

A peer-reviewed version of this preprint was published in PeerJ on 30 October 2014.

[View the peer-reviewed version](https://peerj.com/articles/655) (peerj.com/articles/655), which is the preferred citable publication unless you specifically need to cite this preprint.

Dalby AR, Iqbal M. 2014. A global phylogenetic analysis in order to determine the host species and geography dependent features present in the evolution of avian H9N2 influenza hemagglutinin. PeerJ 2:e655 <https://doi.org/10.7717/peerj.655>

A global phylogenetic analysis in order to determine the host species and geography dependent features present in the evolution of avian H9N2 Influenza Hemagglutinin

A complete phylogenetic analysis of all of the H9N2 hemagglutinin sequences that were collected between 1966 and 2012 was carried out in order to build a picture of the geographical and host specific evolution of the hemagglutinin protein. To improve the quality and applicability of the output data the sequences were divided into subsets based upon location and host species.

The phylogenetic analysis of hemagglutinin reveals that the protein has distinct lineages between China and the Middle East, and that wild birds in both regions retain a distinct form of the H9 molecule, from the same lineage as the ancestral hemagglutinin. The results add further evidence to the hypothesis that the current predominant H9N2 hemagglutinin lineage might have originated in Southern China. The study also shows that there are sampling problems that affect the reliability of this and any similar analysis. This raises questions about the surveillance of H9N2 and the need for wider sampling of the virus in the environment.

The results of this analysis are also consistent with a model where hemagglutinin has predominantly evolved by neutral drift punctuated by occasional selection events. These selective events have produced the current pattern of distinct lineages in the Middle East, Korea and China. This interpretation is in agreement with existing studies that have shown that there is widespread intra-country sequence evolution.

- 1 Andrew R. Dalby
- 2 Faculty of Science and Technology, University of Westminster, 115 New Cavendish Street, Westminster,
- 3 W1W 6UW, UK.
- 4 Munir Iqbal
- 5 Avian Viral Diseases Programme, The Pirbright Institute, Compton Laboratory, Newbury, Berkshire.
- 6 RG20 7NN, UK

- 7 A.Dalby@westminster.ac.uk

8 Introduction

9 The spread of avian influenza viruses (AIV) are a major cause of concern for global animal and
10 public health; these viruses are causing enormous economic losses as well as posing a credible
11 threat for pandemic emergence (Sorrell et al. 2009; Xu et al. 2004; Yu et al. 2011). There has
12 been an increase in the monitoring of disease outbreaks in wild and domestic birds, as well as in
13 other potential hosts such as swine and humans, but we still lack a coordinated global
14 surveillance network (Butler 2012). The Influenza A/H5N1 virus has been the main focus of
15 international monitoring after a series of recent outbreaks, but the emergence of the A/H1N1
16 pandemic virus “swine flu” in 2009 showed that other subtypes also pose a serious threat to
17 human health (Cao et al. 2009). Experiments have been carried out to determine the exact factors
18 of bird to human transmission and of droplet transmission of H9N2 viruses (Sorrell et al. 2009).

19 The H9N2 subtype is a variant of AIV usually associated with low pathogenicity. Due to the
20 lower pathogenicity phenotype of this virus, data collection has been very sporadic. There have
21 been outbreaks of H9N2 in flocks of domestic birds resulting in significant economic loss and
22 with high mortality rates of up to 60% reported during the epizootic of 1998-2001 in Iran (Nili &
23 Asasi 2002). This subtype has also been shown to pass to pigs, ferrets and guinea pigs, as well as
24 to humans in a small number of cases (Butt et al. 2005; Cheng et al. 2011; Lin et al. 2000; Lv et
25 al. 2012; Peiris et al. 1999; Wan et al. 2008; Xu et al. 2004; Yu et al. 2008; Zhang et al. 2009).

26 Antibodies to the virus have also been found in a sero-epidemiological investigation of poultry
27 workers (Pawar et al. 2012; Wang et al. 2009). These cross species infections show that in the
28 future the virus may present a serious threat to human health.

29 The co-circulation of H9N2 with other H5N1, H7N3, H1N1 and H3N2 subtypes has resulted in
30 the emergence of novel reassortant viruses (Monne et al. 2013; Peiris et al. 2001; Sun et al. 2011).

31 The reassorted virus has been shown to possess increased virulence (Iqbal et al. 2009; Marshall et

32 al. 2013). The recent emergence of a novel reassortant H7N9 virus containing internal genes from
33 the H9N2 virus is another of novel AIV in birds which has the capability of infecting humans,
34 with fatal consequences. However in the cases of reassortment human to human transmission has
35 not been demonstrated (Watanabe et al. 2013).

36 There have been a number of recent studies on the evolution of AIV that have incorporated
37 geographical data available from global influenza monitoring (Fusaro et al. 2011; Haase et al.
38 2010; LAM et al. 2012; Wallace et al. 2007). With the growth in the global monitoring efforts,
39 and the widespread use of cheaper DNA sequencing technology there has been a rapid expansion
40 in the number of available sequences. Previous phylogenetic studies of H9N2 hemagglutinin,
41 have focussed on sequences from a single location (Banks et al. 2000; Butt et al. 2010; Kim et al.
42 2006; Li et al. 2005; Song et al. 2011; Xu et al. 2007a). The largest previous phylogeographical
43 study was that of Fusaro *et al.* who surveyed all of the H9N2 viral segments from the Middle East
44 (Fusaro et al. 2011). Fusaro's study defined eight geographical regions covering the Middle East
45 and used maximum likelihood methods to construct the phylogenetic analysis. Bayesian methods
46 were used to evaluate the geographical clade distribution. This work has recently been extended
47 to take a more detailed look at viral evolution in Israel and to show that there have been
48 successive introductions from neighbouring countries (Davidson et al. 2014).

49 A large-scale phylogenetic analysis using eight viral gene segments from 571 complete genomic
50 sequences collected between 1966 and 2009 was carried out by Dong and co-workers (Dong et
51 al. 2011). Geographical details were not the focus of this study as the aim was to establish the
52 lineage structure and genotypes present in H9N2. That analysis revealed 74 lineages and 98
53 genotypes when re-assortment is taken into consideration, but they only identified 7 HA variants.

54 Investigations of the evolution of influenza A H3N2 hemagglutinin using evidence from flu
55 antigen evolution, have shown how the rates of evolution can vary between selection events

56 (epochs) (Koelle et al. 2006; Smith et al. 2004; Thomas & Hertz 2012; Wolf et al. 2006). Wagner
57 used this study to suggest a reconciliation between the selectionist and neutralist views in a
58 network based model (Wagner 2008). Wagner's model also shows that the order in which
59 mutations take place can have an effect on the selection of a group of mutations. In this way
60 random drift is punctuated by selective epochs. Using mathematical models, Bedford *et al.* have
61 shown that in A/H3N2 sequence evolution is constrained by canalisation (Bedford et al. 2012).
62 This agrees with Wagner's hypothesis if there are times when there is limited drift but occasional
63 bursts where the organism escapes the canalisation (Thomas & Hertz 2012; Wolf et al. 2006).

64 It is difficult to assess the performance of phylogenetic models. The methods are sensitive to the
65 distance measures used and this is reflected in the number of sites that can be compared and the
66 number of sequences in the study (Felsenstein & Felsenstein 2004). Clades should be
67 monophyletic if they are produced by divergent evolution (Page & Holmes 2009). If a lineage is
68 geography and host specific then we would expect all of the sequences that form a clade to share
69 the same labels in terms of geography and the species in which they are found. Single base or
70 amino acid changes are ambiguous and could be a product of either divergence or convergence,
71 but larger conserved patterns, especially if they are non-consecutive in the sequence alignments
72 are good indicators of mutations that are responsible for differentiating between clades. These
73 patterns of change can then be examined in terms of their effects on protein structure and
74 function. Ultimately it is the biological function that determines how selection has been
75 responsible for clade differentiation.

76 The current study presents a comprehensive phylogenetic analysis of the H9N2 HA, that includes
77 sequences from all of the geographical regions where H9N2 has been reported. This investigation
78 shows how host species and geographical distribution have shaped the evolution of distinct
79 lineages of H9 HA. This reveals clades that are both geographically and host species dependent.

80 From these analyses new hypotheses can be generated for more specific events such as
81 migrations or intra-species infections.

82 Materials and Methods

83 The complete set of H9N2 hemagglutinin protein sequences were downloaded from the NCBI on
84 the 22nd of May 2013. The search term used for searching the protein database was **H9N2**
85 **hemagglutinin**. The sequences were exported in FASTA format.

86 The complete set of H9N2 hemagglutinin nucleotide sequences were downloaded from the NCBI
87 on the 9th of September 2013. The search term used was **((H9N2 hemagglutinin) NOT**
88 **precursor) NOT partial**. The sequences were downloaded in FASTA format.

89 A dataset for the complete set of Korean H9N2 hemagglutinin sequences was downloaded from
90 the NCBI on the 15th of June 2014. The search term used was **((H9N2 hemagglutinin) NOT H5)**
91 **Korea)**.

92 After removal of long truncations (> 40 amino-acid residues) as well as a group of sequences that
93 were actually from other AIV subtypes, the final protein dataset contained 2045 sequences, the
94 final nucleotide dataset contained 1052 sequences of which the Korean nucleotide subset
95 contained 64 sequences.

96 The edited protein and nucleotide datasets were then broken into sub-groups based on the
97 sequence annotations, using a short text matching program written in Perl. The data was split into
98 subsets based on geographical location and host species.

99 The geographical subsets were based on national boundaries except for China, which was divided
100 into its regions. Where a Chinese region had 10 sequences or less then phylogenetic analysis was
101 not carried out. National borders have been identified as barriers to influenza transmission and so
102 these are appropriate geographical subsets (Wallace & Fitch 2008). A subset was created for the

103 Americas including all of the North American and the single South American (Argentinian)
104 sequence.

105 Species trees were created for chicken, duck (including mallard), quail, pheasant and swine.
106 There are numerous wild bird species and also some environmental samples, but these are often
107 only represented by single cases and so subsets were not created.

108 All of the sequence analysis and editing was carried out in MEGA 5.2 on a Windows 8 computer.
109 Both the protein and gene sequences were aligned using Muscle within the MEGA sequence
110 analysis package, using the default parameters (Edgar 2004; Tamura et al. 2011).

