# Comparative Sequence Analysis of Circulating and Pandemic Flu Viruses in Swine and Humans

## Alicia Lee

CNIC, Chinese Academy of Sciences
Beijing, China
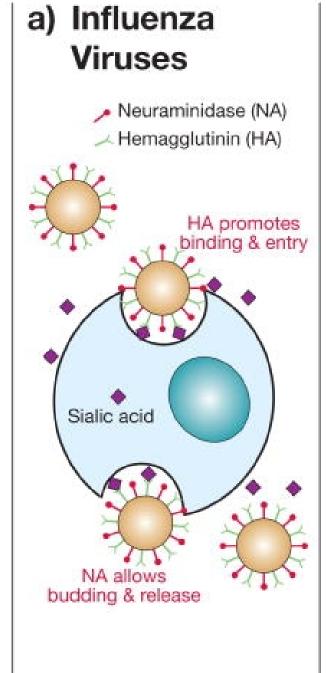
Final Report

August 25, 2010

# Influenza

- Influenza, more commonly known as the flu, is an RNA virus that primarily affects the respiratory system.

- Influenza viruses are named for their surface proteins hemagglutinin (HA) and neuraminidase (NA).

# Proteins

- Hemagglutinin (HA)
- Neuraminidase (NA)
- Receptor Proteins responsible for entry and exit of virus cells
- HA – binding and entry
- NA – budding and release

## a) Influenza Viruses

Neuraminidase (NA)
Hemagglutinin (HA)

HA promotes binding & entry

Sialic acid

NA allows budding & release

# Mutations

- **Antigenic drift** – RNA sequence has mutated the HA or NA protein enough to become immune to previous vaccinations.

- **Antigenic shift** – two or more different influenza viruses affect the same host, leading to a reassortment of genomic segments, creating a new virus against which humans have no pre-existing immunity.
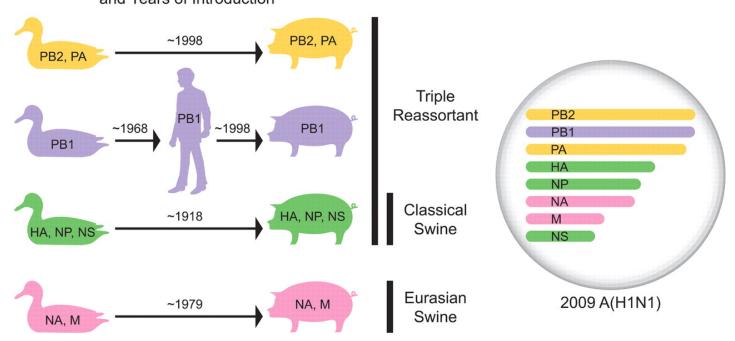
# Pandemic vs Seasonal viruses

- Pandemic strain
  - Infects many people around the world
  - Easily to spread; more contagious
  - 4 pandemics
    - 1918 Spanish Flu (H1N1)
    - 1957 Avian Flu (H2N2)
    - 1968 Hong Kong Flu (H3N2)
    - 2009 Swine Flu (H1N1)

- Seasonal (non-pandemic) strain
  - Localized, more regular occurrences

# 2009 Swine Flu Strain

- The 2009 swine flu strain (SO-IAV) is a 'triple reassortant', with its hemagglutinin gene from the swine lineage.

Gene Segments, Hosts, and Years of Introduction

PB2, PA → ~1998 → PB2, PA

PB1 → ~1968 → PB1 → ~1998 → PB1

Triple Reassortant

HA, NP, NS → ~1918 → HA, NP, NS

Classical Swine

NA, M → ~1979 → NA, M

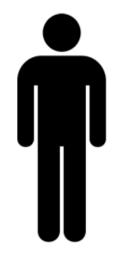Eurasian Swine

PB2
PB1
PA
HA
NP
NA
M
NS

2009 A(H1N1)

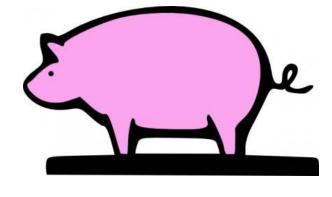- Important to monitor currently circulating swine strains

# Specific Aims

1. Identify signature motifs found in H1N1 pandemic strains affecting humans
   - compare them to sequences found in current swine flu strains circulating among swine using protein sequence analysis

2. Model these mutations using VMD or Chimera
   - visualize protein-ligand interactions and measure binding affinities to human receptor analogs using AutoDock suite tools.

