1    **Appendix S1.**  More details of the genome sequencing procedure and the indices used to

2    measure gene flow.

3

4    **Viral Genome Sequencing.**  All samples were first sequenced using the traditional

5    PCR/Sanger high-throughput sequencing pipeline at the JCVI (Dugan *et al.* 2008).  After

6    sequencing, the sequences were trimmed to remove primer sequence as well as low quality

7    residues, and segments were assembled individually using the small genome assembler Elvira

8    (http://elvira.sourceforge.net/).  Sequencing and assembly difficulties associated with the high

9    level of variability in the sample set (e.g. due to the mixed infection of multiple influenza

10   subtypes) led to the processing of all samples using the next generation sequencing pipeline at

11   JCVI that includes the 454/Roche GS-FLX and the Illumina Genome Analyzer II.  Viral RNA was

12   first reverse transcribed and amplified by multi-segment RT-PCR (M-RTPCR) (Zhou *et al.* 2009).

13   The cDNA for each sample was primed with barcoded random hexamers using the SISPA

14   protocol (Djikeng *et al.* 2008).  One library was prepared for sequencing on the 454/Roche GS-

15   FLX platform using Titanium chemistry while the other was made into a library for sequencing

16   on the Illumina Genome Analyzer II.

17        The sequence reads from the GS-FLX data were sorted by barcode, trimmed, and

18   searched by TBLASTX against custom nucleotide databases of full-length influenza A segments

19   downloaded from GenBank to filter out both chimeric influenza sequences and non-influenza

20   sequences amplified during the random hexamer-primed amplification.  The filtered GS-FLX

21   reads were then binned by segment, and *de novo* assembled using CLC Bio's

22   clc_novo_assemble program.  Because of the short read length of the sequences obtained from

23   the Illumina Genome Analyzer II, these were not subjected to the TBLASTX filtering step.  Both

24   GS-FLX and Illumina reads were then mapped to the selected reference influenza A virus

25   segments using the clc_ref_assemble_long program.  At loci where both GS-FLX and Illumina

26   sequence data agreed on a variation compared to the reference sequence, the latter was

1

1    updated to reflect the difference.  A final mapping of all next generation sequences to the

2    updated reference sequences was then performed.  All sequences generated here have been

3    submitted to GenBank and assigned accession numbers (Table S1).

4

5    **Phylogenetic Analysis**

6    Panoramic phylogenies were estimated using a rapid hill-climbing search method and the

7    GTRGAMMA nucleotide substitution model implemented in RAxML v7.04 (Stamatakis 2006).

8    The tree search was initiated with a random maximum parsimony tree.  200 independent tree

9    searches with different random maximum parsimony starting trees were performed to obtain the

10   phylogeny with the highest likelihood score.  Major lineages of North American wild bird avian

11   influenza viruses were identified and extracted for further analyses.

12          For more accurate phylogenetic analysis of NA-WB-AIV lineages, we used the heuristic

13   search method implemented in PhyML v2.4.5 (Guindon *et al*. 2009).  This allowed us to

14   determine the maximum likelihood phylogenetic tree by optimizing the tree topology and branch

15   lengths on the sequence data.  The tree search was initiated using a BIONJ tree and then

16   employed nearest-neighbor interchange branch-swapping.  All phylogenetic analyses utilized

17   the General Time Reversible (GTR) model of nucleotide substitution which allows variable rates

18   of substitution between the A, T, C and G nucleotides, a proportion of invariable sites (I), and

19   four classes of rate variation among nucleotide sites ($+\Gamma_4$).  The RAxML and PhyML trees were

20   extremely similar in topology.

21          Pseudo-replicates for each data set were generated by resampling the columns of

22   sequence alignment with replacement, and then used in bootstrap analyses to determine

23   phylogenetic robustness.  Because of the very large size of the 'panoramic' data set (>5,000

24   sequences), the neighbor-joining clustering method (as implemented in the PAUP* package –

25   Swofford *et al*. 2003) was used to infer bootstrap trees in this case.  In the case of the North

26   American wild bird AIV lineage, where more phylogenetic accuracy is required, we used the

1    maximum likelihood bootstrap method implemented in PhyML.  More background information

2    about the types of phylogenetic analyses performed here can be found in Lam *et al.* (2010).

3        All AIV gene segments were subject to the independent phylogenetic and

4    phylogeographic analyses, with the exception that the only HA and NA gene segments analyzed

5    were from subtypes H3, H4, N6, and N8.  The other subtypes had HA and NA data sets of

6    insufficient size (<200 sequences or 8 sampling localities) for meaningful phylogeographic

7    analysis.

8

9    **Measuring Gene Flow Between Populations**

10   Three summary statistics – (i) the Fixation index ($F_{ST}$), (ii) modified Slatkin-Maddison's *s* ($\sigma$),

11   and (iii) the rate of state transition (*q*) – were used to determine the level of gene flow between

12   two localities from either the sample of virus sequences in the two localities or from the inferred

13   phylogenetic trees.

