

Perspectives

Genesis of non-coding RNA genes- a sequence connection with protein genes separated by evolutionary time

Nicholas Delihias

Department of Microbiology and Immunology, Renaissance School of Medicine, Stony Brook University, Stony Brook, New York, 11794-5222 USA

Email: Nicholas.delihias@stonybrook.edu

Abstract: A small phylogenetically conserved sequence of 11,231 bp termed FAM247 is repeated in human chromosome 22 by segmental duplications. This sequence forms part of diverse genes that span evolutionary time, the protein genes being the earliest as they are present in zebrafish and/or mice genomes, the long non-coding RNA genes and pseudogenes the most recent as they appear to be present only in the human genome. We propose that the conserved sequence provides a nucleation site for new gene development at evolutionary conserved chromosomal loci where the FAM247 sequences reside. The FAM247 sequence also carries information in its open reading frames that provides protein exon amino acid sequences; one exon plays an integral role in immune system regulation, specifically, the function of ubiquitin specific protease (USP18) in the regulation of interferon. An analysis of this multifaceted sequence and the genesis of genes that contain it are presented.

Keywords: gene evolution; gene formation; long non-coding RNA genes; pseudogenes; *USP18*; *GGT5*

1. Introduction

The genesis of genes has been a major topic of interest for several decades [1, 2]. One mechanism of gene origins is the formation of genes from the duplication of existing genes [1, 3]. This is considered one of the major processes for new gene formation, but it has also been shown that there is a prevalence of gene birth from non-coding DNA via de novo gene formation [4-7]; this pathway also significantly contributes to new gene formation [4, 7].

In this treatise we present and analyze gene development by an evolutionary conserved ancestral sequence. This is a repeat element, previously termed clincRNA [8] and now termed FAM247. The long intergenic non-coding RNA (lincRNA) *FAM247A* gene sequence has been used as a guide to find homologous sequences and heretofore FAM247 is used in place of *FAM247A*. We propose that FAM247 carries information to form nucleation sites for gene development and this is exemplified with the maturation of pseudogenes by addition of chromosomal sequences at specific sites on FAM247. The FAM247 sequence also carries information in terms of its open reading frames. These open reading frames are found to form exon sequences of proteins.

FAM247 forms part of non-coding RNA genes that appear to be human specific. The sequence is also found in protein genes, *USP18* (ubiquitin specific protease) and *GGT5* (gamma glutamyltransferase). Both these genes date back in evolutionary time, *USP18* over 350 million years ago (MYA) and *GGT5* over 90 MYA. Thus, the FAM247 sequence has formed a part of genes through much of vertebrate evolution.

2. Background on conserved linked sequences

FAM247 is present in different segmental duplications or low copy repeats (LCR22) in human chromosome 22 (chr22) as part of phylogenetically conserved linked gene sequences [8]. Figure 1 is a representation of these linked sequences and shows conserved nearest neighbor sequence signatures found in humans. The linked gene sequences are repeated in chr 22 and generate gene families. The signatures are also representative of ancestral primate linked-sequences, e.g., the sequence arrangement in Figure 1b is present in the Rhesus Old World monkey (*Macaca mulatta*) where FAM247 and spacer sequences are linked to *GGT1* on chr10. The spacer sequence (3953 bp) depicted in Figure 1 is also evolutionarily conserved. It is present in *Pan troglodytes* (chimpanzee), *Papio Anubis* (olive baboon), *Pongo abelii* (Sumatran orangutan), and *Macaca mulatta* (Rhesus monkey) genomes but it does not encode genes or form part of genes [8]. That it is evolutionarily conserved indicates it may have a function. FAM247 is the common denominator in Figure 1a and 1b. In Figure 1c, the FAM247 sequence depicted is embedded the *USP18* gene.

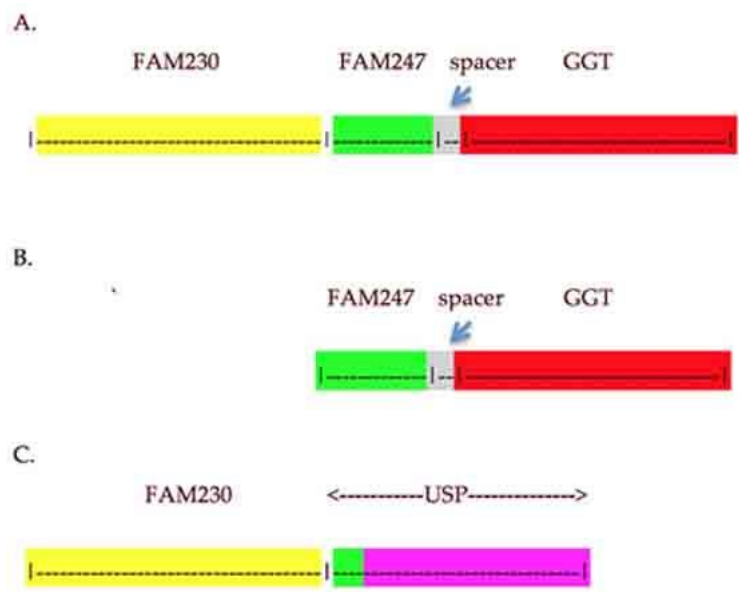


Figure 1. Schematic representation of evolutionarily conserved linked sequences with different colors depicting different sequences, as described in [8]. In Figure 1c, the FAM247 sequence (green highlight) is embedded the *USP18* gene.

Table 1 contains a list of human gene families that are found in repeat units shown in Figure 1 and indicates the sequence or chromosomal locus of origin. For example, GGT represents the locus of origin of *GGT1*, the gamma-glutamyltransferase and gamma-glutamyltransferase light chain genes and their respective pseudogenes; FAM247 is the sequence/locus of origin of *GGT5*.

Table 1. Human genes and gene families found in linked sequences: FAM230-FAM247-GGT, FAM230-USP.

