

1 **DISCOVERY OF UFO PROTEINS: HUMAN-VIRUS CHIMERIC PROTEINS GENERATED**
2 **DURING INFLUENZA VIRUS INFECTION.**

3
4
5 Yixuan Ma^{1,11}, Matthew Angel^{2,11}, Guojun Wang^{9,10,11}, Jessica Sook Yui Ho^{1,11}, Nan Zhao¹, Justine
6 Noel¹, Natasha Moshkina¹, James Gibbs², Jiajie Wei², Brad Rosenberg¹, Jeffrey Johnson³, Max Chang⁴,
7 Zuleyma Peralta⁵, Nevan Krogan³, Christopher Benner⁴, Harm van Bakel⁵, Marta Łuksza³, Benjamin D.
8 Greenbaum⁶, Emily R. Miraldi⁷, Adolfo Garcia-Sastre⁸, Jonathan W. Yewdell^{2,12} and Ivan Marazzi^{9,12*}
9

10 ¹Department of Microbiology, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA

11 ²Laboratory of Viral Diseases, National Institute of Allergy and Infectious Diseases, NIH, Bethesda, MD
12 20892, USA

13 ³Department of Cellular and Molecular Pharmacology, University of California, San Francisco, San
14 Francisco, CA 94158, USA

15 ⁴Department of Medicine, School of Medicine, University of California San Diego, La Jolla, CA 92037,
16 USA

17 ⁵Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York,
18 NY 10029, USA

19 ⁶Tisch Cancer Institute, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA;
20 Department of Medicine, Hematology and Medical Oncology, Icahn School of Medicine at Mount Sinai,
21 New York, NY 10029, USA; Department of Oncological Sciences, Icahn School of Medicine at Mount
22 Sinai, New York, NY 10029, USA; Department of Pathology, Icahn School of Medicine at Mount Sinai,
23 New York, NY 10029, USA

24 ⁷Divisions of Immunobiology and Biomedical Informatics, Cincinnati Children's Hospital, Cincinnati,
25 OH 45229, USA; Department of Pediatrics, University of Cincinnati College of Medicine, Cincinnati, OH
26 45257, USA

27 ⁸Department of Microbiology, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA;
28 Global Health and Emerging Pathogens Institute, Icahn School of Medicine at Mount Sinai, New York,
29 NY 10029, USA; Division of Infectious Diseases, Department of Medicine, Icahn School of Medicine at
30 Mount Sinai, New York, NY 10029, USA

31 ⁹Department of Microbiology, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA;
32 Global Health and Emerging Pathogens Institute, Icahn School of Medicine at Mount Sinai, New York,
33 NY 10029, USA

34 ¹⁰The State Key Laboratory of Reproductive Regulation and Breeding of Grassland Livestock, College of
35 Life Sciences, Inner Mongolia University, Hohhot, 010070, China

36 ¹¹These authors contributed equally

37 ¹²These senior authors contributed equally

38 *Corresponding and Lead author: ivan.marazzi@mssm.edu
39
40

41 **ABSTRACT**

42
43 **Influenza A virus (IAV) is a threat to mankind because it generates yearly epidemics and poorly**
44 **predictable sporadic pandemics with catastrophic potential. IAV has a small RNA genome**
45 **composed of 8 mini-chromosomes (segments) that constitute a 5'UTR followed by a coding region**
46 **and a 3'UTR. Transcription of IAV RNA into mRNA depends on host mRNA, as the viral**
47 **polymerase cleaves 5'm7G-capped nascent transcripts to use as primers to initiate viral mRNA**
48 **synthesis. We hypothesized that captured host transcripts bearing AUG could drive the expression**
49 **of upstream ORFs in the viral segments, a phenomenon that would depend on the translatability of**
50 **the viral 5'UTRs. Here we report the existence of this mechanism, which generates host-virus**
51 **chimeric proteins. We label these proteins as Upstream Flu ORFs (UFO). Depending on the frame,**
52 **two types of host-virus UFO proteins are made: canonical viral proteins with human-derived N**
53 **term extensions or novel uncharacterized proteins. Here we show that both types are made during**

54 **IAV infection. Sequences that enable chimeric protein synthesis are conserved across IAV strains,**
55 **indicating that selection allowed the expansion of the proteome diversity of IAV in infected cells to**
56 **include multiple human-virus proteins.**

57
58 **Keywords:** Influenza, uORFs, viral evolution, chimeric protein
59

60 INTRODUCTION

61
62 Influenza A virus (IAV), of the family *Orthomyxoviridae*, is a highly contagious human and animal
63 pathogen responsible for significant levels of morbidity and mortality worldwide. The virus bears a
64 single-stranded, negative-sense RNA genome that is organized into eight segments (Bouvier and Palese,
65 2008). Viral mRNA transcription and genome replication both occur within the host nucleus, and require
66 the three-subunit viral RNA-dependent RNA polymerase (RdRP) complex comprising of PB1, PB2 and
67 PA proteins (Bouvier and Palese, 2008; Te Velthuis and Fodor, 2016).
68

69 IAV viral mRNA synthesis is primed using 5' methyl-7-guanosine (m7G) capped short RNA sequences
70 cleaved from host RNA polymerase II (RNAPII) dependent transcripts. During this process, named “cap-
71 snatching”, PA cleaves host-capped RNA bound to PB2 to generate 7-20 nucleotide long, capped RNA
72 fragments (Dias et al., 2009). These host-derived fragments are then utilized by PB1 to initiate the
73 transcription of viral mRNAs (Plotch et al., 1981; Reich et al., 2014). Consequently, IAV viral mRNAs
74 are not only genetic hybrids that include both host and viral derived sequences, but also possess diverse 5'
75 sequence heterogeneity (Kopstein et al., 2015; Sikora et al., 2017). Once made, viral mRNA is exported
76 to the cytoplasm and translated by the host machinery.
77

78 Each segment of the IAV genome encodes one major open reading frame (ORF) flanked by 5' and 3'
79 untranslated regions (UTRs). The IAV segments code for eight major structural and non-structural
80 proteins (PB2, PB1, PA, HA, NA, NP, M1, NS1). In addition, IAV utilizes several different mechanisms
81 to expand the coding capacity of the individual segments to generate additional proteins. Segments 7 and
82 8 are spliced to produce M2 and NEP proteins, respectively (Inglis and Brown, 1981; Lamb and Lai,
83 1980; Lamb et al., 1981). PB1-F2 and N40 (Chen et al., 2001; Wise et al., 2009) proteins arise from leaky
84 ribosomal scanning of IAV segment 2. Segment 3 encodes an alternative protein, PA-X, generated by +1-
85 ribosomal frameshift during the translation of PA protein (Jagger et al., 2012). Segment 3 also produces
86 several N-terminally truncated forms of PA due to alternate start codon (AUG) usage (Muramoto et al.,
87 2013). Additional viral proteins, such as M42, might be encoded from alternative spliced mRNAs (Wise
88 et al., 2012).
89

90 Intriguingly, full genome studies on IAV isolates have revealed that the length and sequence context of
91 these accessory proteins varies between IAV strains. These differences are often correlated with altered
92 virulence and/or responses of host cells. For example, PB1-F2 protein derived from the IAV laboratory
93 strain A/H1N1/Puerto Rico/8/1934 induces apoptosis in host cells through the interaction with BAK/BAX
94 (Chen et al., 2001). In contrast, PB1-F2 from a H5N1 IAV strain (A/H5N1/Hong Kong/156/1997) is non-
95 mitochondrial and not pro-apoptotic (Chen et al., 2010). As such, identification of novel accessory
96 proteins and the breadth of their diversity across different strains of IAV may provide insight into viral
97 replication and the interplay with the infected host.
98

99 In this manuscript, we describe the existence of IAV-human protein chimeras. We show that at least three
100 such chimeric proteins are synthesized during influenza virus infection. These proteins are initiated from
101 cap-snatched RNA sequences with upstream AUGs (uAUGs) that initiate translation of the IAV 5' UTR
102 and the downstream viral segment. Through this mechanism, host uAUGs create either viral protein N-
103 terminal extensions and/or the synthesis of novel, heretofore uncharacterized host-viral proteins (UFOs).
104 We show that both types of proteins are expressed in infected cells, as our analyses reveal the existence of

105 HA and NP extensions driven by host RNA and also identify an uvORF in segment 2 that generates a
106 novel, ~77 amino acid long protein (PB1 Upstream Flu ORf; PB1-UFO). Full length PB1-UFO is
107 conserved in more than 90% of isolates of IAV. HA and NP extensions are conserved in 99% of IAV
108 isolates. Overall, our analysis reveal that host-viral protein chimeras are (1) segment-specific, (2)
109 conserved across IAV strains and (3) undergo differential selection pressures according to 5' UTR and
110 coding region constrains, resulting in fixation of N term extensions and novel ORFs that are sampling
111 evolutionary space through genomic overprinting.

112

113 RESULTS

114

115 IAV 5' "UTRs" are potentially translatable

116 IAV transcription is initiated by host RNA cap snatching (**Figure 1A**). This process generates 5' host
117 derived extension of IAV segments. We hypothesized that this mechanism, used to express canonical
118 viral proteins (**Figure 1B**, Outcome 1), could generate upstream host-virus chimeric ORFs with coding
119 potential. Depending on the reading frame, an upstream host derived AUGs may either initiate the
120 synthesis of N-terminal extended viral proteins (**Figure 1B**, Outcome 2) or novel uvORFs that overprint
121 the canonical viral ORF (**Figure 1B**, Outcome 3). These outcomes are contingent on three assumptions:
122 (1) premature stop codons are not present in translation frames of interest (2) viral "UTRs" lacking stop
123 codons are evolutionarily conserved in IAV strains (3) AUGs are present in cap-snatched host sequences
124 enabling translation of host-virus chimeric RNA.

125

126 To address the first two points, we analyzed the nucleotide sequence variability within the 5'UTRs of all
127 eight segments among all H1N1 strains available from the GISAID Database (Shu and McCauley, 2017).
128 5'UTRs are highly conserved within each individual segment, as shown by the positional weight matrices
129 (**Figure 1C**, top panels and **Figure S1**). To determine if viral 5'UTRs can be translated to generate long
130 peptides, we retrieved the most commonly occurring 5'UTR nucleotide sequences per segment (**Figure**
131 **S1**). These sequences were then translated in all three strains *in silico* (**Figure 1C**). This revealed that the
132 UTRs of 3 (PB1, PA and M) of 8 viral segments possess conserved stop codons in-frame and upstream of
133 the major ORF start codon. Thus, 5/8 viral segments have the potential to code for N-terminally extended
134 viral proteins if an upstream AUG is captured from host mRNA. Surprisingly, we also detected the
135 absence of stop codons in the alternate translation reading frames of several viral segments (**Figure 1C**,
136 Segments PB2, PB1, PA, NA and HA). We were thus intrigued with the possibility that these segments
137 encode novel long peptides given an upstream, host-donated AUG in the right context.

138

139 Host RNA bearing upstream AUGs are present in viral mRNA

140 We next determined the abundance of AUGs in host-snatched sequences generated from PR8 infected-
141 A549 cells by RNA-sequencing. Host oligonucleotides with AUG codons constituted approximately
142 ~12% of all cap-snatched sequences, and were present at similar ratios in all eight segments of the virus
143 (**Figure 2A**).

144

145 Do uAUGs result in N-terminal extensions of viral proteins or generate uvORFs *in silico*? We aligned and
146 extended viral derived sequences from sequenced host-virus RNA chimeras to match the reference
147 sequences of the A/H1N1/Puerto Rico/8/1934 IAV. Viral UTRs with uAUGs were then translated *in*
148 *silico* revealing all possible N-terminal protein extension and/or uvORFs in the data set. We define N-
149 terminal protein extensions as ORFs with uAUGs in frame with the canonical ORFs and without a stop
150 codon in the uORF. By contrast, we define ORFs with uAUGs out of frame with the canonical ORF as
151 uvORFs. Putative sequences that would generate a novel ORF but would not contain a stop codon across
152 the whole length of a viral segment were excluded because of the inherent instability of mRNA lacking
153 stop codon (Simms et al., 2017). In an effort to be stringent with our analysis, host sequences that begun
154 with an AUG at the 5' were also removed from our analysis as it is unclear if the ribosome would be able
155 to recognize these as start codons (**Figure 2A**, yellow bars). uvORF length filters were also not applied at

156 this stage of the analysis as we reasoned that the ribosomal complex initiation and assembly should be
157 independent of ORF length (Chew et al., 2016).

158
159 We mapped host-derived sequences in PR8 infected A549 via RNA sequencing. Our analysis revealed
160 that host-derived uAUGs are present in all three translational reading frames, and at similar frequencies in
161 the eight viral segments (**Figure 2B**). As expected (**Figure 1C and S1**), individual viral segments
162 exhibited different propensities to generate N-terminally extended proteins (Orange bars; **Figure 2C**)
163 compared to uvORFs (Blue bars; **Figure 2C**). ~19% of uAUGs in host-derived sequences in PB2, HA,
164 NP, NA and NS segments (**Figure 2C**) initiate N-terminal extensions of the major ORF, but next to none
165 in PB1, PA and M segments. uvORFs are present in all segments at significant frequencies. Given that
166 that viral genes are among the highest expressed RNA in the cells during infection, this suggests that
167 uvORF containing viral RNAs are likely to be present at levels similar to most other host mRNA in the
168 cell.

