

Identical Subsequences of Contiguous Amino Acids in Influenza Virus Hemagglutinin and in Human Proteins

Joel K Weltman*

Faculty of Medicine, Alpert Medical School, Brown University, USA

Abstract

Background: Influenza virus is a significant public health problem throughout the world. Increased insight into the basic biology of the virus may enable the development of more effective anti-influenza preventives and therapeutics.

Methodology: The occurrence of specific amino acid subsequences in H1N1 and H3N2 influenza virus hemagglutinins was used for joint detection of those subsequences in human proteins. Only subsequences consisting of at least 5 contiguous amino acids were considered for further study.

Results: Ten H1N1 hemagglutinin amino acid subsequences and nine H3N2 hemagglutinin amino subsequences were identified as also occurring in proteins of human origin. The length of the subsequences selected for further study, ranged from 5 contiguous amino acids to 8 contiguous amino acids.

Conclusion: The joint occurrence of amino acid subsequences in influenza hemagglutinins and in human proteins may help explain the relatively low efficacy of current anti-influenza vaccines. It is proposed that the identification of the joint subsequences may be useful for the improved design of anti-influenza therapeutics and especially anti-influenza vaccines.

Keywords: Influenza virus; H1N1; H3N2 influenza virus; Anti-influenza vaccines

Introduction

Influenza remains a highly significant global public health problem for which vaccination is an extremely important preventive intervention [1,2]. However, currently used anti-influenza vaccines are only 43-62% effective [3]. Therefore, it is extremely urgent and important for us to increase our understanding of influenza virus biology so that more effective vaccines and therapeutic agents can be designed.

This is a report of the occurrence of amino acid subsequences within sequences of hemagglutinin proteins of influenza virus subtypes H1N1, H3N2 and also within human proteins. Understanding the biology of protein subsequences which occur both in the influenza virus and in the human host may facilitate the design of anti-influenza vaccines and drugs.

This research is a continuation of the previous report of the occurrence of identical tetrapeptide subsequences in influenza hemagglutinins and in human proteins [4]. The present report addresses identical subsequences of at least pentapeptide length.

Materials and Methods

Entire sets of full-length (566 amino acids) hemagglutinin (HA) proteins of influenza H1N1 virus (16,679 sequences) and influenza H3N2 virus (19,641 sequences) were downloaded from the NCBI Influenza Virus Database (<https://www.ncbi.nlm.nih.gov/genomes/FLU/Database/nph-select.cgi?go=database>) on 08 Feb 2019. A consensus sequence of each of the two datasets was determined with JalView 2.10.5 [5].

Computing was performed with 64-bit Anaconda3 Python 3.7.1 (<https://www.anaconda.com/distribution/>). Information entropy (H) (in bits) was computed by the method of Shannon [6].

Viral HA protein 3-state secondary structure (h=alpha helix, e=extended strand and c=random coil) was computed on the RaptorX

website [7]. Numpy arrays of the multiplication products of the secondary protein structures were obtained by multiplying each H1N1 HA secondary structure array (length=566 aa positions) by the corresponding H3N2 HA secondary structure array (length=566 aa positions) according to equations 1, 2 and 3, where the symbol \times represents element-wise multiplication:

$$\text{helix product} = h(\text{H1N1}) \times h(\text{H3N2}) \quad (\text{equation 1})$$

$$\text{extended strand product} = e(\text{H1N1}) \times e(\text{H3N2}) \quad (\text{equation 2})$$

$$\text{random coil product} = c(\text{H1N1}) \times c(\text{H3N2}) \quad (\text{equation 3})$$

Because of a sequence length limitation, the sequences from viral HA helical region aa 385-427 were split into two subsequences, 20 and 23 amino acids in length, prior to BLASTP analysis. The consensus amino acids at positions with multiplication product values ≥ 0.75 in the influenza HA product array were selected as probes for matching subsequences in human proteins on the NCBI BLASTP website.

