 174:962–968. doi: 10.1093/infdis/174.5.962
- [103] R. E. Dornsife, M. H. St Clair, A. T. Huang, T. J. Panella, G. W. Koszalka, C. L. Burns, and D. R. Averett. 1991. Anti-human immunodeficiency virus synergism by zidovudine (3'-azidothymidine) and didanosine (dideoxyinosine) contrasts with their additive inhibition of normal human marrow progenitor cells. *Antimicrob Agents Chemother*, 35:322–328. doi: 10.1128/AAC.35.2.322
- [104] V. A. Johnson, D. P. Merrill, J. A. Videler, T. C. Chou, R. E. Byington, J. J. Eron, R. T. D'Aquila, and M. S. Hirsch. 1991. Two-drug combinations of zidovudine, didanosine, and recombinant interferon-alpha A inhibit replication of zidovudine-resistant human immunodeficiency virus type 1 synergistically in vitro. *J Infect Dis*, 164:646–655. doi: 10.1093/infdis/164.4.646
- [105] S. W. Cox, K. Apéria, J. Albert, and B. Wahren. 1994. Comparison of the sensitivities of primary isolates of HIV type 2 and HIV type 1 to antiviral drugs and drug combinations. *AIDS Res Hum Retroviruses*, 10:1725–1729. doi: 10.1177/095632029300400407
- [106] J. Y. Feng, J. K. Ly, F. Myrick, D. Goodman, K. L. White, E. S. Svarovskaia, K. Borroto-Esoda, and M. D. Miller. 2009. The triple combination of tenofovir, emtricitabine and efavirenz shows synergistic anti-HIV-1 activity in vitro: a mechanism of action study. *Retrovirology*, 6:44. doi: 10.1186/1742-4690-6-44
- [107] B. L. Jilek, M. Zarr, M. E. Sampah, S. A. Rabi, C. K. Bullen, J. Lai, L. Shen, and R. F. Siliciano. 2012. A quantitative basis for antiretroviral therapy for HIV-1 infection. *Nat Med*, 18:446–451. doi: 10.1038/nm.2649
- [108] R. Kulkarni, R. Hluhanich, D. M. McColl, M. D. Miller, and K. L. White. 2014. The combined anti-HIV-1 activities of emtricitabine and tenofovir plus the integrase inhibitor elvitegravir or raltegravir show high levels of synergy in vitro. *Antimicrob Agents Chemother*, 58:6145–6150. doi: 10.1128/AAC.03591-14
- [109] Y. K. Chow, M. S. Hirsch, D. P. Merrill, L. J. Bechtel, J. J. Eron, J. C. Kaplan, and R. T. D'Aquila. 1993. Use of evolutionary limitations of HIV-1 multidrug resistance to optimize therapy. *Nature*, 361:650–654. doi: 10.1038/361650a0
- [110] B. A. Larder, S. D. Kemp, and P. R. Harrigan. 1995. Potential mechanism for

sustained antiretroviral efficacy of AZT-3TC combination therapy. *Science*, 269: 696–699. doi: 10.1126/science.7542804

- [111] A. C. Collier, R. W. Coombs, M. A. Fischl, P. R. Skolnik, D. Northfelt, P. Boutin, C. J. Hooper, L. D. Kaplan, P. A. Volberding, L. G. Davis, D. R. Henrard, S. Weller, and L. Corey. 1993. Combination therapy with zidovudine and didanosine compared with zidovudine alone in HIV-1 infection. *Ann Intern Med*, 119:786–793. doi: 10.7326/0003-4819-119-8-199310150-00003
- [112] J. J. Eron, S. L. Benoit, J. Jemsek, R. D. MacArthur, J. Santana, J. B. Quinn, D. R. Kuritzkes, M. A. Fallon, and M. Rubin. 1995. Treatment with lamivudine, zidovudine, or both in HIV-positive patients with 200 to 500 CD4+ cells per cubic millimeter. *N Engl J Med*, 333:1662–1669. doi: 10.1056/NEJM199512213332502
- [113] A. C. Collier, R. W. Coombs, D. A. Schoenfeld, R. L. Bassett, J. Timpone, A. Baruch, M. Jones, K. Facey, C. Whitacre, V. J. McAuliffe, H. M. Friedman, T. C. Merigan, R. C. Richman, C. Hooper, and L. Corey. 1996. Treatment of human immunodeficiency virus infection with saquinavir, zidovudine, and zalcitabine. *N Engl J Med*, 334:1011–1017. doi: 10.1056/NEJM199604183341602
- [114] S. M. Hammer, D. A. Katzenstein, M. D. Hughes, H. Gundacker, R. T. Schooley, R. H. Haubrich, W. K. Henry, M. M. Lederman, J. P. Phair, M. Niu, M. S. Hirsch, and T. C. Merigan. 1996. A trial comparing nucleoside monotherapy with combination therapy in HIV-infected adults with CD4 cell counts from 200 to 500 per cubic millimeter. *N Engl J Med*, 335:1081–1090. doi: 10.1056/NEJM199610103351501
- [115] L. D. Saravolatz, D. L. Winslow, G. Collins, J. S. Hodges, C. Pettinelli, D. S. Stein, N. Markowitz, R. Reves, M. O. Loveless, L. Crane, M. Thompson, and D. Abrams. 1996. Zidovudine alone or in combination with didanosine or zalcitabine in HIV-infected patients with the acquired immunodeficiency syndrome or fewer than 200 CD4 cells per cubic millimeter. *N Engl J Med*, 335:1099–1106. doi: 10.1056/NEJM199610103351503
- [116] J. Derbyshire, Delta Coordinating Committee, et al. 1996. Delta: a randomised double-blind controlled trial comparing combinations of zidovudine plus didanosine or zalcitabine with zidovudine alone in HIV-infected individuals. *Lancet*, 348:283–291. doi: 10.1016/S0140-6736(96)05387-1
- [117] S. M. Hammer, K. E. Squires, M. D. Hughes, J. M. Grimes, L. M. Demeter, J. S. Currier, J. Eron, Jr, J. E. Feinberg, H. Balfour, Jr, L. R. Deyton, J. A. Chodakewitz, and M. A. Fischl. 1997. A controlled trial of two nucleoside analogues plus indinavir in persons with human immunodeficiency virus infection and CD4 cell counts of 200 per cubic millimeter or less. *N Engl J Med*, 337:725–733. doi: 10.1056/NEJM199709113371101
- [118] R. M. Gulick, J. W. Mellors, D. Havlir, J. J. Eron, C. Gonzalez, D. McMahon, D. D. Richman, F. T. Valentine, L. Jonas, A. Meibohm, E. A. Emini, and J. A.

- Chodakewitz. 1997. Treatment with indinavir, zidovudine, and lamivudine in adults with human immunodeficiency virus infection and prior antiretroviral therapy. *N Engl J Med*, 337:734–739. doi: 10.1056/NEJM199709113371102
- [119] R. D. Moore and R. E. Chaisson. 1999. Natural history of HIV infection in the era of combination antiretroviral therapy. *AIDS*, 13:1933–1942
- [120] Antiretroviral Therapy Cohort Collaboration. 2008. Life expectancy of individuals on combination antiretroviral therapy in high-income countries: a collaborative analysis of 14 cohort studies. *Lancet*, 372:293–299. doi: 10.1016/S0140-6736(08)61113-7
- [121] O. Keiser, P. Taffé, M. Zwahlen, M. Battegay, E. Bernasconi, R. Weber, M. Rickenbach, and S. H. I. V. C. S. . 2004. All cause mortality in the Swiss HIV Cohort Study from 1990 to 2001 in comparison with the Swiss population. *AIDS*, 18: 1835–1843
- [122] A. I. van Sighem, L. A. J. Gras, P. Reiss, K. Brinkman, F. de Wolf, and A. T. H. E. N. A. n. o. c. s. . 2010. Life expectancy of recently diagnosed asymptomatic HIV-infected patients approaches that of uninfected individuals. *AIDS*, 24:1527–1535. doi: 10.1097/QAD.0b013e32833a3946
- [123] F. Nakagawa, M. May, and A. Phillips. 2013. Life expectancy living with HIV: recent estimates and future implications. *Curr Opin Infect Dis*, 26:17–25. doi: 10.1097/QCO.0b013e32835ba6b1
- [124] L. F. Johnson, J. Mossong, R. E. Dorrington, M. Schomaker, C. J. Hoffmann, O. Keiser, M. P. Fox, R. Wood, H. Prozesky, J. Giddy, D. B. Garone, M. Cornell, M. Egger, A. Boulle, and I. E. D. t. E. A. I. D. S. S. A. C. . 2013. Life expectancies of South African adults starting antiretroviral treatment: collaborative analysis of cohort studies. *PLoS Med*, 10:e1001418. doi: 10.1371/journal.pmed.1001418
- [125] S. G. Deeks, S. R. Lewin, and D. V. Havlir. 2013. The end of AIDS: HIV infection as a chronic disease. *Lancet*, 382:1525–1533. doi: 10.1016/S0140-6736(13)61809-7
- [126] A. R. Cillo, M. D. Sobolewski, R. J. Bosch, E. Fyne, M. Piatak, Jr, J. M. Coffin, and J. W. Mellors. 2014. Quantification of HIV-1 latency reversal in resting CD4+ T cells from patients on suppressive antiretroviral therapy. *Proc Natl Acad Sci U S A*, 111:7078–7083. doi: 10.1073/pnas.1402873111
- [127] T. W. Chun, D. Engel, M. M. Berrey, T. Shea, L. Corey, and A. S. Fauci. 1998. Early establishment of a pool of latently infected, resting CD4(+) T cells during primary HIV-1 infection. *Proc Natl Acad Sci U S A*, 95:8869–8873
- [128] J. B. Whitney, A. L. Hill, S. Sanisetty, P. Penaloza-MacMaster, J. Liu, M. Shetty, L. Parenteau, C. Cabral, J. Shields et al. 2014. Rapid seeding of the viral

- reservoir prior to SIV viraemia in rhesus monkeys. *Nature*, 512:74–77. doi: 10.1038/nature13594
- [129] X. Contreras, M. Schweneker, C.-S. Chen, J. M. McCune, S. G. Deeks, J. Martin, and B. M. Peterlin. 2009. Suberoylanilide hydroxamic acid reactivates HIV from latently infected cells. *J Biol Chem*, 284:6782–6789. doi: 10.1074/jbc.M807898200
- [130] S. Xing, C. K. Bullen, N. S. Shroff, L. Shan, H.-C. Yang, J. L. Manucci, S. Bhat, H. Zhang, J. B. Margolick, T. C. Quinn, D. M. Margolis, J. D. Siliciano, and R. F. Siliciano. 2011. Disulfiram reactivates latent HIV-1 in a Bcl-2-transduced primary CD4+ T cell model without inducing global T cell activation. *J Virol*, 85: 6060–6064. doi: 10.1128/JVI.02033-10
- [131] T. W. Chun, L. Carruth, D. Finzi, X. Shen, J. A. DiGiuseppe, H. Taylor, M. Hermankova, K. Chadwick, J. Margolick, T. C. Quinn, Y. H. Kuo, R. Brookmeyer, M. A. Zeiger, P. Barditch-Crovo, and R. F. Siliciano. 1997. Quantification of latent tissue reservoirs and total body viral load in HIV-1 infection. *Nature*, 387: 183–188. doi: 10.1038/387183a0
- [132] D. Baltimore. 1970. RNA-dependent DNA polymerase in virions of RNA tumour viruses. *Nature*, 226:1209–1211. doi: 10.1038/2261209a0
- [133] H. M. Temin and S. Mizutani. 1970. RNA-dependent DNA polymerase in virions of Rous sarcoma virus. *Nature*, 226:1211–1213. doi: 10.1038/2261211a0
- [134] D. P. Grandgenett, A. C. Vora, and R. D. Schiff. 1978. A 32,000-dalton nucleic acid-binding protein from avian retravirus cores possesses DNA endonuclease activity. *Virology*, 89:119–132. doi: 10.1016/0042-6822(78)90046-6
- [135] A. T. Panganiban and H. M. Temin. 1984. The retrovirus pol gene encodes a product required for DNA integration: identification of a retrovirus int locus. *Proc Natl Acad Sci U S A*, 81:7885–7889
- [136] P. Schwartzberg, J. Colicelli, and S. P. Goff. 1984. Construction and analysis of deletion mutations in the pol gene of Moloney murine leukemia virus: a new viral function required for productive infection. *Cell*, 37:1043–1052. doi: 10.1016/0092-8674(84)90439-2
- [137] L. A. Donehower and H. E. Varmus. 1984. A mutant murine leukemia virus with a single missense codon in pol is defective in a function affecting integration. *Proc Natl Acad Sci U S A*, 81:6461–6465
- [138] A. T. Panganiban and H. M. Temin. 1983. The terminal nucleotides of retrovirus DNA are required for integration but not virus production. *Nature*, 306:155–160. doi: 10.1038/306155a0
- [139] F. D. Veronese, R. Rahman, T. D. Copeland, S. Oroszlan, R. C. Gallo, and M. G.

- Sarngadharan. 1987. Immunological and chemical analysis of P6, the carboxyl-terminal fragment of HIV P15. *AIDS Res Hum Retroviruses*, 3:253–264. doi: 10.1089/aid.1987.3.253
- [140] H. G. Göttlinger, T. Dorfman, J. G. Sodroski, and W. A. Haseltine. 1991. Effect of mutations affecting the p6 gag protein on human immunodeficiency virus particle release. *Proc Natl Acad Sci U S A*, 88:3195–3199. doi: 10.1073/pnas.88.8.3195
- [141] B. Strack, A. Calistri, S. Craig, E. Popova, and H. G. Göttlinger. 2003. AIP1/ALIX is a binding partner for HIV-1 p6 and EIAV p9 functioning in virus budding. *Cell*, 114:689–699. doi: 10.1016/S0092-8674(03)00653-6
- [142] B. Crise, L. Buonocore, and J. K. Rose. 1990. CD4 is retained in the endoplasmic reticulum by the human immunodeficiency virus type 1 glycoprotein precursor. *J Virol*, 64:5585–5593
- [143] S. Bour, F. Boulerice, and M. A. Wainberg. 1991. Inhibition of gp160 and CD4 maturation in U937 cells after both defective and productive infections by human immunodeficiency virus type 1. *J Virol*, 65:6387–6396
- [144] J. Sodroski, C. Rosen, F. Wong-Staal, S. Z. Salahuddin, M. Popovic, S. Arya, R. C. Gallo, and W. A. Haseltine. 1985. Trans-acting transcriptional regulation of human T-cell leukemia virus type III long terminal repeat. *Science*, 227:171–173. doi: 10.1126/science.2981427
- [145] J. Sodroski, R. Patarca, C. Rosen, F. Wong-Staal, and W. Haseltine. 1985. Location of the trans-activating region on the genome of human T-cell lymphotropic virus type III. *Science*, 229:74–77. doi: 10.1126/science.2990041
- [146] T. K. Howcroft, K. Strelbel, M. A. Martin, and D. S. Singer. 1993. Repression of MHC class I gene promoter activity by two-exon Tat of HIV. *Science*, 260: 1320–1322. doi: 10.1126/science.8493575
- [147] B. E. Meyer and M. H. Malim. 1994. The HIV-1 Rev trans-activator shuttles between the nucleus and the cytoplasm. *Genes Dev*, 8:1538–1547. doi: 10.1101/gad.8.13.1538
- [148] J. Sodroski, W. C. Goh, C. Rosen, A. Dayton, E. Terwilliger, and W. Haseltine. 1986. A second post-transcriptional trans-activator gene required for HTLV-III replication. *Nature*, 321:412–417. doi: 10.1038/321412a0
- [149] M. B. Feinberg, R. F. Jarrett, A. Aldovini, R. C. Gallo, and F. Wong-Staal. 1986. HTLV-III expression and production involve complex regulation at the levels of splicing and translation of viral RNA. *Cell*, 46:807–817. doi: 10.1016/0092-8674(86)90062-0
- [150] D. M. Knight, F. A. Flomerfelt, and J. Ghrayeb. 1987. Expression of the art/trs

- protein of HIV and study of its role in viral envelope synthesis. *Science*, 236: 837–840. doi: 10.1126/science.3033827
- [151] M. H. Malim, J. Hauber, R. Fenrick, and B. R. Cullen. 1988. Immunodeficiency virus rev trans-activator modulates the expression of the viral regulatory genes. *Nature*, 335:181–183. doi: 10.1038/335181a0
 - [152] D. Gutman and C. J. Goldenberg. 1988. Virus-specific splicing inhibitor in extracts from cells infected with HIV-1. *Science*, 241:1492–1495. doi: 10.1126/science.3047873
 - [153] M. H. Malim, J. Hauber, S. Y. Le, J. V. Maizel, and B. R. Cullen. 1989. The HIV-1 rev trans-activator acts through a structured target sequence to activate nuclear export of unspliced viral mRNA. *Nature*, 338:254–257. doi: 10.1038/338254a0
 - [154] M. H. Malim, S. Bhnlein, J. Hauber, and B. R. Cullen. 1989. Functional dissection of the HIV-1 Rev trans-activator—derivation of a trans-dominant repressor of Rev function. *Cell*, 58:205–214. doi: 10.1016/0092-8674(89)90416-9
 - [155] G. Yu and R. L. Felsted. 1992. Effect of myristoylation on p27 nef subcellular distribution and suppression of HIV-LTR transcription. *Virology*, 187:46–55. doi: 10.1016/0042-6822(92)90293-X
 - [156] J. V. Garcia and A. D. Miller. 1991. Serine phosphorylation-independent downregulation of cell-surface CD4 by nef. *Nature*, 350:508–511. doi: 10.1038/350508a0
 - [157] R. E. Benson, A. Sanfridson, J. S. Ottinger, C. Doyle, and B. R. Cullen. 1993. Downregulation of cell-surface CD4 expression by simian immunodeficiency virus Nef prevents viral super infection. *J Exp Med*, 177:1561–1566. doi: 10.1084/jem.177.6.1561
 - [158] C. Aiken, J. Konner, N. R. Landau, M. E. Lenburg, and D. Trono. 1994. Nef induces CD4 endocytosis: requirement for a critical dileucine motif in the membrane-proximal CD4 cytoplasmic domain. *Cell*, 76:853–864. doi: 10.1016/0092-8674(94)90360-3
 - [159] J. Lama, A. Mangasarian, and D. Trono. 1999. Cell-surface expression of CD4 reduces HIV-1 infectivity by blocking Env incorporation in a Nef- and Vpu-inhibitable manner. *Curr Biol*, 9:622–631. doi: 10.1016/S0960-9822(99)80284-X
 - [160] T. M. Ross, A. E. Oran, and B. R. Cullen. 1999. Inhibition of HIV-1 progeny virion release by cell-surface CD4 is relieved by expression of the viral Nef protein. *Curr Biol*, 9:613–621. doi: 10.1016/S0960-9822(99)80283-8
 - [161] N. Michel, I. Allespach, S. Venzke, O. T. Fackler, and O. T. Keppler. 2005. The Nef protein of human immunodeficiency virus establishes superinfection immunity

