



Tracking the Geographic Spread of Avian Influenza (H5N1) with Multiple Phylogenetic Trees

Journal:	<i>Cladistics</i>
Manuscript ID:	CLA-09-05-0268.R1
Manuscript Type:	Article
Keywords:	Biogeography & Applications, Phylogeography & Biogeography & Applications, Cladistics, Molecular systematics & Genomics, Microbiology, Taxon sampling & Methodology

 scholar**ONE**TM
Manuscript Central

Tracking the Geographic Spread of Avian Influenza (H5N1) with Multiple Phylogenetic Trees

Rasmus Hovmöller ^{1,2}, Boyan Alexandrov ², Jori Hardman ², and Daniel Janies ²

1 Mathematical Biosciences Institute, The Ohio State University, Columbus, OH, USA
2 Department of Biomedical Informatics, The Ohio State University, Columbus, OH,
USA.

For Peer Review

Abstract	3
Introduction.....	5
Influenza viruses and pathogenic H5N1	5
Origin of H5N1	6
Phylogeographic analysis of H5N1	7
Materials and methods	9
Current datasets.....	9
Wallace and Fitch datasets.....	11
Phylogenetic analysis.....	12
Character optimization.....	13
Counting transmission events	15
Results.....	17
Alignment and phylogenetic trees	17
Reanalysis of Wallace and Fitch (2008)	18
Frequency of discovery of transmission routes	19
Visualization	20
Discussion	20
Phylogenetic analysis and H5N1	20
Character optimization in phylogeographic analysis	21
Current state of the H5N1 epidemic	25
Perils of using a single tree	29
Problems with using multiple trees.....	30
Conclusions.....	31
Acknowledgements.....	33
Literature cited.....	34
Figure captions.....	44
Figure 1:.....	44
Figure 2:	44
Figure 3:	45

Abstract

Avian influenza (H5N1) has been of great social and economic importance since it first infected humans in Hong Kong in 1997. A highly pathogenic strain has spread from China and has killed humans in east Asia, west Africa, south Asia, and the Middle East. Recently, several molecular phylogenetic studies have focused on the relationships of various clades of H5N1 and their spread over time, space, and various hosts. These studies examining the geographic spread of H5N1 have based their conclusions on a single tree. This tree often results from the analysis of the genomic segment coding for hemagglutinin (HA) or neuraminidase (NA) proteins and a limited sample of viral isolates. Here we present the first study using multiple candidate trees to estimate geographic transmission routes of H5N1. In addition, we use all high quality HA and NA sequences available to the public as of June 2008. We estimated geographic transmission routes of H5N1 by optimizing multistate characters with states representing different geographical regions over a pool of presumed minimum length trees. We also developed means to visualize our results in Keyhole Markup Language (KML) for virtual globes. We provide these methods as a web application entitled "Routemap" (<http://routemap.osu.edu>). The resulting visualizations are akin to airline route maps but they depict the routes of spread of viral lineages. We compare our results with the results of previous studies. We focus on the sensitivity of results to sampling of tree space, character coding schemes, optimization methods, and taxon sampling. In conclusion, we

1
2
3 find that using one tree and a single character optimization method will ignore many of
4
5 the transmission routes indicated by genetic sequence and geographic data.
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Introduction

Influenza viruses and pathogenic H5N1

Influenza viruses cause disease in humans, wild and domesticated birds and other homeothermic animals. The current classification divides influenza into three major groups: A, B, and C. Within influenza A, lineages are classified with a system of antibodies to the surface proteins of the virions, hemagglutinin and neuraminidase. This leads to the HXNY nomenclature where X includes 16 antigenic subtypes and Y represents nine different antigenic subtypes. Influenza B is only known to infect humans and pinnipeds (seals, sea lions and walruses). Influenza C is only known to infect humans and swine. On the other hand, influenza A, including subtype H5N1, has a much wider host range. Influenza A infects aves as well as many mammalian groups such as artiodactyl (swine and bovids), equids (horses), carnivorans (canids, domestic and wild felids, pinnipeds, mustelids, and viverrids). H5N1 has also been found in blow flies (Diptera: Sarcophagidae) in the vicinity of farms containing infected poultry in Japan (Sawabe et al, 2004).

Influenza viruses have two means of creating genetic diversity. The segmented genome of influenza viruses enables genomic reassortment between strains. When a host is simultaneously infected by more than one strain of influenza virions that carry gene segments from two or more original viruses, termed "reassortants" can be produced. In addition, influenza viruses are negative strand-RNA viruses. These viruses use RNA

1
2
3 polymerases without proofreading functions, resulting in a high mutation rate (Holland et
4
5 al., 1982). Both of these molecular evolutionary processes allow many lineages of
6
7 influenza to escape host immune responses (Smith et al., 2004) and antiviral therapeutics
8
9 (Bright, et al., 2005; Hill et al., 2009).

16 The origins of H5N1 and its variable pathogenicity

17
18
19
20
21 The RNA for hemagglutinin (HA) encodes a polyprotein with two subunits that must be
22 cleaved by host endoproteases to enable the virion to fuse with the host's cell membrane.
23
24 Highly pathogenic strains of H5N1 are characterized by multiple basic amino acids at the
25 cleavage site between the two subunits of the hemagglutinin protein (Kawaoka, 1987;
26 Subbarao, 1998). Pathogenicity is a multigenic trait that is also mediated by features such
27
28 as whether the polymerase basic 2 F2 protein is produced or not (Zamarin et al., 2006)
29
30 The earliest known influenza strain characterized as serotype H5N1 was isolated from
31 birds in Scotland in 1959 (WHO, 2004) and sequenced by De et al., in 1988. Descendents
32 carrying the HA of this strain have been found in wild birds in North America since the
33
34 1970s (Matrosovich et al., 1999; Zhou et al., 1999; Obenauer et al., 2006). However, thus
35 far these isolates of H5 from North America appear to be strains with low pathogenicity
36
37 that infect only birds (USDA, 2007).

38
39
40
41
42 A separate highly pathogenic lineage of H5N1 has caused significant international
43 concern due to their spread across Eurasia and Africa. This pathogenic lineage can be
44 traced back to a 1996 isolate from a goose in China's Guangdong province (Xu et al.,
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 1999). As there was little monitoring of H5N1 in Asia before the human cases of H5N1
4
5 in Hong Kong in 1997, the deep origins of the pathogenic lineages of H5N1 remains
6
7 unclear. Further complicating the deep history of H5N1, there are few isolates of H5N1
8 collected in the period 1960 to 1995 that have been fully sequenced¹ and put into the
9 public domain.
10
11
12
13
14
15
16

17 When an outbreak of H5N1 in 1997 in live bird markets in Hong Kong infected 18
18 people, six of whom died, this was the first indication that H5N1 had not only become
19 very pathogenic in birds but also that it has achieved the ability to infect humans (Xu et
20 al., 1999; Chan, 2002). Since 2003, there have been 438 documented human cases of
21 H5N1 and 60% have died (WHO, 2009). To date, the H5N1 virus has killed people in the
22 Middle East, west Africa, and throughout Asia.
23
24
25
26
27
28
29
30
31
32

33 Phylogeographic analysis of H5N1 34 35 36 37

38 In most phylogeographic studies of H5N1 (e.g., Ducez et al., 2007; Guan et al., 2002;
39 Wallace et al., 2007; Wallace and Fitch, 2008; Wang et al., 2008), transmission routes²
40 have typically been determined by visual inspection of a single phylogenetic tree.
41 Kilpatrick et al., (2006) used a single tree in conjunction with data from observations on
42 migration routes of wild birds and international poultry trade to estimate opportunities for
43 the virus to spread across various borders. Smith et al., (2006) and Salzberg et al., (2007)
44 optimized geographic regions on a phylogenetic tree. These analyses have the advantage
45
46
47
48
49
50
51
52
53
54

55 ¹ Most sequences from this period are of low quality and were not included in this study.
56
57 ² With the phrases "transmission route" and "transmission event", we refer to a viral
58 lineage crossing a political border in a character evolution study.
59
60

1
2
3 that the nodes of the tree are assigned ancestral states, enabling areas of origin to be
4
5 estimated. Janies et al., (2007) and Hill et al., (2009) optimized other features such host
6
7 shifts and key mutations onto a phylogenetic tree for H5N1 and projected the tree into a
8
9 virtual globe. Using virtual globes also allows for animation over time to reconstruct the
10
11 geographic spread of viral lineages. All of these studies present a single tree per genetic
12
13 segment or whole genome for optimization and visualization (Smith et al., 2006; Salzberg
14
15 et al., 2007; and Janies et al., 2007; 2008; Hill et al., 2009).
16
17
18
19
20
21

22 In this paper, we compare our results based on multiple trees and large datasets to the
23
24 work of Wallace et al., (2007) and Wallace and Fitch, (2008). Wallace et al. (2007)
25
26 analyzed 192 HA sequences and Wallace and Fitch (2008) analyzed 482 HA and 430 NA
27
28 sequences. These authors used a variety of optimality criteria for tree search. This
29
30 research group chose to present a single parsimony tree for the 2007 paper and a single
31
32 maximum likelihood tree for the 2008 paper.
33
34
35
36
37

38 Here we examine whether the analysis of a single tree in phylogeographic studies of
39 avian influenza (H5N1) tends to underestimate possible routes of transmission of the
40
41 virus across geographical borders. We also update the HA and NA datasets for H5N1
42
43 with recently released sequence data and address the implication for regional spread of
44
45 influenza.
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 We use combined efficient heuristic search strategies with computing clusters to find
4
5 large sets a set of presumed most parsimonious trees (MPT). We also explore various
6
7 coding schemes for political and geographic boundaries as characters.
8
9

10
11
12 We developed a novel method for visualization of transmission routes of diseases. The
13
14 "Supramap" method of Janies et al., (2007) can produce very complex visualizations,
15 especially for large datasets, and only displays a single tree. In order to address these
16 issues we developed a technique, termed "Routemap", that allows the user to 1) display
17 only the routes of interest and 2) summarize transmission between regions implied by one
18 to many candidate trees. The essential difference between Supramap and Routemap is
19 that Supramap projects one tree whereas Routemap distills information from a pool of
20 trees. Moreover Routemap presents the results of character change studies not as a
21 projected tree but rather as an interactive map with lines representing viral transmission
22 routes between localities. We distribute the visualization method as an easy to use web
23 application called "Routemap" (<http://routemap.osu.edu>). In doing so, we aim that the
24 application will be of general utility in biomedical and natural science applications of
25 phylogeography.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

Materials and methods

Current datasets

50
51
52
53 This study is based on all high quality nucleotide sequences for HA and NA genes from
54 H5N1 available in The National Institutes of Health's GenBank
55
56
57
58
59
60

(<http://www.ncbi.nlm.nih.gov>) as of June 18, 2008. We add a single sequence made available by via two labs³ participating in the Global Initiative for Sharing of All Influenza Data (GISAID, <http://www.gisaid.org>).

The following quality criteria were used to determine sequence quality: Sequences must be >75% complete, must not contain frame breaking insertions or deletions, and must not have a stop codon except at the end. Sequences from viral strains passaged through egg or other cell cultures or laboratory animals several times were excluded. When multiple sequences from the same isolate (as determined by GenBank taxonomic identification number) were available, only the most complete sequence was included. A total of 1646 HA and 1335 NA sequences were included in the final analyses. These sequence alignments are available at <http://routemap.osu.edu/publications>.

The segmented genome of influenza A viruses presents a difficulty for phylogenetic analysis. Since reassortment between diverse strains is reported to be common, two separate genomic segments are unlikely to share a phylogenetic history (Webster et al., 1992; Guan et al., 2002; Hatchette et al., 2004; Chen et al., 2006; Lindstrom et al., 1998, 2004). Reassortment violates some of the fundamental assumptions of the combined analysis (i.e. total evidence) approach such as the presumption of a shared history among

³ Lab for virus isolation: Mona M Aly, Abdel S Arafa from National Laboratory for Quality Control on Poultry Production, Dokki (Cairo), Egypt .

Lab for sequencing: Isabella Monne, Alice Fusaro, Giovanni Cattoli, Ilaria Capua from OIE/FAO Reference Laboratory for Avian Influenza and Newcastle Disease, Istituto Zooprofilattico Sper.le delle Venezie, Legnaro (PD), Italy.

1
2
3 various forms of data reflecting the history of the organismal lineages (Kluge, 1989).
4
5
6 Therefore, we do not combine data from HA and NA into a single dataset.
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Multiple alignment was performed on nucleotide sequences using CLUSTALW (version 2.09 Larkin et al., 2007) using the default settings (gap opening penalty 10, gap extension penalty 0.2, DNA transition weight 0.5). A small number of internal gaps in the alignment were redistributed such that the final alignment did not imply any breaks in the reading frame in ancestral sequences. The 5' and 3' ends of the alignment were trimmed using the start and stop codons as reference. Leading and trailing gaps that remained were replaced with question marks to indicate missing data. Internal gaps were treated as a fifth character state. These datasets are hereafter referred to as "HA1646" and "NA1335" datasets.

Wallace and Fitch datasets

The HA nucleotide sequence dataset analyzed by Wallace and Fitch (2008) was recreated from the list of GenBank accession numbers available as supplemental data to their paper. Following the materials and methods section of Wallace and Fitch (2008), their HA sequence dataset was aligned using CLUSTALW with the default settings without any manual adjustment. Leading and trailing gap markers were replaced with question marks to indicate missing data. Internal gap markers were treated as a fifth character state. This dataset is hereafter referred to as the "HA482WF" dataset.

Phylogenetic analysis

Phylogenetic analysis was performed using TNT (Goloboff et al., 2008) in parallel and sequential modes. All trees were rooted with A/Chicken/Scotland/1959/H5N1 as an outgroup. This study focuses on the highly pathogenic Asian lineage of H5N1 that is typically assumed to have originated with A/Goose/Guangdong/1996. We used the older isolate to check this assumption and polarize geographic character state changes. A preliminary analysis was performed to find the length of presumed most parsimonious trees. The TNT command 'xmult = level 5 replications 100000' was used for the preliminary runs on a single CPU. The argument 'level 5' allows for very thorough searches and was used only to find the presumed minimum length for each dataset. The preliminary runs were interrupted after 24 hours if they failed to produce shorter trees than previously found.

To create a representative sampling of the space of MPTs, the presumed minimum tree length was used to guide a new search to find 1000 independent hits to minimum length.

The TNT command 'xmult = level 0 hits 1000 giveupscore n' was used, where 'n' is the tree length of the most parsimonious trees found by the preliminary run.

To ensure that trees were well sampled, a cumulative series of consensus trees were calculated to check that a stable consensus had been reached. The first consensus tree was based on trees 0-9, the second on the first consensus plus trees 10-19, the third on the

1
2
3 second consensus and trees 20-29 [...] until all trees found from sampling 1000 hits to
4 minimum length have been included. Consensus trees were calculated using the TNT
5 command 'nelson* [starttree] [endtree];' which calculates a strict
6 consensus of all trees in the given range and saves the consensus as the last tree in
7 memory. The TNT command 'tnodes' was used to count the number of nodes that each
8 consensus tree contained to provide a metric of their stability. A consistent number of
9 nodes in a consensus tree indicate that we have collected a pool of candidate trees that
10 well represents topologies at presumed minimum length.
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Character optimization

Transmission events were calculated using a modified version of the method developed by Slatkin and Maddison (1989). To find transmission events, different geographic regions are treated as character states of a single multistate character, and are mapped onto trees obtained from phylogenetic analyses using standard optimization methods (Farris, 1970; Fitch, 1971; Swofford and Maddison, 1987).

To make our results comparable with previous phylogeographic studies of H5N1, we use equal weight between geographical character states. However, if preferred, the multistate character can be coded as a Sankoff character with a transition matrix specifying the relative costs of each type of transition between character states (Sankoff and Rousseau, 1975).

We counted the minimum and maximum number of transmissions between all pairs of geographic character states. When a character state change is observed, this is interpreted as a possible transmission route of the virus. The state changes are directional, so results are expected with different numbers for transformations from state A to state B and state B to state A. Figure 1 illustrates how state changes are counted.

Geographic information for the molecular datasets was extracted from the 'TSeq_orgname' field in GenBank's TinySeq XML records. For each molecular dataset (HA1646, NA1335, and HA482WF), new geographic character matrices were created. Each geographic character matrix contained an entry for a multistate character for each taxon in the nucleotide dataset. The order of taxa was maintained in nucleotide and geographic character datasets as this is important for consistency when using compact tree format and other tree formats that use numerical taxon names in TNT.

Three geographic character data sets of 31 character states were used to find transmission routes. Due to the limitations of the character state space allowed by the tree search software, TNT, we were unable to issue a single set of geographic character states suitable for comparing the HA1646 and NA1335 datasets to some previous studies. However, we were able to compare our results to Wallace and Fitch (2008). The first geographic character set we created, termed "W&F", was used to assign each taxon to a geographic category matching those used by Wallace and Fitch (2008). For taxa collected in localities not matching any of the 28 character states used in W&F, we assigned

1
2
3 character states to these taxa with the most fitting of three new states: "Other China",
4
5 "Other Asia" and "Other World".
6
7

8 We created a second character set termed "Africa" that focuses on relationships and
9 sources of H5N1 in regions of Africa and Europe. A third character set "World" was used
10 to look at broader patterns among regions of the globe. The character set World treats
11 mainland China as a single entity, treats Europe as four regions (north, east, south and
12
13 west), and all of Africa as a single region.
14
15
16
17
18

23 Counting transmission events

24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For each count of transmission events, a geographic character dataset containing a single multistate character and compact tree file containing the MPTs found during the search on the corresponding nucleotide data were loaded into memory. The TNT command change was used to calculate the maximum number of character state changes for pairs of geographic character states. The syntax used was 'change] ./ [n] / [region X] [region Y] '. The ']' indicates that results should be displayed character by character, and the '. / [n] / ' denotes that character number 'n' should be optimized in all trees in memory. For example, the command 'change] ./0/HongKong Fujian' will return the global minimum and maximum number of transmission events from Hong Kong to Fujian, calculated over all MPTs, when only a single multistate character is used. Each direction of change among states (e.g., A to B and B to A) was calculated for all possible pairs of states in the dataset.