The complete set of non-redundant (nr) sequences of human (*Homo sapiens* (tax id: 9606)) proteins on the National Center for Biotechnology Information database (https://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastp&PAGE_TYPE=BlastSearch&LINK_LOC=blasthome) was screened with BLASTP for occurrence of the subsequences that had been identified within the secondary structures of the H1N1 and H3N2 influenza virus HA proteins [8]. In each case,

*Corresponding author: Joel K Weltman, Clinical Professor Emeritus of Medicine, Alpert Medical School, Brown University, Providence, RI 02912, USA, Tel: 4012457588; E-mail: joel_weltman@brown.edu

Received March 26, 2019; Accepted May 04, 2019; Published May 10, 2019

Citation: Weltman JK (2019) Identical Subsequences of Contiguous Amino Acids in Influenza Virus Hemagglutinin and in Human Proteins. J Med Microb Diagn 8: 301. doi:10.4172/2161-0703.1000301

Copyright: © 2019 Weltman JK. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

the 100 most frequent influenza subsequences within the human protein sequences were downloaded. Each subsequence within the downloaded BLASTP results was considered for further analysis if that subsequence contained at least 5 contiguous amino acids.

Results

The downloaded H1N1 and H3N2 HA sequence sets yielded the following full-length consensus amino acid sequences:

H1N1_HA = M K A I L V V L L Y T F A T A N A
 D T L C I G Y H A N N S T D T V D T V L E K N V T V
 T H S V N L L E D K H N G K L C K L R G V A P L H L G K C N I A G W I L G N
 P E C E S L S T A S S W S Y I V E T S S D N G T C Y P G D F I D Y E E L R E Q L S S
 V S S F E R F E I F P K T S S W P N H D S N K G V T A A C P H A G A K S F Y K N L I
 W L V K K G N S Y P K L S K S Y I N D K G K E V L V L W G I H P S T S A D Q Q S
 L Y Q N A D A Y V F V G T S R Y S K K F K P E I A I R P K V R D Q E G R M N Y Y
 W T L V E P G D K I T F E A T G N L V P R Y A F A M E R N A G S G I I S D T
 P V H D C N T T C Q T P K G A I N T S L P F Q N I H P I T I G K C P K Y V K
 S T K L R L A T G L R N V P S I Q S R G L F G A I A G F I E G G W T G M V D G
 W Y G Y H H Q N E Q G S G Y A A D L K S T Q N A I D K I T N K V N S V I E K
 M N T Q F T A V G K E F N H L E K R I E N L N K K V D D G F L D I W T Y N A E L L V
 L L E N E R T L D Y H D S N V K N L Y E K V R S Q L K N N A K E I G N G C F E F Y H
 K C D N T C M E S V K N G T Y D Y P K Y S E E A K L N R E E I D G V K L E S T R I Y Q
 I L A I Y S T V A S S L V L V S L G A I S F W M C S N G S L Q C R I C I

and

H3N2_HA = M K T I I A L S Y I L C L V F A Q K L P G N D N S T A T
 L C L G H H A V P N G T I V K T I T N D R I E V T N A T E L V Q N S S I G E I C
 D S P H Q I L D G E N C T L I D A L L G D P Q C D G F Q N K K W D L F V E R
 S K A Y S N C Y P Y D V P D Y A S L R S L V A S S G T L E F N N E S F N W T
 G V T Q N G T S S A C I R R S N S S F S R L N W L T H L N Y K Y P A L N V T
 M P N N E Q F D K L Y I W G V H H P G T D K D Q I F L Y A Q S S G R I T V S T
 K R S Q Q A V I P N I G S R P R I R D I P S R I S I Y W T I V K P G D I L L I N S T G N
 L I A P R G Y F K I R S G K S S I M R S D A P I G K C K S E C I T P N G S I P N D K
 P F Q N V N R I T Y G A C P R Y V K Q S T L K L A T G M R N V P E K Q T R G I F
 G A I A G F I E N G W E G M V D G W Y G F R H Q N S E G R Q A A D L K S T Q A
 A I D Q I N G K L N R L I G T N E K F H Q I E K E F S E V E G R I Q D L E K Y V E D T

**K I D L W S Y N A E L L V A L E N Q H T I D L T D S E M N K L F E K T K K Q L R E
 N A E D M G N G C F K I Y H K C D N A C I G S I R N G T Y D H N V Y R D E A L N
 N R F Q I K G V E L K S G Y K D W I L W I S F A I S C F L L C V A L L G F I M W A C
 Q K G N I R C N I C I**

The secondary structures of these consensus amino acid sequences are shown in Figure 1. In Figure 1, there is significant correlation between the H1N1 and H3N2 distributions of helices (Spearman rho=0.8062, p=1.1389 × 10⁻¹³⁰), beta strands (Spearman rho=0.6761, p=7.8722 × 10⁻⁷⁷) and random coils (Spearman rho = 0.5764, p = 2.016 × 10⁻⁵¹).