# Modified Aims

1. Identify signature motifs, or highly conserved regions, found in H1N1 pandemic strains affecting humans
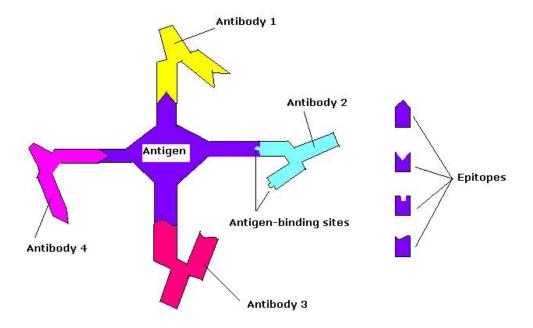   - Compare to currently circulating swine strains and seasonal strains focusing on immune epitopes

**vs**

# Epitopes

- Parts of the virus which are recognized by antibodies or T cells
- Important mutations occur at the epitopes
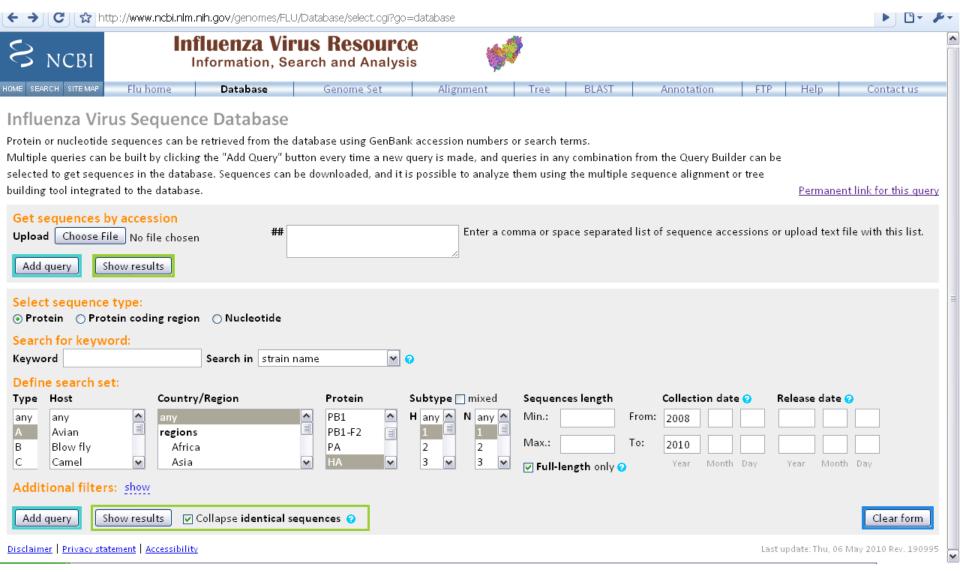
# Aim 1: Find motifs characterized by pandemic strains

- Tools:
  - Influenza Virus Resource (IVR) at NCBI
  - Basic Local Alignment Search Tool (BLAST)
  - MEME suite
    - Multiple Em for Motif Elicitation (MEME)
    - Motif Alignment and Search Tool (MAST)
  - Immune Epitope Database (IEDB)
  - Chimera

# Influenza Virus Resource



Used IVR to find strains to analyze

# Some strains

- Seasonal H1N1 (total of 15)
  - A/Brisbane/59/2007
- Pandemic H1N1 (total of 97)
  - A/California/04/2009 (3LZG)
  - A/New York/1682/2009 2009/04/27 HA
- Previously circulating swine H1N1 strain (1)
  - A/swine/Iowa/15/30(H1N1) (1RUY)
- Currently circulating H1N1 swine strains (26)
  - A/swine/Minnesota/03000/2010(H1N1)
  - A/swine/Beijing/26/2008(H1N1)
  - A/swine/Hong Kong/NS29/2009(H1N1)