1
2
3 From the output of the 'change' command, tables were created containing the global
4 minimum and maximum values for each location pair and direction. When the number of
5 state changes was >0, this result was interpreted as a possible transmission route from
6 one region to another.
7
8
9
10
11
12
13
14

15 Transmission events were calculated for the HA428WF sequence dataset using the W&F
16 character state set. For 1646HA and 1335NA sequence datasets, all three (W&F, Africa,
17
18 and World) character state sets were used.
19
20
21
22
23
24

25 We have created a web-application, Routemap, (<http://routemap.osu.edu>) to produce
26 visualizations of transmission events from the output of the geographic character
27 optimization studies over pools of trees. The required inputs to Routemap are 1) a table of
28 comma-separated values (CSV) containing the taxon label, the geographic place for the
29 isolates, the latitude, and longitude in decimal degrees, and 2) a matrix of aligned
30 nucleotide or amino acid data for the isolates in FASTA format. As an option, the user
31 can input a set of precomputed trees. In the first step, Routemap will produce a TNT file
32 that the user can execute on their machine with the command 'tnt p filename,'.
33
34 The user then will load the TNT output back into the Routemap via the browser, where
35 the underlying application will complete the geographic calculations. Routemap's final
36 output is a Keyhole Markup Language (KML) file compatible with Google Earth
37 (<http://earth.google.com>) and other geographic information systems such as ESRI's
38 ArcGIS Explorer (<http://www.esri.com/software/arcgis/explorer/index.html>). A
39 visualization produced by Routemap is similar to a web-based airline route map. When
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 using a Routemap in a virtual globe the user selects localities of interest and viral
4 transmission events implied by the analysis of phylogenetic trees are displayed as
5 radiating from origins or arriving at destinations. The user can make one or more origin
6 and destination selections to make as complex or as simple of a visualization as needed.
7
8 The lines indicating routes are colored to indicate directional, ambiguous or bidirectional
9 travel. Example datasets, KML files, scripts and instructions for users to make their own
10 Routemaps are available at <http://routemap.osu.edu>.
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Results

Alignment and phylogenetic trees

The alignment of the HA1646 dataset consisted of 1646 taxa and 1716 characters. 937 characters in HA1646 were parsimony-informative. The length of the presumed most parsimonious trees was 7284 steps. A search for 1000 independent hits to the MPT length resulted in 2721 unique trees. The number of nodes in the consensus stabilized at 500 trees at 695 nodes.

The alignment of the NA1335 dataset consisted of 1335 taxa and 1410 characters. 744 characters in NA1335 were parsimony-informative. The length of the presumed most parsimonious trees was 4795 steps. A search for 1000 independent hits to the MPT length resulted in 2649 unique trees. The number of nodes in the consensus stabilized at 860 trees at 540 nodes.

The HA482WF dataset consisted of 482 taxa and 1809 characters. 937 characters in HA482WF were parsimony-informative. Under our search strategy, the length of the presumed most parsimonious trees is 3257 steps. A search for 1000 independent hits to the MPT length resulted in 2470 unique trees. The consensus stabilized at 251 nodes after 140 trees were collected.

Reanalysis of Wallace and Fitch (2008)

Results of the reanalysis of Wallace and Fitch's (2008) HA nucleotide dataset (HA482WF) show that using a single tree, and a single optimization path, misestimates the frequency of transmission events. Moreover the use of a single tree can fail to detect possible transmission events. Figure 2 illustrates a comparison of the transmission events reported in Wallace and Fitch's (2008) with our reanalysis of HA428WF and W&F character data for localities. We used the global maximum number of transmissions from the pool of 2470 MPTs in our reanalysis. The global maximum of implied transmission events corresponds to Wallace and Fitch's use of DELTRAN optimization (Swofford and Maddison, 1987; Swofford, 1990). Thus a transmission route is indicated if it is found in any MPT.

There are 756 possible transmission routes in W&F character state data. We compare the results found by Wallace and Fitch's (2008) optimization of these character data over a single tree they derived from HA482WF sequence dataset to our reanalysis using 2470 MPT for HA482WF. 517 of the possible routes are not found by either study, one transmission route result is found only in Wallace and Fitch (2008), 169 new routes were

discovered by our reanalysis, and 69 routes are suggested by both Wallace and Fitch (2008) and our reanalysis. Detailed results are presented in supplemental data at <http://routemap.osu.edu/publications>.

We also analyzed W&F based on all the MPTs we considered for HA428WF or HA1646 sequence datasets. Of 756 possible routes in W&F, 528 routes were not found when optimizing these character data over MPTs we found for the HA428WF or HA1646 sequence datasets. Four transmission routes are found only in HA428WF, 158 routes are found only in HA1646, and 61 routes occur in both HA428WF and the HA1646 dataset. Thoughtful examination of transmission routes implied by the HA1646 sequence dataset required addition of three character states for viral isolates from regions not included in the HA428WF dataset. 111 additional transmission routes were detected in MPTs implied by HA1646 when we included three localities not considered by Wallace and Fitch, (2008) but considered by other authors (e.g., Duceatz et al., 2006, 2007; Salzberg et al., 2007).

Frequency of discovery of transmission routes

We compared the transmission routes found by Wallace and Fitch (2008) for a single tree to the pool of MPTs we derived from the HA428WF dataset. Wallace and Fitch (2008) did not find 22% of the transmission routes that are found in 95% or more of the MPT that we found for HA428WF. See supplemental data at <http://routemap.osu.edu/publications> for a summary table.

Visualization

The panels in Figure 3 contain visualizations produced with Routemap of viral traffic in Africa, east Asia, and southeast Asia from the HA1646 dataset and the World geographic character coding scheme. Visualizations of the HA1646 and NA1335 dataset under the W&F, World, and Africa geographic character coding schemes, as well as the HA428WF dataset are available at <http://routemap.osu.edu/publications>.

Discussion

Phylogenetic analysis and H5N1

Metadata on pathogens such as date, host, and geographic location for viral isolation are often combined with phylogenetic trees based on pathogen sequence data by researchers and public health officials during an epidemic. In order to interpret metadata in a phylogenetic context, character optimization is a necessary tool. Studies that use metadata and character optimization include investigations of: where H5N1 spreads (Wallace et al., 2007, and Wallace and Fitch, 2008; Salzberg et al., 2007), where and when H5N1 has switched hosts (Janies et al., 2007) or experienced mutations that confer resistance to drugs (Hill et al., 2009).

Character optimization in phylogeographic analysis

Here we demonstrate that when reporting transmission events from character state changes, a single tree and a single way of optimizing the character will only tell a fraction of the story. Wallace et al., (2007) and Wallace and Fitch (2008) used MIGRAPHYLA (HoDac et al., 2007) and PAUP* (Swofford, 2008) to calculate transmission events, interpreted as character state changes using a single multistate character with geographical regions as character states. As Wallace and Fitch (2008) point out: "PAUP's assignments may represent only one of several possible most-parsimonious traces". Thus Wallace and Fitch (2008) inferred transmission events from a single possible optimization out of many, of a 28-state character, on a single tree, out of many possible trees. When we examine the length of this 28-state character on the pool of 2470 MPTs from the HA482WF sequence dataset, the lengths are between 107 and 118 steps. This range of lengths is consistent with the 115 steps reported in Figure S2 in Wallace and Fitch (2008). On the other hand, when we examine the average number of transmission events supported by the MPTs we found for the HA482WF dataset, the numbers of transmission events are very different. Each MPT supports on average 131 different transmission routes in contrast to the 70 unique transmission routes reported by Wallace and Fitch (2008).

The difference in results between TNT and PAUP* is because PAUP* only considers a single optimization path of the character, when calculating steps on polarized character state changes. This is the behavior of PAUP* even in cases where there are multiple equally parsimonious optimization paths for the character and each path may support

1
2
3 different sets of polarized state changes. In contrast, in TNT it is possible to calculate the
4 maximum (or minimum) number of polarized state changes (e.g., from state A to state B)
5 between two states of a multistate character. Moreover, in TNT this calculation can be
6 done over multiple trees and multiple optimization paths per tree. Although the sum of all
7 state changes calculated this way with TNT will exceed the minimum length of the
8 character, it is important to find all possible state changes.
9
10
11
12
13
14
15
16
17
18
19

20 We also discovered an artifact in the case of closely related taxa. The number of
21 observed state changes between two states should not necessarily be interpreted as the
22 number of implied transmission events. In cases where the spread of an infectious agent
23 is rapid, it is very likely that isolates from neighboring regions will be genetically
24 identical or fail to have obtained any parsimony informative mutations. In such cases, the
25 number of observed state changes in the MPTs in a clade with no supported branches
26 (*sensu* the command `collapse 3` in TNT) will only reflect the proportion of taxa from
27 each location, when optimizing over a pool of dichotomous, MPTs. In other words, given
28 a character with two geographic states (location A and location B), the maximum number
29 of transmission events between A and B will be the maximum length of the character in
30 the subtree. For example, if there are 10 taxa from location A and five from location B in
31 a subtree, the maximum number of transmission events will be five, when counting
32 changes from A to B or from B to A.
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52

53 In cases where there are many sequences in a region where there is little or no genetic
54 difference, the artifact we describe above can occur. In other cases, where distinct genetic
55
56
57
58
59
60

lineages have been exchanged between two places, counting the changes in geographic characters is a useful tool to understand viral traffic. Both the HA428WF and HA1646 sequence datasets contain large subclades (e.g., in Africa) of viruses with identical genetic sequences. The recent invasion of H5N1 into Africa without much molecular evolution in HA makes it of interest for future studies to examine these isolates from a multi locus approach as done by Salzberg, et al., (2007) and Janies et al., (2007). These concepts are also applicable to pandemic H1N1, which emerged in early 2009 and has spread rapidly around the world.

Current phylogenetics software (PAUP*, Swofford, 2008; TNT, Goloboff et al., 2008) limit a categorical character to 31 states. These programs are suitable for most analyses of morphological and sequence data in which the data are used to find trees. MESQUITE (Maddison and Maddison, 2007) can summarize character state changes over pools of trees. MESQUITE is thus an alternative for these analyses, but is limited to 56 states and would have to be coupled to a strong tree search program. POY4 (Varón et al., 2008) may provide a solution to the state space limitation as the user can define any alphabet to describe states for features of interest of pathogens and transformation costs among states in a look-up table (e.g., a Sankoff matrix). Heuristic efficiency will be important as an accurate regional map may require hundreds or thousands of geographical character states.

When a widely distributed taxon such as influenza A is studied, geographic character coding requires some abstraction of the globe into regions. We followed the 28-location

1
2
3 state coding scheme of Wallace and Fitch (2008), as well as implementing two alternative
4 coding schemes (i.e. Africa and World) to make our results comparable to other studies.
5
6

7 When comparing the HA1646 and NA1335 sequence datasets to the results of Wallace
8 and Fitch (2008), 275 isolates were not compatible with the 28 states defined by Wallace
9 and Fitch (2008). We coded the isolates that fell outside of the character state space with
10 one of three new character states: Other Asia, Other China, and Other World. This coding
11 scheme allowed our analyses to find transmission routes outside the scope of Wallace and
12 Fitch's (2008) original publication, while retaining comparable results.
13
14
15
16
17
18
19
20
21
22
23
24

25 By including more isolates from all regions sampled by HA1646 and NA1335 sequence
26 datasets, we identified more possible routes of transmission than Wallace and Fitch
27 (2008). For example, when using the same character data as W&F, our results agree with
28 those of an independent study, Wang et al., (2008), in finding a single unambiguous
29 transmission from the Hunan province of China to Indonesia. Wallace and Fitch (2008)
30 suggested a transmission from Guangdong, China to Indonesia. The Guangdong to
31
32 Indonesia route is not observed under any optimization of character data in W&F across
33
34 the MPTs we found for the HA1646 sequence dataset.
35
36
37
38
39
40
41
42
43
44
45

46 In all studies of the phylogeography of influenza there are problems associated with the
47 imprecise naming of viruses. For example in the HA1646 sequence dataset, 18 viruses
48 are labeled "China" which represents an area of 9,598,08 km². This area is comparable to
49 10,180,000 km² for all of Europe, including Russia west of the Ural Mountains. Other
50
51 viral isolates are named for cities. In some cases viral isolates are named for regions.
52
53
54
55
56
57
58
59
60

1
2
3 Many workers, including Wallace have suggested that GenBank records for viral isolates
4
5 should include latitude and longitude records to avoid geographic imprecision (Butler,
6
7 2008). Similarly, there is much room for improvement in metadata standards for naming
8
9 of host taxa (Janies et al., 2007). In the wake of the H1N1 pandemic, GISAID has vastly
10
11 improved metadata standards. High quality metadata will improve research in viral
12
13 character evolution and geography studies.
14
15

17
18
19 We used the results from HA1646 and NA1335 sequence datasets, and character states
20
21 from the World coding scheme to find transmission routes discussed in the literature. The
22
23 major difference between W&F and World is that mainland China is treated as a single
24
25 region in the World character state set, and countries in Europe and Africa are separate
26
27 into regions (discussed above). Consolidation of regions within China was necessary to
28
29 allow all countries (as defined by political borders) in east Asia to fit into the 31 character
30
31 state limit.
32
33
34
35
36
37
38
39

40 Current state of the H5N1 epidemic 41 42 43 44

45 Since 2003, H5N1 has been endemic in Indonesia. All phylogenetic studies including HA
46
47 sequences isolated from H5N1 in Indonesia have found a monophyletic clade of
48
49 Indonesian isolates, independent of optimality criterion (neighbor joining: Guan et al.,
50
51 2004; Smith et al., 2006; Salzberg et al., 2007, Bayesian: Wang et al., 2008; Smith et al.,
52
53 2006, maximum likelihood: Kilpatrick et al., 2006; Wallace and Fitch, 2008; Hill et al.,
54
55 2009, parsimony: Janies et al., 2007; Salzberg et al., 2007; Wallace et al., 2007).

1
2
3
4
5 It is apparent that the H5N1 strains within Indonesia occurred due to a single
6 transmission event from southern China. Janies et al., (2007) use whole genomes to
7 demonstrate this link. Guan et al., (2004) use HA data to suggest the closest relatives to
8 isolates from Indonesia are viruses from Yunnan province in China. Wallace et al.,
9 (2007), and Wallace and Fitch (2008) also us HA data to find that Guangdong province in
10 China is the most likely source. These results stand in contrast to Kilpatrick et al., (2006),
11 who used data from international trade to assert that the source of H5N1 in Indonesia was
12 wild birds imported from Thailand. Our analyses of sequence and geographic data for
13 H5N1 are in disagreement with the assertion of Kilpatrick et al, as we never find a
14 Thailand to Indonesia route. Our analyses of the HA1646 sequence dataset, suggest that
15 the geographic origin of H5N1 for Indonesia was Hunan province in China. Wang et al.,
16 (2008) concur. In our analyses of HA1646 sequence data, the Hunan-Indonesia migration
17 route appears in all MPTs independent of optimization method.
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39 Similarly our analyses of the NA1335 sequence dataset suggest Hunan as a possible
40 source for Indonesia but some of the trees for NA1335 also indicate Hong Kong and the
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Similarly our analyses of the NA1335 sequence dataset suggest Hunan as a possible
source for Indonesia but some of the trees for NA1335 also indicate Hong Kong and the
Guangxi Zhuang Autonomous Region of China as possible sources for Indonesia.
Conflict between genes result could represent a reassortment event in which viruses in
Indonesia have HA and NA genes with distinct histories or could simply reflect
ambiguity in the data.

1
2
3
4
5 Smith et al., (2006), found that all H5N1 from Vietnam, Cambodia and Thailand form a
6 single clade, except for one isolate from Vietnam. We were unable to find this clade
7 structure in any tree found among the most parsimonious trees from the analyses of the
8 HA1646 or NA1335 datasets. Although most isolates from Vietnam, Thailand, Laos,
9 Cambodia and Vietnam are found in a monophyletic clade, isolates from Vietnam and
10 Thailand are always found outside of this clade. This result suggests that the results of
11 Smith et al., (2006) are an artifact of limited taxon sampling.
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

In the pool of trees resulting from searches of the HA1646 sequence data, we see at least three separate introductions of H5N1 to Japan. There were transmissions events from China to Japan in 2003 and 2008. Under certain optimizations, Vietnam is also a possible source of the 2008 isolates in Japan. In contrast, the source of the 2004 outbreak in Japan appears to be South Korea in all trees and all optimizations based on analysis of HA1646 sequence data and the W&F character data.

These comparisons illustrate the danger in excluding data. Wallace et al., (2007) and Wallace and Fitch (2008) both find the China-Japan connection, but excluded isolates collected in 2002-2003 from South Korea despite the fact these data were available in public databases as of 2005. Due to this exclusion, Wallace et al., (2007) and Wallace and Fitch (2008) did not find South Korea as a source of the Japanese outbreak in 2004. Our phylogenetic analysis indicating transmission of H5N1 from South Korea to Japan in

1
2
3 2004 is in agreement with Kilpatrick et al., (2006) who based their inference on trade and
4
5 bird migration data.
6
7
8
9

10 The H5N1 outbreaks in Africa have been investigated using HA sequences (Ducatez et
11
12 al., 2006) as well as full genomes (Ducatez et al., 2007, Salzberg et al., 2007). However,
13
14 these studies were conducted when there was a limited set of putative outgroup(s) for the
15
16 African lineages in the public domain. These studies conclude that the sources of H5N1
17
18 are not southeast Asia, but rather central Asia and Europe. The pools of trees resulting
19
20 from our analysis of HA1646 agree with Ducatez et al., (2006; 2007) that H5N1 has been
21
22 introduced to Nigeria on three separate occasions. Using the "Africa" character state set,
23
24 we find corroboration for the sources to be central Asia or Europe rather than southeast
25
26 Asia. By examining at all trees and all optimizations we found for HA1646 and the W&F
27
28 character state set, we can not rule out the source of one transmission to Africa was the
29
30 Qinghai province of China (Salzberg et al., 2007).

31
32
33
34 Sources for H5N1 in Europe are unclear⁴, as multiple possible sources are indicated for
35
36 most countries. The basis for the ambiguity is the high similarity of the isolates found in
37
38 Europe, Africa, and the Middle East in 2005-2006. This is another empirical example of
39
40 the artifact we described above, where very low genetic variation leads to every region
41
42 being a possible source for its neighbors.