The secondary structure distributions shown in Figure 1 were used to calculate a helix multiplication product array (equation 1), an extended strand multiplication product array (equation 2) and a random coil multiplication product array (equation 3) by multiplying corresponding secondary structure values at each of the H1N1 and H3N2 amino acid positions. These resulting element-wise multiplication product arrays are shown in Figure 2.

The distributions of H1N1 × H3N2 influenza hemagglutinin secondary structure multiplication product domains shown in Figure 2 are much simpler than the pre-multiplication distributions. Each simplified domain that has a multiplication product equal to or greater than 0.75 and that also contains five or more contiguous amino acids is marked in red. The large multiplication product and high contiguity cut off values were both invoked to minimize the effects of randomization on the outcomes of this study. There are five of these red-marked domains with alpha helical structure, two marked domains with extended strand structure and three marked domains with random coil structure. The amino acid composition of each of these 10 regions is next given in Table 1.

The minimum sequence length in Table 1 is 5 contiguous amino acids and the maximum is 43 contiguous amino acids. Each of the influenza virus amino acid sequences indicated in Table 1 was used as a probe in a BLASTP search of human protein sequences. Because of their length of 43 amino acids, the H1N1 and H3N2 subsequences from viral HA helical region 385-427 were each fractionated into subse-

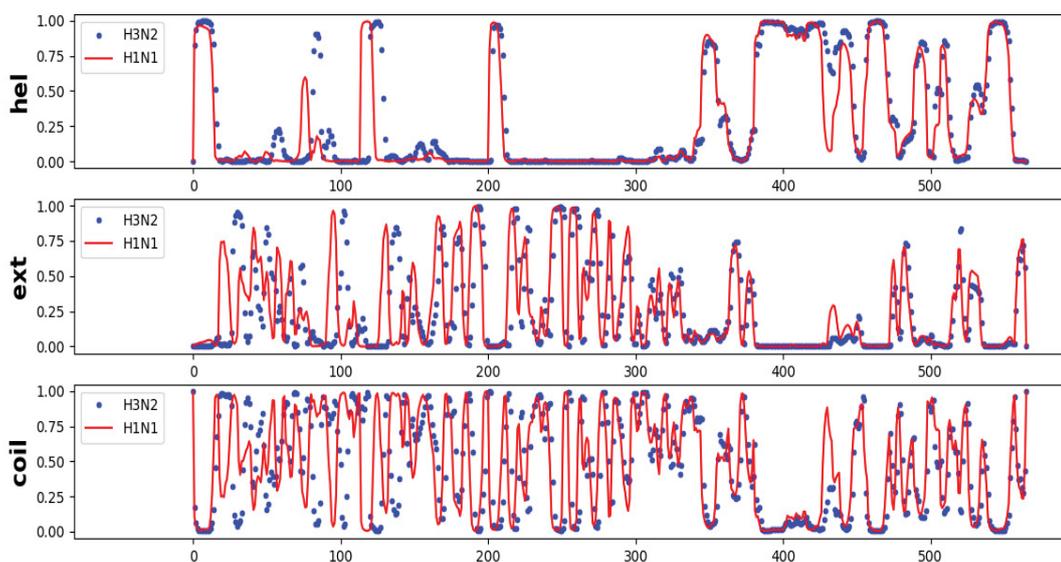


Figure 1: 3-State Secondary Structure of H1N1 and H3N2 Hemagglutinins. hel=alpha helix, ext=extended strand and coil=random coil.

