1 # Illumina-based sequencing framework for accurate detection and
2 mapping of influenza virus defective interfering particle-associated
3 RNAs
4

5 Fadi G. Alnaji[1], Jessica R. Holmes[2,3], Gloria Rendon[2,3], J. Cristobal Vera[1,2], Chris
6 Fields[2,3], Brigitte E. Martin[1], and Christopher B. Brooke*[1,2]
7

8 [1]Department of Microbiology, University of Illinois at Urbana-Champaign
9 [2]Carl R. Woese Institute for Genomic Biology, University of Illinois at Urbana-
10 Champaign
11 [3]High Performance Biological Computing at the Roy J. Carver Biotechnology Center,
12 University of Illinois at Urbana-Champaign
13 * Corresponding author (cbrooke@illinois.edu)
14

## 15 Abstract

16 The mechanisms and consequences of defective interfering particle (DIP) formation
17 during influenza virus infection remain poorly understood. The development of next
18 generation sequencing (NGS) technologies has made it possible to identify large numbers
19 of DIP-associated sequences, providing a powerful tool to better understand their
20 biological relevance. However, NGS approaches pose numerous technical challenges
21 including the precise identification and mapping of deletion junctions in the presence of
22 frequent mutation and base-calling errors, and the potential for numerous experimental
23 and computational artifacts. Here we detail an Illumina-based sequencing framework and
24 bioinformatics pipeline capable of generating highly accurate and reproducible profiles of
25 DIP-associated junction sequences. We use a combination of simulated and experimental
26 control datasets to optimize pipeline performance and demonstrate the absence of
27 significant artifacts. Finally, we use this optimized pipeline to generate a high-resolution
28 profile of DIP-associated junctions produced during influenza virus infection and
29 demonstrate how this data can provide insight into mechanisms of DIP formation. This
30 work highlights the specific challenges associated with NGS-based detection of DIP-
31 associated sequences, and details the computational and experimental controls required
32 for such studies.
33

## 34 Importance

35 Influenza virus defective interfering particles (DIPs) that harbor internal deletions within
36 their genomes occur naturally during infection in humans and cell culture. They have been
37 hypothesized to influence the pathogenicity of the virus; however, their specific function
38 remains elusive. The accurate detection of DIP-associated deletion junctions is crucial for
39 understanding DIP biology but is complicated by an array of technical issue that can bias
40 or confound results. Here we demonstrate a combined experimental and computational
41 framework for detecting DIP-associated deletion junctions using next generation
42 sequencing (NGS). We detail how to validate pipeline performance and provide the
43 bioinformatics pipeline for groups interested in using it. Using this optimized pipeline, we

44  detect hundreds of distinct deletion junctions generated during IAV infection, and use
45  these data to test a long-standing hypothesis concerning the molecular details of DIP
46  formation.
47
48  **INTRODUCTION**
49  Influenza A virus (IAV) DIPs were first described over 60 years ago, and are classically
50  defined by their ability to interfere with the production of wild-type virus(1, 2). This ability
51  has been linked to the ability of DI RNAs to both outcompete wild-type (WT) genomic
52  RNAs for resources and packaging into virions, as well as to more potently stimulate the
53  induction of anti-viral immunity through cytosolic RNA sensors (3–6). DIPs have also been
54  implicated in influencing the outcome of influenza virus infection in humans(7). The
55  specific mechanisms and broader functional consequences of DIP formation during IAV
56  infection remain poorly understood.
57
58  IAV DIPs are characterized by the presence of large internal deletions in one or more
59  genome segments that disrupt essential open reading frames while retaining the
60  sequences required for replication and packaging(5). As such, the mapping of DIP-
61  associated deletions has helped to define the minimum sequences required for genome
62  replication and packaging (8, 9). These deletions are believed to result from a poorly
63  defined process by which the viral RNA-dependent RNA polymerase (RdRp) ceases RNA
64  polymerization at one site of the viral RNA template (donor site), only to resume at another
65  site downstream (acceptor site), resulting in a failure to copy an internal stretch of the WT
66  template (10). Until recently, the ability to characterize these DIP-associated deletion
67  junction sites (breakpoints) has been limited based on the need to clone and Sanger
68  sequence individual DIP-associated RNAs. As a result, the number of individual DIP-
69  associated RNA sequences that have been analyzed has been relatively small, hindering
70  efforts to define the factors that govern DIP deletion formation.
71
72  The advent of next generation sequencing (NGS) has increased the number of individual
73  recombinant sequences that can be identified within a given sample by orders of
74  magnitude. However, the identification and analysis of DIP-associated RNAs by NGS
75  poses new challenges, including the successful alignment of junction-containing (or
76  junction-spanning) reads to the viral reference sequence, the precise definition and
77  localization of DIP-associated deletion breakpoints, and the differentiation of true DIP
78  deletion sequences from the artifactual recombinants that can form during reverse
79  transcription, PCR, and/or sequencing. Without careful optimization and validation, these
80  issues can easily compromise efforts to define the genetic profile of DIP populations.
81
82  Here, we describe the development and validation of an Illumina-based sequencing
83  framework for the identification and analysis of influenza virus DIP-associated deletion
84  junctions. The bioinformatics pipeline combines the Bowtie 2 alignment algorithm with the
85  ViReMa (Virus Recombination Mapper) algorithm developed by Andrew Routh and a
86  collection of additional scripts for data processing and analysis (11, 12). We used
87  simulated NGS datasets and a panel of experimental control samples to optimize and

88  quantify the sensitivity, precision, and reproducibility of our pipeline. Subsequently, we
89  used the optimized pipeline to fine-tune the experimental protocol from sample
90  preparation to RNA sequencing to better detect and map DIP-associated deletions
91  generated during experimental IAV infection. This work highlights the computational and
92  experimental controls needed for Illumina-based NGS studies of viral recombination, and
93  provides an optimized, user-friendly sequencing and bioinformatics pipeline for the
94  identification and analysis of IAV DIP-associated sequences. Higher resolution analysis
95  of these deletion sequences can shed light on both the specific molecular mechanisms
96  of DIP formation, as well as how DIPs may affect the overall behavior of viral populations.
97

## RESULTS

### Overview of the pipeline

100  The sequencing framework we describe here encompasses sample preparation,
101  sequencing, and data analysis (**Fig 1A**). In brief, we generate 8-segment, full-length
102  amplicons from viral samples and sequence these using the Illumina MiSeq sequencing
103  platform. Datasets are quality-filtered and aligned to the viral reference genome using
104  Bowtie 2 in a conservative manner that disallows soft clipping. Thus, reads containing
105  deletion junctions fail to align, and are fed into the ViReMa algorithm to detect DIP-
106  associated deletion junctions. Finally, the identified junctions are mapped to the viral
107  genome and output as a matrix containing the segment name, junction start and end sites,
108  and NGS read support that can easily be analyzed using additional software tools. Below,
109  we outline the approaches we have taken to optimize and validate the various steps in
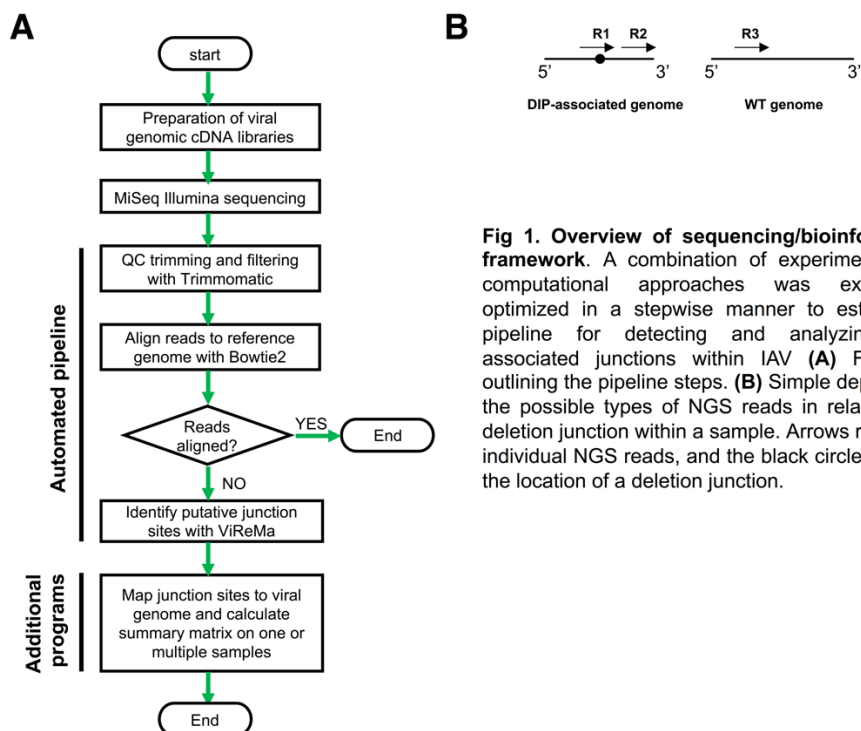110  the process.
111



**Fig 1. Overview of sequencing/bioinformatics framework**. A combination of experimental and computational approaches was extensively optimized in a stepwise manner to establish a pipeline for detecting and analyzing DIP-associated junctions within IAV **(A)** Flowchart outlining the pipeline steps. **(B)** Simple depiction of the possible types of NGS reads in relation to a deletion junction within a sample. Arrows represent individual NGS reads, and the black circle denotes the location of a deletion junction.

