# On the importance of curated databases
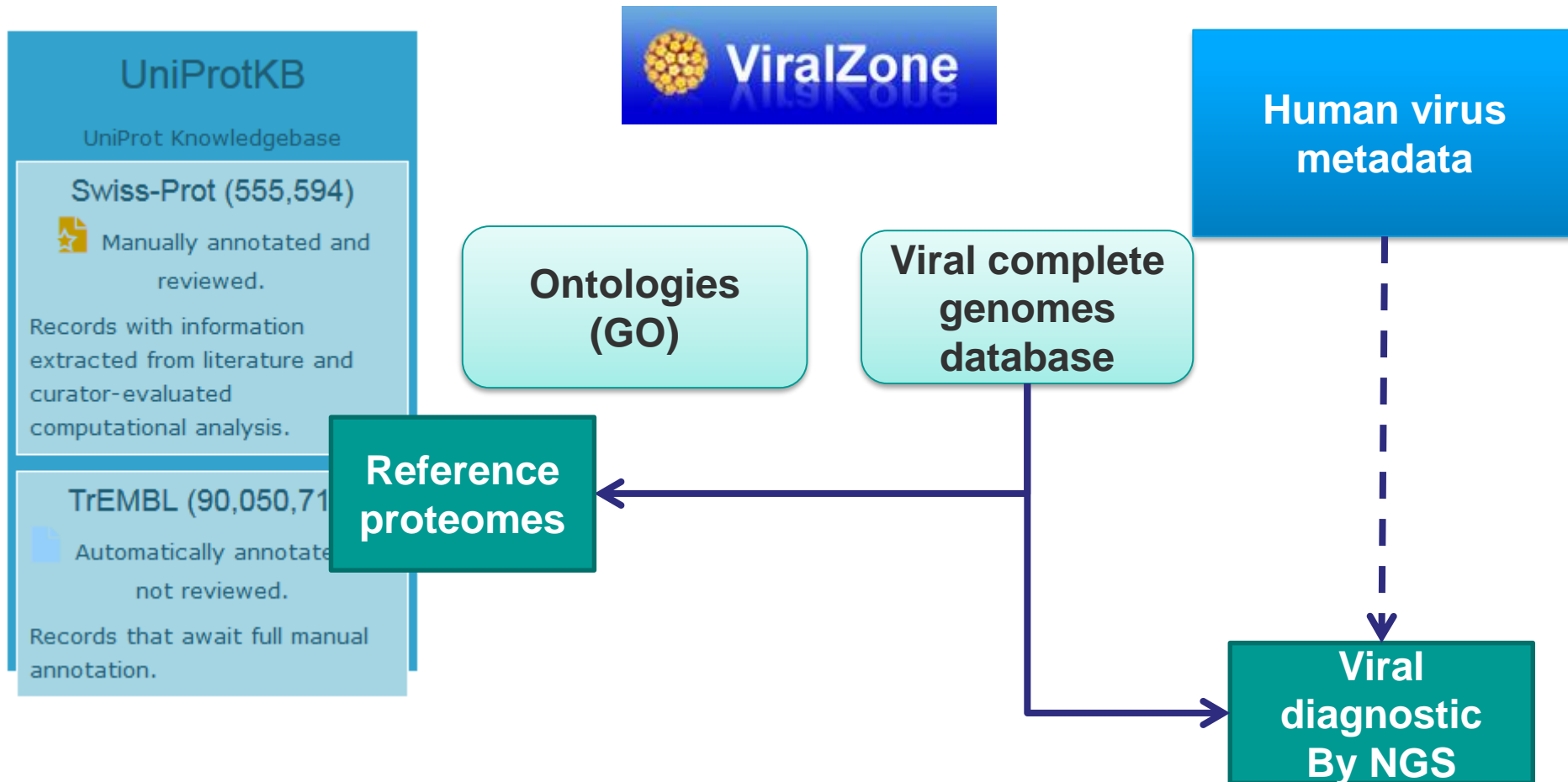
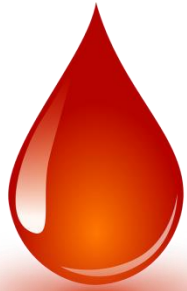**SIB Swiss Institute of Bioinformatics**