111 Model evaluation was carried out for the protein dataset within MEGA. This analysis showed that
112 the Jones-Taylor-Thornton model with gamma distributed rates amongst sites was the best model
113 (supplementary table 1) (Jones et al. 1992). Phylogenetic trees for the protein sequences were
114 then constructed using Maximum Likelihood within MEGA using the JTT + G substitution
115 model. Each model was tested with 500 bootstraps replicates. To make the calculations tractable
116 the heuristic nearest neighbour interchange was used. The initial tree was created with
117 neighbourhood joining.

118 Model evaluation for the corresponding nucleotide dataset showed that the Tamura-Nei
119 substitution model with gamma distributed rates amongst sites was the best performing model
120 (supplementary table 2) (Tamura & Nei 1993). Phylogenetic trees for the gene sequences were
121 constructed using Maximum Likelihood within MEGA using the TN93 + G substitution model.
122 Models were tested with 500 bootstrap replicates. Smaller numbers of sequences in the gene
123 dataset meant that it was possible to calculate the bootstrap values for the chicken trees, which
124 could not be calculated in the protein trees.

125 The amino acids responsible for clade formation were determined manually using the alignment
126 view within MEGA. This view shows only the amino acid changes relative to the conserved
127 sequence, and so it is a matter of scanning the sequences looking for substitutions that correspond
128 to the different clusters.

129 All of the trees have been presented as cladograms rather than as phylograms because the paper
130 uses a cladistic approach for the analysis. The topology (ordering of the clusters) is significant
131 factor for investigation rather than the distances between groups. The removal of this distance
132 data and also the bootstrap data from the figures improves the clarity of the diagrams, but this
133 omitted information is available in the supplementary files that include the full phylogenetic
134 analysis. Results and Discussion

135 **Global Nucleotide Phylogenetic Tree Analysis.**

136 A complete phylogenetic analysis of the nucleotide dataset was carried out. Computational limits
137 on memory and processor speed make it impossible to carry out a complete analysis on the
138 protein dataset (especially as it contains almost twice as many sequences). A condensed view of
139 the principal clades from the global nucleotide tree can be seen in figure 1. One noticeable
140 absence from the global nucleotide tree is the quail sequences. These were omitted during the
141 editing of the sequences to remove truncated sequences where the ends of the sequences were
142 missing. The virus can be broken into three main clusters labelled A, B and the Main Chinese
143 Clade. (The complete tree can be found in supplementary figure 1). There is also a small clade
144 (labelled C) that splits from the root of the tree between clades A and B. This clade contains
145 Chinese sequences that were isolated from chickens between 1999-2002. This topology agrees
146 with those from the existing literature except that there is some disagreement in the identification
147 of the lineages (Li et al. 2003; Peiris et al. 2001; Perk et al. 2006; Yu et al. 2008). Ji used a
148 lineage and sub-lineage nomenclature that identified four lineages and 2 sub-lineages
149 corresponding to Clade A (lineages 9.1, 9.2 and 9.3), Clade B (sub-lineage 9.4.1) and the main
150 Chinese Clade (sub-lineage 9.4.2). Huang and co-workers used a lineage naming system based on
151 the prototypes Y439, TY WS 66, G1 and Y280 (Huang et al. 2010). Y439, and TY WS 66
152 (Turkey, Wisconsin 1966) are both in clade A. G1 is in clade B and Y280 is in the main Chinese

153 clade in China 1 1994-2010. The large study by Dong identified HK/G1/97, BJ/1/94, HK/289/78,
154 HK/AF157/92, KR/96323/96, DE/113/95 and WI/1/66 as prototype sequences for the different
155 HA lineages (Dong et al. 2011). The G1 lineage and WI/1/66 lineages are well established and
156 correspond to Clade B and Clade A respectively. BJ/1/94 is in the same clade as Y280,
157 HK/289/78, HK/AF157/92 and DE/113/95 are not in the current dataset because of truncations,
158 but there are sub-clades for Hong Kong duck sequences and Korean chicken sequences that
159 correspond to HK/289/78 and KR/96323/96 respectively. This suggests that it might be possible
160 to divide Clade a into further sub-clades but given the limited number of sequences sampled there
161 is insufficient evidence to be able to carry this out at present.

162 Within the Main Chinese clade there are a series of nested sub-clades labelled China 1-9 where
163 the annotations follow the correct date order, with the exception of clade 8. In clade 8 there is a
164 single sequence from a chicken in the Shandong region collected in 1999, which is a probable
165 outlier. All of the other sequences within the clade are from 2003 onwards, which would be
166 consistent with the splitting dates for the other sub-clades. It is possible that this is an example of
167 convergent evolution between recent Chinese sequences and an earlier branching of the
168 phylogenetic tree, but a single sequence is insufficient evidence to corroborate this. An alternative
169 explanation is that this results from sequencing error but this is also unlikely given the large
170 number of bases that would have to be incorrectly identified (for an alignment between the
171 Shandong sequence, G1, G9/Y280, Wisconsin 1966 sequences see supplementary figure 2). It is
172 possible that this could be a database error, but again this seems unlikely given the provenance
173 and tracking systems within GenBank. In the absence of database error and a clear evolutionary
174 connection between this lone sequence and the rest of the clade does provide evidence for the
175 inadequate sampling of sequences. This nested structure of the Chinese sequences had previously
176 been reported by Song et al. (Song et al. 2011).

177 The expanded tree (supplementary figure 1) shows that the geographical sampling has not been
178 systematic and that it has been carried out in a sporadic and haphazard manner. There are breaks
179 in the dates of sequence annotations in some of the clades that are homogeneous for location.
180 Date gaps in annotations at specific locations are likely to be a result of inadequate sampling
181 rather than reliable evidence for the loss and subsequent migratory return of the clade. The focus
182 on China because of outbreaks of H9N2 within flocks of domestic birds has resulted in a very
183 dense phylogenetic tree for the Chinese sequences, resulting in an artificially Chinese focused
184 distribution to the phylogenetic trees. There is only a single sequence from South America this is
185 the result of a sampling effect rather than reflecting the actual H9N2 distribution.

186 Clade A (figure 2) corresponds to wild bird infections that are distributed world-wide and include
187 examples from North America, South America, Europe and the Far East. This is an important
188 clade because it also contains the original sequence of the hemagglutinin from the H9N2 subtype
189 that was found in a turkey in Wisconsin in 1966. This clade also contains a number of recent
190 sequences from Korea the US and Europe and so this clade remains extant. The topology of this
191 clade disagrees with that from Kim and co-workers who constructed a tree for regions 1-1104
192 using DNA Star (Kim et al. 2006). The complete sequence for HA is over 1700 bases, and the
193 method used here is Maximum Likelihood, whereas DNA Star uses the less reliable UPGMA tree
194 generation algorithm within ClustalW. In their tree the recent Korean sequences were placed
195 outside clade A and beyond the Y280 sequences as a distinct clade. The bootstrap values are high
196 for this region of the tree (> 99%) and it has geographical consistency with the rest of the clade,
197 and so the topology presented in the current paper is most likely to be correct.

198 Clade B (figure 3) contains mostly Middle Eastern sequences although there are a few Chinese
199 sequences that form a sub-clade. This clade corresponds to the extended G1 lineage used by
200 Fusaro *et al.* and Monne *et al.* (Monne et al. 2013) Iran and Israel are the two countries most

201 strongly represented in this clade. There are also a number of sequences from the Arabian
202 peninsular and a large grouping from the Indian sub-continent.

203 The sub-lineage structure of Clade B shows that the sequences have evolved substantially from
204 the G1 prototype. The G1 prototype sub-clade became extinct in Iran in 2004 and in Israel in
205 2007. This was replaced by a new sub-lineage (labelled 725 sub-clade in figure 3) that appears to
206 have originated in chicken flocks in Iran in 1998, although it is only found in a large number of
207 Iranian samples after 2004. This sub-lineage is similar to that previously identified by Fusaro *et*
208 *al.* and Monne *et al.* (labelled cluster C in their papers) (Fusaro et al. 2011; Monne et al. 2013).
209 There is a sub-clade from the Indian subcontinent that originated in the Punjab/Haryana region in
210 2003. There is a linking clade than originated in Pakistan in 2004 before spreading back to Iran in
211 2009. The final sub-clade seems to have originated in the Arabian Peninsular in 2006 and
212 correspond to cluster B from the studies of Fusaro *et al.* and Monne *et al.* This sub-lineage then
213 spread to Israel and most recently to Egypt. This is in good agreement with the previous studies
214 on the origin of the Egyptian virus (Abdel-Moneim et al. 2012). The phylogenetic tree shows that
215 after the initial Egyptian outbreak in 2010 there has been a marked diversification in the
216 sequences.

217 Clade C (figure 4) is a Chinese clade that has annotated dates between 1999 and 2001. In the
218 phylogenetic analysis this clade falls between Clade A and Clade B. As there are only a small
219 number of sequences within the clade there is no discernable pattern to the distribution of the
220 sequences within the Chinese regions, although it appears to have originated in Guangdong in
221 1999. As there are no more sequences after those from 2002 this clade can be considered extinct.

222 The first Chinese clade contains the G9/Y280 lineage (figure 5). After the banning of live quail
223 from poultry markets the G1 lineage disappeared from Chinese poultry leaving only the G9/Y280
224 lineage (Choi et al. 2004). In the study of Li et al. four lineages were specified G1, TY/WS/66,

225 Y439 and Beijing/1/94, but no sequences from G1 or Y439 were found in their sample (Li et al.
226 2005). Cong and co-workers identified two more lineages within this clade based on antigenic
227 studies and nucleotide phylogenetic trees (Cong et al. 2007). These are represented by prototype
228 sequences from Shanghai 1998 and a swine genotype. These lineages only cover a small number
229 of the possible sub-clades within this first Chinese Clade (figure 5). As no sequences have been
230 sampled from this clade since 2010 it is possible that the G7/Y280 lineage is now extinct, but this
231 hypothesis cannot be confirmed without a longer period of absence of viruses from this clade.

232 This study shows a series of new clades that had not been previously identified and that originate
233 in 1997 in Beijing. This is also the origin of the clade labelled China 2. All of the subsequent
234 nested Chinese clades are related to this Beijing sequence. Some sub-clades such as China 4,
235 China 5 and China 6 also appear to have become extinct and so there is good evidence for
236 successive selective sweeps through the viral population.