114 **Optimization of analysis pipeline using simulated data**
115 All bioinformatic pipelines have the potential to introduce artifacts and biases during data
116 analysis. Therefore, we first aimed to optimize the sensitivity and precision of our
117 bioinformatics pipeline using simulated NGS datasets where we absolutely know the
118 identity and frequency of all DIP-associated deletion sequences present. IAV DIP-
119 associated deletions can be found in nearly all (if not all) genome segments at a wide
120 range of frequencies (13, 14). To mimic this natural variation, we used MetaSim to
121 generate a panel of Illumina MiSeq-based NGS simulated datasets that contain DIP-
122 associated deletions in all genome segments at varying frequencies and locations (see
123 **Table 1, Fig S1**). We used a simple Perl script to randomly generate deletion junctions
124 within the terminal ~600nts of A/California/07/09 (Cal07), since these regions have been
125 shown to be hotspots for DIP-associated deletions(9, 13, 15). We also generated a
126 negative control dataset that lacks deletions to quantify the occurrence of false positives
127 generated by the pipeline. Critically, we introduced a nucleotide substitution frequency of
128 ~1% into these datasets, based on the published Illumina MiSeq empirical error model(16,
129 17). Each dataset comprised ~1million 2x250nts paired-end reads, mirroring the read
130 depth that we expect per sample on a typical sequencing run.
131

| Dataset name | Total paired-end reads | Total junction count | Total junction NGS reads | Total WT NGS individual reads | Junction count | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | | PB2 | PB1 | PA | HA | NP | NA | M | NS |
| Cal07-400 | ~1 million (2X250). Total 2 million individual reads | 400 | 1838898 | 116548 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 |
| Cal07-200 | | 200 | 1774920 | 225080 | 25 | 25 | 25 | 25 | 25 | 25 | 25 | 25 |
| Cal07 | | 0 | 0 | 2000000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

132 Table 1. Description of the simulated datasets used in this study.
133

134 **Optimization of alignment**
135 We first optimized the filtering of reads that contain deletion junctions (**Fig 1B,** R1), from
136 those that don't include junctions (**Fig 1B,** R2,R3). To do this, we aligned all reads to the
137 WT reference genome using Bowtie 2. Reads that successfully align should not contain
138 deletion junctions and are saved for further analysis, while reads that fail to align are fed
139 through the ViReMa algorithm. The performance of this alignment step is highly
140 dependent upon the mismatch penalty scores that are used during alignment. If mismatch
141 penalties are too stringent, reads with random mutations or base calling errors will fail to
142 align and be sent to ViReMa, increasing both the chances of false positives and the total
143 computational time per sample; too lenient, and true junction-spanning reads will
144 successfully align and be excluded from downstream analysis.
145
146 We used a junction-rich simulated dataset (Cal07-400) to test the effects of varying the
147 alignment penalty score on the output of ViReMa (**Fig 2A**). We observed that a penalty
148 score of 0.3 minimized the number of unaligned reads (and thus potential for false
149 positives) without diminishing the number junction-spanning reads detected. This value
150 was used for all subsequent analysis.

151

## Optimization of ViReMa operation

153 We next optimized the sensitivity and precision with which the pipeline detects deletion
154 junctions. The ability of ViReMa to accurately map true junction-containing reads is
155 affected by three factors. The first is the method the algorithm uses to identify breakpoints.
156 ViReMa extracts and aligns a seed sequence of 20-30nts (the default value of 25 was
157 used in this study) from the beginning of each read and begins aligning the downstream
158 nucleotides. If at any point the downstream alignment fails (as would be the case for a
159 deletion breakpoint), ViReMa generates a new seed sequence starting from that location
160 for realignment. Thus, breakpoints cannot be detected if they occur within the terminal 25
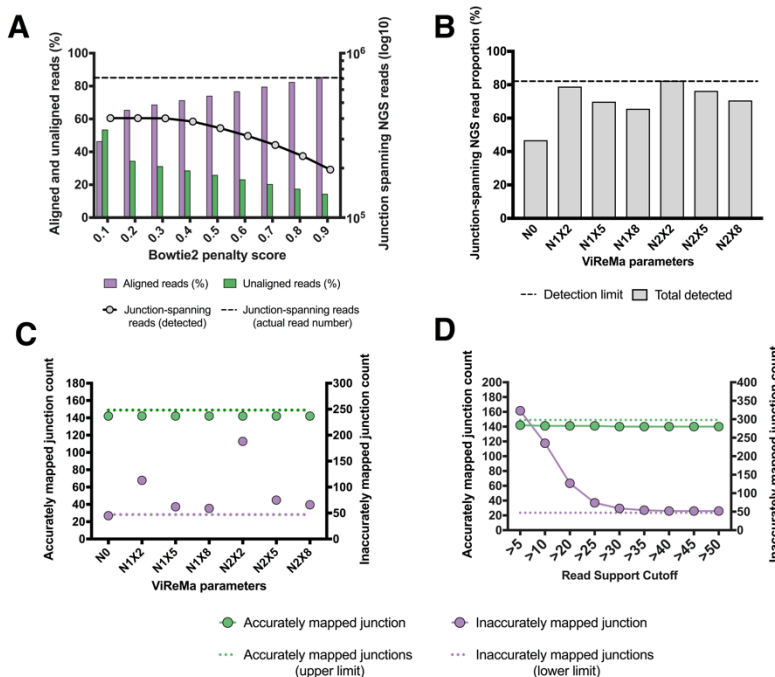161 nts of a read.



**Fig 2. Optimization of bioinformatics pipeline using simulated data. (A)** Quantification of the effects of varying the bowtie 2 penalty score on the number of junction spanning reads detected by ViReMa in the Cal07-400 simulated dataset (black line; dashed line represents the actual number of junction-spanning reads present in the dataset). The percentages of reads that aligned to reference genome (purple) and failed to align (green) are also shown for each penalty score. **(B)** The effects of ViReMa --X and --N parameters on the percentage of junction spanning reads present in the simulated dataset that were successfully detected. Dashed line shows the maximum theoretical sensitivity (~81.5%), based on the ViReMa seed length of 25nts. **(C)** The effects of ViReMa --X and --N parameters on the number of accurately (green) and inaccurately mapped (purple) deletion junctions reported by the pipeline using the Cal07-200 simulated dataset. The maximum possible number of accurate junctions and the minimum number of inaccurate junctions (resulting from junctions adjacent to direct repeat sequences) are shown for comparison **(D)** Effects of varying the minimum read support cutoff (RSC) on junction detection. Analysis performed on the Cal07-200 simulated dataset using N1X8 ViReMa values.

162

163 The second factor is the presence of short direct repeats adjacent to the junction site.
164 These repeats result in a situation where multiple potential breakpoints can give rise to
165 the same final sequence, making precise definition of the true breakpoints impossible (**Fig**
166 **S2 and S3A**). ViReMa deals with these 'fuzzy' regions through the parameter 'Defuzz',
167 which can be set to report the junction either to the 5' end, 3' end, or the middle of the
168 ambiguous region. For consistency's sake, we pushed all fuzzy junctions towards the 3'
169 of the ambiguous region. The effects of direct repeats on breakpoint mapping are
170 impossible to avoid and vary somewhat between IAV genome segments. Importantly,
171 while this effect reduces the precision of breakpoint mapping, it does not affect the ability
172 of the pipeline to determine the actual sequences of DIP-associated RNAs.