- by a dual strategy to downregulate cell-surface CCR5 and CD4. *Curr Biol*, 15: 714–723. doi: 10.1016/j.cub.2005.02.058
- [162] O. Schwartz, V. Maréchal, S. Le Gall, F. Lemonnier, and J. M. Heard. 1996. Endocytosis of major histocompatibility complex class I molecules is induced by the HIV-1 Nef protein. *Nat Med*, 2:338–342. doi: 10.1038/nm0396-338
- [163] K. L. Collins, B. K. Chen, S. A. Kalams, B. D. Walker, and D. Baltimore. 1998. HIV-1 Nef protein protects infected primary cells against killing by cytotoxic T lymphocytes. *Nature*, 391:397–401. doi: 10.1038/34929
- [164] P. Stumptner-Cuvelette, S. Morchoisne, M. Dugast, S. Le Gall, G. Raposo, O. Schwartz, and P. Benaroch. 2001. HIV-1 Nef impairs MHC class II antigen presentation and surface expression. *Proc Natl Acad Sci U S A*, 98:12144–12149. doi: 10.1073/pnas.221256498
- [165] A. D. Blagoveshchenskaya, L. Thomas, S. F. Feliciangeli, C. H. Hung, and G. Thomas. 2002. HIV-1 Nef downregulates MHC-I by a PACS-1- and PI3K-regulated ARF6 endocytic pathway. *Cell*, 111:853–866. doi: 10.1016/S0092-8674(02)01162-5
- [166] X. N. Xu, B. Laffert, G. R. Screaton, M. Kraft, D. Wolf, W. Kolanus, J. Mongkolsapay, A. J. McMichael, and A. S. Baur. 1999. Induction of Fas ligand expression by HIV involves the interaction of Nef with the T cell receptor zeta chain. *J Exp Med*, 189:1489–1496. doi: 10.1084/jem.189.9.1489
- [167] J. A. Schrager and J. W. Marsh. 1999. HIV-1 Nef increases T cell activation in a stimulus-dependent manner. *Proc Natl Acad Sci U S A*, 96:8167–8172. doi: 10.1073/pnas.96.14.8167
- [168] J. K. Wang, E. Kiyokawa, E. Verdin, and D. Trono. 2000. The Nef protein of HIV-1 associates with rafts and primes T cells for activation. *Proc Natl Acad Sci U S A*, 97:394–399. doi: 10.1073/pnas.97.1.394
- [169] A. Simmons, V. Aluvihare, and A. McMichael. 2001. Nef triggers a transcriptional program in T cells imitating single-signal T cell activation and inducing HIV virulence mediators. *Immunity*, 14:763–777. doi: 10.1016/S1074-7613(01)00158-3
- [170] J. A. Schrager, V. Der Minassian, and J. W. Marsh. 2002. HIV Nef increases T cell ERK MAP kinase activity. *J Biol Chem*, 277:6137–6142. doi: 10.1074/jbc.M107322200
- [171] M. Schindler, J. Münch, O. Kutsch, H. Li, M. L. Santiago, F. Bibollet-Ruche, M. C. Müller-Trutwin, F. J. Novembre, M. Peeters, V. Courgnaud, E. Bailes, P. Roques, D. L. Sodora, G. Silvestri, P. M. Sharp, B. H. Hahn, and F. Kirchhoff. 2006. Nef-mediated suppression of T cell activation was lost in a lentiviral lineage that gave rise to HIV-1. *Cell*, 125:1055–1067. doi: 10.1016/j.cell.2006.04.033

- [172] F. Kirchhoff, M. Schindler, A. Specht, N. Arhel, and J. Münch. 2008. Role of Nef in primate lentiviral immunopathogenesis. *Cell Mol Life Sci*, 65:2621–2636. doi: 10.1007/s00018-008-8094-2
- [173] F. Kirchhoff. 2009. Is the high virulence of HIV-1 an unfortunate coincidence of primate lentiviral evolution? *Nat Rev Microbiol*, 7:467–476. doi: 10.1038/nrmicro2111
- [174] A. M. Sheehy, N. C. Gaddis, J. D. Choi, and M. H. Malim. 2002. Isolation of a human gene that inhibits HIV-1 infection and is suppressed by the viral Vif protein. *Nature*, 418:646–650. doi: 10.1038/nature00939
- [175] R. Mariani, D. Chen, B. Schröfelbauer, F. Navarro, R. König, B. Bollman, C. Münk, H. Nymark-McMahon, and N. R. Landau. 2003. Species-specific exclusion of APOBEC3G from HIV-1 virions by Vif. *Cell*, 114:21–31. doi: 10.1016/S0092-8674(03)00515-4
- [176] A. M. Sheehy, N. C. Gaddis, and M. H. Malim. 2003. The antiretroviral enzyme APOBEC3G is degraded by the proteasome in response to HIV-1 Vif. *Nat Med*, 9: 1404–1407. doi: 10.1038/nm945
- [177] M. Marin, K. M. Rose, S. L. Kozak, and D. Kabat. 2003. HIV-1 Vif protein binds the editing enzyme APOBEC3G and induces its degradation. *Nat Med*, 9: 1398–1403. doi: 10.1038/nm946
- [178] X. Yu, Y. Yu, B. Liu, K. Luo, W. Kong, P. Mao, and X.-F. Yu. 2003. Induction of APOBEC3G ubiquitination and degradation by an HIV-1 Vif-Cul5-SCF complex. *Science*, 302:1056–1060. doi: 10.1126/science.1089591
- [179] R. S. Harris, K. N. Bishop, A. M. Sheehy, H. M. Craig, S. K. Petersen-Mahrt, I. N. Watt, M. S. Neuberger, and M. H. Malim. 2003. DNA deamination mediates innate immunity to retroviral infection. *Cell*, 113:803–809. doi: 10.1016/S0092-8674(03)00423-9
- [180] B. Mangeat, P. Turelli, G. Caron, M. Friedli, L. Perrin, and D. Trono. 2003. Broad antiretroviral defence by human APOBEC3G through lethal editing of nascent reverse transcripts. *Nature*, 424:99–103. doi: 10.1038/nature01709
- [181] H. Zhang, B. Yang, R. J. Pomerantz, C. Zhang, S. C. Arunachalam, and L. Gao. 2003. The cytidine deaminase CEM15 induces hypermutation in newly synthesized HIV-1 DNA. *Nature*, 424:94–98. doi: 10.1038/nature01707
- [182] D. Lecossier, F. Bouchonnet, F. Clavel, and A. J. Hance. 2003. Hypermutation of HIV-1 DNA in the absence of the Vif protein. *Science*, 300:1112. doi: 10.1126/science.1083338
- [183] E. A. Cohen, E. F. Terwilliger, J. G. Sodroski, and W. A. Haseltine. 1988. Identifi-

- cation of a protein encoded by the vpu gene of HIV-1. *Nature*, 334:532–534. doi: 10.1038/334532a0
- [184] K. Strelbel, T. Klimkait, and M. A. Martin. 1988. A novel gene of HIV-1, vpu, and its 16-kilodalton product. *Science*, 241:1221–1223. doi: 10.1126/science.3261888
- [185] R. L. Willey, F. Maldarelli, M. A. Martin, and K. Strelbel. 1992. Human immunodeficiency virus type 1 Vpu protein induces rapid degradation of CD4. *J Virol*, 66: 7193–7200
- [186] S. Bour, U. Schubert, and K. Strelbel. 1995. The human immunodeficiency virus type 1 Vpu protein specifically binds to the cytoplasmic domain of CD4: implications for the mechanism of degradation. *J Virol*, 69:1510–1520
- [187] W. L. Marshall, D. C. Diamond, M. M. Kowalski, and R. W. Finberg. 1992. High level of surface CD4 prevents stable human immunodeficiency virus infection of T-cell transfectants. *J Virol*, 66:5492–5499
- [188] M. J. Cortés, F. Wong-Staal, and J. Lama. 2002. Cell surface CD4 interferes with the infectivity of HIV-1 particles released from T cells. *J Biol Chem*, 277: 1770–1779. doi: 10.1074/jbc.M109807200
- [189] S. J. D. Neil, T. Zang, and P. D. Bieniasz. 2008. Tetherin inhibits retrovirus release and is antagonized by HIV-1 Vpu. *Nature*, 451:425–430. doi: 10.1038/nature06553
- [190] K. Strelbel, T. Klimkait, F. Maldarelli, and M. A. Martin. 1989. Molecular and biochemical analyses of human immunodeficiency virus type 1 vpu protein. *J Virol*, 63:3784–3791
- [191] S. H. Hughes, P. R. Shank, D. H. Spector, H. J. Kung, J. M. Bishop, H. E. Varmus, P. K. Vogt, and M. L. Breitman. 1978. Proviruses of avian sarcoma virus are terminally redundant, co-extensive with unintegrated linear DNA and integrated at many sites. *Cell*, 15:1397–1410. doi: 10.1016/0092-8674(78)90064-8
- [192] S. Wain-Hobson, P. Sonigo, O. Danos, S. Cole, and M. Alizon. 1985. Nucleotide sequence of the AIDS virus, LAV. *Cell*, 40:9–17. doi: 10.1016/0092-8674(85)90303-4
- [193] M. A. Muesing, D. H. Smith, C. D. Cabradilla, C. V. Benton, L. A. Lasky, and D. J. Capon. 1985. Nucleic acid structure and expression of the human AIDS/lymphadenopathy retrovirus. *Nature*, 313:450–458. doi: 10.1038/313450a0
- [194] L. Ratner, W. Haseltine, R. Patarca, K. J. Livak, B. Starcich, S. F. Josephs, E. R. Doran, J. A. Rafalski, E. A. Whitehorn, and K. Baumeister. 1985. Complete nucleotide sequence of the AIDS virus, HTLV-III. *Nature*, 313:277–284. doi: 10.1038/313277a0

- [195] R. Sanchez-Pescador, M. D. Power, P. J. Barr, K. S. Steimer, M. M. Stempien, S. L. Brown-Shimer, W. W. Gee, A. Renard, A. Randolph, and J. A. Levy. 1985. Nucleotide sequence and expression of an AIDS-associated retrovirus (ARV-2). *Science*, 227:484–492. doi: 10.1126/science.2578227
- [196] E. C. Holmes, L. Q. Zhang, P. Simmonds, C. A. Ludlam, and A. J. Brown. 1992. Convergent and divergent sequence evolution in the surface envelope glycoprotein of human immunodeficiency virus type 1 within a single infected patient. *Proc Natl Acad Sci U S A*, 89:4835–4839
- [197] S. Bonhoeffer, E. C. Holmes, and M. A. Nowak. 1995. Causes of HIV diversity. *Nature*, 376:125. doi: 10.1038/376125a0
- [198] H. A. Ross and A. G. Rodrigo. 2002. Immune-mediated positive selection drives human immunodeficiency virus type 1 molecular variation and predicts disease duration. *J Virol*, 76:11715–11720. doi: 10.1128/JVI.76.22.11715-11720.2002
- [199] S. M. Wolinsky, B. T. Korber, A. U. Neumann, M. Daniels, K. J. Kunstman, A. J. Whetsell, M. R. Furtado, Y. Cao, D. D. Ho, and J. T. Safrit. 1996. Adaptive evolution of human immunodeficiency virus-type 1 during the natural course of infection. *Science*, 272:537–542. doi: 10.1126/science.272.5261.537
- [200] R. Shankarappa, J. B. Margolick, S. J. Gange, A. G. Rodrigo, D. Upchurch, H. Farzadegan, P. Gupta, C. R. Rinaldo, G. H. Learn, X. He, X. L. Huang, and J. I. Mullins. 1999. Consistent viral evolutionary changes associated with the progression of human immunodeficiency virus type 1 infection. *J Virol*, 73: 10489–10502
- [201] D. D. Ho, A. U. Neumann, A. S. Perelson, W. Chen, J. M. Leonard, and M. Markowitz. 1995. Rapid turnover of plasma virions and CD4 lymphocytes in HIV-1 infection. *Nature*, 373:123–126. doi: 10.1038/373123a0
- [202] X. Wei, S. K. Ghosh, M. E. Taylor, V. A. Johnson, E. A. Emini, P. Deutsch, J. D. Lifson, S. Bonhoeffer, M. A. Nowak, and B. H. Hahn. 1995. Viral dynamics in human immunodeficiency virus type 1 infection. *Nature*, 373:117–122. doi: 10.1038/373117a0
- [203] A. S. Perelson, P. Essunger, Y. Cao, M. Vesalanen, A. Hurley, K. Saksela, M. Markowitz, and D. D. Ho. 1997. Decay characteristics of HIV-1-infected compartments during combination therapy. *Nature*, 387:188–191. doi: 10.1038/387188a0
- [204] D. Finzi, M. Hermankova, T. Pierson, L. M. Carruth, C. Buck, R. E. Chaisson, T. C. Quinn, K. Chadwick, J. Margolick, R. Brookmeyer, J. Gallant, M. Markowitz, D. D. Ho, D. D. Richman, and R. F. Siliciano. 1997. Identification of a reservoir for HIV-1 in patients on highly active antiretroviral therapy. *Science*, 278:1295–1300. doi: 10.1126/science.278.5341.1295

- [205] J. A. Hoxie, B. S. Haggarty, J. L. Rackowski, N. Pillsbury, and J. A. Levy. 1985. Persistent noncytopathic infection of normal human T lymphocytes with AIDS-associated retrovirus. *Science*, 229:1400–1402. doi: 10.1126/science.2994222
- [206] T. Folks, D. M. Powell, M. M. Lightfoote, S. Benn, M. A. Martin, and A. S. Fauci. 1986. Induction of HTLV-III/LAV from a nonvirus-producing T-cell line: implications for latency. *Science*, 231:600–602. doi: 10.1126/science.3003906
- [207] S. M. Berget, C. Moore, and P. A. Sharp. 1977. Spliced segments at the 5' terminus of adenovirus 2 late mRNA. *Proc Natl Acad Sci USA*, 74:3171–3175
- [208] L. T. Chow, R. E. Gelinas, T. R. Broker, and R. J. Roberts. 1977. An amazing sequence arrangement at the 5' ends of adenovirus 2 messenger RNA. *Cell*, 12: 1–8. doi: 10.1016/0092-8674(77)90180-5
- [209] N. Lohse, A.-B. E. Hansen, G. Pedersen, G. Kronborg, J. Gerstoft, H. T. Sørensen, M. Vaeth, and N. Obel. 2007. Survival of persons with and without HIV infection in Denmark, 1995–2005. *Ann Intern Med*, 146:87–95
- [210] D. D. Richman, D. M. Margolis, M. Delaney, W. C. Greene, D. Hazuda, and R. J. Pomerantz. 2009. The challenge of finding a cure for HIV infection. *Science*, 323: 1304–1307. doi: 10.1126/science.1165706
- [211] A. B. Hutchinson, P. G. Farnham, H. D. Dean, D. U. Ekwueme, C. del Rio, L. Kamimoto, and S. E. Kellerman. 2006. The economic burden of HIV in the United States in the era of highly active antiretroviral therapy: evidence of continuing racial and ethnic differences. *J Acquir Immune Defic Syndr*, 43:451–457. doi: 10.1097/01.qai.0000243090.32866.4e
- [212] B. R. Schackman, K. A. Gebo, R. P. Walensky, E. Losina, T. Muccio, P. E. Sax, M. C. Weinstein, G. R. Seage, R. D. Moore, and K. A. Freedberg. 2006. The lifetime cost of current human immunodeficiency virus care in the United States. *Med Care*, 44:990–997. doi: 10.1097/01.mlr.0000228021.89490.2a
- [213] T. Fukuhara, T. Hosoya, S. Shimizu, K. Sumi, T. Oshiro, Y. Yoshinaka, M. Suzuki, N. Yamamoto, L. A. Herzenberg, L. A. Herzenberg, and M. Hagiwara. 2006. Utilization of host SR protein kinases and RNA-splicing machinery during viral replication. *Proc Natl Acad Sci USA*, 103:11329–11333. doi: 10.1073/pnas.0604616103
- [214] D. Mandal, Z. Feng, and C. M. Stoltzfus. 2010. Excessive RNA splicing and inhibition of HIV-1 replication induced by modified U1 small nuclear RNAs. *J Virol*, 84:12790–12800. doi: 10.1128/JVI.01257-10
- [215] Y.-H. Zheng, H.-F. Yu, and B. M. Peterlin. 2003. Human p32 protein relieves a post-transcriptional block to HIV replication in murine cells. *Nat Cell Biol*, 5: 611–618. doi: 10.1038/ncb1000

- [216] C. Gélinas and H. M. Temin. 1986. Nondefective spleen necrosis virus-derived vectors define the upper size limit for packaging reticuloendotheliosis viruses. *Proc Natl Acad Sci USA*, 83:9211–9215
- [217] S. A. Herman and J. M. Coffin. 1987. Efficient packaging of readthrough RNA in ALV: implications for oncogene transduction. *Science*, 236:845–848. doi: 10.1126/science.3033828
- [218] N. H. Shin, D. Hartigan-O'Connor, J. K. Pfeiffer, and A. Telesnitsky. 2000. Replication of lengthened Moloney murine leukemia virus genomes is impaired at multiple stages. *J Virol*, 74:2694–2702
- [219] C. M. Stoltzfus. 2009. Chapter 1. Regulation of HIV-1 alternative RNA splicing and its role in virus replication. *Adv Virus Res*, 74:1–40. doi: 10.1016/S0065-3527(09)74001-1
- [220] S. Y. Kim, R. Byrn, J. Groopman, and D. Baltimore. 1989. Temporal aspects of DNA and RNA synthesis during human immunodeficiency virus infection: evidence for differential gene expression. *J Virol*, 63:3708–3713
- [221] R. J. Pomerantz, D. Trono, M. B. Feinberg, and D. Baltimore. 1990. Cells nonproductively infected with HIV-1 exhibit an aberrant pattern of viral RNA expression: a molecular model for latency. *Cell*, 61:1271–1276. doi: 10.1016/0092-8674(90)90691-7
- [222] A. L. Brass, D. M. Dykxhoorn, Y. Benita, N. Yan, A. Engelman, R. J. Xavier, J. Lieberman, and S. J. Elledge. 2008. Identification of host proteins required for HIV infection through a functional genomic screen. *Science*, 319:921–926. doi: 10.1126/science.1152725
- [223] R. König, Y. Zhou, D. Elleder, T. L. Diamond, G. M. C. Bonamy, J. T. Irelan, C.-Y. Chiang, B. P. Tu, P. D. D. Jesus et al. 2008. Global analysis of host-pathogen interactions that regulate early-stage HIV-1 replication. *Cell*, 135:49–60. doi: 10.1016/j.cell.2008.07.032
- [224] F. D. Bushman, N. Malani, J. Fernandes, I. D'Orso, G. Cagney, T. L. Diamond, H. Zhou, D. J. Hazuda, A. S. Espeseth, R. Knig, S. Bandyopadhyay, T. Ideker, S. P. Goff, N. J. Krogan, A. D. Frankel, J. A. T. Young, and S. K. Chanda. 2009. Host cell factors in HIV replication: meta-analysis of genome-wide studies. *PLoS Pathog*, 5:e1000437. doi: 10.1371/journal.ppat.1000437
- [225] S. Jäger, P. Cimermancic, N. Gulbahce, J. R. Johnson, K. E. McGovern, S. C. Clarke, M. Shales, G. Mercenne, L. Pache et al. 2012. Global landscape of HIV-human protein complexes. *Nature*, 481:365–370. doi: 10.1038/nature10719
- [226] A. Bansal, J. Carlson, J. Yan, O. T. Akinsiku, M. Schaefer, S. Sabbaj, A. Bet, D. N. Levy, S. Heath, J. Tang, R. A. Kaslow, B. D. Walker, T. Ndung'u, P. J.