Gene/gene family	Type	Locus origin
*FAM230A-J	lincRNA	FAM230
FAM247A-D	lincRNA	FAM247
POM121L9P, POM121L10P	pseudogene	FAM247
BCRP3	pseudogene	FAM247
GGT1, GGT2	protein	GGT
GGTLC2	protein	GGT
GGTLC3	protein	GGT
GGT3P	pseudogene	GGT
GGT4P	pseudogene	GGT
GGTLC4P	pseudogene	GGT
GGTLC5P	pseudogene	GGT
GGT5	protein	FAM247
USP18	protein	FAM247

*Some FAM230 family members such as FAM230J do not have the linked sequence signatures.

A description of genes is as follows. *POM121L9P* and *POM121L10P* are member of the POM121 transmembrane nucleoporin like 1 pseudogene family. *BCRP3* is a member of the BCR pseudogene family. *BCR*, a large gene of 137,529 bp is the activator of RhoGEF and GTPase and formerly termed breakpoint cluster region protein. The *BCR* gene is important clinically as it is associated with the production the Philadelphia chromosome in chronic myelogenous leukaemia [9, 10]. *POM121L9P* and *BCRP3* stem from the FAM247 sequence at chromosomal loci where the GGT sequence is found, as represented in Figure 1b. *USP18* is the ubiquitin specific peptidase gene, a member of the deubiquitinating protease family; the protein product plays a major role in interferon regulation [11] and has multiple functions [12].

3. lincRNA gene families

The FAM230 lincRNA and FAM247 lincRNA gene families (<https://www.genenames.org/>) exemplify how segmental duplications or low copy repeats in chromosome 22 are a driving element in the genesis and proliferation of lincRNA gene families. Ten FAM230 and five FAM247 genes are present in chr22 low copy repeats (LCR22) that originate from sequence duplications [8]. FAM230 family genes differ from one another in sequence, transcript sequence and exon number, and in RNA expression in various fetal developing tissues [8, 13, 14]. Their functions are not known. FAM230 sequences are also present in primates, but these sequences are annotated as predicted protein genes or pseudogenes, not as lincRNA genes, e.g., more than eleven genes that contain the FAM230 sequence in chimpanzee and gorilla are annotated as protein genes and two FAM230 sequences in Rhesus monkey and olive baboon are annotated as pseudogenes. An example is *LOC106992440*, which is found in the Rhesus monkey is

annotated as an uncharacterized pseudogene and resides in a chr locus that is homologous to that of human *FAM230D* linked to *USP18* in humans; the Rhesus *LOC106992440* has 58% identity with the human lincRNA *FAM230D*. In this and in other cases, the detection of an experimental transcript that verifies the computational prediction of a protein gene or pseudogene is essential. Evolutionarily, the FAM230 sequence may have originated in the Rhesus monkey or Old World monkeys as the FAM230 sequence is present in the Rhesus genome but is not found in the genome of the Prosimian primate ancestor, *Philippine tarsier*.

The *FAM247* lincRNA gene family may have newly formed in humans as there are few or no differences in gene sequence or in RNA transcript expression in somatic and fetal developing tissues [8, 13]. An homologous sequence to FAM247 is present in chimpanzee and is linked to *GGT2*. It has the full length FAM247 sequence [8], but the chimpanzee sequence has not been annotated as encoding a gene. Segments of sequences homologous to FAM247 are found in other primate genomes (gorilla, orangutan, Rhesus monkey, *Philippine tarsier*), and these FAM247 sequences may date back evolutionarily to the house mouse and zebrafish (discussed below).

4. The FAM247 sequence is present in diverse genes

A significant property of the FAM247 sequence is that it forms part of diverse genes. Sequences homologous to FAM247 form genes that include lincRNA genes, pseudogenes, and protein genes (Figure 2). These genes stem from phylogenetically conserved nearest neighbor gene loci where the FAM247 sequence is linked to adjacent genes that form signatures containing gene families e.g., *FAM230E-FAM247C-GGT3P* present in segmental duplication LCR22A and *FAM230B-FAM247A-GGT2* in LCR22D. Other than the FAM247 lincRNA family genes, which contain the entire 11,231 bp, only segments of FAM247 are found to be part of other genes. The ends of these segments represent sequence breaks, i.e. bp positions ~5958 and ~8000-8200 (Figure 2, see numbers above green highlighted FAM247). These are regions that contain *Alu* sequences at their ends (Figure 2, caption), and they may provide sites for attachment of other sequences to FAM247. FAM247 contains a total of fifteen *Alu* elements.