173 The third factor is the potential for base calling errors or mutations to result in erroneous
174 junction mapping (**Fig S3B**). Even though reported junctions in this category are derived
175 from real junctions, they can be viewed as false positives in that they are reported as

176  distinct junctions that do not actually exist in the viral population. Altogether, these three
177  factors set a ceiling on the maximum number of deletion junctions that can be accurately
178  detected and mapped. Using our simulated datasets, we knew how many deletion
179  junctions actually existed, exactly where they were located, and whether or not they were
180  adjacent to direct repeats (see Materials and Method) that could result in incorrect
181  mapping. This allowed us to systematically optimize the sensitivity and precision of the
182  software pipeline.
183
184  We tested how varying the ViReMa operating parameters affected both junction-spanning
185  read detection and actual junction reporting. We used the Cal07-200 dataset to challenge
186  ViReMa across a range of --N parameter (number of mismatches allowed) and --X
187  parameter (mismatch distance from the putative junction location) values. We first asked
188  how varying the --N and --X parameters influenced the total number of junction-spanning
189  reads detected (**Fig 2B**). We found that using N=0 (--X is irrelevant at this condition)
190  significantly decreased the number of junction-spanning reads detected compared with
191  non-zero --N and --X values. We next asked how increasing the --N and --X values
192  affected the number of accurately and inaccurately mapped junctions reported (**Fig 2C**).
193  We observed a clear correlation between the --X parameter and junction-mapping
194  precision, as increasing the --X value decreased the number of inaccurately mapped
195  junctions. Overall, we found that using N=1 and X=8 reduced inaccurate junction mapping
196  to the minimum amount possible, given the occurrence of direct repeats adjacent to
197  23.5% (47 of 200) of junctions in the dataset.
198
199  We next asked whether setting a minimum read support cutoff (RSC) to report a junction
200  affected the numbers of both accurate and inaccurate junctions that the pipeline identified.
201  Requiring that a given junction be represented within a minimum number of reads can
202  decrease the number of erroneously mapped junctions arising from base calling errors
203  but could also result in some true junctions being lost due to insufficient read coverage.
204  We aligned our simulated Cal07-200 dataset with Bowtie 2 and used the resulting
205  unaligned reads to challenge ViReMa using different RSC values (**Fig 2D**). We found that
206  the number of true junctions reported by the pipeline was very close to the theoretical
207  maximum, with minimal drop-off across the range of RSCs tested. In contrast, we
208  observed that the number of inaccurately reported junctions was highly sensitive to the
209  RSC value used. An RSC of >30 was needed to lower the number of inaccurately reported
210  junctions to the minimal limit (determined by the number of 'fuzzy' junctions with adjacent
211  direct repeats in the dataset).
212
213  Altogether, these data highlight the importance of optimizing RSC values and the ViReMa
214  --N and --X parameters for maximizing the sensitivity of junction detection while
215  minimizing the number of false positives. We set our default values at RSC>30, --N=1,
216  and --X=8 for subsequent analysis.
217
218  **Validation of sequencing pipeline**
219  After optimizing the bioinformatics component of our pipeline using simulated datasets,

220  we examined the ability of the pipeline to detect DIP-associated deletions within complex
221  viral populations from experimental samples. Our overall strategy was based on the
222  universal, eight segment RT-PCR approach pioneered by Zhou et al. (18). Critically, there
223  are a number of steps within the library preparation and sequencing steps that have the
224  potential to introduce artifacts that can compromise junction detection and analysis. In
225  particular, we were concerned about the potential for recombination during reverse
226  transcription, PCR, and/or sequencing to generate junctions that will be called by the
227  pipeline (19, 20). To address this, we prepared several control sample libraries,
228  sequenced them on the MiSeq, and ran the results through our optimized pipeline.

230  To quantify false positive generation during the PCR and/or sequencing steps, we
231  constructed libraries without using actual viral RNA or reverse transcriptase. To do this,
232  we generated an equimolar ratio mixture of full length PCR amplicons from each of the
233  eight IAV genome segments, using reverse genetics plasmids encoding the gene
234  segments from A/Puerto Rico/8/1934 (PR8) as templates. These amplicons were gel
235  purified to ensure correct, full-length size, and then used as template for the universal
236  amplification PCR and subsequent library preparation. Our analysis pipeline detected no
237  breakpoints in this control, indicating that none of the steps in our pipeline from PCR
238  onwards were significant sources of false positive signals.

240  We next sequenced a recombinant Cal07 stock that was grown under low MOI conditions
241  to minimize the frequency of DIPs (21). We performed two independent RNA extractions
242  and reverse transcription reactions on this stock to serve as technical replicates (named
243  Par1 and Par2). ViReMa detected 6 and 7 DIP-associated deletion junctions from Par1
244  and Par2, respectively, with junction-spanning reads representing ~0.1-0.2% of the total
245  reads **(Fig 3A)**. The majority of these reads were derived from a single shared deletion
246  junction in HA (indicated by the following nomenclature: 615_1132_HA). 4 other DIP
247  junctions were shared between replicates, each with low NGS read depth (ranging
248  between 19 and 94). Two unshared junctions in Par1 and one in Par2 were actually
249  reported in both replicates but failed to reach the level of detection in one replicate.

251  The significant overlap in the specific junctions that were reported from the two replicates
252  suggested that these junctions were produced by the viral polymerase (and were thus
253  bona fide DIP-associated sequences) rather than by the reverse transcriptase. However,
254  the generation of the same junction in independent RT reactions could also indicate the
255  existence of strong hotspots for RT recombination. To more directly address the potential
256  contribution of RT-derived recombinants, we performed two independent experiments.
257  First, we compared the junctions detected in HA segment libraries generated from Par 1
258  using two different RT enzymes, Invitrogen Superscript III and Agilent AccuScript.
259  Second, we performed *in vitro* transcription of a plasmid-derived Cal07 HA segment using
260  T7 RNA polymerase, which then was used as a template for RT-PCR to produce the
261  amplicon library for sequencing. The IAV polymerase was not involved in this control; thus
262  any deletions detected will have been generated by T7 polymerase or the RT enzyme.

264    The junctions reported from libraries generated by the two RT enzymes had significant
265    overlap and were both dominated by 615_1132_HA (**Fig 3B**). In contrast, we detected
266    none of the Par1-derived junctions in the library generated from T7-transcribed HA (**Fig
267    3B**). Although the read depth coverage was comparable to Par1, 615_1132_HA was
268    completely absent, and the three junctions that were detected had minimal read support
269    and were not seen in virus-derived libraries. Altogether, these results suggest that the
270    formation of deletion junctions during the reverse transcription reaction is rare, and that
271    the Par1-derived junctions we observed are most likely derived from true DIPs present
272    within our viral stock, despite the stock having been prepared at low MOI. This highlights
273    the difficulty in producing a completely DIP-free virus preparation.
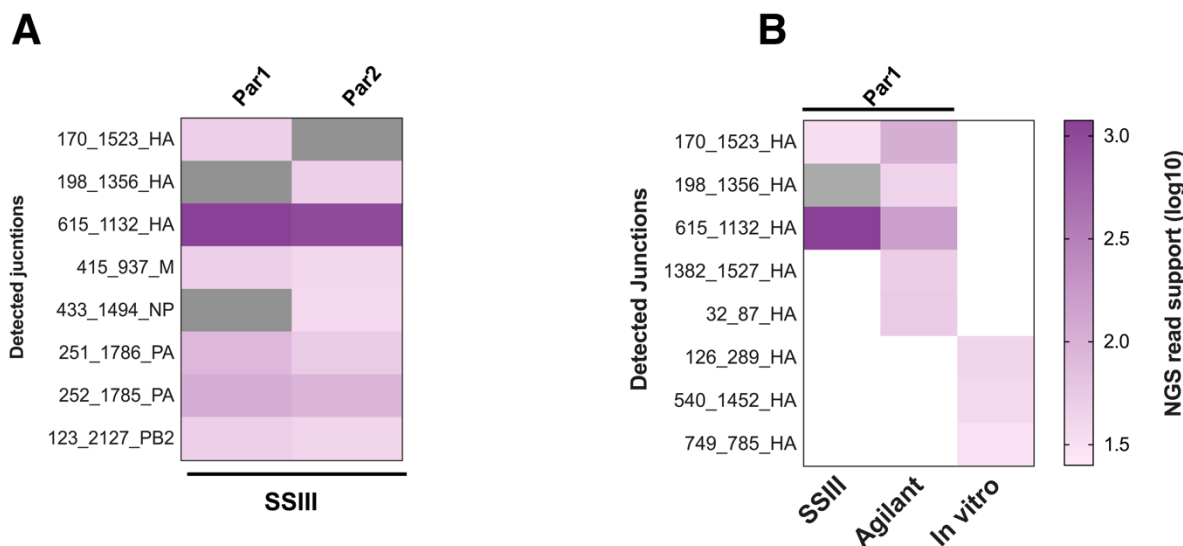274



**Fig 3. DIP-associated deletion junctions present in virus working stock.** We performed two independent RNA extractions, RT reactions, PCR amplifications, and library preparations from a single recombinant Cal07 working stock grown at low MOI (Par1 and Par2). **(A)** Comparison of deletion junctions detected in Par1 and Par2 samples. Purple blocks represent instances where junction was detected but read support was below RSC **(B)** Comparison of HA segment junctions detected in two libraries generated from independent RT reactions using two different RT enzymes, along with a library generated from in vitro T7-transcribed viral RNA. White blocks denote no detection

275
276

### Generation of DIP-enriched populations through high MOI passage

278    To test the ability of the pipeline to detect real DIP-associated RNAs, we enriched for
279    DIPs through serial undiluted passage of Cal07 in MDCK cells. We confirmed the
280    presence of DIPs by amplifying full-length genomic cDNA at each passage and examining
281    the size distribution of PCR products by gel electrophoresis **(Fig 4A)**, as previously
282    described (21). The gradual disappearance of the polymerase segments, which form the
283    majority of DIPs, and the appearance of a smear below the shortest IAV segment (NS
284    ~0.9kb) were consistent with the accumulation of DIPs over passage. Based on these
285    results, we picked P1, P3, and P6 as representative samples for sequencing.