43
44
45
46
47
48
49
50
51
52
53
54
55 ⁴ One important exception is the case of an infected eagle (*Nisaetus nipalensis* Hodgson,
56
57 1836) that was smuggled from Thailand to Belgium but did not cause any further spread
58
59 (Van Borm et al., 2005).
60

Perils of using a single tree

A single tree is likely to miss several transmission routes found frequently among the MPTs, as illustrated in Fig. 4. All of the publications on phylogeographic analysis of H5N1 we have reviewed here use a single tree for the results they report. This practice is independent of whether they use Slatkin and Maddison (1989) transmission event analysis (Wallace et al., 2007; Wallace and Fitch, 2008) or other optimization methods [POY3 (Wheeler et al., 2005) command -diagnose in Janies et al., 2007; MESQUITE command trace:parsimony optimzation in Salzberg et al., 2007] or observe transmission events by inspection (e.g., Duceatz et al., 2006, 2007; Wang et al., 2008).

The use of one tree is also independent of which optimality criteria researchers use in tree search. In some cases (e.g., Smith et al., 2006; Wang et al., 2008) only a single neighbor-joining tree is produced. Janies et al., (2007) created a binary tree on which to optimize based on a strict consensus of MPT using the TNT command randtree*. In other cases, information from multiple trees is represented as a consensus tree. A consensus tree can be a majority rules consensus from bootstrap resampling under maximum likelihood (Kilpatrick et al., 2006; Smith et al., 2006), parsimony (Guan et al., 2004; Salzberg et al., 2007;), or a distribution of trees from a Bayesian analysis (Smith et al., 2006).

There are several problems with the use of a consensus tree in the context of character optimization. First, a consensus tree does not represent all MPTs. Also the consensus

1
2
3 topology is never one of the most parsimonious trees. Furthermore, trees can be
4
5 constructed that do not change the consensus tree, but still are not found among the
6
7 MPTs. Finally, a set of most parsimonious trees can be found, where the strict consensus
8
9 tree would imply transmission routes that are never found in any of the individual MPTs,
10
11 (as illustrated in Figure 4). It is thus inappropriate to optimize character state changes on
12
13 a consensus tree.
14
15

17
18
19 Under parsimony, each tree in the pool of MPTs and each optimization of the character
20
21 states are just as optimal as the next one. The transmission routes found in only a few
22
23 trees and under certain optimizations are just as supported by the optimality criterion as
24
25 those trees and optimization found frequently. However, under a model-based optimality
26
27 criterion, probability calculations for the existence of transmission routes could be
28
29 calculated. A distribution of trees from a Bayesian analysis, where each tree has a
30
31 posterior probability could be used to provide confidence values for transitions between
32
33 the character states in a multistate character for geography optimized on trees of various
34
35 posterior probability. Alternatively the proportional likelihood can be calculated for
36
37 ancestral geographic states. One could employ a method such as described to study the
38
39 evolution morphological characters under likelihood (Lewis, 2001) and is implemented in
40
41 MESQUITE (Maddison and Maddison, 2003).
42
43
44
45
46
47
48

49 **Problems with using multiple trees** 50

51
52
53
54 When examining character state changes over a pool of MPTs, different types of
55
56 optimization need to be taken into account. If one examines the global minimum of state
57
58
59
60

1
2
3 changes over a pool of trees, they will find only those state changes that appear in all
4
5 trees, under all optimizations. When dealing with large numbers of trees, a tree will often
6
7 be found where a specific state change does not occur. For example, one has a subtree
8
9 with only unsupported branches and 100 taxa from location A, 100 from location B and a
10
11 single taxon from location C. When looking at this subtree in the MPTs, one can find a
12
13 tree and an optimization pathway that does not have a transmission event between
14
15 location A and B, making the global minimum number of transmission events from A to
16
17 B equal zero.
18
19
20
21
22
23
24

25 In contrast, use of the global minimum number of polarized state changes dramatically
26 changes the number of indicated transmission routes for a set of trees when compared to
27
28 use of the maximum number of state changes. For example, when considering the
29 maximum number of directional state changes on the HA1646 dataset, we find a total of
30
31 267 different transmission routes. This results contrasts with a result of 20 transmission
32
33 routes found when looking at the global minimum number of directional state changes. In
34
35 summary, the global minimum is a severe test of putative transmission routes as the route
36
37 has to appear in every tree in the pool of MPTs to register a global minimum > zero. The
38
39 global maximum is useful to find the range of possible transmission routes a set of MPTs
40
41 implies, irrespective of how frequently a route is found.
42
43
44
45
46
47
48

49 Conclusions 50 51 52 53 54

55 Phylogenetic tree search and character evolution analyses have been valuable tools in
56 many fields including surveillance of infectious diseases. In this context, parsimony is
57
58
59
60

sometimes used for evaluation of trees and almost always used for character optimization. The multiple equally optimal trees produced by parsimony reflect of the natural ambiguity of the sequence data. Character optimization under parsimony can find the global minimum and maximum values for state changes in a pool of trees reflecting various degrees of ambiguity underlying transmission routes for organismal lineages.

With efficient character optimization and tree search heuristics it is possible to find and analyze very large pools of trees for datasets with thousands of terminal taxa. Moreover, because results are sensitive to tree choice, any tree search method that considers multiple trees has an advantage over methods that consider a single tree, such as neighbor joining.

The methods of character optimization and visualization used and developed here are also applicable to sets of trees produced by maximum likelihood and Bayesian methods.

Maximum likelihood methods are commonly used for tree search and sometimes preferred for the operational reason that they support one tree slightly more than near suboptimal candidate trees. While Bayesian methods produce a distribution of trees, conclusions are typically discussed in relationship to the single tree presented (often a majority rules consensus and sometimes a maximum *a posteriori* tree). As we have demonstrated here, a single tree, as well as consensus trees are insufficient for estimating transmission routes of infectious disease. In addition to considering multiple trees, when using Slatkin and Maddison's (1989) method of geographic analysis, it is important to take into account sensitivity to different strategies for character state optimization.

1
2
3 It is our intention that the methods described in this paper should be general and useful
4
5 for phylogeographic analysis of any type of problem above and beyond infectious
6
7 diseases. We welcome everyone to our website <http://routemap.osu.edu> to test our
8
9 methods and create KML visualization of your own data.
10
11
12
13
14
15

Acknowledgements

16
17
18
19
20
21 We thank the Pablo Goloboff, Steve Farris, Kevin Nixon, and the Willi Hennig Society
22
23 for making TNT available. Pablo Goloboff, Universidad Nacional de Tucumán and
24
25 Diego Pol, CONICET, Museo Paleontológico Egidio Feruglio provided insights and
26
27 discussions on character analysis. We also thank two anonymous reviewers for their
28
29 constructive criticisms.
30
31
32
33
34

35
36 We acknowledge the Department of Biomedical Informatics and The Mathematical
37
38 Biosciences Institute of The Ohio State University (OSU) for space, funding, and
39
40 administrative support. We thank the Medical Center Information Services team of OSU
41
42 and the Mathematical Biosciences Institute for hosting computing clusters used in this
43
44 study. We acknowledge this material is based upon work supported by, or in part by, the
45
46 U.S. Army Research Laboratory and Office under grant number W911NF-05-1-0271.
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Literature cited

Bright, R.A., Medina, M., Xu, X., Perez-Oronoz, G., Wallis, T.R., Davis, X., Povinelli, L., Cox, N., Klimov, A. , 2005. Incidence of adamantane resistance among influenza A (H3N2) viruses isolated worldwide from 1994 to 2005: A cause for concern. 2005. Lancet 366,1175-1181.

Butler, D. 2008. Politically correct names given to flu viruses. Nature News online 23 April 2008. doi:10.1038/452923a .

Chan, P. Outbreak of Avian Influenza A(H5N1) Virus infection in Hong Kong in 1997 (2002). Clin. Infect. Dis. 34 (Suppl 2), S58–64.

Chen, H., Bahri, S., Chen, Y., Cheung, C.L., Duan, L., Fan, X.H., Guan, Y. Guo, C.T., Hassan, S.S., Huang, K., Leung, Y.H.C., Li, K.S., Lu, H.R., Naipospos, T.S.P., Nguyen, T.D., Peiris, J.S.M., Qin, K., Rayner, J.M., Smith, G.J.D., Vijaykrishna, D., Wang, J., Webster, R.G., Wu, W.L., Xia, N.S., Xu, K.M., Yuen, K.Y., Zhang, J.X., Zhang, L.J. 2006. Establishment of multiple sublineages of H5N1 influenza virus in Asia: Implications for pandemic control. PNAS. 103, 2845-2850.

De, B. K., Brownlee, G. G., Kendal, A. P., Shaw, M. W. 1988. Complete sequence of cDNA clone of the hemagglutinin gene of influenza A/Chicken/Scotland/59 (H5N1)

virus: comparison with contemporary North American and European strains. Nucleic Acids Res. 16, 4181-4182.

Ducatez, M.F., Ammerlaan, W., de Landtsheer, S., Fouchier, R.A.M., Muller, C.P., Niesters, H.G.M., Olinger, C.M., Osterhaus, A.D.M.E., Owoade, A.A. 2006. Avian Flu: Multiple introductions of H5N1 in Nigeria. Nature 442, 37.

Ducatez, M. F., Ammerlaan, W., De Landtsheer, S., Fouchier, R. A. M., Muller, C. P., Olinger, C. M., Osterhaus, A. D. M. E., Ouedraogo, J. B., Owoade, A. A., Sow, A., Tahita, M. C., Tarnagda, Z. 2007. Molecular and antigenic evolution and geographical spread of H5N1 highly pathogenic avian influenza viruses in western Africa. J. Gen. Virol. 88, 2297-2306.

Farris, J.S. Methods for computing Wagner trees. 1970. Syst. Zool. 19, 83-92.

Fitch, W.M. 1971. Toward defining the course of evolution: minimum change for a specific tree topology. Syst. Zool. 20, 406-416.

Goloboff, P. A., Farris, J.S., Nixon, K.C. 2008. TNT, a free program for phylogenetic analysis. Cladistics 24, 774-786

GISAID, 2009. Sample provided by Aly, M. M., Arafa, A. S. National laboratory for quality control on poultry production, Dokki (Cairo), Egypt. Sequence submitted by

1
2
3 Monne, I., Fusaro, A., Capua, I., Cattoli, G. OIE/FAO Reference Laboratory for avian
4 influenza and Newcastle disease, Istituto Zooprofilattico Sperimentale delle Venezie, Legnaro
5 (PD), Italy.
6
7
8
9
10
11
12

13 Guan, Y., Dyrting, K.C., Ellis, T. M., Kong, K. F., Peiris, Y., Shortridge, K. F., Sit, T.,
14
15 Zhang, L. J. 2002. H5N1 influenza viruses isolated from geese in southeastern China:
16 evidence for genetic reassortment and interspecies transmission to ducks. *Virology* 292,
17
18 16-23.
19
20
21
22
23

24 Guan, Y., Auewarakul, P., Buranathai, C., Chaiseng, A., Chen, H., Duan, L.,
25
26 Estoepangestie, A.T.S., Hanh, N.T.H., Li, K.S., Long, H.T., Nguyen, T.D., Peiris, J.S.M.,
27
28 Poon, L.L.M., Puthavathana, P., Rahardjo, A.P., Shortridge, K.F., Smith, G.J.D., Wang,
29
30 J., Webby, R.J., Webster, R.G., Xu, K.M., Yuen, K.Y. 2004. Genesis of a highly
31 pathogenic and potentially pandemic H5N1 influenza virus in eastern Asia. *Nature* 430,
32
33 209-213.
34
35
36
37
38
39
40

41 Hatchette, T.F., Baker, A., Johnson, C., Pryor, S. P., Walker, D., Webster, R.G. 2004.
42
43 Influenza A viruses in feral Canadian ducks: extensive reassortment in nature. *J. Gen.*
44
45 *Virol.* 85, 2327-2337.
46
47
48
49

50 Hill, A., Guralnick, R. Wilson, M., Habib, F. Janies, D. 2009. Evolution of drug
51 resistance in multiple distinct lineages of H5N1 avian influenza. *Infection, Genetics, and*
52
53 *Evolution.* 9, 169-178.
54
55
56
57
58
59
60

1
2
3
4
5 HoDac, H., Fitch, W.M., Lathrop, R.H., Wallace, R.G. 2007. MigraPhyla: Statistical
6 analysis of migration events through a phylogeny. Version 1.0b manual.
7
8
9

10 <http://pd.bio.uci.edu/ee/WallaceR/MigraPhyla.html>
11
12
13
14

15 Holland, J., Grabau, E., Horodyski, F., Nichol, S., Spindler, K., VandePol, S. 1982. Rapid
16 evolution of RNA genomes. *Science* 215, 1577-1585.
17
18

19 Huelsenbeck, J.P and Ronquist, F. 2001. MRBAYES: Bayesian inference of phylogenetic
20 trees. *Bioinformatics* 17, 754-755.
21
22

23
24
25
26
27 Janies, D., Hill, A.W., Habib, F., Guralnick, R., Waltari, E., Wheeler, W.C. 2007.
28 Genomic analysis and geographic visualization of the spread of avian influenza (H5N1).
29
30 Syst. Biol. 56, 321-329.
31
32
33
34
35

36 Kawaoka Y., Alexander D.J., Nestorowicz A., Webster R.G. 1987. Molecular analyses of
37 the hemagglutinin genes of H5 influenza viruses: origin of a virulent turkey strain.
38
39 Virology 158, 218-227.
40
41
42

43
44
45 Kilpatrick, A. M., Chmura, A.A., Gibbons, D.W., Fleischer, R.C., Marra, P.P., Daszak, P.
46 2006. Predicting the global spread of H5N1 avian influenza. *PNAS*. 103, 19368-19373
47
48

49
50
51
52 Kluge, A.G., 1989. A concern for evidence and a phylogenetic hypothesis of
53 relationships among Epicrates (Boidae, Serpentes). *Syst. Zool* 38, 7-25.
54
55
56
57
58
59
60

Larkin, M.A., Blackshields G., Brown N.P., Chenna R., McGettigan P.A., McWilliam H., Valentin F., Wallace I.M., Wilm A., Lopez R., Thompson J.D., Gibson T.J., Higgins D.G. 2007. Clustal W and Clustal X version 2.0. Bioinformatics 23, 2947-2948.

Lewis, P.O. 2001. A likelihood approach to inferring phylogeny from discrete morphological characters. Syst. Biol. 50:913-925.

Lindstrom, S.E., Hiromoto, Y., Nerome, R., Omoe, K., Sugita, S., Yamazaki, Y., Takahashi, T., Nerome, K. 1998. Phylogenetic analysis of the entire genome of influenza A (H3N2) viruses from Japan: evidence for genetic reassortment of the six internal genes. J. Virol. 72, 8021-8031.

Lindstrom, S.E., Cox, N.J., Klimov, A. 2004. Genetic analysis of human H2N2 and early H3N2 influenza viruses, 1957-1972: evidence for genetic divergence and multiple reassortment events. Virology 328, 101-119.

Maddison, W. P. and D.R. Maddison. 2008. Mesquite: a modular system for evolutionary analysis. <http://mesquiteproject.org>

Matrosovich, M., Kawaoka, Y., Webster, R., Zhou, N. 1999. The surface glycoproteins of H5 influenza viruses isolated from humans, chickens, and wild aquatic birds have distinguishable properties. J. Virol. 73, 1146-1155.

1
2
3
4
5 Obenauer, J.C., Denson, J., Fan, Y., Finkelstein, D.B., Hoffmann, E., Krauss, S., Ma, J.,
6
7 Mehta, P.K., Mukatira, S., Naeve, C.W., Rakestraw, K.M., Su, X., Wang, J., Webster,
8
9 R.G., Xu, X., Zhang, Z., Zheng, J. 2006. Large-scale sequence analysis of avian influenza
10 isolates. *Science* 311, 1576-1580.
11
12
13
14
15
16

17 Salzberg, S.L., Ali, A.S.M., Aly, M.M., Brown, I.H., Capua, I., Cattoli, G., Couacy-
18 Hymann, E., De Mia, G.M., Dung, D.H., Ghedin, E., Guercio, A., Janies, D.A., Joannis,
19 T., Kingsford, C., Osmani, A., Padalino, I., Saad, M.D., Saviæ, V., Sengamalay, N.A.,
20 Spiro, D.J., Yingst, S., Zaborsky, J., Zorman-Rojs, O. 2007. Genome analysis linking
21 recent European and African influenza (H5N1) viruses. *Emerg. Infect. Dis.* 13, 713-718.
22
23
24
25
26
27
28
29
30
31

32 Sankoff, D. and Rousseau, P. 1975. Locating the vertices of a Steiner tree in arbitrary
33 space. *Mathematical Programming*, p. 240246.
34
35
36
37

38 Sawabe, K., Hayashi, T., Hoshino, K., Hotta, A., Isawa, H., Kobayashi, M., Kurahashi,
39 H., Saito, T., Sasaki, T., Tanabayashi, K., Tsuda, Y., Yamada, A. 2006. Detection and
40 isolation of highly pathogenic H5N1 avian influenza A viruses from blow flies collected
41 in the vicinity of an infected poultry farm in Kyoto, Japan, 2004. *Am. J. Trop. Med. Hyg.*
42 75, 327-332.
43
44
45
46
47
48
49
50
51
52

53 Slatkin, M. and Maddison, W. P. 1989. A cladistic measure of gene flow inferred from
54 the phylogeny of alleles. *Genetics* 123, 603-613.
55
56
57
58
59
60

1
2
3
4
5 Smith, D., Lapedes, A., de Jong, J., Bestebroer, T., Rimmelzwaan, G., Osterhaus, A.,
6
7
8 Fouchier, R. 2004. Mapping the Antigenic and Genetic Evolution of Influenza Virus.
9
10 Science 305, 371-376.
11
12
13
14

15 Smith, G.J.D., Bui, N.A., Chen, H., Cheung, C.L., Dao, T.V., de Jong, M.D., Farrar, J.,
16
17 Guan, Y., Hassan, S.S., Hien, T.T., Leung, Y.H.C., Li, K.S., Naipospos, T.S.P., Nguyen,
18
19 T.D., Nguyen, T.V., Nguyen, V.C., Peiris, J.S.M., Poon, L.L.M., Rayner, J.M., Usman,
20
21 T.B., Vijaykrishna, D., Webster, R.G., Zhang, J.X., Zhang, L.J. 2006. Evolution and
22
23 adaptation of H5N1 influenza virus in avian and human hosts in Indonesia and Vietnam.
24
25 Virology 350, 258-268.
26
27
28
29
30
31