169
170 **Host-derived uAUGs drive the translation of uvORFs and N-terminal extensions during infection.**

171 We next sought to determine if viral N-terminal extensions or uvORFs are translated during infection. In
172 silico analyses suggest that, as a function of the frame (F), the probable N-terminal extensions in PB2,
173 HA, NP, NA, NS segments are very consistent in lengths within the given individual segment (**Figure**
174 **3A**). Extensions ranged from 9 – 17 amino acids, with the longest occurring in the NP segment (**Figure**
175 **3A**). In contrast, the lengths of uvORFs in the PB2, PB1, PA, HA and M segments hovered at or below 20
176 amino acids (**Figure 3B**). Most importantly, we found long conserved uvORFs in PB2, PB1, PA and HA
177 segments, ranging from 40+ residues (HA) to nearly 80 residues (PB1).

178
179 If snatched host uAUGs initiate translation of viral 5'UTRs, these sequences should be enriched in
180 translating ribosomes in infected cells. Using RNA-seq and Ribo-seq, we mapped the 5' end of ribosome
181 footprint sequences from harringtonine-treated PR8-infected cells to the viral genome. This revealed an
182 accumulation of reads ~12nt upstream of the IAV canonical start codons in all eight segments (**Figure**
183 **3C**). Notably, we observed a large number of reads in the host-derived portion of the 5' UTR, consistent
184 with ribosome initiation. Furthermore, host sequences demonstrated a 7.5-fold enrichment in the Ribo-
185 seq data set vs. the host primer sequences present in the RNA-seq data set of poly A containing IAV
186 mRNA (**Figure 3D**).

187
188 If ribosomes initiate on host derived AUGs, many of the 5' sequences will be too short to extend from the
189 ribosome, given the brevity of their snatched caps, making P-site phasing problematic by standard
190 Riboseq analysis. We therefore used the location of AUGs within the primers to identify the reading
191 frame being translated. With few exceptions, initiation occurred evenly in all three reading frames. AUG
192 codons tended to aggregate closer to the transcriptional start site, despite being depressed at the -4
193 position in all segments. Frequencies also tended to be lower towards the 5' end of the primer (**Figure**
194 **S2A**). This phenomenon is also observed when the frequency of primers containing AUG is compared to
195 primer length (**Figure S2B**).

196
197 Finally, to verify that chimeric proteins are indeed translated, using targeted proteomic analysis we
198 evaluated the presence of UTR-derived and chimeric peptides in PR8 IAV infected cells. We
199 unequivocally identified peptides that originate from predicted N-terminal extensions of the NP and HA
200 and the long uvORF present in PB1 segment (**Figure 3E**).

201
202 Most host-derived sequences were from protein coding genes (**Figure 3F**), with similar distributions
203 between the three segments (**Figure 3G, top panels**). Host caps were derived from different genes
204 (**Figure 3G, bottom panels**) and were predominantly obtained from high expressing mRNAs (**Figure**
205 **S2C**).

206

207 **PB1-UFO is a host-virus chimeric protein expressed during infection**

208 We were especially intrigued by the ~77-amino acid uvORF present in the 5' UTR of IAV segment 2
209 (encoding PB1). This is one of the longest, conserved non-canonical ORFs in IAVs (**Figure 4A**). We
210 designate this protein PB1-UFO and proceeded to characterize it.

211
212 To evaluate the physiological role of PB1-UFO during infection we used reverse genetics to generate
213 wild-type control (CTRL-3 and CTRL+9) and PB1-UFO-deficient mutant (KO-3 and KO+9) IAVs in the
214 H1N1/Puerto Rico/8/1934 (PR8) strain background (**Figure 4B**). As PB1-UFO is predicted to translate
215 from an alternative -2 reading frame, we could make single nucleotide substitutions to introduce
216 premature stop codons for PB1-UFO without modifying the amino acid sequence of PB1 (**Figure 4B**).
217 PB1-UFO truncating mutations (KO-3 and KO+9) were introduced at either the -3 or +9 nucleotides
218 relative to the PB1 start ATG codon. We included viruses with mutations that did not disrupt either the
219 PB1-UFO or PB1 reading as controls (CTRL-3 or CTRL+9). Mutant and control viruses all yielded
220 stocks with similar particle counts as indicated by HA titers (**Figure S3A and S3B**). All viruses
221 demonstrate similar growth in MDCK cells (**Figure S3A and S3B**) at high (40°C; **Figure S3A**) or
222 physiological (37°C; **Figure S3B**) temperatures, demonstrating that PB1-UFO is not required for
223 replication under these conditions.

224
225 PB1-UFO-deficient viruses were also able to replicate normally in the lungs of infected BALB/c mice as
226 measured by virus yields after intranasal infection. Mice infected with PB1-UFO-deficient or control
227 viruses displayed weight loss (**Figures 4C, middle panels**) over a range of infecting doses. Survival
228 curves of mice infected with PB1-UFO-deficient viruses and controls revealed similar minimum lethal
229 doses (MLD₅₀) for both the KO-3 and KO+9 mutant viruses (**Figure 4C, column 1 and 3**) when
230 compared to CTRL-3 and CTRL+9 viruses (**Figure 4C, column 2 and 4**). These data indicate that PB1-
231 UFO is non-essential for virulence in mice, a phenomenon that is shared by most non-canonical protein
232 from IAV. Since mutations in accessory IAV proteins, which cause subtle differences in pathogenesis,
233 often display molecular phenotypes, we therefore isolated RNA from the lungs of mice infected with 100
234 PFUs of KO-3, CTRL-3, KO+9 or CTRL+9 viruses in biological replicates for RNA-seq analysis at days
235 three (n=2 per condition) and six post infection (n=3 per condition; **Figure S4A-S4D**).

236
237 RNA-seq showed that viral RNA (vRNA) was transcribed at similar levels between PB1-UFO mutant and
238 control viruses (**Figure S4A**), consistent with the minimal difference in viral lung titers (**Figure 4C, top**
239 **panels**). Importantly, mutant and control viruses exhibited a distinct transcriptome signature at day six
240 but not day three post infection (**Figure S4B and S4C**), as observed through differential gene expression
241 analysis (**Figure S4C**). Our analysis also suggested that the two PB1-UFO mutant viruses behaved
242 similarly to each other during infection (**Figure S4B**; Comparison m3 v p9), supporting the conclusion
243 that the difference in gene expression between control and mutant viruses are due to loss of PB1-UFO,
244 and not just alterations in viral RNA sequences. Similar trends in gene expression differences are present
245 in the top 32 differentially expressed genes (**Figure S4C**). Genes differentially expressed at day 6 post
246 infection are predominantly related to angiogenesis and protein folding (**Figure S4D**).

247
248 To check whether PB1-UFO expression could be detected by the immune system, we inserted the
249 SIINFEKL (SIIN) model MHC class I peptide (SIIN) into the 5' UTR upstream of the native PB1
250 initiation codon and in frame with PB1-UFO in a recombinant PR8 virus (**Figure 4D**). SIIN is efficiently
251 processed by the class I pathway and generates a high affinity complex with the mouse Kb class I
252 molecule that can be detected on the cell surface with high specificity and sensitivity by the 25-D1.16
253 mAb (Porgador et al., 1997) (**Figure S3C**). HEK293Kb cells infected with the PB1-UFO (SIIN) virus
254 demonstrated increased staining in flow cytometry relative to control uninfected cells or cells infected
255 with wildtype PR8 virus (**Figure S3D**). Extending this finding, mouse DC2.4 cells infected with PB1-
256 UFO (SIIN) PR8 activated transgenic OT-I CD8⁺ T cells (highly specific for Kb-SIIN (Hogquist et al.,
257 1994)) as determined by upregulation of CD25 and CD69 and also induced OT-I T-cell proliferation

258 when compared to the negative controls. (**Figure 4E**). As a positive control, we used a recombinant IAV
259 expressing SIIN(PB1-Ub-SIIN) at high levels (Wei et al., 2019). These data confirm that PB1-UFO is
260 translated, expressed during infection and that T cell immunosurveillance extends to peptides encoded by
261 uvORFs.

262

263 **Conservation of host-IAV proteins across strains and time**

264 To establish the contribution of the PB1-UFO to viral fitness we examined its conservation among all
265 H1N1, H3N2 and H5N1 IAV strains deposited in public sequence databases (Shu and McCauley, 2017).
266 Due to its high mutation rate, IAV evolution occurs extremely rapidly, and conservation of the ORF
267 provides strong evidence for its contribution to IAV transmission in its natural hosts. The PB1-UFO is
268 highly conserved across these three virus subtypes, all of which encode proteins of similar length and
269 amino acid composition (**Figure 5A**). Truncating mutations occurred relatively infrequently, at 8% of
270 H1N1 sequences, 3% of H5N1 sequences, and not present in H3N2 isolates.

271

272 We next designed a statistical model (**Figure S5A**) to query whether the conserved length of PB1-UFO
273 occurs more frequently than expected by chance. To increase statistical power, we focused only on
274 sequences derived from the H3N2 subtype of viruses as they had the most abundant number of full-length,
275 unique PB1 segment sequences. Over 92% of H3N2 PB1-segment sequences encode a 77-amino acid
276 PB1-UFO (**Figure 5B and 5C**). This is highly significant given that the random mutation model predicts
277 an average ORF of ~19 amino acids (**Figure S5B**). Likewise, analyses on synonymous mutations showed
278 that PB1-UFO is highly likely to maintain a long amino acid sequence pattern, implying that the
279 maintenance of longer sequences is due to protein function rather than the random production of short
280 peptides (**Figure 5D**).

281

282 **Evolutionary analyses on chimeric protein maintenance**

283 In an effort to quantify and infer selection of PB1-UFO in H3N2 strain of influenza virus over time we
284 used a frequency propagator model (Luksza and Lassig, 2014; Strelkova and Lässig, 2012) (**Figure 6A-**
285 **6B**). We compared the likelihoods of non-synonymous and/or synonymous mutations to reach fixation in
286 the PB1-UFO coding sequence of the IAV 5' UTR (R1, **Figure 6C**) to the corresponding likelihood of
287 synonymous mutations, which should evolve at near neutrality, occurring in the main PB1 coding
288 sequence (R3, **Figure 6C**). A similar analysis was done using the nucleotide sequences where PB1-UFO
289 and PB1 ORF overlap (R2, **Figure 6D**). Our analysis suggests that nucleotide mutations are overall
290 strongly repressed in the IAV 5'UTR, consistent with its role in priming viral transcription and viral
291 packaging (**Figure 6C**, Black line). Despite this, the fixation probabilities of synonymous mutations
292 occurring in PB1-UFO (**Figure 6C**, red line; Frequency propagator ratio =0.263+/-0.094) were 2 fold
293 increased over that of non-synonymous mutations (**Figure 6C**, blue line; Frequency propagator ratio
294 =0.134+/-0.040). This suggests that mutations that preserved the PB1-UFO peptide sequence are better
295 tolerated within the viral UTR. In contrast, when the nucleotide sequence of PB1-UFO overlapped PB1
296 (R2, **Figure 6D**), there was no difference in fixation rates between synonymous or non-synonymous
297 mutations (Frequency propagator ratio = 0.924 +/-0.650 (synonymous; red lines) and 1.121+/-0.225 (non-
298 synonymous; blue lines) (**Figure 6D**), indicating that changes to amino acid sequence are more tolerated
299 in the C-terminal of PB1-UFO. Taken together, our analyses suggest that while selection is heterogeneous
300 across the PB1-UFO frame, there is a positive selection pressure to maintain both the PB1-UFO protein
301 length, and N-terminal sequences.

302

303 Similarly, HA- and NA- UFO extensions are conserved more than 99% of H3N2 and are under positive
304 selection (**Figure 6E**). Overall, our result indicates that mutations that disrupt the amino acid sequences
305 of PB1-UFO and/or viral extensions are not well tolerated. IAV thus maintains the capacity, throughout
306 strains and time, to encode for chimeric proteins. This indicates a role for such host-dependent viral
307 protein diversity in viral tropism and life cycle.

308

309 **DISCUSSION**

310

311 In this manuscript, we describe the existence of a novel mechanism employed by IAV to generate hitherto
312 uncharacterized host-virus chimeric proteins. This mechanism employs the generation of host-virus
313 chimeric RNAs that are translated into chimeric proteins. We show that during IAV infection, two
314 classes of chimeric proteins are made: (1) viral proteins with host-encoded N-terminal, and (2) chimeric
315 host-virus proteins with novel open reading frames, which we termed uvORF proteins. We show that
316 these gene products are expressed in infected cells, surveilled by CD8+ T cells and modulate the antiviral
317 response.