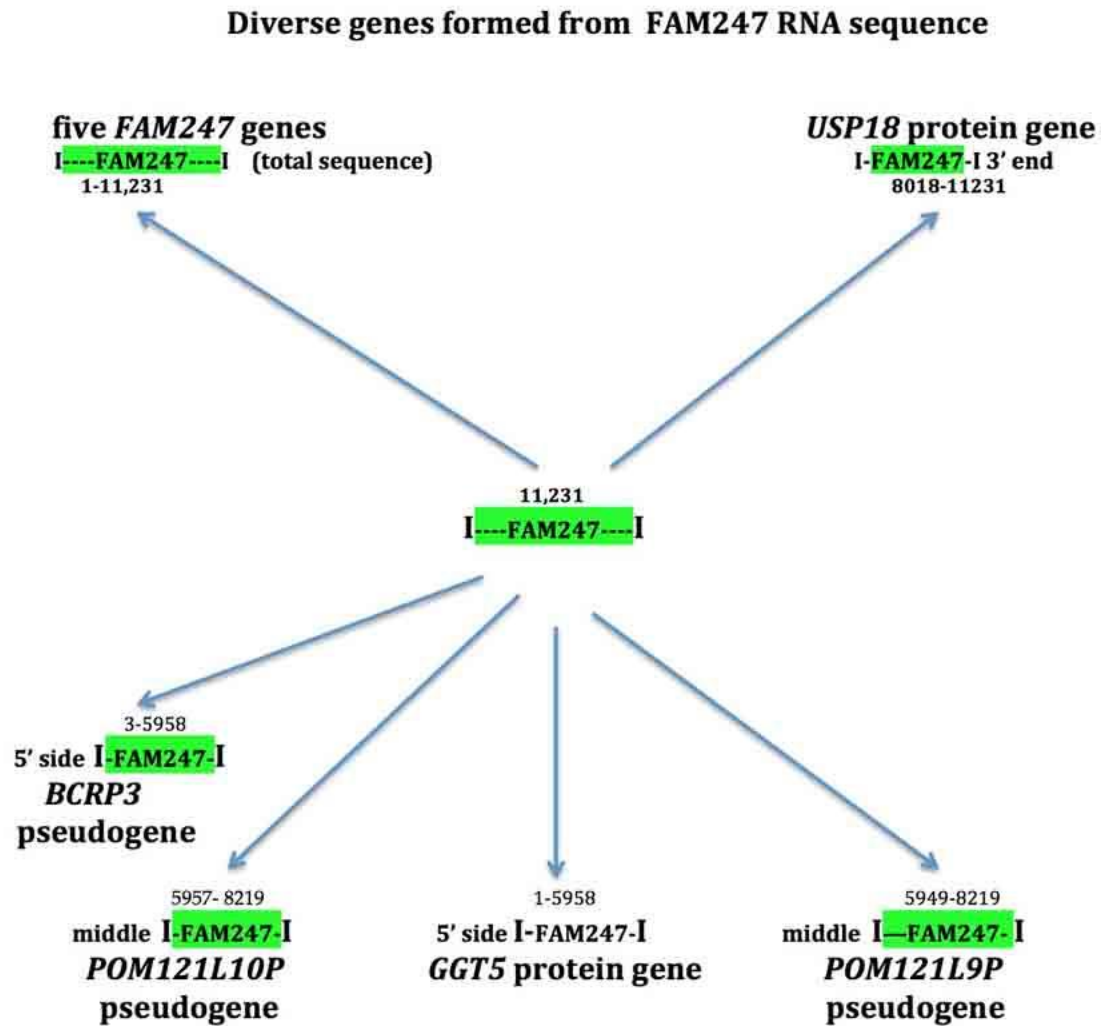


Figure 2. Protein genes, lincRNA genes and pseudogenes that stem from FAM247 sequence and contain different sections of the FAM247 sequence (shown in bp position numbers above FAM247 highlighted in green). Analysis of FAM247 by the RepeatMasker program:

RepeatMasker <http://www.repeatmasker.org/cgi-bin/WEBRepeatMasker> shows the presence of *Alu* sequences in the FAM247 sequence at bp positions 6007-6285 and at 8063-8374, regions close to breaks. Region 8063-8374 bp has seven *Alu* elements in tandem repeats.

4.1. *USP18*

A comparison of *USP18* chromosomal coordinates at loci in different species shows that the position is evolutionarily conserved relative to adjacent genes (Figure 3). This provides nearest neighbor signatures where evolutionary history and origins of *USP18* can be assessed. Two neighbor genes, *PEX26* (peroxisomal biogenesis factor 26) and *TUBA8* (tubulin alpha 8) are in homologous loci that show a conserved orientation with respect to each other in the chromosomes of mice [*Mus musculus* (house mouse)] and primates (Figure 3a-c). In zebrafish (*Danio rerio*, a member the Cyprinidae family of freshwater fish), the tubulin gene (termed tuba8l4 tubulin, alpha 8 like 4) appears to have moved to a different chromosome and developed into two genes, *TUBAa* and *TUBAb*; this results in *PEX28* and *USP18* as immediate neighbor genes in zebrafish chr4 (Figure 3d). The nearest neighbor history of the *USP18* gene locus and the display of evolutionary conservation of gene positions relative to each other are consistent with a common lineage of the *USP18* gene.