Geneva, Switzerland

SIB
Swiss Institute of
Bioinformatics

# Digitalisation of virus knowledge
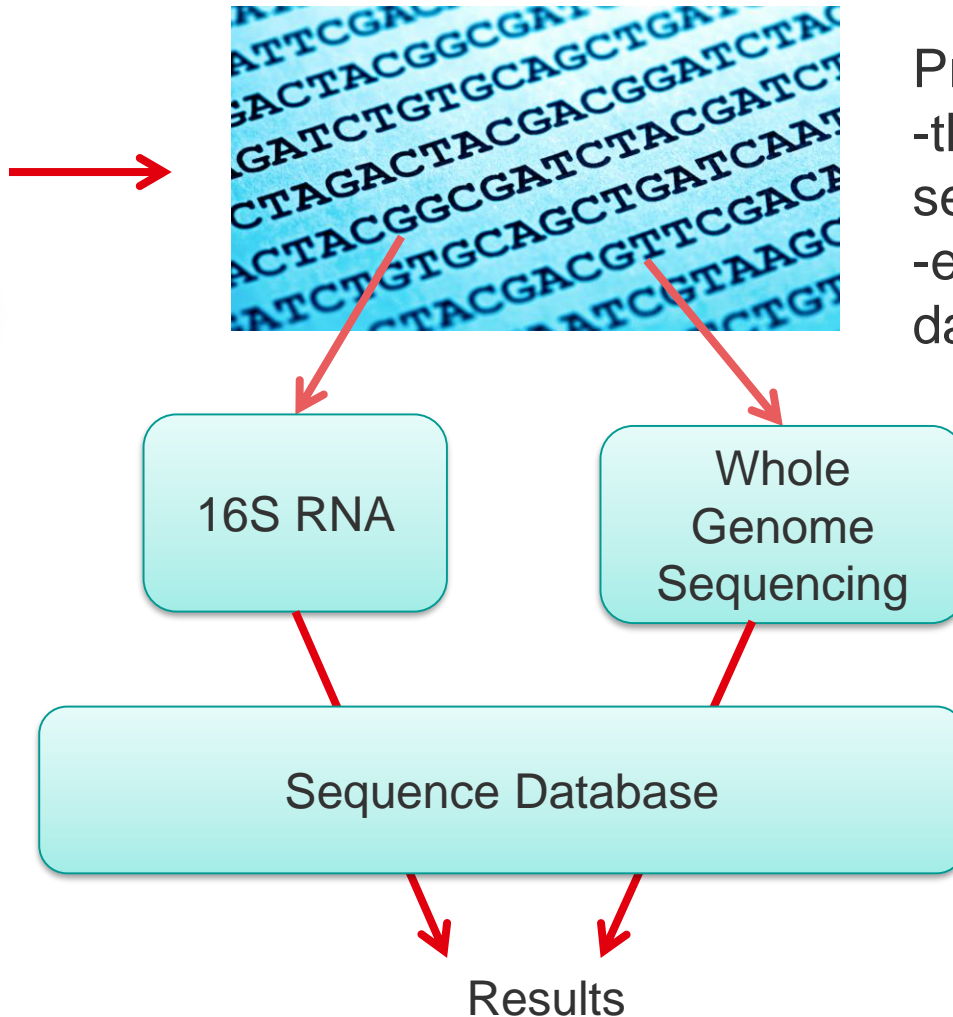## SIB, SwissProt group

# Detection by next generation sequencing

Patient
sample

Present challenge:
-the analysis step (speed,
sensivity)
-efficient references
databases.

