### **Untangling the cellular origins of viruses**

Computing the early history of biochemistry and life

#### Gustavo Caetano-Anollés

University of Illinois, Urbana-Champaign, Illinois, USA

**DNA** habitats and its **RNA** inhabitants 3 - 5 July 2014 Salzburg - Austria Viruses, Mobile Genetic Elements, Viroids, Introns, Ribozymes and other RNAgents

## **Retrodictive exploration**

3D STRUCTURE

#### **QUESTIONS:**

What is the role of viruses in cellular evolution? Is there a truly universal tree of life? Are viruses monophyletic or polyphyletic? How are viral groups related to each other? Is there a preferential direction of gene transfer?

Ideographic framework historical and retrodictive Structure provides a window into function and evolution

### The macromolecular world is modular





Modules carry deep evolutionary signal





Evolutionary conservation increases with hierarchical complexity

### **Domains exhibit levels of structural abstraction**





## **Modularity increases in evolution**



**Domain diversity (occurrence)** and **reuse (abundance)** in proteomes carry **strong and deep phylogenetic signal** 



Module diversity and reuse increase in evolution, even if domains are universal Kim and Caetano-Anollés (2012) BMC Evol Biol 12: 13



## **Retrodictive exploration**



## **Structural phylogenomics**

COMPUTING THE HISTORY OF THE PROTEIN WORLD:

Genomic demography and phylogeny reconstruction





## **Structural phylogenomics**

COMPUTING THE HISTORY OF THE PROTEIN WORLD:

**Recent applications** 

#### **Coevolutionary history of the ribosome**

Harish and Caetano-Anollés (2012) *PLoS ONE* 7: e32776 Caetano-Anollés and Sun (2014) *Front Genet* 5: 127.



#### Origin of the genetic code in protein flexibility

Caetano-Anollés et al. (2013) *PLoS ONE* 8: e72225 Debes et al. (2013) *PLoS Comput Biol* 9: e1002861



## **Origin of viruses**



Caetano-Anollés 2014©

# Three main lines of thought about three main events of emergence

Nasir et al. (2012) Mobile Genet Elements 12:156

**V** : Origin of viruses

**C** : Origin of cells

L : Origin of diversified cells\*



Last universal cellular ancestor [LUCELLA]

## The root of the 'universal' tree of life



Trees of proteomes and evolutionary PCA analyses derived from fold superfamilies (FSFs) place the root of the tree in VIRUSES with large-to-medium size genomes Nasir et al. (2012) BMC Evol Biol 12:156



## The root of the 'universal' tree of life



#### The ancient cellular origin of viruses supports the 'reductive' evolution hypothesis and redefines urancestors of life

Nasir et al. (2012) Mobile Genet Elements 12:156





LUCA: Last universal common ancestor **LUCELLA:** Last universal cellular ancestor

### The root of the 'universal' tree of life







Two organic-walled microfossil size ranges: **5-90 μm** (Pilbara Craton, Western Australia) **50-300 μm** (Moodies Group, South Africa)

Sugitani et al. (2007) *Precamb Res* 158: 228-226 Sugitani et al. (2009) *Astrobiology* 9:603-615 Javaux et al. (2010) *Nature* 463:934-938 Wacey et al. (2011) *Nature Geosci* 4:698-702

# Expanded proteome dataset

2,715 proteomes from all 6 viral groups [NCBI RefSeq Viral Resource]

- $_{\circ}$  1,125 dsDNA
- $\circ$  453 ssRNA
- $_{\circ}$  122 dsRNA
- 806 ssRNA(+)
- 95 ssRNA(-)
- 114 Retrotranscribing viruses
- **1,496 proteomes from cellular organisms** [SUPERFAMILY database]
  - 114 Archaea, 1,062 Bacteria, and 320 Eukarya



- Includes protein domains with low sequence identity (< 15%)</p>
- > Common conserved 3D cores and biochemical properties
- > Homologous as defined by the SCOP evolutionary theory
- More conserved and better suited to study remote relationships



## **Sharing patterns**





4,211 proteomes from cells and viruses 1,993 FSF domains (E < 0.0001)