237 As discussed previously the individual unusual sequence from 1999 in China 8 sub-clade is
238 difficult to explain. It might suggest that there is considerable sequence diversity within these
239 nested clades, but because of the low number of individuals carrying a particular sequence it
240 might take a long-time for a sequence to become fixed sufficiently within the population to be
241 found by sporadic sampling. If this is true then all of the clades probably have a much earlier
242 origin and there is a long period before first detection. This causes some concern for tracking the
243 evolution of potential pandemic strains, as they might be circulating quite widely before they are
244 detected for the first time.

245 **Clade Analysis of the Geographical Subsets**

246 From the global phylogenetic tree three interesting geographical subsets were identified; the tree
247 for Korea because it is homogeneous to clade A, and the trees from Iran and Israel as they show

248 successive waves of sequence evolution in clade B. The main Chinese clade has many interesting
249 features but these are difficult if not impossible to untangle and coherently explain because of the
250 limitations of sampling within the data, where there are only small numbers from some regions
251 and then large numbers from others.

252 Both the protein and nucleotide phylogenetic trees for Korea are shown in figure 6. There is a fair
253 agreement between the two trees but once again this emphasizes the problems of sampling and
254 the possible effects this can have on phylogenetic reconstruction. From the nucleotide tree it is
255 tempting to define a clade that contains only wild birds and that suddenly developed in 2005, but
256 this clade is not present in the protein tree, and so it cannot be definitively assigned. There are
257 also some differences in the topology surrounding the swine flu case (marked with a red
258 diamond). Previous studies have investigated the effect of vaccination, which initially
259 suppressed the number of cases that were seen in 2008 before an increase in the number of cases
260 in 2009 (Park et al. 2011). From the trees presented here the period following the introduction of
261 vaccination does correspond to a period of diversification. All of the Korean sequences are from
262 clade A, which is the longest circulating lineage and previously it had not shown a wide diversity
263 of sequences. Vaccination is likely to have had an impact on hemagglutinin evolution, which can
264 be seen in the 2009-2010 clade (Lee & Song 2013).

265 The protein and nucleotide phylogenetic trees for Israel and Iran are shown in Figure 7. In the
266 Iranian tree Clade 1 are duck sequences from Clade A. This is why this clade has the deepest
267 origin within the tree, but it does not contain the earliest sequences. This shows that wild birds
268 can introduce other lineages to a geographical region where another lineage currently dominates.

269 Clades 2 and 5 in the Iranian tree are the oldest clade from the G1-like lineage in these regions.

270 The G1 lineage appears to have originated in Hong Kong in 1997. The initial G1 clade (Clade 2)
271 becomes extinct in 2003 in the nucleotide tree and in 2005 in the protein tree indicating a

272 selective sweep. This event coincides with the loss of sister Clades 1 and 2 in the Israeli protein
273 tree and Clade 1 in the nucleotide tree. These clades were assigned to cluster A by Fusaro *et al.*
274 and Monne *et al.* The second G1-like clade (clade 5, labelled sub-clade 725 in the global
275 phylogenetic tree, Fusaro/Monne Cluster D) continued to be found in Iran until 2009. This has no
276 equivalent in the Israeli trees, which seem to have inherited a G1-like lineage which originated in
277 the Indian sub-continent and then circulated in the Arabian peninsular (see the global
278 phylogenetic tree, supplementary figure1 and figure 3, Fusaro and Monne cluster B) before being
279 found in Israel in 2006/2007 (clade 3 in the nucleotide and protein trees) and in Iran in 2009
280 (clade 4 in the protein and nucleotide trees). A recent paper has shown that there have in fact been
281 several introductions of H9N2 to Israel from Jordan, and the Arabian Peninsula (Davidson et al.
282 2014). This newly introduced G1-like sub-lineage seems to have had a selective advantage over
283 the existing viral G1-like sub-lineages but this can only be confirmed by future sampling.
284 Subsequently this new sub-lineage has also spread to Egypt from Israel.

285 **Clade Analysis of the Host Specific Subsets**

286 The interesting host species subsets are those for ducks, quail and swine. Ducks are important as
287 a possible carrier of the virus between geographical regions and in acting as a reservoir species.
288 Quail have also been associated with acting as a host to allow re-assortment and viral evolution.
289 Finally swine are important because of the over-lapping glycosylation site specificity of swine
290 and human viruses, which suggests that they may act as an intermediate for transmission to
291 humans (Guo et al. 2005).

292 There are problems in interpreting the quail trees because of the absence of the Shantou
293 sequences from the nucleotide tree because of partial sequencing (figure 8). Shantou was the
294 main geographical location for the outbreak of H9N2 in quail between 2000 and 2005 (Xu et al.
295 2007a; Xu et al. 2007b). Live quail were banned from Chinese wet markets after a study had

296 shown the link to the virus. In the literature it was found that by 2003 this had resulted in the G1-
297 like lineage disappearing from China but the data here show that it was still present in quail in the
298 Shantou region until 2005 (Choi et al. 2004). The trees show that quail are hosts for all of the
299 main lineages clade A, clade B (G1-like lineage) and the main Chinese clade. This would support
300 the theory that quail have been the host species responsible for the diversification of the H9N2
301 sub-type, especially given that the original sequence of G-like lineage that is the prototype
302 sequence for clade B was originally found in a Hong Kong quail. (Hossain et al. 2008). This also
303 fits with the epochal model of viral evolution proposed by Wagner (Wagner 2008). There are
304 periods of neutral drift, which are punctuated by selection events. Here the selection event is the
305 formation of a new lineage within a different host species after a long period of neutral drift
306 within clade A.

307 Like quail ducks provide a host for the clade A and main Chinese clade viral lineages (figure 9).
308 There is only a single example of a clade B sequence in ducks, and so they might not be very
309 effective carriers of this virus or this might be explained by their geographical exposure to that
310 lineage (Perez et al. 2003). Ducks are of concern as a host species as they can contribute to the
311 spread of this lineage over a wider geographical region. In the global phylogenetic tree sequences
312 from ducks are often found clustered with those from quail showing that there is frequent transfer
313 between the two hosts.

314 Most of the swine cases have been within the main Chinese clade of sequences (figure 10). There
315 is a single sequence from Korea in 2004 that belongs to clade A. This is important in guiding how
316 we monitor disease outbreaks because it shows that the most geographically widely distributed
317 clade can also produce infection in pigs. The swine flu epidemic of H1N1 originated in Mexico
318 whereas China had been the main focus of surveillance, because of the frequency and previous
319 circulation of the virus. This is also true of H9N2 monitoring which is focused on Southeast

320 Asia, but shows that more widespread monitoring is important in the early detection of
321 pandemics.

322 There is no obvious geographical or temporal pattern in infection in pigs and there are multiple
323 transmissions between birds and pigs that produce a tree with many different sub-clades usually
324 with only a small number of members. This is consistent with their not being a widespread
325 circulation of the virus within pigs, which would give a larger homogeneous cluster of swine
326 viruses from different times and locations due to pig to pig transmission. Previously five amino
327 acids changes had been identified as being swine specific, but none of these were conserved
328 within the swine sequences or even within a single swine clade (Xu et al. 2004). There are a few
329 cases of the S145N mutation that has been identified to change the antibody epitope but this
330 again is sporadic (Ping et al. 2008).

331 **Identifying the Amino Acid Changes Responsible for Clade Formation**

332 The x-ray crystal structures are available for the H9 hemagglutinin and so it is possible to map
333 the amino acid changes responsible for differentiating between the different clades onto the
334 structure (Ha et al. 2001; Ha et al. 2002). The region between amino acids 128 and 275 makes up
335 the receptor subdomain. This domain is responsible for binding to the cellular membrane as part
336 of viral invasion. The stem domain is made up of the first 60 amino acids of the N-terminus and
337 the final 275 amino acids in the C-terminus of which the last 221 are cleaved by proteolysis of
338 the precursor protein at a conserved arginine to produce a second peptide chain. In between these
339 domains is the remains of a catalytic domain – the vestigial enzyme domain (Ha et al. 2002). The
340 amino acids that are specific to the four most distinct clades are given in table 1.

341 The presence of a four key amino acids has been shown to be essential for droplet transmission of
342 the virus H183, A189, E190 and L226 that correspond to residues H191, A197, E198, and L234
343 in the H9 numbering (Sorrell et al. 2009).

344 There are 17 amino acids that distinguish between Clade A and the consensus sequence, only
345 three of these are in the receptor domain and so the majority are in the stem domains. Three
346 mutations are in the enzymatic domain of which the most significant of which is the replacement
347 of a serine at residue 127 with an asparagine as this is on the boundary of the domain and this
348 creates another potential glycosylation site. The receptor domain changes are Q164H, R180E
349 and T206A. Of these the replacement of the basic arginine group by an acidic glutamic acid is the
350 most interesting, because of the change in polarity, none of the amino acid changes affected either
351 the glycosylation sites or altered the residues that were identified as key to viral droplet
352 transmission. Only T8A had been previously identified as an important change when it was
353 shown to be involved in host specificity (Perez et al. 2003)

354 There are 26 amino acid changes that distinguish Clade B. Eleven of these changes can be found
355 in the receptor subdomain. Many of the changes are between leucine, isoleucine and valine.
356 These are conservative changes that preserve the hydrophobicity of the amino acid but change the
357 steric interactions. Another significant proportion of the changes is substitutions to threonine
358 from serine or valine. Three of the serine to threonine changes occur in the receptor domain and
359 this might reflect an altered binding affinity for a larger binding partner. Position A168 had been
360 identified previously as under positive selection (Fusaro et al. 2011). This is supported by the
361 mutation to leucine in this clade. Of the key amino acid changes required for droplet infection
362 only the N191H mutation is clade specific.

363 The amino acid changes responsible for differentiating between clades gives some support to
364 previous studies that have tried to identify residues that are under selection (Fusaro et al. 2011).

365 However most of the clade specific changes have not been previously identified in the literature
366 on the evolution of antigenicity (Kaverin et al. 2004; Skehel & Wiley 2000). The amino acids
367 responsible for glycosylation are conserved throughout the phylogenetic trees, although some of
368 the amino acid changes introduce new asparagine residues, which could be new sites for
369 glycosylation (Guo et al. 2000; Zhang et al. 2004). There are a very large number of sequences
370 with the L234 substitution required for droplet infection but these are not specific to any of the
371 major clades and this shows that the mutation has arisen multiple times.