16S RNA

Whole
Genome
Sequencing

Sequence Database

Results

# GenBank is a repository database



- Unexpected information you can find in
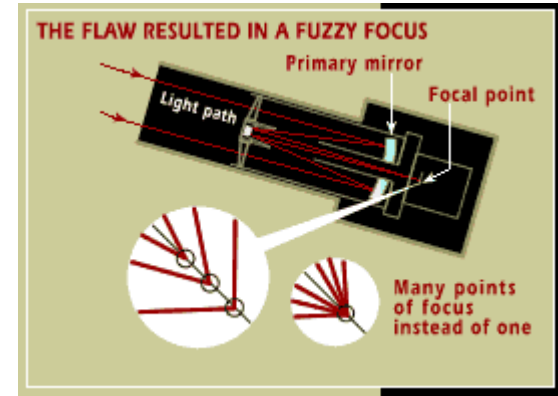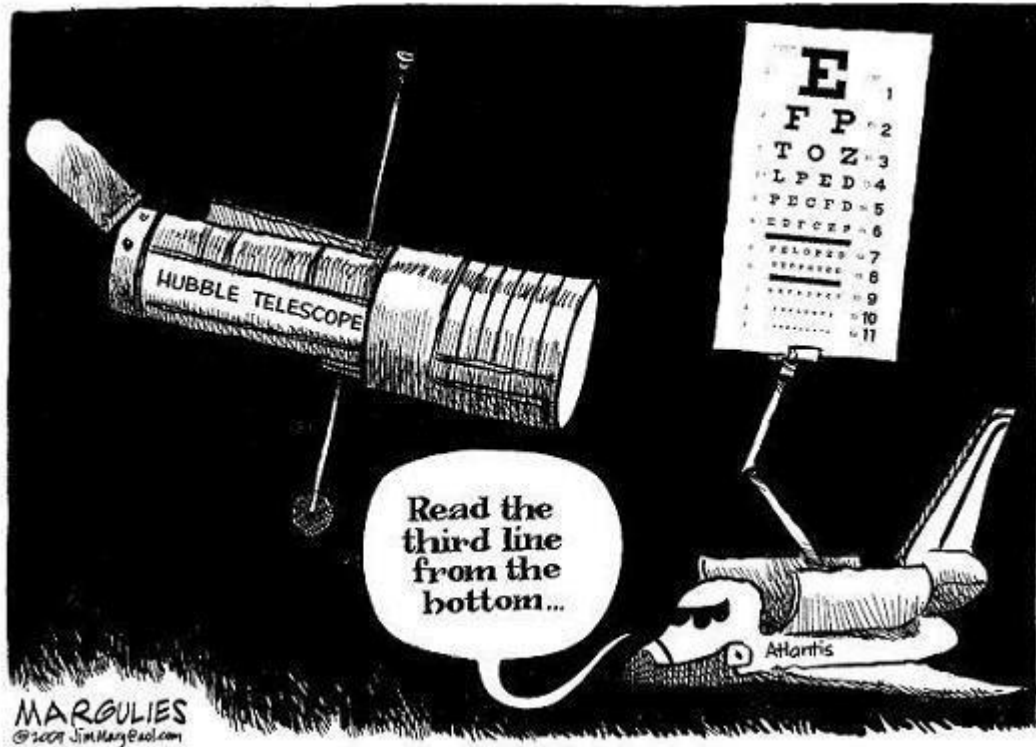*INSDC*=Genbank/EMBL/DDBJ:

```
FT   source          1..124
FT                   /db_xref="taxon:4097"
FT                   /organelle="plastid:chloroplast"
FT                   /organism="Nicotiana tabacum"
FT                   /isolate="Cuban cahibo cigar, gift from
FT                                        President Fidel Castro"
```

```
FT   CDS          complement(45959..47332)
FT               /db_xref="SPTREMBL:Q9UZ71"
FT               /note="PAB2386"
FT               /transl_table=11
FT               /product="4-AMINOBUTYRATE qui se dilate AMINOTRANSFERASE
FT               (EC 2.6.1.19)"
FT               /protein_id="CAB50188.1"
FT               /translation="MDYPRIVVNPPGPKAKELIEREKRVLSTGIGVKLFPLVPKRGFGP
FT               FIEDVDGNVFIDFLAGAAAASTGYSHPKLVKAVKEQVELIQHSMIGYTHSERAIRVAEK
FT               LVKISPIKNSKVLFGLSGSDAVDMAIKVSKFSTRRPWILAFIGAYHGQTLGATSVASFQ
FT               VSQKRGYSPLMPNVFWVPYPNPYRNPWGINGYEEPQELVNRVVEYLEDYVFSHVVPPDE
FT               VAAFFAEPIQGDAGIVVPPENFFKELKKLLDEHGILLVMDEVQTGIGRTGKWFASEWFE
FT               VKPDMIIFGKGVASGMGLSGVIGREDIMDITSGSALLTPAANPVISAAADATLEIIEEE
FT               NLLKNAIEVGSFIMKRLNELKEQFDIIGDVRGKGLMIGVEIVKENGRPDPEMTGKICWR
FT               AFELGLILPSYGMFGNVIRITPPLVLTKEVAEKGLEIIEKAIKDAIAGKVERKVVTWH"
```