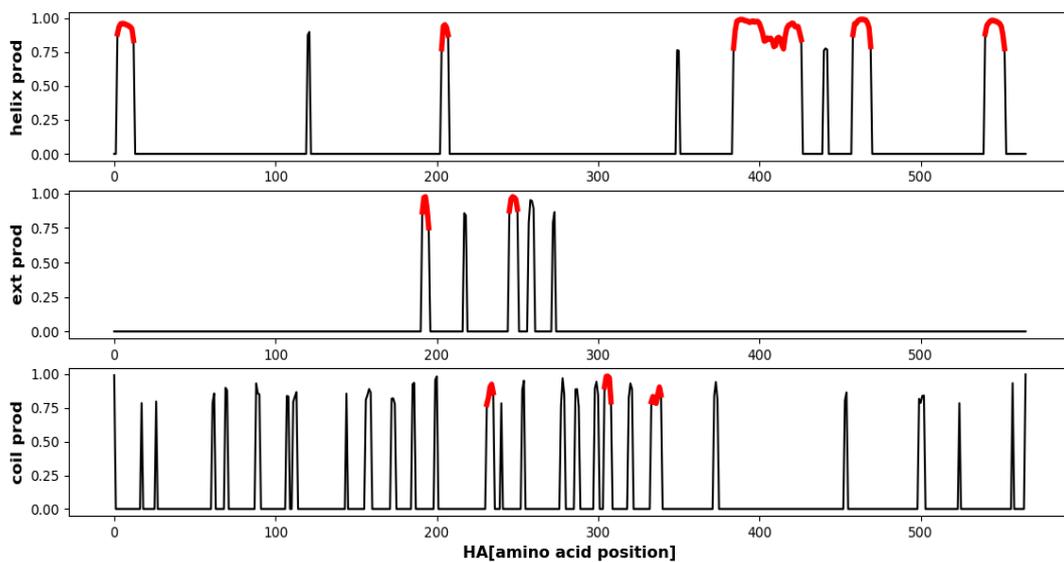


Figure 2: Distribution of H1N1*H3N2 Influenza Hemagglutinin Secondary Structure Multiplication Product Arrays. Element-wise multiplication product arrays are represented by solid black lines. Multiplication product arrays with values ≥ 0.75 and which consist of at least five contiguous amino acids are colored red.

<p>Alpha helix (aa:3-13) len=11 ALVVLLYTFE THIALSYLCL</p>
<p>Alpha helix (aa:204-208) len=5 DQQSL DKDQI</p>
<p>Alpha helix (aa:385-427) len=43 TQNAIDKITNKVNSVIEKMNTQFTAVGKEFNHLEKRIENLNKK STQAAIDQINGKLNRLIGKTNEKFHQIEKEFSEVEGRIQDLEK</p>
<p>Alpha helix (aa:459-470) len=12 VKNLYEKVRSQL EMNKLFEKTKKQ</p>
<p>Alpha helix (aa:541-553) len=13 LVIVVSLGAIISFW FLLCVALLGFIMW</p>
<p>Extended strand (aa:192-196) len=5 VLWGI KLYIW</p>
<p>Extended strand (aa:246-251) len=6 YYWTLV ISIWY</p>
<p>Random coil (aa:232-236) len=5 AIRPK NIGSR</p>
<p>Random coil (aa:305-309) len=5 TSLPF PNDKP</p>
<p>Random coil (aa:334-340) len=7 GLRNVPS TGMRNPV</p>

Table 1: Influenza Hemagglutinin (HA) Consensus Amino Acid Sequences for Use as BLASTP Probes for Identical Subsequences in Human Proteins. H1N1 sequences are colored red and H3N2 sequences are colored blue. Amino acid positions refer to positions within the intact influenza HA protein molecule.

quences of length 20 and 23 amino acids to enable the BLASTP search of proteins of human origin. The longest subsequences each consist of 43 contiguous amino acids and reside in helical domains of the HA molecules. There are 4 pairs of subsequences with the minimum length of 5 contiguous amino acids. Two of these minimum length pairs have random coil structure, one of these pairs consist of extended strands

and one pair consists of alpha helices. The subsequences of greatest length (11, 12, 13 and 43 amino acids) reside in alpha helical domains of the consensus sequences.