372 **Conclusions**

373 The sequence databases are growing at a hyperbolic rate, and with next generation sequencing
374 this level of growth is likely to continue for the foreseeable future. There are two challenges for
375 dealing with this data. The first is computational that requires improved algorithms and
376 implementations especially as we are moving to more computationally intensive methods of
377 analysis. The second is the quality of the data collection itself. Currently data collection is not
378 systematic and this seriously affects the reliability of the models that can be built. Sampling is
379 badly skewed to certain geographical locations, China being a prime example, while others are
380 ignored (Africa and South America). Surveillance of wild birds is particularly problematic.
381 Where there have been international efforts such as in the European Union study of the spread of
382 H5N1 in wild birds, even this was incomplete with no data from Spain and Eastern Europe, and
383 only partial data for France and Germany (Hesterberg et al. 2009).

384 There are also problems with partial and truncated sequences, as these often have to be excluded
385 from analysis. This is becoming less of a problem as sequencing methods become cheaper and so
386 complete sequences become more widely available.

387 Another area where there needs to be a significant improvement is in the quality of the sequence
388 annotations in the databases. For effective phylogeographic analysis it is important that future
389 data should be annotated with as much geographic data as is possible, this must include GPS
390 coordinates and further GIS (geographic information system) information to include habitat and
391 urbanisation measurements would be ideal (Scotch et al. 2011; Yasué et al. 2006).

392 In this paper there are clearly three different principal clades; Clade A – Wisconsin like, Clade B
393 – G1-like and the Main Chinese Cluster – Y280-like, but it is not clear when a new lineage has
394 arisen and when they are no longer “like” the prototype sequences. At the sub-lineage level the

395 assignment of cluster and clades strongly depends on the sampling of the sequences and this has
396 produced some conflicts between different assignments in the literature. Geography is a much
397 stronger determinant of lineage rather than avian host, which seems to provide a much weaker
398 barrier to spread of the virus. However there is a definite barrier between bird species and pigs as
399 hosts.

400 The clade analysis has provided insight into the functional and structural evolution of the protein.
401 There is only a limited overlap between the residues identified in this study as important in clade
402 differentiation and those identified as significant in the existing literature. There is therefore a
403 need for further investigation of the functionality of these newly identified amino acid changes.

404 **Acknowledgments**

405 AD would like to thank Dr Lorna Tinworth for her careful reading of the manuscript.

406 **References**

- 407 Abdel-Moneim AS, Afifi MA, and El-Kady MF. 2012. Isolation and mutation trend
408 analysis of influenza A virus subtype H9N2 in Egypt. *Virology* 9:173.
- 409 Banks J, Speidel E, Harris P, and Alexander D. 2000. Phylogenetic analysis of
410 influenza A viruses of H9 haemagglutinin subtype. *Avian pathology* 29:353-
411 359.
- 412 Bedford T, Rambaut A, and Pascual M. 2012. Canalization of the evolutionary
413 trajectory of the human influenza virus. *BMC biology* 10:38.
- 414 Butler D. 2012. Flu surveillance lacking. *Nature* 483:520-522.
- 415 Butt AM, Siddique S, Idrees M, and Tong Y. 2010. Avian influenza A (H9N2):
416 computational molecular analysis and phylogenetic characterization of viral
417 surface proteins isolated between 1997 and 2009 from the human population. *J*
418 *Virology* 7:319.
- 419 Butt KM, Smith GJ, Chen H, Zhang LJ, Leung YH, Xu KM, Lim W, Webster RG, Yuen KY,
420 Peiris JS, and Guan Y. 2005. Human infection with an avian H9N2 influenza A
421 virus in Hong Kong in 2003. *J Clin Microbiol* 43:5760-5767.
- 422 Cao B, Li X-W, Mao Y, Wang J, Lu H-Z, Chen Y-S, Liang Z-A, Liang L, Zhang S-J, and
423 Zhang B. 2009. Clinical features of the initial cases of 2009 pandemic
424 influenza A (H1N1) virus infection in China. *New England Journal of Medicine*
425 361:2507-2517.
- 426 Cheng VC, Chan JF, Wen X, Wu W, Que T, Chen H, Chan K, and Yuen K. 2011.
427 Infection of immunocompromised patients by avian H9N2 influenza A virus.
428 *Journal of Infection* 62:394-399.
- 429 Choi Y, Ozaki H, Webby R, Webster R, Peiris J, Poon L, Butt C, Leung Y, and Guan Y.
430 2004. Continuing evolution of H9N2 influenza viruses in Southeastern China. *J*
431 *Virology* 78:8609-8614.
- 432 Cong YL, Pu J, Liu QF, Wang S, Zhang GZ, Zhang XL, Fan WX, Brown EG, and Liu JH.
433 2007. Antigenic and genetic characterization of H9N2 swine influenza viruses
434 in China. *J Gen Virol* 88:2035-2041.
- 435 Davidson I, Fusaro A, Heidari A, Monne I, and Cattoli G. 2014. Molecular evolution of
436 H9N2 avian influenza viruses in Israel. *Virus genes* 48:457-463.
- 437 Dong G, Luo J, Zhang H, Wang C, Duan M, Deliberto TJ, Nolte DL, Ji G, and He H.
438 2011. Phylogenetic diversity and genotypical complexity of H9N2 influenza A
439 viruses revealed by genomic sequence analysis. *PloS one* 6:e17212.
- 440 Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high
441 throughput. *Nucleic acids research* 32:1792-1797.
- 442 Felsenstein J, and Felsenstein J. 2004. *Inferring phylogenies*: Sinauer Associates
443 Sunderland.
- 444 Fusaro A, Monne I, Salviato A, Valastro V, Schivo A, Amarin NM, Gonzalez C, Ismail
445 MM, Al-Ankari A-R, and Al-Blawi MH. 2011. Phylogeography and evolutionary
446 history of reassortant H9N2 viruses with potential human health implications.
447 *J Virology* 85:8413-8421.
- 448 Guo YJ, Krauss S, Senne DA, Mo IP, Lo KS, Xiong XP, Norwood M, Shortridge KF,
449 Webster RG, and Guan Y. 2000. Characterization of the pathogenicity of
450 members of the newly established H9N2 influenza virus lineages in Asia.
451 *Virology* 267:279-288.
- 452 Guo YJ, Wen LY, Zhang Y, Wan M, Guo JF, Li Z, and Shu YL. 2005. [Do pigs play a role
453 in human infection with avian influenza A H9N2 viruses]. *Zhonghua Shi Yan*
454 *He Lin Chuang Bing Du Xue Za Zhi* 19:106-109.

- 455 Ha Y, Stevens DJ, Skehel JJ, and Wiley DC. 2001. X-ray structures of H5 avian and H9
456 swine influenza virus hemagglutinins bound to avian and human receptor
457 analogs. *Proceedings of the National Academy of Sciences* 98:11181-11186.
- 458 Ha Y, Stevens DJ, Skehel JJ, and Wiley DC. 2002. H5 avian and H9 swine influenza
459 virus haemagglutinin structures: possible origin of influenza subtypes. *The*
460 *EMBO journal* 21:865-875.
- 461 Haase M, Starick E, Fereidouni S, Strebelow G, Grund C, Seeland A, Scheuner C,
462 Cieslik D, Smietanka K, and Minta Z. 2010. Possible sources and spreading
463 routes of highly pathogenic avian influenza virus subtype H5N1 infections in
464 poultry and wild birds in Central Europe in 2007 inferred through likelihood
465 analyses. *Infection, Genetics and Evolution* 10:1075-1084.
- 466 Hesterberg U, Harris K, Stroud D, Guberti V, Busani L, Pittman M, Piazza V, Cook A,
467 and Brown I. 2009. Avian influenza surveillance in wild birds in the European
468 Union in 2006. *Influenza Other Respir Viruses* 3:1-14.
- 469 Hossain MJ, Hickman D, and Perez DR. 2008. Evidence of expanded host range and
470 mammalian-associated genetic changes in a duck H9N2 influenza virus
471 following adaptation in quail and chickens. *PLoS One* 3:e3170.
- 472 Huang Y, Hu B, Wen X, Cao S, Gavrilov BK, Du Q, Khan MI, and Zhang X. 2010.
473 Diversified reassortant H9N2 avian influenza viruses in chicken flocks in
474 northern and eastern China. *Virus Res* 151:26-32.
- 475 Iqbal M, Yaqub T, Reddy K, and McCauley JW. 2009. Novel genotypes of H9N2
476 influenza A viruses isolated from poultry in Pakistan containing NS genes
477 similar to highly pathogenic H7N3 and H5N1 viruses. *PLoS One* 4:e5788.
- 478 Jones DT, Taylor WR, and Thornton JM. 1992. The rapid generation of mutation data
479 matrices from protein sequences. *Computer applications in the biosciences:*
480 *CABIOS* 8:275-282.
- 481 Kaverin NV, Rudneva IA, Ilyushina NA, Lipatov AS, Krauss S, and Webster RG. 2004.
482 Structural differences among hemagglutinins of influenza A virus subtypes are
483 reflected in their antigenic architecture: analysis of H9 escape mutants. *J Virol*
484 78:240-249.
- 485 Kim JA, Cho SH, Kim HS, and Seo SH. 2006. H9N2 influenza viruses isolated from
486 poultry in Korean live bird markets continuously evolve and cause the severe
487 clinical signs in layers. *Vet Microbiol* 118:169-176.
- 488 Koelle K, Cobey S, Grenfell B, and Pascual M. 2006. Epochal evolution shapes the
489 phylodynamics of interpandemic influenza A (H3N2) in humans. *Science*
490 314:1898-1903.
- 491 LAM TTY, HON CC, Lemey P, Pybus OG, Shi M, Tun HM, Li J, Jiang J, Holmes EC, and
492 LEUNG FCC. 2012. Phylodynamics of H5N1 avian influenza virus in Indonesia.
493 *Molecular ecology* 21:3062-3077.
- 494 Lee DH, and Song CS. 2013. H9N2 avian influenza virus in Korea: evolution and
495 vaccination. *Clin Exp Vaccine Res* 2:26-33.
- 496 Li C, Yu K, Tian G, Yu D, Liu L, Jing B, Ping J, and Chen H. 2005. Evolution of H9N2
497 influenza viruses from domestic poultry in Mainland China. *Virology* 340:70-
498 83.
- 499 Li KS, Xu KM, Peiris JS, Poon LL, Yu KZ, Yuen KY, Shortridge KF, Webster RG, and
500 Guan Y. 2003. Characterization of H9 subtype influenza viruses from the ducks
501 of southern China: a candidate for the next influenza pandemic in humans? *J*
502 *Virol* 77:6988-6994.
- 503 Lin YP, Shaw M, Gregory V, Cameron K, Lim W, Klimov A, Subbarao K, Guan Y, Krauss
504 S, Shortridge K, Webster R, Cox N, and Hay A. 2000. Avian-to-human
505 transmission of H9N2 subtype influenza A viruses: relationship between H9N2
506 and H5N1 human isolates. *Proc Natl Acad Sci U S A* 97:9654-9658.