- Goulder, D. Heckerman, E. Hunter, and P. A. Goepfert. 2010. CD8 T cell response and evolutionary pressure to HIV-1 cryptic epitopes derived from antisense transcription. *J Exp Med*, 207:51–59. doi: 10.1084/jem.20092060
- [227] N. Bakkour, Y.-L. Lin, S. Maire, L. Ayadi, F. Mahuteau-Betzer, C. H. Nguyen, C. Mettling, P. Portales, D. Grierson, B. Chabot, P. Jeanteur, C. Branlant, P. Corbeau, and J. Tazi. 2007. Small-molecule inhibition of HIV pre-mRNA splicing as a novel antiretroviral therapy to overcome drug resistance. *PLoS Pathog*, 3: 1530–1539. doi: 10.1371/journal.ppat.0030159
- [228] M. B. Asparuhova, G. Marti, S. Liu, F. Serhan, D. Trono, and D. Schmperli. 2007. Inhibition of HIV-1 multiplication by a modified U7 snRNA inducing Tat and Rev exon skipping. *J Gene Med*, 9:323–334. doi: 10.1002/jgm.1027
- [229] T. O. Tange, T. H. Jensen, and J. Kjems. 1996. In vitro interaction between human immunodeficiency virus type 1 Rev protein and splicing factor ASF/SF2-associated protein, p32. *J Biol Chem*, 271:10066–10072. doi: 10.1074/jbc.271.17.10066
- [230] R. Berro, K. Kehn, C. de la Fuente, A. Pumfery, R. Adair, J. Wade, A. M. Colberg-Poley, J. Hiscott, and F. Kashanchi. 2006. Acetylated Tat regulates human immunodeficiency virus type 1 splicing through its interaction with the splicing regulator p32. *J Virol*, 80:3189–3204. doi: 10.1128/JVI.80.7.3189-3204.2006
- [231] J. Bohne, A. Schambach, and D. Zychlinski. 2007. New way of regulating alternative splicing in retroviruses: the promoter makes a difference. *J Virol*, 81: 3652–3656. doi: 10.1128/JVI.02105-06
- [232] J. A. Jablonski, A. L. Amelio, M. Giacca, and M. Caputi. 2010. The transcriptional transactivator Tat selectively regulates viral splicing. *Nucleic Acids Res*, 38:1249–1260. doi: 10.1093/nar/gkp1105
- [233] M. Kuramitsu, C. Hashizume, N. Yamamoto, A. Azuma, M. Kamata, N. Yamamoto, Y. Tanaka, and Y. Aida. 2005. A novel role for Vpr of human immunodeficiency virus type 1 as a regulator of the splicing of cellular pre-mRNA. *Microbes Infect*, 7:1150–1160. doi: 10.1016/j.micinf.2005.03.022
- [234] C. Hashizume, M. Kuramitsu, X. Zhang, T. Kurosawa, M. Kamata, and Y. Aida. 2007. Human immunodeficiency virus type 1 Vpr interacts with spliceosomal protein SAP145 to mediate cellular pre-mRNA splicing inhibition. *Microbes Infect*, 9:490–497. doi: 10.1016/j.micinf.2007.01.013
- [235] N. M. Kopelman, D. Lancet, and I. Yanai. 2005. Alternative splicing and gene duplication are inversely correlated evolutionary mechanisms. *Nat Genet*, 37: 588–589. doi: 10.1038/ng1575
- [236] Y. Xing and C. Lee. 2005. Evidence of functional selection pressure for alternative

splicing events that accelerate evolution of protein subsequences. *Proc Natl Acad Sci USA*, 102:13526–13531. doi: 10.1073/pnas.0501213102

- [237] Z. Su, J. Wang, J. Yu, X. Huang, and X. Gu. 2006. Evolution of alternative splicing after gene duplication. *Genome Res*, 16:182–189. doi: 10.1101/gr.4197006
- [238] F. L. Watson, R. Pttmann-Holgado, F. Thomas, D. L. Lamar, M. Hughes, M. Kondo, V. I. Rebel, and D. Schmucker. 2005. Extensive diversity of Ig-superfamily proteins in the immune system of insects. *Science*, 309:1874–1878. doi: 10.1126/science.1116887
- [239] V. W. Pollard and M. H. Malim. 1998. The HIV-1 Rev protein. *Annu Rev Microbiol*, 52:491–532. doi: 10.1146/annurev.micro.52.1.491
- [240] D. F. Purcell and M. A. Martin. 1993. Alternative splicing of human immunodeficiency virus type 1 mRNA modulates viral protein expression, replication, and infectivity. *J Virol*, 67:6365–6378
- [241] M. E. Klotman, S. Kim, A. Buchbinder, A. DeRossi, D. Baltimore, and F. Wong-Staal. 1991. Kinetics of expression of multiply spliced RNA in early human immunodeficiency virus type 1 infection of lymphocytes and monocytes. *Proc Natl Acad Sci USA*, 88:5011–5015
- [242] D. M. Benko, S. Schwartz, G. N. Pavlakis, and B. K. Felber. 1990. A novel human immunodeficiency virus type 1 protein, tev, shares sequences with tat, env, and rev proteins. *J Virol*, 64:2505–2518
- [243] K. Fujita, J. Silver, and K. Peden. 1992. Changes in both gp120 and gp41 can account for increased growth potential and expanded host range of human immunodeficiency virus type 1. *J Virol*, 66:4445–4451
- [244] R. M. McAllister, M. B. Gardner, A. E. Greene, C. Bradt, W. W. Nichols, and B. H. Landing. 1971. Cultivation in vitro of cells derived from a human osteosarcoma. *Cancer*, 27:397–402
- [245] T. Kwan, D. Benovoy, C. Dias, S. Gurd, D. Serre, H. Zuzan, T. A. Clark, A. Schweitzer, M. K. Staples, H. Wang, J. E. Blume, T. J. Hudson, R. Sladek, and J. Majewski. 2007. Heritability of alternative splicing in the human genome. *Genome Res*, 17:1210–1218. doi: 10.1101/gr.6281007
- [246] J. Hull, S. Campino, K. Rowlands, M.-S. Chan, R. R. Copley, M. S. Taylor, K. Rockett, G. Elvidge, B. Keating, J. Knight, and D. Kwiatkowski. 2007. Identification of common genetic variation that modulates alternative splicing. *PLoS Genet*, 3: e99. doi: 10.1371/journal.pgen.0030099
- [247] E. T. Wang, R. Sandberg, S. Luo, I. Khrebtukova, L. Zhang, C. Mayr, S. F.

- Kingsmore, G. P. Schroth, and C. B. Burge. 2008. Alternative isoform regulation in human tissue transcriptomes. *Nature*, 456:470–476. doi: 10.1038/nature07509
- [248] Y. Barash, J. A. Calarco, W. Gao, Q. Pan, X. Wang, O. Shai, B. J. Blencowe, and B. J. Frey. 2010. Deciphering the splicing code. *Nature*, 465:53–59. doi: 10.1038/nature09000
- [249] M. T. Vahey, M. E. Nau, L. L. Jagodzinski, J. Yalley-Ogunro, M. Taubman, N. L. Michael, and M. G. Lewis. 2002. Impact of viral infection on the gene expression profiles of proliferating normal human peripheral blood mononuclear cells infected with HIV type 1 RF. *AIDS Res Hum Retroviruses*, 18:179–192. doi: 10.1089/08892220252781239
- [250] A. B. van 't Wout, G. K. Lehrman, S. A. Mikheeva, G. C. O'Keeffe, M. G. Katze, R. E. Bumgarner, G. K. Geiss, and J. I. Mullins. 2003. Cellular gene expression upon human immunodeficiency virus type 1 infection of CD4(+)T-cell lines. *J Virol*, 77:1392–1402. doi: 10.1128/JVI.77.2.1392-1402.2003
- [251] R. Mitchell, C.-Y. Chiang, C. Berry, and F. Bushman. 2003. Global analysis of cellular transcription following infection with an HIV-based vector. *Mol Ther*, 8: 674–687. doi: 10.1016/S1525-0016(03)00215-6
- [252] M. Rotger, K. K. Dang, J. Fellay, E. L. Heinzen, S. Feng, P. Descombes, K. V. Shianna, D. Ge, H. F. Gnethard, D. B. Goldstein, A. Telenti, S. H. C. Study, and C. for HIV/AIDS Vaccine Immunology. 2010. Genome-wide mRNA expression correlates of viral control in CD4+ T-cells from HIV-1-infected individuals. *PLoS Pathog*, 6:e1000781. doi: 10.1371/journal.ppat.1000781
- [253] S. T. Chang, P. Sova, X. Peng, J. Weiss, G. L. Law, R. E. Palermo, and M. G. Katze. 2011. Next-generation sequencing reveals HIV-1-mediated suppression of T cell activation and RNA processing and regulation of noncoding RNA expression in a CD4+ T cell line. *MBio*, 2. doi: 10.1128/mBio.00134-11
- [254] R. Tewhey, J. B. Warner, M. Nakano, B. Libby, M. Medkova, P. H. David, S. K. Kotsopoulos, M. L. Samuels, J. B. Hutchison, J. W. Larson, E. J. Topol, M. P. Weiner, O. Harismendy, J. Olson, D. R. Link, and K. A. Frazer. 2009. Microdroplet-based PCR enrichment for large-scale targeted sequencing. *Nat Biotechnol*, 27: 1025–1031. doi: 10.1038/nbt.1583
- [255] J. C. Marioni, C. E. Mason, S. M. Mane, M. Stephens, and Y. Gilad. 2008. RNA-seq: an assessment of technical reproducibility and comparison with gene expression arrays. *Genome Res*, 18:1509–1517. doi: 10.1101/gr.079558.108
- [256] R. Morin, M. Bainbridge, A. Fejes, M. Hirst, M. Krzywinski, T. Pugh, H. McDonald, R. Varhol, S. Jones, and M. Marra. 2008. Profiling the HeLa S3 transcriptome using randomly primed cDNA and massively parallel short-read sequencing. *Biotechniques*, 45:81–94. doi: 10.2144/000112900

- [257] J. Eid, A. Fehr, J. Gray, K. Luong, J. Lyle, G. Otto, P. Peluso, D. Rank, P. Baybayan et al. 2009. Real-time DNA sequencing from single polymerase molecules. *Science*, 323:133–138. doi: 10.1126/science.1162986
- [258] M. Schaefer, M. Brown, W. Kilembe, S. Allen, Y. Guo, E. Hunter, and E. Paxinos. Single-molecule complete HIV-1 genome sequencing from 2 linked transmission pairs. In *Conference on Retroviruses and Opportunistic Infections*, 2012
- [259] E. Buratti and F. E. Baralle. 2004. Influence of RNA secondary structure on the pre-mRNA splicing process. *Mol Cell Biol*, 24:10505–10514. doi: 10.1128/MCB.24.24.10505-10514.2004
- [260] J. A. Jablonski, E. Buratti, C. Stuani, and M. Caputi. 2008. The secondary structure of the human immunodeficiency virus type 1 transcript modulates viral splicing and infectivity. *J Virol*, 82:8038–8050. doi: 10.1128/JVI.00721-08
- [261] P. J. Shepard and K. J. Hertel. 2008. Conserved RNA secondary structures promote alternative splicing. *RNA*, 14:1463–1469. doi: 10.1261/rna.1069408
- [262] M. Alló, V. Buggiano, J. P. Fededa, E. Petrillo, I. Schor, M. de la Mata, E. Agirre, M. Plass, E. Eyras, S. A. Elela, R. Klinck, B. Chabot, and A. R. Kornblihtt. 2009. Control of alternative splicing through siRNA-mediated transcriptional gene silencing. *Nat Struct Mol Biol*, 16:717–724. doi: 10.1038/nsmb.1620
- [263] H. Tilgner, C. Nikolaou, S. Althammer, M. Sammeth, M. Beato, J. Valcresel, and R. Guig. 2009. Nucleosome positioning as a determinant of exon recognition. *Nat Struct Mol Biol*, 16:996–1001. doi: 10.1038/nsmb.1658
- [264] S. Schwartz, E. Meshorer, and G. Ast. 2009. Chromatin organization marks exon-intron structure. *Nat Struct Mol Biol*, 16:990–995. doi: 10.1038/nsmb.1659
- [265] T. L. Crabb, B. J. Lam, and K. J. Hertel. 2010. Retention of spliceosomal components along ligated exons ensures efficient removal of multiple introns. *RNA*, 16:1786–1796. doi: 10.1261/rna.2186510
- [266] K. Takahara, U. Schwarze, Y. Imamura, G. G. Hoffman, H. Toriello, L. T. Smith, P. H. Byers, and D. S. Greenspan. 2002. Order of intron removal influences multiple splice outcomes, including a two-exon skip, in a COL5A1 acceptor-site mutation that results in abnormal pro-alpha1(V) N-propeptides and Ehlers-Danlos syndrome type I. *Am J Hum Genet*, 71:451–465. doi: 10.1086/342099
- [267] M. de la Mata, C. Lafaille, and A. R. Kornblihtt. 2010. First come, first served revisited: factors affecting the same alternative splicing event have different effects on the relative rates of intron removal. *RNA*, 16:904–912. doi: 10.1261/rna.1993510
- [268] A. M. Zahler, K. M. Neugebauer, W. S. Lane, and M. B. Roth. 1993. Distinct

- functions of SR proteins in alternative pre-mRNA splicing. *Science*, 260:219–222. doi: 10.1126/science.8385799
- [269] C. W. Smith and J. Valcárcel. 2000. Alternative pre-mRNA splicing: the logic of combinatorial control. *Trends Biochem Sci*, 25:381–388. doi: 10.1016/S0968-0004(00)01604-2
- [270] J. Ule, G. Stefani, A. Mele, M. Ruggiu, X. Wang, B. Taneri, T. Gaasterland, B. J. Blencowe, and R. B. Darnell. 2006. An RNA map predicting Nova-dependent splicing regulation. *Nature*, 444:580–586. doi: 10.1038/nature05304
- [271] H. Y. Xiong, Y. Barash, and B. J. Frey. 2011. Bayesian prediction of tissue-regulated splicing using RNA sequence and cellular context. *Bioinformatics*, 27: 2554–2562. doi: 10.1093/bioinformatics/btr444
- [272] J. T. Witten and J. Ule. 2011. Understanding splicing regulation through RNA splicing maps. *Trends Genet*, 27:89–97. doi: 10.1016/j.tig.2010.12.001
- [273] M. M. O'Reilly, M. T. McNally, and K. L. Beemon. 1995. Two strong 5' splice sites and competing, suboptimal 3' splice sites involved in alternative splicing of human immunodeficiency virus type 1 RNA. *Virology*, 213:373–385. doi: 10.1006/viro.1995.0010
- [274] B. A. Amendt, D. Hesslein, L. J. Chang, and C. M. Stoltzfus. 1994. Presence of negative and positive cis-acting RNA splicing elements within and flanking the first tat coding exon of human immunodeficiency virus type 1. *Mol Cell Biol*, 14: 3960–3970. doi: 10.1128/MCB.14.6.3960
- [275] J. D. Levengood, C. Rollins, C. H. J. Mishler, C. A. Johnson, G. Miner, P. Rajan, B. M. Znosko, and B. S. Tolbert. 2012. Solution structure of the HIV-1 exon splicing silencer 3. *J Mol Biol*, 415:680–698. doi: 10.1016/j.jmb.2011.11.034
- [276] M. Caputi, M. Freund, S. Kammler, C. Asang, and H. Schaal. 2004. A bidirectional SF2/ASF- and SRp40-dependent splicing enhancer regulates human immunodeficiency virus type 1 rev, env, vpu, and nef gene expression. *J Virol*, 78: 6517–6526. doi: 10.1128/JVI.78.12.6517-6526.2004
- [277] C. Asang, I. Hauber, and H. Schaal. 2008. Insights into the selective activation of alternatively used splice acceptors by the human immunodeficiency virus type-1 bidirectional splicing enhancer. *Nucleic Acids Res*, 36:1450–1463. doi: 10.1093/nar/gkm1147
- [278] T. O. Tange, C. K. Damgaard, S. Guth, J. Valcrcel, and J. Kjems. 2001. The hnRNP A1 protein regulates HIV-1 tat splicing via a novel intron silencer element. *EMBO J*, 20:5748–5758. doi: 10.1093/emboj/20.20.5748
- [279] A. Tranell, E. M. Feny, and S. Schwartz. 2010. Serine- and arginine-rich proteins