32 Subbarao, K., Klimov, A., Katz, J., Regnery, H., Lim, W., Hall, H., Perdue, M., Swayne, D.,
33
34 Bender, C., Huang, J., Hemphill, M., Rowe, T., Shaw, M., Xu, X., Fukuda, K., Cox, N. 1998.
35
36 Characterization of an avian influenza a (H5N1) virus isolated from a child with a fatal
37
38 respiratory illness. Science 279, 393-396.
39
40
41
42

43 Swofford, D. L., 1990. PAUP: Phylogenetic Analysis Using Parsimony, Version 3.0.
44
45 Manual. Illinois Natural History Survey, Champaign, Ill.
46
47

48
49
50 Swofford, D.L. 2008. PAUP* Phylogenetic Analysis Using Parsimony (*and Other
51
52 Methods). Version 4. Sinauer Associates, Sunderland, Massachusetts.
53
54
55
56
57
58
59
60

1
2
3 Swofford, D.L. and Maddison, W.P. 1987. Reconstructing ancestral character states
4
5 under Wagner parsimony. *Math. Biosci.* 87, 199-229.
6
7
8
9

10 USDA, 2007. Release No. 0296.06. Avian influenza. Low pathogenic H5N1 vs. highly
11
12 pathogenic H5N1. Latest update, 23 July 2007. <http://www.usda.gov/2006/08/0296.xml>
13
14
15

16
17 Van Borm, S., Boschmans, M., Decaestecker, M., Dupont, G., Hanquet, G., Lambrecht,
18
19 B., Snacken, R., Thomas, I., van den Berg, T. 2005. Highly pathogenic H5N1 influenza
20
21 virus in smuggled Thai eagles, Belgium. *Emerg. Infect. Dis.* 11, 702-705.
22
23

24
25
26 Wallace, R.G., HoDac, H., Lathrop, R.H., Fitch, W.M. 2007. A statistical
27
28 phylogeography of influenza A H5N1. *PNAS* 104, 4473-4478.
29
30

31
32
33 Wallace, R.G. and Fitch, W.M. 2008. Influenza A H5N1 immigration is filtered out at
34
35 some international borders. *PLoS ONE* 3, e1697. doi:10.1371/journal.pone.0001697.
36
37

38
39
40 Wang, J., Bahl, J., Chen, H., Smith, G.J.D., Duan, L., Guan, Y., Peiris, J. S. M.,
41
42 Vijaykrishna, D., Webster, R. G., Zhang, J. X. 2008. Identification of the progenitors of
43
44 Indonesian and Vietnamese avian influenza A (H5N1) viruses from southern China. *J.*
45
46 *Virol.* 82, 3405-3414.
47
48
49

50
51
52 Webster, R.G., Bean, W.J., Chambers, T.M., Gorman, O.T., Kawaoka, Y. 1992.
53
54 Evolution and ecology of influenza A viruses. *Microbiol. Rev.* 56, 152-179.
55
56

1
2
3
4
5 Wheeler,W. C., Delaet, J., Gladstein, D., and Varón, A., 2005. POY (version 3.012).
6
7 Phylogeny reconstruction via optimization of DNAand other data.
8
9
10 <http://research.amnh.org/scicomp/projects/poy.php>
11
12
13
14
15 WHO, 2004. Avian influenza A(H5N1)- update 31: Situation (poultry) in Asia: need for a
16 long-term response, comparison with previous outbreaks. Reported to WHO. 2 March
17
18 2004. http://www.who.int/csr/don/2004_03_02/en/
19
20
21
22
23
24 WHO, 2009. Cumulative Number of Confirmed Human Cases of Avian Influenza
25 A/(H5N1) Reported to WHO. 11 August 2009.
26
27 http://www.who.int/csr/disease/avian_influenza/country/cases_table_2009_08_11/en/index.html
28
29
30
31
32
33
34
35
36 Varón, A., Vinh, L. S., Bomash, I., Wheeler, W. C. 2009. POY 4.1.1. American Museum
37 of Natural History. <http://research.amnh.org/scicomp/projects/poy.php>
38
39
40
41
42
43 Xu X., Cox, N.J., Guo, Y., Subbarao K. 1999. Genetic characterization of the pathogenic
44 influenza A/Goose/Guangdong/1/96 (H5N1) virus: similarity of its hemagglutinin gene to
45 those of H5N1 viruses from the 1997 outbreaks in Hong Kong. Virology 261, 15-19.
46
47
48
49
50
51
52 Zamarin, D., Ortigoza, M.B., Palese, P. 2006. Influenza A Virus PB1-F2 Protein
53 Contributes to Viral Pathogenesis in Mice. Journal of Virology 80: 7976-7983.
54
55
56
57
58
59
60

1
2
3 Zhou, N. N., Shortridge, K.F., Claas, E.C.J., Krauss, S.L., Webster, R.G. 1999. Rapid
4
5 evolution of H5N1 influenza viruses in chickens in Hong Kong. J. Virol. 73, 3366-3374.
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Figure captions**Figure 1.**

The same tree is shown with two different optimizations of a geographic character. The tree on the left shows one transmission from location A to location B and one from location B to location A. The tree on the right has two transmission events from location A to location B, and zero transmission events from location B to location A. The global maximum number of transmissions are thus: location A to location B: 2; location B to location A: 1.

Figure 2.

A comparison of transmission routes found by Wallace and Fitch (2008) and our reanalysis of the WF482HA dataset. Each cell in the grid represents a transmission route between two locations in the 28 state geographical character set used in the original publication. Cells colored black indicate transmission routes that were not found in any analysis. Cells colored yellow indicate that our reanalysis agrees with Wallace and Fitch (2008). Cells colored cyan indicates that the transmission route was not found in our reanalysis. Cells colored red indicate transmission routes found in the reanalysis, but not in Wallace and Fitch (2008).

1
2
3 The rows represent outgoing locations and the columns incoming. For example, the cell
4
5 in the Vietnam column and the Guangdong row represents a transmission from
6
7 Guangdong to Vietnam, while the cell in the Guangdong column and the Vietnam row
8
9 represents a transmission from Vietnam to Guangdong
10
11
12
13
14
15

16 **Figure 3.**
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Migration routes indicated by the HA1646 dataset.

A. Transmissions into Nigeria, using the Africa character set. The orange lines represent transmission routes are bidirectional or ambiguous. In H5N1 in Africa, this result is due to very low genetic variation between the H5N1 in Nigeria and neighboring regions.

B. Transmissions into Indonesia, using the "W&F" character set. Analyses indicate only one source of H5N1 in Indonesia, Hunan province in China.

C. Transmission into Japan, using the "World" character set. Both China and South Korea are sources for the H5N1 in Japan. The direction of the transmission route to Vietnam is either bidirectional or ambiguous.

Figure 4:

1
2
3 The consensus tree can imply transmission routes never found among the most
4
5 parsimonious trees.
6
7
8
9

10 When an unordered multistate character with states 0, A, B and C is optimized on trees a.
11
12 and b., no optimal optimization path has the directional state change C to A. However,
13
14 considering optimal optimization paths in the strict consensus tree, one path has the node
15
16 above the outgroup assigned to state B and thus two changes from B to C, as well as two
17
18 changes from C to A are required.
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Cladistics

location A 1

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

location A 2

location B 1

location A 3

location B 2

Cladistics

■ → □ : 2
□ → ■ : 0

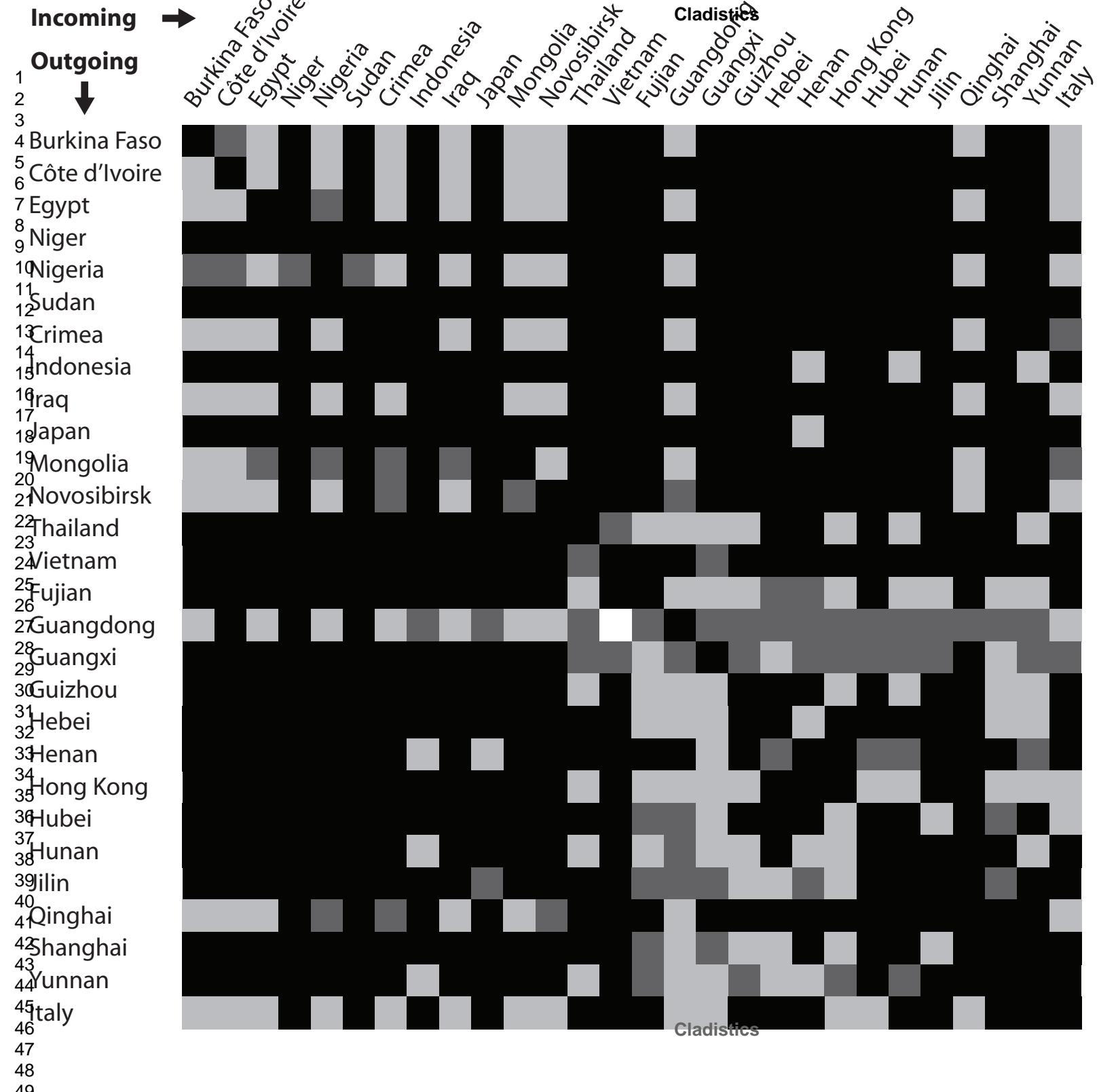
location A 3

location B 2

location B 1

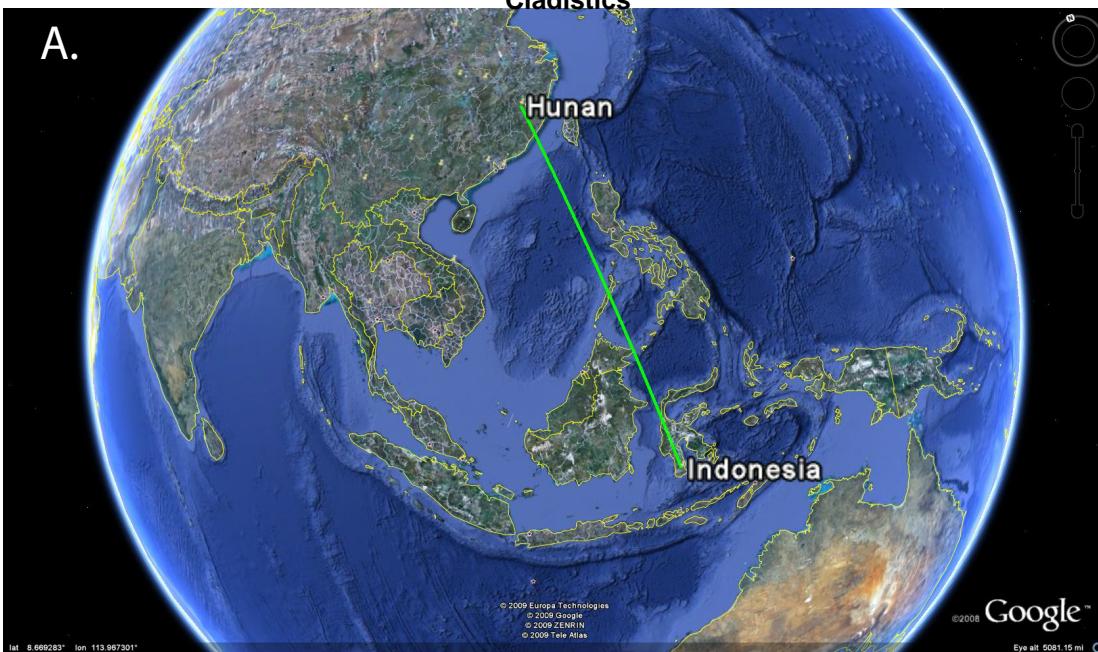
location A 2

location A 1



Cladistics

A.

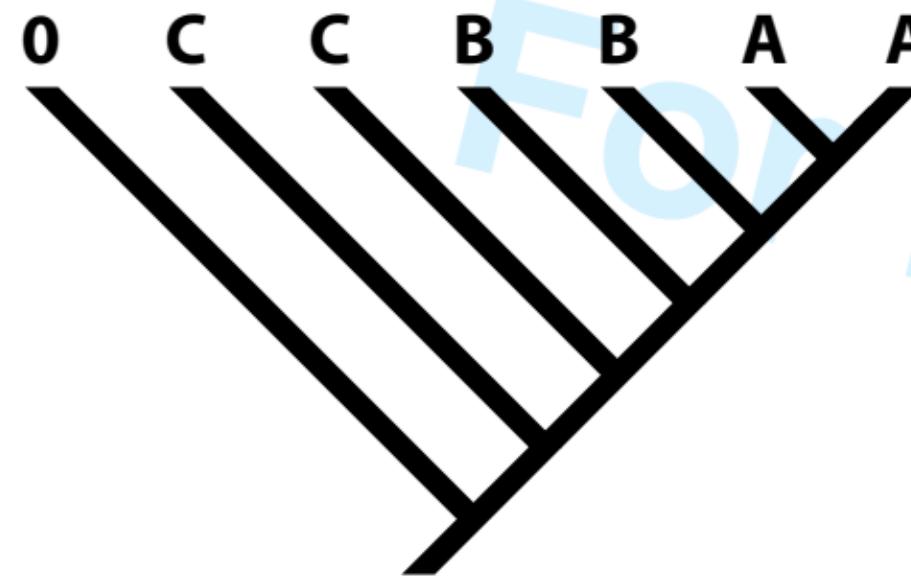


B.

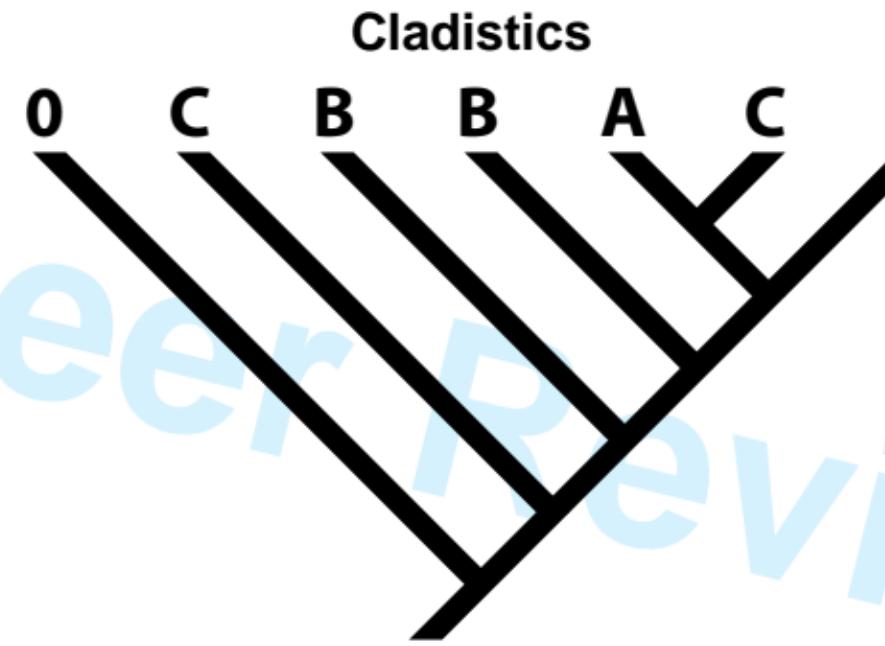


C.

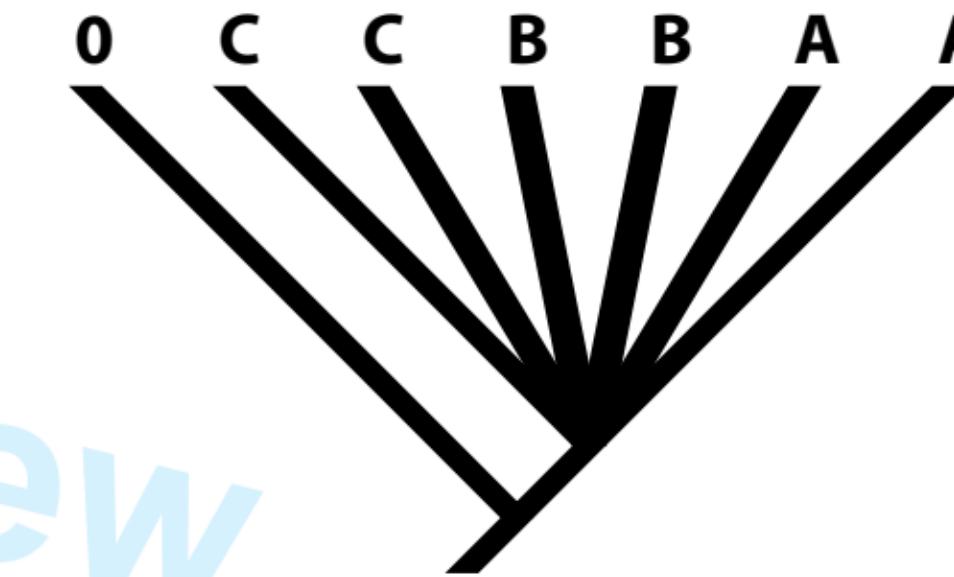




a. Max changes C to A: 0



b. Max changes C to A: 0
Cladistics



c. Strict consensus of a. and b.
Max changes C to A: 2

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 Tracking the Geographic Spread of Avian Influenza (H5N1) with Multiple Phylogenetic Trees

| Rasmus Hovmöller ^{1,2}, Boyan Alexandrov ², Jori Hardman ², and Daniel Janies ²

Deleted: A

1 Mathematical Biosciences Institute, The Ohio State University, Columbus, OH, USA
2 Department of Biomedical Informatics, The Ohio State University, Columbus, OH,
USA.