The H1N1 and H3N2 influenza HA subsequences shown above in Table 1 were next used as BLASTP probes for detection of matching subsequences in human proteins. Specific examples of the detected human: influenza virus matching subsequences are presented below. The observed length of matching contiguous amino acid components of influenza virus hemagglutinins and human proteins is presented as contiguous length. H1N1 sequences and associated information are coloured red and H3N2 sequences and associated information are coloured blue.

helix aa: 3-13 **contiguous length=6, contiguous length=8**

> XP_011542508.1 desumoylating isopeptidase 2 isoform X3 [Homo sapiens]

Influenza Virus ILVVLL-YT

Human ILVVLLSYT

> 6N4X_A Chain A, Metabotropic glutamate receptor 5 [Homo sapiens]

6N4X_B Chain B, Metabotropic glutamate receptor 5 [Homo sapiens]

Influenza Virus THIALSYLCL

Human THIALSYIFCL

helix aa:204-208 **contiguous length = 5, contiguous length = 5**

> XP_016874588.1 glutamate receptor-interacting protein 1 isoform X2 [Homo sapiens]

Influenza Virus DQQSL

Human DQQSL

> NP_872336.2 myelin regulatory factor-like protein [Homo sapiens]

Q96LU7.2 RecName: Full=Myelin regulatory factor-like protein
> XP_016874456.1 myelin regulatory factor-like protein isoform X1 [Homo sapiens]
> XP_016874457.1 myelin regulatory factor-like protein isoform X2 [Homo sapiens]
> XP_016874458.1 myelin regulatory factor-like protein isoform X3 [Homo sapiens]
> XP_016874459.1 myelin regulatory factor-like protein isoform X4 [Homo sapiens]

> XP_011537354.1 myelin regulatory factor-like protein isoform X5 [Homo sapiens]

Influenza Virus DKDQI
Human DKDQI

helix aa:385-427 contiguous length = 6, contiguous length = 6

> XP_011514027.1 alpha-aminoacidic semialdehyde synthase, mitochondrial isoform X1 [Homo sapiens]

Influenza Virus ITNKVNSV
Human ITNKVNMV

> 6BSZ_A Chain A, Human mGlu8 Receptor complexed with glutamate

6BSZ_B Chain B, Human mGlu8 Receptor complexed with glutamate

Influenza Virus AIDQIN
Human AIDQIN

helix aa:459-470 contiguous length = 5, contiguous length = 7

> XP_005255680.1 protein FAM234A isoform X3 [Homo sapiens]

XP_016879252.1 protein FAM234A isoform X3 [Homo sapiens]

XP_016879253.1 protein FAM234A isoform X3 [Homo sapiens]

Length = 543

Influenza Virus VKNLYEKV
Human VKGLYEKV

> NP_001018126.1 caveolae-associated protein 4 [Homo sapiens]

Q5BXX8.2 RecName: Full=Caveolae-associated protein 4; AltName: Full = Muscle-related coiled-coil protein; AltName: Full=Muscle-restricted

coiled-coil protein ACA62935.1 muscle-restricted coiled-coil protein [Homo sapiens]

Influenza Virus NKLFEKTKK
Human NKLFEKTRK

helix aa: 541-553 contiguous length = 5, contiguous length = 7

>NP_000584.2 antigen peptide transporter 1 isoform 1 [Homo sapiens]

Q03518.2 RecName: Full=Antigen peptide transporter 1; Short=APT1; AltName:

Full=ATP-binding cassette sub-family B member 2; AltName:

Full=Peptide supply factor 1; AltName: Full=Peptide transporter

PSF1; Short=PSF-1; AltName: Full=Peptide transporter

TAP1; AltName: Full=Peptide transporter involved in antigen

processing 1; AltName: Full=Really interesting new gene 4 protein

CAA40741.1 peptide transporter [Homo sapiens]

AAH14081.1 Transporter 1, ATP-binding cassette, sub-family B (MDR/TAP) [Homo sapiens]

EAX03647.1 transporter 1, ATP-binding cassette, sub-family B (MDR/TAP), isoform CRA_b [Homo sapiens]

