



Fusaro, A. et al. (2016) Unexpected interfarm transmission dynamics during a highly pathogenic avian influenza epidemic. *Journal of Virology*, 90(14), pp. 6401-6411. (doi:10.1128/jvi.00538-16)

This is the author's final accepted version.

There may be differences between this version and the published version. You are advised to consult the publisher's version if you wish to cite from it.

<http://eprints.gla.ac.uk/121343/>

Deposited on: 04 August 2016

1 **Unexpected Inter-farm Transmission Dynamics During a Highly Pathogenic**
2 **Avian Influenza Epidemic**

3

4 Alice Fusaro^{a#}, Luca Tassoni^a, Adelaide Milani^a, Joseph Hughes^b, Annalisa Salviato^a, Pablo R.
5 Murcia^b, Paola Massi^c, Gianpiero Zamperin^a, Lebara Bonfanti^a, Stefano Marangon^a, Giovanni
6 Cattoli^{a*}, Isabella Monne^a

7

8 Istituto Zooprofilattico Sperimentale delle Venezie, Legnaro (PD), Italy^a; MRC-University of
9 Glasgow Center for Virus Research, Glasgow, United Kingdom^b; Istituto Zooprofilattico
10 Sperimentale della Lombardia e dell'Emilia Romagna, Brescia, Italy^c

11

12

13 Running Head: Influenza Transmission Dynamics

14

15 #Address correspondence to Alice Fusaro, afusaro@izsvenezie.it

16 * Present address: Joint FAO/IAEA Division of Nuclear Techniques in Food and Agriculture,
17 Department of Nuclear Sciences and Applications, International Atomic Energy Agency,
18 Seibersdorf, Austria

19

20 Abstract word counts: 240

21 Text word counts: 5865

22

23 **ABSTRACT**

24 Next Generation Sequencing technology is now being increasingly applied to study the within and
25 between host population dynamics of viruses. However, information on avian influenza virus
26 evolution and transmission during a naturally occurring epidemic is still limited. Here, we use deep
27 sequencing data obtained from clinical samples collected from five industrial holdings and a
28 backyard farm infected during the 2013 highly pathogenic avian influenza (HPAI) H7N7 epidemic
29 in Italy to unravel i) the epidemic virus population diversity, ii) the evolution of virus pathogenicity,
30 and iii) the pathways of viral transmission between different holdings and sheds. We show a high
31 level of genetic diversity of the HPAI H7N7 viruses within a single farm as a consequence of
32 separate bottlenecks and founder effects. In particular, we identified the co-circulation in the index
33 case of two viral strains showing a different insertion at the Hemagglutinin cleavage site, as well as
34 nine nucleotide differences at the consensus level and 92 minority variants. To assess inter-farm
35 transmission, we combined epidemiological and genetic data and identified the index case as the
36 major source of the virus, suggesting the spread of different viral haplotypes from the index farm to
37 the other industrial holdings, probably at different time points. Our results revealed inter-farm
38 transmission dynamics that the epidemiological data alone could not unravel and demonstrated that
39 delay in the disease detection and stamping out was the major cause of the emergence and the
40 spread of the HPAI strain.

41

42

43 **IMPORTANCE**

44 The within and between host evolutionary dynamics of a highly pathogenic avian influenza (HPAI)
45 strain during a naturally occurring epidemic is currently poorly understood. Here, we perform for
46 the first time an in-depth sequence analysis of all the samples collected during a HPAI epidemic
47 and demonstrate the importance to complement outbreak investigations with genetic data to
48 reconstruct the transmission dynamics of the viruses and to evaluate the within and between farms

49 genetic diversity of the viral population. We show that the evolutionary transition from the low
50 pathogenic to the highly pathogenic form occurred within the first infected flock where we
51 identified haplotypes with hemagglutinin cleavage site of different lengths. We also identify the
52 index case as the major source of virus, indicating that prompt application of depopulation measures
53 is essential to limit virus spread to other farms.

54

55 **INTRODUCTION**

56 Today, Next Generation Sequencing (NGS) techniques allow the investigation of viral
57 population dynamics at any level (from within host to the epidemiological scale) with high
58 resolution. In addition, NGS can be used to identify low frequency variants, which may be selected
59 for and transmitted to other hosts. Avian influenza viruses (AIVs) exist in the host as populations of
60 genetically related variants (1). The rate at which genetic diversity is generated within the host, the
61 competitive replication ability of each variant, and the occurrence of genetic drift and of bottleneck
62 events are some of the processes that drive virus evolution.

63 NGS has been applied to avian influenza virus i) to characterize the emergence of mutations in
64 the viral subpopulations associated to an increased virulence (2, 3) or to adaptation to new hosts, (4,
65 5) ii) to study genetic bottlenecks upon transmission events (6, 7); iii) to investigate the dynamics of
66 virus evolution during outbreaks in poultry (8); and iv) to identify co-infection with different
67 subtypes (9). However, application of high throughput sequencing for the exploration of avian
68 influenza virus evolution and transmission during a naturally occurring epidemic is still limited,
69 making the interpretation of genomic data collected from outbreaks far from straightforward.
70 Between August 13th and September 3rd of 2013, thirteen years after the last highly pathogenic
71 avian influenza (HPAI) outbreak, Italy experienced a new avian influenza epidemic caused by a
72 HPAI virus of the H7N7 subtype, which infected five industrial poultry holdings, four of which
73 belonged to a large vertically integrated layer company, and one backyard flock (10). Detailed
74 information on these outbreaks has been provided in a previous study (10). The epidemiological
75 investigation indicated that the contact between free-range hens and wild waterfowl in the first
76 affected holding may have favoured the introduction of a low pathogenic avian influenza (LPAI)
77 virus, which rapidly mutated into a HP form within the infected sheds (10) through the acquisition
78 of multiple basic amino acids at the hemagglutinin (HA) cleavage site, which is considered as being
79 the major molecular determinant of an HPAI virus (11).

80 Here we used NGS to unravel the virus population diversity and the evolution of virus
81 pathogenicity within the affected poultry farms. We also determined the transmission pathways of
82 the H7N7 virus between different holdings and sheds during the course of the epidemic by
83 combining deep sequencing and epidemiological data.

84

85

86 **MATERIALS AND METHODS**

87

88 **Viruses**

89 Fourteen positive clinical samples (organs and swabs) were collected between August 13th and
90 September 3rd 2013 from each infected shed of the five industrial farms and a backyard flock,
91 counting for all the sheds infected during the epidemic (10). Epidemiological information, including
92 collection date, sample type (swabs, organs), farm and shed of origin, number of birds present in
93 each farm at the time of the forfeiture and depopulation date, is available in Table 1.

94 The viral RNA copy numbers (Table 1) were determined for each sample using a quantitative real-
95 time RT-PCR assay with a standard curve targeting the M gene of influenza A, using the published
96 probes and primers from Spackman et al. (12).