286
287 We further confirmed the presence of DIPs by plotting the read coverage of the aligned
288 reads from passages 1, 3, and 6 (**Fig 4B**). These coverage plots clearly reveal the
289 characteristic pattern of DIP-rich populations, with much lower depth of read alignment in
290 the middle portion of the segment compared with the termini. As expected, the number of
291 DIP-associated deletion junctions detected by the pipeline also increased across
292 passages, reaching the highest level at passage 6 (**Fig 4C**). To confirm that these
293 junction-containing sequences were derived from virion rather than cellular RNA, we
294 measured the number of reads that aligned to the host (canine) genome in our samples.
295 We found very few reads derived from the canine genome in all the passages, compared
296 with about 40% of the reads from RNA extracted from infected cells (**Fig S4**).
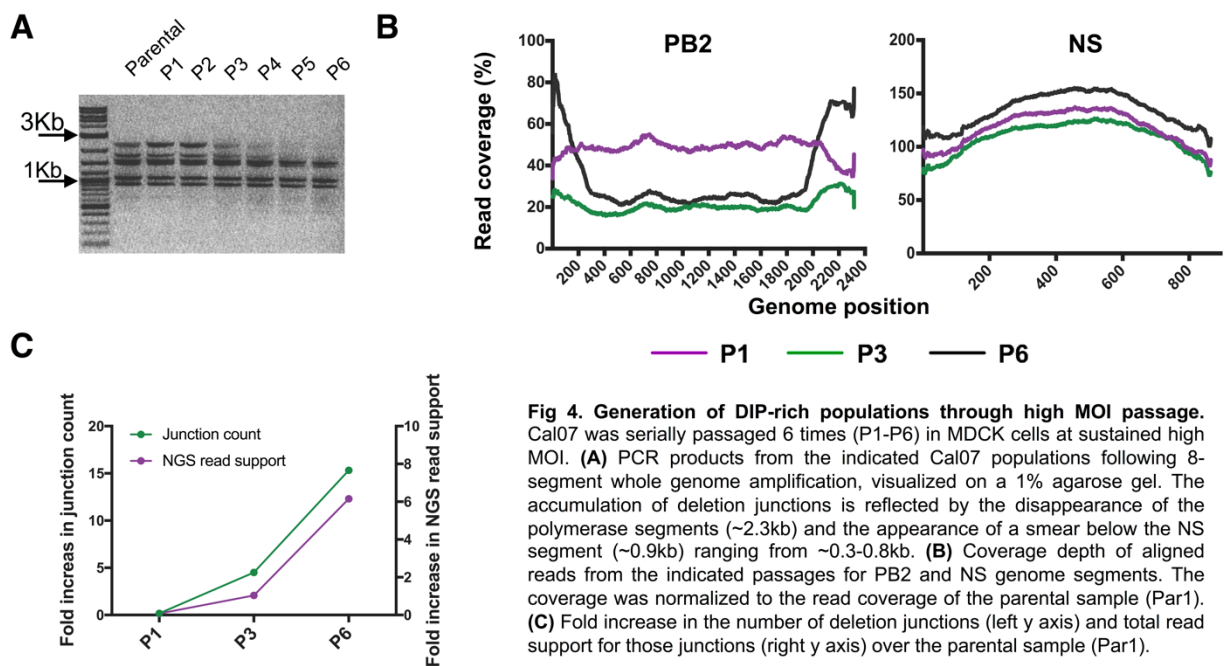297



**Fig 4. Generation of DIP-rich populations through high MOI passage.**
Cal07 was serially passaged 6 times (P1-P6) in MDCK cells at sustained high MOI. **(A)** PCR products from the indicated Cal07 populations following 8-segment whole genome amplification, visualized on a 1% agarose gel. The accumulation of deletion junctions is reflected by the disappearance of the polymerase segments (~2.3kb) and the appearance of a smear below the NS segment (~0.9kb) ranging from ~0.3-0.8kb. **(B)** Coverage depth of aligned reads from the indicated passages for PB2 and NS genome segments. The coverage was normalized to the read coverage of the parental sample (Par1). **(C)** Fold increase in the number of deletion junctions (left y axis) and total read support for those junctions (right y axis) over the parental sample (Par1).

298
299
300 **Reproducibility of pipeline performance**
301 Multiple steps in the combined experimental/computational pipeline could introduce
302 stochasticity into the pipeline performance, thus diminishing overall consistency and
303 reproducibility of output. To examine the reproducibility of our pipeline's performance, we
304 sequenced two separate extractions of a single P6 population (Hereafter known as L1-
305 P6-Rep1 and L1-P6-Rep2, where L refers to lineage) and compared the pipeline outputs
306 between the two replicates (**Fig S5**). We found that the normalized read support values
307 of individual junctions were highly correlated between the two replicate samples, whether
308 the replicates were sequenced on the same MiSeq flowcell (Spearman R = 0.92) or
309 separate ones (Spearman R = 0.91). Thus, the combined steps from RNA extraction to
310 sequence analysis introduce minimal noise into the pipeline output, and pipeline
311 performance is highly reproducible between experiments.
312
313 **Optimization of minimum read support cutoffs**

314   Our experiments using simulated datasets revealed the importance of setting minimal
315   RSCs for maximizing the accuracy of pipeline performance, and suggested that the
316   optimal RSC may differ between datasets. We next attempted to optimize RSC values for
317   our experimental dataset where we did not actually know the precise location and number
318   of junctions present in the population (as we did with our simulated datasets). To quantify
319   precision in junction detection for our experimental dataset, we assumed that base calling
320   errors and mutations that result in inaccurate junction reporting would be stochastic and
321   thus read support for these inaccurate junctions would be highly variable between
322   technical replicates. In contrast, read support for real junctions should be consistent
323   between replicates.
324
325   We assessed the effects of varying the RSC on the degree of correlation between
326   junctions identified in L1-P6-Rep1 and L1-P6-Rep2. We varied the RSC values from 1 to
327   50 for each individual genome segment, and examined the effect on the number of
328   reported junctions (**Fig 5**). We observed a similar pattern to that observed for our
329   simulated data, where raising the RSC to 10 or higher resulted in a large drop-off in the
330   number of reported junctions. We next determined the RSC value that yielded the highest
331   degree of correlation between the two replicates. We identified distinct optimal RSC cutoff
332   values for each segment: 20, 20, 30, 30, and 15 for PB2, PB1, PA, HA, and NA,
333   respectively. The average of these values was used as an RSC for the remaining
334   segments where no enough junctions were detected to perform the correlation test (see
335   below).
336
337   We do not expect these values to be universal, as they likely are influenced by a number
338   of factors that will vary between individual sequencing runs. Also, for different
339   applications, it may be beneficial to lower the RSC to improve detection sensitivity at the
340   cost of precision. Thus, we suggest running two technical replicates with each NGS run
341   to be used as reference to establish optimal per-segment RSC values for that run.
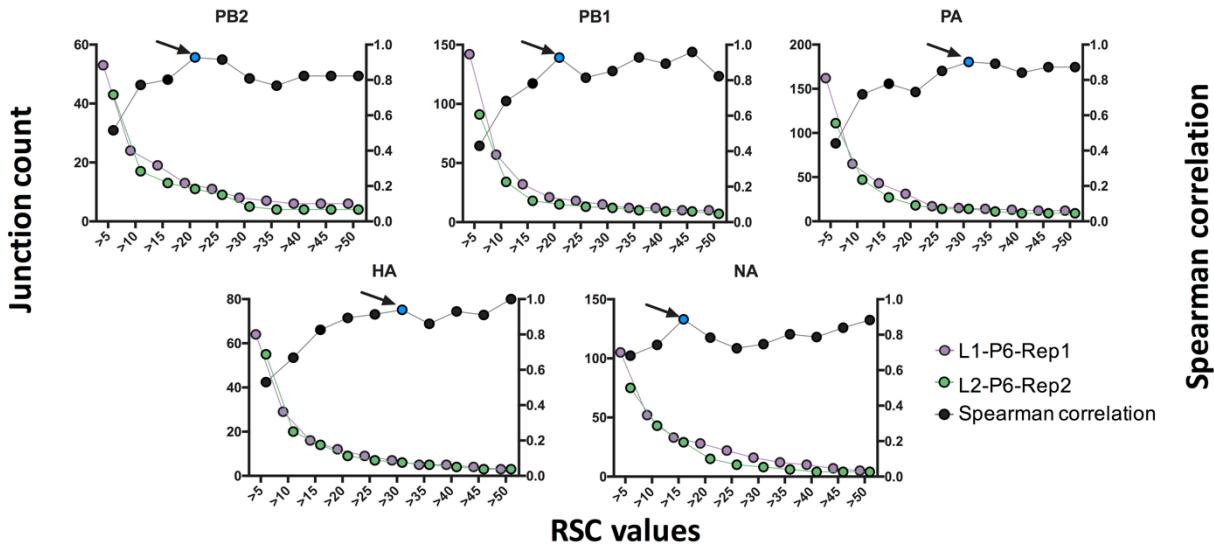342