For Peer Review

Abstract.....	3
Introduction.....	5
Influenza viruses and pathogenic H5N1	5
Origin of H5N1	6
Phylogeographic analysis of H5N1	7
Materials and methods.....	9
Current datasets	9
Wallace and Fitch datasets.....	11
Phylogenetic analysis.....	12
Character optimization.....	13
Counting transmission events	15
Results	17
Alignment and phylogenetic trees	17
Reanalysis of Wallace and Fitch (2008).....	18
Frequency of discovery of transmission routes.....	19
Visualization	20
Discussion.....	20
Phylogenetic analysis and H5N1.....	20
Character optimization in phylogeographic analysis.....	21
Current state of the H5N1 epidemic	25
Perils of using a single tree	29
Problems with using multiple trees	30
Conclusions	31
Acknowledgements	33
Literature cited.....	34
Figure captions	44
Figure 1:.....	44
Figure 2:.....	44
Figure 3:.....	45

Deleted: 29
Deleted: 30
Deleted: 31
Deleted: 33
Deleted: 34
Deleted: 44
Deleted: 44
Deleted: 44
Deleted: 45

Abstract

Avian influenza (H5N1) has been of great social and economic importance since it first infected humans in Hong Kong in 1997. A highly pathogenic strain has spread from China and has killed humans in east Asia, west Africa, south Asia, and the Middle East.

Deleted: out of

Deleted: East

Deleted: West

Deleted: South

Deleted: Several recent

Deleted:,

Deleted: conclusion

Recently, several molecular phylogenetic studies have focused on the relationships of various clades of H5N1 and their spread over time, space, and various hosts. These

Deleted:,

studies examining the geographic spread of H5N1 have based their conclusions on a

single tree. This tree often results from the analysis of the genomic segment coding for

Deleted: dataset in terms

hemagglutinin (HA) or neuraminidase (NA) proteins and a limited sample of viral

Deleted: sampling of

Deleted: of the virus.

Deleted: -

Deleted:,

isolates. Here we present the first study using multiple candidate trees to estimate

geographic transmission routes of H5N1. In addition, we use all high quality HA and NA

sequences available to the public as of June 2008. We estimated geographic transmission

Deleted:,

routes of H5N1 by optimizing multistate characters with states representing different

geographical regions over a pool of presumed minimum length trees. We also developed

Deleted: develop

means to visualize our results in Keyhole Markup Language (KML) for virtual globes.

Deleted: ,

We provide these methods as a web application entitled "Routemap"

Deleted: , in virtual globes.

(<http://routemap.osu.edu>). The resulting visualizations are akin to airline route maps but

they depict the routes of spread of viral lineages. We compare our results with the results

of previous studies. We focus on the sensitivity of results to sampling of tree space,

character coding schemes, optimization methods, and taxon sampling. In conclusion, we

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2 find that using one tree and a single character optimization method will ignore many of
3
4 the transmission routes indicated by genetic sequence and geographic data.
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Deleted: the

For Peer Review

Introduction

Influenza viruses and pathogenic H5N1

Influenza viruses cause disease in humans, wild and domesticated birds and other homeothermic animals. The current classification divides influenza into three major groups: A, B, and C. Within influenza A, lineages are classified with a system of antibodies to the surface proteins of the virions, hemagglutinin and neuraminidase. This leads to the HXNY nomenclature where X includes 16 antigenic subtypes and Y represents nine different antigenic subtypes. Influenza B is only known to infect humans and pinnipeds (seals, sea lions and walruses). Influenza C is only known to infect humans and swine. On the other hand, influenza A, including subtype H5N1, has a much wider host range. Influenza A infects aves as well as many mammalian groups such as artiodactyl (swine and bovids), equids (horses), carnivorans (canids, domestic and wild felids, pinnipeds, mustelids, and viverrids). H5N1 has also been found in blow flies (Diptera: Sarcophagidae) in the vicinity of farms containing infected poultry in Japan (Sawabe et al, 2004).

Influenza viruses have two means of creating genetic diversity. The segmented genome of influenza viruses enables genomic reassortment between strains. When a host is simultaneously infected by more than one strain of influenza virions that carry gene segments from two or more original viruses, termed "reassortants" can be produced. In addition, influenza viruses are negative strand-RNA viruses. These viruses use RNA

Deleted: 9

Deleted: ,

Deleted: .

Deleted: artiodactyls (

Deleted:).

Deleted: felids (domestic and wild),
pinnipeds, artiodactyls

Deleted: .

Deleted: farms

Deleted: reassortments

Deleted: A

Deleted: can produce reassortant

Deleted: several

Deleted: . Influenza viruses

Deleted: , and on

1
2 polymerases without proofreading functions, resulting in a high mutation rate (Holland et
3 al., 1982). Both of these molecular evolutionary processes allow many lineages of
4 influenza to escape host immune responses (Smith et al., 2004) and antiviral therapeutics
5 (Bright, et al., 2005; Hill et al., 2009).

Deleted: ability

Formatted: Font: Times New

The origins of H5N1 and its variable pathogenicity.

17
18 The RNA for hemagglutinin (HA) encodes a polyprotein with two subunits that must be
19 cleaved by host endoproteases to enable the virion to fuse with the host's cell membrane.
20

Deleted: Origin

Deleted: ¶

The earliest known influenza virus characterized as serotype H5N1 was isolated from birds in Scotland in 1959 (WHO, 2004),

Deleted: first sequenced by De et al. in 1988. H5N1 had been found occasionally in wild birds in North America since the 1970s. However, thus far these isolates of H5N1 from North America appear to be strains with low

Deleted: and only infect birds (USDA, 2007).

Deleted: virions

21
22 Highly pathogenic strains of H5N1 are characterized by multiple basic amino acids at the
23 cleavage site between the two subunits of the hemagglutinin protein (Kawaoka, 1987;
24

25 Subbarao, 1998). Pathogenicity is a multigenic trait that is also mediated by features such
26 as whether the polymerase basic 2 F2 protein is produced or not (Zamarin et al., 2006)

27
28 The earliest known influenza strain characterized as serotype H5N1 was isolated from
29 birds in Scotland in 1959 (WHO, 2004) and sequenced by De et al., in 1988. Descendents
30 carrying the HA of this strain have been found in wild birds in North America since the
31 1970s (Matrosovich et al., 1999; Zhou et al., 1999; Obenauer et al., 2006). However, thus
32 far these isolates of H5 from North America appear to be strains with low pathogenicity
33 that infect only birds (USDA, 2007).

34
35 A separate highly pathogenic lineage of H5N1 has caused significant international

Deleted: lineages

Deleted: that have

36
37 concern due to their spread across Eurasia and Africa. This pathogenic lineage can be
38 traced back to a 1996 isolate from a goose in China's Guangdong province (Xu et al.,
39

Deleted: an

1
2 1999). As there was little monitoring of H5N1 in Asia before the human cases of H5N1
3
4 in Hong Kong in 1997, the deep origins of the pathogenic lineages of H5N1 remains
5
6 unclear. Further complicating the deep history of H5N1, there are few isolates of H5N1
7
8 collected in the period 1960 to 1995 that have been fully sequenced¹ and put into the
9
10 public domain.

Deleted: no

11
12
13
14 When an outbreak of H5N1 in 1997 in live bird markets in Hong Kong infected 18
15
16 people, six of whom died, this was the first indication that H5N1 had not only become
17
18 very pathogenic in birds but also that it has achieved the ability to infect humans (Xu et
19
20 al., 1999; Chan, 2002). Since 2003, there have been 438 documented human cases of
21
22 H5N1 and 60% have died (WHO, 2009). To date, the H5N1 virus has killed people in the
23
24 Middle East, west Africa, and throughout Asia.

Deleted: a 1997

Deleted: which six

Deleted: cross the species barrier to

Deleted: 394

Deleted: 63% of infected persons

27 Phylogeographic analysis of H5N1

28
29
30
31 In most phylogeographic studies of H5N1 (e.g. Ducatez et al., 2007; Guan et al., 2002;
32
33 Wallace et al., 2007; Wallace and Fitch, 2008; Wang et al., 2008), transmission routes²
34
35 have typically been determined by visual inspection of a single phylogenetic tree.
36
37 Kilpatrick et al. (2006) used a single tree in conjunction with data from observations on
38
39 migration routes of wild birds and international poultry trade to estimate opportunities for
40
41 the virus to spread across various borders. Smith et al. (2006) and Salzberg et al. (2007)
42
43 optimized geographic regions on a phylogenetic tree. These analyses have the advantage

Deleted:

Deleted: are

Deleted: or a majority rules consensus tree based on bootstrap resampling or Bayesian posterior probabilities (Ronquist and Huelsenbeck, 2001).

Deleted:

Deleted:

Deleted:

44
45
46 ¹ Most sequences from this period are of low quality and were not included in this study.

Deleted: The isolates in question all fail our

47
48 ² With the phrases "transmission route" and "transmission event", we refer to a viral
49 lineage crossing a political border in a character evolution study.

Deleted: requirements,

Deleted: are therefore

1
2 that the nodes of the tree are assigned ancestral states, enabling areas of origin to be
3
4 estimated. Janies et al. (2007) and Hill et al., (2009) optimized other features such host
5 shifts and key mutations onto a phylogenetic tree for H5N1 and projected the tree into a
6 virtual globe. Using virtual globes also allows for animation over time to reconstruct the
7 geographic spread of viral lineages. All of these studies present a single tree per genetic
8 segment or whole genome for optimization and visualization (Smith et al., 2006; Salzberg
9 et al., 2007; and Janies et al., 2007; 2008; Hill et al., 2009).

Deleted: determined.
Deleted: .
Deleted: as
Deleted: and host shifts
Deleted: and animated
Deleted: This analysis combines temporal, genotypic, phenotypic, and
Deleted: studies into one context.

17
18 In this paper, we compare our results based on multiple trees and large datasets to the
19 work of Wallace et al., (2007) and Wallace and Fitch, (2008). Wallace et al. (2007)
20 analyzed 192 HA sequences and Wallace and Fitch (2008) analyzed 482 HA and 430 NA
21 sequences. These authors used a variety of optimality criteria for tree search. This
22 research group chose to present a single parsimony tree for the 2007 paper and a single
23 maximum likelihood tree for the 2008 paper.

Deleted: .
Deleted: They choose

32
33 Here we examine whether the analysis of a single tree in phylogeographic studies of
34 avian influenza (H5N1) tends to underestimate possible routes of transmission of the
35 virus across geographical borders. We also update the HA and NA datasets for H5N1
36 with recently released sequence data and address the implication for regional spread of
37 influenza.

42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2 We use combined efficient heuristic search strategies with computing clusters to find
3 large sets a set of presumed most parsimonious trees (MPT). We also explore various
4 coding schemes for political and geographic boundaries as characters.
5
6
7

Deleted: computationally

Deleted: for

Deleted:)

Deleted: models for

Deleted: the regional

8
9
10 We developed a novel method for visualization of transmission routes of diseases. The
11 "Supramap" method of Janies et al., (2007) can produce very complex visualizations,
12 especially for large datasets, and only displays a single tree. In order to address these
13 issues we developed a technique, termed "Routemap", that allows the user to 1) display
14 only the routes of interest and 2) summarize transmission between regions implied by one
15 to many candidate trees. The essential difference between Supramap and Routemap is
16 that Supramap projects one tree whereas Routemap distills information from a pool of
17 trees. Moreover Routemap presents the results of character change studies not as a
18 projected tree but rather as an interactive map with lines representing viral transmission
19 routes between localities. We distribute the visualization method as an easy to use web
20 application called "Routemap" (<http://routemap.osu.edu>). In doing so, we aim that the
21 application will be of general utility in biomedical and natural science applications of
22 phylogeography.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38

Deleted: methodological implications
are expected to

Materials and methods

Current datasets

45 This study is based on all high quality nucleotide sequences for HA and NA genes from
46 H5N1 available in The National Institutes of Health's GenBank
47
48

Deleted: (NIH)

(<http://www.ncbi.nlm.nih.gov>) as of June 18, 2008. We add a single sequence made available by via two labs³ participating in the Global Initiative for Sharing of All Influenza Data (GISAID, <http://www.gisaid.org>).

Deleted: The
Deleted: Avian
Deleted: :
Deleted:

The following quality criteria were used to determine sequence quality: Sequences must be >75% complete, must not contain frame breaking insertions or deletions, and must not have a stop codon except at the end. Sequences from viral strains passaged through egg or other cell cultures or laboratory animals several times were excluded. When multiple sequences from the same isolate (as determined by GenBank taxonomic identification number) were available, only the most complete sequence was included. A total of 1646 HA and 1335 NA sequences were included in the final analyses. These sequence alignments are available at <http://routemap.osu.edu/publications>.

Deleted: maintained in
Deleted:, eggs, and cell cultures
Deleted: id
Deleted:

The segmented genome of influenza A viruses presents a difficulty for phylogenetic analysis. Since reassortment between diverse strains is reported to be common, two separate genomic segments are unlikely to share a phylogenetic history (Webster et al., 1992; Guan et al., 2002; Hatchette et al., 2004; Chen et al., 2006; Lindstrom et al., 1998, 2004). Reassortment violates some of the fundamental assumptions of the combined analysis (i.e. total evidence) approach such as the presumption of a shared history among

Deleted: genes
Deleted: , or
Deleted: -
Deleted:,

³ Lab for virus isolation: Mona M Aly, Abdel S Arafa from National Laboratory for Quality Control on Poultry Production, Dokki (Cairo), Egypt.

Lab for sequencing: Isabella Monne, Alice Fusaro, Giovanni Cattoli, Ilaria Capua from OIE/FAO Reference Laboratory for Avian Influenza and Newcastle Disease, Istituto Zooprofilattico Sperimentale delle Venezie, Legnano (PD), Italy.

1
2 various forms of data reflecting the history of the organismal lineages (Kluge, 1989).
3
4
5
6
7
8
9

10 Therefore, we do not combine data from HA and NA into a single dataset.
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

Multiple alignment was performed on nucleotide sequences using CLUSTALW (version
2.09 Larkin et al., 2007) using the default settings (gap opening penalty 10, gap extension
penalty 0.2, DNA transition weight 0.5). A small number of internal gaps in the
alignment were redistributed such that the final alignment did not imply any breaks in the
reading frame in ancestral sequences. The 5' and 3' ends of the alignment were trimmed
using the start and stop codons as reference. Leading and trailing gaps that remained were
replaced with question marks to indicate missing data. Internal gaps were treated as a
fifth character state. These datasets are hereafter referred to as "HA1646" and "NA1335"
datasets.

Deleted: ClustalW

Deleted: (

Deleted: '

Formatted: Font: Times New

Formatted: Font: Times New

Deleted: 5'.

Deleted: "

Deleted: "

Wallace and Fitch datasets

The HA nucleotide sequence dataset analyzed by Wallace and Fitch (2008) was recreated
from the list of GenBank accession numbers available as supplemental data to their
paper. Following the materials and methods section of Wallace and Fitch (2008), their
HA sequence dataset was aligned using CLUSTALW with the default settings without
any manual adjustment. Leading and trailing gap markers were replaced with question
marks to indicate missing data. Internal gap markers were treated as a fifth character
state. This dataset is hereafter referred to as the "HA482WF" dataset.

Deleted:)

Deleted: ClustalW (Larkin et al. 2007)

Deleted: characters

Deleted: datasets are

1 2 Phylogenetic analysis 3

4
5
6 Phylogenetic analysis was performed using TNT (Goloboff et al., 2008) in parallel and
7 sequential modes. All trees were rooted with A/Chicken/Scotland/1959/H5N1 as an
8 outgroup. This study focuses on the highly pathogenic Asian lineage of H5N1 that is
9 typically assumed to have originated with A/Goose/Guangdong/1996. We used the older
10 isolate to check this assumption and polarize geographic character state changes. A
11 preliminary analysis was performed to find the length of presumed most parsimonious
12 trees. The TNT command 'xmult = level 5 replications 100000' was used
13 for the preliminary runs on a single CPU. The argument 'level 5' allows for very
14 thorough searches and was used only to find the presumed minimum length for each
15 dataset. The preliminary runs were interrupted after 24 hours if they failed to produce
16 shorter trees than previously found.

17 Deleted: .

18 Deleted: Although this

19 Deleted: flu

20 Deleted: can be traced back

21 Deleted: , we use

22 Deleted: isolates

23 Deleted: root the phylogenies

24 Deleted: the optimizations of

25 Deleted: states.

26 Deleted: after 24 hours

27
28 To create a representative sampling of the space of MPTs, the presumed minimum tree
29 length was used to guide a new search to find 1000 independent hits to minimum length.
30
31 The TNT command 'xmult = level 0 hits 1000 giveupscore n' was
32 used, where 'n' is the tree length of the most parsimonious trees found by the preliminary
33 run.

34
35 To ensure that trees were well sampled, a cumulative series of consensus trees were
36 calculated to check that a stable consensus had been reached. The first consensus tree was
37 based on trees 0-9, the second on the first consensus plus trees 10-19, the third on the
38

39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2 second consensus and trees 20-29 [...] until all trees found from sampling 1000 hits to
3 minimum length have been included. Consensus trees were calculated using the TNT
4 command 'nelsen* [startree] [endtree];' which calculates a strict
5 consensus of all trees in the given range and saves the consensus as the last tree in
6 memory. The TNT command 'tnodes' was used to count the number of nodes that each
7 consensus tree contained to provide a metric of their stability. A consistent number of
8 nodes in a consensus tree indicate that we have collected a pool of candidate trees that
9 well represents topologies at presumed minimum length.