318

319 **Chimeric UFO proteins: Novel proteins and N-terminal extensions**

320 In human genes, there is increasing evidence that upstream start codons (uAUGs) in the 5' UTR initiate
321 translation of short ORFs (Calvo et al., 2009; Wang and Rothnagel, 2004). uAUGs/uORFs are thought to
322 be mainly important in regulating expression of downstream ORFs by controlling ribosomal scanning
323 efficiencies (Calvo et al., 2009). However, there is evidence that suggests that some uORFs encode
324 biologically active peptides that contribute to evolutionary fitness (Andrews and Rothnagel, 2014;
325 Combier et al., 2008; Wen et al., 2009).

326

327 We have characterized the N-terminal HA and NP extensions, as well as PB1-UFO, because they could
328 be identified unambiguously by analyzing chimeric host-virus or UTR derived peptide through mass
329 spectrometry (whose sequence do not exist in 'conventional' human and viral proteome databases). The
330 expression of other uvORF proteins remains to be determined. It is important to recognize that other
331 uvORFs or N terminal extensions, like the chimeric HA and NP described here, might be difficult to
332 detect due to their N-term heterogeneity and partial overlapping sequences with the canonical protein.

333

334 In fact, we showed that based on the length of host snatched sequences and the viral UTRs, the chimeric-
335 protein extension bear N-termini consisting essentially of hypervariable peptides encoded by host-derived
336 RNA. In the cell, proteins containing variable sequences ("quasi protein-species") can be generated in
337 expressed proteins under normal conditions. This occurs through natural errors in protein synthesis as the
338 translation apparatus is tuned to optimize the occurrence of semi-random amino acid substitutions.
339 Translational fidelity is adaptive, maintained by cell and tissue type, and likely functions to cushion stress
340 (Ribas de Pouplana et al., 2014). For instance, conditions of oxidative stress alters the specificity of
341 Methionine(Met)-amino acyl synthetase, increasing Met-charging of non-Met tRNA to increase the Met
342 content of proteins (Netzer et al., 2009). This presumably protects them from oxidative damage (Levine et
343 al., 1996). Instead of relying wholly on adaptive mis-translation, IAV uvORFs appear to have a built-in
344 mechanism to diversify their proteome during infection. One intriguing hypothesis is that usage of
345 human-derived protein appendices might 'confound' MHC-class I surveillance.

346

347 In a similar vein, uORFs are particularly common in host cell mRNAs encoding regulatory and stress-
348 responsive proteins (Bondke Persson et al., 2015; Starck et al., 2016; Young and Wek, 2016), suggesting
349 that these genetic elements respond to changes in the cell's environment. Stress, leads to global
350 translational repression and preferential usage of uAUGs. This results in pervasive translation of human
351 uORFs as documented in cancer cells (Sendoel et al., 2017) and activated T cells (Starck et al., 2016). In
352 line with this, our data indicate that uAUGs are particularly abundant in high expressing genes that are
353 cap-snatched by IAV. By generating human-viral mRNA chimeras during infection IAV may be co-
354 opting the altered host mRNA expression to drive the expression of its newly expanded proteome.

355

356 **Evolutionary considerations: Overprinting and the mis-naming of UTRs**

357 Genetic overprinting typically occurs when a pre-existing reading frame acquires mutations that enable
358 translation in alternative reading frames while maintaining function of the ancestral frame. This is an
359 important mechanism to create new proteins, especially in the context of compact genomes (viral,

360 prokaryotic, eukaryotic organelles) with little coding capacity (Keese and Gibbs, 1992; Kovacs et al.,
361 2010; Poulin et al., 2003; Sabath et al., 2012).

362
363 Alternative reading frames created by overprinting can be translated by two mechanisms. One way is via
364 leaky scanning ribosomes that bypass the canonical AUG and decode a downstream out-of-frame
365 initiation codon. Important viral virulence factors, like PB1-F2 from influenza virus, are generated by
366 such a mechanism (Chen et al., 2001). Alternative reading frames may also be translated via ribosome
367 frameshift, in which ribosomes slip and skips (either forward or backward) one or two nucleotides to shift
368 to a new reading frame. IAV uses this process to create PA-X (Jagger et al., 2012). HIV also uses this
369 process to express and regulate the expression of Gag and Gag-Pol proteins, which are encoded by the
370 same ORFs (Fernandes et al., 2016).

371
372 While genetic overprinting could be selectively advantageous for some organisms, maintaining
373 overlapping ORFs requires additional regulatory mechanisms (e.g. regulating dynamic expression of
374 multiple proteins upon stimulation and/or ways to stop expression of one ORF in favor of the second).
375 These limitations, along with the fact that a mutation in one ORF will also often affect a second ORF,
376 ends up imposing too many constraints, thus limiting the functional evolutionary space that pathogens
377 require to sample as a mean to adapt to hosts.

378
379 PB1-UFO represents a unique product of overprinting because it is encoded by sequences from two
380 organisms: virus and host, with host sequences providing translatability to viral UTR sequences. Our
381 analyses suggest that PB1-UFO is undergoing stabilizing selection in the 5'UTR, where divergent forms
382 of the protein, generated by non-synonymous mutations, appear to be preferentially removed from the
383 population. This implies that PB1-UFO support viral fitness, as we would not otherwise expect
384 differences in fixation probabilities of synonymous or non-synonymous mutations occurring in the IAV
385 5'UTR.

386 387 **A new player in the host-pathogen arms race**

388 The capacity of a pathogen to overcome host barriers and establish infection is based on the expression of
389 pathogen-derived proteins. To understand how a pathogen antagonizes the host and establishes infection
390 we need to have a clear understanding of what protein a pathogen encodes, how they function, and the
391 manner in which they contribute to virulence. The current dogma about many life-threatening pathogens
392 is that they encode a small, finite number of proteins because of their limited genomes. RNA viruses,
393 including IAV, are a prime example of this paradigm. We now show that there is another level of
394 complexity to this equation.

395 396 **AUTHOR CONTRIBUTIONS**

397
398 Conceptualization, A.G.-S., J.W.Y. and I.M.; Methodology, I.M., J.W.Y., Y.M., M.A., G.W., J.H.;
399 Formal Analysis, Y.M., M.A., G.W., J.H., N.Z., J.N., N.M., J.G., J.W., J.J., M.C., Z.P., H.v.B., M.L.,
400 E.R.M.; Investigation, Y.M., M.A., G.W., J.H., N.Z., J.N., N.M., J.J., M.C., Z.P., H.v.B., M.L., E.R.M.;
401 Resources, Y.M., M.A., J.J., M.C., H.v.B., E.R.M., A.G.-S.; Writing – Original Draft, I.M., J.W.Y.,
402 Writing – Review & Editing, I.M., J.W.Y., Y.M., J.H., M.A., G.W., H.v.B., M.L., B.D.G., E.R.M., A.G.-
403 S.; Visualization, Y.M., M.A., G.W., J.H., J.J., M.C., Z.P., E.R.M.; Funding Acquisition, A.G.-S., I.M.;
404 Supervision, I.M.

405 406 **ACKNOWLEDGMENTS**

407
408 We thank the Genomics and Mouse facility at Icahn School of Medicine at Mount Sinai, the Global
409 Health and Emerging Pathogens Institute (GHEPI) at Mount Sinai, and the entire Marazzi Lab team. I.M.
410 is supported by Burroughs Wellcome Fund 1017892 and by Chan Zuckerberg Initiative 2018-191895.

411 I.M. and H.v.B. are supported by NIH grant R01AI113186. A.G.-S. and I.M. are supported by the NIH
412 grant U19AI135972 FLUOMICS.

413

414 **DECLARATION OF INTERESTS**

415

416 The authors declare no competing interests.

417 **FIGURE LEGENDS**

418

419 **Figure 1. IAV 5' UTR are conserved and translatable in all three reading frames.** (A) Schematic of
420 viral cap-snatching occurring during IAV infection. (B) Presence of upstream AUGs in host-derived
421 segments of viral mRNA may drive the formation of viral protein extensions or novel host-viral chimeric
422 proteins. (C; top panels) Predicted peptide sequences derived upon translation of all three ribosome
423 reading frames in the 5'UTR. (C; lower panels) Sequence conservation analysis of IAV H1N1 strain
424 5'UTR within individual viral segments.

425

426 **Figure 2. Upstream AUGs are present in host derived viral RNAs.** (A) Percentage of cap snatched
427 sequences bearing uAUGs in the first codon (yellow) and bearing uAUGs at other positions (blue). (B)
428 Incorporation of host transcript sequences increases the diversity of putative alternative start codons. For
429 each viral segment, the frequency and position of alternative start codons is shown relative to native start
430 of the viral genes. For each reading frame, the frequency and location of the first in-frame stop codon are
431 indicated. (C) Percentage of N-terminal protein extension (orange) and uvORFs (blue) derived from cap
432 snatched sequences bearing uAUGs in the 8 viral segments.

433

434 **Figure 3. IAV 5'UTRs are translated.** (A) Length distribution of N-terminal protein extension in
435 individual segments. Each dot represents a protein predicted to be translated from a unique host-virus
436 chimeric RNA. F represents the reading frame. (B) Length distribution of uvORFs of individual segments.
437 Each dot represents a protein predicted to be translated from a unique host-virus chimeric RNA. F
438 represents the reading frame. (C) Ribosome profiling of harringtonine-treated A549 cells infected with
439 A/Puerto Rico/8/1934 (H1N1). (D) Frequency of primers containing AUG codons in harringtonine Ribo-
440 seq and RNA-seq datasets. (E) Schematic shows peptides identified in HA, NP N-terminal extension and
441 PB1-UFO via Mass Spectrometry analysis. (F) Source of CS events contributing to PB1-UFO. (G)
442 Transcript biotypes contributing CS sequences for vmRNAs leading to uORFs and those that do not for
443 segments PB1, HA and NP (upper panels). Relative CS abundance for the top-100 genes contributing to
444 uORFs. Names and CS event read counts are shown for the top 5 genes (lower panels).

445

446 **Figure 4. PB1-UFO.** (A) Nucleotide and amino acid sequence of PB1-UFO protein. PB1-UFO peptide
447 detected in mass spectrometry is highlighted in yellow. (B) Schematic of mutations used to construct
448 PB1-UFO-null and control viruses in the PR8 virus background. (C) Viral titers (top panels), percentage
449 of body weight loss (middle panels) and survival curves (lower panels) of mice infected with differing
450 concentrations of the indicated viruses. MLD50 of each virus is indicated at the bottom. (D) Nucleotide
451 and amino acid sequence of PB1 and PB1-UFO N-terminal sequences with the SIINFEKL peptide
452 insertion in 5' UTR to generate the PB1-UFO (SIIN) virus. (E) Infected DC2.4 co-cultured OT-I
453 activation assay. CD69 and CD25 expression at 24 hours post co-culture, and cell proliferation assay at 48
454 hours post co-culture.

455

456 **Figure 5. Bioinformatics analysis on conservation of PB1-UFO protein sequences.** (A) Top five most
457 common PB1-UFO protein sequences in three Influenza A strains, H1N1, H3N2 and H5N1. (B) Density
458 plot of predicted length of H3N2 PB1-UFO protein sequences. Over 92% of sequences are predicted to
459 generate a protein of 77aa (medium blue), ~3% are shorter than 77aa (light blue), ~1.5% are longer than
460 77aa (dark blue), rest of sequences are predicted not to generate PB1-UFO protein (grey). (C) P value
461 distribution/volcano plot of H3N2 PB1-UFO protein sequence length. Each dot represents the difference
462 between observed length and expected length of each individual sequence. (D) The line plot shows the
463 number of synonymous mutations in frame of WT H3N2 PB1 (x-axis) that mutate stop codons in frame
464 of H3N2 PB1-UFO (y-axis).