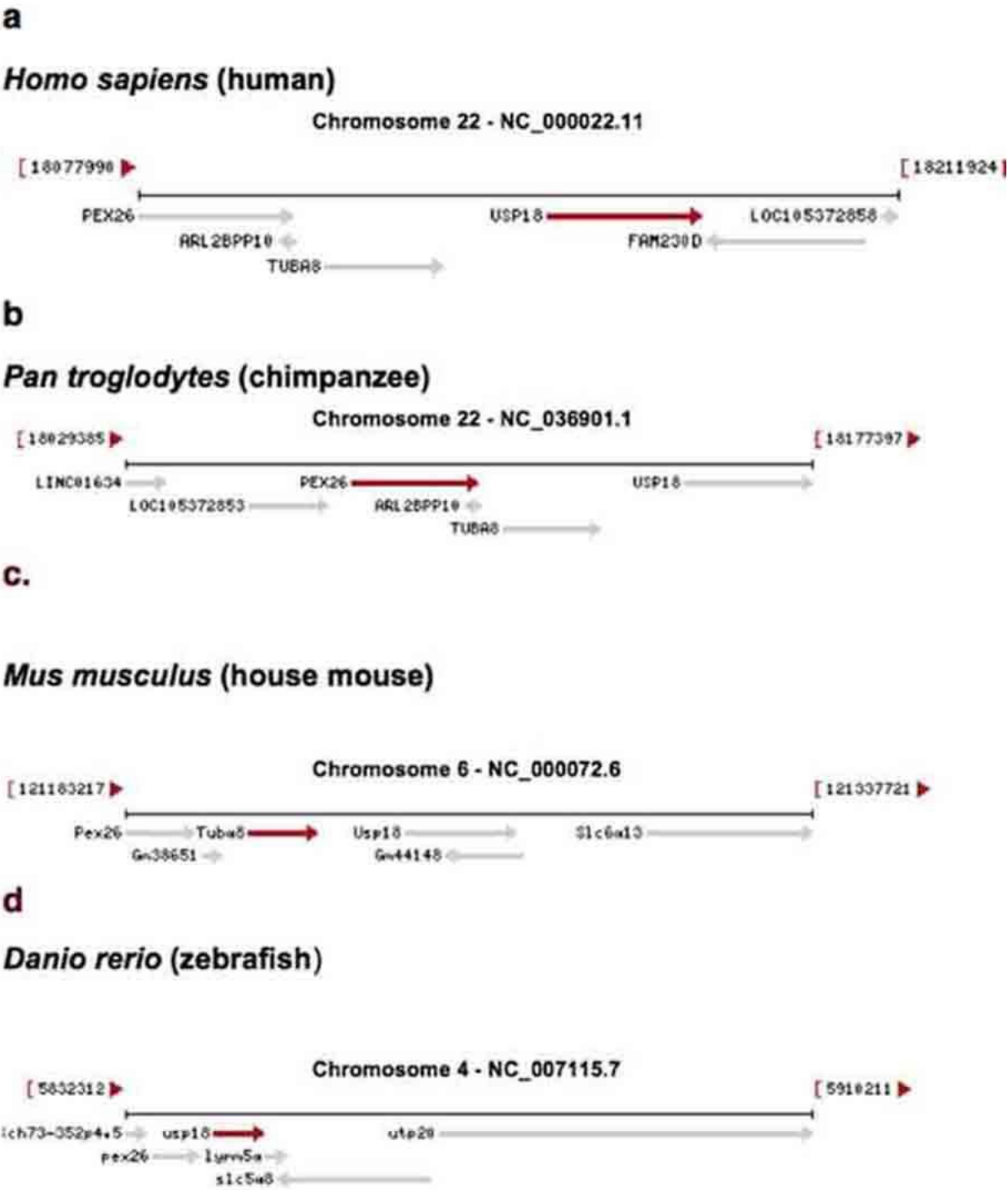


Figure 3. Genes that are adjacent to *USP18* are found in different species and are shown in the gene maps. Drawings of gene arrangement are taken directly from the NCBI website: <https://www.ncbi.nlm.nih.gov/gene> [15].

To analyze the phylogenetic relatedness of *USP18* gene nt and aa sequences, sequences were aligned from zebrafish, the house mouse, three primate species and humans. The resultant percent sequence identities mimic evolutionary distances between species (Table 2) with a linear change in nt and aa sequence with time between the primate and mouse species (not shown). The pattern shows a continuum

of gene nt and protein aa sequence change with evolutionary time and is consistent with a common lineage of the *USP18* gene that dates to an ancestor of zebrafish, more than about 350 million years ago (MYA). This parallels the nearest neighbor gene history of *USP18*.

Table 2. *USP18* gene and protein sequence identities and evolutionary time between species

Species	USP18 gene nt sequence % identity	USP18 aa sequence % identity	Evolutionary age (MYA)*
human	100%	100%	0 MYA
chimpanzee	99%	99%	6 MYA
Rhesus monkey	92%	94%	25 MYA
<i>Philippine tarsier</i>	66%	80%	50 MYA
mouse	51%	71%	90 MYA
zebrafish	39%	31%	350 MYA

*approximate age in million years ago (MYA) [16]

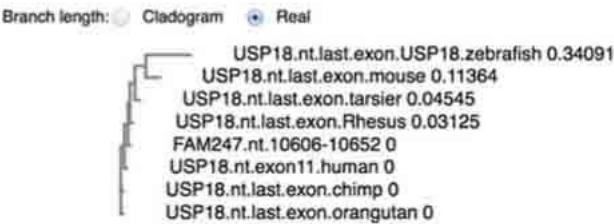
4.2. *USP18* exon 11

Both human exon 11, which encodes the last 14 aa (the carboxy terminal end) of the USP18 peptidase, and the 3' UTR of the USP18 mRNA sequence are provided by the FAM247 sequence [8]. The identity between the FAM247 nt sequence and the human/primate exon11 nt sequences is 100%, with the exception of that of Philippine tarsier (Figure 4). The sequence of the carboxy terminal exon is more stable than that of the sequence of entire gene (compare with Table 2). The identities of the *USP18* 3'UTR sequences from various species compared to FAM247 (Table 3) shows the 3' UTR sequence is also conserved in primates, but to a lesser extent than that of that of exon 11 and is more similar to the *USP18* gene.

A.

Phylogenetic Tree

This is a Neighbour-joining tree without distance corrections.



B.

Nucleotide sequence identities of last exon of *USP18* compared to FAM247

Source of nt sequence	Percent identity relative to FAM247
FAM247.nt.10606-10652	100
USP18.nt.exon11.human	100
USP18.nt.last.exon.chimp	100
USP18.nt.last.exon.orangutan	100
USP18.nt.last.exon.Rhesus.monkey	100
USP18.nt.last.exon.Philippine.tarsier	91
USP18.nt.last.exon.mouse	75
Usp18.nt.last .exon.zebrafish	52.3

C.

Alignment of nt sequences from last exon of *USP18* and FAM247

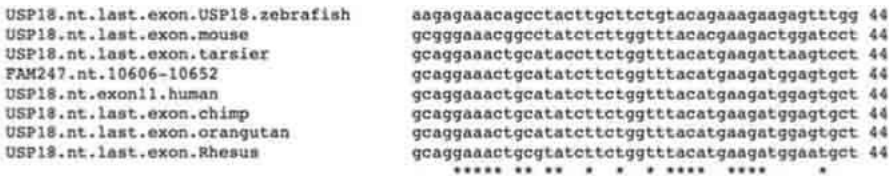


Figure 4. Alignment of the USP18 terminal exon nt sequences from seven species compared with the FAM247 sequence. Data obtained using the EBI Clustal Omega sequence alignment and phylogeny programs. The EMBL-EBI Clustal Omega Multiple Sequence Alignment program [17] at website: <http://www.ebi.ac.uk/Tools/msa/clustalo/> was used for nt sequence alignment. A. Phylogenetic tree of USP18 terminal exon sequences from seven species and the FAM247 sequence. B. The percent identities created using Clustal 2.1. C. Alignment of nt sequences.

The sequence similarity of 52% between FAM247 and zebrafish exon 11 (Figure 4B), the presence of a number of invariant nt positions (Figure 4C) and the similarity with the 3'UTR sequence (53%) (Table 3) suggests that this part of the FAM247 sequence was present in the *USP18* sequence of zebrafish. The invariant nt residues of exon 11, e.g., positions nt 5-9 (Figure 4C) may relate to the functional importance of the *USP18* carboxy terminal end in its role in the regulation of the immune system by *USP18* [11, 12]. These invariant nt positions may give a picture, albeit a small picture of what the ancient FAM247 looked like. *Alu* elements precede the homologous sequence of Exon 11 in FAM247, but do not overlap it.