22 Character optimization

26 Transmission events were calculated using a modified version of the method developed
27 by Slatkin and Maddison (1989). To find transmission events, different geographic
28 regions are treated as character states of a single multistate character, and are mapped
29 onto trees obtained from phylogenetic analyses using standard optimization methods
30 (Farris, 1970; Fitch, 1971; Swofford and Maddison, 1987).

Deleted: The method assumes an island biogeographical model, where each location is equally accessible from every other location.

38 To make our results comparable with previous phylogeographic studies of H5N1, we use
39 equal weight between geographical character states. However, if preferred, the multistate
40 character can be coded as a Sankoff character with a transition matrix specifying the
41 relative costs of each type of transition between character states (Sankoff and Rousseau,
42 1975).

We counted the minimum and maximum number of transmissions between all pairs of geographic character states. When a character state change is observed, this is interpreted as a possible transmission route of the virus. The state changes are directional, so results are expected with different numbers for transformations from state A to state B and state B to state A. Figure 1 illustrates how state changes are counted.

Geographic information for the molecular datasets was extracted from the

'TSeq_orgname' field in GenBank's TinySeq XML records. For each molecular dataset (HA1646, NA1335, and HA482WF), new geographic character matrices were created. Each geographic character matrix contained an entry for a multistate character for each taxon in the nucleotide dataset. The order of taxa was maintained in nucleotide and geographic character datasets as this is important for consistency when using compact tree format and other tree formats that use numerical taxon names in TNT.

- Deleted:** the GenBank
- Deleted:** record.
- Deleted:** and
- Deleted:** containing
- Deleted:** kept the same

Three geographic character data sets of 31 character states were used to find transmission

routes. Due to the limitations of the character state space allowed by the tree search

software, TNT, we were unable to issue a single set of geographic character states

suitable for comparing the HA1646 and NA1335 datasets to some previous studies.

- Deleted:** phylogenetic
- Deleted:** all

However, we were able to compare our results to Wallace and Fitch (2008). The first geographic character set we created, termed "W&F", was used to assign each taxon to a geographic category matching those used by Wallace and Fitch (2008). For taxa collected in localities not matching any of the 28 character states used in W&F, we assigned

1 character states to these taxa with the most fitting of three new states: "Other China",
2 "Other Asia" and "Other World".

Deleted: OtherChina", "OtherAsia

3
4 We created a second character set termed "Africa" that focuses on relationships and
5 sources of H5N1 in regions of Africa and Europe. A third character set "World" was used
6 to look at broader patterns among regions of the globe. The character set World treats
7 mainland China as a single entity, treats Europe as four regions (north, east, south and
8 west), and all of Africa as a single region.

Deleted: OtherWorld". The

Deleted: "

Deleted: "

Deleted: North, East, South

Deleted: West

Counting transmission events

For each count of transmission events, a geographic character dataset containing a single
multistate character and compact tree file containing the MPTs found during the search
on the corresponding nucleotide data were loaded into memory. The TNT command
change was used to calculate the maximum number of character state changes for pairs
of geographic character states. The syntax used was 'change]./[n]/[region X]
[region Y]'. The ']' indicates that results should be displayed character by character,
and the './[n]' denotes that character number 'n' should be optimized in all trees in
memory. For example, the command 'change]./0/HongKong Fujian' will
return the global minimum and maximum number of transmission events from Hong
Kong to Fujian, calculated over all MPTs, when only a single multistate character is used.

Deleted: .

Each direction of change among states (e.g., A to B and B to A) was calculated for all
possible pairs of states in the dataset.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2 From the output of the 'change' command, tables were created containing the global
3 minimum and maximum values for each location pair and direction. When the number of
4 state changes was >0, this result was interpreted as a possible transmission route from
5 one region to another.
6
7
8
9

10
11
12 Transmission events were calculated for the HA428WF sequence dataset using the W&F
13 character state set. For 1646HA and 1335NA sequence datasets, all three (W&F, Africa,
14 and World) character state sets were used.
15
16

17
18
19 We have created a web-application, Routemap, (<http://routemap.osu.edu>) to produce
20 visualizations of transmission events from the output of the geographic character
21 optimization studies over pools of trees. The required inputs to Routemap are 1) a table of
22 comma-separated values (CSV) containing the taxon label, the geographic place for the
23 isolates, the latitude, and longitude in decimal degrees, and 2) a matrix of aligned
24 nucleotide or amino acid data for the isolates in FASTA format. As an option, the user
25 can input a set of precomputed trees. In the first step, Routemap will produce a TNT file
26 that the user can execute on their machine with the command 'tnt p filename,'.
27 The user then will load the TNT output back into the Routemap via the browser, where
28 the underlying application will complete the geographic calculations. Routemap's final
29 output is a Keyhole Markup Language (KML) file compatible with Google Earth
30 (<http://earth.google.com>) and other geographic information systems such as ESRI's
31 ArcGIS Explorer (<http://www.esri.com/software/arcgis/explorer/index.html>).
32 A visualization produced by Routemap is similar to a web-based airline route map. When
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Deleted: wrote
Deleted: computer program in Java
Deleted: www.java.com

Deleted: input
Deleted: the program is two tables
Deleted:). One table contains
Deleted: potential
Deleted: sources
Deleted: disease agent in rows
Deleted: geographic destination for the disease agent in columns
Deleted: each cell contains the number of character transformation events between two localities discovered by the character optimization process represented by an integer. Another table contains the latitudes and longitudes of all the source and destination localities in the first table expressed
Deleted: Java program's

Deleted: ESRI's
Deleted:)

1
2 using a Routemap in a virtual globe the user selects localities of interest and viral
3

Deleted: NASA's Worldwind
(<http://worldwind.arc.nasa.gov>). The KML displays possible

4 transmission events implied by the analysis of phylogenetic trees are displayed as
5

Deleted: a source or going to a

6 radiating from origins or arriving at destinations. The user can make one or more origin
7

Deleted: . Colors are used

8 and destination selections to make as complex or as simple of a visualization as needed.
9

Deleted: the direction of travel

10 The lines indicating routes are colored to indicate directional, ambiguous or bidirectional
11 travel. Example datasets, KML files, scripts and instructions for users to make their own
12

Deleted: Datasets

13 Routemaps are available at <http://routemap.osu.edu>.
14

Deleted: people.mbi.ohio-state.edu/rhovmoller/

Deleted: /.

17 Results

21 Alignment and phylogenetic trees

25 The alignment of the HA1646 dataset consisted of 1646 taxa and 1716 characters. 937
26 characters in HA1646 were parsimony-informative. The length of the presumed most
27 parsimonious trees was 7284 steps. A search for 1000 independent hits to the MPT length
28 resulted in 2721 unique trees. The number of nodes in the consensus stabilized at 500
29 trees at 695 nodes.
30
31
32
33
34

37 The alignment of the NA1335 dataset consisted of 1335 taxa and 1410 characters. 744
38 characters in NA1335 were parsimony-informative. The length of the presumed most
39 parsimonious trees was 4795 steps. A search for 1000 independent hits to the MPT length
40 resulted in 2649 unique trees. The number of nodes in the consensus stabilized at 860
41 trees at 540 nodes.
42
43
44
45
46
47
48
49

The HA482WF dataset consisted of 482 taxa and 1809 characters. 937 characters in HA482WF were parsimony-informative. Under our search strategy, the length of the presumed most parsimonious trees is 3257 steps. A search for 1000 independent hits to the MPT length resulted in 2470 unique trees. The consensus stabilized at 251 nodes after 140 trees were collected.

Reanalysis of Wallace and Fitch (2008)

Results of the reanalysis of Wallace and Fitch's (2008) HA nucleotide dataset (HA482WF) show that using a single tree, and a single optimization path, misestimates the frequency of transmission events. Moreover the use of a single tree can fail to detect possible transmission events. Figure 2 illustrates a comparison of the transmission events reported in Wallace and Fitch's (2008) with our reanalysis of HA428WF and W&F character data for localities. We used the global maximum number of transmissions from the pool of 2470 MPTs in our reanalysis. The global maximum of implied transmission events corresponds to Wallace and Fitch's use of DELTRAN optimization (Swofford and Maddison, 1987; Swofford, 1990). Thus a transmission route is indicated if it is found in any MPT.

Deleted: and in some cases fails

Deleted: among regions.

Deleted: most closely

Deleted:) optimization.

There are 756 possible transmission routes in W&F character state data. We compare the results found by Wallace and Fitch's (2008) optimization of these character data over a single tree they derived from HA482WF sequence dataset to our reanalysis using 2470 MPT for HA482WF. 517 of the possible routes are not found by either study, one transmission route result is found only in Wallace and Fitch (2008), 169 new routes were

discovered by our reanalysis, and 69 routes are suggested by both Wallace and Fitch (2008) and our reanalysis. Detailed results are presented in supplemental data at <http://routemap.osu.edu/publications>.

Deleted: people.mbi.ohio-state.edu/rhovmoller/

Deleted: /.

We also analyzed W&F based on all the MPTs we considered for HA428WF or HA1646 sequence datasets. Of 756 possible routes in W&F, 528 routes were not found when optimizing these character data over MPTs we found for the HA428WF or HA1646 sequence datasets. Four transmission routes are found only in HA428WF, 158 routes are found only in HA1646, and 61 routes occur in both HA428WF and the HA1646 dataset. Thoughtful examination of transmission routes implied by the HA1646 sequence dataset required addition of three character states for viral isolates from regions not included in the HA428WF dataset. 111 additional transmission routes were detected in MPTs implied by HA1646 when we included three localities not considered by Wallace and Fitch,

(2008) but considered by other authors (e.g., Ducatez et al., 2006; Salzberg et al., 2007).

Deleted: .

Frequency of discovery of transmission routes

We compared the transmission routes found by Wallace and Fitch (2008) for a single tree to the pool of MPTs we derived from the HA428WF dataset. Wallace and Fitch (2008) did not find 22% of the transmission routes that are found in 95% or more of the MPT that we found for HA428WF. See supplemental data at <http://routemap.osu.edu/publications> for a summary table.

1 2 Visualization 3

4
5 The panels in Figure 3 contain visualizations produced with Routemap of viral traffic in
6 Africa, east Asia, and southeast Asia from the HA1646 dataset and the World geographic
7 character coding scheme. Visualizations of the HA1646 and NA1335 dataset under the
8 W&F, World, and Africa geographic character coding schemes, as well as the HA428WF
9 dataset are available at <http://routemap.osu.edu/publications>.
10
11
12
13
14
15
16
17
18
19

Deleted: Figure 3. contains a visual summary of the transmission routes concerning Indonesia, South Korea, Japan and Africa using our system for visualization.¶

20 Discussion 21

22 Phylogenetic analysis and H5N1 23

24
25 Metadata on pathogens such as date, host, and geographic location for viral isolation are
26 often combined with phylogenetic trees based on pathogen sequence data by researchers
27 and public health officials during an epidemic. In order to interpret metadata in a
28 phylogenetic context, character optimization is a necessary tool. Studies that use
29 metadata and character optimization include investigations of where H5N1 spreads
30 (Wallace et al., 2007, and Wallace and Fitch, 2008; Salzberg et al., 2007), where and
31 when H5N1 has switched hosts (Janies et al., 2007) or experienced mutations that confer
32 resistance to drugs (Hill et al., 2009).
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Deleted: during

Deleted: examinations

Deleted: investigations on

Deleted:), and to answer questions as to where and when

Deleted: evolved in H5N1

Character optimization in phylogeographic analysis

Here we demonstrate that when reporting transmission events from character state changes, a single tree and a single way of optimizing the character will only tell a fraction of the story. Wallace et al. (2007) and Wallace and Fitch (2008) used MIGRAPHYLA

Deleted: ¶

(HoDac et al., 2007) and PAUP* (Swofford, 2008) to calculate transmission events,

Deleted: .

Deleted: MigraPhyla

interpreted as character state changes using a single multistate character with

geographical regions as character states. As Wallace and Fitch (2008) point out: "PAUP's

Deleted: outline

assignments may represent only one of several possible most-parsimonious traces". Thus

Wallace and Fitch (2008) inferred transmission events from a single possible

optimization out of many of a 28-state character on a single tree out of many possible

trees. When we examine the length of this 28-state character on the pool of 2470 MPTs

from the HA482WF sequence dataset, the lengths are between 107 and 118 steps. This

range of lengths is consistent with the 115 steps reported in Figure S2 in Wallace and

Fitch (2008). On the other hand, when we examine the average number of transmission

Deleted: ¶

events supported by the MPTs we found for the HA482WF dataset, the numbers of

transmission events are very different. Each MPT supports on average 131 different

transmission routes in contrast to the 70 unique transmission routes reported by Wallace

and Fitch (2008).

The difference in results between TNT and PAUP* is because PAUP* only considers a single optimization path of the character, when calculating steps on polarized character state changes. This is the behavior of PAUP* even in cases where there are multiple equally parsimonious optimization paths for the character and each path may support

Deleted: This difference in results is because PAUP* only considers a single optimization path of the character, when calculating steps on polarized character state changes. In contrast, in TNT it is possible to calculate the maximum (or minimum) state changes between two specified states of a multi state character, from multiple trees using multiple optimization paths. Although the sum of all state changes calculated this way will exceed the minimum length of the character, it is a more efficient method of finding possible directional state changes. It is evident that when reporting transmission events from character state changes, a single tree and a single way of optimizing the character will only tell a fraction of the story. We also discovered an artifact in the case of closely related taxa.

1
2 different sets of polarized state changes. In contrast, in TNT it is possible to calculate the
3 maximum (or minimum) number of polarized state changes (e.g., from state A to state B)
4 between two states of a multistate character. Moreover, in TNT this calculation can be
5 done over multiple trees and multiple optimization paths per tree. Although the sum of all
6 state changes calculated this way with TNT will exceed the minimum length of the
7 character, it is important to find all possible state changes.
8
9
10
11
12
13

14
15
16 We also discovered an artifact in the case of closely related taxa. The number of
17 observed state changes between two states should not necessarily be interpreted as the
18 number of implied transmission events. In cases where the spread of an infectious agent
19 is rapid, it is very likely that isolates from neighboring regions will be genetically
20 identical or fail to have obtained any parsimony informative mutations. In such cases, the
21 number of observed state changes in the MPTs in a clade with no supported branches
22 (*sensu* the command collapse 3 in TNT) will only reflect the proportion of taxa from
23 each location, when optimizing over a pool of dichotomous, MPTs. In other words, given
24 a character with two geographic states (location A and location B), the maximum number
25 of transmission events between A and B will be the maximum length of the character in
26 the subtree. For example, if there are 10 taxa from location A and five from location B in
27 a subtree, the maximum number of transmission events will be five, when counting
28 changes from A to B or from B to A.

Deleted: there are rapid spreads

Deleted: only carry

Deleted: non-

Deleted: differences.

Formatted: Font: Italic

Deleted: . Given

Deleted: place

Deleted: place

Deleted: either

Deleted: the other way around.

45 In cases where there are many sequences in a region where there is little or no genetic
46 difference, the artifact we describe above can occur. In other cases, where distinct genetic
47
48
49
50
51
52
53
54
55
56
57
58
59
60

lineages have been exchanged between two places, counting the changes in geographic characters is a useful tool to understand viral traffic. Both the HA428WF and HA1646 sequence datasets contain large subclades (e.g., in Africa) of viruses with identical genetic sequences. The recent invasion of H5N1 into Africa without much molecular evolution in HA makes it of interest for future studies to examine these isolates from a multi locus approach as done by Salzberg, et al., (2007) and Janies et al., (2007). These concepts are also applicable to pandemic H1N1, which emerged in early 2009 and has spread rapidly around the world.

Deleted: of
Deleted: can be
Deleted: ¶
Deleted: .

Deleted: in (

Current phylogenetics software (PAUP*, Swofford, 2008; TNT, Goloboff et al., 2008) limit a categorical character to 31 states. These programs are suitable for most analyses of morphological and sequence data in which the data are used to find trees. MESQUITE (Maddison and Maddison, 2007) can summarize character state changes over pools of trees. MESQUITE is thus an alternative for these analyses, but is limited to 56 states and would have to be coupled to a strong tree search program. POY4 (Varón et al., 2008)

Deleted: is
Deleted: Mesquite
Deleted:
Deleted: and
Deleted: tool that can be explored
Deleted: . POY

may provide a solution to the state space limitation as the user can define any alphabet to describe states for features of interest of pathogens and transformation costs among states in a look-up table (e.g., a Sankoff matrix). Heuristic efficiency will be important as an accurate regional map may require hundreds or thousands of geographical character states.

Deleted:
Deleted:
Deleted:; Sankoff and Rousseau, 1975
Deleted: This is an active area for future research, as trading off coarseness with which the world is divided into a grid to save computational space and time is a challenge for many modelers across the sciences (Lindley, 2009).

When a widely distributed taxon such as influenza A is studied, geographic character coding requires some abstraction of the globe into regions. We followed the 28-location

Deleted:

1
2 state coding scheme of Wallace and Fitch (2008), as well as implementing two alternative
3 coding schemes (i.e. Africa and World) to make our results comparable to other studies.
4

Deleted: 'Africa'

Deleted: 'World'

5 When comparing the HA1646 and NA1335 sequence datasets to the results of Wallace
6 and Fitch (2008), 275 isolates were not compatible with the 28 states defined by Wallace
7 and Fitch (2008). We coded the isolates that fell outside of the character state space with
8 one of three new character states: Other Asia, Other China, and Other World. This coding
9 scheme allowed our analyses to find transmission routes outside the scope of Wallace and
10 Fitch's (2008) original publication, while retaining comparable results.
11

Deleted: other

Deleted: other

Deleted: other

Deleted: allows

12 By including more isolates from all regions sampled by HA1646 and NA1335 sequence
13 datasets, we identified more possible routes of transmission than Wallace and Fitch
14 (2008). For example, when using the same character data as W&F, our results agree with
15 those of an independent study, Wang et al. (2008), in finding a single unambiguous
16 transmission from the Hunan province of China to Indonesia. Wallace and Fitch (2008)
17 suggested a transmission from Guangdong, China to Indonesia. The Guangdong to
18 Indonesia route is not observed under any optimization of character data in W&F across
19 the MPTs we found for the HA1646 sequence dataset.
20

Deleted: found

Deleted: routes, as implied by the data,

Deleted: .