14       **(i) $F_{ST}$.**  $F_{ST}$ is derived from the *F*-statistics developed by Sewall Wright to study

15   population structure (Wright 1942).  It was later modified and generalized for genetic sequence

16   data sampled in different populations (Nei 1982; Lynch & Crease 1990).  Our study employed a

17   $F_{ST}$ similar to that generalized by Hudson et al. (1992) using the GTR substitution model to

18   estimate pairwise genetic distances ($\delta_{xy}$) between two nucleotide sequences (*x* and *y*).  $F_{ST}$

19   between the populations sampled at geographical states *i* and *j* is defined as,

20   $$F_{ST} = 1 - \frac{H_w}{H_b}$$

21   $H_w$ is the average within-population variation in states *i* and *j*, which is defined as,

22   $$H_w = \frac{\left(\frac{2}{n_i(n_i - 1)}\sum_{ix<iy} \delta_{ix,iy}\right) + \left(\frac{2}{n_j(n_j - 1)}\sum_{jx<jy} \delta_{jx,jy}\right)}{2}$$

23   where $n_i$ and $n_j$ are the total numbers of sequences sampled from the populations at

24   geographical states *i* and *j*, respectively.  *ix* and *iy* are the $x^{th}$ and $y^{th}$ sequences in the

1    population sampled in geographical state $i$, while $jx$ and $jy$ are the $x^{th}$ and $y^{th}$ sequences in the

2    population sampled in geographical state $j$.

3    $H_b$ is the between-population variation of state $i$ and $j$, which is defined as,

4    $$H_b = \frac{\sum_{ix,jy} \delta_{ix,jy}}{n_i n_j}$$

5    This is the average pairwise comparison of sequences sampled in the populations of

6    geographical states $i$ and $j$. Extreme $F_{ST}$ estimates of 0 and 1 indicate no and complete

7    population subdivision, respectively. $F_{ST}$ was calculated for each pair of geographical states in

8    the NA-WB-AIV data sets, and was plotted against the spatial distance (km) between each pair

9    of geographical states (Fig. S2). Geographical states with too few sampled sequences ($n \leq 5$ or

10   less than 1% of all sequences) were excluded. The PAUP* program (Swofford 2003) was used

11   to estimate the GTR+I+$\Gamma_4$ genetic distance between sequences. A PERL script implementing

12   the calculation of $F_{ST}$ from the genetic distances is available upon request.

13   **(ii) Modified Slatkin-Maddison's s**. Slatkin and Maddison proposed and demonstrated

14   that the minimum number of geographical state change ($s$) (between two geographical states)

15   observed in the phylogeny provides an estimate of the level of gene flow between these

16   populations (Slatkin & Maddison 1989; Hudson *et al.* 1992),

17   $$s \approx Nm$$

18   where $N$ is the population size and $m$ is the migration rate between the populations.

19   In situations where the sampling frequencies are highly variable among localities, the

20   number of sequences sampled in two geographical locations is the limiting factor for the

21   minimum number of the geographical state changes ($s$) that can be observed in the phylogeny.

22   To overcome this problem, we used a modified Slatkin-Maddison's $s$, denoted 'σ', to infer the

23   level gene flow between two geographical states (e.g. $A$ and $B$). This is defined as,

24   $$\sigma = \frac{s}{s_p}$$

1 where $s$ is number of changes between state $A$ and state $B$ in the observed phylogeny, and $s_P$ is

2 the number of changes between state $A$ and state $B$ occurring in the phylogeny simulated

3 assuming panmixis. In other words, $s_P$ determines the number of state changes that could be

4 observed with the current phylogenetic structure and number of samples available in different

5 geographical states. This scaling turned Slatkin-Maddison's $s$ into a measure of the relative

6 level of gene flow compared to the scenario of the completely unrestricted gene flow (i.e.

7 panmixis). Extreme σ estimates of 0 and 1 (they may possibly but rarely exceed 1) indicate no

8 gene flow and panmixis, respectively.

9      The values of $s$ between different geographical states were estimated from each

10 maximum likelihood gene phylogeny using the Fitch parsimony method (Fitch 1971)

11 (implemented in the JAVA BEAST library). Polytomies were resolved randomly 100 times and

12 average $s$ values were obtained. To obtain $s_P$, each pair of geographical states was randomly

13 shuffled at the tree tips for 1,000 iterations, and ancestral geographical states were again

14 reconstructed using parsimony. This generated a distribution of $s_P$ under panmixis, and the

15 average $s$ is divided by the mean $s_P$ to give σ. To account for topological uncertainty, 200

16 bootstrap maximum likelihood trees were analyzed in the manner described above, and the

17 upper and lower 95% estimates were taken as the resultant uncertainty of σ. A JAVA computer

18 program implementing this procedure is available upon request.