- 55 and 75 (SRp55 and SRp75) induce production of HIV-1 vpr mRNA by inhibiting the 5'-splice site of exon 3. *J Biol Chem*, 285:31537–31547. doi: 10.1074/jbc.M109.077453
- [280] C. M. Stoltzfus and J. M. Madsen. 2006. Role of viral splicing elements and cellular RNA binding proteins in regulation of HIV-1 alternative RNA splicing. *Curr HIV Res*, 4:43–55. doi: 10.2174/157016206775197655
- [281] P. Legrain and M. Rosbash. 1989. Some cis- and trans-acting mutants for splicing target pre-mRNA to the cytoplasm. *Cell*, 57:573–583. doi: 10.1016/0092-8674(89)90127-X
- [282] U. Fischer, S. Meyer, M. Teufel, C. Heckel, R. Lhrmann, and G. Rautmann. 1994. Evidence that HIV-1 Rev directly promotes the nuclear export of unspliced RNA. *EMBO J*, 13:4105–4112
- [283] K. A. Jones and B. M. Peterlin. 1994. Control of RNA initiation and elongation at the HIV-1 promoter. *Annu Rev Biochem*, 63:717–743. doi: 10.1146/annurev.bi.63.070194.003441
- [284] T. W. McCloskey, M. Ott, E. Tribble, S. A. Khan, S. Teichberg, M. O. Paul, S. Pahwa, E. Verdin, and N. Chirmule. 1997. Dual role of HIV Tat in regulation of apoptosis in T cells. *J Immunol*, 158:1014–1019
- [285] G. R. Campbell, E. Pasquier, J. Watkins, V. Bourgarel-Rey, V. Peyrot, D. Esquieu, P. Barbier, J. de Mareuil, D. Braguer, P. Kaleebu, D. L. Yirrell, and E. P. Loret. 2004. The glutamine-rich region of the HIV-1 Tat protein is involved in T-cell apoptosis. *J Biol Chem*, 279:48197–48204. doi: 10.1074/jbc.M406195200
- [286] H. B. Miller, T. J. Robinson, R. Gordn, A. J. Hartemink, and M. A. Garcia-Blanco. 2011. Identification of Tat-SF1 cellular targets by exon array analysis reveals dual roles in transcription and splicing. *RNA*, 17:665–674. doi: 10.1261/rna.2462011
- [287] M. E. Rogel, L. I. Wu, and M. Emerman. 1995. The human immunodeficiency virus type 1 vpr gene prevents cell proliferation during chronic infection. *J Virol*, 69:882–888
- [288] R. A. Fouchier, B. E. Meyer, J. H. Simon, U. Fischer, A. V. Albright, F. Gonzlez-Scarano, and M. H. Malim. 1998. Interaction of the human immunodeficiency virus type 1 Vpr protein with the nuclear pore complex. *J Virol*, 72:6004–6013
- [289] A. K. Gubitz, W. Feng, and G. Dreyfuss. 2004. The SMN complex. *Exp Cell Res*, 296:51–56. doi: 10.1016/j.yexcr.2004.03.022
- [290] C. Trapnell, B. A. Williams, G. Pertea, A. Mortazavi, G. Kwan, M. J. van Baren, S. L. Salzberg, B. J. Wold, and L. Pachter. 2010. Transcript assembly and quan-

- tification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol*, 28:511–515. doi: 10.1038/nbt.1621
- [291] M. F. Rogers, J. Thomas, A. S. Reddy, and A. Ben-Hur. 2012. SpliceGrapher: detecting patterns of alternative splicing from RNA-Seq data in the context of gene models and EST data. *Genome Biol*, 13:R4. doi: 10.1186/gb-2012-13-1-r4
- [292] A. Ryo, Y. Suzuki, K. Ichiyama, T. Wakatsuki, N. Kondoh, A. Hada, M. Yamamoto, and N. Yamamoto. 1999. Serial analysis of gene expression in HIV-1-infected T cell lines. *FEBS Lett*, 462:182–186. doi: 10.1016/S0014-5793(99)01526-4
- [293] G. Lefebvre, S. Desfarges, F. Uyttebroeck, M. Muoz, N. Beerenswinkel, J. Rougemont, A. Telenti, and A. Ciuffi. 2011. Analysis of HIV-1 expression level and sense of transcription by high-throughput sequencing of the infected cell. *J Virol*, 85:6205–6211. doi: 10.1128/JVI.00252-11
- [294] C. Y. Ou, S. Kwok, S. W. Mitchell, D. H. Mack, J. J. Sninsky, J. W. Krebs, P. Feorino, D. Warfield, and G. Schochetman. 1988. DNA amplification for direct detection of HIV-1 in DNA of peripheral blood mononuclear cells. *Science*, 239:295–297. doi: 10.1126/science.3336784
- [295] T. W. Chun, D. Finzi, J. Margolick, K. Chadwick, D. Schwartz, and R. F. Siliciano. 1995. In vivo fate of HIV-1-infected T cells: quantitative analysis of the transition to stable latency. *Nat Med*, 1:1284–1290
- [296] R. T. Davey, N. Bhat, C. Yoder, T. W. Chun, J. A. Metcalf, R. Dewar, V. Natarajan, R. A. Lempicki, J. W. Adelsberger et al. 1999. HIV-1 and T cell dynamics after interruption of highly active antiretroviral therapy (HAART) in patients with a history of sustained viral suppression. *Proc Natl Acad Sci U S A*, 96:15109–15114
- [297] D. Finzi, J. Blankson, J. D. Siliciano, J. B. Margolick, K. Chadwick, T. Pierson, K. Smith, J. Lisziewicz, F. Lori, C. Flexner, T. C. Quinn, R. E. Chaisson, E. Rosenberg, B. Walker, S. Gange, J. Gallant, and R. F. Siliciano. 1999. Latent infection of CD4+ T cells provides a mechanism for lifelong persistence of HIV-1, even in patients on effective combination therapy. *Nat Med*, 5:512–517. doi: 10.1038/8394
- [298] J. D. Siliciano, J. Kajdas, D. Finzi, T. C. Quinn, K. Chadwick, J. B. Margolick, C. Kovacs, S. J. Gange, and R. F. Siliciano. 2003. Long-term follow-up studies confirm the stability of the latent reservoir for HIV-1 in resting CD4+ T cells. *Nat Med*, 9:727–728. doi: 10.1038/nm880
- [299] L. S. Weinberger, R. D. Dar, and M. L. Simpson. 2008. Transient-mediated fate determination in a transcriptional circuit of HIV. *Nat Genet*, 40:466–470. doi: 10.1038/ng.116
- [300] A. Singh, B. Razooky, C. D. Cox, M. L. Simpson, and L. S. Weinberger. 2010. Transcriptional bursting from the HIV-1 promoter is a significant source of

stochastic noise in HIV-1 gene expression. *Biophys J*, 98:L32–L34. doi: 10.1016/j.bpj.2010.03.001

- [301] B. S. Razooky and L. S. Weinberger. 2011. Mapping the architecture of the HIV-1 Tat circuit: A decision-making circuit that lacks bistability and exploits stochastic noise. *Methods*, 53:68–77. doi: 10.1016/j.ymeth.2010.12.006
- [302] H. J. Muller. 1930. Types of visible variations induced by X-rays in *Drosophila*. *J Genet*, 22:299–334
- [303] M. Gaszner and G. Felsenfeld. 2006. Insulators: exploiting transcriptional and epigenetic mechanisms. *Nat Rev Genet*, 7:703–713. doi: 10.1038/nrg1925
- [304] A. Jordan, P. Defechereux, and E. Verdin. 2001. The site of HIV-1 integration in the human genome determines basal transcriptional activity and response to Tat transactivation. *EMBO J*, 20:1726–1738. doi: 10.1093/emboj/20.7.1726
- [305] A. Jordan, D. Bisgrove, and E. Verdin. 2003. HIV reproducibly establishes a latent infection after acute infection of T cells in vitro. *EMBO J*, 22:1868–1877. doi: 10.1093/emboj/cdg188
- [306] R. Pearson, Y. K. Kim, J. Hokello, K. Lassen, J. Friedman, M. Tyagi, and J. Karn. 2008. Epigenetic silencing of human immunodeficiency virus (HIV) transcription by formation of restrictive chromatin structures at the viral long terminal repeat drives the progressive entry of HIV into latency. *J Virol*, 82:12291–12303. doi: 10.1128/JVI.01383-08
- [307] F. Romerio, M. N. Gabriel, and D. M. Margolis. 1997. Repression of human immunodeficiency virus type 1 through the novel cooperation of human factors YY1 and LSF. *J Virol*, 71:9375–9382
- [308] J. J. Coull, F. Romerio, J. M. Sun, J. L. Volker, K. M. Galvin, J. R. Davie, Y. Shi, U. Hansen, and D. M. Margolis. 2000. The human factors YY1 and LSF repress the human immunodeficiency virus type 1 long terminal repeat via recruitment of histone deacetylase 1. *J Virol*, 74:6790–6799. doi: 10.1128/JVI.74.15.6790-6799.2000
- [309] G. He and D. M. Margolis. 2002. Counterregulation of chromatin deacetylation and histone deacetylase occupancy at the integrated promoter of human immunodeficiency virus type 1 (HIV-1) by the HIV-1 repressor YY1 and HIV-1 activator Tat. *Mol Cell Biol*, 22:2965–2973. doi: 10.1128/MCB.22.9.2965-2973.2002
- [310] M. K. Lewinski, D. Bisgrove, P. Shinn, H. Chen, C. Hoffmann, S. Hannenhalli, E. Verdin, C. C. Berry, J. R. Ecker, and F. D. Bushman. 2005. Genome-wide analysis of chromosomal features repressing human immunodeficiency virus transcription. *J Virol*, 79:6610–6619. doi: 10.1128/JVI.79.11.6610-6619.2005

- [311] L. Shan, H.-C. Yang, S. A. Rabi, H. C. Bravo, N. S. Shroff, R. A. Irizarry, H. Zhang, J. B. Margolick, J. D. Siliciano, and R. F. Siliciano. 2011. Influence of host gene transcription level and orientation on HIV-1 latency in a primary-cell model. *J Virol*, 85:5384–5393. doi: 10.1128/JVI.02536-10
- [312] M. J. Pace, E. H. Graf, L. M. Agosto, A. M. Mexas, F. Male, T. Brady, F. D. Bushman, and U. O'Doherty. 2012. Directly infected resting CD4+ T cells can produce HIV Gag without spreading infection in a model of HIV latency. *PLoS Pathog*, 8:e1002818. doi: 10.1371/journal.ppat.1002818
- [313] T. Lenasi, X. Contreras, and B. M. Peterlin. 2008. Transcriptional interference antagonizes proviral gene expression to promote HIV latency. *Cell Host Microbe*, 4:123–133. doi: 10.1016/j.chom.2008.05.016
- [314] Y. Han, Y. B. Lin, W. An, J. Xu, H.-C. Yang, K. O'Connell, D. Dordai, J. D. Boeke, J. D. Siliciano, and R. F. Siliciano. 2008. Orientation-dependent regulation of integrated HIV-1 expression by host gene transcriptional readthrough. *Cell Host Microbe*, 4:134–146. doi: 10.1016/j.chom.2008.06.008
- [315] L. Shan, K. Deng, N. S. Shroff, C. M. Durand, S. A. Rabi, H.-C. Yang, H. Zhang, J. B. Margolick, J. N. Blankson, and R. F. Siliciano. 2012. Stimulation of HIV-1-specific cytolytic T lymphocytes facilitates elimination of latent viral reservoir after virus reactivation. *Immunity*, 36:491–501. doi: 10.1016/j.jimmuni.2012.01.014
- [316] D. Boehm, V. Calvanese, R. D. Dar, S. Xing, S. Schroeder, L. Martins, K. Aull, P.-C. Li, V. Planell, J. E. Bradner, M.-M. Zhou, R. F. Siliciano, L. Weinberger, E. Verdin, and M. Ott. 2013. BET bromodomain-targeting compounds reactivate HIV from latency via a Tat-independent mechanism. *Cell Cycle*, 12:452–462. doi: 10.4161/cc.23309
- [317] A. Savarino, A. Mai, S. Norelli, S. E. Daker, S. Valente, D. Rotili, L. Altucci, A. T. Palamara, and E. Garaci. 2009. “Shock and kill” effects of class I-selective histone deacetylase inhibitors in combination with the glutathione synthesis inhibitor buthionine sulfoximine in cell line models for HIV-1 quiescence. *Retrovirology*, 6: 52. doi: 10.1186/1742-4690-6-52
- [318] N. M. Archin, A. L. Liberty, A. D. Kashuba, S. K. Choudhary, J. D. Kuruc, A. M. Crooks, D. C. Parker, E. M. Anderson, M. F. Kearney, M. C. Strain, D. D. Richman, M. G. Hudgens, R. J. Bosch, J. M. Coffin, J. J. Eron, D. J. Hazuda, and D. M. Margolis. 2012. Administration of vorinostat disrupts HIV-1 latency in patients on antiretroviral therapy. *Nature*, 487:482–485. doi: 10.1038/nature11286
- [319] A. Bosque and V. Planell. 2009. Induction of HIV-1 latency and reactivation in primary memory CD4+ T cells. *Blood*, 113:58–65. doi: 10.1182/blood-2008-07-168393
- [320] A. Bosque and V. Planell. 2011. Studies of HIV-1 latency in an ex vivo model

- that uses primary central memory T cells. *Methods*, 53:54–61. doi: 10.1016/jymeth.2010.10.002
- [321] X. Wu, Y. Li, B. Crise, and S. M. Burgess. 2003. Transcription start regions in the human genome are favored targets for MLV integration. *Science*, 300:1749–1751. doi: 10.1126/science.1083413
- [322] R. S. Mitchell, B. F. Beitzel, A. R. W. Schroder, P. Shinn, H. Chen, C. C. Berry, J. R. Ecker, and F. D. Bushman. 2004. Retroviral DNA integration: ASLV, HIV, and MLV show distinct target site preferences. *PLoS Biol*, 2:e234. doi: 10.1371/journal.pbio.0020234
- [323] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2012
- [324] S. Sherrill-Mix, M. K. Lewinski, M. Famiglietti, A. Bosque, N. Malani, K. E. Ocwieja, C. C. Berry, D. Looney, L. Shan, L. M. Agosto, M. J. Pace, R. F. Siliciano, U. O'Doherty, J. Guatelli, V. Planelles, and F. D. Bushman. 2013. HIV latency and integration site placement in five cell-based models. *Retrovirology*, 10:90. doi: 10.1186/1742-4690-10-90
- [325] C. Berry, S. Hannenhalli, J. Leipzig, and F. D. Bushman. 2006. Selection of target sites for mobile DNA integration in the human genome. *PLoS Comput Biol*, 2:e157. doi: 10.1371/journal.pcbi.0020157
- [326] G. P. Wang, A. Ciuffi, J. Leipzig, C. C. Berry, and F. D. Bushman. 2007. HIV integration site selection: analysis by massively parallel pyrosequencing reveals association with epigenetic modifications. *Genome Res*, 17:1186–1194. doi: 10.1101/gr.6286907
- [327] H. Mochizuki, J. P. Schwartz, K. Tanaka, R. O. Brady, and J. Reiser. 1998. High-titer human immunodeficiency virus type 1-based vector systems for gene delivery into nondividing cells. *J Virol*, 72:8873–8883
- [328] Y. Han, K. Lassen, D. Monie, A. R. Sedaghat, S. Shimoji, X. Liu, T. C. Pierson, J. B. Margolick, R. F. Siliciano, and J. D. Siliciano. 2004. Resting CD4+ T cells from human immunodeficiency virus type 1 (HIV-1)-infected individuals carry integrated HIV-1 genomes within actively transcribed host genes. *J Virol*, 78:6122–6133. doi: 10.1128/JVI.78.12.6122-6133.2004
- [329] G. Plesa, J. Dai, C. Baytop, J. L. Riley, C. H. June, and U. O'Doherty. 2007. Addition of deoxynucleosides enhances human immunodeficiency virus type 1 integration and 2LTR formation in resting CD4+ T cells. *J Virol*, 81:13938–13942. doi: 10.1128/JVI.01745-07
- [330] N. Malani. hiReadsProcessor R package. URL <http://github.com/malnirav/hiReadsProcessor>

- [331] W. J. Kent. 2002. BLAT—the BLAST-like alignment tool. *Genome Res*, 12:656–664. doi: 10.1101/gr.229202
- [332] C. C. Berry, K. Ocwieja, N. Malani, and F. D. Bushman. 2014. Comparing DNA integration site clusters with Scan Statistics. *Bioinformatics*, 30:1493–1500. doi: 10.1093/bioinformatics/btu035
- [333] J. Ernst and M. Kellis. 2010. Discovery and characterization of chromatin states for systematic annotation of the human genome. *Nat Biotechnol*, 28:817–825. doi: 10.1038/nbt.1662
- [334] A. S. Hinrichs, D. Karolchik, R. Baertsch, G. P. Barber, G. Bejerano, H. Clawson, M. Diekhans, T. S. Furey, R. A. Harte et al. 2006. The UCSC Genome Browser Database: update 2006. *Nucleic Acids Res*, 34:D590–D598. doi: 10.1093/nar/gkj144
- [335] K. R. Rosenbloom, C. A. Sloan, V. S. Malladi, T. R. Dreszer, K. Learned, V. M. Kirkup, M. C. Wong, M. Maddren, R. Fang et al. 2013. ENCODE data in the UCSC Genome Browser: year 5 update. *Nucleic Acids Res*, 41:D56–D63. doi: 10.1093/nar/gks1172
- [336] J. Han, S.-G. Park, J.-B. Bae, J. Choi, J.-M. Lyu, S. H. Park, H. S. Kim, Y.-J. Kim, S. Kim, and T.-Y. Kim. 2012. The characteristics of genome-wide DNA methylation in naïve CD4+ T cells of patients with psoriasis or atopic dermatitis. *Biochem Biophys Res Commun*, 422:157–163. doi: 10.1016/j.bbrc.2012.04.128
- [337] L. R. Meyer, A. S. Zweig, A. S. Hinrichs, D. Karolchik, R. M. Kuhn, M. Wong, C. A. Sloan, K. R. Rosenbloom, G. Roe et al. 2013. The UCSC Genome Browser database: extensions and updates 2013. *Nucleic Acids Res*, 41:D64–D69. doi: 10.1093/nar/gks1048
- [338] Z. Wang, C. Zang, J. A. Rosenfeld, D. E. Schones, A. Barski, S. Cuddapah, K. Cui, T.-Y. Roh, W. Peng, M. Q. Zhang, and K. Zhao. 2008. Combinatorial patterns of histone acetylations and methylations in the human genome. *Nat Genet*, 40: 897–903. doi: 10.1038/ng.154
- [339] A. Barski, S. Cuddapah, K. Cui, T.-Y. Roh, D. E. Schones, Z. Wang, G. Wei, I. Chepelev, and K. Zhao. 2007. High-resolution profiling of histone methylations in the human genome. *Cell*, 129:823–837. doi: 10.1016/j.cell.2007.05.009
- [340] Z. Wang, C. Zang, K. Cui, D. E. Schones, A. Barski, W. Peng, and K. Zhao. 2009. Genome-wide mapping of HATs and HDACs reveals distinct functions in active and inactive genes. *Cell*, 138:1019–1031. doi: 10.1016/j.cell.2009.06.049
- [341] D. E. Schones, K. Cui, S. Cuddapah, T.-Y. Roh, A. Barski, Z. Wang, G. Wei, and K. Zhao. 2008. Dynamic regulation of nucleosome positioning in the human genome. *Cell*, 132:887–898. doi: 10.1016/j.cell.2008.02.022