### **Spread of viral FSFs in cellular proteomes**

#### f index \*

| FSF Description                                     | Molecular Function      | Detailed Function             | Distribution in Viruses                             | Archaea (/) | Bacteria (/) | Eukarya (/) |
|---|-------------------------|-------------------------------|---|-------------|--------------|-------------|
| Inhibitor of apoptosis (IAP) repeat                 | Intracellular processes | Cell cycle, Apoptasis         | dsDNA, ssRNA(+)                                     | 0.00        | 0.00         | 0.73        |
| ATP-dependent DNA ligase DNA-binding domain         | Regulation              | DNA-binding                   | dsDNA   | 0.98        | 0.12         | 0.99        |
| Ribonuclease H-like                                 | Metabolism              | Nucleotide m/tr               | dsDNA, ssDNA, ssRNA[-], retrotranscribing           | 1.00        | 1.00         | 1.00        |
| DNA/RNA polymerases                                 | Information             | <b>DNA</b> replication/repair | dsDNA, dsRNA, ssRNA(+), ssRNA(-), retrotranscribing | 1.00        | 0.98         | 1.00        |
| Translation initiation factor 2 beta, alF2beta, N-t | Information             | Translation                   | diDNA   | 1.00        | 0.00         | 1.00        |
| RING/U-box  | Information             | <b>ONA</b> replication/repair | dsDNA, ssRNA(-)                                     | 0.04        | 0.11         | 1.00        |
| dUTPase-like  | Metabolism              | Nucleotide m/tr               | dsDNA, retrotranscribing                            | 0.98        | 0.89         | 0.91        |
| Nucleic acid-binding proteins                       | Information             | <b>DNA</b> replication/repair | dsDNA, ssDNA  | 1.00        | 1.00         | 1.00        |
| R1 subunit of ribonucleotide reductase, N-termin    | r Metabolism            | Nucleotide m/tr               | dsDNA   | 0.61        | 0.78         | 0.96        |
| DNA ligase/mRNA capping enzyme, catalytic don       | r Information           | <b>DNA</b> replication/repair | dsDNA   | 1.00        | 0.99         | 1.00        |
| Zinc-binding domain of translation initiation fact  | General                 | Ion binding                   | dsDNA   | 1.00        | 0.00         | 0.99        |
| P-loop containing nucleoside triphosphate hydro     | General                 | Small molecule binding        | (dsDNA, ssDNA, dsRNA, ssRNA(+)                      | 1.00        | 1.00         | 1.00        |
| Ferritin-like                                       | Intracellular processes | ion m/tr                      | dsDNA   | 0.99        | 0.99         | 1.00        |
| PFL-like glycyl radical enzymes                     | Metabolism              | Other enzymes                 | dsDNA   | 1.00        | 0.98         | 0.98        |
| Protein kinase-like (PK-like)                       | Regulation              | Kinases/phosphatases          | dsDNA, dsRNA, netrotranscribing                     | 1.00        | 0.96         | 1.00        |
| Cytidine deaminase-like                             | Metabolism              | Other enzymes                 | diDNA   | 0.96        | 0.99         | 1.00        |
| Ribosomal protein SS domain 2-like                  | Information             | Translation                   | dsDNA   | 1.00        | 1.00         | 1.00        |
| Cryptochrome/photolyase FAD-binding domain          | General                 | General                       | dsDNA   | 0.32        | 0.53         | 0.83        |
| FAD-linked reductases, C-terminal domain            | Metabolism              | Redox                         | dsDNA   | 0.61        | 0.85         | 1.00        |
| ADP-ribosylation                                    | Metabolism              | Secondary metabolism          | dsDNA, ssDNA.                                       | 0.47        | 0.33         | 0.98        |
| RPB5-like RNA polymerase subunit                    | Information             | Transcription                 | dsDNA   | 0.98        | 0.00         | 0.93        |
| TNF-like  | Extracellular processes | Immune response               | dsDNA   | 0.00        | 0.05         | 0.35        |
| GroES-like  | Intracellular processes | Protein modification          | diDNA   | 0.80        | 0.98         | 1.00        |
| Translation proteins                                | Information             | Translation                   | diDNA   | 1.00        | 1.00         | 1.00        |

\*Percentage of proteomes encoding an FSF divided by the total number of proteomes (in a relative 0-1 scale)