# The importance of database



Whatever the quality of the samples, sequencing and bioinformatics,
a flawed database can blur the results

# There was no database for virus complete genomes



Virus??

# Viral genomes

# Avoid RefSeq to assess full length criteria

```
LOCUS       NC_028891               22947 bp    DNA     linear   VRL 05-JAN-2016
DEFINITION  Hawaiian green turtle herpesvirus thymidine kinase (UL23),
            membrane-associated protein (UL24), minor capsid protein (UL25),
            capsid maturation protease (UL26), virion scaffolding protein
            (UL26.5), virion membrane glycoprotein B (gB), DNA
            cleavage/packaging protein (UL28), single-stranded DNA-binding
            protein (UL29), DNA polymerase catalytic subunit (pol), nuclear
            phosphoprotein (UL31), DNA cleavage/packaging (UL32), DNA
            cleavage/packaging protein (UL33), membrane-associated
            phosphoprotein (UL34), and basic phosphorylated capsid protein
            (UL35) genes, complete cds; and very large tegument protein (UL36)
            gene, partial cds.
ACCESSION   NC_028891
VERSION     NC_028891.1
DBLINK      BioProject: PRJNA307765
KEYWORDS    RefSeq.
SOURCE      Hawaiian green turtle herpesvirus
  ORGANISM  Hawaiian green turtle herpesvirus

COMMENT     PROVISIONAL REFSEQ: This record has not yet been subject to final
            NCBI review. The reference sequence is identical to AF035003.
            COMPLETENESS: full length.
```

?

# Refseq provisional

| Mogiana tick virus | | | |
|---|---|---|---|
| Mogiana tick virus | (2963 nt) | NC_034222 | proteins: 1 |
| Mogiana tick virus | (2728 nt) | NC_034224 | proteins: 1 |
| Mogiana tick virus | (2629 nt) | NC_034225 | proteins: 1 |
| Mogiana tick virus | (2705 nt) | NC_034223 | proteins: 1 |

```
COMMENT     PROVISIONAL REFSEQ: This record has not yet been subject to final
            NCBI review. The reference sequence is identical to AF035003.
            COMPLETENESS: full length.
```

**Characterisation of divergent flavivirus NS3 and NS5 protein**
**detected in Rhipicephalus microplus ticks from Brazil.**

Maruyama SR[1], Castro-Jorge LA[1], Ribeiro JM[1], Gardinassi LG[2], Garcia GR[1], Brandão LG[1], Rodrigu
EP[1], Ferreira BR[5], Fonseca BA[1], Miranda-Santos IK[1].