ALQ33804.1 transporter 1 ATP-binding cassette sub-family B isoform 1, partial [Homo sapiens]

ARB08462.1 TAP1 [Homo sapiens]

ARB08463.1 TAP1 [Homo sapiens]

ARB08464.1 TAP1 [Homo sapiens]

ARB08465.1 TAP1 [Homo sapiens]

ARB08473.1 TAP1 [Homo sapiens]

ARB08474.1 TAP1 [Homo sapiens]

ARB08477.1 TAP1 [Homo sapiens]

ARB08481.1 TAP1 [Homo sapiens]

ARB08482.1 TAP1 [Homo sapiens]

ARB08484.1 TAP1 [Homo sapiens]

prf||1703418A RING4 gene

Influenza Virus LVLVV--SLG--AISF

Human LVLVVLSSLGEMAIPE

> AAH15195.1 Fc fragment of IgE, high affinity I, receptor for; alpha polypeptide [Homo sapiens]

Influenza Virus LLCVALL

Human LLCVALL

extended aa:192-196 contiguous length = 5, contiguous length < 5

> AAI42999.1 Asparagine-linked glycosylation 11, alpha-1,2-mannosyltransferase homolog (yeast) [Homo sapiens]

Influenza Virus VLWGI

Human VLWGI

extended aa:246-251 contiguous length = 5, contiguous length = 5

> NP_036550.1 fascin-2 isoform 1 [Homo sapiens]

O14926.1 RecName: Full=Fascin-2; AltName: Full=Retinal fascin AAB86481.1 retinal fascin [Homo sapiens]

AAC18604.1 retinal fascin [Homo sapiens]

AAI26296.1 Fascin homolog 2, actin-bundling protein, retinal (Strongylocentrotus

purpuratus) [Homo sapiens]

AAI30331.1 Fascin homolog 2, actin-bundling protein, retinal (Strongylocentrotus purpuratus) [Homo sapiens]

Influenza Virus YWTLV

Human **YWTLV**
> CBX25812.1 cytochrome c oxidase subunit III, partial (mitochondrion) [Homo sapiens subsp. 'Denisova']

Influenza Virus **ISIIYW**
Human **ISIIYW**
coil aa:232-236 contiguous length = 5, contiguous length = 5
> NP_001034794.1 trophinin isoform 5 [Homo sapiens]
XP_006724663.1 trophinin isoform X1 [Homo sapiens]
XP_011529110.1 trophinin isoform X1 [Homo sapiens]
XP_011529111.1 trophinin isoform X1 [Homo sapiens]
XP_011529113.1 trophinin isoform X1 [Homo sapiens]
XP_011529114.1 trophinin isoform X1 [Homo sapiens]
XP_011529115.1 trophinin isoform X1 [Homo sapiens]
XP_016885256.1 trophinin isoform X1 [Homo sapiens]
Q12816.3 RecName: Full=Trophinin; AltName: Full=MAGE-D3 antigen
BAA83066.4 KIAA1114 protein [Homo sapiens]

Influenza Virus **AIRPK**
Human **AIRPK**
> XP_011525670.1 zinc finger protein 541 isoform X1 [Homo sapiens]
XP_011525671.1 zinc finger protein 541 isoform X1 [Homo sapiens]
XP_011525672.1 zinc finger protein 541 isoform X1 [Homo sapiens]
XP_011525673.1 zinc finger protein 541 isoform X1 [Homo sapiens]

Influenza Virus **NIGSR**
Human **NIGSR**
coil aa:305-309 contiguous length = 5, contiguous length = 5
> AAO15766.1 a disintegrin-like and metalloprotease with thrombospondin type
1 motifs 20 [Homo sapiens]

Influenza Virus **TSLPF**
Human **TSLPF**
> XP_016865452.1 G-protein coupled receptor 98 isoform X1 [Homo sapiens]

Influenza Virus **PNDKP**
Human **PNDKP**
coil aa:334-340 contiguous length = 6, contiguous length = 5
> EAW73283.1 BCL2-interacting killer (apoptosis-inducing), isoform CRA_b [Homo sapiens]