465

466 **Figure 6. Evolutionary analysis on H3N2 PB1-UFO protein sequences.**

467 (A) Strain tree of H3N2 IAV viruses. Mutations occurring in the N-terminal PB1-UFO frame overlapping
468 the viral 5'UTR (region1, R1, top panel; yellow region) are indicated as color dots. (B) Same as in A, but
469 for mutations occurring in the C-terminal PB1-UFO frame overlapping the PB1 (region 2, R2, top panel;
470 yellow region). (C) Frequency propagator ratio of the indicated classes of mutations occurring in the N-
471 terminal PB1-UFO frame overlapping the viral 5'UTR (region1, R1; yellow). Fixation probabilities were
472 compared to those of synonymous mutations occurring in the region of PB1 that does not overlap PB1-
473 UFO (region 3, R3; blue). Error bars indicate sampling uncertainties. $g(X) < 1$: negative selection, $g(X)$
474 ≈ 1 : weak/heterogeneous selection; $g(X) > 1$: positive selection. (D) Same as in (C), but for C-terminal
475 sequences of PB1-UFO frame overlapping the PB1-frame (region 2, R2; yellow). (E) Percentage of
476 observed HA and NP N-terminal extension protein sequences.
477

478 **Figure S1. Viral 5'UTRs are conserved (Related to Figure 1).** Multiple sequence alignments of H1N1
479 IAV 5'UTRs per segment. The overall distribution of each unique nucleotide sequence is indicated on the
480 left, and the consensus sequence of each UTR is indicated below each alignment.
481

482 **Figure S2. IAV 5'UTRs are associated with active ribosomes (Related to Figure 2).** (A) Frequency of
483 AUG codons by position relative to the viral transcription initiation site. (B) Frequency of AUGs at each
484 position in primers compared to length for harringtonine-arrested ribo-seq and RNA-seq datasets.
485 Expected frequency shown in green. (C) Log (counts per million) of transcripts that are cap-snatched by
486 IAV or control transcripts.
487

488 **Figure S3. Controls for Figure 4 (Related to Figure 4).** (A) and (B) shows growth properties of viruses
489 in MDCK cells. Cells were infected with viruses at MOI of 0.001 and incubated at 40°C (A) and 37°C (B).
490 (C) Schematic of SIINFEKL mechanism of action. (D) SIINFEKL expression from 293Kb cells infected
491 with PB1-UFO (SIIN) virus.
492

493 **Figure S4. RNA-Seq analysis on PB1-UFO mutant and control virus infected mice (Related to**
494 **Figure 4).**

495 (A) Viral RNA levels in PB1-UFO mutant and control viruses in the lungs of infected mice at days 3 and
496 days 6 post infection. (B) Two factor model analyses of RNA sequencing data of PB1-UFO mutant and
497 control viruses infected mouse lungs at days 3 and days 6 post infection. (C) Heatmap showing the top 32
498 differentially expressed genes ($FDR < 0.1$, $|\text{Log}_2\text{FC}| > 1$) when comparing PB1-UFO mutant and control
499 virus infected lungs at day 6 post infection. (D) Gene ontology of genes predicted to be differentially
500 expressed during infection of control or PB1-UFO deficient viruses in mice.
501

502 **Figure S5. Controls related to Figure 5 (Related to Figure 5).**

503 (A) Schematic of predicted sequence length model of PB1-UFO proteins. (B) Density plot showing the
504 expected lengths of H3N2 PB1-UFO proteins, based on random codon-shuffled sequences.
505
506
507
508
509
510
511

512 MATERIALS AND METHODS

513

514 Cells

515 Human embryonic kidney 293T cells, Madin-Darby canine kidney (MDCK) cells, and Human lung
516 carcinoma epithelial A549 cells were maintained in Dulbecco's modified Eagle's medium (DMEM;
517 Corning) containing 10% newborn calf serum (FBS; Peak Serum) and antimicrobial drugs. Human. All
518 cells were maintained at 37 °C with 5% CO₂.

519

520 Viruses

521 Using plasmid-based reverse genetics (Fodor et al., 1999), we generated recombinant influenza viruses
522 with PB1 mutations by using the A/Puerto Rico/8/1934 (PR8) strain as the backbone. A wild-type (WT)
523 recombinant (PR8 WT) was generated, as well as four PB1 substitution mutants.

524 The first mutant, bearing a premature stop codon in the PB1-UFO protein at the position of three
525 nucleotides before the start of PB1 open reading frame, was named as (PB1-UFO KO-3). The second
526 mutant, which preserved expression of full length PB1-UFO protein even with a point mutation at the
527 position of three nucleotides before the start of PB1 open reading frame, was named as (PB1-UFO Ctrl-3).

528 This virus acted as a control of PB1-UFO KO-3. The third mutant containing a stop codon in the PB1-
529 UFO protein at the position of nine nucleotides after the start of PB1 open reading frame was named as
530 (PB1-UFO KO+9). The fourth mutant, which preserved the expression of full length PB1-UFO protein
531 even with a point mutation at the position of nine nucleotides after the start of PB1 open reading frame
532 was named as (PB1-UFO Ctrl+9). This virus acted as a control of PB1-UFO KO+9. Mutations were
533 confirmed by sequencing both plasmids and viruses. The stock virus titers were the average of three
534 independent experiments.

535

536 Growth kinetics of Viruses in Cell Culture

537 MDCK cells were infected with viruses at a multiplicity of infection (MOI) of 0.001, incubated for one
538 hour at 37 °C, washed twice, and then cultured with Opti-MEM and TPCK-treated trypsin at 40°C and
539 37°C for 72 h. Supernatants were collected at the indicated time points. Hemagglutination titer (HA) were
540 tested in 0.5% turkey red blood cells and virus titers were determined by plaque assay in MDCK cells.

541

542 Mouse studies

543 All mice procedures were performed following protocols approved by the Icahn School of Medicine at
544 Mount Sinai Institutional Animal Care and Use Committee (IACUC). All the animal studies were carried
545 out in strict accordance with the recommendations in the Guide for the Care and Use of Laboratory
546 Animals of the National Research Council. Eight-week-old female BALB/c mice were obtained from
547 Jackson Laboratories (Bar Harbor, ME). Mice were anesthetized by intraperitoneal injection of a mixture
548 of ketamine and xylazine before infection.

549 Groups of five mice were inoculated intranasally. with 100, 50, 25, 10, or 5 PFU of virus. Mice were
550 monitored daily for clinical signs of illness and weight loss. Upon reaching 75% of initial body weight,
551 animals were humanely euthanized with carbon dioxide (CO₂) as per the IACUC protocol.

552 Groups of five mice were intranasally (i.n.) infected with 100 plaque-forming unites (PFU) of viruses in a
553 volume of 50 µl, two and three mice were euthanized on 3 and 6 days post-inoculation (d.p.i.),
554 respectively. The middle lobe of the lung was collected for total RNA extraction, and the post-caval lobes
555 of the lung was collected to determine virus titers by plaque assay on MDCK cells.

556

557 RNA sequencing

558 After adaptor removal with cutadapt (Martin, 2011) and base-quality trimming to remove 3' read
559 sequences if more than 20 bases with Q <20 were present, paired-end reads were mapped to the mouse
560 (mm10) reference genome with STAR (Dobin et al., 2013), and gene-count summaries were generated
561 with featureCounts (Liao et al., 2013). DESeq2 (Love et al., 2014) was used to variance-normalize the

562 data before a 2-factor model (gene ~ ConditionTime + Mutant) was applied to identify differentially
563 expressed genes. RNA-seq raw data are deposited in GEO under accession GSE128519.

564

565 **Proteomic Strategy**

566 Mass spectrometry was performed using purified lysates obtained from PR8 IAV infected A549
567 cells. Targeted identification of chimeric proteins was conducted using datasets derived from the entire
568 human and IAV reference sequence merged with the set of predicted IAV uORFs and viral protein
569 extensions. Common contaminants were filtered out and missing values in the data matrix were attributed
570 an intensity score of 0.

571

572 **Ribo-seq analysis**

573 Ribosome footprint reads were trimmed with cutadapt (Martin, 2011), and aligned to the human (hg38)
574 and A/Puerto Rico/8/1938 (H1N1) genomes with STAR (Dobin et al., 2013). The 5' end mapping was
575 then performed for all reads aligning to the influenza genome. Host-derived transcriptional primer
576 sequences were extracted from reads with partially mapping to the 5' end of each segment. Analysis of
577 AUG composition was performed using custom in-house Perl scripts which are available upon request.

578

579 **Antigen Expression and T cell immunosurveillance Assays**

580 HEK293T cells stably expressing mouse K^b MHC-I (HEK293K^b) were infected with influenza A viruses.
581 At 18 hours post infection, cells were stained with Alexa 647-labelled MAb 25D-1.16 (anti-K^b-SIIN) to
582 measure surface expression of K^b-SIIN complexes flow cytometry. For T-cell activation assays, OT-I T-
583 cells were harvested from the spleen and lymph nodes of OT-1 transgenic mice and purified on the
584 AutoMACS with the CD8a+ T Cell Isolation Kit (Milteny, Germany), and stained with CellTrace Violet
585 (Thermo Fisher, Waltham, MA) DC2.4 cells were infected with influenza A viruses for 18 hours, and
586 then co-cultured with OT-I T-cells. T-cells were stained with anti-CD25 and anti-CD28 labeled antibodies
587 at 24 hours post co-culture for activation assays. T-cell proliferation assays were conducted at 48 hours
588 post infection by measuring CellTrace Violet staining by flow cytometry.

589

590 **Computational analyses**

591 **Sequence data set**

592 Our study is based on a data set of 26,472 human influenza A/H3N2 sequences available from the
593 GISAID database (Shu and McCauley, 2017), which contain 6,244 unique PB1 strains. We included only
594 full length sequences using a custom script that is available upon request.

595

596 **Random sequence model**

597 We constructed codon usage matrix for each of individual nucleotide sequence. Using the codon usage
598 table, a protein sequence in open reading frame is used as input to generate multiple random nucleotide
599 sequences. We then translate random nucleotide sequences to protein sequences in frame which may
600 generate PB1-UFO protein. Using a custom script, we calculate the average stop codon positions of
601 random PB1-UFO protein sequences as its expected value. By comparing with its expected value, we
602 determine the likelihood that the translated PB1-UFO sequence was obtained randomly and its deviation
603 from the expected PB1-UFO length.

604

605 **Strain tree reconstruction**

606 Our analysis is based on an ensemble of strain trees obtained from the PB1 sequence data set. Such trees
607 describe the genealogy of influenza strains resulting from an evolutionary process under
608 selection (Strelkova and Lässig, 2012). The tree ensemble is obtained with FastTree (Price et al., 2010),
609 which very time-efficiently reconstructs maximum-likelihood phylogenies. We use a general time-
610 reversible model. We refine the tree topology with RAXML (Stamatakis, 2014). Given the output topology,
611 we reconstruct maximum-likelihood sequences and timing of internal nodes with the TreeTime package
612 (Sagulenko et al., 2018).

613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634

Mapping of mutations

Maximum likelihood maps point mutations between directly related strains onto the branches of the tree. A mutation on a given branch marks an origination event of a single-nucleotide polymorphism, that is, the appearance of a nucleotide difference between the strains descending from that branch and its ancestral lineage. A reconstructed strain tree with all mapped mutations, which are partitioned into two classes: (a) synonymous mutations, (b) nonsynonymous mutations. These mutations are the basis of our fitness model (Luksza and Lassig, 2014).

Frequency propagator ratio analysis

Our analysis is based on a set of codons in PB1-UFO coding sequence overlapping the IAV 5'UTR (R1), the overlapping sequence between PB1-UFO and PB1 ORF (R2), and those in the main PB1 coding sequence (R3) respectively. To quantify selection on a class of mutations, we use the frequency propagator ratio

$$g(X) = \frac{G(X)}{G_0(X)}$$

where $G(X)$ is the likelihood that a mutation in a given class reaches frequency X , and $G_0(X)$ is the likelihood for synonymous mutations occurring in the main PB1 coding sequence, which should evolve near neutrality. To predict the evolutionary direction of a given subset of codons, we compare fixation probabilities (or probabilities of reaching high frequencies) of mutations in that region with those in the null-class region R3.

635 **REFERENCES**

636

637 Andrews, S.J., and Rothnagel, J.A. (2014). Emerging evidence for functional peptides encoded by short
638 open reading frames. *Nat Rev Genet* 15, 193-204.

639 Bondke Persson, A., Benko, E., Staudacher, J.J., Kasim, M., Persson, P.B., Ujvari, S.J., Fählng, M.,
640 Ostareck-Lederer, A., Ostareck, D.H., Naarmann-de Vries, I.S., *et al.* (2015). Hypoxia-induced gene
641 expression results from selective mRNA partitioning to the endoplasmic reticulum. *Nucleic Acids*
642 *Research* 43, 3219-3236.

643 Bouvier, N.M., and Palese, P. (2008). The biology of influenza viruses. *Vaccine* 26 *Suppl* 4, D49-53.

644 Calvo, S.E., Pagliarini, D.J., and Mootha, V.K. (2009). Upstream open reading frames cause widespread
645 reduction of protein expression and are polymorphic among humans. *Proc Natl Acad Sci U S A* 106,
646 7507-7512.

647 Chen, C.J., Chen, G.W., Wang, C.H., Huang, C.H., Wang, Y.C., and Shih, S.R. (2010). Differential
648 localization and function of PB1-F2 derived from different strains of influenza A virus. *J Virol* 84, 10051-
649 10062.

650 Chen, W., Calvo, P.A., Malide, D., Gibbs, J., Schubert, U., Bacik, I., Basta, S., O'Neill, R., Schickli, J.,
651 Palese, P., *et al.* (2001). A novel influenza A virus mitochondrial protein that induces cell death. *Nat Med*
652 7, 1306-1312.

653 Chew, G.L., Pauli, A., and Schier, A.F. (2016). Conservation of uORF repressiveness and sequence
654 features in mouse, human and zebrafish. *Nat Commun* 7, 11663.

655 Combier, J.P., de Billy, F., Gamas, P., Niebel, A., and Rivas, S. (2008). Trans-regulation of the
656 expression of the transcription factor MtHAP2-1 by a uORF controls root nodule development. *Genes*
657 *Dev* 22, 1549-1559.

658 Dias, A., Bouvier, D., Crepin, T., McCarthy, A.A., Hart, D.J., Baudin, F., Cusack, S., and Ruigrok, R.W.
659 (2009). The cap-snatching endonuclease of influenza virus polymerase resides in the PA subunit. *Nature*
660 458, 914-918.