Deleted: MPT

21 In all studies of the phylogeography of influenza there are problems associated with the
22 imprecise naming of viruses. For example in the HA1646 sequence dataset, 18 viruses
23 are labeled "China" which represents an area of 9,598,08 km². This area is comparable to
24 10,180,000 km² for all of Europe, including Russia west of the Ural Mountains. Other
25 viral isolates are named for cities. In some cases viral isolates are named for regions.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Deleted: "

Deleted: "

Deleted: is

Deleted: mountains.

1
2 Many workers, including Wallace have suggested that GenBank records for viral isolates
3 should include latitude and longitude records to avoid geographic imprecision (Butler,
4 2008). Similarly, there is much room for improvement in metadata standards for naming
5 of host taxa (Janies et al., 2007). In the wake of the H1N1 pandemic, GISAID has vastly
6 improved metadata standards. High quality metadata will improve research in viral
7 character evolution and geography studies.

8
9
10
11
12 We used the results from HA1646 and NA1335 sequence datasets, and character states
13 from the World coding scheme to find transmission routes discussed in the literature. The
14 major difference between W&F and World is that mainland China is treated as a single
15 region in the World character state set, and countries in Europe and Africa are separate
16 into regions (discussed above). Consolidation of regions within China was necessary to
17 allow all countries (as defined by political borders) in east Asia to fit into the 31 character
18 state limit.
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

Deleted: 'world'

Deleted: states

34 Current state of the H5N1 epidemic

35
36
37
38 Since 2003, H5N1 has been endemic in Indonesia. All phylogenetic studies including HA
39 sequences isolated from H5N1 in Indonesia have found a monophyletic clade of
40 Indonesian isolates, independent of optimality criterion (neighbor joining: Guan et al.,
41 2004; Smith et al., 2006; Salzberg et al., 2007, Bayesian: Wang et al., 2008; Smith et al.,
42 2006, maximum likelihood: Kilpatrick et al., 2006; Wallace and Fitch, 2008; Hill et al.,
43 2009, parsimony: Janies et al., 2007; Salzberg et al., 2007; Wallace et al., 2007).

44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 It is apparent that the H5N1 strains within Indonesia occurred due to a single
5 transmission event from southern China. Janies et al., (2007) use whole genomes to
6 demonstrate this link. Guan et al. (2004) use HA data to suggest the closest relatives to
7 isolates from Indonesia are viruses from Yunnan province in China. Wallace et al.
8 (2007), and Wallace and Fitch (2008) also us HA data to find that Guangdong province in
9 China is the most likely source. These results stand in contrast to Kilpatrick et al., (2006),
10 who used data from international trade to assert that the source of H5N1 in Indonesia was
11 wild birds imported from Thailand. Our analyses of sequence and geographic data for
12 H5N1 are in disagreement with the assertion of Kilpatrick et al. as we never find a
13 Thailand to Indonesia route. Our analyses of the HA1646 sequence dataset, suggest that
14 the geographic origin of H5N1 for Indonesia was Hunan province in China. Wang et al.
15 (2008) concur. In our analyses of HA1646 sequence data, the Hunan-Indonesia migration
16 route appears in all MPTs independent of optimization method.

Deleted: endemic

Deleted: singe

Deleted: While

Deleted:

Deleted: China's

Deleted:,

Deleted:

Deleted:) as well as

Deleted: use

Deleted: indicate

Deleted: to be

Deleted: Pathogen

Deleted: is

Deleted: this

Deleted: as

Deleted: ¶
¶

Deleted: source

Deleted: into

Deleted: appears to be

Deleted:

32
33 Similarly our analyses of the NA1335 sequence dataset suggest Hunan as a possible
34 source for Indonesia but some of the trees for NA1335 also indicate Hong Kong and the
35 Guangxi Zhuang Autonomous Region of China as possible sources for Indonesia.

Deleted: In

Deleted:, all MPTs contain

Deleted:. Some

Deleted:,

Deleted:,

Deleted: This

Deleted: viruses

Deleted: be

Deleted: a reflection of

36
37
38 Conflict between genes result could represent a reassortment event in which viruses in
39 Indonesia have HA and NA genes with distinct histories or could simply reflect
40 ambiguity in the data.

Smith et al. (2006) found that all H5N1 from Vietnam, Cambodia and Thailand form a single clade, except for one isolate from Vietnam. We were unable to find this clade structure in any tree found among the most parsimonious trees from the analyses of the HA1646 or NA1335 datasets. Although most isolates from Vietnam, Thailand, Laos, Cambodia and Vietnam are found in a monophyletic clade, isolates from Vietnam and Thailand are always found outside of this clade. This result suggests that the results of Smith et al. (2006) are an artifact of limited taxon sampling.

In the pool of trees resulting from searches of the HA1646 sequence data, we see at least three separate introductions of H5N1 to Japan. There were transmission events from China to Japan in 2003 and 2008. Under certain optimizations, Vietnam is also a possible source of the 2008 isolates in Japan. In contrast, the source of the 2004 outbreak in Japan appears to be South Korea in all trees and all optimizations based on analysis of HA1646 sequence data and the W&F character data.

These comparisons illustrate the danger in excluding data. Wallace et al. (2007) and Wallace and Fitch (2008) both find the China-Japan connection, but excluded isolates collected in 2002-2003 from South Korea despite the fact these data were available in public databases as of 2005. Due to this exclusion, Wallace et al. (2007) and Wallace and Fitch (2008) did not find South Korea as a source of the Japanese outbreak in 2004. Our phylogenetic analysis indicating transmission of H5N1 from South Korea to Japan in

- Deleted:**
- Deleted:**) and Guan et al. (2004)
- Deleted:** isolates
- Deleted:** southeast Asia (
- Deleted:**, Laos, Malaysia, Myanmar, Singapore,
- Deleted:** and Vietnam)
- Deleted:** monophyletic clade. In contrast we are
- Deleted:** support for
- Deleted:** resulting
- Deleted:** searches on HA6146
- Deleted:** sequence data. It is true that most isolates from southeast Asia are found in a single clade in both the HA and NA consensus trees, but we find in all trees a few
- Deleted:** Malaysia
- Deleted:** Vietnam
- Deleted:** Laos placed external to the main clade comprised of
- Deleted:** southeast Asia. This leads us to believe
- Deleted:** findings
- Deleted:**
- Deleted:**) and Guan et al. (2004
- Deleted:**
- Deleted:** mainland
- Deleted:** from
- Deleted:** 'W&F'
- Deleted:**

- Deleted:**

1
2 2004 is in agreement with Kilpatrick et al. (2006) who based their inference on trade and
3 bird migration data.
4
5

Deleted: .

6
7 The H5N1 outbreaks in Africa have been investigated using HA sequences (Ducatez et
8
9 al., 2006) as well as full genomes (Ducatez et al., 2007, Salzberg et al., 2007). However,
10
11

12 these studies were conducted when there was a limited set of putative outgroup(s) for the
13

Deleted: both

Deleted: non-

Deleted: isolates

Deleted: Both

14 African lineages in the public domain. These studies conclude that the sources of H5N1
15

16 are not southeast Asia, but rather central Asia and Europe. The pools of trees resulting
17

18 from our analysis of HA1646 agree with Ducatez et al. (2006; 2007) that H5N1 has been
19

Deleted: .

20 introduced to Nigeria on three separate occasions. Using the "Africa" character state set,
21

Deleted: 'Africa'

22 we find corroboration for the sources to be central Asia or Europe rather than southeast
23

Deleted: of not southeast origin but

Deleted: .

24 Asia. By examining at all trees and all optimizations we found for HA1646 and the W&F
25

26 character state set, we can not rule out the source of one transmission to Africa was the
27

28 Qinghai province of China (Salzberg et al., 2007).
29
30

31 Sources for H5N1 in Europe are unclear⁴, as multiple possible sources are indicated for
32

33 most countries. The basis for the ambiguity is the high similarity of the isolates found in
34

35 Europe, Africa, and the Middle East in 2005-2006. This is another empirical example of
36

Deleted: middle east

37 the artifact we described above, where very low genetic variation leads to every region
38

Deleted: an

39 being a possible source for its neighbors.
40
41

42
43
44
45
46
47 ⁴ One important exception is the case of an infected eagle (*Nisaetus nipalensis* Hodgson,
48 1836) that was smuggled from Thailand to Belgium but did not cause any further spread
49 (Van Borm et al., 2005).
50
51
52
53
54
55
56
57
58
59
60

1 2 Perils of using a single tree 3

4
5
6 A single tree is likely to miss several transmission routes found frequently among the
7 MPTs, as illustrated in Fig. 4. [All of the publications on phylogeographic analysis of](#)
8 [H5N1 we have reviewed here use a single tree for the results they report. This practice is](#)
9 [independent of whether they use Slatkin and Maddison \(1989\) transmission event](#)
10 [analysis \(Wallace et al., 2007; Wallace and Fitch, 2008\) or other optimization methods](#)
11 [\[POY3 \(Wheeler et al., 2005\) command -diagnose in Janies et al., 2007; MESQUITE](#)
12 [command trace:parsimony optimization in Salzberg et al., 2007\] or observe](#)
13 [transmission events by inspection \(e.g., Ducatez et al., 2006, 2007; Wang et al., 2008\).](#)
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Deleted: All of the publications on phylogeographic analysis of H5N1 we have reviewed here use single tree for the reported results. This practices is independent of whether they use Slatkin and Maddison (1989) transmission event analysis (Wallace et al., 2007; Wallace and Fitch, 2008) or observe transmission events by inspection (e.g. Ducatez et al., 2006, 2007; Wang et al., 2008), and which optimality criteria they use tree search. In some cases (e.g. Salzberg et al., 2007; Smith et al., 2006; Wang et al., 2008), only a single neighbor-joining tree is produced.

The use of one tree is also independent of which optimality criteria researchers use in tree search. In some cases (e.g., Smith et al., 2006; Wang et al., 2008) only a single neighbor-joining tree is produced. Janies et al., (2007) created a binary tree on which to optimize based on a strict consensus of MPT using the TNT command randtree*. In other cases, information from multiple trees is represented as a consensus tree. A consensus tree can be a majority rules consensus from bootstrap resampling under maximum likelihood (Kilpatrick et al., 2006; Smith et al., 2006), parsimony (Guan et al., 2004; Salzberg et al., 2007;) or a distribution of trees from a Bayesian analysis (Smith et al., 2006).
Deleted:).

Deleted: In other cases, information from multiple trees are condensed into a consensus tree, this

There are several problems with the use of a consensus tree in the context of character optimization. First, a consensus tree does not represent all MPTs. Also the consensus

Deleted: trees

Deleted: a

1 topology is never one of the most parsimonious trees. Furthermore, trees can be
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Deleted: a shape that was not found during the tree search itself.

constructed that do not change the consensus tree, but still are not found among the

MPTs. Finally, a set of most parsimonious trees can be found, where the strict consensus tree would imply transmission routes that are never found in any of the individual MPTs, (as illustrated in Figure 4). It is thus inappropriate to optimize character state changes on a consensus tree.

Under parsimony, each tree in the pool of MPTs and each optimization of the character states are just as optimal as the next one. The transmission routes found in only a few trees and under certain optimizations are just as supported by the optimality criterion as those trees and optimization found frequently. However, under a model-based optimality criterion, probability calculations for the existence of transmission routes could be calculated. A distribution of trees from a Bayesian analysis, where each tree has a posterior probability could be used to provide confidence values for transitions between the character states in a multistate character for geography optimized on trees of various posterior probability. Alternatively the proportional likelihood can be calculated for ancestral geographic states. One could employ a method such as described to study the evolution morphological characters under likelihood (Lewis, 2001) and is implemented in MESQUITE (Maddison and Maddison, 2003).

Deleted: Each

Deleted: is

Deleted: Similar to the fallacy of interpreting percentage in majority rules consensus as support for clades, the rarer

Deleted:,

Deleted: in

Problems with using multiple trees

When examining character state changes over a pool of MPTs, different types of

optimization need to be taken into account. If one examines the global minimum of state

Deleted: When we look at

1 changes over a pool of trees, they will find only those state changes that appear in all
 2 trees under all optimizations. When dealing with large numbers of trees, a tree will often
 3 be found where a specific state change does not occur. For example, one has a subtree
 4 with only unsupported branches and 100 taxa from location A, 100 from location B and a
 5 single taxon from location C. When looking at this subtree in the MPTs, one can find a
 6 tree and an optimization pathway that does not have a transmission event between
 7 location A and B, making the global minimum number of transmission events from A to
 8 B equal zero.

Deleted: (transmissions)
Deleted: we
Deleted: and short branch lengths, for nearly every type of state change
Deleted: can
Deleted: when looking at the minimum number of state changes. ¶
 ¶ Say you have
Deleted: you

Deleted: 0

19 In contrast, use of the global minimum number of polarized state changes dramatically
 20 changes the number of indicated transmission routes for a set of trees when compared to
 21 use of the maximum number of state changes. For example, when considering the
 22 maximum number of directional state changes on the HA1646 dataset, we find a total of
 23 267 different transmission routes. This results contrasts with a result of 20 transmission
 24 routes found when looking at the global minimum number of directional state changes. In
 25 summary, the global minimum is a severe test of putative transmission routes as the route
 26 has to appear in every tree in the pool of MPTs to register a global minimum > zero. The
 27 global maximum is useful to find the range of possible transmission routes a set of MPTs
 28 implies, irrespective of how frequently a route is found.

Deleted: When optimizing the 'world' geographic character state set on HA1646 under global minimum number of state changes we find only 20 supported transmission routes (directional changes from one state to another). In contrast, when looking at the same character state set and pool of trees, but looking for any transmission in any tree when maximizing the number of state changes we find 267 routes. Wallace and Fitch (2008) found a total of 70 transmission routes, when using DELTRAN optimization (Swofford and Maddison, 1987). ¶

41 Conclusions

42 Phylogenetic tree search and character evolution analyses have been valuable tools in
 43 many fields, including surveillance of infectious diseases. In this context, parsimony is

Deleted:, geographic,
Deleted:. Here we demonstrate that
Deleted: analysis,

1 sometimes used for evaluation of trees and almost always used for character
2

3 optimization. The multiple equally optimal trees produced by parsimony reflect of the
4 natural ambiguity of the sequence data. Character optimization under parsimony can find
5 the global minimum and maximum values for state changes in a pool of trees reflecting
6 various degrees of ambiguity underlying transmission routes for organismal lineages.
7

8 With efficient character optimization and tree search heuristics it is possible to find and
9 analyze very large pools of trees for datasets with thousands of terminal taxa. Moreover,
10 because results are sensitive to tree choice, any tree search method that considers multiple
11 trees has an advantage over methods that consider a single tree, such as neighbor joining.
12

13 The methods of character optimization and visualization used and developed here are also
14 applicable to sets of trees produced by maximum likelihood and Bayesian methods.
15

16 Maximum likelihood methods are commonly used for tree search and sometimes
17 preferred for the operational reason that they support one tree slightly more than near
18 suboptimal candidate trees. While Bayesian methods produce a distribution of trees,
19 conclusions are typically discussed in relationship to the single tree presented (often a
20 majority rules consensus and sometimes a maximum *a posteriori* tree). As we have
21 demonstrated here, a single tree, as well as consensus trees are insufficient for estimating
22 transmission routes of infectious disease. In addition to considering multiple trees, when
23 using Slatkin and Maddison's (1989) method of geographic analysis, it is important to
24 take into account sensitivity to different strategies for character state optimization.
25

26 **Deleted:**, and tree search algorithms that employ heuristics enabling analysis of large datasets are key technologies for infectious disease surveillance and research. Parsimony is computationally fast and with tree search heuristics can find very large pools of candidate trees that form a representative sample of the space of possible topologies and the hypotheses they support.¶
27

28 **Deleted:** with respect to the questions of interest.
29

30 **Deleted:**.

31 **Deleted:** we use

1
2 It is our intention that the methods described in this paper should be general and useful
3
4 for phylogeographic analysis of any type of problem above and beyond infectious
5
6 diseases. We welcome everyone to our website <http://routemap.osu.edu> to test our
7
8 methods and create KML visualization of your own data.

13 **Acknowledgements**

18 We thank the Pablo Goloboff, Steve Farris, Kevin Nixon, and the Willi Hennig Society
19 for making TNT available. Pablo Goloboff, Universidad Nacional de Tucumán and
20 Diego Pol, CONICET, Museo Paleontológico Egidio Feruglio provided insights and
21 discussions on character analysis. We also thank two anonymous reviewers for their
22 constructive criticisms.

Deleted: authors

Deleted: Tucuman

Deleted: for

Deleted:

30 We acknowledge the Department of Biomedical Informatics and The Mathematical
31 Biosciences Institute of The Ohio State University (OSU) for space, funding, and
32 administrative support. We thank the Medical Center Information Services team of OSU
33 and the Mathematical Biosciences Institute for hosting computing clusters used in this
34 study. We acknowledge this material is based upon work supported by, or in part by, the
35 U.S. Army Research Laboratory and Office under grant number W911NF-05-1-0271.

Deleted:). We thank The

Deleted: Medical Center Information Services team of OSU and the Mathematical Biosciences Institute hosts for hosting computing clusters used in this study. We acknowledge this material is based upon work supported by, or in part by, the U.S. Army Research Laboratory and Office under grant number W911NF-05-1-0271. ¶

Deleted: .**Literature cited**

Bright, R.A., Medina, M., Xu, X., Perez-Oronoz, G., Wallis, T.R., Davis, X., Povinelli, L., Cox, N., Klimov, A. , 2005. Incidence of adamantane resistance among influenza A (H3N2) viruses isolated worldwide from 1994 to 2005: A cause for concern. 2005. Lancet 366,1175-1181.

Butler, D. 2008. Politically correct names given to flu viruses. Nature News online 23 April 2008. doi:10.1038/452923a .

Chan, P. Outbreak of Avian Influenza A(H5N1) Virus infection in Hong Kong in 1997 (2002). Clin. Infect. Dis. 34 (Suppl 2), S58–64.