**VP2-3 gene, complete cds,**
**..**



Flavivirus polyprotein

Signal peptidase ↓  Golgi protease ↓  NS3 protease ▽

# Refseq is for annotation



Reference
Puerto rico 1934

# Representing the diversity calls for many references

# Toward detecting all viruses

- Virology have focused on pathogens

- Recent genomic explorations of human samples have revealed dozens of previously unrecognized viruses.

# ViralZone complete genome dataset

| Influenza databases | HIV database | HPV database | Adenovirus database | HCV database |
|---|---|---|---|---|

**ViPR**
Virus Pathogen Resource

| HBV database |
|---|

GenBank

Extraction and
Curation by families

**For Eukaryotic viruses: 70,352 complete genomes**
**317,979 sequences**

# Toward an automatic detection of virus complete genomes

Indidual criteria
for each 121 virus families

GenBank

Virus
Complete
genomes

Testing negative selection by families;
-genome length
-CDS all complete
-Number of CDS



100nm

Variola virus
360nm

Herpesvirus
200nm

Rabies
180X80nm

Measles
150nm

HIV-1
120nm

SARS
120nm

Influenza virus
100nm

Adenovirus
90nm

Rotavirus
80nm

Ebolavirus
80x970nm

Papillomavirus
60nm

Dengue virus,
Zika virus
50nm

Hepatitis C virus
50nm

Hepatitis B virus
42nm

Hepatitis A virus,
Poliovirus
30nm

Parvovirus
20nm

# Trimming the branches

223 complete sequences

54 clusters at 95% identity

# Manual vs clustering regarding virus variability

# Adding annotation to reference sequences



GenBank TaxID

Species

genotype

Isolate name

GenBank accession

Segment

**Virus sequence**

# Genotyping within a virus species



Global distribution of HCV genotypes

World Health Organisation 2009

key
- Genotype 1
- Genotype 2
- Genotype 3
- Genotype 4
- Genotype 5
- Genotype 6

# Genotype references data

We are gathering data from various sources: WHO, CDC, publications, Book, virus databases…

Status: in construction

Already done: HAV, HBV, HCV, HSV-1, HEV, HRV, HPV, measles virus, pegivirus, rotavirus, rubella virus, TTV, VZV, West Nile virus, Yellow fever virus, Zika virus