**Fig 5. Determination of optimal read support cutoffs for experimental data.** Plots showing the numbers of deletion junctions reported in the indicated genome segments for two technical library preparation and sequencing replicates generated from a single DIP-rich viral population (L1-P6-Rep1 and L2-P6-Rep2; left y axis). Black dots represent the results of Spearman correlation tests between the replicates at each RSC condition (right y axis). Blue dots indicate the point with the highest degree of correlation and minimum decrease of junction count for each genome segment.

343
344
345    **Analysis of DIP-associated deletion profiles**
346    We next examined the overall diversity of DIP-associated deletion junctions within the P6
347    populations from the two independent lineages (L1-P6-Rep1 and L2-P6), and found
348    dozens of distinct deletion junctions scattered across the viral genome in both lineages
349    (**Fig 6A**). Junctions were not evenly distributed across the genome segments, as few to
350    no junctions were detected in the NP, M, or NS segments. Within each segment, the read
351    support for individual junctions varied significantly (**Fig 6B**). When we compared the
352    deletion junction repertoires between the two passage lineages, we observed that a
353    significant fraction of the detected junctions was shared between the two, and that these
354    shared junctions exhibited a high degree of correlation in terms of read support (**Fig 6C**).
355    These data suggest that specific DIP-associated deletions may be consistently formed by
356    Cal07. While there was substantial diversity in terms of the number of distinct deletion
357    junctions present, when we plotted the locations of those these junctions within the
358    genome segments, we observed that they were largely confined within clear hotspots
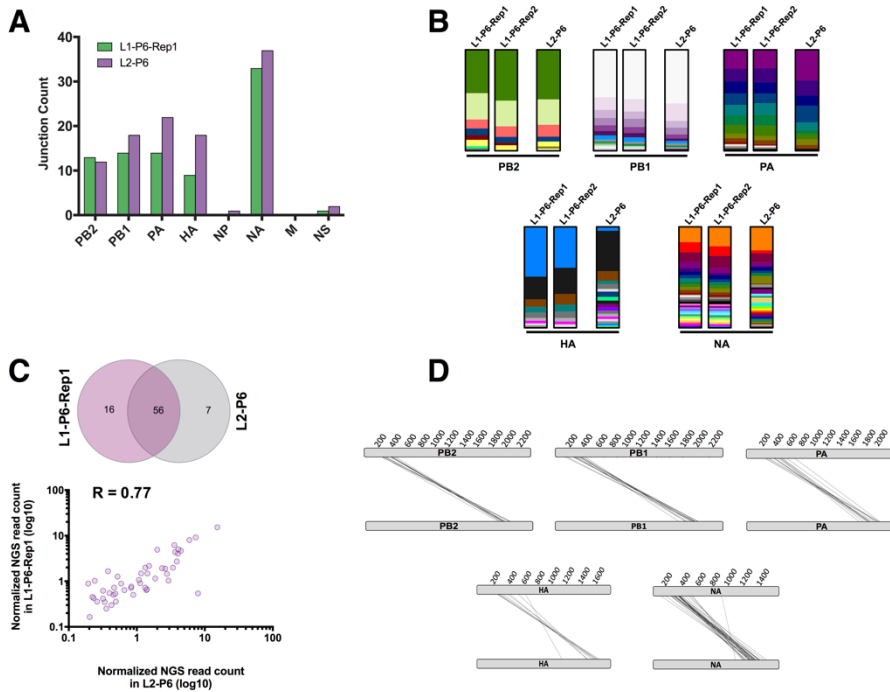359    towards the termini of the segments with few exceptions (**Fig 6D**).
360

**Fig 6. DIP-associated junctions analysis. (A)** The total number of detected junctions in the individual genome segments for the two independent passage 6 populations. **(B)** Stacked column charts showing the proportional abundance of each deletion junction per segment between lineage 1 (including both technical replicates) and lineage 2 at passage 6. Each color bar represents a unique junction within each segment, whose height reflects the relative NGS read support, normalized to the total number of NGS junction-spanning reads for the indicated segment. **(C)** Comparison between L1-P6-REP1 and L2-P6 in relation to the number of shared junctions (upper panel) and the correlation between their NGS reads (lower panel) **(D)** Parallel coordinates diagrams showing the deletion junctions in P6-Rep1 mapped to their actual respective locations on each segment. Each individual junction is represented by a black line that connects the donor and acceptor sites of the breakpoint.

### Effect of varying template input on pipeline performance

We next asked whether the amount of cDNA template that goes into the library preparation affects the sensitivity and stochasticity of junction detection by the pipeline. We serially diluted both the amount of viral RNA template used in the RT reaction and the amount of cDNA template used in the PCR and compared pipeline outputs from the DIP-rich L1-P6-Rep1 population. We first tested the correlation of detected DIP-associated junctions between a limited number of dilutions ranging from 1:3 to 1:15. We observed that the correlation of read support values between specific junctions across dilutions was more consistent when cDNA was diluted, rather than RNA, suggesting that RNA dilution may increase the stochasticity of downstream PCR amplification (**Fig S6A**).

Based on this, we performed whole genome PCR using a dilution series of L1-P6-Rep1-derived cDNA (spanning roughly $4*10^8$ to $4*10^6$ NP genome equivalents per PCR) as template (**Fig 7A**). We observed that there is an optimal amount of input cDNA template for maximizing junction detection. Diluting the input cDNA 1:120 (corresponding to ~$4*10^6$ NP genome equivalents) increased the number of detected junctions over 4-fold compared with undiluted input. Although the number of DIP-associated junctions was increased, the distribution of junctions across segments and their mapped locations were consistent with our earlier results (**Fig S6B and Fig 6**).