97

98 **Generation of viral sequence data**

99 Total RNA was purified from the 14 infected clinical samples using the Nucleospin RNA kit
100 (Macherey–Nagel, Duren, Germany). Complete influenza A virus genomes were amplified with the
101 SuperScript III One-Step RT-PCR system with Platinum Taq High Fidelity (Invitrogen, Carlsbad,
102 CA) using one pair of primers complementary to the conserved elements of the influenza A virus
103 promoter as described in (13). PCR products were visualized on a 0.7% agarose gel. Sequencing
104 libraries were obtained using Nextera DNA XT Sample preparation kit (Illumina) following the
105 manufacturer’s instructions and quantified using the Qubit dsDNA High Sensitivity kit (Invitrogen,

106 USA). The average fragment length was determined using the Agilent High Sensitivity Bioanalyzer
107 Kit. Finally the indexed libraries were pooled in equimolar concentrations and sequenced in
108 multiplex for 250 bp paired-end on Illumina MiSeq, according to the manufacturer's instructions.

109

110 **Quality trimming, assembly and SNP detection**

111 Illumina MiSeq reads were inspected using FASTQC to assess the quality of data. Fastq files were
112 cleaned with PRINSEQ and Trim Galore to remove low quality bases at the 5' and 3' end of each
113 read and to exclude reads with a Phred quality score below 30 and shorter than 80 nucleotides. The
114 filtered, trimmed reads were aligned to the eight gene segments of A/chicken/Italy/13VIR4727-
115 11/2013, for which the consensus genome were previously obtained using Sanger method (data not
116 shown), using BWA-MEM v.0.7.5a (<http://arxiv.org/abs/1303.3997v2>). The BAM alignment files
117 were parsed using the diversiTools program (<http://josephhughes.github.io/btctools/>) to determine
118 the average base-calling error probability and to identify the frequency of polymorphisms at each
119 site relative to the reference used for the alignment. In order to minimize artefacts introduced
120 through RT-PCR and sequencing errors, for all the analysis conducted throughout this study we
121 considered only polymorphisms with a frequency above 2% identified in positions with a minimum
122 coverage of 500. This choice was based on the comparison of data obtained from two technical
123 replicates of three samples (4541-8, 4541-9, 4541-34), sequenced on two different Illumina
124 sequencing machines (MiSeq), starting from two separate libraries obtained from the same
125 extracted RNA. This threshold should guarantee the exclusion of 99.6% of the errors from our deep
126 sequencing data (Fig. 1). For each replicate, only the aligned genome with the highest coverage was
127 used in the following analyses.

128 For each gene, we calculated the number of synonymous and non-synonymous polymorphisms
129 present either at a consensus level or as subpopulations and normalized to the number of
130 synonymous and non-synonymous sites in the coding regions. Significant differences between the

131 frequencies of the two types of mutation in the different genes were calculated using a two-way
132 ANOVA. A value of $P < 0.05$ was considered significant.

133

134 **Genetic distance, entropy and transmission tree**

135 We computed the genetic distance between the complete genome of all pairs of individuals (S1 and
136 S2) using the following formula: $d = \frac{1}{N} \sum_{i=1}^N (|f_{A_{iS1}} - f_{A_{iS2}}| + |f_{C_{iS1}} - f_{C_{iS2}}| + |f_{T_{iS1}} - f_{T_{iS2}}| +$
137 $|f_{G_{iS1}} - f_{G_{iS2}}|)^2$, where $f_{A_{iS1}}$, $f_{C_{iS1}}$, $f_{T_{iS1}}$, $f_{G_{iS1}}$ are the frequencies of nucleotide A, C, T and G at
138 position i in the two samples and N is the length of the sequence. This matrix was used to compute
139 a neighbour-joining phylogenetic tree using the web server T-REX (14). In addition, we combined
140 the distance matrix and the collection dates to reconstruct the transmission tree of the H7N7 during
141 the Italian outbreak, using SeqTrack (15), a graph-based approach particularly suitable to infer
142 maximum parsimony genealogies of viruses in densely sampled disease outbreak. The *adegenet*
143 (16) and *igraph* packages (17) for the R software were used to perform the analysis and to draw the
144 network.

145 To measure the complexity of the viral populations within a sample, we calculated the Shannon
146 entropy of each sample and each gene using the following equation:

$$E = -\frac{1}{N} \sum_{i=1}^N (f_{iA} \ln f_{iA} + f_{iG} \ln f_{iG} + f_{iT} \ln f_{iT} + f_{iC} \ln f_{iC})$$

147 where f_i is the frequency of the nucleotide A, T, G or C at position i and N is the total length of the
148 gene segment (average entropy per gene) or of the genome (average entropy per sample). Only
149 nucleotides with a frequency above the 2% threshold identified in positions with a minimum
150 coverage of 500 were included in this calculation. We used one-way ANOVA to determine
151 significant differences between the entropies of each gene for each sample. A value of $P < 0.05$ was
152 considered significant.

153

154 **Phylogenetic analyses**

155 Consensus sequences of the complete genome of the 14 samples were aligned using MAFFT v. 7
156 (18) and compared with the most related sequences available in GenBank and in GISAID (accessed
157 on May 2015). In addition, representative H7 viruses circulating in wild and domestic birds in
158 Europe and H7 viruses responsible of important epidemics were included in the alignment.
159 Maximum likelihood (ML) phylogenetic trees were obtained for each gene segment using the best-
160 fit general time reversible (GTR) model of nucleotide substitution with gamma-distributed rate
161 variation among sites (with four rate categories, Γ 4) available in RAxML-MPI v.8.1.7 (19). To
162 assess the robustness of individual nodes of the phylogeny, one hundred bootstrap replicates were
163 performed. Phylogenetic trees were visualized with the program FigTree v1.4
164 (<http://tree.bio.ed.ac.uk/software/figtree/>).
165 The eight gene segments of the influenza virus genome were manually concatenated and the
166 alignment was used to construct a phylogenetic network using the Median Joining method
167 implemented in the program NETWORK 4.5 (<http://www.fluxus-engineering.com>) (20). This
168 method uses a parsimony approach to reconstruct the relationships between highly similar
169 sequences, and allows the creation of “median vectors”, which represents unsampled sequences,
170 that are used to connect the existing genotypes in the most parsimonious way. The parameter
171 *epsilon* was set to 10 and the transition to transversion ratio to 3:1. A bootstrap resampling process
172 (1000 replicates) using a distance-based method (NeighborNet) implemented in SplitsTree4
173 v.4.14.2 (21) was used to assess the robustness of the network edges.