### **Biases in functional** preferences of viral FSFs

Extracellular processes underrepresented in Archaea Intracellular processes and general metabolism overrepresented in Eukarya



### Viruses enhance molecular biodiversity

#### **FSFs shared with viruses are widely distributed in the world of proteomes, especially with ssDNA, dsRNA and ssRNA (-) viruses**



## **67 virus-specific FSFs**

- Three times more abundant than Archaea-specific FSFs
- Under-represented (HGT and sampling biases)
- Harbor pathogenic and immunological roles
- Include major viral capsid proteins
- Interesting drug targets
- Present in all six viral subgroups
- Could not have originated in cells!









### **Virus-host relationships and FSF sharing**

## While viruses do not cross superkingdom barriers and have specific host preferences, they harbor a common structural core



Loss of viral lineages in Archaea and Bacteria?
Late appearance of RNA viruses?

Caetano-Anollés 2014©

#### 34 FSFs common to viruses associated with the three superkingdoms

| 1   | FSF Description                             | Molecular Function      | Detailed Function             | Distribution  | A(f) \$( | n I        | 00   |
|-----|---|-------------------------|-------------------------------|---|----------|------------|------|
| 2   | Ribonuclease H-like                         | Metabolism              | Nucleotide m/tr               | dsONA, ssDNA, ssRNA(-), retrotranscribing           | 1.00 1.0 | 00 1       | 1.00 |
| 3   | ONA/RNA polymerases                         | information             | DNA replication/repair        | dsONA, dsRNA, ssRNA(+), ssRNA(-), retrotranscribing | 1.00 0.5 | 98 3       | 1.00 |
| .4  | dUTPase-like                                | Metabolism              | Nucleotide m/tr               | dsDNA, retrotranscribing                            | 0.98 0.1 | 89 0       | 0.91 |
| 5   | Nucleic acid-binding proteins               | Information             | <b>DNA replication/repair</b> | diONA, siDNA  | 1.00 1.0 | 00 3       | 1.00 |
| 6   | P-loop containing nucleoside triphosphate h | General                 | Small molecule binding        | dsDNA, ssDNA, dsRNA, ssRNA(+)                       | 1.00 1.0 | 00 1       | 1.00 |
| 7   | (Phosphotyrosine protein) phosphatases il   | Regulation              | Kinases/phosphatases          | diONA   | 0.51 0.5 | 51 3       | 1.00 |
| - 8 | PIN domain-like                             | Other                   | Unknown function              | diONA   | 1.00 0.5 | 99 7       | 1.00 |
| 9   | N-terminal nucleophile aminohydrolases (Nt  | Metabolism              | Other enzymes                 | acha All viral around                               | 1.0      |            | 00   |
| 10  | Chaperone J-domain                          | Intracellular processes | Protein modification          | eona All viral groups                               | 0.4      |            | 20   |
| 11  | Putative DNA-binding domain                 | Regulation              | DNA-binding                   | diONA   | 1.0      | <b>Л</b>   | 99   |
| 12  | Winged helix DNA binding domain             | Regulation              | ONA-binding                   | dsONA   | 1.0      |            | 00   |
| 13  | Uracil-ONA glycosylase-like                 | Information             | <b>DNA replication/repair</b> | diONA   | 1.0      | 5          | 99   |
| 14  | Nucleotide-diphospho-sugar transferases     | Metabolism              | Transferases                  | dsONA, dsRNA  | 1.0      | ~          | 90   |
| 15  | vWA-like                                    | Extracellular processes | Cell adhesion                 | diONA   | 0.9      | 0          | 00   |
| 16  | 5-adenosyl-L-methionine-dependent methylt   | Metabolism              | Transferases                  | dsONA, dsRNA, ssRNA(+)                              | 1.0      |            | 90   |
| 17  | Thioredoxin-like                            | Metabolism              | Redox                         | diONA   | 0.9      |            | 00   |
| 1.8 | ONA clamp                                   | information             | <b>DNA</b> replication/repair | (BONA   | 1.0      |            | 00   |
| 19  | Metallo-dependent phosphatases              | intracellular processes | Proteases                     | dsDNA viruses:                                      | 1.0      |            | 00   |
| 20  | ONA breaking-rejoining enzymes              | information             | <b>DNA</b> replication/repair | dsDNA   | 0.8      | $\bigcirc$ | 99   |
| 21  | ATPase domain of HSP90 chaperone/DNA to     | Intracellular processes | Protein modification          | elona Common  | 1.0      |            | 00   |
| 22  | Rad51 N-terminal domain-like                | Information             | DNA replication/repair        | (KONA   | 1.0      |            | 83   |
| 23  | NAD(P)-binding Rossmann-fold domains        | General                 | Small molecule binding        | denominator!  | 1.0      | Y          | 00   |
| 24  | Radical SAM enzymes                         | 1                       |                               | dsONA   | 1.0      |            | 00   |
| 25  | Restriction endonuclease-like               | Many FS                 | SFs are                       | diONA   | 1.0      | 0          | 00   |
| 26  | TPR-like                                    |                         |                               | diONA   | 0.8      |            | 00   |
| 27  | UDP-Glycosyltransferase/glycogen phosphor   | sunerfol                | ds with                       | dsONA, dsRNA  | 1.0      |            | 00   |
| 28  | Thymidylate synthase-complementing protein  | Jupenio                 |                               | diONA   | 0.5      |            | 25   |
| 29  | Concanavalin A-like lectins/glucanases      | contral f               | iunctions                     | diONA, diRNA  | 0.4      | D          | 00   |
| 30  | ARM repeat                                  | Central I               | unctions                      | dsONA, ssRNA(+)                                     | 0.8      |            | 00   |
| 31  | beta-beta-alpha zinc fingers                | Regulation              | DNA-binding                   | diONA   | 0.3      |            | 00   |
| 32  | Bacterial hemolysins                        | Metabolism              | Other enzymes                 | dsONA, ssDNA, dsRNA                                 | 0.3      |            | 43   |
| 33  | Trypsin-like serine proteases               | Intracellular processes | Proteases                     | dsONA, ssRNA(+)                                     | 0.5      |            | 97   |
| 34  | Fibronectin type III                        | Extracellular processes | Cell adhesion                 | dsONA   | 0.25 0.5 | 51 6       | 0.99 |
| 35  | Pectin lyase-like                           | Metabolism              | Polysaccharide m/tr           | dsONA, ssRNA(+)                                     | 0.68 0.1 | 72:0       | 0.89 |