# Protein Sequence Alignment

Positions from 121 till 180

```
        Consensus sequence          QLSSVS SFERFEIFPKAS SWPNHDTTRGVTAACPHAGAKSFYRNLI WLVKKGNSYPKLSK
ACQ76318  A/California/04/2009(H1N1)  . . . . . . . . . . . . . . . . . T . . . . . . . S N K . . . . . . . . . . . . . . . . . . . . K . . . . . . . . . .
ACR10223  A/Hamburg/4/2009(H1N1)      . . . . . . . . . . . . . . . . . T . . . . . . . S N K . . . . . . . . . . . . . . . . . . . . K . . . . . . . . . .
AAB52905  A/swine/Iowa/15/1930(H1N1)  . . . . . . . K . . . . . T . . . . . . E . . . . . . . . . . . . . Y . . S . . . . . . L . . . . . E . . . . . . .
ADE28750  A/Brisbane/59/2007(H1N1)    . . . . . . . . . . . . . . . . . E . . . - . V T . . S . . S . . S . N . E S . . . . . . L . . T G . N G L . . N . . .
ADD96946  A/Hamburg/INS92/2009(H1N1)  . . . . . . . . . K . . . . T . . . . . . . . S N K . . . . . . . . . . . . . . . . . . . . K . . . . . . . . . .
ADC32526  A/swine/Argentina/SAGiles-31215/2009(H1N1) . . . . T . . . . . . . S N K . . . . . . . . . . . . . . . . . . . . K . . . . . . . . . .
ACK57777  A/swine/Beijing/26/2008(H1N1)  . . . T . . . . . . . . . . . . . . T . . . . . . . . . T . V . . . S . . S . . V N . . . . . L . . I . . . . . .
ACK57767  A/swine/Fujian/58/2008(H1N1)   . . . T . . . . . . . . . . . . . . T . . . . . . . . . T . V . . . S . . S . . V N . . . . . L . . I . . . . . .
ADG08528  A/swine/Hong Kong/3001/2009(H1N1)  . . . T . . . . . . . . . . . . . . T . . . . . . . . . T . V . . . S . . S . . . N . . . . . L . . I . . . . . .
ADG08198  A/swine/Hong Kong/NS613/2009(H1N1)  . . . T . . . . . . . . . . . . . . T . . K . . . . K . T . V . . . S . S . . . . . . . . . . I . Q . E . . . . N .
BAI83394  A/swine/Ratchaburi/NIAH101942/2008(H1N1)  . . . . . . . . . . . . . . . R E . . . . . . . E . D . . . . . . . . Y . . . N . . . . . . . . . . . . . . . . .
ACK57757  A/swine/Shandong/128/2008(H1N1)  . . . T . . . . . . . . . . . . . . T . . . . . . . . . T . V . . . S . . S . . V N . . . . . L . . I . . . . . .
```

Positions from 121 till 180

```
        Consensus sequence           QLSSVS SFERFEIFPKES SWPNHTVTKGVSASCSHNGKSSFYRNLLWLTGKNGLYPNLSK
ABR15918  A/Auckland/582/2000(H1N1)   . . . . . . . . . . . . . . . . . . . . . . . . . . . . - . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
AAP34323  A/Beijing/262/1995(H1N1)    . . . . . . . . . . . . . . . . . . . . . . K . . . - . . T . . . . . . . . . . . . . . . . . E . . . . . . . . . . N
ABW23335  A/California/09/2006(H1N1)  . . . . . . . . . . . . . . . . . . . . . . . - . . . . . . . . . . . . E . . . . . . . . . . . . . . . . . . . . . X .
ABP49338  A/California/10/1978(H1N1)  . . . . . . . . . . . . . . . . . R . . . K . N . . R . T . . . . K . . . . . . . . . . . . E . . . S . . . . . . .
ACD47246  A/Hawaii/44/2007(H1N1)      . . . . . . . . . . . . . . . . . . . . . . . . - . . . . . . E . . . . K . . . . . . . . . . . . . . . . . . . . .
ABQ44416  A/Memphis/1/1987(H1N1)      . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . T . . . K . . . . . . . . . . E . . . S . . . . . .
ABW23319  A/Minnesota/02/2007(H1N1)   . . . . . . . . . . . . . . . . . . . . . . . . - . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
ADC45736  A/Nagasaki/07N020/2008 (H1N1)  . . . . . . . . . . . . . . . . . . . . . . . . . . . . . E . . . . . . . . . . . . . . . . . . . . . . . . .
AAP34324  A/New Caledonia/20/1999(H1N1)  . . . . . . . . . . . . . . . . . . . . . . . - . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
ACK99034  A/Norway/50/2008(H1N1)      . . . . . . . . . . . . . . . . . . . . . . . . - . . . . . . E . . . . . . . . . . . . . . . . . . . . . . . . .
ACF41834  A/Puerto Rico/8/1934(H1N1)  . . . . . . . . . . . . . . . . . . . . N T N - . T . A . . . E . . . . . . . . . . . . . . . . . E . E . S . . . . K N
ACV49666  A/Siena/4/1987(H1N1)        . . . . . . . . . . . . . . . . . . . . . . K . . . . . . . . T . A . . . K . R . . . . . . . . . . . E . . . . . . .
AAK70450  A/Switzerland/5389/95 (H1N1)  . . . . . . . . . . . . . . . . . . . . . . . . . . T . . . . . . . . . . . K . . . . . . . E . . . . . . . . .
AAP34322  A/Texas/36/1991(H1N1)       . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . T T . . . . . . . . . . . . . . . . K . . . . . . V . .
```