174

175 **Nucleotide sequence accession numbers**

176 MiSeq sequences were submitted to the NCBI Sequence Read Archive (SRA,
177 <http://www.ncbi.nlm.nih.gov/Traces/sra/>) under accession numbers SRR3036850, SRR3036852,
178 SRR3036854, SRR3036856, SRR3036860, SRR3036864, SRR3036910, SRR3036911,
179 SRR3036914, SRR3036916, SRR3036917, SRR3036919, SRR3036920, SRR3036945. Consensus

180 sequences of the 14 H7N7 viruses were submitted to GISAID under accession numbers EPI677984
181 to EPI678095.

182

183

184 **RESULTS**

185 **Phylogenetic analysis of consensus sequences**

186 To investigate influenza virus variation during the HPAI H7N7 epidemic, we sequenced the eight
187 genomic segments for all the clinical samples received from each infected farm. The highest
188 number of positive samples (8) was submitted from the three infected sheds (shed 2, 4 and 5) of the
189 index case, while only one sample per infected shed was received from the remaining five outbreak
190 sites, for a total of one or two samples per farm. Farms are labelled from 1 to 6, according to the
191 collection date of the samples. Details of location, date of sample collection, farm characteristics,
192 sample type and mean depth of coverage are reported in Table 1.

193

194 Our maximum likelihood phylogenetic analyses of the consensus sequences show that the fourteen
195 HPAI H7N7 viruses form a distinct genetic group, defined by high bootstrap values (>96%) and
196 long branches in all the eight phylogenies, suggesting the occurrence of a single viral introduction
197 (Fig 2). This group includes also the sequences of the complete genome available for one of the
198 three poultry workers involved in the depopulation, who developed conjunctivitis due to HPAI
199 H7N7 infection, suggesting a direct transmission of the virus from poultry to human (22). In the HA
200 and NA phylogenetic trees, the Italian H7N7 HPAI cluster with H7 viruses collected in Europe
201 between 2009 and 2014. In particular, the HA gene segment of the Italian samples show the highest
202 similarity (99.1-99.3%) with an LPAI H7N7 virus collected from a wild bird in Italy in 2014, for
203 which only the HA sequence is available (Fig 2), while the NA gene segment display the highest
204 identity (99-99.1%) with an H7N7 virus collected from chicken in the Netherlands (phylogenetic
205 tree is available upon request). In the phylogenies of the internal gene segments the Italian samples

206 group with viruses of different subtypes circulating mainly among wild birds in Eurasian countries
207 (phylogenetic trees are available upon request).

208

209 **High genetic variability of the first infected flock**

210 Surprisingly, molecular analysis of the eight viruses collected from the index case shows the co-
211 circulation of two highly pathogenic strains with a different insertion at the HA cleavage site
212 compared to a H7 LP virus. Specifically, sequences of the two viruses from shed 5 (4541-9 and
213 4541-34) show an insertion of 6 nucleotides, while the remaining samples identified in sheds 2 and
214 4 possess a longer cleavage site with a nine nucleotide insertion (Fig. 3).

215 To better understand the evolution of the pathogenicity of the H7N7 viruses within the first infected
216 flock, we focused our analysis on the deep sequencing data of the HA cleavage site. The sequencing
217 coverage in this genetic region ranges from 4445 for the sample 4541-7 to 23511 for the sample
218 4541-34. We did not identify any reads showing the cleavage site typical of a LPAI strain. 99.9% of
219 the reads of the two samples from shed 5 (named for clarity V+6) possess a cleavage site with an
220 insertion of six nucleotides, with only a few reads containing an insertion of three, five and nine
221 nucleotides (Table 2). 99.7% to 99.9% of the reads of viruses from shed 2 (named V+9) have an
222 insertion of nine nucleotides, with only a few minority variants showing an insertion of six, seven
223 or eight nucleotides (Table 2). On the other hand, in one of the samples from shed 4 (4541-33) we
224 identified a mixed population with both type of cleavage sites displaying an insertion of nine
225 (95.7%) and six (4.1%) nucleotides.

226 Similarly to the samples from shed 2, the majority (from 99.9% to 100%) of the viral population of
227 the subsequent outbreaks possesses the longer cleavage site, suggesting that this variant (V+9) may
228 have a higher fitness advantage (Fig. 3).

229 Besides the cleavage site, the samples V+9 collected from shed 2 and 4 of the first infected farm
230 can be distinguished from the two samples from shed 5 by nine nucleotide signatures (HA G471A,
231 PB2 A347G and T1891G, PA A347G and T1891G, NP G219A and C316A, NS1 G353A and

232 G378A; Fig. 3), which resulted in three amino acid changes (PA Q116R, PA C631G, NS1 R118K).
233 These signatures are maintained in all samples identified in the subsequent outbreaks, suggesting
234 that only viruses from sheds 2 and 4 of the index case were transmitted to the other five farms (Fig
235 3). In addition, we identified one non-synonymous mutation at position 130 of the M2 gene,
236 responsible of the amino acid substitution D44N, which is shared between the V+6 viruses and the
237 samples 4527-11 from shed 2 of farm 1, 4603 from farm 2, 4678 from farm 3 and 5091 from farm 5
238 (Fig 3). However, whether this mutation emerged by chance in the 4 viruses or arose in the shed 2
239 virus of the index case and was then transmitted to the other outbreaks or was acquired by the V+9
240 samples through a reassortment event cannot be assessed.

241 To determine whether the shed 5 viruses (V+6) were the progenitors of the variant V+9, we
242 examined the presence of the nine signature mutations (Fig 3) as minority variants in the analysed
243 samples. None of the mutations typical of the V+9 viruses were already present in shed 5 viruses
244 (V+6) with a frequency higher than 2% (the frequency threshold used in this study, see the
245 Materials and Methods section for details). Similarly, none of the mutations characteristic of V+6
246 (Fig 3) was identified in the subpopulations of the V+9 samples, except for the virus from shed 4 of
247 the index case (4541-33), which, besides the shorter cleavage site, possessed subpopulations
248 containing all the mutations distinctive of V+6 variant, with a frequency ranging from 3% to 9%,
249 confirming the presence of a mixed population (V+6 and V+9).

250

251 **Genetic diversity of H7N7 viruses**

252 Overall, we observed mutations at 185 sites (excluding the HA cleavage site) distributed among the
253 eight gene segments, of which 111 are non-synonymous and 74 synonymous. Specifically, a total of
254 35 consensus-level nucleotide substitutions are recovered along the entire genome, defining 11
255 different genomes (named from A to K in Fig 3), five of which identified within the first infected
256 farm (A to E). The PB2 gene, with a total of ten nucleotide variants (8 synonymous and 2 non-
257 synonymous), is the segment showing the highest number of mutations at the consensus level. The

258 nucleotide distance among the fourteen viruses ranged from 0-0.1% for the PA, HA and NA genes
259 to 0-0.2% for the PB2, PB1, NP genes and 0-0.4% for the NS gene. Notably, 13 out of 35 mutations
260 distributed along twelve proteins (HA, NA, PB2, PB1, PB1-F2, PA, PA-X, NP, M1, M2, NS1 and
261 NS2) are non-synonymous, with the PA protein showing the highest number of amino acid
262 variations (4) (Fig 3).