19      To study rates of gene flow within and between migratory flyways, transitions in

20 geographical state, e.g. Minnesota ↔ Alaska, (estimated from the phylogeny using the

21 parsimony method described above) were categorized as transitions in flyways, e.g. MF ↔ PF,

22 depending on which flyway the geographical state belongs to (e.g. Minnesota belongs to

23 Mississippi Flyway (MF), Alaska belongs to Pacific Flyway (PF), see Table S2). Hence, $s$ and

24 $s_P$ were counted for each type of transition among four flyways (i.e. MF↔PF, MF↔MF, MF↔CF,

25 MF↔AF, PF↔PF, PF↔CF, PF↔AF, CF↔CF, CF↔AF and AF↔AF).

1     **(iii) Rate of state transition ($q$) by Pagel's maximum likelihood method.**  Pagel

2 (1994) developed a maximum likelihood method to study trait evolution (as discrete characters)

3 along the phylogenetic tree (Pagel 1994).  We employed this method and treated the

4 geographical states as the trait at every tip.  Formally, this method maximizes

5 $L = f(D|T, M)$

6 where $L$ is the likelihood of the data set $D$ (i.e. the trait states of the taxa) given the phylogeny $T$,

7 and the model of evolution, $M$, of the trait states.  In the analysis of gene flow between each pair

8 of different geographical states, we estimated the parameters for the state transition rates ($q$) in

9 the model.  In total, there were 120 reversible transition rate parameters (e.g. $q_{MN \leftrightarrow TX}$, $q_{ND \leftrightarrow OR}$,

10 $q_{SD \leftrightarrow WA}$) of interest given the 16 geographical states in the study.  The huge number of

11 parameters made the model statistically untraceable.  To reduce the large number of

12 parameters, we simplified the state transition model by reducing the rate parameters for two

13 states to rate parameters for two flyways.  This constituted the flyway-specific rate model (FRM)

14 of AIV gene flow.  In the FRM there were 10 reversible transition rate parameters for four

15 flyways (and 6 parameters for the three flyway model (3-FRM), which combines MF and CF):

16 $q_{CF \leftrightarrow AF}$ , $q_{CF \leftrightarrow PF}$ , $q_{CF \leftrightarrow MF}$ , $q_{CF \leftrightarrow CF}$ , $q_{AF \leftrightarrow MF}$ , $q_{AF \leftrightarrow PF}$ , $q_{AF \leftrightarrow AF}$ , $q_{PF \leftrightarrow PF}$ , $q_{PF \leftrightarrow MF}$ , $q_{MF \leftrightarrow MF}$.  These parameters

17 were estimated from the phylogeny with the highest likelihood score, using the discrete

18 maximum likelihood method described by Pagel (1994) as implemented in the APE library

19 (Paradis *et al.* 2004) running in the R package version 2.11.1.  Various starting values of

20 transition rates were attempted (20, 10, 5, 1, 0.1 and 0.01), and the estimates with the best

21 likelihoods were kept.  The flyway-specific rate estimates are shown in Table S3.

22

23 **References**

24 Djikeng A., Halpin R., Kuzmickas R., Depasse J., Feldblyum J., Sengamalay N., Afonso C.,

25     Zhang X., Anderson N.G., Ghedin E. & Spiro D.J. (2008) Viral genome sequencing by

26     random priming methods. *BMC Genomics*, 9, 5

Dugan V.G., Chen R., Spiro D.J., Sengamalay N., Zaborsky J., Ghedin E., Nolting J., Swayne D.E., Runstadler J.A., Happ G.M., Senne D.A., Wang R., Slemons R.D., Holmes E.C. & Taubenberger J.K. (2008) The evolutionary genetics and emergence of avian influenza viruses in wild birds. *PLoS Pathog*, 4, e1000076

Fitch W.M. (1971) Toward defining the course of evolution: minimum change for a specified tree topology. *Syst Zool*, 20, 406-16

Hudson R.R., Slatkin M. & Maddison W.P. (1992) Estimation of levels of gene flow from DNA sequence data. *Genetics*, 132, 583-9

Lam T.T., Hon C.C. & Tang J.W. (2010) Use of phylogenetics in the molecular epidemiology and evolutionary studies of viral infections. *Critical Rev Clin Lab Sci*, 47, 5-49

Lynch M. & Crease T.J. (1990) The analysis of population survey data on DNA sequence variation. *Mol Biol Evol*, 7, 377-94

Nei M. (1982) Evolution of human races at the gene level. *Prog Clin Biol Res*, 103, 167-81

Pagel M. (1994) Detecting correlated evolution on phylogenies: a general method for the comparative analysis of discrete characters. *Proc R Soc Lond B*, 255, 37-45

Paradis E., Claude J. & Strimmer K. (2004) APE: Analyses of phylogenetics and evolution in R language. *Bioinformatics*, 20, 289-90

Slatkin M. & Maddison W.P. (1989) A cladistic measure of gene flow inferred from the phylogenies of alleles. *Genetics*, 123, 603-13

Swofford D.L. (2003) PAUP*. Phylogenetic Analysis Using Parsimony (*and Other Methods). In. Sinauer Associates, Sunderland, Massachusetts.

Wright S. (1942) Isolation by distance. *Genetics*, 28, 114–38

Zhou B., Donnelly M.E., Scholes D.T., St George K., Hatta M., Kawaoka Y. & Wentworth D.E. (2009) Single-reaction genomic amplification accelerates sequencing and vaccine production for classical and swine origin human influenza a viruses. *J Virol*, 83, 10309-13.

1