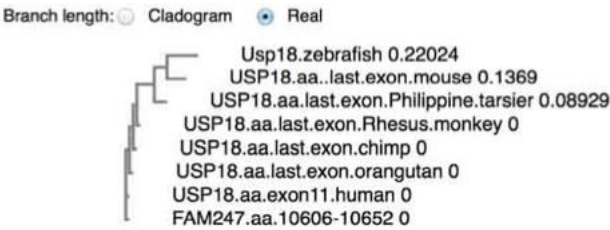
Table 3. Nucleotide sequence identities of 3'UTR of *USP18* from different species

Source of nt sequence	% identity relative to FAM247 3' end
FAM247A 3' end nt 10653-11231	100
USP18 3'UTR human	99.8
USP18 3'UTR chimp	98.6
USP18 3'UTR Rhesus monkey	90.2
USP18 3'UTR Philippine tarsier	71.9
USP18.3'UTR mouse	49.5
Usp18.3'UTR zebrafish	53.1

Figure 5 shows USP18 aa sequence percent identity, sequence alignment and a phylogenetic tree produced from an alignment of *USP18* terminal exon aa sequences from different species with the translated aa sequence of FAM247. Eight of the 14 amino acid residues that form the terminal exon are totally conserved from primates to zebrafish, together with FAM247 translated aa sequence. (Figure 5c). The USP18 carboxy terminal peptide sequence interacts with the INFAR2 interferon receptor and is an important regulator of IFN signaling [11]; in addition, the carboxyl end sequence functions in deISGylation [18, 19, 12]. A mutation in L365 in the exon 11 sequence ³⁵⁹QETAYLL₃₆₅VYMKMEC₃₇₂ abolishes deISGylation and INAFR2 binding [19]; L₃₆₅ is one of the evolutionary conserved amino acids of exon 11 (Figure 5). The mutation may alter the protein conformation necessary for function. On the other hand, the high number aa residues conserved relative to the FAM247 translated aa sequence adds to the suggestion the FAM247 sequence was present in zebrafish *USP18*.

Phylogenetic Tree

This is a Neighbour-joining tree without distance corrections.



Amino acid sequence identities of last exon of *USP18*

Source of aa sequence	Percent identity relative to FAM247
FAM247.aa.10606-10652	100
USP18.aa.exon11.human	100
USP18.aa.last.exon.chimp	100
USP18.aa.last.exon.orangutan	100
USP18.aa.last.exon.Rhesus.monkey	100
USP18.aa.last.exon.Philippine.tarsier	78.6
USP18.aa..last.exon.mouse	64.3
Usp18.zebrafish	57.1

Amino acid sequence of last exon of *USP18*

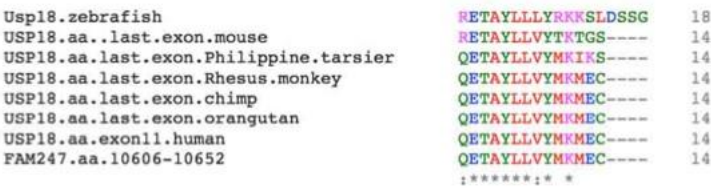


Figure 5. Alignment of the USP18 terminal exon amino acid sequences from seven species compared with the FAM247 translated amino acid sequence. Data obtained using the EBI Clustal Omega sequence alignment and phylogeny programs. The EMBL-EBI Clustal Omega Multiple Sequence Alignment program [17] at website: <http://www.ebi.ac.uk/Tools/msa/clustalo/> was used for aa sequence alignment. Top. Phylogenetic tree of USP18 terminal exon aa sequences from seven species and the FAM247 aa sequence. Middle. The percent identities created using Clustal 2.1. Bottom. Alignment of aa sequences.

4.3. *GGT5*

The human *GGT5* protein gene resides in chromosomal segmental duplication LCR22G and is linked to pseudogene *POM121L9P* with a spacer sequence and the pseudogene *GGTLC4P* situated between *GGT5* and *POM121L9P* (Figure 6) [8]. The *GGT5* nearest gene/sequence arrangement is more complex than that of the signatures shown in Figure 1b. *GGT5* is an anomaly as its sequence does not stem from a GGT locus, as other GGT family members do, but from the chromosomal site containing the FAM247 sequence [8]. *GGT5* carries a sequence homologous to the 5' half of the FAM247A sequence, bp positions 1-5958 (Figures 2 and 6) and *POM121L9P* contains part of the 3' half of FAM247 (bp 5949-8219). The *GGTLC4P* pseudogene derives its sequence from GGT (Figure 6).

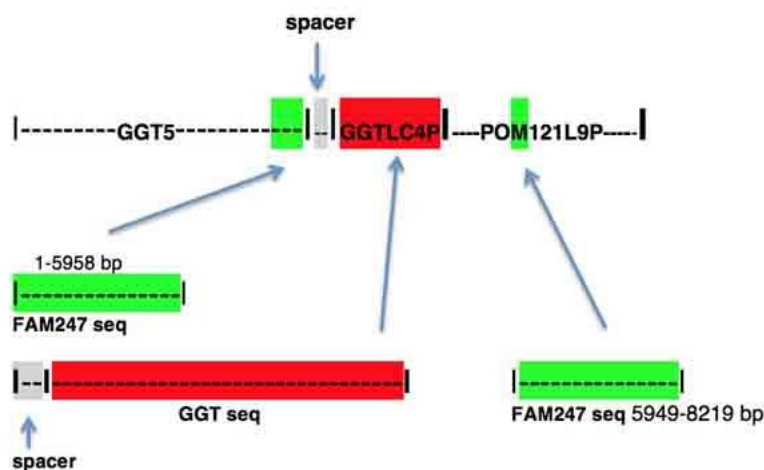


Figure 6. The genes linked to *GGT5* in LCR22G with nearest neighbor arrangements (top schematic) and the source of sequences found in human linked genes *GGT5*-*GGTLC4P*-*POM121L9P*.

GGT5 and *POM121L9P* appear to have formed at very different evolutionary times. FAM247 is part of *GGT5* genes in non-human primates, including *Philippine tarsier*. In addition, FAM247 provides the sequence for exon 1 of *GGT5*. There is a significant similarity between the FAM247 nt sequence and that of the mouse *GGT5* exon 1 (Figure 7A). There is not enough evidence to suggest that the mouse *GGT5* contains the entire 5' half of the FAM247 sequence but alignment of the mouse exon 1 nt sequence with FAM247 shows that a significant number of nucleotides are invariant (Figure 7B). Although there is invariance in 50 out of 173 nt between the FAM247 sequence and zebrafish *GGT5* exon 1, the zebrafish exon sequence shows significant differences, which makes it difficult to further assess a sequence similarity. The exon 1 data are consistent with the formation of the *GGT5* gene with the FAM247 sequence that occurred before the evolutionary appearance of primates and appearing in mice or an ancestor to mice.