661 Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M., and
662 Gingeras, T.R. (2013). STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 29, 15-21.

663 Fernandes, J.D., Faust, T.B., Strauli, N.B., Smith, C., Crosby, D.C., Nakamura, R.L., Hernandez, R.D.,
664 and Frankel, A.D. (2016). Functional Segregation of Overlapping Genes in HIV. *Cell* 167, 1762-
665 1773.e1712.

666 Fodor, E., Devenish, L., Engelhardt, O.G., Palese, P., Brownlee, G.G., and Garcia-Sastre, A. (1999).
667 Rescue of influenza A virus from recombinant DNA. *J Virol* 73, 9679-9682.

668 Hogquist, K.A., Jameson, S.C., Heath, W.R., Howard, J.L., Bevan, M.J., and Carbone, F.R. (1994). T cell
669 receptor antagonist peptides induce positive selection. *Cell* 76, 17-27.

670 Inglis, S.C., and Brown, C.M. (1981). Spliced and unspliced RNAs encoded by virion RNA segment 7 of
671 influenza virus. *Nucleic Acids Res* 9, 2727-2740.

- 672 Jagger, B.W., Wise, H.M., Kash, J.C., Walters, K.A., Wills, N.M., Xiao, Y.L., Dunfee, R.L.,
673 Schwartzman, L.M., Ozinsky, A., Bell, G.L., *et al.* (2012). An overlapping protein-coding region in
674 influenza A virus segment 3 modulates the host response. *Science* *337*, 199-204.
- 675 Keese, P.K., and Gibbs, A. (1992). Origins of genes: "big bang" or continuous creation?
676 *Proceedings of the National Academy of Sciences* *89*, 9489.
- 677 Koppstein, D., Ashour, J., and Bartel, D.P. (2015). Sequencing the cap-snatching repertoire of H1N1
678 influenza provides insight into the mechanism of viral transcription initiation. *Nucleic Acids Res* *43*,
679 5052-5064.
- 680 Kovacs, E., Tompa, P., Liliom, K., and Kalmar, L. (2010). Dual coding in alternative reading frames
681 correlates with intrinsic protein disorder. *Proceedings of the National Academy of Sciences* *107*, 5429.
- 682 Lamb, R.A., and Lai, C.J. (1980). Sequence of interrupted and uninterrupted mRNAs and cloned DNA
683 coding for the two overlapping nonstructural proteins of influenza virus. *Cell* *21*, 475-485.
- 684 Lamb, R.A., Lai, C.J., and Choppin, P.W. (1981). Sequences of mRNAs derived from genome RNA
685 segment 7 of influenza virus: colinear and interrupted mRNAs code for overlapping proteins. *Proc Natl*
686 *Acad Sci U S A* *78*, 4170-4174.
- 687 Levine, R.L., Mosoni, L., Berlett, B.S., and Stadtman, E.R. (1996). Methionine residues as endogenous
688 antioxidants in proteins. *Proc Natl Acad Sci U S A* *93*, 15036-15040.
- 689 Liao, Y., Smyth, G.K., and Shi, W. (2013). featureCounts: an efficient general purpose program for
690 assigning sequence reads to genomic features. *Bioinformatics* *30*, 923-930.
- 691 Love, M.I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and dispersion for
692 RNA-seq data with DESeq2. *Genome Biol* *15*, 550.
- 693 Luksza, M., and Lassig, M. (2014). A predictive fitness model for influenza. *Nature* *507*, 57-61.
- 694 Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing reads. *2011* *17*,
695 3.
- 696 Muramoto, Y., Noda, T., Kawakami, E., Akkina, R., and Kawaoka, Y. (2013). Identification of novel
697 influenza A virus proteins translated from PA mRNA. *J Virol* *87*, 2455-2462.
- 698 Netzer, N., Goodenbour, J.M., David, A., Dittmar, K.A., Jones, R.B., Schneider, J.R., Boone, D., Eves,
699 E.M., Rosner, M.R., Gibbs, J.S., *et al.* (2009). Innate immune and chemically triggered oxidative stress
700 modifies translational fidelity. *Nature* *462*, 522-526.
- 701 Plotch, S.J., Bouloy, M., Ulmanen, I., and Krug, R.M. (1981). A unique cap(m7GpppXm)-dependent
702 influenza virion endonuclease cleaves capped RNAs to generate the primers that initiate viral RNA
703 transcription. *Cell* *23*, 847-858.
- 704 Porgador, A., Yewdell, J.W., Deng, Y., Bennink, J.R., and Germain, R.N. (1997). Localization,
705 quantitation, and in situ detection of specific peptide-MHC class I complexes using a monoclonal
706 antibody. *Immunity* *6*, 715-726.

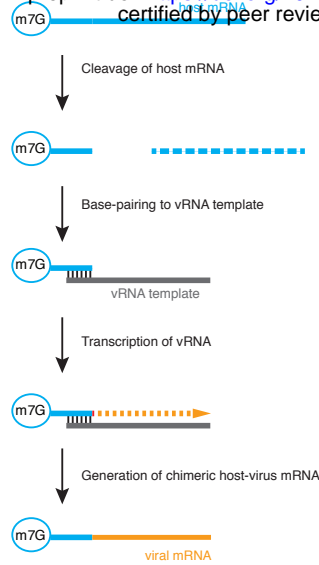
- 707 Poulin, F., Brueschke, A., and Sonenberg, N. (2003). Gene Fusion and Overlapping Reading Frames in
708 the Mammalian Genes for 4E-BP3 and MASK. *Journal of Biological Chemistry* 278, 52290-52297.
- 709 Price, M.N., Dehal, P.S., and Arkin, A.P. (2010). FastTree 2 – Approximately Maximum-Likelihood
710 Trees for Large Alignments. *PLOS ONE* 5, e9490.
- 711 Reich, S., Guilligay, D., Pflug, A., Malet, H., Berger, I., Crepin, T., Hart, D., Lunardi, T., Nanao, M.,
712 Ruigrok, R.W., *et al.* (2014). Structural insight into cap-snatching and RNA synthesis by influenza
713 polymerase. *Nature* 516, 361-366.
- 714 Ribas de Pouplana, L., Santos, M.A., Zhu, J.H., Farabaugh, P.J., and Javid, B. (2014). Protein
715 mistranslation: friend or foe? *Trends Biochem Sci* 39, 355-362.
- 716 Sabath, N., Wagner, A., and Karlin, D. (2012). Evolution of Viral Proteins Originated De Novo by
717 Overprinting. *Molecular Biology and Evolution* 29, 3767-3780.
- 718 Sagulenko, P., Puller, V., and Neher, R.A. (2018). TreeTime: Maximum-likelihood phylodynamic
719 analysis. *Virus Evolution* 4.
- 720 Sendoel, A., Dunn, J.G., Rodriguez, E.H., Naik, S., Gomez, N.C., Hurwitz, B., Levorse, J., Dill, B.D.,
721 Schramek, D., Molina, H., *et al.* (2017). Translation from unconventional 5' start sites drives tumour
722 initiation. *Nature* 541, 494-499.
- 723 Shu, Y., and McCauley, J. (2017). GISAID: Global initiative on sharing all influenza data - from vision to
724 reality. *Euro Surveill* 22.
- 725 Sikora, D., Rocheleau, L., Brown, E.G., and Pelchat, M. (2017). Influenza A virus cap-snatches host
726 RNAs based on their abundance early after infection. *Virology* 509, 167-177.
- 727 Simms, C.L., Yan, L.L., and Zaher, H.S. (2017). Ribosome Collision Is Critical for Quality Control
728 during No-Go Decay. *Mol Cell* 68, 361-373 e365.
- 729 Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
730 phylogenies. *Bioinformatics* 30, 1312-1313.
- 731 Starck, S.R., Tsai, J.C., Chen, K., Shodiya, M., Wang, L., Yahiro, K., Martins-Green, M., Shastri, N., and
732 Walter, P. (2016). Translation from the 5' untranslated region shapes the integrated stress response.
733 *Science* 351, aad3867.
- 734 Strelkova, N., and Lässig, M. (2012). Clonal Interference in the Evolution of Influenza. *Genetics* 192,
735 671.
- 736 Te Velthuis, A.J., and Fodor, E. (2016). Influenza virus RNA polymerase: insights into the mechanisms
737 of viral RNA synthesis. *Nat Rev Microbiol* 14, 479-493.
- 738 Wang, X.Q., and Rothnagel, J.A. (2004). 5'-untranslated regions with multiple upstream AUG codons can
739 support low-level translation via leaky scanning and reinitiation. *Nucleic Acids Res* 32, 1382-1391.
- 740 Wei, J., Kishton, R.J., Angel, M., Conn, C.S., Dalla-Venezia, N., Marcel, V., Vincent, A., Catez, F., Ferre,
741 S., Ayadi, L., *et al.* (2019). Ribosomal Proteins Regulate MHC Class I Peptide Generation for
742 Immunosurveillance. *Mol Cell*.

- 743 Wen, Y., Liu, Y., Xu, Y., Zhao, Y., Hua, R., Wang, K., Sun, M., Li, Y., Yang, S., Zhang, X.J., *et al.*
744 (2009). Loss-of-function mutations of an inhibitory upstream ORF in the human hairless transcript cause
745 Marie Unna hereditary hypotrichosis. *Nat Genet* 41, 228-233.
- 746 Wise, H.M., Foeglein, A., Sun, J., Dalton, R.M., Patel, S., Howard, W., Anderson, E.C., Barclay, W.S.,
747 and Digard, P. (2009). A complicated message: Identification of a novel PB1-related protein translated
748 from influenza A virus segment 2 mRNA. *J Virol* 83, 8021-8031.
- 749 Wise, H.M., Hutchinson, E.C., Jagger, B.W., Stuart, A.D., Kang, Z.H., Robb, N., Schwartzman, L.M.,
750 Kash, J.C., Fodor, E., Firth, A.E., *et al.* (2012). Identification of a novel splice variant form of the
751 influenza A virus M2 ion channel with an antigenically distinct ectodomain. *PLoS Pathog* 8, e1002998.
- 752 Young, S.K., and Wek, R.C. (2016). Upstream Open Reading Frames Differentially Regulate Gene-
753 specific Translation in the Integrated Stress Response. *Journal of Biological Chemistry* 291, 16927-16935.
754