- 507 Lv J, Wei B, Yang Y, Yao M, Cai Y, Gao Y, Xia X, Zhao X, Liu Z, Li X, Wang H, Yang H,
508 Roesler U, Miao Z, and Chai T. 2012. Experimental transmission in guinea pigs
509 of H9N2 avian influenza viruses from indoor air of chicken houses. *Virus Res*
510 170:102-108.
- 511 Marshall N, Priyamvada L, Ende Z, Steel J, and Lowen AC. 2013. Influenza virus
512 reassortment occurs with high frequency in the absence of segment
513 mismatch. *PLoS pathogens* 9:e1003421.
- 514 Monne I, Hussein HA, Fusaro A, Valastro V, Hamoud MM, Khalefa RA, Dardir SN,
515 Radwan MI, Capua I, and Cattoli G. 2013. H9N2 influenza A virus circulates in
516 H5N1 endemically infected poultry population in Egypt. *Influenza Other Respir*
517 *Viruses* 7:240-243.
- 518 Nili H, and Asasi K. 2002. Natural cases and an experimental study of H9N2 avian
519 influenza in commercial broiler chickens of Iran. *Avian Pathol* 31:247-252.
- 520 Page RD, and Holmes EC. 2009. *Molecular evolution: a phylogenetic approach*: John
521 Wiley & Sons.
- 522 Park KJ, Kwon H-i, Song M-S, Pascua PNQ, Baek YH, Lee JH, Jang H-L, Lim J-Y, Mo I-P,
523 and Moon H-J. 2011. Rapid evolution of low-pathogenic H9N2 avian influenza
524 viruses following poultry vaccination programmes. *Journal of General Virology*
525 92:36-50.
- 526 Pawar SD, Tandale BV, Raut CG, Parkhi SS, Barde TD, Gurav YK, Kode SS, and Mishra
527 AC. 2012. Avian influenza H9N2 seroprevalence among poultry workers in
528 Pune, India, 2010. *PLoS One* 7:e36374.
- 529 Peiris JS, Guan Y, Markwell D, Ghose P, Webster RG, and Shortridge KF. 2001.
530 Cocirculation of avian H9N2 and contemporary "human" H3N2 influenza A
531 viruses in pigs in southeastern China: potential for genetic reassortment? *J*
532 *Virology* 75:9679-9686.
- 533 Peiris M, Yuen KY, Leung CW, Chan KH, Ip PL, Lai RW, Orr WK, and Shortridge KF.
534 1999. Human infection with influenza H9N2. *Lancet* 354:916-917.
- 535 Perez D, Webby R, Hoffmann E, and Webster R. 2003. Land-based birds as potential
536 disseminators of avian/mammalian reassortant influenza A viruses. *Avian Dis*
537 47:1114-1117.
- 538 Perk S, Banet-Noach C, Shihmanter E, Pokamunski S, Pirak M, Lipkind M, and Panshin
539 A. 2006. Genetic characterization of the H9N2 influenza viruses circulated in
540 the poultry population in Israel. *Comp Immunol Microbiol Infect Dis* 29:207-
541 223.
- 542 Ping J, Li C, Deng G, Jiang Y, Tian G, Zhang S, Bu Z, and Chen H. 2008. Single-amino-
543 acid mutation in the HA alters the recognition of H9N2 influenza virus by a
544 monoclonal antibody. *Biochem Biophys Res Commun* 371:168-171.
- 545 Scotch M, Sarkar IN, Mei C, Leaman R, Cheung K-H, Ortiz P, Singraur A, and Gonzalez
546 G. 2011. Enhancing phylogeography by improving geographical information
547 from GenBank. *Journal of biomedical informatics* 44:S44-S47.
- 548 Skehel JJ, and Wiley DC. 2000. Receptor binding and membrane fusion in virus entry:
549 the influenza hemagglutinin. *Annual review of biochemistry* 69:531-569.
- 550 Smith DJ, Lapedes AS, de Jong JC, Bestebroer TM, Rimmelzwaan GF, Osterhaus AD,
551 and Fouchier RA. 2004. Mapping the antigenic and genetic evolution of
552 influenza virus. *Science* 305:371-376.
- 553 Song Xf, Han P, and Chen YPP. 2011. Genetic variation of the hemagglutinin of avian
554 influenza virus H9N2. *J Med Virol* 83:838-846.
- 555 Sorrell EM, Wan H, Araya Y, Song H, and Perez DR. 2009. Minimal molecular
556 constraints for respiratory droplet transmission of an avian-human H9N2
557 influenza A virus. *Proc Natl Acad Sci U S A* 106:7565-7570.
- 558 Sun Y, Qin K, Wang J, Pu J, Tang Q, Hu Y, Bi Y, Zhao X, Yang H, and Shu Y. 2011. High
559 genetic compatibility and increased pathogenicity of reassortants derived

- 560 from avian H9N2 and pandemic H1N1/2009 influenza viruses. *Proceedings of*
561 *the National Academy of Sciences* 108:4164-4169.
- 562 Tamura K, and Nei M. 1993. Estimation of the number of nucleotide substitutions in
563 the control region of mitochondrial DNA in humans and chimpanzees. *Mol Biol*
564 *Evol* 10:512-526.
- 565 Tamura K, Peterson D, Peterson N, Stecher G, Nei M, and Kumar S. 2011. MEGA5:
566 molecular evolutionary genetics analysis using maximum likelihood,
567 evolutionary distance, and maximum parsimony methods. *Mol Biol Evol*
568 28:2731-2739.
- 569 Thomas PG, and Hertz T. 2012. Constrained evolution drives limited influenza
570 diversity. *BMC biology* 10:43.
- 571 Wagner A. 2008. Neutralism and selectionism: a network-based reconciliation.
572 *Nature Reviews Genetics* 9:965-974.
- 573 Wallace RG, and Fitch WM. 2008. Influenza A H5N1 immigration is filtered out at
574 some international borders. *PLoS One* 3:e1697.
- 575 Wallace RG, HoDac H, Lathrop RH, and Fitch WM. 2007. A statistical phylogeography
576 of influenza A H5N1. *Proceedings of the National Academy of Sciences*
577 104:4473-4478.
- 578 Wan H, Sorrell EM, Song H, Hossain MJ, Ramirez-Nieto G, Monne I, Stevens J, Cattoli
579 G, Capua I, Chen LM, Donis RO, Busch J, Paulson JC, Brockwell C, Webby R,
580 Blanco J, Al-Natour MQ, and Perez DR. 2008. Replication and transmission of
581 H9N2 influenza viruses in ferrets: evaluation of pandemic potential. *PLoS One*
582 3:e2923.
- 583 Wang M, Fu CX, and Zheng BJ. 2009. Antibodies against H5 and H9 avian influenza
584 among poultry workers in China. *N Engl J Med* 360:2583-2584.
- 585 Watanabe T, Kiso M, Fukuyama S, Nakajima N, Imai M, Yamada S, Murakami S,
586 Yamayoshi S, Iwatsuki-Horimoto K, and Sakoda Y. 2013. Characterization of
587 H7N9 influenza A viruses isolated from humans. *Nature* 501:551-555.
- 588 Wolf YI, Viboud C, Holmes EC, Koonin EV, and Lipman DJ. 2006. Long intervals of
589 stasis punctuated by bursts of positive selection in the seasonal evolution of
590 influenza A virus. *Biol Direct* 1:357-360.
- 591 Xu C, Fan W, Wei R, and Zhao H. 2004. Isolation and identification of swine influenza
592 recombinant A/Swine/Shandong/1/2003(H9N2) virus. *Microbes Infect* 6:919-
593 925.
- 594 Xu KM, Li KS, Smith GJ, Li JW, Tai H, Zhang JX, Webster RG, Peiris JS, Chen H, and
595 Guan Y. 2007a. Evolution and molecular epidemiology of H9N2 influenza A
596 viruses from quail in southern China, 2000 to 2005. *J Virol* 81:2635-2645.
- 597 Xu KM, Smith GJ, Bahl J, Duan L, Tai H, Vijaykrishna D, Wang J, Zhang JX, Li KS, Fan
598 XH, Webster RG, Chen H, Peiris JS, and Guan Y. 2007b. The genesis and
599 evolution of H9N2 influenza viruses in poultry from southern China, 2000 to
600 2005. *J Virol* 81:10389-10401.
- 601 Yasué M, Feare CJ, Bennun L, and Fiedler W. 2006. The epidemiology of H5N1 avian
602 influenza in wild birds: why we need better ecological data. *BioScience*
603 56:923-929.
- 604 Yu H, Hua RH, Wei TC, Zhou YJ, Tian ZJ, Li GX, Liu TQ, and Tong GZ. 2008. Isolation
605 and genetic characterization of avian origin H9N2 influenza viruses from pigs
606 in China. *Vet Microbiol* 131:82-92.
- 607 Yu H, Zhou Y-J, Li G-X, Ma J-H, Yan L-P, Wang B, Yang F-R, Huang M, and Tong G-Z.
608 2011. Genetic diversity of H9N2 influenza viruses from pigs in China: a
609 potential threat to human health? *Vet Microbiol* 149:254-261.
- 610 Zhang M, Gaschen B, Blay W, Foley B, Haigwood N, Kuiken C, and Korber B. 2004.
611 Tracking global patterns of N-linked glycosylation site variation in highly

612 variable viral glycoproteins: HIV, SIV, and HCV envelopes and influenza
613 hemagglutinin. *Glycobiology* 14:1229-1246.
614 Zhang P, Tang Y, Liu X, Liu W, Zhang X, Liu H, Peng D, Gao S, Wu Y, Zhang L, Lu S,
615 and Liu X. 2009. A novel genotype H9N2 influenza virus possessing human
616 H5N1 internal genomes has been circulating in poultry in eastern China since
617 1998. *J Virol* 83:8428-8438.