Positions from 121 till 180

```
        Consensus sequence              QLSSVS SFERFEIFPKTS SWPNHDSNKGVTAACPHAGAKSFYKNLIWLVKKGNSYPKLSK
ACQ84467  A/New York/1682/2009(H1N1)      . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
ADH01958  A/Aalborg/INS133/2009(H1N1)     . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
ADA83041  A/Abakan/02/2009(H1N1)          . . . . . . . . . . . . . . . . . . . . . . . . . . . . X . . . . . . . . . . . . . . . . . . . . . . . . . . . .
ADG21140  A/Afghanistan/N09840/2009(H1N1) . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
ACX31945  A/Aichi/198/2009(H1N1)          . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
ACV67245  A/Alabama/02/2009(H1N1)         . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
ADC32423  A/Ancona/97/2009(H1N1)
```

# MEME

# MEME results

**E-value**: the probability of finding an equally well-conserved pattern in random sequences
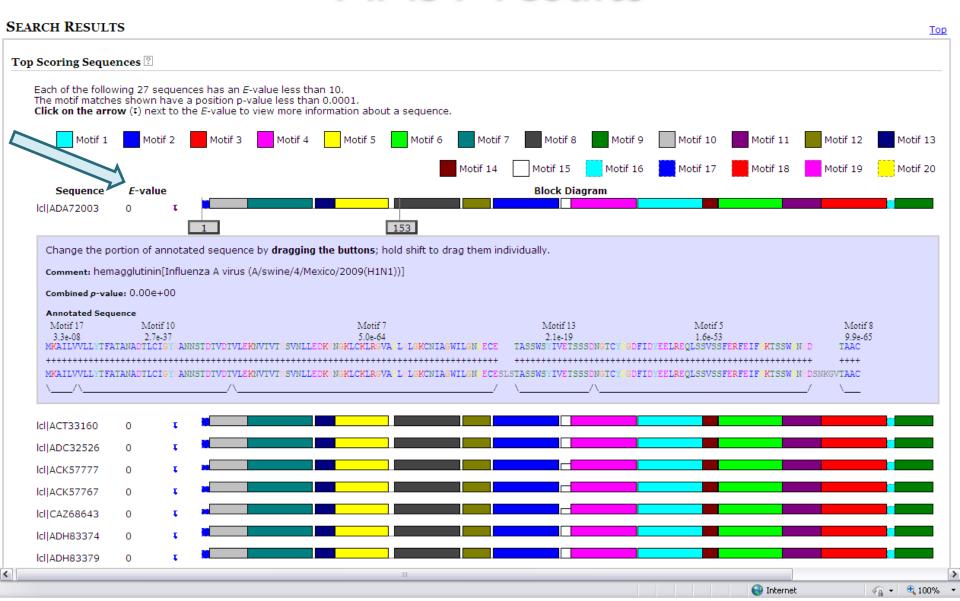