263 Besides these consensus-level variant sites, our deep sequencing analysis identifies 209 minority
264 variants in 151 sites (97 non-synonymous and 54 synonymous) with a frequency ranging from 2%
265 to 49.8% (Fig 4). The virus collected from shed 4 of the index case (4541-33), which displayed a
266 mixed population of V+6 and V+9, and the sample 4603 collected from farm 2, comprise the
267 highest number of minority variants (respectively, 40 and 41). On the contrary, we did not detect
268 any subpopulations in the samples 4541-34 and 4541-9. No correlation between the number of
269 variants and the type of samples used for the analysis (pool, organs or swab) was observed (Pearson
270 test, p-value=0.254; r=0.33).

271 We measured the complexity of the viral population of each sample using Shannon entropy
272 (represented by the size of the circles in Fig 5). In the first infected flock, entropy measures
273 fluctuate considerably: the lowest values are observed for the two viruses from the shed 5 (V+6)
274 (Wilcoxon rank-sum test, p-values range from 9.37×10^{-14} to 4.7×10^{-3}), suggesting that these
275 samples (4541-9 and 4541-34) had recently experienced a narrow bottleneck and had not recovered
276 from the loss of complexity. Conversely, viruses from shed 2 show intermediate values of entropy,
277 while samples 4541-33 from shed 4 of the index case, 4603 from farm 2 and 4678 from farm 3
278 displayed entropy levels significantly higher compared to the other samples (Wilcoxon rank-sum
279 test, p-values range from 9.37×10^{-14} to 1.49×10^{-3}), consistent with the high genetic diversity
280 observed across their genomes. There is no significant Pearson correlation between within-host
281 virus diversity and viral RNA content (p-value=0.487; r=-0.2; the number of RNA copies are
282 reported in Table 1). Thus, the significantly different entropies between the analysed samples may
283 be simply a bias associated to the time elapsed between infection and sampling, or alternatively

284 they may be due to the occurrence of random or selective bottleneck events of different intensity
285 within separate sheds or farms.
286 To evaluate which could be the major force - selective bottleneck or random founder effect -
287 driving the virus evolution, we compared the relative diversity changes on a gene-by-gene basis
288 (Table S1). We would not expect selective bottlenecks to affect all the genes in the same way and
289 thus the entropy and the number of non-synonymous mutations found in the genes should vary. To
290 this aim we calculated the mean entropy and the number of synonymous and non-synonymous
291 normalised mutations separately for all the genes of each sample (Table S1). There was no
292 significant difference in the entropy between the genes (one-way ANOVA, $p=0.196$) or between the
293 number of non-synonymous and synonymous mutations (two-way ANOVA, $p=0.249$), suggesting
294 that the reduction in diversity observed in the analysed samples was probably due to founder effects
295 rather than selective bottlenecks.

296

297 **Minority variants transmitted between sheds and farms**

298 Focusing our analysis of the first infected flock, we observed that only a few mutations were shared
299 at a shed and farm level, while the majority of the minor changes were unique to individual
300 samples. Specifically, at the shed level we detected 44 minority changes in viruses from shed 2, of
301 which 22 are found in individual samples and not shared with others, and 48 in viruses from shed 4,
302 of which 37 are identified in single samples. Similarly, at the farm scale we counted 92 mutations,
303 of which 12 are shared between 2-4 samples, while 59 minority variants were identified in single
304 individuals (Fig 4).

305 Interestingly, five of these variants, shared between viruses of the first infected farm, are fixed in
306 the viral population of at least one sample (variants highlighted with black arrows in Fig 4) and
307 three of them were also transmitted or independently acquired by viruses collected from the other
308 premises. Four of these are non-synonymous mutations fixed in the viral population, which cause
309 changes at the protein level (NS M119T, M2 D44N, PA V100I, PB2 K574R) (Fig 3).

310 We detected only seven minority variants (HA 1351A, M 942G, M 955G, PA 1251 G, PA 1748A,
311 PB2 981G and NA 390A) shared between two or three farms. Interestingly, five of them result in
312 amino acid mutations (HA D451N, M2 D85G, PA R583Q, PB2 G327G, NA M130I). These non-
313 synonymous mutations may be advantageous variants, associated with changes in viral fitness or
314 due to adaptive evolution of the virus, or alternatively they may be neutral or deleterious
315 polymorphisms, which occurred because of random genetic drift or hitchhiking.

316

317 **Transmission dynamics of the H7N7 virus**

318 To assess the inter-farm transmission, a Median Joining phylogenetic network was inferred using
319 the concatenated consensus sequences of the eight gene segments of the 14 analysed viruses (Fig 6).
320 Within the first infected farm we identified five sequence genotypes (grey circles): one within shed
321 5, two within shed 2, and two in shed 4. Viruses from sheds 2 and 4 appear to be at the origin of the
322 infection to the other farms, although one or two median vectors (red circles), which represent the
323 lost ancestral sequences, separate them from viruses of the other holdings, except for the sample
324 5051-3 from farm 6, which appears to be a direct descendant of shed 2 viruses (bootstrap value
325 85.5). Sequences from farms 2 to 6 grouped within two main clusters which shared a common
326 ancestor (c1 and c2): c1 includes viruses collected from farms 2 (4603), 3 (4778) and 5 (5091),
327 while c2 contains virus sequences from farms 4 (4774) and 6 (5051-1). Sequences within these two
328 clusters are separated by 6 to 10 nucleotide differences, whereas 9-13 differences are observed
329 between viruses of the two clusters. Therefore, the high number of mutations and median vectors
330 identified between the analysed samples and the low number of viruses available for the analysis
331 makes the relationship between sequences hard to determine and we cannot exclude that sampling
332 bias may have affected the results. Our deep sequencing data may contribute to better understand
333 this relationship. To this aim, first we inferred a neighbour-joining phylogenetic tree based on the
334 distance matrix calculated from our NGS data, which confirmed the clustering identified by our
335 network analysis (Fig 4). Then we used the distance matrix and the collection dates to reconstruct a

336 transmission tree using the graph-based algorithm SeqTrack. This approach, which considers the
337 sampled viruses as a fraction of the genealogy, is particularly suitable to infer the transmission
338 pathway during disease outbreaks, where one strain can be the ancestor of another strain (Fig 5).
339 Despite 21 days passing from the first to the last outbreak, the inferred genealogy suggests that all
340 but one of the outbreaks descend directly from shed 2 (V+9) of the index case. The only exception
341 is represented by the virus (5091) collected from the backyard farm on September 2 (farm 5), which
342 appears to have been infected directly by farm 3.