### A possible early origin of dsDNA viruses

dsDNA viruses can be as large as cells



### dsDNA viruses harbor the largest number and the most shared FSFs

| Subgroup     | Total | Unique | ssDNA | retroviruses | dsRNA | ssRNA |
|--------------|-------|--------|-------|--------------|-------|-------|
| dsDNA        | 587   | 530    | 7     | 11           | 9     | 22    |
| ssDNA        | 14    | 4      | -     | 0            | 0     | 0     |
| retroviruses | 27    | 12     | -     | -            | 0     | 1     |
| dsRNA        | 28    | 11     | -     | -            | -     | 1     |
| ssRNA        | 68    | 38     | -     | -            | -     | -     |

Nasir et al. (2014) Frontiers Microbiology 5:194



Rooted tree of proteomes built from an FSF domain census Kim and Caetano-Anollés (2011) *BMC Evol Biol* 11: 140.



Venn taxonomic groups

### History of the viral and cellular FSF repertoire

Superkingdom-specific FSFs shared with viruses (AV, BV, and EV) appear after those that are not shared





### **Extreme reductive tendencies in viruses**



Bacteria



Archaea







## **Diversification of life**





### **Some conclusions**



#### **Early cellular origins of viruses**

- Viruses originated from proto-cells that coexisted with cellular ancestors ~3.4 Gy ago
- Early origin of virocell lineages by reductive evolution
- Ancient virocells had DNA replicons
- DNA viruses are ancient and monophyletic



#### Late rise of viral novelties

- Late appearance of capsids and parasitic life cycles
- Viruses enhance cellular molecular diversity
- Early origins of archaeal DNA viruses and spread of novelties to RNA viruses in other superkingdoms

## Acknowledgements



Arshan Nasir

Collaborators:









PARIS

Questions

NSF

Kyung Mo Kim (Daejeon) Patrick Forterre (Paris)

## **Thank you!**