- Input: Pandemic, Swine and Seasonal strains
- 17 motifs found
- Width: 6-50
- Zero or one occurrence of motifs per sequence
- 12 sequences:
  - 4 swine
  - 4 seasonal
  - 4 pandemic
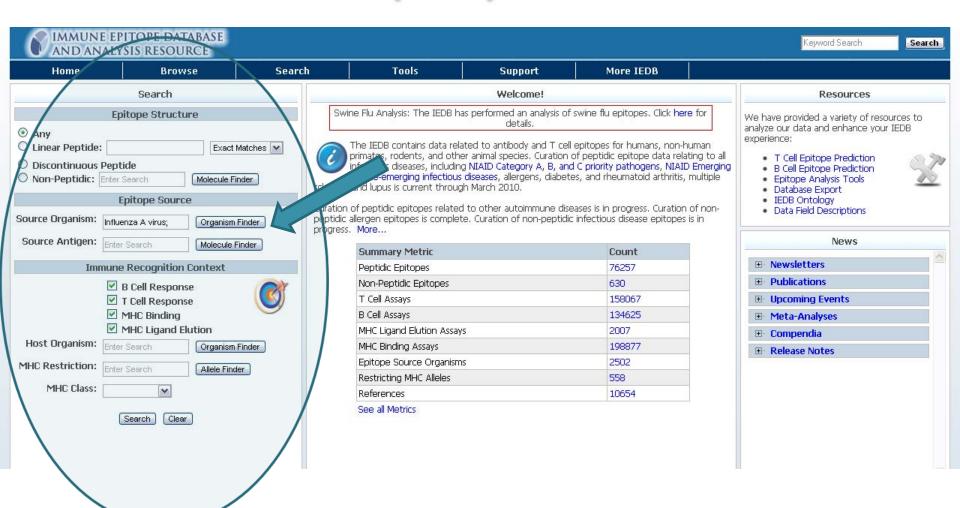
# MEME Swine and Pandemic motifs



Pandemic 97_0-1

Swine 27_0-1

# MAST results



Pandemic 97 0-1 vs Swine 27

# Immune Epitope Database

| Home | Browse | Search | Tools | Support | More IEDB |

Epitope

2391 item(s) found, displaying 1 to 25 (Click the column headers to adjust the sorting)

« previous  1  2  3  4  5  6  7  8  9 ... 95  96  next »    Go To »  1

Export all results: 📗 (full)

| Epitope ID ↑ | Structure | Source Antigen | Source Organism |
|---|---|---|---|
| 133 | AAFEDLRVLSFIRG | nucleoprotein | Influenza A virus |
| 134 | AAFEDLRVLSFIRGTKVSPR | Nucleoprotein | Influenza A virus |
| 360 | AAPIEHIASM | Polymerase acidic protein | Influenza A virus (A/Wilson-Smith/1933(H1N1)) |
| 570 | ACKRGPGSGFFSRLN | Hemagglutinin (1 more) | Influenza A virus (1 more) |
| 695 | ADKRITEMI | Polymerase basic protein 2 | Influenza A virus (A/Ann Arbor/6/1960(H2N2)) |
| 714 | ADLKSTQAAIDQING | Hemagglutinin precursor | Influenza A virus (A/X-31(H3N2)) |
| 729 | ADMSIGVTV | RNA-directed RNA polymerase catalytic subunit | Influenza A virus (A/Ann Arbor/6/1960(H2N2)) |
| 730 | ADMSIGVTVI | RNA-directed RNA polymerase catalytic subunit | Influenza A virus (A/Ann Arbor/6/1960(H2N2)) |
| 754 | ADQKSTQNAI | Hemagglutinin precursor | Influenza A virus (A/Wilson-Smith/1933(H1N1)) |
| 798 | ADYEELREQLSSVSSFERFE | Hemagglutinin precursor | Influenza A virus |
| 825 | AEAIIVAMV | Polymerase basic protein 2 | Influenza A virus (A/Ann Arbor/6/1960(H2N2)) |
| 826 | AEAIIVAMVF | Polymerase basic protein 2 | Influenza A virus (A/Ann Arbor/6/1960(H2N2)) |
| 857 | AEDMGNGCF | Hemagglutinin | Influenza A virus (A/Moscow/343/2003(H3N2)) |
| 984 | AEIEDLIFL | nucleocapsid protein (2 more) | Influenza A virus (A/Viet Nam/1194/2004(H5N1)) (2 more) |
| 985 | AEIEDLIFS | Nucleoprotein | Influenza A virus (A/Bilthoven/4791/81(H3N2)) |
| 1021 | AEKPKFLPDL | Polymerase acidic protein | Influenza A virus (A/Ann Arbor/6/1960(H2N2)) |
| 1055 | AELLVALEN | hemagglutinin | Influenza A virus (A/Memphis/102/1972(H3N2)) |