Further dilution of input template beyond 1:120 resulted in a decrease in sensitivity. Importantly, dilution across the range tested did not result in a failure to detect any of the junctions reported in the undiluted sample. We also observed that the correlation of read support values between specific junctions across dilutions tracked closely with the sensitivity (**Fig 7B**). Altogether, these observations indicate that optimization of the cDNA

388    template input amount can significantly improve the sensitivity of DIP-associated junction
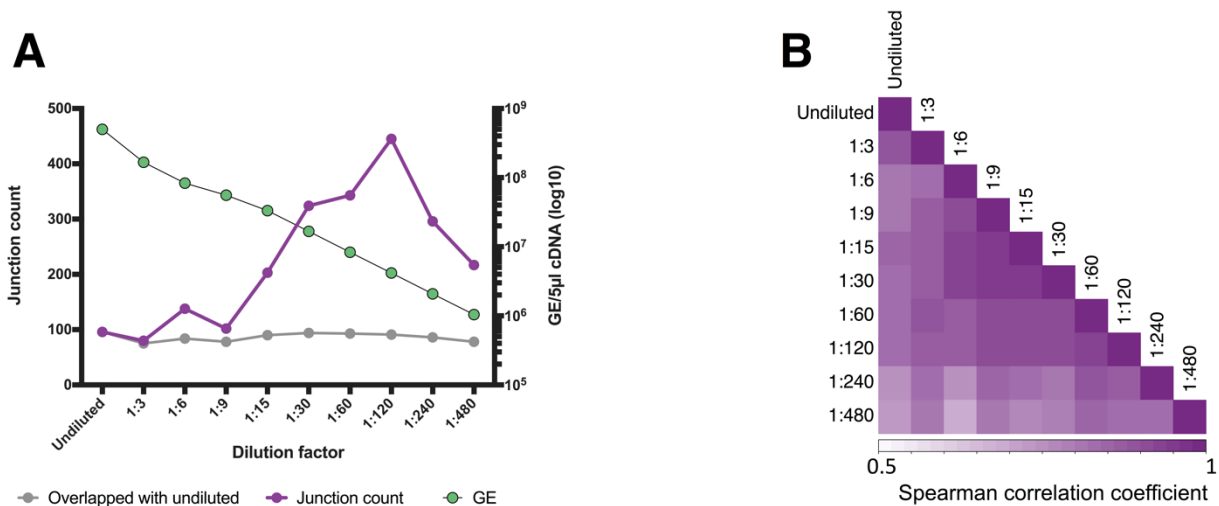389    detection.
390



**Fig 7. Effects of viral template input on the detection of DIP-associated junctions.** We serially diluted cDNA generated from the L1-P6-Rep1 sample, and compared sequencing results between libraries generated with these dilutions as templates. **(A)** For each dilution, the total numbers of detected junctions (purple) are shown, along with the number of specific junctions detected that were also detected in the undiluted sample (grey). The copy number of viral cDNA molecules included in downstream PCR and library preparation for each dilution was determined by RT-qPCR (green; right y axis). **(B)** Read support values for all deletion junctions common across the diluted and undiluted samples were normalized to the total number of deletion junction-spanning reads for each sample and used to perform a Spearman correlation between all pairs of samples using R cor function.

391
392

393    **Lack of association between direct repeats and junction formation**
394    Direct repeat sequences (detailed in **Fig S2 and S3A**) are common across the IAV
395    genome and have previously been hypothesized to contribute to DIP-associated deletion
396    formation by promoting viral polymerase slippage (10, 15). We leveraged the large
397    number of DIP-associated deletion junctions that we identified in this study to test this
398    hypothesis. We asked whether the deletion junctions in the DIP-enriched sample L1-P6-
399    Rep1 were found more frequently adjacent to direct repeats than would be expected if
400    the junctions were located randomly in the viral genome. We compared the frequency of
401    deletion junctions associated with direct repeats between the L1-P6-Rep1 and L2-P6
402    populations (where all deletions are formed by the viral RdRp) and the Cal07-200
403    simulated dataset, where all deletions are randomly localized (**Fig 8**). The frequency of
404    direct repeats of varying lengths at junction sites in the real viral populations was not
405    significantly different than that seen in the simulated data, indicating that direct repeat
406    sequences are not enriched at DIP-associated junctions and arguing against a significant
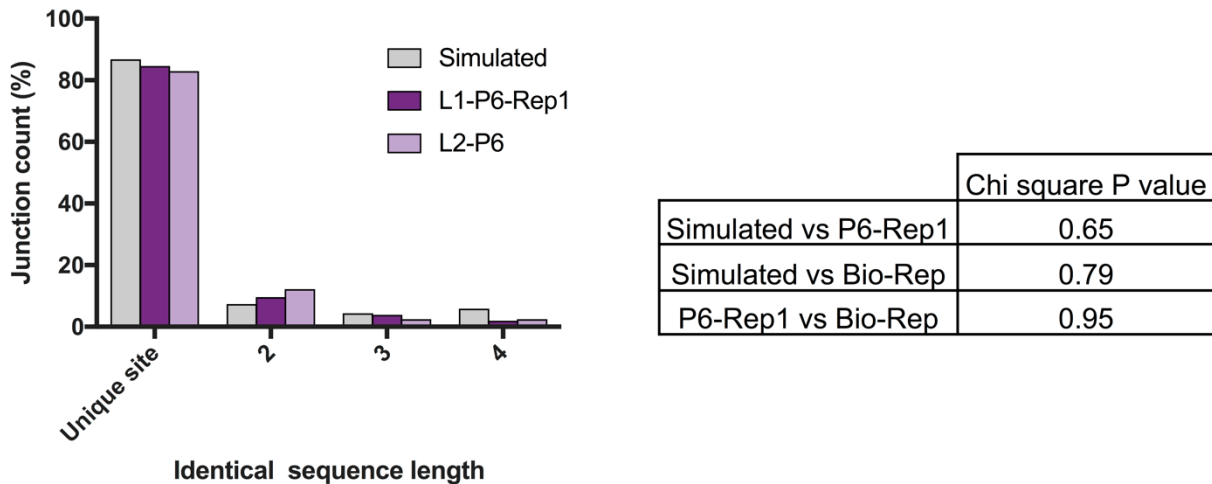407    role for direct repeats in DIP formation.
408

**Fig 8. Direct repeat sequences are not over-represented at DIP-associated deletion junctions.** The percentages of deletion junctions within the polymerase segments that occurred at unique sites or at sites with direct nucleotide repeats with lengths 2-4nts was compared between L1-P6-Rep1, L2-P6, and the Cal07-200 simulated dataset. The number of junctions was plotted and compared by Chi Square. The table shows the Chi Square P values between every possible pair of the samples.

409

410

411 **DISCUSSION**

412 Sensitive and accurate detection of DIP-associated sequences within viral populations is
413 critical for defining how DIPs form and function during IAV infection. Here we outline a
414 pipeline to detect DIP-associated junctions within viral populations using Illumina-based
415 short read sequencing, and validate its performance using a combination of simulated
416 and experimental control datasets.

417

418 Our primary goal was to develop and optimize a reasonably simple and straightforward
419 sequencing framework that accounts for the potential artifacts that can potentially
420 confound NGS-based DIP detection efforts. We chose the Illumina sequencing platform
421 because it is widely available, easy to use, conducive to sample multiplexing, and
422 because it has a relatively low rate of base-calling errors. One concern that we had initially
423 was that recombination during reverse transcription, PCR, or sequencing might make
424 identification of bona fide DIP-associated junctions a challenge. Two recently developed
425 technologies, CirSeq and ClickSeq, largely eliminate this issue, but also significantly
426 increase the amount of labor involved in library preparation (22, 23). We observed that
427 the occurrence of non-viral recombination that occurs during our library preparation and
428 sequencing procedures was vanishingly small, and can effectively be ignored. Thus, while
429 both Cirseq and ClickSeq are enormously useful in certain circumstances, our data

430  indicates that such methods are not required to generate highly accurate and sensitive
431  profiles of IAV DIPs.
432
433  A significant shortcoming of the method we detail here is that the measured read support
434  for individual deletion junctions does not necessarily reflect the actual frequencies of
435  these deletions within the viral population. This is due to both biasing of PCR amplification
436  towards shorter products, as well as the uneven distribution of read coverage across the
437  viral genome. For situations where the accurate measurement of individual DIP genotype
438  frequencies is critical, we recommend pairing a cDNA barcoding method such as primer
439  ID (24, 25) with a platform capable of long-read sequencing, such as PacBio or Oxford
440  Nanopore (26). Alternatively, direct sequencing of viral RNA using the Oxford Nanopore
441  platform may also prove to be useful for accurate measurement of junction frequencies
442  (27).
443
444  When we used our pipeline to examine DIP-enriched viral populations generated through
445  serial high-MOI passage, we detected dozens of distinct DIP-associated deletion
446  junctions, revealing a high degree of diversity within the DIP population. Although the
447  majority of these DIP-associated junctions were derived from the polymerase segments
448  as expected, we also detected a substantial proportion of deletions within the HA and NA
449  segments, but not the NP, M, and NS segments. The non-random distribution of junctions
450  across the genome segments mirrors what has been reported elsewhere, and highlights
451  how little we know about the specific molecular mechanisms that regulate DIP formation.
452
453  We hope that the approach detailed here, and the associated bioinformatics pipeline
454  prove useful to other groups interested in defective interfering particle biology. Our
455  approach is optimized for influenza virus sequences; however, the approaches and
456  controls detailed here can easily be adapted to other RNA virus systems.
457
458  **MATERIALS AND METHODS**
459  **Viruses and Cells**
460  Madin-Darby canine kidney cells (MDCK; obtained from Dr. Jonathan Yewdell) and
461  human embroyonic kidney 293 cells (293T; obtained from Dr. Joanna Shisler) were grown
462  in Minimal Essential Medium (MEM) + GlutaMAX (Gibco), supplemented with 8.3% fetal
463  bovine serum (Seradigm), at 37°C and 5% $CO_2$. Recombinant A/California/07/09 (Cal07)
464  virus was rescued via the standard 8-plasmid reverse genetics approach. Briefly, 60-90%
465  confluent 293T cells were transfected with 500ng of the following plasmids (pDZ::PB2,
466  pDZ::PB1, pDZ::PA, pDZ::HA, pDZ::NP, pDZ::NA, pDZ::M, pDZ::NS) using JetPRIME
467  (Polyplus) according to the manufacturer's instructions. Cal07 reverse genetics plasmids
468  were originally obtained as A/California/04/2009-encoding plasmids from Dr. Jonathan
469  Yewdell. We introduced A660G and A335G substitutions into the HA and NP plasmids,
470  respectively, to convert them to match the amino acid sequence of A/California/07/2009
471  HA and NP (NCBI accession# CY121680, CY121683). A seed stock was prepared by
472  amplifying a plaque isolate from the rescue supernatants. Virus working stocks were
473  generated by infecting MDCK cells with seed stock at an MOI of 0.001 TCID50/cell and

474    collecting and clarifying supernatants at 48 hpi.