A. Percent identity of nucleotide sequences from *GGT5* exon 1 compared to FAM247

Source of exon 1 nt sequence	Percent identity relative to FAM247
FAM247 nt. 392-567	100
Exon 1 <i>GGT5</i> human	88
Exon 1 <i>GGT5</i> Philippine tarsier	80
Exon 1 <i>GGT5</i> mouse	69
Exon 1 <i>GGT5b</i> zebrafish	46

B. Alignment of nucleotide sequences from *GGT5* exon 1 and FAM247

exon1.GGT5b.zebrafish	atggccaaa-----tctcagtcgaggcgatgctgcttttgtttactcgtttagt	50
exon1.GGT5.mouse	atgggttgggttcacaggccacggctcgtcctggtcctgc-----tgggtgtaggtcta	54
FAM247.392-567	atggccaggactacggagccatgggtgacctggtcctgctggggctggggctggggctg	60
exon1.Philippine.tarsier	atggcctggggctgcaggcccatcagcctggtcctgctggggctggggctggggctg	60
exon1.GGT5.human	atggccggggctacggggccacggctcagcctagtcctgc-----tgggtgtggggctg	54
	***** * * *	
exon1.GGT5b.zebrafish	gtgcactgctgcattatctgcatttgcatactt-----ttcagcaaacagaa	98
exon1.GGT5.mouse	g-gtctgggttatcgttgtgttggctgcggctccttctcctcgtcaagcctcttgggtcc	113
FAM247.392-567	g-cgctggctgtcattgtgctggctgtggtcctctctcgcacacaggcccatgtgaccc	119
exon1.Philippine.tarsier	g-cactagtcatttctgctggctgtggtcctcctcgcacacaggcccatgtgaccc	119
exon1.GGT5.human	g-cgctggctgtcattgtgctggctgtggtcctctcgcacacaggcccatgtgaccc	113
	* * * * *	
exon1.GGT5b.zebrafish	atgcgactttactcgggcgcgtgtttctgcggactctctcatgtgctcggacatcggcag	158
exon1.GGT5.mouse	cgggtccttcacgcgtgctgcggtagcggctgactccaagatctgctcggatattggacg	173
FAM247.392-567	-cgggcctttgcccacgcgcgtgttgcgtgactccaaggtcttctcaaatattgtacg	178
exon1.Philippine.tarsier	ccaggcctttgcccacgcgcgtgttgcgtgactccaaggtctgctcggacattggacg	179
exon1.GGT5.human	ccaggcctttgcccacgcgcgtgttgcggcggactccaaggtctgctcggatattggacg	173
	* * * * *	
exon1.GGT5b.zebrafish	g	159
exon1.GGT5.mouse	-	173
FAM247.392-567	-	178
exon1.Philippine.tarsier	g	180
exon1.GGT5.human	-	173

Figure 7. A. The percent identity of FAM247 sequence with that of *GGT5* exon 1 sequences from four species. Data were obtained using Clustal 2.1. **B.** Alignment of the *GGT5* exon nucleotide sequences from the four species, compared with the FAM247 sequence, positions 192-567. Data obtained using the EBI Clustal Omega sequence alignment and phylogeny programs. The EMBL-EBI Clustal Omega Multiple Sequence Alignment program [17] at website: <http://www.ebi.ac.uk/Tools/msa/clustalo/>

4.4. *POM121L9P*

POM121L9P has a very different sequence from the other *POM121L* family pseudogenes and it is a unique gene. A schematic of the compositional make-up of the *POM121L9P* gene shows that the pseudogene contains most of the sequence homologous to the putative parent gene, protein gene *POM121L1* (LOC101929738 putative POM121-like protein 1, 2379 bp) on its 5' side, and the *BCRP1* pseudogene sequence (that is homologous to the 3' section of the *BCR* gene that includes *BCR* terminal exons 19-23) on its 3' side (Figure 8). FAM247 may have formed a nucleation site for addition of these motifs, which are copies of sequences from different regions of the genome. The sequence motifs are found attached onto 5' and 3' ends of FAM247 at bp positions where there are *Alu* sequences (FAM247 bp position 5949 and 8219). Also, the complete *POM121L-1* sequence has an *Alu* sequence at bp positions 2309 -2379 (the end of *POM121L-1* is position 2304, the attachment site with FAM247). *Alu* sequences may have facilitated the addition of *POM121L-1* to FAM247. The *BCR* sequence addition to *POM121L9P* is more complex as there is an undefined sequence between the two (bp position 4479-5779 on *POM121L9P*), and there are no *Alu* sequences detected in the *BCR* sequence at the junction site. The human *POM121L9P* pseudogene RNA transcript is highly expressed in somatic testis tissue, and there is broad expression of circular RNAs in developing fetal tissues with major expression in lung and adrenal

tissues [13, 14]. Its functions are not known but should be of interest in view of the strong RNA expression levels.

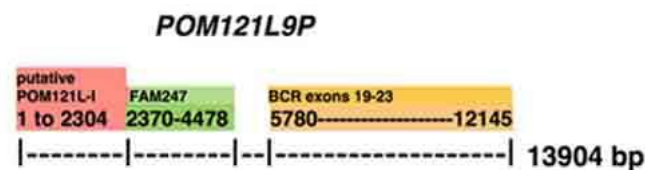


Figure 8. A schematic of the compositional make-up of the pseudogene *POM121L9P*. The numbers under the motifs shown represent the bp positions on the *POM121L9P* sequence. The FAM247 sequence that forms part of *POM121L9P* consists of FAM247 positions 5949-8219, where there are *Alu* sequences at both ends.