Figure 1

Phylogenetic Overview

A compressed view of the complete phylogenetic tree with the major clades shown in a condensed format.

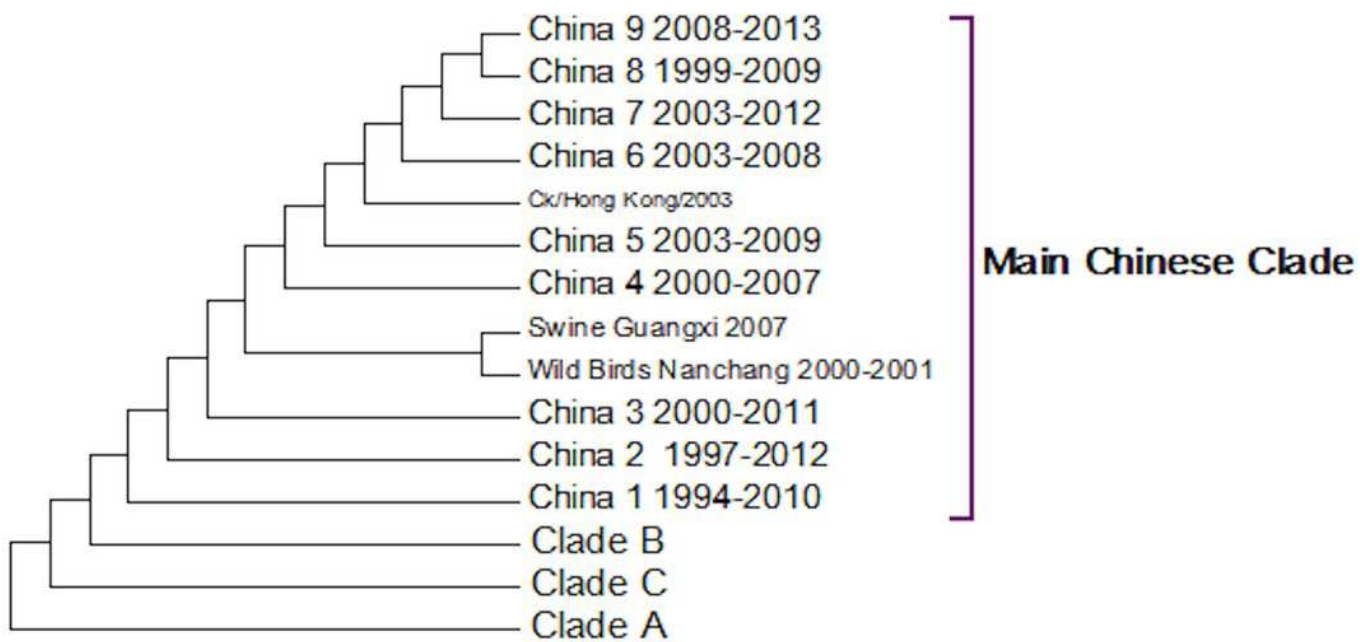


Figure 2

Clade A

This is the original lineage of H9N2 that was first isolated in a turkey in Wisconsin in 1966.

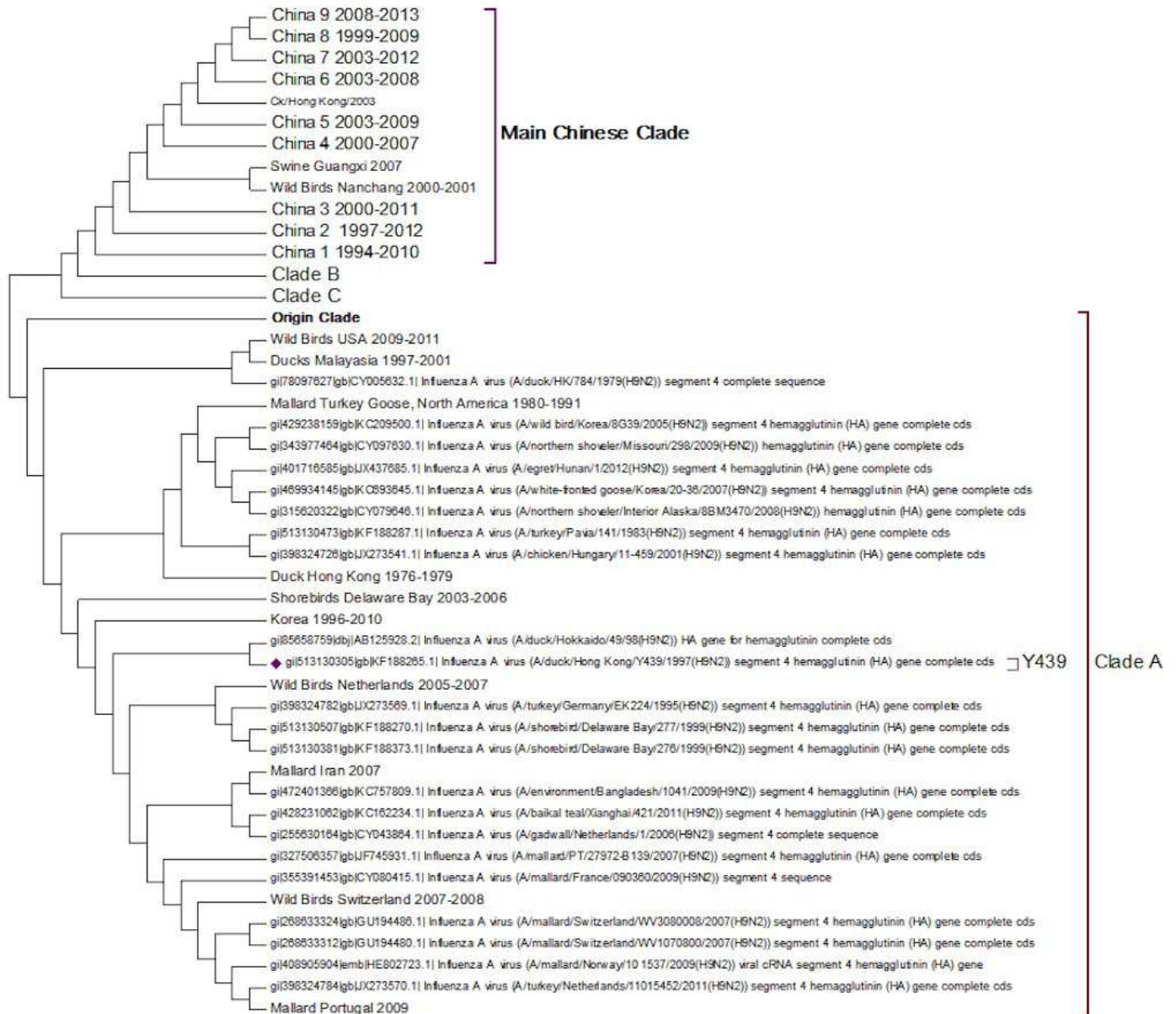


Figure 3

Clade B

This is also known as the G1 lineage.