Example: WHO measles

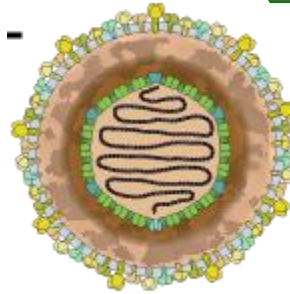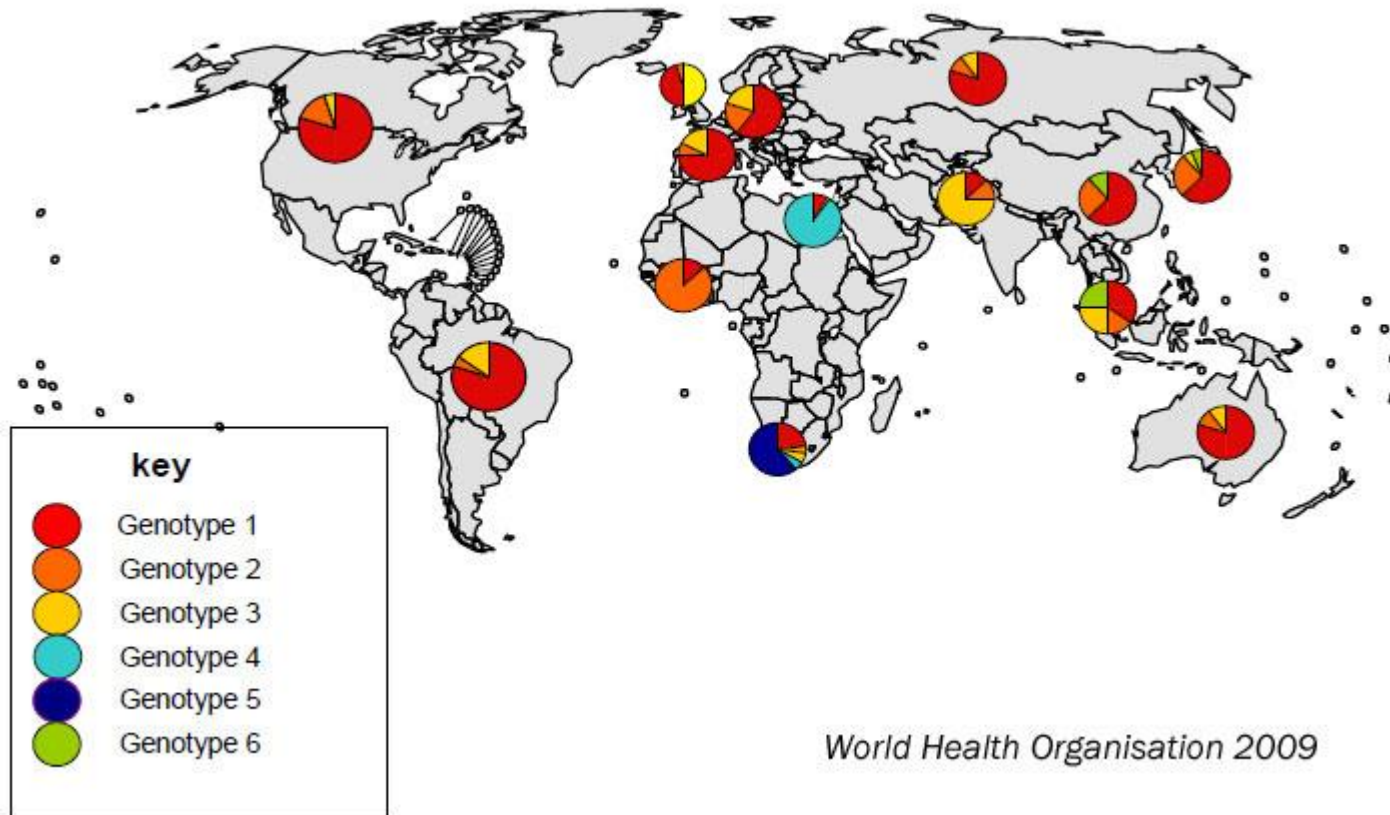| Genotype – Génotype | Last observed* – Dernière observation* | Reference strain – Souche de référence | GenBank H | sp Genbank N |
|---|---|---|---|---|
| A | 2008 | MVi/Maryland.USA/0.54 | U03669 | U01987 |
| B1ᵃ | 1983 | MVi/Yaounde.CMR/12.83 | AF079552 | U01998 |
| B2 | 2011 | MVi/Libreville.GAB/0.84 | L46753 | U01994 |
| B3 | 2011 | MVi/New York.USA/0.94 | L46752 | L46753 |
| | | MVi/Ibadan.NGA/0.97/1 | AJ239133 | AJ232203 |
| C1ᵃ | 1992 | MVi/Tokyo.JPN/0.84 | AY047365 | AY043459 |
| C2 | 2007 | MVi/Maryland.USA/0.77 | M81898 | M89921 |
| | | MVi/Erlangen.DEU/0.90 | Z80808 | X84872 |
| D1ᵃ | 1986 | MVi/Bristol.GBR/0.74 | Z80805 | D01005 |
| D2 | 2005 | MVi/Johannesburg.ZAF/0.88/1 | AF085498 | U64582 |
| D3 | 2004 | MVi/Illinois.USA/0.89/1 | M81895 | U01977 |
| D4 | 2011 | MVi/Montreal.CAN/0.89 | AF079554 | U01976 |
| D5 | 2010 | MVi/Palau/0.93 | L46757 | L46758 |
| | | MVi/Bangkok.THA/0.93/1 | AF009575 | AF07955 |
| D6 | 2007 | MVi/New Jersey.USA/0.94/1 | L46749 | L46750 |
| D7 | 2007 | MVi/Victoria.AUS/16.85 | AF247202 | AF243450 |
| | | MVi/Illinois.USA/50.99 | AY043461 | AY037020 |
| D8 | 2011 | MVi/Manchester.GBR/30.94 | U29285 | AF280803 |
| D9 | 2011 | MVi/Victoria.AUS/12.99 | AY127853 | AF481485 |
| D10 | 2005 | MVi/Kampala.UGA/51.01/1 | AY923213 | AY923185 |
| D11 | 2011 | MVi/Menglian.Yunnan.CHN/47.09 | GU440576 | GU440571 |
| Eᵃ | 1987 | MVi/Goettingen.DEU/0.71 | Z80797 | X84879 |
| Fᵃ | 1994 | MVs/Madrid.ESP/0.94 (SSPE) | Z80830 | X84865 |
| G1ᵃ | 1983 | MVi/Berkeley.USA/0.83 | AF079553 | U01974 |
| G2 | 2004 | MVi/Amsterdam.NLD/49.97 | AF171231 | AF171232 |
| G3 | 2011 | MVi/Gresik.IDN/17.02 | AY184218 | AY184217 |
| H1 | 2011 | MVi/Hunan.CHN/0.93/7 | AF045201 | AF045212 |
| H2 | 2003 | MVi/Beijing.CHN/0.94/1 | AF045203 | AF045217 |