- [342] F. Hsu, W. J. Kent, H. Clawson, R. M. Kuhn, M. Diekhans, and D. Haussler. 2006. The UCSC Known Genes. *Bioinformatics*, 22:1036–1046. doi: 10.1093/bioinformatics/btl048
- [343] J. Friedman, T. Hastie, and R. Tibshirani. 2010. Regularization paths for generalized linear models via coordinate descent. *J Stat Softw*, 33:1–22
- [344] I. H. Greger, F. Demarchi, M. Giacca, and N. J. Proudfoot. 1998. Transcriptional interference perturbs the binding of Sp1 to the HIV-1 promoter. *Nucleic Acids Res*, 26:1294–1301
- [345] A. De Marco, C. Biancotto, A. Knezevich, P. Maiuri, C. Vardabasso, and A. Marcello. 2008. Intragenic transcriptional cis-activation of the human immunodeficiency virus 1 does not result in allele-specific inhibition of the endogenous gene. *Retrovirology*, 5:98. doi: 10.1186/1742-4690-5-98
- [346] J. S. Waye and H. F. Willard. 1987. Nucleotide sequence heterogeneity of alpha satellite repetitive DNA: a survey of alphoid sequences from different human chromosomes. *Nucleic Acids Res*, 15:7549–7569. doi: 10.1093/nar/15.18.7549
- [347] J. Jurka, V. V. Kapitonov, A. Pavlicek, P. Klonowski, O. Kohany, and J. Walichiewicz. 2005. Repbase Update, a database of eukaryotic repetitive elements. *Cytogenet Genome Res*, 110:462–467. doi: 10.1159/000084979
- [348] E. Verdin, P. Paras, and C. Van Lint. 1993. Chromatin disruption in the promoter of human immunodeficiency virus type 1 during transcriptional activation. *EMBO J*, 12:3249–3259
- [349] C. Van Lint, S. Emiliani, M. Ott, and E. Verdin. 1996. Transcriptional activation and chromatin remodeling of the HIV-1 promoter in response to histone acetylation. *EMBO J*, 15:1112–1120
- [350] K. G. Lassen, K. X. Ramyar, J. R. Bailey, Y. Zhou, and R. F. Siliciano. 2006. Nuclear retention of multiply spliced HIV-1 RNA in resting CD4+ T cells. *PLoS Pathog*, 2:e68. doi: 10.1371/journal.ppat.0020068
- [351] M. Dieudonné, P. Maiuri, C. Biancotto, A. Knezevich, A. Kula, M. Lusic, and A. Marcello. 2009. Transcriptional competence of the integrated HIV-1 provirus at the nuclear periphery. *EMBO J*, 28:2231–2243. doi: 10.1038/emboj.2009.141
- [352] R. F. Siliciano and W. C. Greene. 2011. HIV Latency. *Cold Spring Harb Perspect Med*, 1:a007096. doi: 10.1101/cshperspect.a007096
- [353] M. Lusic, B. Marini, H. Ali, B. Lucic, R. Luzzati, and M. Giacca. 2013. Proximity to PML nuclear bodies regulates HIV-1 latency in CD4+ T cells. *Cell Host Microbe*, 13:665–677. doi: 10.1016/j.chom.2013.05.006

- [354] Q. Pan, O. Shai, L. J. Lee, B. J. Frey, and B. J. Blencowe. 2008. Deep surveying of alternative splicing complexity in the human transcriptome by high-throughput sequencing. *Nature Genetics*, 40:1413–1415. doi: 10.1038/ng.259
- [355] F. Pagani, M. Raponi, and F. E. Baralle. 2005. Synonymous mutations in CFTR exon 12 affect splicing and are not neutral in evolution. *Proc Natl Acad Sci U S A*, 102:6368–6372. doi: 10.1073/pnas.0502288102
- [356] C. Wang, Y. Mitsuya, B. Gharizadeh, M. Ronaghi, and S. R. W. 2007. Characterization of mutation spectra with ultra-deep pyrosequencing: Application to HIV-1 drug resistance. *Genome Research*, 17:1195–1201. doi: 10.1101/gr.6468307
- [357] K. Wang, R. Wernersson, and S. Brunak. 2011. The strength of intron donor splice sites in human genes displays a bell-shaped pattern. *Bioinformatics*, 27: 3079–3084. doi: 10.1093/bioinformatics/btr532
- [358] C. Carrera, M. Pinilla, L. Pérez-Alvarez, and M. M. Thomson. 2010. Identification of unusual and novel HIV type 1 spliced transcripts generated in vivo. *AIDS Res Hum Retroviruses*, 26:815–820. doi: 10.1089/aid.2010.0011
- [359] M. Lützelberger, L. S. Reinert, A. T. Das, B. Berkhout, and J. Kjems. 2006. A novel splice donor site in the gag-pol gene is required for HIV-1 RNA stability. *J Biol Chem*, 281:18644–18651. doi: 10.1074/jbc.M513698200
- [360] J. Salfeld, H. G. Gtlinger, R. A. Sia, R. E. Park, J. G. Sodroski, and W. A. Haseltine. 1990. A tripartite HIV-1 tat-env-rev fusion protein. *EMBO J*, 9:965–970
- [361] S. Schwartz, B. K. Felber, D. M. Benko, E. M. Fenyö, and G. N. Pavlakis. 1990. Cloning and functional analysis of multiply spliced mRNA species of human immunodeficiency virus type 1. *J Virol*, 64:2519–2529
- [362] J. Smith, A. Azad, and N. Deacon. 1992. Identification of two novel human immunodeficiency virus type 1 splice acceptor sites in infected T cell lines. *J Gen Virol*, 73 (Pt 7):1825–1828
- [363] J. A. Jablonski and M. Caputi. 2009. Role of cellular RNA processing factors in human immunodeficiency virus type 1 mRNA metabolism, replication, and infectivity. *J Virol*, 83:981–992. doi: 10.1128/JVI.01801-08
- [364] A. Tranell, S. Tingsborg, E. M. Feny, and S. Schwartz. 2011. Inhibition of splicing by serine-arginine rich protein 55 (SRp55) causes the appearance of partially spliced HIV-1 mRNAs in the cytoplasm. *Virus Res*, 157:82–91. doi: 10.1016/j.virusres.2011.02.010
- [365] H. Zhou, M. Xu, Q. Huang, A. T. Gates, X. D. Zhang, J. C. Castle, E. Stec, M. Ferrer, B. Strulovici, D. J. Hazuda, and A. S. Espeseth. 2008. Genome-scale RNAi screen

- for host factors required for HIV replication. *Cell Host Microbe*, 4:495–504. doi: 10.1016/j.chom.2008.10.004
- [366] Y. Zhu, G. Chen, F. Lv, X. Wang, X. Ji, Y. Xu, J. Sun, L. Wu, Y.-T. Zheng, and G. Gao. 2011. Zinc-finger antiviral protein inhibits HIV-1 infection by selectively targeting multiply spliced viral mRNAs for degradation. *Proc Natl Acad Sci U S A*, 108:15834–15839. doi: 10.1073/pnas.1101676108
- [367] M. J. Saltarelli, E. Hadziyannis, C. E. Hart, J. V. Harrison, B. K. Felber, T. J. Spira, and G. N. Pavlakis. 1996. Analysis of human immunodeficiency virus type 1 mRNA splicing patterns during disease progression in peripheral blood mononuclear cells from infected individuals. *AIDS Res Hum Retroviruses*, 12: 1443–1456. doi: 10.1089/aid.1996.12.1443
- [368] E. Delgado, C. Carrera, P. Nebreda, A. Fernndez-Garca, M. Pinilla, V. Garca, L. Prez-lvarez, and M. M. Thomson. 2012. Identification of new splice sites used for generation of rev transcripts in human immunodeficiency virus type 1 subtype C primary isolates. *PLoS One*, 7:e30574. doi: 10.1371/journal.pone.0030574
- [369] P. Grabowski. 2011. Alternative splicing takes shape during neuronal development. *Curr Opin Genet Dev*, 21:388–394. doi: 10.1016/j.gde.2011.03.005
- [370] M. Llorian and C. W. J. Smith. 2011. Decoding muscle alternative splicing. *Curr Opin Genet Dev*, 21:380–387. doi: 10.1016/j.gde.2011.03.006
- [371] J. Y. Ip, A. Tong, Q. Pan, J. D. Topp, B. J. Blencowe, and K. W. Lynch. 2007. Global analysis of alternative splicing during T-cell activation. *RNA*, 13:563–572. doi: 10.1261/rna.457207
- [372] J. D. Topp, J. Jackson, A. A. Melton, and K. W. Lynch. 2008. A cell-based screen for splicing regulators identifies hnRNP LL as a distinct signal-induced repressor of CD45 variable exon 4. *RNA*, 14:2038–2049. doi: 10.1261/rna.1212008
- [373] S. Sonza, H. P. Mutimer, K. O'Brien, P. Ellery, J. L. Howard, J. H. Axelrod, N. J. Deacon, S. M. Crowe, and D. F. J. Purcell. 2002. Selectively reduced tat mRNA heralds the decline in productive human immunodeficiency virus type 1 infection in monocyte-derived macrophages. *J Virol*, 76:12611–12621
- [374] D. Dowling, S. Nasr-Esfahani, C. H. Tan, K. O'Brien, J. L. Howard, D. A. Jans, D. F. j Purcell, C. M. Stoltzfus, and S. Sonza. 2008. HIV-1 infection induces changes in expression of cellular splicing factors that regulate alternative viral splicing and virus production in macrophages. *Retrovirology*, 5:18. doi: 10.1186/1742-4690-5-18
- [375] R. Collman, J. W. Balliet, S. A. Gregory, H. Friedman, D. L. Kolson, N. Nathanson, and A. Srinivasan. 1992. An infectious molecular clone of an unusual macrophage-

- tropic and highly cytopathic strain of human immunodeficiency virus type 1. *J Virol*, 66:7517–7521
- [376] H. Deng, R. Liu, W. Ellmeier, S. Choe, D. Unutmaz, M. Burkhardt, P. Di Marzio, S. Marmon, R. E. Sutton, C. M. Hill, C. B. Davis, S. C. Peiper, T. J. Schall, D. R. Littman, and N. R. Landau. 1996. Identification of a major co-receptor for primary isolates of HIV-1. *Nature*, 381:661–666. doi: 10.1038/381661a0
- [377] N. R. Landau and D. R. Littman. 1992. Packaging system for rapid production of murine leukemia virus vectors with variable tropism. *J Virol*, 66:5110–5113
- [378] X. Wei, J. M. Decker, S. Wang, H. Hui, J. C. Kappes, X. Wu, J. F. Salazar-Gonzalez, M. G. Salazar, J. M. Kilby, M. S. Saag, N. L. Komarova, M. A. Nowak, B. H. Hahn, P. D. Kwong, and G. M. Shaw. 2003. Antibody neutralization and escape by HIV-1. *Nature*, 422:307–312. doi: 10.1038/nature01470
- [379] N. Srinivasakumar, N. Chazal, C. Helga-Maria, S. Prasad, M. L. Hammarskjöld, and D. Rekosh. 1997. The effect of viral regulatory protein expression on gene delivery by human immunodeficiency virus type 1 vectors produced in stable packaging cell lines. *J Virol*, 71:5841–5848
- [380] D. C. Shugars, M. S. Smith, D. H. Glueck, P. V. Nantermet, F. Seillier-Moiseiwitsch, and R. Swanstrom. 1993. Analysis of human immunodeficiency virus type 1 nef gene sequences present in vivo. *J Virol*, 67:4639–4650
- [381] K. J. Travers, C.-S. Chin, D. R. Rank, J. S. Eid, and S. W. Turner. 2010. A flexible and efficient template format for circular consensus sequencing and SNP detection. *Nucleic Acids Res*, 38:e159. doi: 10.1093/nar/gkq543
- [382] T. A. Thanaraj and F. Clark. 2001. Human GC-AG alternative intron isoforms with weak donor sites show enhanced consensus at acceptor exon positions. *Nucleic Acids Res*, 29:2581–2593. doi: 10.1093/nar/29.12.2581
- [383] M. Aebei, H. Hornig, R. A. Padgett, J. Reiser, and C. Weissmann. 1986. Sequence requirements for splicing of higher eukaryotic nuclear pre-mRNA. *Cell*, 47: 555–565
- [384] M. Burset, I. A. Seledtsov, and V. V. Solovyev. 2000. Analysis of canonical and non-canonical splice sites in mammalian genomes. *Nucleic Acids Res*, 28:4364–4375. doi: 10.1093/nar/28.21.4364
- [385] M. Burset, I. A. Seledtsov, and V. V. Solovyev. 2001. SpliceDB: database of canonical and non-canonical mammalian splice sites. *Nucleic Acids Res*, 29: 255–259. doi: 10.1093/nar/29.1.255
- [386] N. Sheth, X. Roca, M. L. Hastings, T. Roeder, A. R. Krainer, and R. Sachidanan-

- dam. 2006. Comprehensive splice-site analysis using comparative genomics. *Nucleic Acids Res*, 34:3955–3967. doi: 10.1093/nar/gkl556
- [387] J. C. Guatelli, T. R. Gingeras, and D. D. Richman. 1990. Alternative splice acceptor utilization during human immunodeficiency virus type 1 infection of cultured cells. *J Virol*, 64:4093–4098
- [388] C. Kuiken, B. Foley, T. Leitner, C. Apetrei, B. Hahn, I. Mizrachi, J. Mullins, A. Rambaut, S. Wolinsky, and B. Korber. 2010. HIV Sequence Compendium 2010. Theoretical Biology and Biophysics Group, Los Alamos National Laboratory, New Mexico. URL <http://www.hiv.lanl.gov/content/sequence/HIV/COMPENDIUM/2010compendium.html>
- [389] C. Burge, T. Tuschl, and P. Sharp. 1999. Splicing of precursors to mRNAs by the spliceosomes. *Cold Spring Harbor Monograph Archive*, 37. doi: 10.1101/087969589.37.525
- [390] T. E. M. Abbink and B. Berkhout. 2008. RNA structure modulates splicing efficiency at the human immunodeficiency virus type 1 major splice donor. *J Virol*, 82:3090–3098. doi: 10.1128/JVI.01479-07
- [391] K. Verhoef, P. S. Bilodeau, J. L. van Wamel, J. Kjems, C. M. Stoltzfus, and B. Berkhout. 2001. Repair of a Rev-minus human immunodeficiency virus type 1 mutant by activation of a cryptic splice site. *J Virol*, 75:3495–3500. doi: 10.1128/JVI.75.7.3495-3500.2001
- [392] A. M. Zahler, C. K. Damgaard, J. Kjems, and M. Caputi. 2004. SC35 and heterogeneous nuclear ribonucleoprotein A/B proteins bind to a juxtaposed exonic splicing enhancer/exonic splicing silencer element to regulate HIV-1 tat exon 2 splicing. *J Biol Chem*, 279:10077–10084. doi: 10.1074/jbc.M312743200
- [393] S. K. Arya, C. Guo, S. F. Josephs, and F. Wong-Staal. 1985. Trans-activator gene of human T-lymphotropic virus type III (HTLV-III). *Science*, 229:69–73
- [394] K. E. Ocwieja, S. Sherrill-Mix, R. Mukherjee, R. Custers-Allen, P. David, M. Brown, S. Wang, D. R. Link, J. Olson, K. Travers, E. Schadt, and F. D. Bushman. 2012. Dynamic regulation of HIV-1 mRNA populations analyzed by single-molecule enrichment and long-read sequencing. *Nucleic Acids Res*, 40: 10345–10355. doi: 10.1093/nar/gks753
- [395] J. He, S. Choe, R. Walker, P. Di Marzio, D. O. Morgan, and N. R. Landau. 1995. Human immunodeficiency virus type 1 viral protein R (Vpr) arrests cells in the G2 phase of the cell cycle by inhibiting p34cdc2 activity. *J Virol*, 69:6705–6711
- [396] J. B. Jowett, V. Planelles, B. Poon, N. P. Shah, M. L. Chen, and I. S. Chen. 1995. The human immunodeficiency virus type 1 vpr gene arrests infected T cells in the G2 + M phase of the cell cycle. *J Virol*, 69:6304–6313