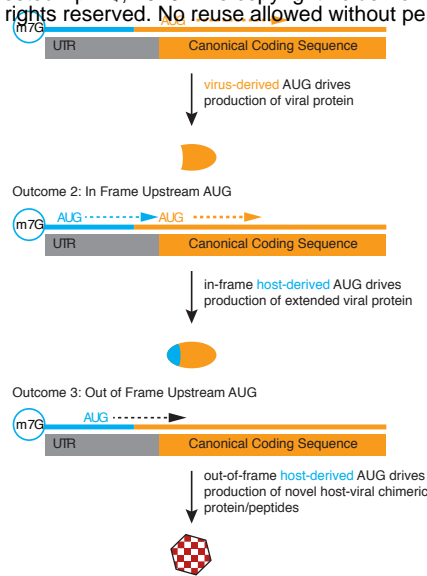
FIGURE 1

A

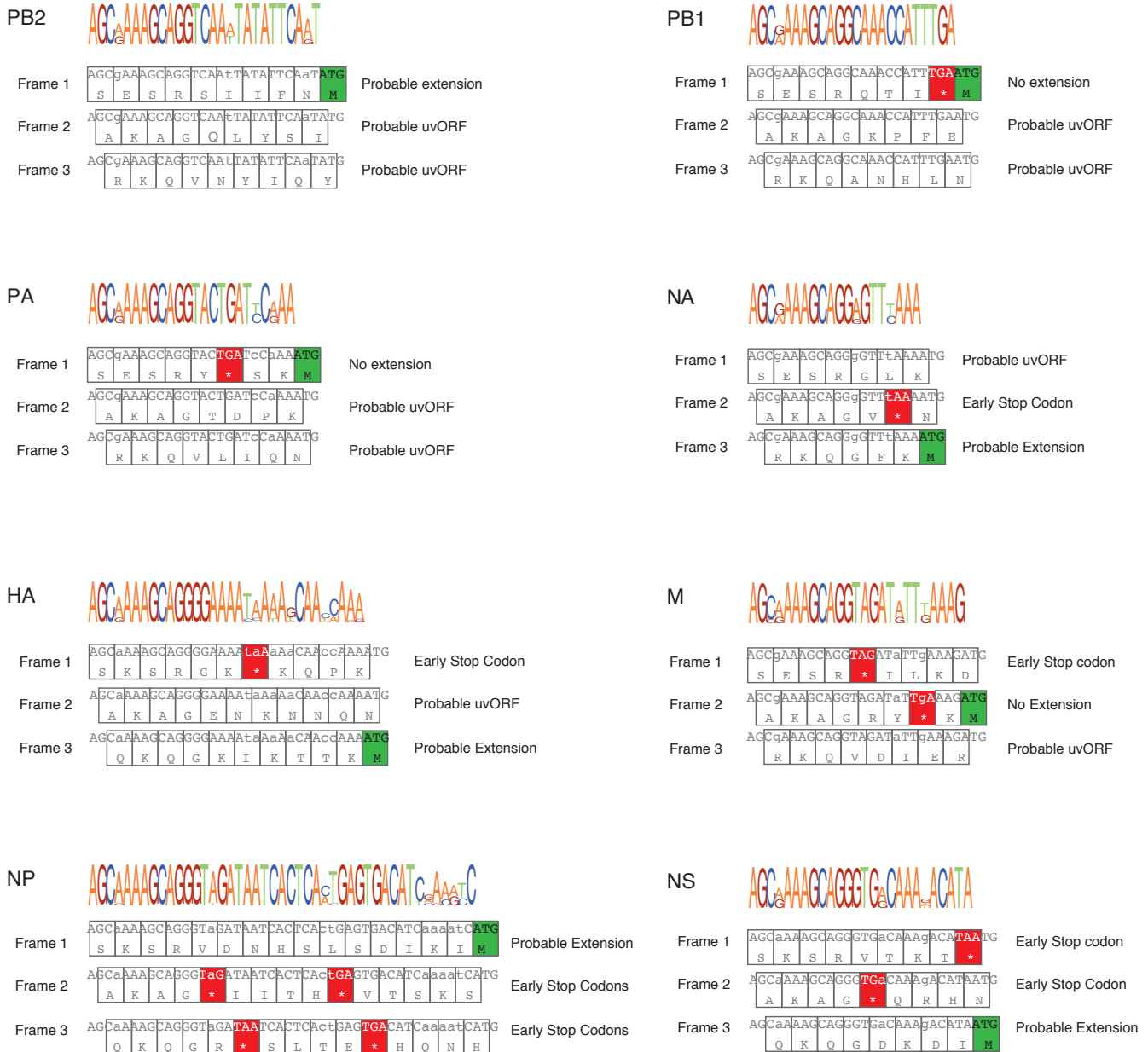
bioRxiv preprint doi: <https://doi.org/10.1101/597617>; this version posted April 8, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



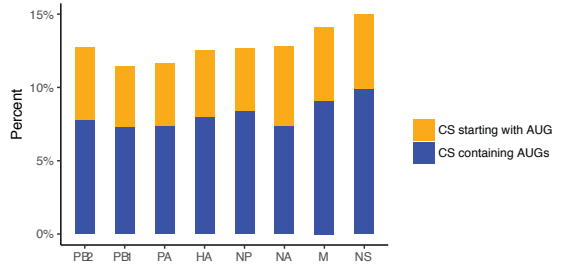
B



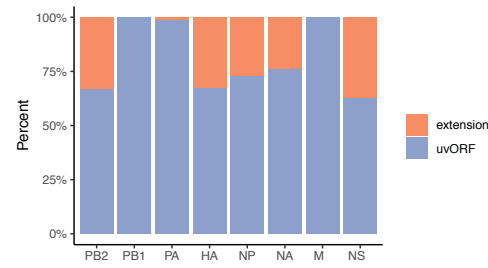
C



A



C



B

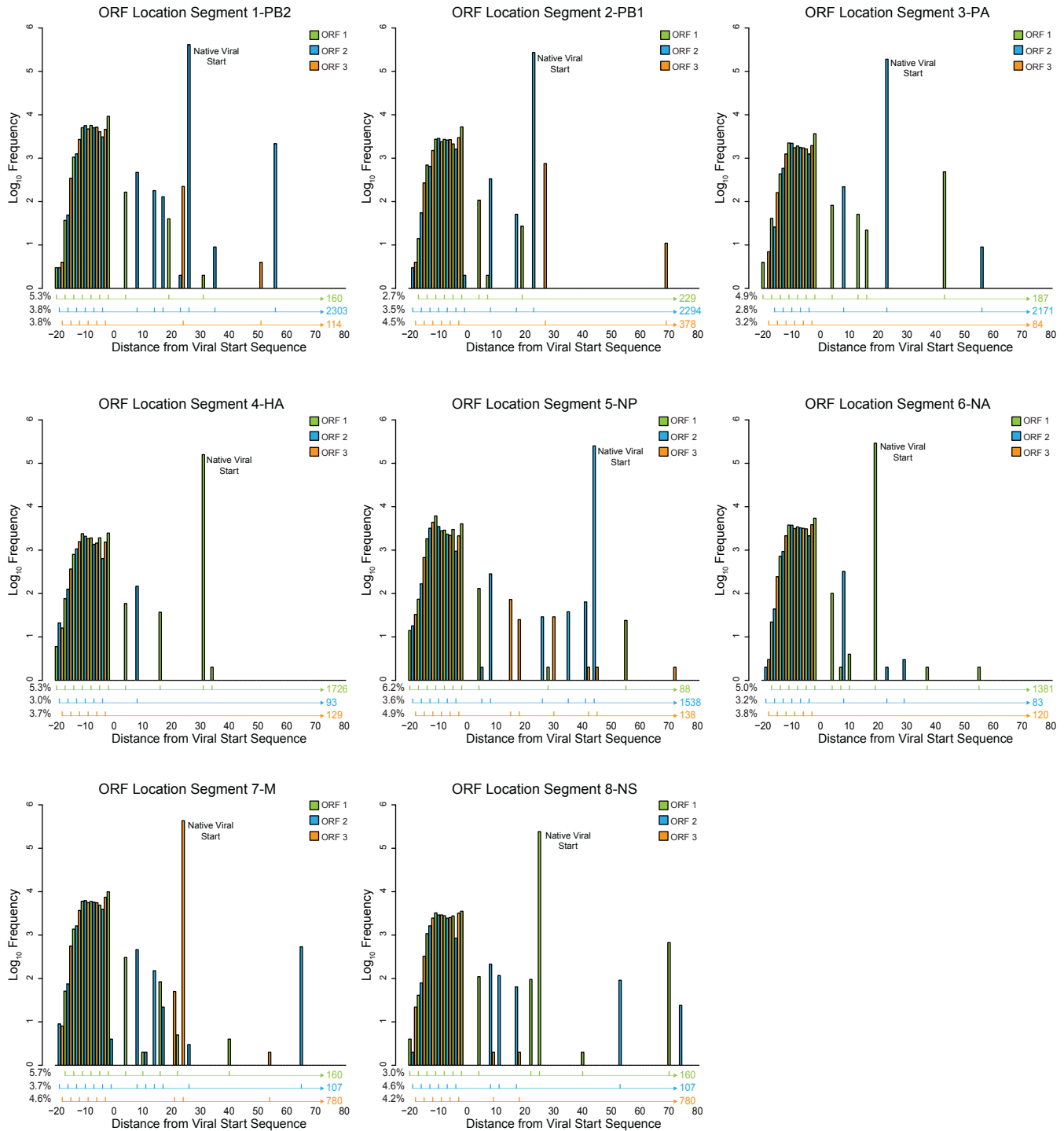


FIGURE 3

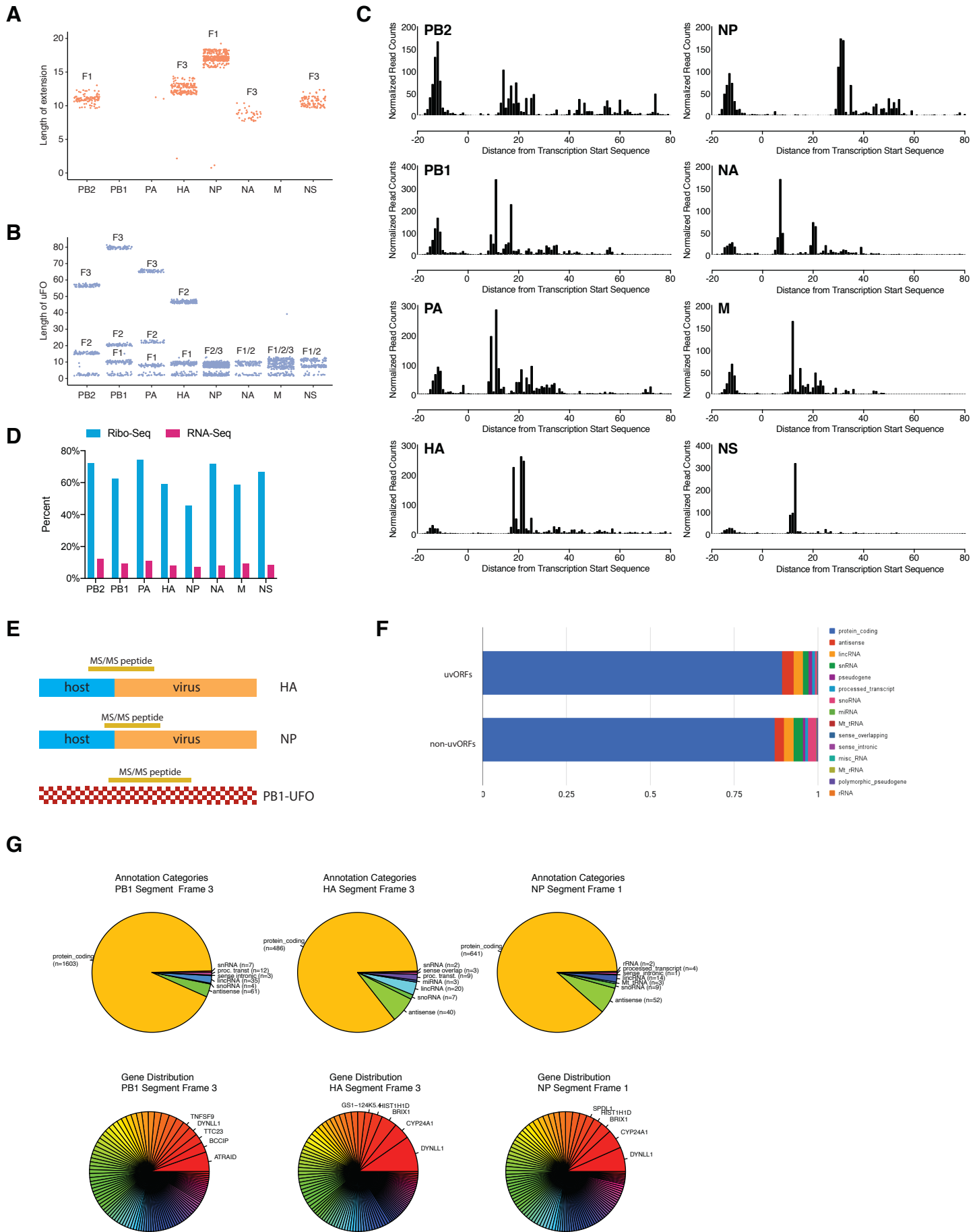


FIGURE 4

A

PB1-UFO translation

```

1 CGAAGCAGGCAAAAC CATTGAATGATGT CAATCCGACCTTACT TTCTTAAAGTGGC AGCACAAATGCTAT
1 R K Q A N H L N G C Q S D L T F L K S A S T K C Y
76 AGCACAACTTCCC TTACTGGAGACC TCCTTACAGCCATGG GACAGGAACAGGATA CACCATGGATACTGT
26 K H N F P L Y W R P S L Q P W D R N R I H H G Y C
151 CACAGGACACATCA GTAATCAGAAAGGG AAGATGGACACAAA CACCGAACTGGAGC ACCGCAACTCAACCC
51 Q Q D T S V L R K G K K N D N K H R N W S T A T Q P
226 GATTGA
76 D *
    
```

■ Peptide Found via Mass Spectrometry
* Stop Codon

B

WT
mRNA AGCGAAAGCAGGCAAAACCATTTGAATGGATGTCAATCCGACCTTACTTTTC
PB1 M D V N P T L L F
PB1-UFO R K Q A N H L N G C Q S D L T F

KO-3
mRNA AGCGAAAGCAGGCAAAACCATTTGAATGGATGTCAATCCGACCTTACTTTTC
PB1 M D V N P T L L F
PB1-UFO R K Q A N H .