# Adding annotation to reference sequences

## Virus sequence

>LC190490; **species=**Rotavirus A; **taxid=**28875; **genotype=**unknown; **segment=**1; **isolate=**Isolate LC190490;

AGTTGTTGATCTGTGTGAATCAGACTGCGACAGTTCGAGTTTGAAGCGAAAGCTAGCAACAGTATCAACA
GGTTTTATTTTGGATTTGGAAACGAGAGTTTCTGGTCATGAAAAACCCAAAAAAGAAATCCGGAGGATTC
CGGATTGTCAATATGCTAAAACGCGGAGTAGCCCGTGTGAGCCCCTTTGGGGGCTTGAAGAGGCTGCCAG
CCGGACTTCTGCTGGGTCATGGGCCCATCAGGATGGTCTTGGCAATTCTAGCCTTTTTGAGATTCACGGC
AATCAAGCCATCACTGGGTCTCATCAATAGATGGGGTTCAGTGGGGAAAAAGAGGCTATGGAAATAATA
AAGAAGTTCAAGAAAGATCTGGCTGCCATGCTGAGAATAATCAATGCTAGGAAGGAGAAGAAGAGACGAG
GCGCAGATACTAGTGTCGGAATTGTTGGCCTCCTGCTGACCACAGCTATGGCAGCGGAGGTCACTAGACG
TGGGAGTGCATACTATATGTACTTGGACAGAAACGATGCTGGGGAGGCCATATCTTTTCCAACCACATTG
GGGATGAATAAGTGTTATA

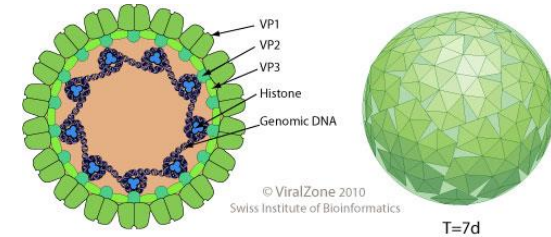# Reads can be assigned to a species level

# In search of the best database

Classical approach

Manual Db
(11,256)

Complete
genomes
Database
(16,980)

Source: ViralZone
(247,326)

Can be used in existing pipelines

Research approach

ORF
database
(102,280)