343 However, based on the number of shared mutations between the analysed sequences, we may
344 speculate further scenarios. For example, sample 4603 from farm 2 shared two fixed mutations with
345 samples 4678 (farm 3) and 5091 (farm 5) (group c2 of the network analysis), thus a transmission
346 event from farm 2 to farm 3 cannot be excluded. Similarly, viruses 4774 and 5051-1 share 3 unique
347 minority variants and 1 unique fixed mutation, making a transmission event between these two
348 farms highly plausible. In addition, samples 4774 and 5051-1 share 1 fixed mutations and 3
349 minority variants (group c1 of the network analysis), and in turn they share 2 fixed mutations with
350 the sample 5051-3. Although viruses 5051-1 and 5051-3 were collected from two different sheds of
351 farm 6, we observed a relatively high nucleotide distance between them. Specifically, they show 7
352 and 14 nucleotide differences at the population and subpopulation level, respectively, although all
353 the consensus level mutations were present as minority variants in the other sample (Fig 4). Thus,
354 the occurrence of two separate introductions in farm 6 from the index case and/or farm 4 cannot be
355 excluded.

356 Overall, these analyses indicate that shed 2 of the index case is the major source of the virus. An
357 early strain (c1) appears to have spread from the first infected flock to farm 2 (19 August) and 3 (21
358 August) and then from farm 3 to the backyard farm 5 (2 September). Since farms 2 and 3 belong to
359 different companies (circle outlines in Fig 5) and are located 50 Km apart (map in Fig 6), it is more
360 plausible that viruses with similar genetic characteristics were transmitted from the index case to
361 both holdings. A later spread with a slightly different strain (c2) may have occurred from the first

362 infected flock to farm 4 (27 August) and 6 (3 September). These two farms are located in the same
363 area, with a distance of 3 Km, and belong to the same layer company as the first infected holding,
364 thus an exchange of virus between them cannot be ruled out.

365

366

367 **DISCUSSION**

368 Acquisition of a virulent phenotype by H7 avian influenza viruses may have devastating
369 consequences to the poultry industry and in some instance can create major human health issues,
370 including the risk of generating a new pandemic strain (23). Despite the identification of multiple
371 basic amino acids at the HA cleavage site as one of the most important molecular markers of virus
372 pathogenicity, the mechanisms underlying the emergence, spread and evolution of HPAI during an
373 epidemic are poorly understood and limited to few studies (3, 24). Here we performed for the first
374 time a deep sequencing analysis of all the samples collected during a HPAI epidemic to evaluate the
375 transmission dynamics and the within and between farms genetic diversity of the viral population.
376 We showed that the fourteen H7N7 Italian samples collected from six different farms form a cluster
377 distinct to other Eurasian sequences for all the eight gene segments, suggesting the occurrence in
378 the poultry population of a single viral introduction. The high similarity of the HA gene segment
379 with a virus collected from a wild bird in Italy and the contact between free-range hens and wild
380 waterfowl in the first infected farm (10), indicates that the LPAI progenitor strain may have been
381 introduced from the wild bird population into the first infected holding, where it rapidly mutated
382 into a HP form.

383 Despite our phylogenies suggesting a single viral introduction, we observed a high genetic
384 variability of H7N7 between the different sheds of the first infected flock. In particular, at the
385 consensus level, viruses collected from shed 5 possessed a shorter HA cleavage site and nine
386 nucleotide differences compared to the viruses from sheds 2. None of the fixed mutations had ever
387 been described in previous HPAI H7 outbreaks or recognized as associated to a specific phenotypic

388 effect. Further studies will be necessary to evaluate their possible impact on virus fitness, host range
389 and virulence. This number of nucleotide substitutions (9) is not compatible with the occurrence of
390 different introductions, when a higher number of mutations is usually observed (25). Moreover, we
391 noticed that the highest genetic distance (mean \pm standard error) between the two groups of H7N7
392 viruses (V+6 and V+9) ranged from 0.1% \pm 0.1% for the HA, NA and PA genes to 0.4% \pm 0.2% for
393 the NS gene, while, the overall mean distance among the sequences included in our phylogenies
394 ranged from 1.8% \pm 0.2% for the M gene to 5.5% \pm 0.3% for the NA gene. This evidence indicates
395 that the Italian H7N7 sequences are significantly closer to each other than any other random
396 sequences in their tree seems to be, which supports the hypothesis that the two variants V+6 and
397 V+9 had very likely derived from a single introduction. Hence, this high genetic variability can be
398 explained by i) a rapid evolution of the virus following some bottleneck events or a strong selection,
399 ii) independent evolution of the same virus within two separate sheds, or iii) the establishment of
400 the infection starting from two different seeding variants from the same progenitor viral population
401 comprising a cloud of diverse viruses. Nevertheless, our analysis of the mutation spectra of viral
402 populations suggests that the two variants arose as a consequence of a founder event or a narrow
403 population bottleneck. Indeed, the haplotype V+6, circulating in shed 5, was not identified in the
404 viral subpopulations of shed 2 and similarly haplotype V+9, identified in shed 2, was not detected
405 as a minority population in shed 5 animals. In addition, at the HA cleavage site of viruses from
406 sheds 2 (V+9) and 5 (V+6) we identified only a total of 16 and 3 reads with an insertion
407 respectively of six and nine nucleotides.

408 Entropy values obtained for the two viruses from shed 5 further supports this hypothesis. Samples
409 founded by few viral particles should have low entropy, since the strong bottleneck/founder effect
410 drastically reduce the diversity of the viral population. On the other hand, samples that experienced
411 relatively loose bottlenecks should display higher entropy. Therefore, the low entropy values of the
412 viruses from shed 5 may indicate that they had recently experienced a narrow bottleneck/founder
413 effect or that they had been subjected to a strong selection which had reduced the within-host

414 diversity and fixed adaptive mutations. To distinguish between selective and random
415 bottlenecks/founder effects we compared the relative diversity changes on a gene-by-gene basis.
416 We showed that there were no significant differences between the entropy values and the number of
417 non-synonymous mutations for the different genes, suggesting that founder effects caused by the
418 transmission bottlenecks are a major driving force of virus evolution during this epidemic.
419 On the other hand, viruses from shed 2 show intermediate entropy values, suggesting that i) they
420 were founded by a larger seeding population, ii) they experienced a high-level of replication, or iii)
421 they had circulated within the shed for a longer period of time. The latter suggestion is supported by
422 the identification of H7-specific antibodies in animals from this shed, but not in animals from sheds
423 4 and 5 (10), while the second hypothesis may be supported by the high number of dead birds found
424 in shed 2 compared to the other sheds, considering the virulence of the two variants were equal
425 (intravenous pathogenicity index of 3 for both variants, data not shown).
426 However, we cannot exclude that difference in the within-host genetic diversity between the
427 analysed samples could be associated to a different time of sampling since infection. Unfortunately,
428 the lack of information on the exact time of entrance of the virus in each farm and shed makes it
429 impossible to exclude this possible bias and to ascertain the process responsible of the reduction of
430 the genetic variability, which may be caused both by bottlenecks of different sizes that occurred
431 around the same time since infection or by similar bottlenecks that occurred at different time points
432 relative to sampling.
433 Surely, sequences of early viruses might have helped us to provide a better characterization of the
434 evolution of this strain within the index case. Indeed, the identification of H7-specific antibodies in
435 animals from shed 2 and from the outer sheds 1 and 7, where no viruses were isolated (10),
436 indicates that the virus had been circulating undetected within the farm before its identification,
437 likely with a low pathogenic phenotype.