In an homologous nearest neighbor gene arrangement present in chimpanzee chr22 that are annotated as glutathione hydrolase light chain 2 gene (*LOC749018*) and putative POM121-like protein 1 gene, *LOC112206778*; these are linked to *GGT5* through the spacer sequence (Figure 9). Thus the human pseudogenes *GGTLC4P* and *POM121L9P* sequences are annotated as protein genes in the homologous chromosomal loci of chimpanzee. This is another example of ncRNA genes in humans annotated as protein genes in non-human primates and is of significance, but isolation of protein products from the chimpanzee genes is essential to confirm this.

Of importance is that 69% of the *POM121L9P* sequence is present in the chimpanzee genome with a 98% identity, at the genomic region where the comparable chromosomal locus resides in chimpanzee. There are no FAM247 or *POM121L9P* sequences that have been found linked to *GGT5* in Rhesus. It appears that the development of the *POM121L9P* sequence began in the chimpanzee, with but a partial sequence.

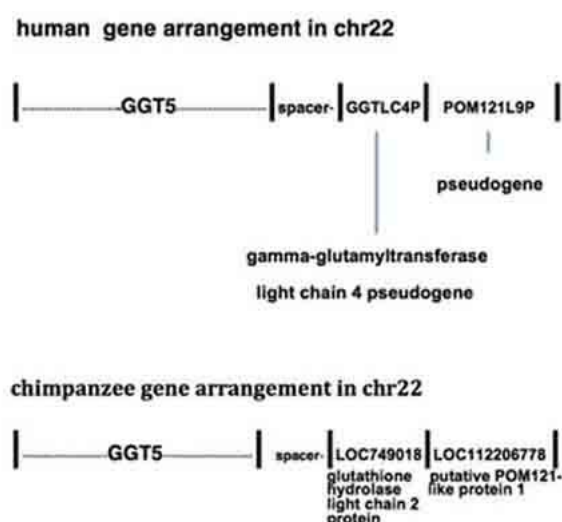


Figure 9. Nearest neighbor gene arrangements in human and chimpanzee chromosomal loci where the *GGT5* gene resides.

4.5. *BCRP3* and *POM121L10P*

Human *BCRP3* and *POM121L10P* are linked to *GGT1* in the gene/sequence arrangement, *GGT1*-spacer-*BCRP3*-*POM121L10P*, which is present in chr22 LCR22H. *FAM247* forms part of the two pseudogenes: *BCRP3*, which has the *FAM247* positions 33-5958 and *POM121L10*, positions 5957- 8219 (Figure 2). Thus parts of the 5' and 3' regions of *FAM247* are found in genes found on these linked genes, which is similar to the presence of *FAM247* in genes *GGT5* and *POM121L9P*.

BCRP3 is a member of the *BCRP* pseudogene family consisting of eight pseudogenes, all of which contain the homologous sequence of the 3' end sequence of *BCR* protein gene. However, *BCRP3* differs as it contains additional sequence motifs (Figure 10) and is the only *BCRP* family member that contains the *FAM247* sequence. The *BCRP3* gene appears to have a unique sequence. The compositional make-up of *BCRP3* shows that its 5' side has the *FAM247* sequence, which is followed by a 4255 bp segment of the immunoglobulin lambda locus (*IGL*) and the 3' end of the *BCR* sequence (Figure 10). The *IGL* sequence is homologous to the *IGL* locus V segments and three C segments, which are known to not encode immunoglobulin proteins. The *IGL* sequence has an *Alu* sequence at the junction with *FAM247*, which may relate to the process of attachment of *IGL* to *FAM247*. In terms of RNA expression, the pseudogene shows broad expression of linear RNA in 27 normal somatic tissues and a broad expression of circular RNA in developing fetal tissues [13, 14].

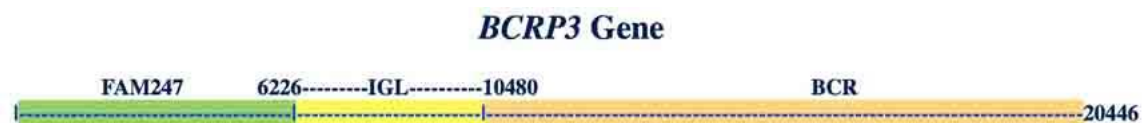


Figure 10. Sequence motif of the *BCRP3* gene. Positions 6226 -10480 of *BCRP3* span the *IGL* insert. The total length of *BCRP3* is 20446 bp.

The *POM121L10P* sequence is linked to *BCRP3* on chr 22. It also contains the *FAM247* sequence (Figure 2). *POM121L10P* is compositionally made up of nearly the entire sequence of the related pseudogene *POM121L1P*, but has a 1062 bp sequence at its 3' end that consists of a copy of the 3' end of the *BCR* gene. *POM121L10P* also appears to be a unique gene construct. The *POM121L10P* linear RNA transcript is strongly expressed in testes; circular RNAs are broadly expressed fetal tissues. [13, 14]. Thus both this gene and *BCRP3* show a robust RNA expression. It should be pointed out that there are additional *POM121LP* pseudogene family members that carry the *FAM247* sequence but are not addressed here.

In the rhesus genome, *GGT1* is linked to the spacer sequence and followed by the *FAM247* sequence, which is similar to that of the human *GGT1* gene/sequence arrangement (Figure 1B). Rhesus gene *LOC107000612*, annotated as a "breakpoint cluster region protein-like protein" is situated close to *GGT1*. This is part of the homologous chromosomal region where the pseudogene *BCRP3* resides in the human genome; approximately 78% of the human *BCRP3* sequence is present in the Rhesus genome at this locus. The *BCRP3* sequence has not been detected in the early primate *Philippine tarsier*. The earliest appearance of the *BCRP3* sequence is in the Rhesus species and the sequence appeared to have matured into a

pseudogene in humans. There was a large chromosomal expansion of the rhesus monkey genome between genes *GGT1* and *GGT5*. The chromosomal length between genes *GGT1* and *GGT5* in *Philippine tarsier* is 2872 bp; in the rhesus monkey it is 216,200 bp. Thus there is a 75-fold sequence expansion between *GGT1* and *GGT5* in rhesus. Segments of the *BCRP3* gene may have formed with this chromosomal expansion. This may account for the source of the *BCRP3* sequence in Rhesus, but the sequence is not found in *Philippine tarsier*.