Source: GenBank
(1,423,048)

Better representation of virus variability
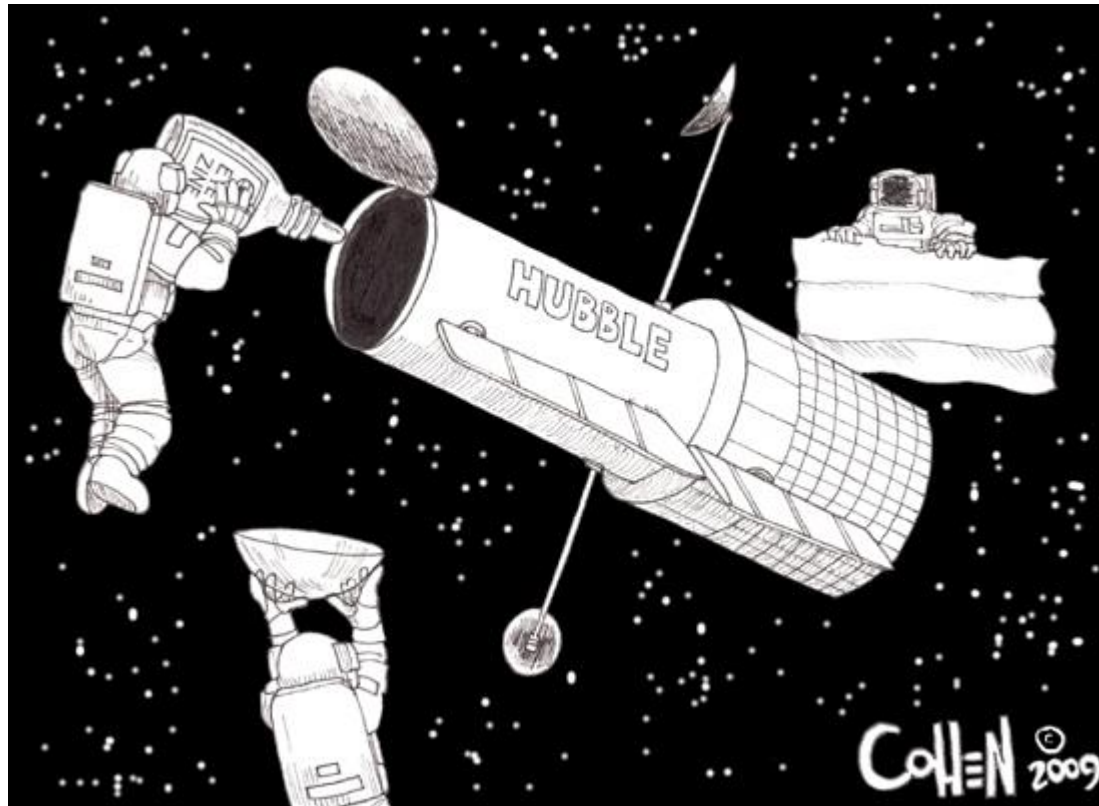Haplotypes

# Getting a «16S RNA» for viruses



All paramyxoviridae can be identified by their L polymerase polymorphism

# Conclusion

- Databases are a **key elemen**t to identify/characterize microbial organisms

- Clinical metagenomics needs **dedicated databases**

- **Sequence curation** with clinical focus facilitates the intrepretation of results

# Thank you for your attention

# Acknowledgements

**Swiss-Prot and Vital-IT**
Ioannis Xenarios

**Swiss-Prot Group**
-Alan Bridge
-Patrick Masson
-Chantal Hulo
-Edouard de Castro
-Andrea Auchinchloss

**Vital-IT**
-Anne Gleizes
-Nicolas Guex
-Christian Iseli
-Thomas Junier

**SIB**
Jacques Fellay
Valérie Barbié
Aitana Lebrand

**Collaborations**
Laurent Kaiser
Samuel Cordey
Florian Laubscher