CTRL-3
mRNA AGCGAAAGCAGGCAAAACCATTCGAATGGATGTCAATCCGACCTTACTTTTC
PB1 M D V N P T L L F
PB1-UFO R K Q A N H S N G C Q S D L T F

KO+9
mRNA AGCGAAAGCAGGCAAAACCATTTGAATGGATGTAAATCCGACCTTACTTTTC
PB1 M D V N P T L L F
PB1-UFO R K Q A N H L N G C

CTRL+9
mRNA AGCGAAAGCAGGCAAAACCATTTGAATGGATGTGAATCCGACCTTACTTTTC
PB1 M D V N P T L L F
PB1-UFO R K Q A N H L N G C E S D L T F

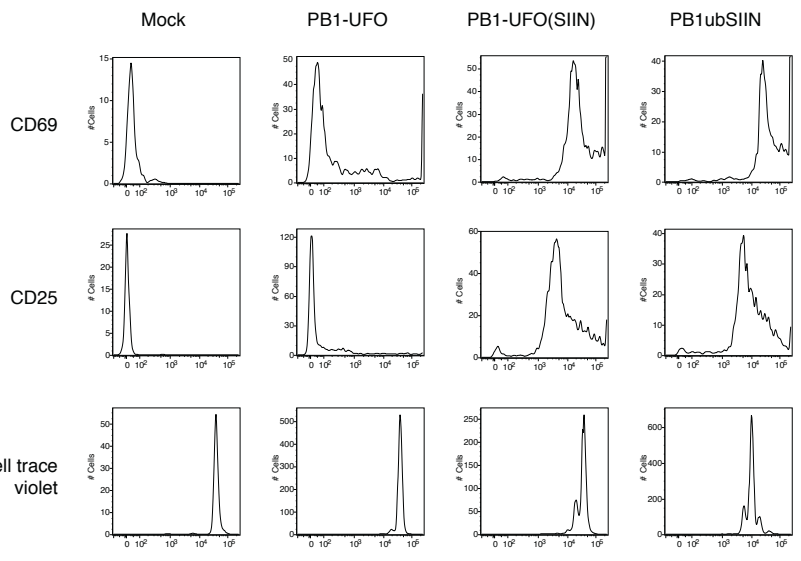
D

```

mRNA AGCGAAAGCAGGCAAAAGGATTTGA-----ATGGATGTCAATCCGACCTTA...
PB1 M D V N P T L...
PB1-UFO R K Q A N H L N-----G C Q S D L...

mRNA AGCGAAAGCAGGCAAAAGGATTTGAGTATAATCAACTTTGAAAACTGAATGGATGTCAATCCGACCTTA...
PB1 M D V N P T L...
PB1-UFO SIIN R K Q A N H L S I I N F E K L N G C Q S D L...
    
```

E



C

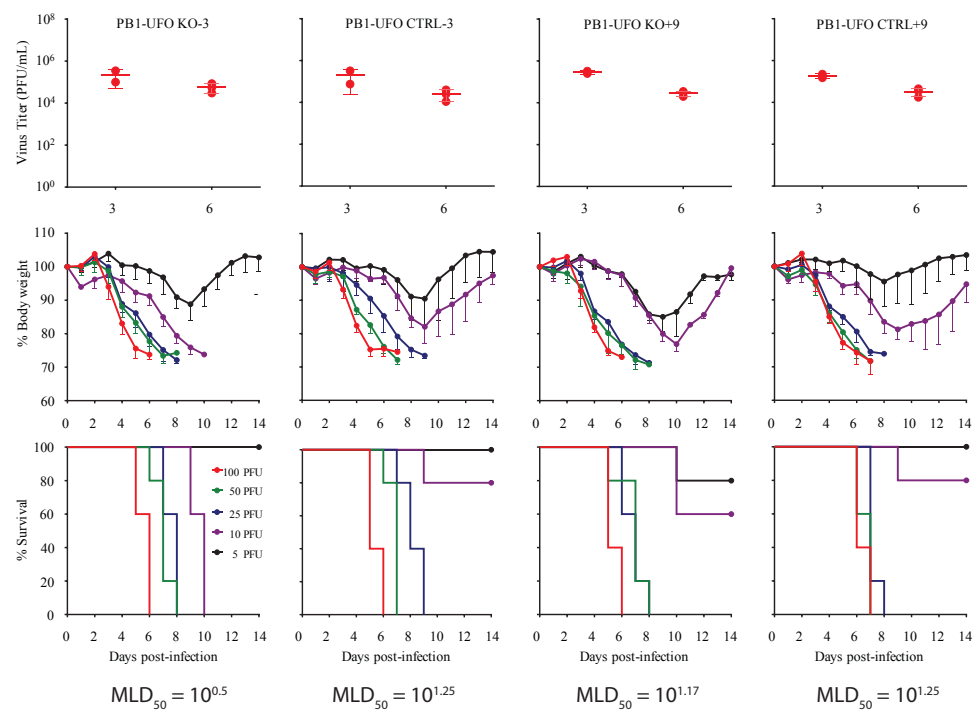
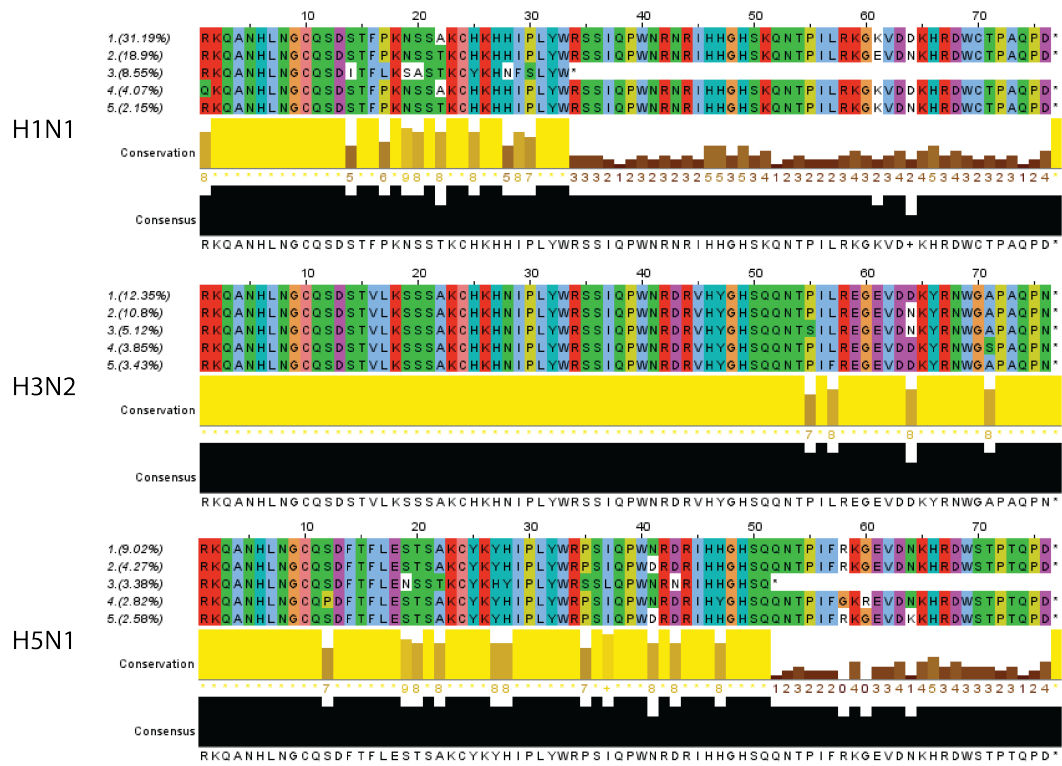
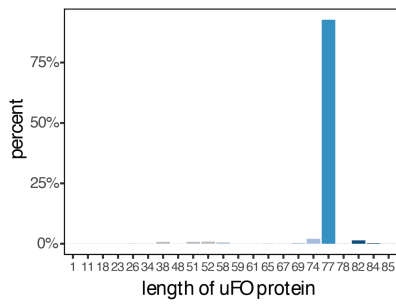