5. Conclusions

Both the FAM247 lincRNA gene family and the various pseudogenes appear to have the repeat FAM247 sequence as a foundation for gene development, however, the mechanism of formation and the compositional make-up between the lincRNA genes and pseudogenes greatly differs. The FAM247 family (as well as the FAM230 lincRNA gene family) formed by gene duplication and family members display sequences that are “variations on a theme”. Although pseudogenes *BCRP3*, *POM121L9P* and *POM121L10P* contain duplications of part of or entire portions of parent protein genes, they formed differently by a process of addition of large unrelated genomic sequences to the FAM247 sequence, with the resultant formation of unique pseudogene sequences. *Alu* elements are present in FAM247 at sites of attachment; these may contribute to the process of sequence addition. As these pseudogenes are unique with large sequences unrelated to the parent protein gene, the question is whether they should be called pseudogenes. How *USP18* and *GGT5* protein genes developed is not known but a putative ancient FAM247 sequence was likely involved. A separate but important aspect of the FAM247 sequence in cellular and molecular functions is that it contributes the amino acid sequence for protein exons, the first exon of *GGT5* and last exon of *USP18*. The functions of the carboxyl terminal aa sequence of *USP18* are of major significance because of the role in the regulation of the immune system.

Funding: This research received no external funding

Conflicts of Interest: The author declares no conflict of interest.

References

1. Ohno, S. Gene duplication and the uniqueness of vertebrate genomes circa 1970-1999. *Semin Cell Dev Biol.* **1999**, *10*, 517-522.
2. Jacob, F. Evolution and tinkering. *Science* **1977**, *196*, 1161-1166
3. Wang, W.; Yu, H.; Long, M. Duplication-degeneration as a mechanism of gene fission and the origin of new genes in *Drosophila* species. *Nat Genet.* **2004**, *36*, 523-527.
4. Carvunis, A.R.; Rolland, T.; Wapinski, I.; Calderwood, M.A.; Yildirim, M.A.; Simonis, N.; Charleatoux, B.; Hidalgo, C.A.; Barbette, J.; Santhanam, B.; et al. Proto-genes and de novo gene birth. *Nature* **2012**, *487*, 370-374.
5. McLysaght, A.; Guerzoni, D. New genes from non-coding sequence: the role of de novo protein-coding genes in eukaryotic evolutionary innovation. *Philos Trans R Soc Lond B Biol Sci.* **2015**, *370*, 20140332.
6. Schlotterer, C. Genes from scratch—the evolutionary fate of de novo genes. *Trends Genet.* **2015**, *31*, :215-219.

7. Van Oss, S.B.; Carvunis, A.R.; De novo gene birth. *PLoS Genet.* **2019**, *15*(5):e1008160.
8. Delihias, N. Formation of human long intergenic non-coding RNA genes, pseudogenes, and protein genes: Ancestral sequences are key players. *PLoS One*, **2020**, *15*(3):e0230236.
9. Nowell, P.; Hungerford, D.; A minute chromosome in human chronic granulocytic leukemia. *Science*, **1960**, *132*, 1497.
10. de Klein, A.; van Kessel, A.G.; Grosveld, G.; Bartram, C.R.; Hagemeijer, A.; Bootsma, D.; Spurr, N.K.; Heisterkamp, N.; Groffen, J.; Stephenson, J.R. A cellular oncogene is translocated to the Philadelphia chromosome in chronic myelocytic leukaemia. *Nature*, **1982**, *300*, 765–767.
11. Arimoto, KI.; Löchte, S.; Stoner, S.A.; Burkart, C.; Zhang, Y.; Miyauchi, S.; Wilmes, S.; Fan, J.B.; Heinisch, J.J.; Li, Z. STAT2 is an essential adaptor in USP18-mediated suppression of type I interferon signaling. *Nat. Struct. Mol. Biol.* **2017**, *24*, 279–289.
12. Honke, N.; Shaabani, N.; Zhang, D.E.; Hardt, C.; Lang, K.S. Multiple functions of USP18. *Cell Death Dis.* **2016**, *7*(11):e2444.
13. Szabo, L.; Morey, R.; Palpant, N.J.; Wang, P.L.; Afari, N.; Jiang, C.; Parast, M.M.; Murry, C.; Laurent, L.C.; Salzman, J. Statistically based splicing detection reveals neural enrichment and tissue-specific induction of circular RNA during human fetal development. *Genome Biol.* **2015**, *16*(1):126.
14. Fagerberg, L.; Hallström, B.M.; Oksvold, P.; Kampf, C.; Djureinovic, D.; Odeberg, J.; Habuka, M.; Tahmasebpour, S.; Danielsson, A.; Edlund, K.; et al. Analysis of the human tissue-specific expression by genome-wide integration of transcriptomics and antibody based proteomics. *Mol. Cell Proteomics* **2014**, *13*, 397–406.
15. O'Leary, N.A.; Wright, M.W.; Brister, J.R.; Ciufu, S.; Haddad, D.; McVeigh, R.; Rajput B.; Robbertse, B.; Smith-White, B.; Ako-Adjei, D.; et al. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **2016**, *44*, D733–D745.
16. Siepel, A. Phylogenomics of primates and their ancestral populations. *Genome Res.* **2009**, *19*, 1929–1941.
17. Madeira, F.; Park, Y.M.; Lee, J.; Buso, N.; Gur, T.; Madhusoodanan, N.; Basutkar, P.; Tivey, A.R.N.; Potter, S.C.; Finn, R.D. The EMBL-EBI Search and Sequence Analysis Tools APIs in **2019**, *Nucleic Acids Res.* **2019**, *47*(W1):W636–W641.
18. Malakhov, M.P.; Malakhova, O.A.; Kim, K.I.; Ritchie, K.J.; Zhang, D.E. Protein ISGylation Modulates the JAK-STAT Signaling Pathway. *J Biol Chem.* **2002**, *277*, 9976–9981.
19. Dauphinee, S.M.; Richer, E.; Eva, M.M.; McIntosh, F.; Paquet, M.; Dangoor, D.; Burkart, C.; Zhang, D.E.; Gruenheid, S.; Gros, P. Contribution of increased ISG15, ISGylation and deregulated type I IFN signaling in Usp18 mutant mice during the course of bacterial infections. *Genes Immun.* **2014**, *15*, 282–292.