- [397] W. C. Goh, M. E. Rogel, C. M. Kinsey, S. F. Michael, P. N. Fultz, M. A. Nowak, B. H. Hahn, and M. Emerman. 1998. HIV-1 Vpr increases viral expression by manipulation of the cell cycle: a mechanism for selection of Vpr in vivo. *Nat Med*, 4:65–71. doi: 10.1038/nm0198-065
- [398] R. A. Marciniak and P. A. Sharp. 1991. HIV-1 Tat protein promotes formation of more-processive elongation complexes. *EMBO J*, 10:4189–4196
- [399] P. Wei, M. E. Garber, S. M. Fang, W. H. Fischer, and K. A. Jones. 1998. A novel CDK9-associated C-type cyclin interacts directly with HIV-1 Tat and mediates its high-affinity, loop-specific binding to TAR RNA. *Cell*, 92:451–462. doi: 10.1016/S0092-8674(00)80939-3
- [400] S. Kanazawa, T. Okamoto, and B. M. Peterlin. 2000. Tat competes with CIITA for the binding to P-TEFb and blocks the expression of MHC class II genes in HIV infection. *Immunity*, 12:61–70. doi: 10.1016/S1074-7613(00)80159-4
- [401] M. Barboric, J. H. N. Yik, N. Czudnochowski, Z. Yang, R. Chen, X. Contreras, M. Geyer, B. Matija Peterlin, and Q. Zhou. 2007. Tat competes with HEXIM1 to increase the active pool of P-TEFb for HIV-1 transcription. *Nucleic Acids Res*, 35: 2003–2012. doi: 10.1093/nar/gkm063
- [402] S. K. O'Brien, H. Cao, R. Nathans, A. Ali, and T. M. Rana. 2010. P-TEFb kinase complex phosphorylates histone H1 to regulate expression of cellular and HIV-1 genes. *J Biol Chem*, 285:29713–29720. doi: 10.1074/jbc.M110.125997
- [403] L. Muniz, S. Egloff, B. Ughy, B. E. Jády, and T. Kiss. 2010. Controlling cellular P-TEFb activity by the HIV-1 transcriptional transactivator Tat. *PLoS Pathog*, 6: e1001152. doi: 10.1371/journal.ppat.1001152
- [404] J. Corbeil, D. Sheeter, D. Genini, S. Rought, L. Leoni, P. Du, M. Ferguson, D. R. Masys, J. B. Welsh, J. L. Fink, R. Sasik, D. Huang, J. Drenkow, D. D. Richman, and T. Gingeras. 2001. Temporal gene regulation during HIV-1 infection of human CD4+ T cells. *Genome Res*, 11:1198–1204. doi: 10.1101/gr.180201
- [405] C. H. Woelk, F. Ottone, C. R. Plotkin, P. Du, C. D. Royer, S. E. Rought, J. Lozach, R. Sasik, R. S. Kornbluth, D. D. Richman, and J. Corbeil. 2004. Interferon gene expression following HIV type 1 infection of monocyte-derived macrophages. *AIDS Res Hum Retroviruses*, 20:1210–1222. doi: 10.1089/0889222042545009
- [406] M. D. Hyrcza, C. Kovacs, M. Loutfy, R. Halpenny, L. Heisler, S. Yang, O. Wilkins, M. Ostrowski, and S. D. Der. 2007. Distinct transcriptional profiles in ex vivo CD4+ and CD8+ T cells are established early in human immunodeficiency virus type 1 infection and are characterized by a chronic interferon response as well as extensive transcriptional changes in CD8+ T cells. *J Virol*, 81:3477–3486. doi: 10.1128/JVI.01552-06

- [407] J. Q. Wu, D. E. Dwyer, W. B. Dyer, Y. H. Yang, B. Wang, and N. K. Saksena. 2008. Transcriptional profiles in CD8+ T cells from HIV+ progressors on HAART are characterized by coordinated up-regulation of oxidative phosphorylation enzymes and interferon responses. *Virology*, 380:124–135. doi: 10.1016/j.virol.2008.06.039
- [408] A. J. Smith, Q. Li, S. W. Wietgrefe, T. W. Schacker, C. S. Reilly, and A. T. Haase. 2010. Host genes associated with HIV-1 replication in lymphatic tissue. *J Immunol*, 185:5417–5424. doi: 10.4049/jimmunol.1002197
- [409] M. Imbeault, K. Giguère, M. Ouellet, and M. J. Tremblay. 2012. Exon level transcriptomic profiling of HIV-1-infected CD4(+) T cells reveals virus-induced genes and host environment favorable for viral replication. *PLoS Pathog*, 8: e1002861. doi: 10.1371/journal.ppat.1002861
- [410] P. Mohammadi, S. Desfarges, I. Bartha, B. Joos, N. Zanger, M. Muoz, H. F. Gnethard, N. Beerenswinkel, A. Telenti, and A. Ciuffi. 2013. 24 hours in the life of HIV-1 in a T cell line. *PLoS Pathog*, 9:e1003161. doi: 10.1371/journal.ppat.1003161
- [411] X. Peng, P. Sova, R. R. Green, M. J. Thomas, M. J. Korth, S. Proll, J. Xu, Y. Cheng, K. Yi, L. Chen, Z. Peng, J. Wang, R. E. Palermo, and M. G. Katze. 2014. Deep sequencing of HIV-infected cells: insights into nascent transcription and host-directed therapy. *J Virol*, 88:8768–8782. doi: 10.1128/JVI.00768-14
- [412] C. de la Fuente, F. Santiago, L. Deng, C. Eadie, I. Zilberman, K. Kehn, A. Madrukuri, S. Baylor, K. Wu, C. G. Lee, A. Pumfrey, and F. Kashanchi. 2002. Gene expression profile of HIV-1 Tat expressing cells: a close interplay between proliferative and differentiation signals. *BMC Biochem*, 3:14. doi: 10.1186/1471-2091-3-14
- [413] B. Langmead, C. Trapnell, M. Pop, and S. L. Salzberg. 2009. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biology*, 10:R25. doi: 10.1186/gb-2009-10-3-r25
- [414] G. R. Grant, M. H. Farkas, A. D. Pizarro, N. F. Lahens, J. Schug, B. P. Brunk, C. J. Stoeckert, J. B. Hogenesch, and E. A. Pierce. 2011. Comparative analysis of RNA-Seq alignment algorithms and the RNA-Seq unified mapper (RUM). *Bioinformatics*, 27:2518–2528. doi: 10.1093/bioinformatics/btr427
- [415] Q. Li, A. J. Smith, T. W. Schacker, J. V. Carlis, L. Duan, C. S. Reilly, and A. T. Haase. 2009. Microarray analysis of lymphatic tissue reveals stage-specific, gene expression signatures in HIV-1 infection. *J Immunol*, 183:1975–1982. doi: 10.4049/jimmunol.0803222
- [416] A. Subramanian, P. Tamayo, V. K. Mootha, S. Mukherjee, B. L. Ebert, M. A. Gillette, A. Paulovich, S. L. Pomeroy, T. R. Golub, E. S. Lander, and J. P. Mesirov. 2005. Gene set enrichment analysis: a knowledge-based approach for interpreting

- genome-wide expression profiles. *Proc Natl Acad Sci U S A*, 102:15545–15550. doi: 10.1073/pnas.0506580102
- [417] W. J. Kent, C. W. Sugnet, T. S. Furey, K. M. Roskin, T. H. Pringle, A. M. Zahler, and D. Haussler. 2002. The human genome browser at UCSC. *Genome Res*, 12: 996–1006. doi: 10.1101/gr.229102
- [418] H. Li, B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis, R. Durbin, and . G. P. D. P. S. . 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25:2078–2079. doi: 10.1093/bioinformatics/btp352
- [419] R. P. Subramanian, J. H. Wildschutte, C. Russo, and J. M. Coffin. 2011. Identification, characterization, and comparative genomic distribution of the HERV-K (HML-2) group of human endogenous retroviruses. *Retrovirology*, 8:90. doi: 10.1186/1742-4690-8-90
- [420] G. La Mantia, D. Maglione, G. Pengue, A. Di Cristofano, A. Simeone, L. Lanfrancone, and L. Lania. 1991. Identification and characterization of novel human endogenous retroviral sequences preferentially expressed in undifferentiated embryonal carcinoma cells. *Nucleic Acids Res*, 19:1513–1520
- [421] G. La Mantia, B. Majello, A. Di Cristofano, M. Strazzullo, G. Minchiotti, and L. Lania. 1992. Identification of regulatory elements within the minimal promoter region of the human endogenous ERV9 proviruses: accurate transcription initiation is controlled by an Inr-like element. *Nucleic Acids Res*, 20:4129–4136. doi: 10.1093/nar/20.16.4129
- [422] K. E. Plant, S. J. Routledge, and N. J. Proudfoot. 2001. Intergenic transcription in the human beta-globin gene cluster. *Mol Cell Biol*, 21:6507–6514. doi: 10.1128/MCB.21.19.6507-6514.2001
- [423] J. Ling, W. Pi, R. Bollag, S. Zeng, M. Keskinetepe, H. Saliman, S. Krantz, B. Whitney, and D. Tuan. 2002. The solitary long terminal repeats of ERV-9 endogenous retrovirus are conserved during primate evolution and possess enhancer activities in embryonic and hematopoietic cells. *J Virol*, 76:2410–2423. doi: 10.1128/jvi.76.5.2410-2423.2002
- [424] X. Yu, X. Zhu, W. Pi, J. Ling, L. Ko, Y. Takeda, and D. Tuan. 2005. The long terminal repeat (LTR) of ERV-9 human endogenous retrovirus binds to NF-Y in the assembly of an active LTR enhancer complex NF-Y/MZF1/GATA-2. *J Biol Chem*, 280:35184–35194. doi: 10.1074/jbc.M508138200
- [425] R. C. Edgar. 2004. MUSCLE: a multiple sequence alignment method with reduced time and space complexity. *BMC Bioinformatics*, 5:113. doi: 10.1186/1471-2105-5-113

- [426] M. Rotger, J. Dalmau, A. Rauch, P. McLaren, S. E. Bosingher, R. Martinez, N. G. Sandler, A. Roque, J. Liebner et al. 2011. Comparative transcriptomics of extreme phenotypes of human HIV-1 infection and SIV infection in sooty mangabey and rhesus macaque. *J Clin Invest*, 121:2391–2400. doi: 10.1172/JCI45235
- [427] K. Breuer, A. K. Foroushani, M. R. Laird, C. Chen, A. Sribnaia, R. Lo, G. L. Winsor, R. E. W. Hancock, F. S. L. Brinkman, and D. J. Lynn. 2013. InnateDB: systems biology of innate immunity and beyond—recent updates and continuing curation. *Nucleic Acids Res*, 41:D1228–D1233. doi: 10.1093/nar/gks1147
- [428] I. Rusinova, S. Forster, S. Yu, A. Kannan, M. Masse, H. Cumming, R. Chapman, and P. J. Hertzog. 2013. Interferome v2.0: an updated database of annotated interferon-regulated genes. *Nucleic Acids Res*, 41:D1040–D1046. doi: 10.1093/nar/gks1215
- [429] S. T. Chang, M. J. Thomas, P. Sova, R. R. Green, R. E. Palermo, and M. G. Katze. 2013. Next-generation sequencing of small RNAs from HIV-infected cells identifies phased microRNA expression patterns and candidate novel microRNAs differentially expressed upon infection. *MBio*, 4:e00549–e00512. doi: 10.1128/mBio.00549-12
- [430] Z. Kalender Atak, K. De Keersmaecker, V. Gianfelicci, E. Geerdens, R. Vandepoel, D. Pauwels, M. Porcu, I. Lahortiga, V. Brys, W. G. Dirks, H. Quentmeier, J. Cloos, H. Cuppens, A. Uyttebroeck, P. Vandenberghe, J. Cools, and S. Aerts. 2012. High accuracy mutation detection in leukemia on a selected panel of cancer genes. *PLoS One*, 7:e38463. doi: 10.1371/journal.pone.0038463
- [431] E. S. Patel and L.-J. Chang. 2012. Synergistic effects of interleukin-7 and pre-T cell receptor signaling in human T cell development. *J Biol Chem*, 287: 33826–33835. doi: 10.1074/jbc.M112.380113
- [432] M. Imbeault, M. Ouellet, and M. J. Tremblay. 2009. Microarray study reveals that HIV-1 induces rapid type-I interferon-dependent p53 mRNA up-regulation in human primary CD4+ T cells. *Retrovirology*, 6:5. doi: 10.1186/1742-4690-6-5
- [433] S. Iwase, Y. Furukawa, J. Kikuchi, M. Nagai, Y. Terui, M. Nakamura, and H. Yamada. 1997. Modulation of E2F activity is linked to interferon-induced growth suppression of hematopoietic cells. *J Biol Chem*, 272:12406–12414. doi: 10.1074/jbc.272.19.12406
- [434] R. W. Johnstone, J. A. Kerry, and J. A. Trapani. 1998. The human interferon-inducible protein, IFI 16, is a repressor of transcription. *J Biol Chem*, 273: 17172–17177. doi: 10.1074/jbc.273.27.17172
- [435] B. R. Williams. 1999. PKR; a sentinel kinase for cellular stress. *Oncogene*, 18: 6112–6120. doi: 10.1038/sj.onc.1203127

- [436] C. V. Ramana, N. Grammatikakis, M. Chernov, H. Nguyen, K. C. Goh, B. R. Williams, and G. R. Stark. 2000. Regulation of c-myc expression by IFN-gamma through Stat1-dependent and -independent pathways. *EMBO J*, 19:263–272. doi: 10.1093/emboj/19.2.263
- [437] S.-L. Liang, D. Quirk, and A. Zhou. 2006. RNase L: its biological roles and regulation. *IUBMB Life*, 58:508–514. doi: 10.1080/15216540600838232
- [438] F. Maldarelli, C. Xiang, G. Chamoun, and S. L. Zeichner. 1998. The expression of the essential nuclear splicing factor SC35 is altered by human immunodeficiency virus infection. *Virus Res*, 53:39–51
- [439] A. Monette, L. Ajamian, M. López-Lastra, and A. J. Mouland. 2009. Human immunodeficiency virus type 1 (HIV-1) induces the cytoplasmic retention of heterogeneous nuclear ribonucleoprotein A1 by disrupting nuclear import: implications for HIV-1 gene expression. *J Biol Chem*, 284:31350–31362. doi: 10.1074/jbc.M109.048736
- [440] R. Contreras-Galindo, P. López, R. Vélez, and Y. Yamamura. 2007. HIV-1 infection increases the expression of human endogenous retroviruses type K (HERV-K) in vitro. *AIDS Res Hum Retroviruses*, 23:116–122. doi: 10.1089/aid.2006.0117
- [441] R. Contreras-Galindo, M. H. Kaplan, S. He, A. C. Contreras-Galindo, M. J. Gonzalez-Hernandez, F. Kappes, D. Dube, S. M. Chan, D. Robinson, F. Meng, M. Dai, S. D. Gitlin, A. M. Chinnaiyan, G. S. Omenn, and D. M. Markovitz. 2013. HIV infection reveals widespread expansion of novel centromeric human endogenous retroviruses. *Genome Res*, 23:1505–1513. doi: 10.1101/gr.144303.112
- [442] N. Bhardwaj, F. Maldarelli, J. Mellors, and J. M. Coffin. 2014. HIV-1 infection leads to increased transcription of human endogenous retrovirus HERV-K (HML-2) proviruses in vivo but not to increased virion production. *J Virol*, 88:11108–11120. doi: 10.1128/JVI.01623-14
- [443] R. B. Jones, H. Song, Y. Xu, K. E. Garrison, A. A. Buzdin, N. Anwar, D. V. Hunter, S. Mujib, V. Mihajlovic, E. Martin, E. Lee, M. Kuciak, R. A. S. Raposo, A. Bozorgzad, D. A. Meiklejohn, L. C. Ndhlovu, D. F. Nixon, and M. A. Ostrowski. 2013. LINE-1 retrotransposable element DNA accumulates in HIV-1-infected cells. *J Virol*, 87:13307–13320. doi: 10.1128/JVI.02257-13
- [444] P. Medstrand and D. L. Mager. 1998. Human-specific integrations of the HERV-K endogenous retrovirus family. *J Virol*, 72:9782–9787
- [445] C. Macfarlane and P. Simmonds. 2004. Allelic variation of HERV-K(HML-2) endogenous retroviral elements in human populations. *J Mol Evol*, 59:642–656. doi: 10.1007/s00239-004-2656-1
- [446] K. Büscher, U. Trefzer, M. Hofmann, W. Sterry, R. Kurth, and J. Denner. 2005.