# BLAST

# Multiple alignment

# Creating FASTA files

- Created FASTA files from BLAST results to get MEME results

# MEME



MEME_ImEplessthan200  (pandemic)

# MEME

# PSPM –Position Specific Scoring Matrix

- ## Protein alphabet:
  - ○ ACDEFGHIKLMNPQRSTVWY

**Data Formats** [?]

View the motif in ⦿ PSPM Format [?]   ○ PSSM Format [?]   ○ BLOCKS Format [?]   ○ FASTA Format [?]   ○ Raw Format [?]   or ○ Hide

```
letter-probability matrix: alength= 20 w= 8 nsites= 160 E= 1.6e-347
0.125000 0.000000 0.000000 0.006250 0.137500 0.000000 0.000000 0.000000 0.150000 0.187500
0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.125000 0.131250 0.000000 0.137500
0.000000 0.000000 0.000000 0.350000 0.000000 0.062500 0.000000 0.000000 0.125000 0.000000
0.000000 0.000000 0.000000 0.000000 0.000000 0.331250 0.131250 0.000000 0.000000 0.000000
0.131250 0.000000 0.000000 0.137500 0.000000 0.006250 0.000000 0.000000 0.187500 0.000000
0.000000 0.062500 0.000000 0.000000 0.137500 0.268750 0.000000 0.068750 0.000000 0.000000
0.131250 0.000000 0.000000 0.000000 0.275000 0.000000 0.000000 0.000000 0.000000 0.206250
0.000000 0.000000 0.000000 0.000000 0.000000 0.256250 0.000000 0.000000 0.131250 0.000000
0.000000 0.131250 0.000000 0.137500 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000
0.006250 0.000000 0.131250 0.000000 0.137500 0.000000 0.000000 0.062500 0.000000 0.393750
0.000000 0.000000 0.000000 0.137500 0.000000 0.000000 0.000000 0.325000 0.125000 0.068750
0.000000 0.131250 0.193750 0.000000 0.012500 0.006250 0.000000 0.000000 0.000000 0.000000
0.000000 0.000000 0.000000 0.000000 0.137500 0.000000 0.262500 0.000000 0.068750 0.000000
0.000000 0.325000 0.000000 0.137500 0.000000 0.000000 0.000000 0.000000 0.068750 0.000000
0.131250 0.000000 0.312500 0.000000 0.000000 0.068750 0.000000 0.000000 0.000000 0.343750
0.000000 0.000000 0.137500 0.000000 0.000000 0.000000 0.006250 0.000000 0.000000 0.000000
```

# MAST



MAST_ImEp200vspandemic97

# Aim 2: Model mutations known to occur in pandemic strains into the circulating swine flu strain

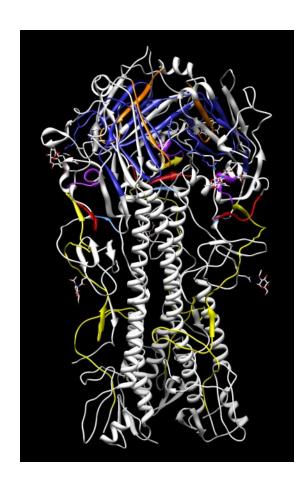- Did some work on Aim 2 but could not carry it further since Aim 1 was not accomplished

# Chimera



1RVX aligned with 3LZG and1RUY



1RVT LSTc



1RUY_mappedMEME swine motifs



3LZG_mapped MEME pandemic94 motifs

# Results

- Found motifs using MEME
  - Changed number of expected repetitions
  - Changed length
  - Categorized by pandemic, seasonal and swine
  - Immune Epitopes
    - Overall MEME E-values showed that most motifs are statistically significant
- Submitted MEME results to MAST
  - Most pandemic motifs were found in swine and seasonal strains as well (and vice versa)
    - E-values were very small numbers signifying that the MEME motifs found in pandemic strains were conserved in swine and seasonal hemagglutinin strains as well

# Results cont…

- Learned: Most hemagglutinin strains are conserved.
  - Which may explain why most papers focus more on specific amino acid changes instead of motifs
  - Difficult to tell whether a motif is significant (may result from functional, structural or evolutionary relationships between sequences)

- Next steps
  - Motifs found by MEME need to be analyzed further
    - background check through literature
    - Use other resources available through the internet to check the significance of the motifs

# Significance

- Knowing signature motifs that correlate to potential virulence for cross-species infection of swine flu viruses into humans would be very useful for **diagnostic assay development**.