**Generation of DIP stocks through high MOI passage**

476    Confluent MDCK cells in 96-well plates were infected with IAV Cal07 at an MOI of 5
477    TCID50/cell. Supernatants (200µl total per well) were harvested at 24 hpi (passage 1)
478    and pooled. 100µl/well of this pooled supernatant was used to infect a 96 well plate of
479    fresh MDCK cells to generate the next passage. This process was repeated 6 times to
480    produce passages 1-6 in two independent lineages (1 96-well plate per lineage).

**IAV Genome amplification**

482    Viral RNA was extracted from 140µl of cell culture supernatant using the QIAamp viral
483    RNA kit (Qiagen) and eluted in 60µl distilled H2O (dH2O). For cDNA reactions, 3µL of
484    RNA     was     mixed     with     1µL     (2µM)     MBTUni-12     primer     (5'-
485    ACGCGTGATCAGCRAAAGCAGG-3') + 1µL (10µM) dNTPs + 8µL dH2O. The mixture
486    was incubated for 5 minutes at 65°C and then placed on ice for 2 min. Subsequently, the
487    mixture was removed from ice and the following was added: 1µL SuperScript III RT
488    (Invitrogen), 4µL of 5X First-Strand Buffer (Comes with SSIII kit), 1µL of DTT, 1µL RNase-
489    in (Invitrogen). The reaction was incubated at 45°C for 50 min, followed by a 15 min
490    incubation at 70°C for inactivation. 5µL of cDNA product was mixed with the following for
491    PCR     amplification:          2.5µL     (10µM)     MBTUni-12_4R     primer     (5'-
492    ACGCGTGATCAGCRAAAGCAGG-3'),     2.5µL     (10µM)     MBTUni-13     primer     (5'-
493    ACGCGTGATCAGTAGAAACAAGG-3'), 0.5µL Phusion polymerase (NEB), 10µL - 5x HF
494    buffer, 1µL (10mM dNTPs mix), and 28.5µL dH2O. The PCR reaction conditions used:
495    98°C (30 s) followed by 25 cycles of 98°C (10 s), 57°C (30 s) and 72°C (1:30 min), a
496    terminal extension of 72°C (5 min), and a final 10°C hold. PCR products were purified
497    using the PureLink PCR purification kit (Invitrogen) with the <300nt cutoff option and
498    eluted in 30µL dH2O. There was no difference in deletion junction detection when we
499    purified the PCR products with the lower cutoff option (data not shown).
500
**NGS library preparation**

502    We started with ~20ng of the PCR products in a volume of 50µl. The Covaris M220
503    sonicator (Covaris) was used to fragment the DNA. Three different conditions were used
504    to generate different average fragment lengths of 300, 500, 700 base pairs (bp): (I) 300
505    bp = Peak Power 50, Duty Factor 20 and Cycles/Burst 200 for 2:40 min, (II) 500 bp =
506    Peak Power 50, Duty Factor 10 and Cycles/Burst 200 for 1:30 min, and (III) ~600 bp =
507    fragment length Peak Power 50, Duty Factor 10 and Cycles/Burst 200 for 1 min. In our
508    hands, the fragmentation length did not have any effect on our sequencing results (data
509    not shown). For the sake of consistency, we used the 300 bp fragmentation length. To
510    confirm the PCR products, we visualized the amplicons on a Fragment Analyzer (AATI)
511    with the DNF-486 high sensitivity NGS kit before and after fragmentation. Next, we used
512    KAPA Hyper Prep kit (Roche) to construct the libraries according to the manual. To
513    eliminate the possibility of index hopping (or index switching), we used the TruSeq Unique
514    Dual Indexes (UDI) from Illumina. The Adapter ligation step was carried out with 5µl of
515    Truseq UDIs diluted 1:10 with 10nM Tris. For maximum efficiency we increased the
516    ligation time to 30mins. We then performed 3 cycles of PCR with the Kapa library

517 amplification primers diluted 1:5 in water followed by a cleanup step with 40$\mu$l of AxyPrep
518 Mag PCR beads (Thermofisher). We then mixed the libraries at an equimolar ratio and
519 carried out a qPCR to accurately quantitate the library pool and maximize the number of
520 clusters in the sequencing flowcell. A size selection step was not needed. Finally, the
521 pooled libraries were sequenced with paired-ends 2x250nt reads on an Illumina MiSeq
522 using V2 chemistry. The fastq files were generated and demultiplexed with the bcl2fastq
523 v2.20 Conversion Software (Illumina).
524

525 **Simulated Datasets**
526 All the simulated datasets used in this study were generated by MetaSim (v0.9.1) (28),
527 a genomic and metagenomics simulator. Several reference library sequences composed
528 of WT reference sequences of IAV Cal07 or PR8 (see Table S1 for NCBI accession
529 numbers), mixed with a defined DIP sequence population - generated randomly within
530 the first and last 600 nts of all the segments - were used in Metasim for data simulation.
531 The configurations were fixed across all datasets to maintain the preferable conditions.
532 The reference sequences were fragmented into 350 nts fragments length with a standard
533 deviation of +/- 50 and were simulated into ~1 million 2x250nts paired-end reads per
534 sample, with a total mutation rate of ~1% based on the published Illumina empirical error
535 model, and corresponds to substitutions as the indel error rate is negligible within Illumina
536 MiSeq. One dataset was simulated with no DIP sequences as a control sample for any
537 computational artifacts. Metasim generated two FASTA files of 1 million reads per file per
538 sample (~2 million single-end reads = 1 million paired-end reads), which subsequently
539 were used for the optimization process.

540 **Sequencing analysis of DIP-associated junctions**
541 The raw sequencing reads were quality-filtered by Trimmomatic (v0.36) (Parameters:
542 `ILLUMINACLIP:TruSeq3-PE-2.fa:2:15:10`            `SLIDINGWINDOW:3:20`
543 `LEADING:28 TRAILING:28`) (29) and any reads shorter than 75nts were removed from
544 the datasets. The paired reads were concatenated into one file and treated as single-end
545 when aligned end-to-end to the WT reference sequences using Bowtie2 (v2.3.1)
546 (Parameters: --score-min `L,0,-0.3`). Subsequently, the algorithm ViReMa (v0.10) was
547 used to analyze the remaining un-aligned reads (putative junction-spanning reads)
548 (Parameters `-DeDup --MicroInDel_Length 20 --Defuzz 3 --N 1 --X 8`).
549 Next, the DIP-associated deletion junctions and their read support were extracted from
550 ViReMa output files and sorted per segment, using an in-house Perl scrip, for data
551 analysis and visualization. To detect any MDCK genome leakage, the datasets were
552 aligned against the dog genome (assembly CanFam3.1). All scripts are available at
553 *https://github.com/BROOKELAB/Influenza-virus-DI-identification-pipeline*.
554

555 **Quantification of sensitivity and precision**
556 To calculate the actual number of junction-spanning reads in Fig 2A, reads that derived
557 from DIP-associated sequences were counted by their FASTA headers, which contain
558 the source of each read, produced by MetaSim. To calculate the maximum theoretical
559 sensitivity of ViReMa (Fig 2B) based on seed length of 25nts and two allowed mutations
560 (--N = 2), the number of mutations was subtracted from the seed length, which on its turn