438 Moreover, the identification of three human infections during this epidemic highlights the need for
439 constant monitoring during avian influenza outbreaks for the emergence in poultry of amino acid
440 signatures associated with interspecies transmission to provide early warning of pandemic potential.
441

442 Our analysis of the transmission dynamics indicates that only one of the two variants (V+9),
443 probably the one with the highest fitness advantage, was transmitted from the index case to the
444 other farms. Four out of the six infected farms (farms 1, 2, 4, 6) belong to one large vertically
445 integrated layer company (Fig 5), therefore virus dissemination might have occurred through shared
446 equipment, human-mediated mechanical transport, and also through infected workers, as H7N7
447 virus was diagnosed for three humans involved in the control of the epidemic (22). The low number
448 of shared mutations between farms (seven) suggests that the transmission depended on the
449 dissemination of a few viral particles. However, the high frequency threshold (2%) used in this
450 study to identify the minority variants and the scarce number of analysed samples for each farm
451 need to be taken into consideration.

452 In the farms for which it was possible to sequence more than one sample (eight for farm 1 and two
453 for farm 6), we identified the co-circulation in the same premise of different related variants and the
454 possible occurrence of multiple introductions in the same holdings (i.e. farm 6), which can be
455 detected only through the sequencing of a larger number of samples. Moreover, the high number of
456 median vectors identified between the analysed samples in our phylogenetic network reveals
457 missing ancestral sequences from our analyses, which might have been detected with increased
458 sampling. As a consequence, increasing the number of viruses sampled from each farm and also
459 from the environment could increase the resolution of our inter-farm transmission dynamic.

460 We identify farm 1 as the major source for the spread of the virus to the other four industrial
461 holdings, while the rural farm (farm 5) appears to have received the virus from the turkey farm
462 (farm 3). Interestingly, this finding allowed the National authorities to demonstrate the occurrence
463 of uncontrolled movements of birds from the infected turkey flock (farm 3), underlining the

464 importance of genetic data to complement the outbreak investigation. Despite 21 days elapsing
465 from the index case (August 13) to the last outbreak (September 3), the late depopulation date of the
466 first infected flock (August 27) and the ability of the avian virus to persist in the environment (26),
467 might explain the virus spread between these two holdings (1 and 6).

468 In addition, results of our analysis of the transmission dynamics suggests that, despite farm 2 being
469 located in close proximity to farms 4 and 6, transmission links are absent between these two
470 premises. On the contrary, the virus sampled from this farm appears to be more related to the virus
471 from farms 1 and 3, located, respectively, 38 km and 36 km from farm 2.

472 This finding suggests that multiple introductions of different viral haplotypes occurred from farm 1
473 to the other farms, probably at different time points and with different transmission modes, i.e.
474 neighbourhood spread (i.e. farm 1 and 3), human-mediated transport among farms of the same
475 company (i.e. farm 1 and 2 or farm 1 and 4). These different means of viral diffusion have been
476 observed also during other HPAI epidemics (24, 27) suggesting that long distance transmission
477 events may play an important role for the virus dissemination into new areas.

478 Overall this study shows that analysis of deep sequencing data can complement epidemiological
479 investigations, providing important insights and revealing unexpected dynamics on the inter-farm
480 transmission network. Specifically, we demonstrated that the delay in the disease detection and
481 stamping out in the index case might have been the major cause of the emergence and the spread of
482 the HPAI strain. Epidemiological investigations did not recognize the central role of the first
483 infected flock in the diffusion of the virus to most of the farms, and suggested an epidemiological
484 link between farms 2, 5 and 6, which has not been confirmed by our data. In addition the
485 epidemiological data alone was not sufficient to trace back the source of the virus detected in the
486 rural farm (farm 5), which we demonstrated to be linked to the turkey farm (farm 3).

487 Moreover, we show that a farm can harbour a high level of heterogeneity, potentially caused either
488 by separate bottlenecks and founder effects in the different sheds, or by multiple viral introductions
489 from different sources. Hence, the importance during the control activities to collect and analyse

490 several samples from each infected farm to provide a complete picture of the evolutionary process
491 during an avian influenza epidemic.

492

493

494 **FUNDING INFORMATION**

495 This work was financially supported by the European projects Epi-SEQ (research project supported
496 under the 2nd Joint Call for Transnational Research Projects by EMIDA ERA-NET [FP7 project
497 no. 219235]) and PREDEMICS (research project supported by the European Community's Seventh
498 Framework Programme [FP7/2007-2013] under grant agreement n. 278433). PRM and JH are
499 supported by the Medical Research Council of the United Kingdom (grant number G0801822).

500

501

502 **ACKNOWLEDGEMENTS**

503 We acknowledge the authors, originating and submitting laboratories of the sequences from
504 GISAID's EpiFlu Database used in this study. A detailed list is found in supplementary table S2.

505

506 **REFERENCES**

- 507 1. **Manrubia SC, Escarmís C, Domingo E, Lázaro E.** 2005. High mutation
508 rates, bottlenecks, and robustness of RNA viral quasispecies. *Gene* **347**:273–
509 282.
- 510 2. **Iqbal M, Reddy KB, Brookes SM, Essen SC, Brown IH, McCauley JW.**
511 2014. Virus Pathotype and Deep Sequencing of the HA Gene of a Low
512 Pathogenicity H7N1 Avian Influenza Virus Causing Mortality in Turkeys.
513 *PLoS One* **9**:e87076.