Chen, H., Bahri, S., Chen, Y., Cheung, C.L., Duan, L., Fan, X.H., Guan, Y., Guo, C.T., Hassan, S.S., Huang, K., Leung, Y.H.C., Li, K.S., Lu, H.R., Naipospos, T.S.P., Nguyen, T.D., Peiris, J.S.M., Qin, K., Rayner, J.M., Smith, G.J.D., Vijaykrishna, D., Wang, J., Webster, R.G., Wu, W.L., Xia, N.S., Xu, K.M., Yuen, K.Y., Zhang, J.X., Zhang, L.J. 2006. Establishment of multiple sublineages of H5N1 influenza virus in Asia: Implications for pandemic control. PNAS. 103, 2845-2850.

De, B. K., Brownlee, G. G., Kendal, A. P., Shaw, M. W. 1988. Complete sequence of cDNA clone of the hemagglutinin gene of influenza A/Chicken/Scotland/59 (H5N1)

virus: comparison with contemporary North American and European strains. Nucleic Acids Res. 16, 4181-4182.

Ducatez, M.F., Ammerlaan, W., de Landtsheer, S., Fouchier, R.A.M., Muller, C.P., Niesters, H.G.M., Olinger, C.M., Osterhaus, A.D.M.E., Owoade, A.A. 2006. Avian Flu: Multiple introductions of H5N1 in Nigeria. Nature 442, 37.

Ducatez, M. F., Ammerlaan, W., De Landtsheer, S., Fouchier, R. A. M., Muller, C. P., Olinger, C. M., Osterhaus, A. D. M. E., Ouedraogo, J. B., Owoade, A. A., Sow, A., Tahita, M. C., Tarnagda, Z. 2007. Molecular and antigenic evolution and geographical spread of H5N1 highly pathogenic avian influenza viruses in western Africa. J. Gen. Virol. 88, 2297-2306.

Farris, J.S. Methods for computing Wagner trees. 1970. Syst. Zool. 19, 83-92.

Fitch, W.M. 1971. Toward defining the course of evolution: minimum change for a specific tree topology. Syst. Zool. 20, 406-416.

Goloboff, P. A., Farris, J.S., Nixon, K.C. 2008. TNT, a free program for phylogenetic analysis. Cladistics 24, 774-786

GISAID, 2009. Sample provided by Aly, M. M., Arafa, A. S. National laboratory for quality control on poultry production, Dokki (Cairo), Egypt. Sequence submitted by

1
2 Monne, I., Fusaro, A., Capua, I., Cattoli, G. OIE/FAO Reference Laboratory for avian
3 influenza and Newcastle disease, Istituto Zooprofilattico Sperimentale delle Venezie, Legnaro
4 (PD), Italy.
5
6
7
8
9

10 Guan, Y., Dyrting, K.C., Ellis, T. M., Kong, K. F., Peiris, Y., Shortridge, K. F., Sit, T.,
11 Zhang, L. J. 2002. H5N1 influenza viruses isolated from geese in southeastern China:
12 evidence for genetic reassortment and interspecies transmission to ducks. *Virology* 292,
13 16-23.
14
15
16
17

18
19 Guan, Y., Auewarakul, P., Buranathai, C., Chaiseng, A., Chen, H., Duan, L.,
20 Estoepangestie, A.T.S., Hanh, N.T.H., Li, K.S., Long, H.T., Nguyen, T.D., Peiris, J.S.M.,
21 Poon, L.L.M., Puthavathana, P., Rahardjo, A.P., Shortridge, K.F., Smith, G.J.D., Wang,
22 J., Webby, R.J., Webster, R.G., Xu, K.M., Yuen, K.Y. 2004. Genesis of a highly
23 pathogenic and potentially pandemic H5N1 influenza virus in eastern Asia. *Nature* 430,
24 209-213.
25
26
27
28
29
30
31

32
33 Hatchette, T.F., Baker, A., Johnson, C., Pryor, S. P., Walker, D., Webster, R.G. 2004.
34 Influenza A viruses in feral Canadian ducks: extensive reassortment in nature. *J. Gen.*
35
36 Virol.

37 85, 2327-2337.

38
39 Hill, A., Guralnick, R. Wilson, M., Habib, F. Janies, D. 2009. Evolution of drug
40 resistance in multiple distinct lineages of H5N1 avian influenza. *Infection, Genetics, and*
41
42 *Evolution*. 9, 169-178.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 HoDac, H., Fitch, W.M., Lathrop, R.H., Wallace, R.G. 2007. MigratPhyla: Statistical
5 analysis of migration events through a phylogeny. Version 1.0b manual.
6
7 <http://pd.bio.uci.edu/ee/WallaceR/MigratPhyla.html>
8
9

10
11 Holland, J., Grabau, E., Horodyski, F., Nichol, S., Spindler, K., VandePol, S. 1982. Rapid
12 evolution of RNA genomes. *Science* 215, 1577-1585.
13
14

15 Huelsenbeck, J.P and Ronquist, F. 2001. MRBAYES: Bayesian inference of phylogenetic
16 trees. *Bioinformatics* 17, 754-755.
17
18

19 Janies, D., Hill, A.W., Habib, F., Guralnick, R., Waltari, E., Wheeler, W.C. 2007.
20 Genomic analysis and geographic visualization of the spread of avian influenza (H5N1).
21
22 *Syst. Biol.* 56, 321-329.
23
24

25 Kawaoka Y., Alexander D.J., Nestorowicz A., Webster R.G. 1987. Molecular analyses of
26 the hemagglutinin genes of H5 influenza viruses: origin of a virulent turkey strain.
27
28 *Virology* 158, 218-227.
29
30

31 Kilpatrick, A. M., Chmura, A.A., Gibbons, D.W., Fleischer, R.C., Marra, P.P., Daszak, P.
32
33 2006. Predicting the global spread of H5N1 avian influenza. *PNAS*. 103, 19368-19373
34
35

36 Kluge, A.G., 1989. A concern for evidence and a phylogenetic hypothesis of
37 relationships among Epicrates (Boidae, Serpentes). *Syst. Zool* 38, 7-25.
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Larkin M.A., Blackshields G., Brown N.P., Chenna R., McGettigan P.A., McWilliam H., Valentin F., Wallace I.M., Wilm A., Lopez R., Thompson J.D., Gibson T.J., Higgins D.G. 2007. Clustal W and Clustal X version 2.0. Bioinformatics 23, 2947-2948.

Formatted: Don't adjust space between Latin and Asian text, Don't adjust space between Asian text and numbers

Lewis, P.O. 2001. A likelihood approach to inferring phylogeny from discrete morphological characters. Syst. Biol. 50:913-925.

Deleted: Lindley, D. 2009. Calculating the Future. Communications of the ACM 52, 9-11. ¶ ¶

Lindstrom, S.E., Hiromoto, Y., Nerome, R., Omoe, K., Sugita, S., Yamazaki, Y., Takahashi, T., Nerome, K. 1998. Phylogenetic analysis of the entire genome of influenza A (H3N2) viruses from Japan: evidence for genetic reassortment of the six internal genes. J. Virol. 72, 8021-8031.

Lindstrom, S.E., Cox, N.J., Klimov, A. 2004. Genetic analysis of human H2N2 and early H3N2 influenza viruses, 1957-1972: evidence for genetic divergence and multiple reassortment events. Virology 328, 101-119.

Maddison, W. P. and D.R. Maddison. 2008. Mesquite: a modular system for evolutionary analysis. <http://mesquiteproject.org>

Deleted: <http://mesquiteproject.org>
Formatted: Font: Courier New, 10 pt

Matrosovich, M., Kawaoka, Y., Webster, R., Zhou, N. 1999. The surface glycoproteins of H5 influenza viruses isolated from humans, chickens, and wild aquatic birds have distinguishable properties. J. Virol. 73, 1146-1155.

1
2
3
4 Obenauer, J.C., Denson, J., Fan, Y., Finkelstein, D.B., Hoffmann, E., Krauss, S., Ma, J.,
5
6 Mehta, P.K., Mukatira, S., Naeve, C.W., Rakestraw, K.M., Su, X., Wang, J., Webster,
7
8 R.G., Xu, X., Zhang, Z., Zheng, J. 2006. Large-scale sequence analysis of avian influenza
9 isolates. *Science* 311, 1576-1580.
10
11
12
13

14 Salzberg, S.L., Ali, A.S.M., Aly, M.M., Brown, I.H., Capua, I., Cattoli, G., Couacy-
15 Hymann, E., De Mia, G.M., Dung, D.H., Ghedin, E., Guercio, A., Janies, D.A., Joannis,
16 T., Kingsford, C., Osmani, A., Padalino, I., Saad, M.D., Saviæ, V., Sengamalay, N.A.,
17 Spiro, D.J., Yingst, S., Zaborsky, J., Zorman-Rojs, O. 2007. Genome analysis linking
18 recent European and African influenza (H5N1) viruses. *Emerg. Infect. Dis.* 13, 713-718.
19
20
21
22
23
24
25
26

27 Sankoff, D. and Rousseau, P. 1975. Locating the vertices of a Steiner tree in arbitrary
28 space. *Mathematical Programming*, p. 240246.
29
30
31
32

33 Sawabe, K., Hayashi, T., Hoshino, K., Hotta, A., Isawa, H., Kobayashi, M., Kurahashi,
34 H., Saito, T., Sasaki, T., Tanabayashi, K., Tsuda, Y., Yamada, A. 2006. Detection and
35 isolation of highly pathogenic H5N1 avian influenza A viruses from blow flies collected
36 in the vicinity of an infected poultry farm in Kyoto, Japan, 2004. *Am. J. Trop. Med. Hyg.*
37
38
39 75, 327-332.
40
41
42
43
44

45 Slatkin, M. and Maddison, W. P. 1989. A cladistic measure of gene flow inferred from
46 the phylogeny of alleles. *Genetics* 123, 603-613.
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 Smith, D., Lapedes, A., de Jong, J., Bestebroer, T., Rimmelzwaan, G., Osterhaus, A.,
5 Fouchier, R. 2004. Mapping the Antigenic and Genetic Evolution of Influenza Virus.
6 Science 305, 371-376.
7
8
9

10
11
12 Smith, G.J.D., Bui, N.A., Chen, H., Cheung, C.L., Dao, T.V., de Jong, M.D., Farrar, J.,
13 Guan, Y., Hassan, S.S., Hien, T.T., Leung, Y.H.C., Li, K.S., Naipospos, T.S.P., Nguyen,
14 T.D., Nguyen, T.V., Nguyen, V.C., Peiris, J.S.M., Poon, L.L.M., Rayner, J.M., Usman,
15 T.B., Vijaykrishna, D., Webster, R.G., Zhang, J.X., Zhang, L.J. 2006. Evolution and
16 adaptation of H5N1 influenza virus in avian and human hosts in Indonesia and Vietnam.
17
18 Virology 350, 258-268.
19
20
21
22
23

24
25 Subbarao, K., Klimov, A., Katz, J., Regnery, H., Lim, W., Hall, H., Perdue, M., Swayne, D.,
26 Bender, C., Huang, J., Hemphill, M., Rowe, T., Shaw, M., Xu, X., Fukuda, K., Cox, N. 1998.
27 Characterization of an avian influenza a (H5N1) virus isolated from a child with a fatal
28 respiratory illness. Science 279, 393-396.
29
30
31
32
33

34
35 Swofford, D. L., 1990. PAUP: Phylogenetic Analysis Using Parsimony, Version 3.0.
36
37 Manual. Illinois Natural History Survey, Champaign, Ill.
38
39
40
41
42

43 Swofford, D.L. 2008. PAUP* Phylogenetic Analysis Using Parsimony (*and Other
44 Methods). Version 4. Sinauer Associates, Sunderland, Massachusetts.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2 Swofford, D.L. and Maddison, W.P. 1987. Reconstructing ancestral character states
3
4 under Wagner parsimony. *Math. Biosci.* 87, 199-229.
5
6
7
8
9

10 USDA, 2007. Release No. 0296.06. Avian influenza. Low pathogenic H5N1 vs. highly
11 pathogenic H5N1. Latest update, 23 July 2007. <http://www.usda.gov/2006/08/0296.xml>
12
13
14
15

16 Van Borm, S., Boschmans, M., Decaestecker, M., Dupont, G., Hanquet, G., Lambrecht,
17 B., Snacken, R., Thomas, I., van den Berg, T. 2005. Highly pathogenic H5N1 influenza
18 virus in smuggled Thai eagles, Belgium. *Emerg. Infect. Dis.* 11, 702-705.
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

30 Wallace, R.G., HoDac, H., Lathrop, R.H., Fitch, W.M. 2007. A statistical
31 phylogeography of influenza A H5N1. *PNAS* 104, 4473-4478.
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

30 Wallace, R.G. and Fitch, W.M. 2008. Influenza A H5N1 immigration is filtered out at
31 some international borders. *PLoS ONE* 3, e1697. doi:10.1371/journal.pone.0001697.
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

30 Wang, J., Bahl, J., Chen, H., Smith, G.J.D., Duan, L., Guan, Y., Peiris, J. S. M.,
31 Vijaykrishna, D., Webster, R. G., Zhang, J. X. 2008. Identification of the progenitors of
32 Indonesian and Vietnamese avian influenza A (H5N1) viruses from southern China. *J.*
33 *Virol.* 82, 3405-3414.
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

30 Webster, R.G., Bean, W.J., Chambers, T.M., Gorman, O.T., Kawaoka, Y. 1992.
31
32 Evolution and ecology of influenza A viruses. *Microbiol. Rev.* 56, 152-179.
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 Wheeler, W. C., Delaet, J., Gladstein, D., and Varón, A., 2005. POY (version 3.012).
5
6
7
8
9

10
11 Phylogeny reconstruction via optimization of DNA and other data.
12
13
14 <http://research.amnh.org/scicomp/projects/poy.php>
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

WHO, 2004. Avian influenza A(H5N1)- update 31: Situation (poultry) in Asia: need for a long-term response, comparison with previous outbreaks. Reported to WHO. 2 March 2004. http://www.who.int/csr/don/2004_03_02/en/

WHO, 2009. Cumulative Number of Confirmed Human Cases of Avian Influenza

A/(H5N1) Reported to WHO. 11 August 2009.

http://www.who.int/csr/disease/avian_influenza/country/cases_table_2009_08_11/en/index.html

Varón, A., Vinh, L. S., Bomash, I., Wheeler, W. C. 2009. POY 4.1.1. American Museum of Natural History. <http://research.amnh.org/scicomp/projects/poy.php>

Xu X., Cox, N.J., Guo, Y., Subbarao K. 1999. Genetic characterization of the pathogenic influenza A/Goose/Guangdong/1/96 (H5N1) virus: similarity of its hemagglutinin gene to those of H5N1 viruses from the 1997 outbreaks in Hong Kong. *Virology* 261, 15-19.

Zamarin, D., Ortigoza, M.B., Palese, P. 2006. Influenza A Virus PB1-F2 Protein Contributes to Viral Pathogenesis in Mice. *Journal of Virology* 80: 7976-7983.

Deleted: 14 January

Deleted: 01_14

Deleted: ¶
World Health Organization/World Organisation for Animal Health/Food and Agriculture Organization H5N1 Evolution Working Group. Toward a unified nomenclature system for highly pathogenic avian influenza virus (H5N1). 2008. *Emerg. Infect. Dis.* Available from <http://www.cdc.gov/EID/content/14/7/e1.htm> doi: 10.3201/eid1407.071681¶

Deleted: ¶

Deleted:
Figure captions: ¶
¶ **Figure 1:** ¶

Formatted: Swedish (Sweden)

1
2 Zhou, N. N., Shortridge, K.F., Claas, E.C.J., Krauss, S.L., Webster, R.G. 1999. Rapid
3
4 evolution of H5N1 influenza viruses in chickens in Hong Kong. J. Virol. 73, 3366-3374.

5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Figure captions

Figure 1.

The same tree is shown with two different optimizations of a geographic character. The tree on the left shows one transmission from location A to location B and one from location B to location A. The tree on the right has two transmission events from location A to location B, and zero transmission events from location B to location A. The global maximum number of transmissions are thus: location A to location B: 2; location B to location A: 1.

Deleted: first

Deleted: second

Deleted: transmissions

Deleted: none

Deleted: hence

Deleted: ,

Deleted: :

Figure 2.

A comparison of transmission routes found by Wallace and Fitch (2008) and our reanalysis of the WF482HA dataset. Each cell in the grid represents a transmission route between two locations in the 28 state geographical character set used in the original publication. Cells colored black indicate transmission routes that were not found in any analysis. Cells colored yellow indicate that our reanalysis agrees with Wallace and Fitch (2008). Cells colored cyan indicates that the transmission route was not found in our reanalysis. Cells colored red indicate transmission routes found in the reanalysis, but not in Wallace and Fitch (2008).

Deleted: &

Deleted: Black indicates

Deleted: A

Deleted: cell indicated

Deleted: the original,

Deleted: and

Deleted: that it is

Deleted: the original study.

The rows represent outgoing locations and the columns incoming. For example, the cell in the Vietnam column and the Guangdong row represents a transmission from Guangdong to Vietnam, while the cell in the Guangdong column and the Vietnam row represents a transmission from Vietnam to Guangdong

Figure 3.

Migration routes indicated by the HA1646 dataset.

A. Transmissions into Nigeria, using the Africa character set. The orange lines represent transmission routes are bidirectional or ambiguous. In H5N1 in Africa, this result is due to very low genetic variation between the H5N1 in Nigeria and neighboring regions.

B. Transmissions into Indonesia, using the "W&F" character set. Analyses indicate only one source of H5N1 in Indonesia, Hunan province in China.

C. Transmission into Japan, using the "World" character set. Both China and South Korea are sources for the H5N1 in Japan. The direction of the transmission route to Vietnam is either bidirectional or ambiguous.

Deleted::

Deleted: 'Africa'

Deleted: represents

Deleted: that

Deleted: either

Deleted:,

Deleted: where the direction is indeterminable.

Deleted: case, this

Deleted: an artefact of

Deleted: Neighboring

Deleted:

Deleted: 'Wallace & Fitch'

Deleted: the H5N1in

Deleted: 'World'

Deleted:

Deleted: not possible to determine.

Figure 4:

1
2 The consensus tree can imply transmission routes never found among the most
3 parsimonious trees.

4
5
6
7
8 When an unordered multistate character with states 0, A, B and C is optimized on trees a,
9 and b., no optimal optimization path has the directional state change C to A. However,
10 considering optimal optimization paths in the strict consensus tree, one path has the node
11 above the outgroup assigned to state B and thus two changes from B to C, as well as two
12 changes from C to A are required.

13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review