561     was multiplied by 2 to account for both termini ((25-2)*2=46). Subsequently, this number
562     was subtracted from the possible cutting site of a 250nts read and divided by the total
563     number of cutting sites and multiplied by 100 ((249-46)/249*100)=81.5%). To calculate
564     the number of accurately and inaccurately mapped junctions in Fig 2C,D, the seed
565     sequences of the Cal07-200 dataset were used against ViReMa with --N set to 0, and the
566     remaining parameters were kept the same. These sequences were generated initially to
567     establish the seed for MetaSim to simulate sequencing, therefore their lengths are varied
568     between ~350-1800. The long ones were trimmed to <1000nts, so ViReMa would take
569     them as reads (the maximum default read length that ViReMa could take is ~1024nts)
570     and, critically, the junction locations were maintained. The junctions that occurred within
571     the first or last 25nts were removed (4 junction sequences). Finally, the junctions that
572     accurately mapped were counted, which found to be 149 versus 47 inaccurately mapped
573     junctions.
574

575     **Correlation analysis**
576     For the correlation tests, the NGS read support count for each DIP-associated junction
577     was normalized to the total detected junction-spanning reads of every sample. Next, the
578     correlation was calculated based on Spearman rank correlation using either R (cor
579     function) or an online tool at:
580     http://www.biostathandbook.com/spearman.html
581

582     **ACKNOWLEDGMENTS**

589

590     **REFERENCES**

591     1.   Von Magnus P. 1954. Incomplete forms of influenza virus. Adv Virus Res 2:59–79.

592     2.   von MAGNUS P. 1951. Propagation of the PR8 strain of influenza A virus in chick
593         embryos. II. The formation of incomplete virus following inoculation of large doses of
594         seed virus. Acta Pathol Microbiol Scand 28:278–293.

595     3.   Rezelj VV, Levi LI, Vignuzzi M. 2018. The defective component of viral populations.
596         Curr Opin Virol 33:74–80.

597     4.   Baum A, Sachidanandam R, García-Sastre A. 2010. Preference of RIG-I for short
598         viral RNA molecules in infected cells revealed by next-generation sequencing. Proc
599         Natl Acad Sci U S A 107:16303–16308.

5.  Nayak DP, Chambers TM, Akkina RK. 1985. Defective-interfering (DI) RNAs of influenza viruses: origin, structure, expression, and interference. Curr Top Microbiol Immunol 114:103–151.

6.  Brooke CB. 2017. Population Diversity and Collective Interactions during Influenza Virus Infection. J Virol 91.

7.  Vasilijevic J, Zamarreño N, Oliveros JC, Rodriguez-Frandsen A, Gómez G, Rodriguez G, Pérez-Ruiz M, Rey S, Barba I, Pozo F, Casas I, Nieto A, Falcón A. 2017. Reduced accumulation of defective viral genomes contributes to severe outcome in influenza virus infected patients. PLOS Pathog 13:e1006650.

8.  Sherry L, Punovuori K, Wallace LE, Prangley E, DeFries S, Jackson D. 2016. Identification of cis-acting packaging signals in the coding regions of the influenza B virus HA gene segment. J Gen Virol 97:306–315.

9.  Hutchinson EC, von Kirchbach JC, Gog JR, Digard P. 2010. Genome packaging in influenza A virus. J Gen Virol 91:313–328.

10. Jennings PA, Finch JT, Winter G, Robertson JS. 1983. Does the higher order structure of the influenza virus ribonucleoprotein guide sequence rearrangements in influenza viral RNA? Cell 34:619–627.

11. Routh A, Johnson JE. 2014. Discovery of functional genomic motifs in viruses with ViReMa-a Virus Recombination Mapper-for analysis of next-generation sequencing data. Nucleic Acids Res 42:e11.

12. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. Nat Methods 9:357–359.

13. Jennings PA, Finch JT, Winter G, Robertson JS. 1983. Does the higher order structure of the influenza virus ribonucleoprotein guide sequence rearrangements in influenza viral RNA? Cell 34:619–627.

14. Janda JM, Davis AR, Nayak DP, De BK. 1979. Diversity and generation of defective interfering influenza virus particles. Virology 95:48–58.

15. Saira K, Lin X, DePasse JV, Halpin R, Twaddle A, Stockwell T, Angus B, Cozzi-Lepri A, Delfino M, Dugan V, Dwyer DE, Freiberg M, Horban A, Losso M, Lynfield R, Wentworth DN, Holmes EC, Davey R, Wentworth DE, Ghedin E, INSIGHT FLU002 Study Group, INSIGHT FLU003 Study Group. 2013. Sequence analysis of in vivo defective interfering-like RNA of influenza A H1N1 pandemic virus. J Virol 87:8064–8074.

633    16.    Loman NJ, Misra RV, Dallman TJ, Constantinidou C, Gharbia SE, Wain J, Pallen
634           MJ. 2012. Performance comparison of benchtop high-throughput sequencing
635           platforms. Nat Biotechnol 30:434–439.

636    17.    Schirmer M, Ijaz UZ, D'Amore R, Hall N, Sloan WT, Quince C. 2015. Insight into
637           biases and sequencing errors for amplicon sequencing with the Illumina MiSeq
638           platform. Nucleic Acids Res 43:e37.

639    18.    Zhou B, Donnelly ME, Scholes DT, St George K, Hatta M, Kawaoka Y,
640           Wentworth DE. 2009. Single-reaction genomic amplification accelerates sequencing
641           and vaccine production for classical and Swine origin human influenza a viruses. J
642           Virol 83:10309–10313.

643    19.    Görzer I, Guelly C, Trajanoski S, Puchhammer-Stöckl E. 2010. The impact of
644           PCR-generated recombination on diversity estimation of mixed viral populations by
645           deep sequencing. J Virol Methods 169:248–252.

646    20.    Lahr DJG, Katz LA. 2009. Reducing the impact of PCR-mediated recombination
647           in molecular evolution and environmental studies using a new-generation high-
648           fidelity DNA polymerase. BioTechniques 47:857–866.

649    21.    Xue J, Chambers BS, Hensley SE, López CB. 2016. Propagation and
650           Characterization of Influenza Virus Stocks That Lack High Levels of Defective Viral
651           Genomes and Hemagglutinin Mutations. Front Microbiol 7:326.

652    22.    Routh A, Head SR, Ordoukhanian P, Johnson JE. 2015. ClickSeq:
653           Fragmentation-Free Next-Generation Sequencing via Click Ligation of Adaptors to
654           Stochastically Terminated 3'-Azido cDNAs. J Mol Biol 427:2610–2616.

655    23.    Acevedo A, Andino R. 2014. Library preparation for highly accurate population
656           sequencing of RNA viruses. Nat Protoc 9:1760–1769.

657    24.    Kosik I, Ince WL, Gentles LE, Oler AJ, Kosikova M, Angel M, Magadán JG, Xie
658           H, Brooke CB, Yewdell JW. 2018. Influenza A virus hemagglutinin glycosylation
659           compensates for antibody escape fitness costs. PLOS Pathog 14:e1006796.

660    25.    Jabara CB, Jones CD, Roach J, Anderson JA, Swanstrom R. 2011. Accurate
661           sampling and deep sequencing of the HIV-1 protease gene using a Primer ID. Proc
662           Natl Acad Sci U S A 108:20166–20171.

663    26.    Jaworski E, Routh A. 2017. Parallel ClickSeq and Nanopore sequencing
664           elucidates the rapid evolution of defective-interfering RNAs in Flock House virus.
665           PLOS Pathog 13:e1006365.

666    27.    Garalde DR, Snell EA, Jachimowicz D, Sipos B, Lloyd JH, Bruce M, Pantic N,
667           Admassu T, James P, Warland A, Jordan M, Ciccone J, Serra S, Keenan J, Martin

668   S, McNeill L, Wallace EJ, Jayasinghe L, Wright C, Blasco J, Young S, Brocklebank
669   D, Juul S, Clarke J, Heron AJ, Turner DJ. 2018. Highly parallel direct RNA
670   sequencing on an array of nanopores. Nat Methods 15:201–206.

671   28.   Richter DC, Ott F, Auch AF, Schmid R, Huson DH. 2008. MetaSim—A
672   Sequencing Simulator for Genomics and Metagenomics. PLOS ONE 3:e3373.

673   29.   Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for
674   Illumina sequence data. Bioinforma Oxf Engl 30:2114–2120.

675   30.   Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis
676   G, Durbin R, 1000 Genome Project Data Processing Subgroup. 2009. The
677   Sequence Alignment/Map format and SAMtools. Bioinforma Oxf Engl 25:2078–
678   2079.

679   31.   Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing
680   genomic features. Bioinforma Oxf Engl 26:841–842.

681