- It provides **markers** for monitoring the circulating swine flu viruses and whether the changes occurring would give rise to **potentially pandemic** viruses.

- Further research can also be done on **other proteins** of the virus to see if having the same signature motifs do indeed suggest potential for virulence.

# References

- Altschul, S.F., Gish, W., Miller, W., Myers, E.W. & Lipman, D.J. (1990) "Basic local alignment search tool." J. Mol. Biol. 215:403-410.

- Bailey, T.L. and Elkan, C. (1994) Fitting a mixture model by expectation maximization to discover motifs in biopolymers Proceedings of the Second International Conference on Intelligent Systems for Molecular Biology, August In Altman, R.B., Brutlag, D.L., Karp, P.D., Lathrop, R.H., Searls, D.B. (Eds.). Menlo Park, CA AAAI Press pp. 28–36 .

- Bailey, T.L. and Gribskov, M. (1998) 'Combining evidence using P-values: application to sequence homology searches *Bioinformatics*, **14**, 48–54.

- Bailey TL, Williams N, Misleh C, Li WW (2006) MEME: discovering and analyzing DNA and protein sequence motifs. Nucleic Acids Res 34: W369–373.

- Bao Y., P. Bolotov, D. Dernovoy, B. Kiryutin, L. Zaslavsky, T. Tatusova, J. Ostell, and D. Lipman. The Influenza Virus Resource at the National Center for Biotechnology Information. J. Virol. 2008 Jan;82(2):596-601.

- Influenza A Virus From Wikipedia, the free encyclopedia; http://en.wikipedia.org/wik i/Inf luenza_A_virus

- Jameel, Shahid, The 2009 influenza pandemic. Current Science, Vol. 98, No. 3, 10 Feb 2010. <http://www.ias.ac.in/currsci/10feb2010/306.pdf >

- Vita R, Zarebski L, Greenbaum JA, Emami H, Hoof I, Salimi N, Damle R, Sette A, Peters B. The immune epitope database 2.0. Nucleic Acids Res. 2010 Jan;38(Database issue):D854-62. Epub 2009 Nov 11.

- Garten R. J., Davis C. T., Russell C. A., Shu B., Lindstrom S., Balish A., Sessions W. M., Xu X., Skepner E., Deyde V., Okomo-Adhiambo M., Gubareva L., Barnes J., Smith C. B., Emery S. L., Hillman M. J., Rivailler P., Smagala J., de Graaf M., Burke D. F., Fouchier R. A., Pappas C., Alpuche-Aranda C. M., Lopez-Gatell H., Olivera H., Lopez I., Myers C. A., Faix D., Blair P. J., Yu C., Keene K. M., Dotson P. D., Jr., Boxrud D., Sambol A. R., Abid S. H., George K. St., Bannerman T., Moore A. L., Stringer D. J., Blevins P., Demmler-Harrison G. J., Ginsberg M., Kriner P., Waterman S., Smole S., Guevara H. F., Belongia E. A., Clark P. A., Beatrice S. T., Donis R., Katz J., Finelli L., Bridges C. B., Shaw M., Jernigan D. B., Uyeki T. M., Smith D. J., Klimov A. I., Cox N. J. 2009. Antigenic and genetic characteristics of swine-origin 2009 A(H1N1) influenza viruses circulating in humans. Science 325:197–201

- Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng EC, Ferrin TE. UCSF Chimera--a visualization system for exploratory research and analysis. J Comput Chem. 2004 Oct;25(13):1605-12.

- http://www.cdc.gov/flu/swineflu/key_facts.htm

- Immune Epitope Database and Analysis Resource: www.immuneepitope.org

# Acknowledgements

- UCSD: National Biomedical Computation Resource (NBCR)
  - Dr. Wilfred Li
- CNIC, CAS
  - Dr. Kai Nan
  - Dr. Jianjun Yu
  - Guangyuan Liu
  - Haiyan Xu
  - Wei Chen
- PRIME
  - Dr. Gabriele Wiesenhausen
  - Dr. Peter Arzberger
  - Teri Simas
  - Jim Galvin
  - Tricia Taylor-Oliveira
- National Science Foundation, IOSE-0710726