- 514 3. **Monne I, Fusaro A, Nelson MI, Bonfanti L, Mulatti P, Hughes J, Murcia**
515 **PR, Schivo A, Valastro V, Moreno A, Holmes EC, Cattoli G.** 2014.
516 Emergence of a Highly Pathogenic Avian Influenza Virus from a Low-
517 Pathogenic Progenitor. *J Virol* **88**:4375–4388.
- 518 4. **Jonges M, Welkers MRA, Jeeninga RE, Meijer A, Schneeberger P,**
519 **Fouchier RAM, de Jong MD, Koopmans M.** 2014. Emergence of the
520 virulence-associated PB2 E627K substitution in a fatal human case of highly
521 pathogenic avian influenza virus A(H7N7) infection as determined by Illumina
522 ultra-deep sequencing. *J Virol* **88**:1694–1702.
- 523 5. **Poole DS, Yú S, Cai Y, Dinis JM, Müller MA, Jordan I, Friedrich TC,**
524 **Kuhn JH, Mehle A.** 2014. Influenza A virus polymerase is a site for adaptive
525 changes during experimental evolution in bat cells. *J Virol* **88**:12572–12585.
- 526 6. **Varble A, Albrecht RA, Backes S, Crumiller M, Bouvier NM, Sachs D,**
527 **García-Sastre A, tenOever BR.** 2014. Influenza A Virus Transmission
528 Bottlenecks Are Defined by Infection Route and Recipient Host. *Cell Host*
529 *Microbe* **16**:691–700.
- 530 7. **Wilker PR, Dinis JM, Starrett G, Imai M, Hatta M, Nelson CW, O'Connor**
531 **DH, Hughes AL, Neumann G, Kawaoka Y, Friedrich TC.** 2013. Selection
532 on haemagglutinin imposes a bottleneck during mammalian transmission of
533 reassortant H5N1 influenza viruses. *Nat Commun* **4**:2636.
- 534 8. **Fusaro A, Tassoni L, Hughes J, Milani A, Salviato A, Schivo A, Murcia**
535 **PR, Bonfanti L, Cattoli G, Monne I.** 2015. Evolutionary trajectories of two

- 536 distinct avian influenza epidemics: Parallelisms and divergences. *Infect Genet*
537 *Evol* **34**:457–466.
- 538 9. **Yu X, Jin T, Cui Y, Pu X, Li J, Xu J, Liu G, Jia H, Liu D, Song S, Yu Y,**
539 **Xie L, Huang R, Ding H, Kou Y, Zhou Y, Wang Y, Xu X, Yin Y, Wang J,**
540 **Guo C, Yang X, Hu L, Wu X, Wang H, Liu J, Zhao G, Zhou J, Pan J, Gao**
541 **GF, Yang R, Wang J.** 2014. Influenza H7N9 and H9N2 viruses: coexistence in
542 poultry linked to human H7N9 infection and genome characteristics. *J Virol*
543 **88**:3423–3431.
- 544 10. **Bonfanti L, Monne I, Tamba M, Santucci U, Massi P, Patregnani T, Loli**
545 **Piccolomini L, Natalini S, Ferri G, Cattoli G, Marangon S.** 2014. Highly
546 pathogenic H7N7 avian influenza in Italy. *Vet Rec* **174**:382–382.
- 547 11. **Bosch FX, Garten W, Klenk HD, Rott R.** 1981. Proteolytic cleavage of
548 influenza virus hemagglutinins: primary structure of the connecting peptide
549 between HA1 and HA2 determines proteolytic cleavability and pathogenicity of
550 Avian influenza viruses. *Virology* **113**:725–735.
- 551 12. **Spackman E, Senne D a, Myers TJ, Bulaga LL, Garber LP, Perdue ML,**
552 **Lohman K, Daum LT, Suarez DL.** 2002. Development of a Real-Time
553 Reverse Transcriptase PCR Assay for Type A Influenza Virus and the Avian
554 H5 and H7 Hemagglutinin Subtypes Development of a Real-Time Reverse
555 Transcriptase PCR Assay for Type A Influenza Virus and the Avian H5 and H7
556 Hemagglutini. *J Clin Microbiol* **40**:3256–3260.
- 557 13. **Zhou B, Donnelly ME, Scholes DT, St. George K, Hatta M, Kawaoka Y,**

- 558 **Wentworth DE.** 2009. Single-Reaction Genomic Amplification Accelerates
559 Sequencing and Vaccine Production for Classical and Swine Origin Human
560 Influenza A Viruses. *J Virol* **83**:10309–10313.
- 561 14. **Boc A, Diallo AB, Makarenkov V.** 2012. T-REX: a web server for inferring,
562 validating and visualizing phylogenetic trees and networks. *Nucleic Acids Res*
563 **40**:W573–W579.
- 564 15. **Jombart T, Eggo RM, Dodd PJ, Balloux F.** 2011. Reconstructing disease
565 outbreaks from genetic data: a graph approach. *Heredity (Edinb)* **106**:383–390.
- 566 16. **Jombart T.** 2008. adegenet: a R package for the multivariate analysis of
567 genetic markers. *Bioinformatics* **24**:1403–1405.
- 568 17. **Gabor C, Tamas N.** 2006. The igraph software package for complex network
569 research. *InterJournal*.
- 570 18. **Katoh K, Standley DM.** 2013. MAFFT Multiple Sequence Alignment
571 Software Version 7: Improvements in Performance and Usability. *Mol Biol*
572 *Evol* **30**:772–780.
- 573 19. **Stamatakis A.** 2014. RAxML version 8: a tool for phylogenetic analysis and
574 post-analysis of large phylogenies. *Bioinformatics* **30**:1312–1313.
- 575 20. **Bandelt HJ, Forster P, Röhl A.** 1999. Median-joining networks for inferring
576 intraspecific phylogenies. *Mol Biol Evol* **16**:37–48.
- 577 21. **Huson DH, Bryant D.** 2006. Application of phylogenetic networks in
578 evolutionary studies. *Mol Biol Evol* **23**:254–67.
- 579 22. **Puzelli S, Rossini G, Facchini M, Vaccari G, Di Trani L, Di Martino A,**

- 580 **Gaibani P, Vocale C, Cattoli G, Bennett M, McCauley JW, Rezza G, Moro**
581 **ML, Rangoni R, Finarelli AC, Landini MP, Castrucci MR, Donatelli I.**
582 2014. Human infection with highly pathogenic A(H7N7) avian influenza virus,
583 Italy, 2013. *Emerg Infect Dis* **20**:1745–1749.
- 584 23. **Capua I, Marangon S.** 2006. Control of Avian Influenza in Poultry. *Emerg*
585 *Infect Dis* **12**:1319–1324.
- 586 24. **Bataille A, van der Meer F, Stegeman A, Koch G.** 2011. Evolutionary
587 analysis of inter-farm transmission dynamics in a highly pathogenic avian
588 influenza epidemic. *PLoS Pathog* **7**:e1002094.
- 589 25. **Bouwstra R, Koch G, Heutink R, Harders F, van der Spek A, Elbers A,**
590 **Bossers A.** 2015. Phylogenetic analysis of highly pathogenic avian influenza
591 A(H5N8) virus outbreak strains provides evidence for four separate
592 introductions and one between-poultry farm transmission in the Netherlands,
593 November 2014. *Eurosurveillance* **20**:21174.
- 594 26. **Brown JD, Swayne DE, Cooper RJ, Burns RE, Stallknecht DE.** 2007.
595 Persistence of H5 and H7 avian influenza viruses in water. *Avian Dis* **51**:285–
596 289.
- 597 27. **Souris M, Gonzalez J-P, Shanmugasundaram J, Corvest V, Kittayapong P.**
598 2010. Retrospective space-time analysis of H5N1 Avian Influenza emergence
599 in Thailand. *Int J Health Geogr* **9**:3.