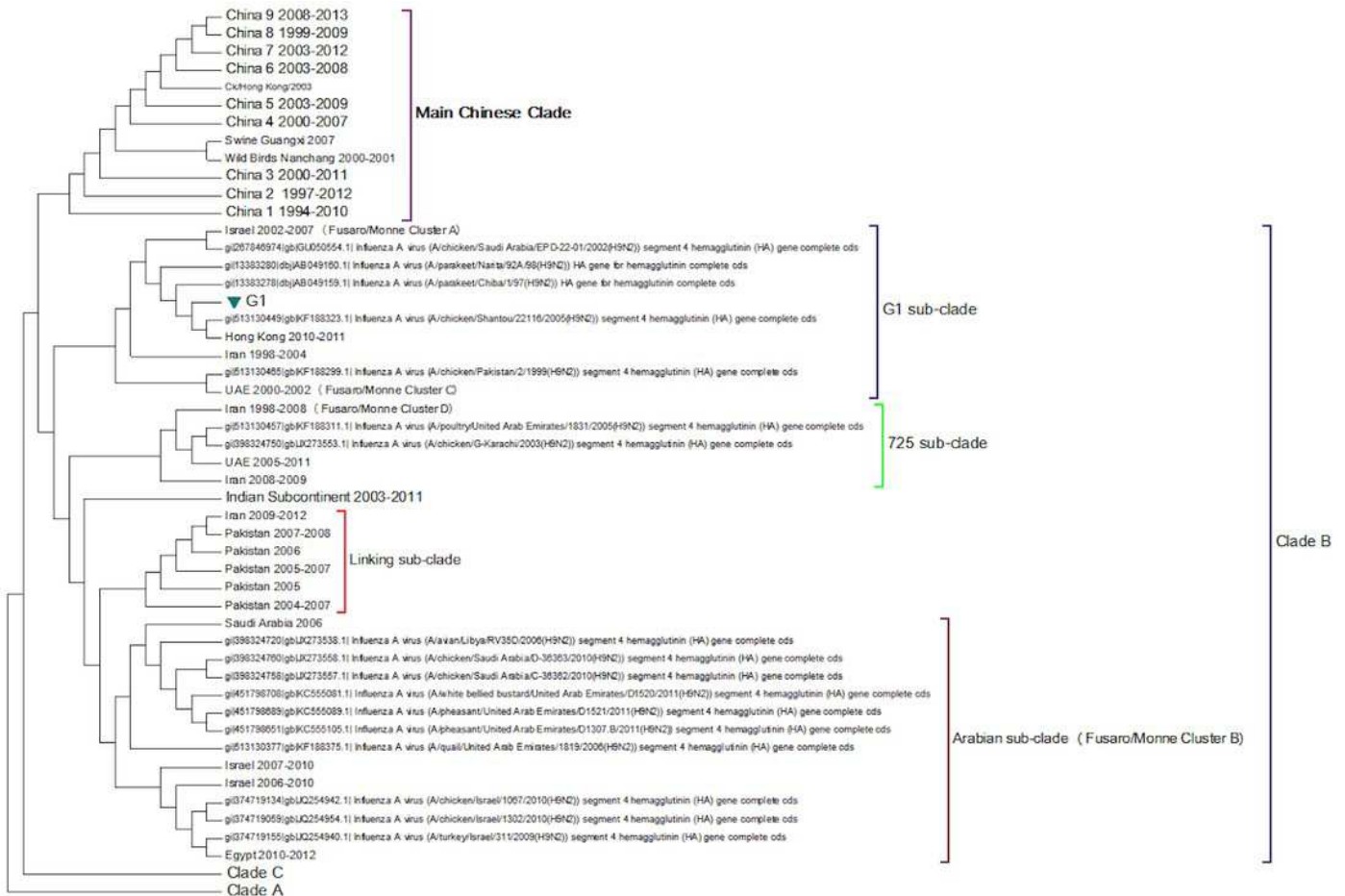


Figure 4

Clade C

This is a small clade found between clades A and B

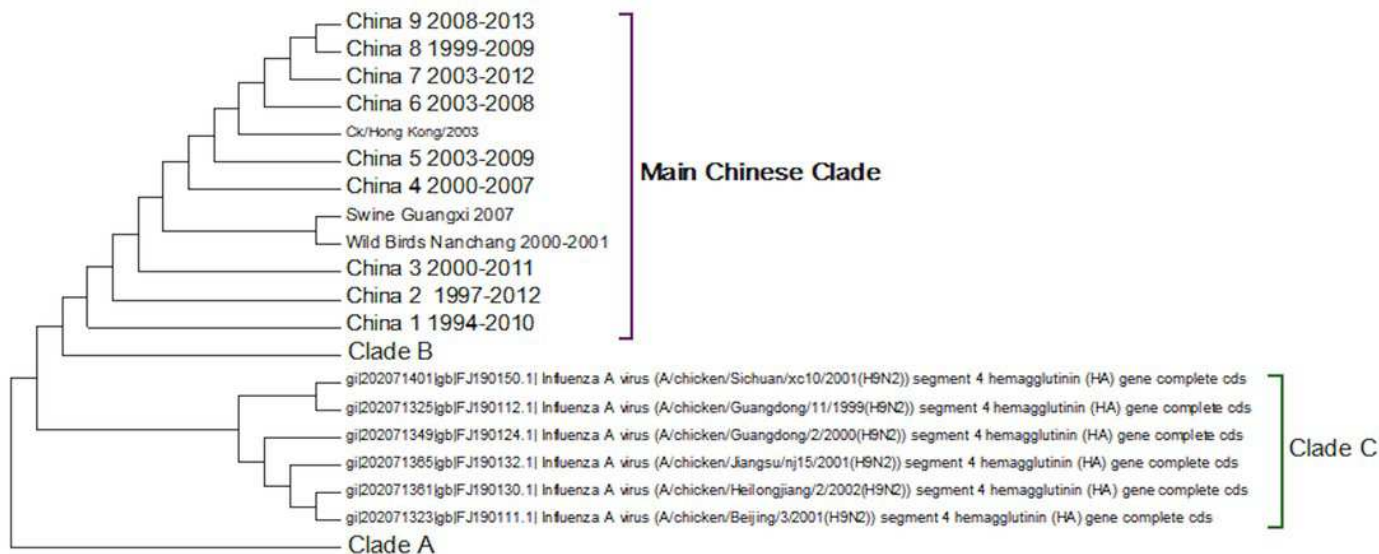
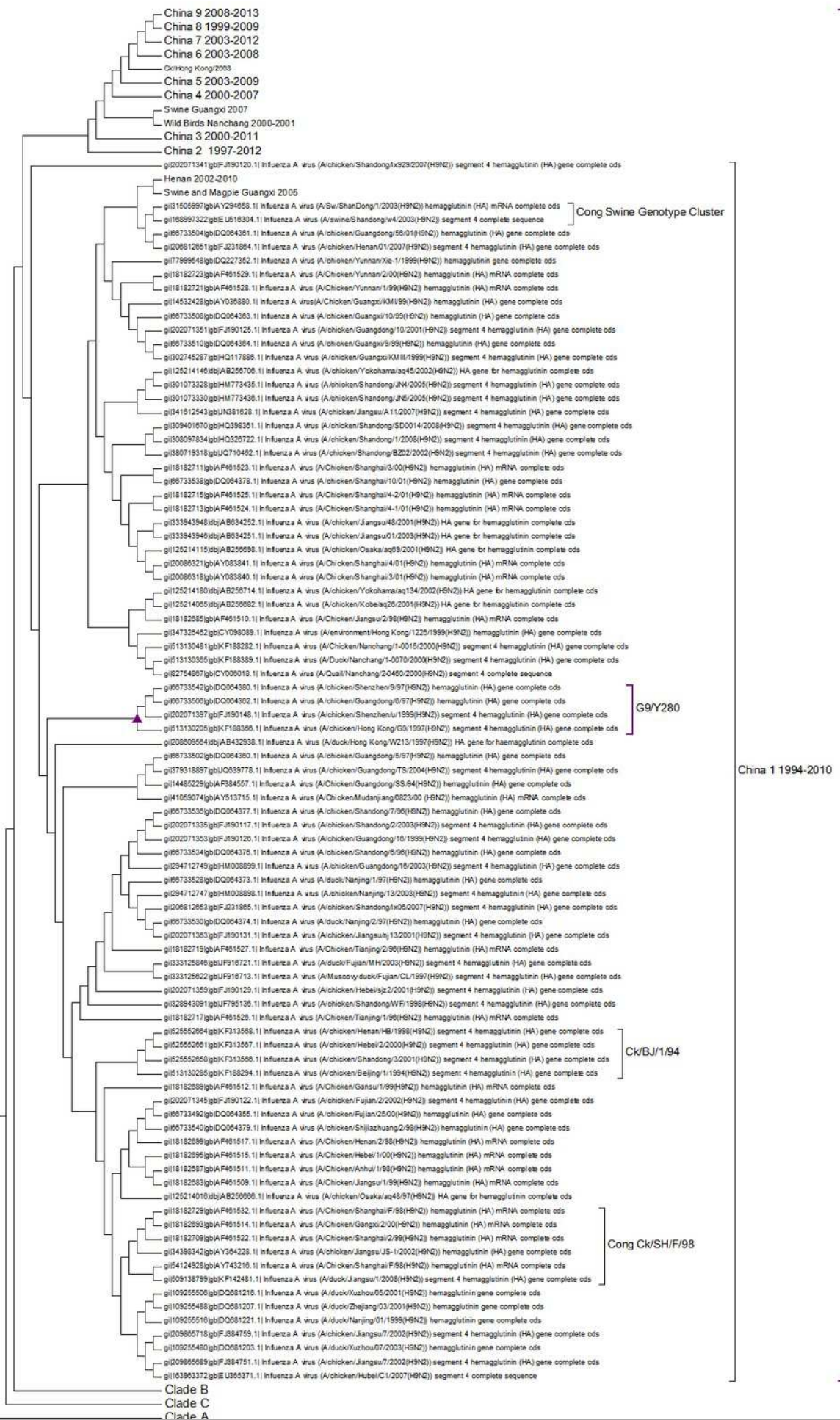


Figure 5

The main Chinese Clade

This contains the Y280 lineage but this has divided into a series of deeply rooted nested clades.



Cong Swine Genotype Cluster

G9/Y280

China 1 1994-2010

Main Chinese Clade

Ck/BJ/194

Cong Ck/SH/F/98

Clade B
Clade C
Clade A

Figure 6

The Korean nucleotide and amino acid phylogenetic trees.

A) The nucleotide phylogenetic tree. B) The amino acid phylogenetic tree.

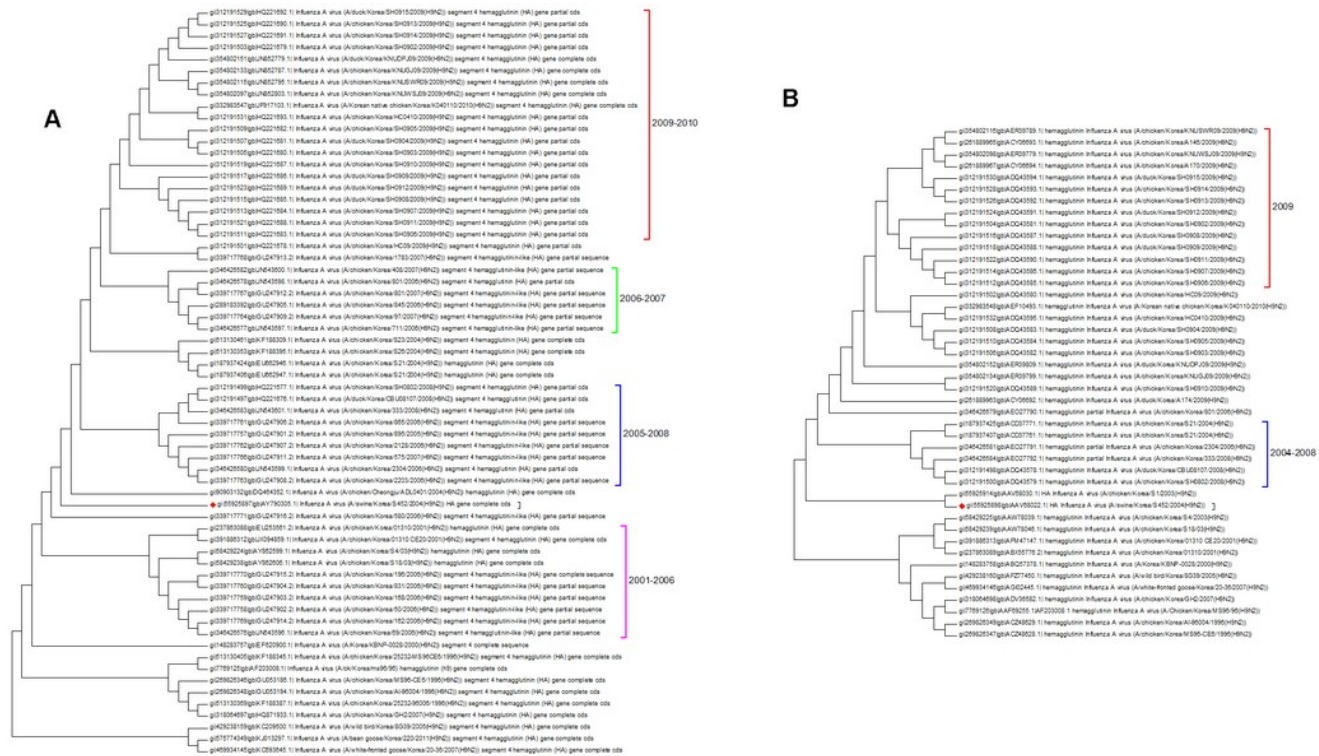


Figure 7

The Iranian and Israeli nucleotide and amino acid phylogenetic trees.