Expression of human endogenous retrovirus K in melanomas and melanoma cell lines. *Cancer Res*, 65:4172–4180. doi: 10.1158/0008-5472.CAN-04-2983

- [447] G. Howard, R. Eiges, F. Gaudet, R. Jaenisch, and A. Eden. 2008. Activation and transposition of endogenous retroviral elements in hypomethylation induced tumors in mice. *Oncogene*, 27:404–408. doi: 10.1038/sj.onc.1210631
- [448] R. C. Iskow, M. T. McCabe, R. E. Mills, S. Torene, W. S. Pittard, A. F. Neuwald, E. G. Van Meir, P. M. Vertino, and S. E. Devine. 2010. Natural mutagenesis of human genomes by endogenous retrotransposons. *Cell*, 141:1253–1261. doi: 10.1016/j.cell.2010.05.020
- [449] E. Lee, R. Iskow, L. Yang, O. Gokcumen, P. Haseley, L. J. Luquette, 3rd, J. G. Lohr, C. C. Harris, L. Ding, R. K. Wilson, D. A. Wheeler, R. A. Gibbs, R. Kucherlapati, C. Lee, P. V. Kharchenko, P. J. Park, and C. G. A. R. N. . 2012. Landscape of somatic retrotransposition in human cancers. *Science*, 337:967–971. doi: 10.1126/science.1222077
- [450] S. W. Criscione, Y. Zhang, W. Thompson, J. M. Sedivy, and N. Neretti. 2014. Transcriptional landscape of repetitive elements in normal and cancer human cells. *BMC Genomics*, 15:583. doi: 10.1186/1471-2164-15-583
- [451] A. W. Whisnant, H. P. Bogerd, O. Flores, P. Ho, J. G. Powers, N. Sharova, M. Stevenson, C.-H. Chen, and B. R. Cullen. 2013. In-depth analysis of the interaction of HIV-1 with cellular microRNA biogenesis and effector mechanisms. *MBio*, 4:e000193. doi: 10.1128/mBio.00193-13
- [452] N. F. Lahens, I. H. Kavakli, R. Zhang, K. Hayer, M. B. Black, H. Dueck, A. Pizarro, J. Kim, R. Irizarry, R. S. Thomas, G. R. Grant, and J. B. Hogenesch. 2014. IVT-seq reveals extreme bias in RNA sequencing. *Genome Biol*, 15:R86. doi: 10.1186/gb-2014-15-6-r86
- [453] R. D. Hockett, J. M. Kilby, C. A. Derdeyn, M. S. Saag, M. Sillers, K. Squires, S. Chiz, M. A. Nowak, G. M. Shaw, and R. P. Bucy. 1999. Constant mean viral copy number per infected cell in tissues regardless of high, low, or undetectable plasma HIV RNA. *J Exp Med*, 189:1545–1554. doi: 10.1084/jem.189.10.1545
- [454] R. J. De Boer, R. M. Ribeiro, and A. S. Perelson. 2010. Current estimates for HIV-1 production imply rapid viral clearance in lymphoid tissues. *PLoS Comput Biol*, 6:e1000906. doi: 10.1371/journal.pcbi.1000906
- [455] T. Ikeda, J. Shibata, K. Yoshimura, A. Koito, and S. Matsushita. 2007. Recurrent HIV-1 integration at the BACH2 locus in resting CD4+ T cell populations during effective highly active antiretroviral therapy. *J Infect Dis*, 195:716–725. doi: 10.1086/510915
- [456] T. A. Wagner, S. McLaughlin, K. Garg, C. Y. K. Cheung, B. B. Larsen, S. Styrcak,

- H. C. Huang, P. T. Edlefsen, J. I. Mullins, and L. M. Frenkel. 2014. Proliferation of cells with HIV integrated into cancer genes contributes to persistent infection. *Science*, 345:570–573. doi: 10.1126/science.1256304
- [457] F. Maldarelli, X. Wu, L. Su, F. R. Simonetti, W. Shao, S. Hill, J. Spindler, A. L. Ferris, J. W. Mellors, M. F. Kearney, J. M. Coffin, and S. H. Hughes. 2014. Specific HIV integration sites are linked to clonal expansion and persistence of infected cells. *Science*, 345:179–183. doi: 10.1126/science.1254194
- [458] L. B. Cohn, I. T. Silva, T. Y. Oliveira, R. A. Rosales, E. H. Parrish, G. H. Learn, B. H. Hahn, J. L. Czartoski, M. J. McElrath, C. Lehmann, F. Klein, M. Caskey, B. D. Walker, J. D. Siliciano, R. F. Siliciano, M. Jankovic, and M. C. Nussenzweig. 2015. HIV-1 integration landscape during latent and active infection. *Cell*, 160: 420–432. doi: 10.1016/j.cell.2015.01.020
- [459] A. R. W. Schröder, P. Shinn, H. Chen, C. Berry, J. R. Ecker, and F. Bushman. 2002. HIV-1 integration in the human genome favors active genes and local hotspots. *Cell*, 110:521–529. doi: 10.1016/S0092-8674(02)00864-4
- [460] T. Brady, Y. N. Lee, K. Ronen, N. Malani, C. C. Berry, P. D. Bieniasz, and F. D. Bushman. 2009. Integration target site selection by a resurrected human endogenous retrovirus. *Genes Dev*, 23:633–642. doi: 10.1101/gad.1762309
- [461] B. Marini, A. Kertesz-Farkas, H. Ali, B. Lucic, K. Lisek, L. Manganaro, S. Pongor, R. Luzzati, A. Recchia, F. Mavilio, M. Giacca, and M. Lusic. 2015. Nuclear architecture dictates HIV-1 integration site selection. *Nature*. doi: 10.1038/nature14226
- [462] M. Cavazzana-Calvo, E. Payen, O. Negre, G. Wang, K. Hehir, F. Fusil, J. Down, M. Denaro, T. Brady et al. 2010. Transfusion independence and HMGA2 activation after gene therapy of human β -thalassaemia. *Nature*, 467:318–322. doi: 10.1038/nature09328
- [463] S. Hacein-Bey-Abina, A. Garrigue, G. P. Wang, J. Soulier, A. Lim, E. Morillon, E. Clappier, L. Caccavelli, E. Delabesse et al. 2008. Insertional oncogenesis in 4 patients after retrovirus-mediated gene therapy of SCID-X1. *J Clin Invest*, 118: 3132–3142. doi: 10.1172/JCI35700
- [464] A. Moiani, Y. Paleari, D. Sartori, R. Mezzadra, A. Miccio, C. Cattoglio, F. Cochiarella, M. R. Lidonni, G. Ferrari, and F. Mavilio. 2012. Lentiviral vector integration in the human genome induces alternative splicing and generates aberrant transcripts. *J Clin Invest*, 122:1653–1666. doi: 10.1172/JCI61852
- [465] D. Cesana, J. Sgualdino, L. Rudilosso, S. Merella, L. Naldini, and E. Montini. 2012. Whole transcriptome characterization of aberrant splicing events induced by lentiviral vector integrations. *J Clin Invest*, 122:1667–1676. doi: 10.1172/JCI62189

- [466] S. Pääbo, D. M. Irwin, and A. C. Wilson. 1990. DNA damage promotes jumping between templates during enzymatic amplification. *J Biol Chem*, 265:4718–4721
- [467] S. J. Odelberg, R. B. Weiss, A. Hata, and R. White. 1995. Template-switching during DNA synthesis by *Thermus aquaticus* DNA polymerase I. *Nucleic Acids Res*, 23:2049–2057. doi: 10.1093/nar/23.11.2049
- [468] X.-C. Zeng and S.-X. Wang. 2002. Evidence that BmTXK beta-BmKCT cDNA from Chinese scorpion *Buthus martensi Karsch* is an artifact generated in the reverse transcription process. *FEBS Lett*, 520:183–4; author reply 185
- [469] B. Tasic, C. E. Nabholz, K. K. Baldwin, Y. Kim, E. H. Rueckert, S. A. Ribich, P. Cramer, Q. Wu, R. Axel, and T. Maniatis. 2002. Promoter choice determines splice site selection in protocadherin alpha and gamma pre-mRNA splicing. *Mol Cell*, 10:21–33
- [470] M. Geiszt, K. Lekstrom, and T. L. Leto. 2004. Analysis of mRNA transcripts from the NAD(P)H oxidase 1 (Nox1) gene. Evidence against production of the NADPH oxidase homolog-1 short (NOH-1S) transcript variant. *J Biol Chem*, 279: 51661–51668. doi: 10.1074/jbc.M409325200
- [471] J. Cocquet, A. Chong, G. Zhang, and R. A. Veitia. 2006. Reverse transcriptase template switching and false alternative transcripts. *Genomics*, 88:127–131. doi: 10.1016/j.ygeno.2005.12.013
- [472] C. J. McManus, J. D. Coolon, M. O. Duff, J. Eipper-Mains, B. R. Graveley, and P. J. Wittkopp. 2010. Regulatory divergence in *Drosophila* revealed by mRNA-seq. *Genome Res*, 20:816–825. doi: 10.1101/gr.102491.109
- [473] B. Cogné, R. Snyder, P. Lindenbaum, J.-B. Dupont, R. Redon, P. Moullier, and A. Leger. 2014. NGS library preparation may generate artifactual integration sites of AAV vectors. *Nat Med*, 20:577–578. doi: 10.1038/nm.3578
- [474] E. Gilboa, S. W. Mitra, S. Goff, and D. Baltimore. 1979. A detailed model of reverse transcription and tests of crucial aspects. *Cell*, 18:93–100. doi: 10.1016/0092-8674(79)90357-X
- [475] G. X. Luo and J. Taylor. 1990. Template switching by reverse transcriptase during DNA synthesis. *J Virol*, 64:4321–4328
- [476] J. Houseley and D. Tollervey. 2010. Apparent non-canonical trans-splicing is generated by reverse transcriptase in vitro. *PLoS One*, 5:e12271. doi: 10.1371/journal.pone.0012271
- [477] A. Meyerhans, J. P. Vartanian, and S. Wain-Hobson. 1990. DNA recombination during PCR. *Nucleic Acids Res*, 18:1687–1691

- [478] D. J. G. Lahr and L. A. Katz. 2009. Reducing the impact of PCR-mediated recombination in molecular evolution and environmental studies using a new-generation high-fidelity DNA polymerase. *Biotechniques*, 47:857–866. doi: 10.2144/000113219
- [479] W. Al-Ahmadi, L. Al-Haj, F. A. Al-Mohanna, R. H. Silverman, and K. S. A. Khabar. 2009. RNase L downmodulation of the RNA-binding protein, HuR, and cellular growth. *Oncogene*, 28:1782–1791. doi: 10.1038/onc.2009.16
- [480] R. B. Jones, K. E. Garrison, S. Mujib, V. Mihajlovic, N. Aidarus, D. V. Hunter, E. Martin, V. M. John, W. Zhan et al. 2012. HERV-K-specific T cells eliminate diverse HIV-1/2 and SIV primary isolates. *J Clin Invest*, 122:4473–4489. doi: 10.1172/JCI64560
- [481] K. Boller, O. Janssen, H. Schuldes, R. R. Tönjes, and R. Kurth. 1997. Characterization of the antibody response specific for the human endogenous retrovirus HTDV/HERV-K. *J Virol*, 71:4581–4588
- [482] K. E. Garrison, R. B. Jones, D. A. Meiklejohn, N. Anwar, L. C. Ndhlovu, J. M. Chapman, A. L. Erickson, A. Agrawal, G. Spotts, F. M. Hecht, S. Rakoff-Nahoum, J. Lenz, M. A. Ostrowski, and D. F. Nixon. 2007. T cell responses to human endogenous retroviruses in HIV-1 infection. *PLoS Pathog*, 3:e165. doi: 10.1371/journal.ppat.0030165
- [483] R. Tandon, D. SenGupta, L. C. Ndhlovu, R. G. S. Vieira, R. B. Jones, V. A. York, V. A. Vieira, E. R. Sharp, A. A. Wiznia, M. A. Ostrowski, M. G. Rosenberg, and D. F. Nixon. 2011. Identification of human endogenous retrovirus-specific T cell responses in vertically HIV-1-infected subjects. *J Virol*, 85:11526–11531. doi: 10.1128/JVI.05418-11
- [484] D. SenGupta, R. Tandon, R. G. S. Vieira, L. C. Ndhlovu, R. Lown-Hecht, C. E. Ormsby, L. Loh, R. B. Jones, K. E. Garrison, J. N. Martin, V. A. York, G. Spotts, G. Reyes-Terán, M. A. Ostrowski, F. M. Hecht, S. G. Deeks, and D. F. Nixon. 2011. Strong human endogenous retrovirus-specific T cell responses are associated with control of HIV-1 in chronic infection. *J Virol*, 85:6977–6985. doi: 10.1128/JVI.00179-11
- [485] W. Pi, Z. Yang, J. Wang, L. Ruan, X. Yu, J. Ling, S. Krantz, C. Isales, S. J. Conway, S. Lin, and D. Tuan. 2004. The LTR enhancer of ERV-9 human endogenous retrovirus is active in oocytes and progenitor cells in transgenic zebrafish and humans. *Proc Natl Acad Sci U S A*, 101:805–810. doi: 10.1073/pnas.0307698100
- [486] X. H.-F. Zhang and L. A. Chasin. 2006. Comparison of multiple vertebrate genomes reveals the birth and evolution of human exons. *Proc Natl Acad Sci USA*, 103:13427–13432. doi: 10.1073/pnas.0603042103
- [487] F. A. Santoni, J. Guerra, and J. Luban. 2012. HERV-H RNA is abundant in

- human embryonic stem cells and a precise marker for pluripotency. *Retrovirology*, 9:111. doi: 10.1186/1742-4690-9-111
- [488] N. V. Fuchs, S. Loewer, G. Q. Daley, Z. Izsvák, J. Löwer, and R. Löwer. 2013. Human endogenous retrovirus K (HML-2) RNA and protein expression is a marker for human embryonic and induced pluripotent stem cells. *Retrovirology*, 10:115. doi: 10.1186/1742-4690-10-115
- [489] A. Fort, K. Hashimoto, D. Yamada, M. Salimullah, C. A. Keya, A. Saxena, A. Bonetti, I. Voineagu, N. Bertin et al. 2014. Deep transcriptome profiling of mammalian stem cells supports a regulatory role for retrotransposons in pluripotency maintenance. *Nat Genet*, 46:558–566. doi: 10.1038/ng.2965
- [490] J. Wang, G. Xie, M. Singh, A. T. Ghanbarian, T. Raskó, A. Szvetnik, H. Cai, D. Besser, A. Prigione, N. V. Fuchs, G. G. Schumann, W. Chen, M. C. Lorincz, Z. Ivics, L. D. Hurst, and Z. Izsvák. 2014. Primate-specific endogenous retrovirus-driven transcription defines naive-like stem cells. *Nature*, 516:405–409. doi: 10.1038/nature13804
- [491] B. Joos, M. Fischer, H. Kuster, S. K. Pillai, J. K. Wong, J. Böni, B. Hirscherl, R. Weber, A. Trkola, H. F. Günthard, and S. H. I. V. C. S. . 2008. HIV rebounds from latently infected cells, rather than from continuing low-level replication. *Proc Natl Acad Sci U S A*, 105:16725–16730. doi: 10.1073/pnas.0804192105
- [492] T. P. Brennan, J. O. Woods, A. R. Sedaghat, J. D. Siliciano, R. F. Siliciano, and C. O. Wilke. 2009. Analysis of human immunodeficiency virus type 1 viremia and provirus in resting CD4+ T cells reveals a novel source of residual viremia in patients on antiretroviral therapy. *J Virol*, 83:8470–8481. doi: 10.1128/JVI.02568-08
- [493] T. A. Wagner, J. L. McKernan, N. H. Tobin, K. A. Tapia, J. I. Mullins, and L. M. Frenkel. 2013. An increasing proportion of monotypic HIV-1 DNA sequences during antiretroviral treatment suggests proliferation of HIV-infected cells. *J Virol*, 87:1770–1778. doi: 10.1128/JVI.01985-12
- [494] M. F. Kearney, J. Spindler, W. Shao, S. Yu, E. M. Anderson, A. O’Shea, C. Rehm, C. Poethke, N. Kovacs, J. W. Mellors, J. M. Coffin, and F. Maldarelli. 2014. Lack of detectable HIV-1 molecular evolution during suppressive antiretroviral therapy. *PLoS Pathog*, 10:e1004010. doi: 10.1371/journal.ppat.1004010
- [495] C. J. L. Murray, K. F. Ortblad, C. Guinovart, S. S. Lim, T. M. Wolock, D. A. Roberts, E. A. Dansereau, N. Graetz, R. M. Barber et al. 2014. Global, regional, and national incidence and mortality for HIV, tuberculosis, and malaria during 1990–2013: a systematic analysis for the Global Burden of Disease Study 2013. *Lancet*, 384:1005–1070. doi: 10.1016/S0140-6736(14)60844-8
- [496] K. A. Sollis, P. W. Smit, S. Fiscus, N. Ford, M. Vitoria, S. Essajee, D. Barnett,

- B. Cheng, S. M. Crowe, T. Denny, A. Landay, W. Stevens, V. Habiyambere, J. Perrins, and R. W. Peeling. 2014. Systematic review of the performance of HIV viral load technologies on plasma samples. *PLoS One*, 9:e85869. doi: 10.1371/journal.pone.0085869
- [497] C. Liu, M. Mauk, R. Gross, F. D. Bushman, P. H. Edelstein, R. G. Collman, and H. H. Bau. 2013. Membrane-based, sedimentation-assisted plasma separator for point-of-care applications. *Anal Chem*, 85:10463–10470. doi: 10.1021/ac402459h
- [498] K. A. Curtis, D. L. Rudolph, I. Nejad, J. Singleton, A. Beddoe, B. Weigl, P. LaBarre, and S. M. Owen. 2012. Isothermal amplification using a chemical heating device for point-of-care detection of HIV-1. *PLoS One*, 7:e31432. doi: 10.1371/journal.pone.0031432
- [499] T. Notomi, H. Okayama, H. Masubuchi, T. Yonekawa, K. Watanabe, N. Amino, and T. Hase. 2000. Loop-mediated isothermal amplification of DNA. *Nucleic Acids Res*, 28:E63
- [500] K. A. Curtis, D. L. Rudolph, and S. M. Owen. 2008. Rapid detection of HIV-1 by reverse-transcription, loop-mediated isothermal amplification (RT-LAMP). *J Virol Methods*, 151:264–270. doi: 10.1016/j.jviromet.2008.04.011
- [501] K. A. Curtis, D. L. Rudolph, and S. M. Owen. 2009. Sequence-specific detection method for reverse transcription, loop-mediated isothermal amplification of HIV-1. *J Med Virol*, 81:966–972. doi: 10.1002/jmv.21490
- [502] Y. Zeng, X. Zhang, K. Nie, X. Ding, B. Z. Ring, L. Xu, L. Dai, X. Li, W. Ren, L. Shi, and X. Ma. 2014. Rapid quantitative detection of Human immunodeficiency virus type 1 by a reverse transcription-loop-mediated isothermal amplification assay. *Gene*, 541:123–128. doi: 10.1016/j.gene.2014.03.015
- [503] N. Hosaka, N. Ndembí, A. Ishizaki, S. Kageyama, K. Numazaki, and H. Ichimura. 2009. Rapid detection of human immunodeficiency virus type 1 group M by a reverse transcription-loop-mediated isothermal amplification assay. *J Virol Methods*, 157:195–199. doi: 10.1016/j.jviromet.2009.01.004
- [504] K. A. Curtis, P. L. Niedzwiedz, A. S. Youngpairoj, D. L. Rudolph, and S. M. Owen. 2014. Real-time detection of HIV-2 by reverse transcription-loop-mediated isothermal amplification. *J Clin Microbiol*, 52:2674–2676. doi: 10.1128/JCM.00935-14
- [505] C. Kuiken, H. Yoon, W. Abfalsterer, B. Gaschen, C. Lo, and B. Korber. 2013. Viral genome analysis and knowledge management. *Methods Mol Biol*, 939:253–261. doi: 10.1007/978-1-62703-107-3_16
- [506] M. Manak, S. Sina, B. Anekella, I. Hewlett, E. Sanders-Buell, V. Ragupathy, J. Kim, M. Vermeulen, S. L. Stramer, E. Sabino, P. Grabarczyk, N. Michael, S. Peel,