600

601

Table 1. Epidemiological information of the 14 samples collected during the HPAI H7N7 outbreak (TS= tracheal swabs).

Farm	Sample	RNA copies/ μ l	Mean depth of coverage	Sample type	Farm type	Collection date	Province	Number of birds	Depopulation date
1 shed 2	4527-11	1,24E+05	19354	Pool of 10 TS	Laying hen (industrial farm)	13 Aug 2013	Ferrara	128000	27 Aug 2013
	4527-12	9,22E+07	36772	Pool of 10 TS					
	4541-7	1,04E+06	24696	Organ pool*					
	4541-32	2,49E+07	53292	Kidney					
1 shed 4	4541-8	5,24E+05	34018	Organ pool*					
	4541-33	4,98E+03	42661	Kidney					
1 shed 5	4541-9	3,93E+04	23390	Organ pool*					
	4541-34	4,32E+04	58893	Kidney					
2	4603-1	2,89E+05	43810	Pool of 10 TS	Laying hen (industrial farm)	19 Aug 2013	Bologna	584900	8 Sept 2013
3	4678	7,88E+05	19893	Organ pool*	Meat turkey (industrial farm)	21 Aug 2013	Ferrara	19850	27 Aug 2013
4	4774	2,30E+08	31804	Organ pool**	Laying hen (industrial farm)	27 Aug 2013	Bologna	121705	8 Sept 2013
5	5091	1,21E+07	24510	Organ pool*	Backyard flock	2 Sept 2013	Ferrara	3	5 Sept 2013
6	5051-1	5,27E+07	46615	Trachea	Pullets (industrial farm)	3 Sept 2013	Bologna	98200	8 Sept 2013
	5051-3	1,02E+08	48562	Trachea					

* pool of organs from 2 animals

** pool of organs from 3 animals

Table 2. Number of reads showing a 0 to 9 nucleotide insertion (compared to the sequence of a H7 LPAI strain, CCAAAGAGAAGA) at the HA cleavage site of the eight samples collected from three different sheds of the index case

N. nt insertion	SHED 5		SHED 4		SHED 2			
	4541-34	4541-9	4541-8	4541-33	4527-11	4527-12	4541-7	4541-32
0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0
3	0	1	0	0	0	0	0	0
4	0	0	0	1	0	0	0	0
5	0	1	0	0	0	0	0	0
6	23509	14861	0	591	11	0	1	4
7	0	0	4	18	5	3	5	27
8	0	0	0	0	1	0	0	5
9	2	1	16587	13700	6660	13725	4439	22929

FIGURE LEGENDS

Fig 1. Distribution of nucleotide frequency differences between three technical replicates. For each genome position with coverage >500 the frequency differences between the four bases (A, C, T and G) were obtained from the comparison of the replicates of the three samples: 4541-8 in yellow, 4541-9 in violet, 4541-34 in blue. The y-axis represents the percentage of nucleotide positions where the highest frequency differences fall within the ranges 0-0.1%, 0.1-0.25%, 0.25-0.5%, 0.5-1%, 1-2% and >2% (x-axis). Frequency differences higher than 2% were observed in only 0.3%-0.4% of all the analysed positions (11501 to 13308) for all the replicates. Thus a 2% threshold allows the exclusion of 99.6% of the possible errors.

Fig 2. ML phylogenetic tree of the HA gene segment of 172 H7 avian influenza viruses. HPAI H7N7 viruses collected during Italian epidemic are coloured according to the farm of collection: grey for farm 1, purple for farm 2, light blue for farm 3, yellow for farm 4, green for farm 5 and orange for farm 6. The numbers at nodes represent bootstrap values (>70%), while branch lengths

are scaled according to the numbers of nucleotide substitutions per site. The tree is mid-point rooted for clarity only.

Fig 3. Consensus level nucleotide and amino acid differences among the complete genome of the 14 Italian H7N7 viruses. Each sample (column) is coloured according to the farm of collection: grey for farm 1, purple for farm 2, light blue for farm 3, yellow for farm 4, green for farm 5 and orange for farm 6. The farm and shed of belonging (i.e. 1-shed 5, corresponds to farm 1, shed 5) and the sample type is indicated above the sample name. The nucleotide (NT) differences identified between each sample and the viruses from shed 5 of the index case (samples 4541-9 and 4541-34, column 1 and 2) are reported. Amino acid mutations (AA) are highlighted in red, while silent mutations are in black. The 11 different genomes identified during this epidemic are indicated in the last row (A to K).

Fig 4. Heat-map of the nucleotide frequency. The horizontal axis represents the samples, coloured according to the farm of collection, while the vertical axis display only the variable nucleotide positions showing nucleotide differences compared to the samples 4541-9 and 4541-34. The colour scale represents the nucleotide frequency according to the scale bar at the top of the figure. White spaces represent positions for which deep sequencing data were not available (coverage <500). Black arrows indicate the variants that are fixed in the viral population of at least one sample of the first infected farm. The dendrogram above the heatmap represents the neighbour-joining tree obtained from the distance matrix calculated from the deep sequencing data.

Fig 5. Transmission tree obtained from deep sequencing data. Each circle represents an individual sample, coloured according to the farm of collection. The size of the circles is proportional to the mean entropy value. The vertical axis represents the time of collection of each sample (samples in the same row belong to the same farm). Circle outlines are assigned accordingly

to the owner of the farm as shown in the figure legend. Connecting arrows correspond to the results obtained from SeqTrack, while dashed lines are alternative hypotheses of transmission events formulated based on the number of shared mutations. Numbers over the lines are the genetic distance calculated from the deep sequencing data between the samples. Coloured area represents genetic groups identified based on the number of shared mutations and the results of both the neighbour-joining phylogenetic tree (Fig 4) and the network analysis (Fig 6).

Fig 6. Median-joining phylogenetic network. A) The network was constructed from the consensus sequences of the eight concatenated gene segments. Each unique sequence genotype is represented by a circle sized relatively to its frequency in the dataset. Numbers next to the circles correspond to the samples showing that particular genotype, while the number within the circle represents the shed where the genotype was identified. Genotypes are coloured according to farm. Branches represent the shortest trees and black circles represent the number of nucleotide mutations that separate each node. Median vectors are shown as red circles. The violet and yellow shading represent the two identified genetic groups C1 and C2. Numbers at each branch represent bootstrap values. B) The map shows the geographic position of the six infected farms.