FIGURE 5

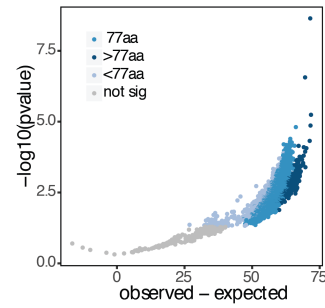
A



B



C



D

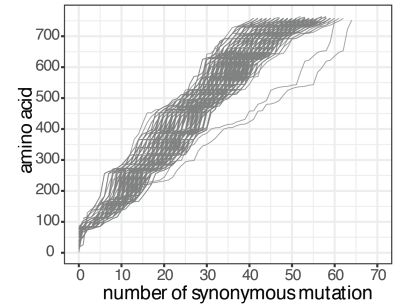


FIGURE 6

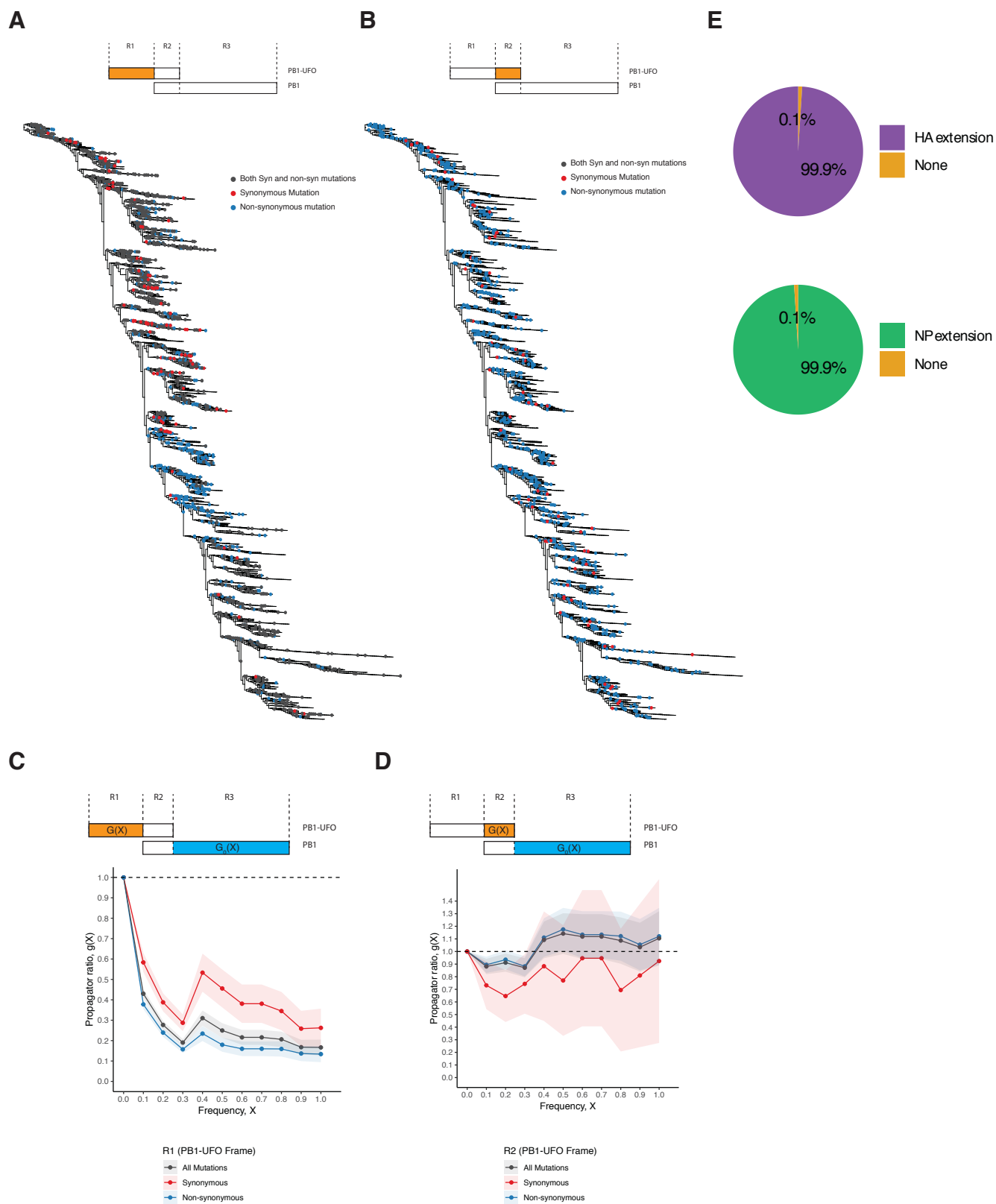


FIGURE S1

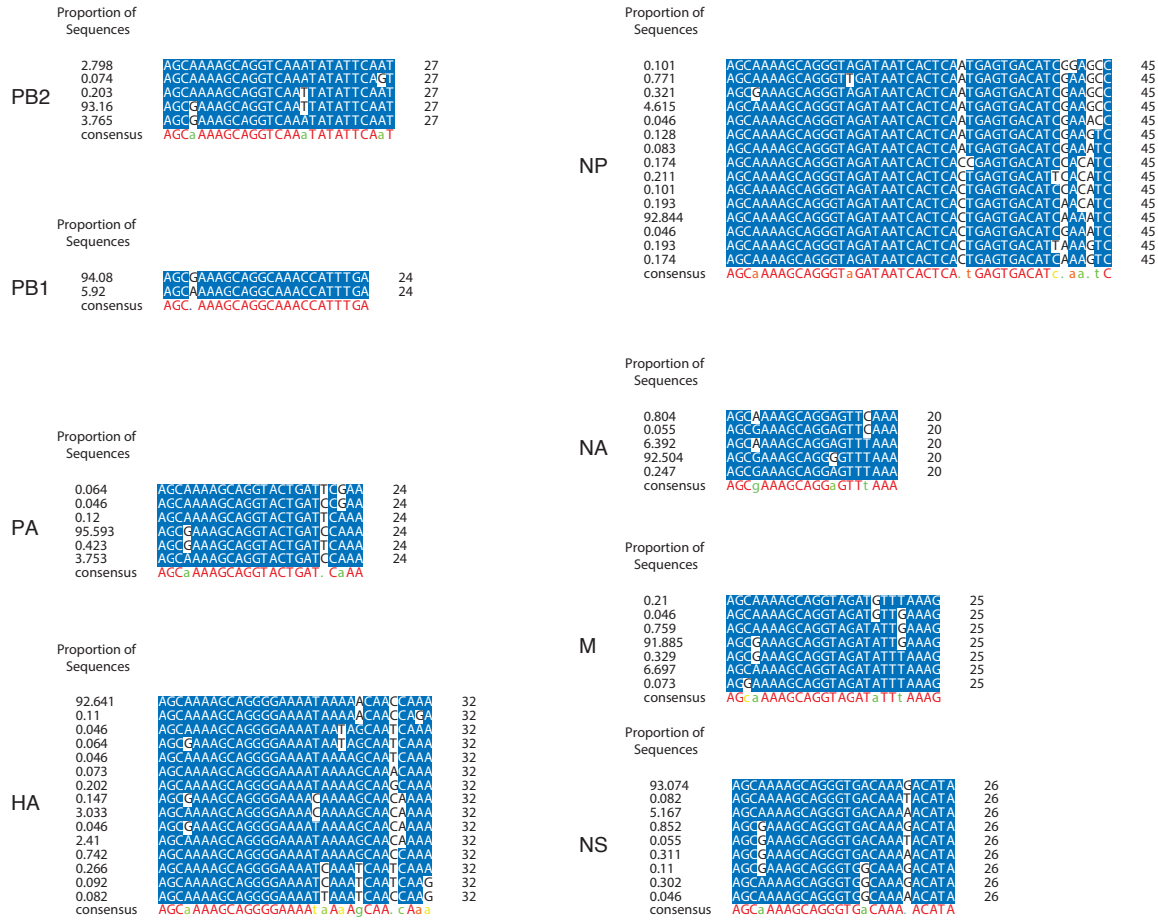
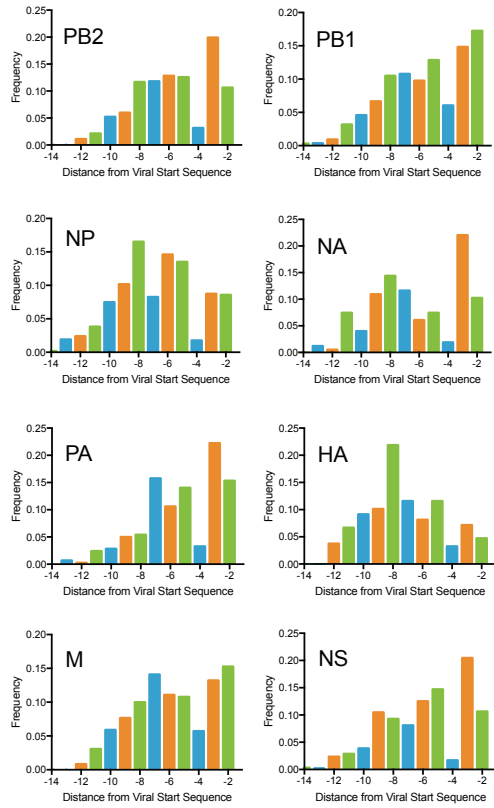
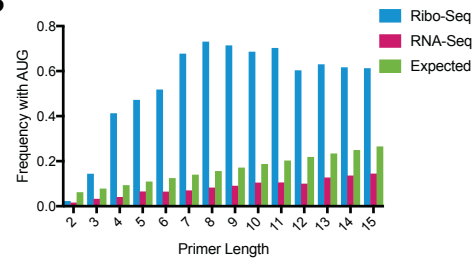


FIGURE S2

A



B



C

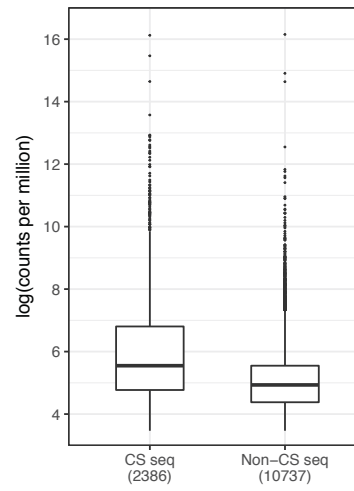


FIGURE S3

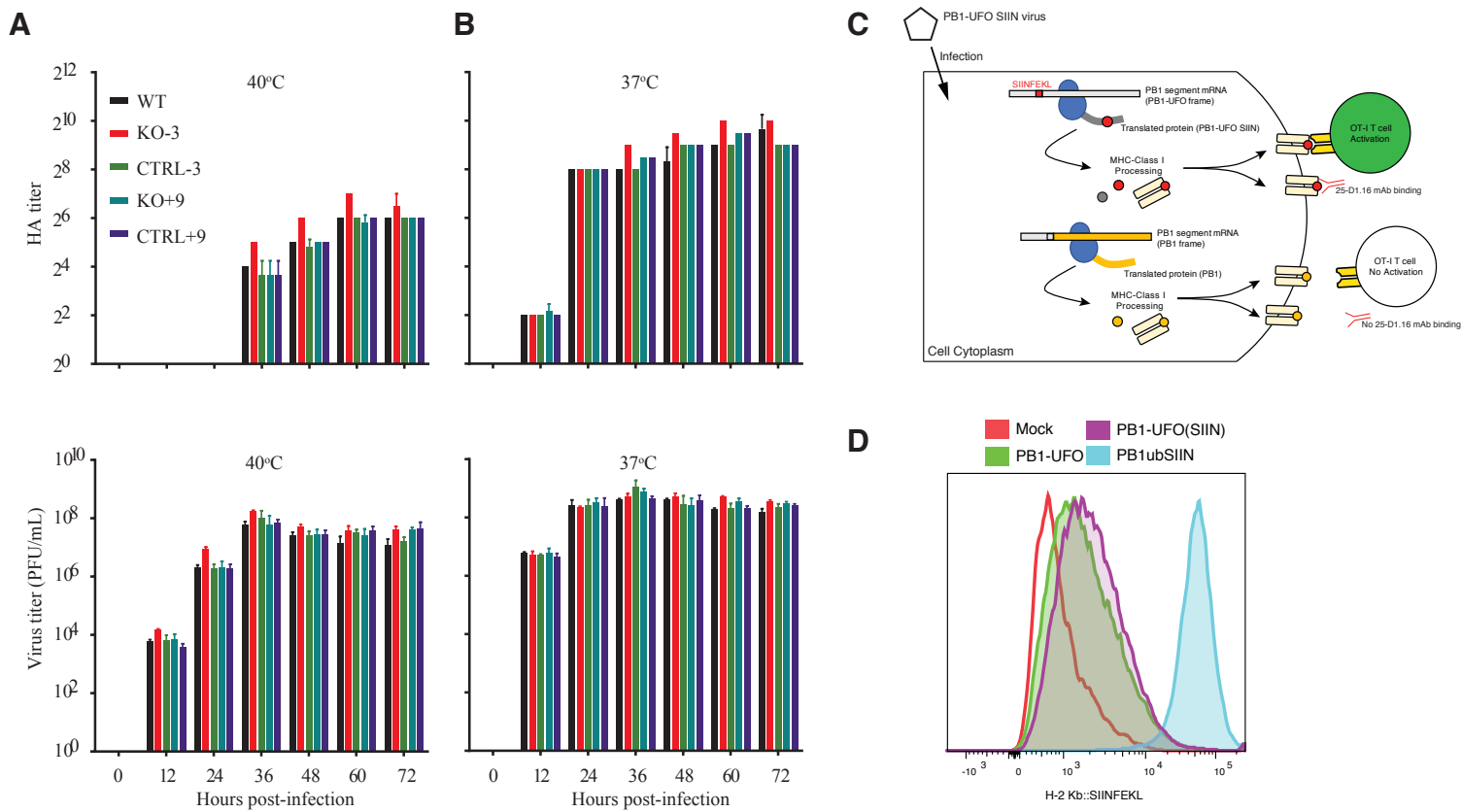


FIGURE S4

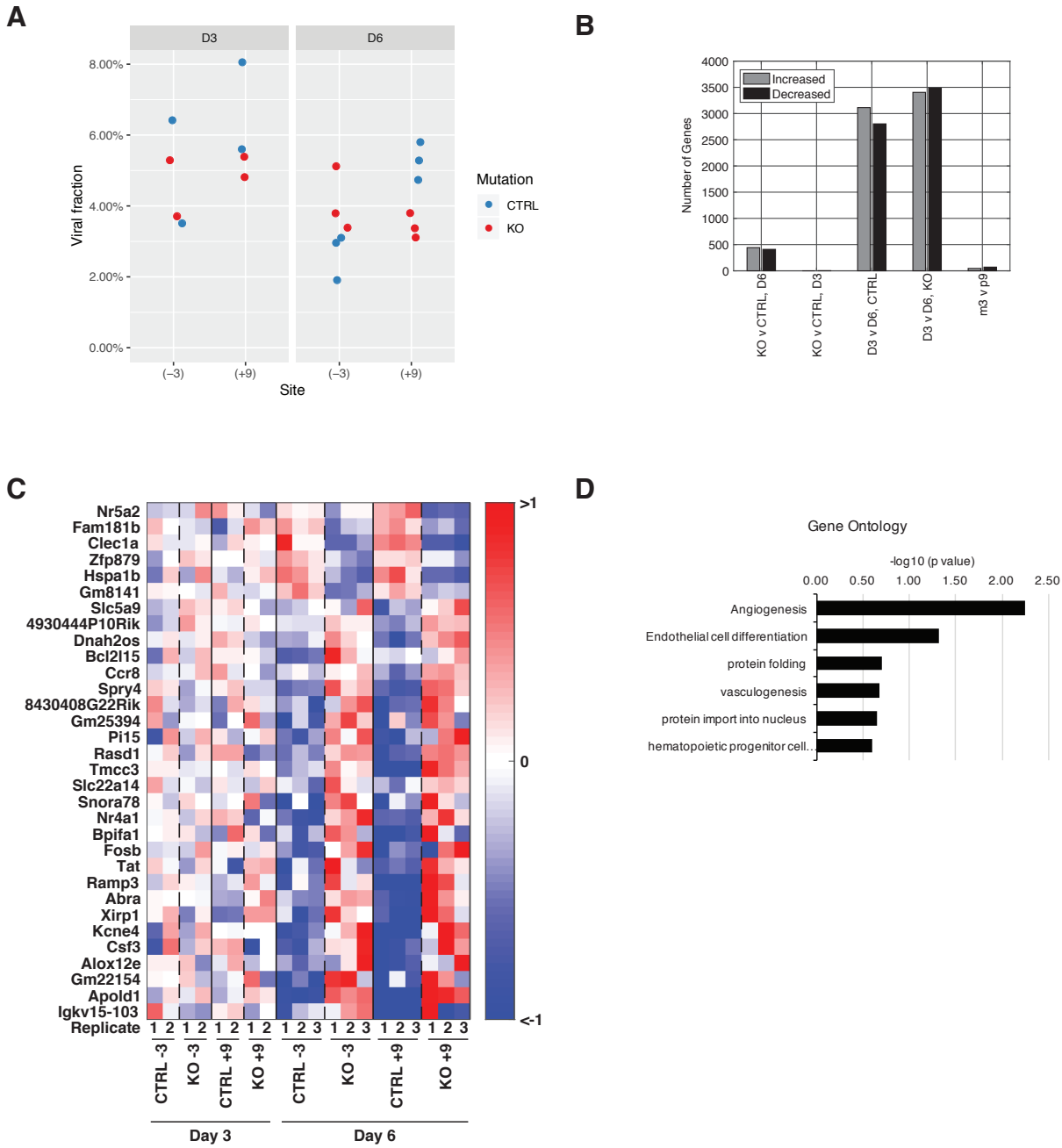
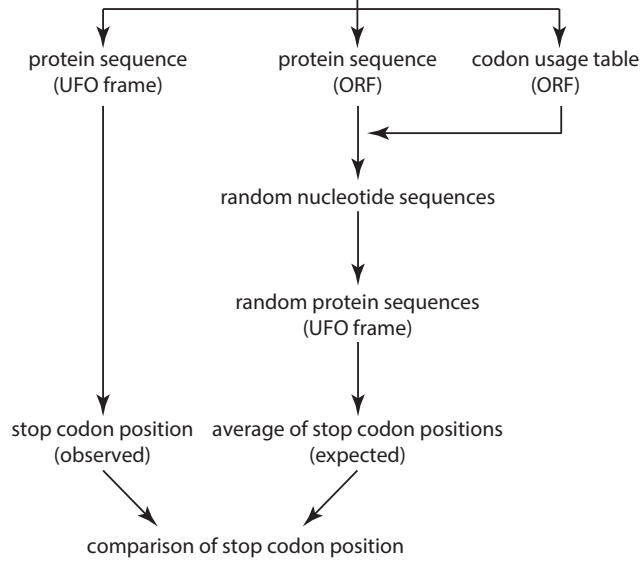


FIGURE S5

A bioRxiv preprint doi: <https://doi.org/10.1101/597617>; this version posted April 8, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



B