- P. Garrett, S. Tovanabutra, M. P. Busch, and M. Schito. 2012. Pilot studies for development of an HIV subtype panel for surveillance of global diversity. *AIDS Res Hum Retroviruses*, 28:594–606. doi: 10.1089/AID.2011.0271
- [507] J. Louwagie, F. E. McCutchan, M. Peeters, T. P. Brennan, E. Sanders-Buell, G. A. Eddy, G. van der Groen, K. Fransen, G. M. Gershay-Damet, and R. Deley. 1993. Phylogenetic analysis of gag genes from 70 international HIV-1 isolates provides evidence for multiple genotypes. *AIDS*, 7:769–780
- [508] L. Buonaguro, M. L. Tornesello, and F. M. Buonaguro. 2007. Human immunodeficiency virus type 1 subtype distribution in the worldwide epidemic: pathogenetic and therapeutic implications. *J Virol*, 81:10209–10219. doi: 10.1128/JVI.00872-07
- [509] N. F. Parrish, F. Gao, H. Li, E. E. Giorgi, H. J. Barbian, E. H. Parrish, L. Zajic, S. S. Iyer, J. M. Decker et al. 2013. Phenotypic properties of transmitted founder HIV-1. *Proc Natl Acad Sci U S A*, 110:6626–6633. doi: 10.1073/pnas.1304288110
- [510] A. G. Abimiku, T. L. Stern, A. Zwandor, P. D. Markham, C. Calef, S. Kyari, W. C. Saxinger, R. C. Gallo, M. Robert-Guroff, and M. S. Reitz. 1994. Subgroup G HIV type 1 isolates from Nigeria. *AIDS Res Hum Retroviruses*, 10:1581–1583
- [511] Zhang, Chung, and Oldenburg. 1999. A simple statistical parameter for use in evaluation and validation of high throughput screening assays. *J Biomol Screen*, 4:67–73. doi: 10.1177/108705719900400206
- [512] C. Liu, E. Geva, M. Mauk, X. Qiu, W. R. Abrams, D. Malamud, K. Curtis, S. M. Owen, and H. H. Bau. 2011. An isothermal amplification reactor with an integrated isolation membrane for point-of-care detection of infectious diseases. *Analyst*, 136:2069–2076. doi: 10.1039/c1an00007a
- [513] C. A. Spina, J. Anderson, N. M. Archin, A. Bosque, J. Chan, M. Famiglietti, W. C. Greene, A. Kashuba, S. R. Lewin et al. 2013. An in-depth comparison of latent HIV-1 reactivation in multiple cell model systems and resting CD4+ T cells from aviremic patients. *PLoS Pathog*, 9:e1003834. doi: 10.1371/journal.ppat.1003834
- [514] G. Lehrman, I. B. Hogue, S. Palmer, C. Jennings, C. A. Spina, A. Wiegand, A. L. Landay, R. W. Coombs, D. D. Richman, J. W. Mellors, J. M. Coffin, R. J. Bosch, and D. M. Margolis. 2005. Depletion of latent HIV-1 infection in vivo: a proof-of-concept study. *Lancet*, 366:549–555. doi: 10.1016/S0140-6736(05)67098-5
- [515] N. M. Archin, M. Cheema, D. Parker, A. Wiegand, R. J. Bosch, J. M. Coffin, J. Eron, M. Cohen, and D. M. Margolis. 2010. Antiretroviral intensification and valproic acid lack sustained effect on residual HIV-1 viremia or resting CD4+ cell infection. *PLoS One*, 5:e9390. doi: 10.1371/journal.pone.0009390
- [516] A. L. Hill, D. I. S. Rosenbloom, F. Fu, M. A. Nowak, and R. F. Siliciano. 2014.

Predicting the outcomes of treatment to eradicate the latent reservoir for HIV-1. *Proc Natl Acad Sci U S A*, 111:13475–13480. doi: 10.1073/pnas.1406663111

- [517] ENCODE Project Consortium. 2012. An integrated encyclopedia of DNA elements in the human genome. *Nature*, 489:57–74. doi: 10.1038/nature11247
- [518] T. Barrett, S. E. Wilhite, P. Ledoux, C. Evangelista, I. F. Kim, M. Tomashevsky, K. A. Marshall, K. H. Phillippy, P. M. Sherman, M. Holko, A. Yefanov, H. Lee, N. Zhang, C. L. Robertson, N. Serova, S. Davis, and A. Soboleva. 2013. NCBI GEO: archive for functional genomics data sets—update. *Nucleic Acids Res*, 41: D991–D995. doi: 10.1093/nar/gks1193
- [519] D. Karolchik, G. P. Barber, J. Casper, H. Clawson, M. S. Cline, M. Diekhans, T. R. Dreszer, P. A. Fujita, L. Guruvadoo et al. 2014. The UCSC Genome Browser database: 2014 update. *Nucleic Acids Res*, 42:D764–D770. doi: 10.1093/nar/gkt1168
- [520] M. Goldman, B. Craft, T. Swatloski, M. Cline, O. Morozova, M. Diekhans, D. Haussler, and J. Zhu. 2015. The UCSC Cancer Genomics Browser: update 2015. *Nucleic Acids Res*, 43:D812–D817. doi: 10.1093/nar/gku1073
- [521] F. Cunningham, M. R. Amode, D. Barrell, K. Beal, K. Billis, S. Brent, D. Carvalho-Silva, P. Clapham, G. Coates et al. 2015. Ensembl 2015. *Nucleic Acids Res*, 43: D662–D669. doi: 10.1093/nar/gku1010
- [522] M. L. Metzker. 2010. Sequencing technologies - the next generation. *Nat Rev Genet*, 11:31–46. doi: 10.1038/nrg2626
- [523] E. R. Mardis. 2011. A decade's perspective on DNA sequencing technology. *Nature*, 470:198–203. doi: 10.1038/nature09796
- [524] K. Wetterstrand. 2015. DNA Sequencing Costs: Data from the NHGRI Genome Sequencing Program (GSP). URL www.genome.gov/sequencingcosts
- [525] B. F. Keele, E. E. Giorgi, J. F. Salazar-Gonzalez, J. M. Decker, K. T. Pham, M. G. Salazar, C. Sun, T. Grayson, S. Wang et al. 2008. Identification and characterization of transmitted and early founder virus envelopes in primary HIV-1 infection. *Proc Natl Acad Sci USA*, 105:7552–7557. doi: 10.1073/pnas.0802203105
- [526] J. F. Salazar-Gonzalez, M. G. Salazar, B. F. Keele, G. H. Learn, E. E. Giorgi, H. Li, J. M. Decker, S. Wang, J. Baalwa et al. 2009. Genetic identity, biological phenotype, and evolutionary pathways of transmitted/founder viruses in acute and early HIV-1 infection. *J Exp Med*, 206:1273–1289. doi: 10.1084/jem.20090378
- [527] P. J. Kanki, D. J. Hamel, J. L. Sankalé, C. Hsieh, I. Thior, F. Barin, S. A. Woodcock, A. Guèye-Ndiaye, E. Zhang, M. Montano, T. Siby, R. Marlink, I. NDoye, M. E.

- Essex, and S. MBoup. 1999. Human immunodeficiency virus type 1 subtypes differ in disease progression. *J Infect Dis*, 179:68–73. doi: 10.1086/314557
- [528] P. Kaleebu, N. French, C. Mahe, D. Yirrell, C. Watera, F. Lyagoba, J. Nakiyingi, A. Rutebemberwa, D. Morgan, J. Weber, C. Gilks, and J. Whitworth. 2002. Effect of human immunodeficiency virus (HIV) type 1 envelope subtypes A and D on disease progression in a large cohort of HIV-1-positive persons in Uganda. *J Infect Dis*, 185:1244–1250. doi: 10.1086/340130
- [529] J. M. Baeten, B. Chohan, L. Lavreys, V. Chohan, R. S. McClelland, L. Certain, K. Mandaliya, W. Jaoko, and J. Overbaugh. 2007. HIV-1 subtype D infection is associated with faster disease progression than subtype A in spite of similar plasma HIV-1 loads. *J Infect Dis*, 195:1177–1180. doi: 10.1086/512682
- [530] N. Kiwanuka, O. Laeyendecker, M. Robb, G. Kigozi, M. Arroyo, F. McCutchan, L. A. Eller, M. Eller, F. Makumbi, D. Birx, F. Wabwire-Mangen, D. Serwadda, N. K. Sewankambo, T. C. Quinn, M. Wawer, and R. Gray. 2008. Effect of human immunodeficiency virus Type 1 (HIV-1) subtype on disease progression in persons from Rakai, Uganda, with incident HIV-1 infection. *J Infect Dis*, 197:707–713. doi: 10.1086/527416
- [531] B. Renjifo, P. Gilbert, B. Chaplin, G. Msamanga, D. Mwakagile, W. Fawzi, M. Essex, T. V. , and H. I. V. S. Group. 2004. Preferential in-utero transmission of HIV-1 subtype C as compared to HIV-1 subtype A or D. *AIDS*, 18:1629–1636
- [532] G. C. John-Stewart, R. W. Nduati, C. M. Rousseau, D. A. Mbori-Ngacha, B. A. Richardson, S. Rainwater, D. D. Panteleeff, and J. Overbaugh. 2005. Subtype C Is associated with increased vaginal shedding of HIV-1. *J Infect Dis*, 192:492–496. doi: 10.1086/431514
- [533] W. Huang, S. H. Eshleman, J. Toma, S. Fransen, E. Stawiski, E. E. Paxinos, J. M. Whitcomb, A. M. Young, D. Donnell, F. Mmiro, P. Musoke, L. A. Guay, J. B. Jackson, N. T. Parkin, and C. J. Petropoulos. 2007. Coreceptor tropism in human immunodeficiency virus type 1 subtype D: high prevalence of CXCR4 tropism and heterogeneous composition of viral populations. *J Virol*, 81:7885–7893. doi: 10.1128/JVI.00218-07
- [534] J. Snoeck, R. Kantor, R. W. Shafer, K. Van Laethem, K. Deforche, A. P. Carvalho, B. Wynhoven, M. A. Soares, P. Cane et al. 2006. Discordances between interpretation algorithms for genotypic resistance to protease and reverse transcriptase inhibitors of human immunodeficiency virus are subtype dependent. *Antimicrob Agents Chemother*, 50:694–701. doi: 10.1128/AAC.50.2.694-701.2006
- [535] P. J. Easterbrook, M. Smith, J. Mullen, S. O’Shea, I. Chrystie, A. de Ruiter, I. D. Tatt, A. M. Geretti, and M. Zuckerman. 2010. Impact of HIV-1 viral subtype on disease progression and response to antiretroviral therapy. *J Int AIDS Soc*, 13:4. doi: 10.1186/1758-2652-13-4

- [536] A. U. Scherrer, B. Ledergerber, V. von Wyl, J. Böni, S. Yerly, T. Klimkait, P. Bürgisser, A. Rauch, B. Hirscherl, M. Cavassini, L. Elzi, P. L. Vernazza, E. Bernasconi, L. Held, H. F. Günthard, and S. H. I. V. C. S. . 2011. Improved virological outcome in White patients infected with HIV-1 non-B subtypes compared to subtype B. *Clin Infect Dis*, 53:1143–1152. doi: 10.1093/cid/cir669
- [537] C. Liu, M. M. Sadik, M. G. Mauk, P. H. Edelstein, F. D. Bushman, R. Gross, and H. H. Bau. 2014. Nuclemeter: a reaction-diffusion based method for quantifying nucleic acids undergoing enzymatic amplification. *Sci Rep*, 4:7335. doi: 10.1038/srep07335
- [538] M. G. Mauk, C. Liu, M. Sadik, and H. H. Bau. 2015. Microfluidic devices for nucleic acid (NA) isolation, isothermal NA amplification, and real-time detection. *Methods Mol Biol*, 1256:15–40. doi: 10.1007/978-1-4939-2172-0_2
- [539] A. Piantadosi, B. Chohan, V. Chohan, R. S. McClelland, and J. Overbaugh. 2007. Chronic HIV-1 infection frequently fails to protect against superinfection. *PLoS Pathog*, 3:e177. doi: 10.1371/journal.ppat.0030177
- [540] R. L. R. Powell, M. M. Urbanski, S. Burda, T. Kinge, and P. N. Nyambi. 2009. High frequency of HIV-1 dual infections among HIV-positive individuals in Cameroon, West Central Africa. *J Acquir Immune Defic Syndr*, 50:84–92. doi: 10.1097/QAI.0b013e31818d5a40
- [541] K. Ronen, C. O. McCoy, F. A. Matsen, D. F. Boyd, S. Emery, K. Odem-Davis, W. Jaoko, K. Mandaliya, R. S. McClelland, B. A. Richardson, and J. Overbaugh. 2013. HIV-1 superinfection occurs less frequently than initial infection in a cohort of high-risk Kenyan women. *PLoS Pathog*, 9:e1003593. doi: 10.1371/journal.ppat.1003593
- [542] A. D. Redd, D. Ssemwanga, J. Vandepitte, S. K. Wendel, N. Ndembu, J. Bukenya, S. Nakubulwa, H. Grosskurth, C. M. Parry, C. Martens, D. Bruno, S. F. Porcella, T. C. Quinn, and P. Kaleebu. 2014. Rates of HIV-1 superinfection and primary HIV-1 infection are similar in female sex workers in Uganda. *AIDS*, 28:2147–2152. doi: 10.1097/QAD.0000000000000365
- [543] A. D. Redd, C. E. Mullis, D. Serwadda, X. Kong, C. Martens, S. M. Ricklefs, A. A. R. Tobian, C. Xiao, M. K. Grabowski, F. Nalugoda, G. Kigozi, O. Laeyendecker, J. Kagaayi, N. Sewankambo, R. H. Gray, S. F. Porcella, M. J. Wawer, and T. C. Quinn. 2012. The rates of HIV superinfection and primary HIV incidence in a general population in Rakai, Uganda. *J Infect Dis*, 206:267–274. doi: 10.1093/infdis/jis325
- [544] S. Jost, M.-C. Bernard, L. Kaiser, S. Yerly, B. Hirscherl, A. Samri, B. Autran, L.-E. Goh, and L. Perrin. 2002. A patient with HIV-1 superinfection. *N Engl J Med*, 347:731–736. doi: 10.1056/NEJMoa020263

- [545] G. Fang, B. Weiser, C. Kuiken, S. M. Philpott, S. Rowland-Jones, F. Plummer, J. Kiprani, B. Shi, R. Kaul, J. Bwayo, O. Anzala, and H. Burger. 2004. Recombination following superinfection by HIV-1. *AIDS*, 18:153–159
- [546] G. Blick, R. M. Kagan, E. Coakley, C. Petropoulos, L. Maroldo, P. Greiger-Zanlungo, S. Gretz, and T. Garton. 2007. The probable source of both the primary multidrug-resistant (MDR) HIV-1 strain found in a patient with rapid progression to AIDS and a second recombinant MDR strain found in a chronically HIV-1-infected patient. *J Infect Dis*, 195:1250–1259. doi: 10.1086/512240
- [547] G. S. Gottlieb, D. C. Nickle, M. A. Jensen, K. G. Wong, R. A. Kaslow, J. C. Shepherd, J. B. Margolick, and J. I. Mullins. 2007. HIV type 1 superinfection with a dual-tropic virus and rapid progression to AIDS: a case report. *Clin Infect Dis*, 45: 501–509. doi: 10.1086/520024
- [548] H. Streeck, B. Li, A. F. Y. Poon, A. Schneidewind, A. D. Gladden, K. A. Power, D. Daskalakis, S. Bazner, R. Zuniga, C. Brander, E. S. Rosenberg, S. D. W. Frost, M. Altfeld, and T. M. Allen. 2008. Immune-driven recombination and loss of control after HIV superinfection. *J Exp Med*, 205:1789–1796. doi: 10.1084/jem.20080281
- [549] O. Clerc, S. Colombo, S. Yerly, A. Telenti, and M. Cavassini. 2010. HIV-1 elite controllers: beware of super-infections. *J Clin Virol*, 47:376–378. doi: 10.1016/j.jcv.2010.01.013
- [550] D. M. Smith, J. K. Wong, G. K. Hightower, C. C. Ignacio, K. K. Koelsch, C. J. Petropoulos, D. D. Richman, and S. J. Little. 2005. HIV drug resistance acquired through superinfection. *AIDS*, 19:1251–1256
- [551] M. Pernas, C. Casado, R. Fuentes, M. J. Pérez-Elías, and C. López-Galíndez. 2006. A dual superinfection and recombination within HIV-1 subtype B 12 years after primoinfection. *J Acquir Immune Defic Syndr*, 42:12–18. doi: 10.1097/01.qai.0000214810.65292.73
- [552] D. L. Robertson, P. M. Sharp, F. E. McCutchan, and B. H. Hahn. 1995. Recombination in HIV-1. *Nature*, 374:124–126. doi: 10.1038/374124b0
- [553] M. H. Malim and M. Emerman. 2001. HIV-1 sequence variation: drift, shift, and attenuation. *Cell*, 104:469–472. doi: 10.1016/S0092-8674(01)00234-3
- [554] S. K. Gire, A. Goba, K. G. Andersen, R. S. G. Sealfon, D. J. Park, L. Kanneh, S. Jalloh, M. Momoh, M. Fullah et al. 2014. Genomic surveillance elucidates Ebola virus origin and transmission during the 2014 outbreak. *Science*, 345: 1369–1372. doi: 10.1126/science.1259657
- [555] WHO Ebola Response Team. 2014. Ebola virus disease in West Africa—the first 9

months of the epidemic and forward projections. *N Engl J Med*, 371:1481–1495. doi: 10.1056/NEJMoa1411100

- [556] World Health Organization. 2015. Ebola situation report: 13 May 2014. URL <http://apps.who.int/ebola/en/current-situation/ebola-situation-report-13-may-2015>
- [557] G. Chowell and H. Nishiura. 2014. Transmission dynamics and control of Ebola virus disease (EVD): a review. *BMC Med*, 12:196. doi: 10.1186/s12916-014-0196-0
- [558] A. S. Fauci. 2014. Ebola—underscoring the global disparities in health care resources. *N Engl J Med*, 371:1084–1086. doi: 10.1056/NEJMp1409494
- [559] World Health Organization. 2015. Interim guidance on the use of rapid Ebola antigen detection tests. URL <http://www.who.int/csr/resources/publications/ebola/ebola-antigen-detection/en/>
- [560] Y. Kurosaki, A. Takada, H. Ebihara, A. Grolla, N. Kamo, H. Feldmann, Y. Kawaoka, and J. Yasuda. 2007. Rapid and simple detection of Ebola virus by reverse transcription-loop-mediated isothermal amplification. *J Virol Methods*, 141:78–83. doi: 10.1016/j.jviromet.2006.11.031
- [561] T. Hoenen, D. Safronetz, A. Groseth, K. R. Wollenberg, O. A. Koita, B. Diarra, I. S. Fall, F. C. Haidara, F. Diallo et al. 2015. Mutation rate and genotype variation of Ebola virus from Mali case sequences. *Science*, 348:117–119. doi: 10.1126/science.aaa5646