A) The Iranian nucleotide phylogenetic tree, B) The Iranian amino acid phylogenetic tree, C) The Israeli nucleotide phylogenetic tree, D) The Israeli amino acid phylogenetic tree.

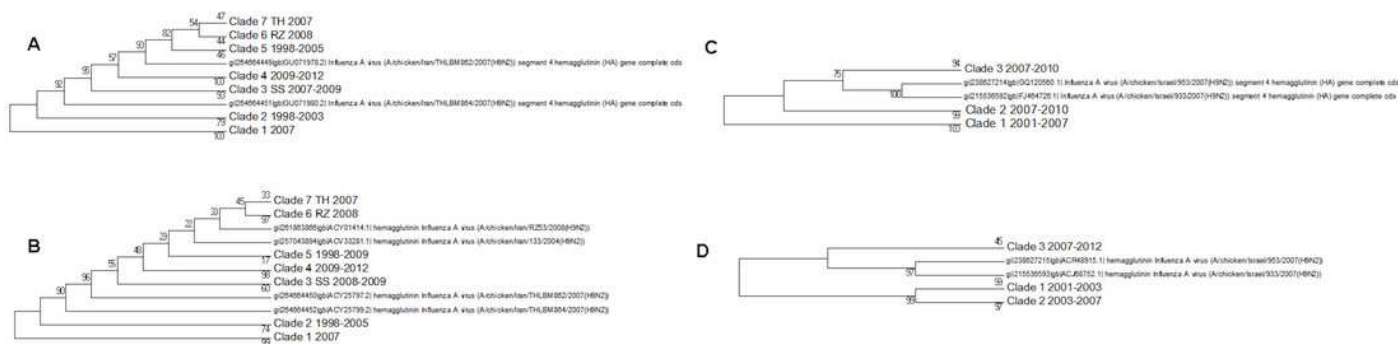


Figure 8

The qual nucleotide and amino acid phylogenetic trees

A) The nucleotide phylogenetic tree, B) The amino acid phylogenetic tree. Numbers on the internal branches are the bootstrap values.



Figure 9

The duck nucleotide phylogenetic tree

Numbers on the internal branches are the bootstrap values.

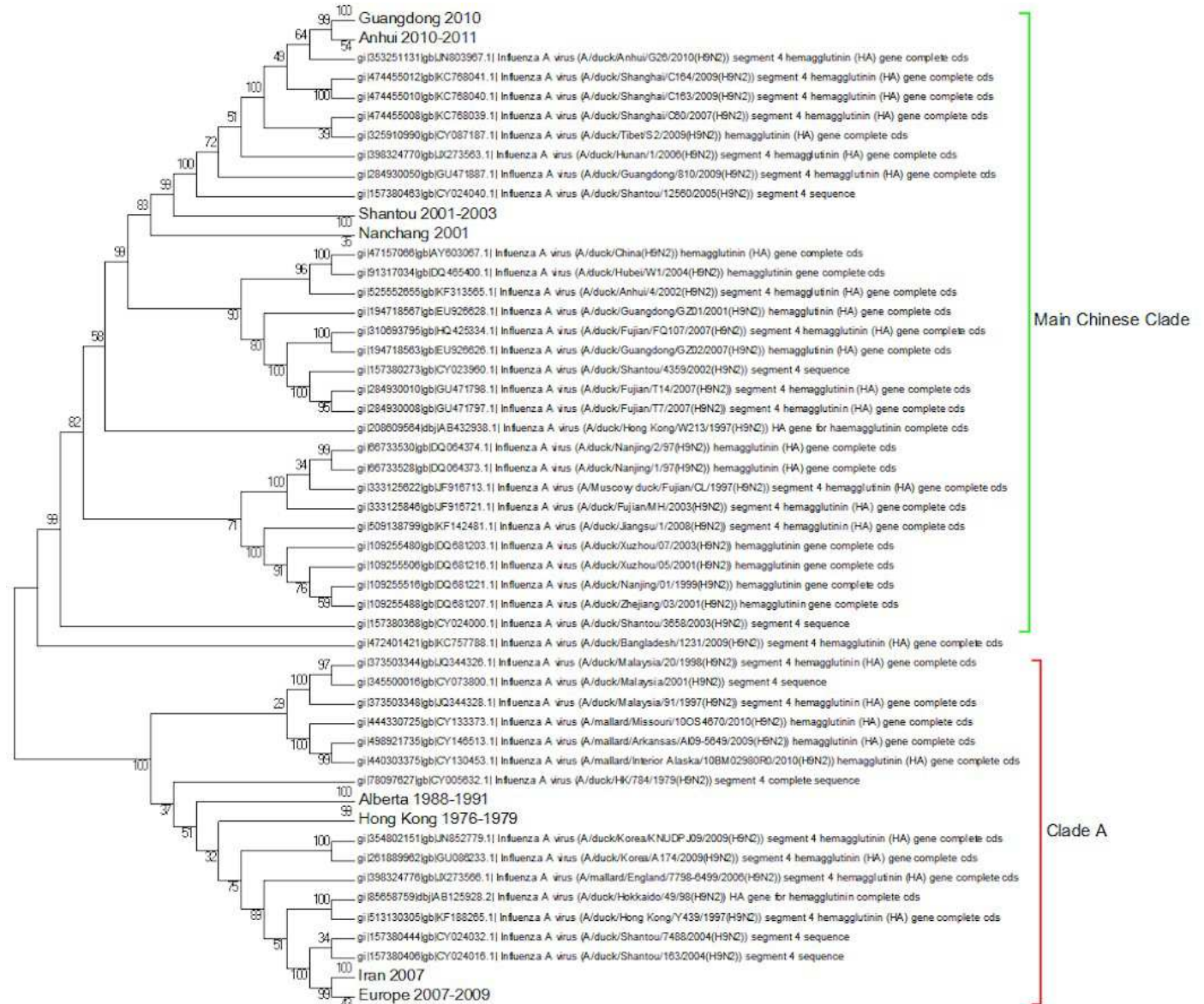


Figure 10

The swine amino acid phylogenetic tree.

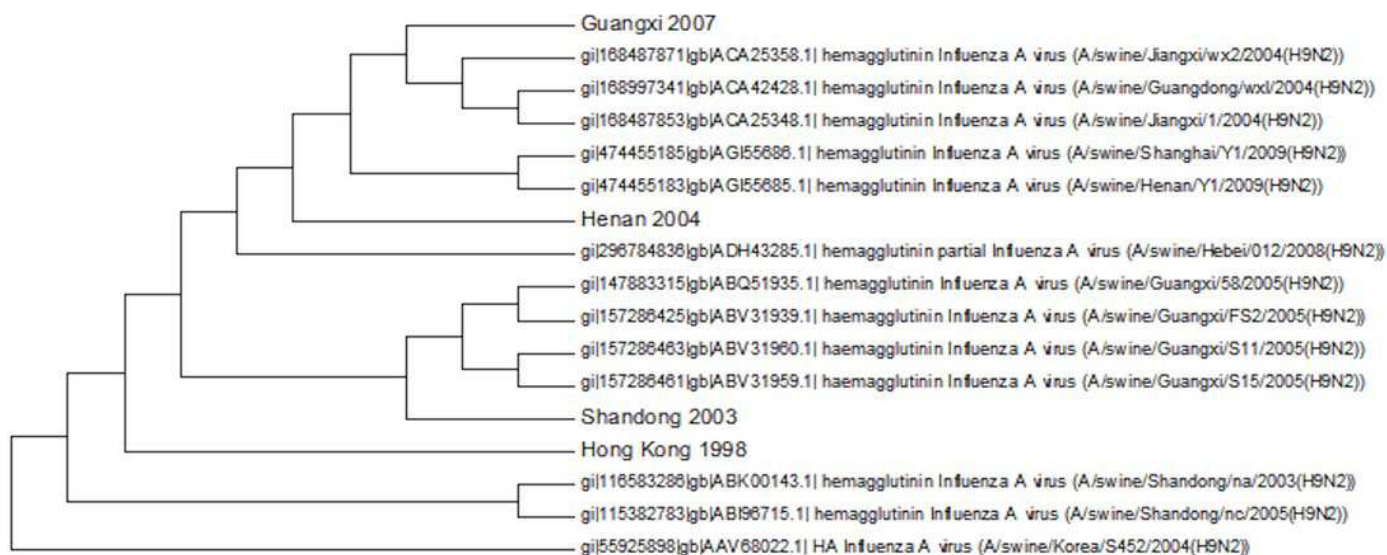


Table 1 (on next page)

The amino acid differences between clades A and B with respect to the main Chinese clade.

Numbering is for the complete H9 hemagglutinin protein. Changes are from the conserved sequence (main Chinese clade) to the clade specific sequence.

Numbering is for the complete H9 hemagglutinin protein. Changes are from the conserved sequence (main Chinese clade) to the clade specific sequence.

Clade A	Clade B
V4 -> T	V3 -> T
T8 -> A	V15 ->T or A
T38 -> I	N40 -> T
A47 -> T	L107 -> T
I79 -> V	S121 -> T
R92 -> K	S143-> T
L122 -> F	S147 -> T
S127 -> N	S158 -> N
Q164 -> H	N167 -> G
R180 -> E	A168 -> L
T206 -> A	M187 -> V
R335 -> A	N191 -> H
S337 -> D	T204 -> I
K381 -> E	R205 -> N
V394 -> I	I217 -> L
K473 -> N	D239 -> N
N496 -> D	R294 -> K
	T299 -> S
	V306 -> I
	N313 -> T
	V318 -> I
	V333 -> I
	S353-> P
	I429 -> V
	V469 -> M
	I537 -> L