Influenza Virus **LRNVPS**

Human **LRNVPS**
>NP_078922.1 protein zyg-11 homolog B [Homo sapiens]
Q9C0D3.2 RecName: Full=Protein zyg-11 homolog B
BAB21821.2 KIAA1730 protein [Homo sapiens]

Influenza Virus **TGMRNVP**
Human **TGMRNHP**

Each of the 20 molecular probes of influenza HA protein origin detected matching identical subsequences of contiguous amino acids of human origin except **KLYIYW**, a pentameric subsequence of contiguous amino acids found in H3N2 hemagglutinins (aa: 192-196) but which was not detected in human proteins. The largest detected human: influenza matching subsequence pair was **THIALSYI** which is a contiguous octameric component of the amino acid 3-13 structural domain of the H3N2 hemagglutinin.

Discussion

This research is a continuation of the previous report of the occurrence of identical tetrapeptide subsequences in influenza hemagglutinins and in human proteins [4]. The present report extends the analysis to identical subsequences of at least pentapeptide length. The nineteen amino acid subsequences reported here as being identical in humans and influenza virus consist of amino acid polymers of contiguous amino acids between 5 and 8 in polymer length. The cut off for this study was increased to 5 amino acids because randomization of secondary structure has been reported to predominate in amino acid sequences below, but not above, pentameric length [9,10].

The methodology used is based upon the similar secondary structural features of H1N1 and H3N2 HA proteins (Figures 1 and 2). Structural features common to subtype H1N1 and H3N2 HA proteins were determined and were used to identify influenza virus HA subsequences for use as probes for subsequences in human proteins. Use of those influenza virus hemagglutinin subsequence probes enabled the detection of human subsequences identical those of the virus.

The occurrence of identical host-virus amino acid subsequences may reduce immunogenicity of the influenza virus in the infected host by favouring immunological tolerance to the organism [11]. Such tolerance would increase the susceptibility of the host to influenza infection and would reduce protective effectiveness of anti-influenza vaccines.

Conclusion

It is reported here that the influenza virus hemagglutinin and various proteins of human origin contain identical subsequences of amino acids. These observed subsequences are 5-8 amino acids in length and thus may have significant biological effects such as induction of immune tolerance to potentially protective epitopes of the influenza virus.

There is a need for an anti-influenza vaccine of greater efficacy. It is proposed that influenza hemagglutinin proteins and peptides lacking human subsequences should be prepared and tested for effective antiviral immunogenicity.

References

1. Iuliano AD, Roguski KM, Chang HH, Muscatello DJ, Palekar R, et al. (2018) Estimates of global seasonal influenza-associated respiratory mortality: A modelling study. Lancet 391: 1285-1300.
2. World Health Organization (2017) WHO Preferred Product Characteristics for Next-Generation Influenza Vaccines.

3. Doyle JD, Chung JR, Kim SS, Gaglani M, Raiyani C, et al. (2019) Interim estimates of 2018-19 seasonal influenza vaccine effectiveness - United States, February 2019. *MMWR Morb Mortal Wkly Rep* 68: 135-139.
4. Weltman JK (2018) Shannon entropy screening of influenza hemagglutinin for tetrapeptides with exact homology to human proteins. *J Med Microb Diagn* 7: 284-286.
5. Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ (2009) Jalview Version 2- a multiple sequence alignment editor and analysis workbench. *Bioinformatics* 25: 1189-1191.
6. Shannon CE (1948) A mathematical theory of communication. *Bell Syst Tech J* 27: 379-423.
7. Källberg M, Wang H, Wang S, Peng J, Wang Z, et al. (2012) Template-based protein structure modeling using the RaptorX web server. *Nat Protoc* 7: 1511-1522.
8. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, et al. (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 25: 3389-3402.
9. Figureau A, Soto MA, Toha J (2003) A pentapeptide-based method for protein secondary structure prediction. *Protein Eng* 16: 103-107.
10. Turjanski P, Ferreiro DU (2018) On the natural structure of amino acid patterns in families of protein sequences.
11. Chirino AJ, Ary ML, Marshall SA (2004) Minimizing the immunogenicity of protein therapeutics. *Drug Discov Today* 9: 82-90.