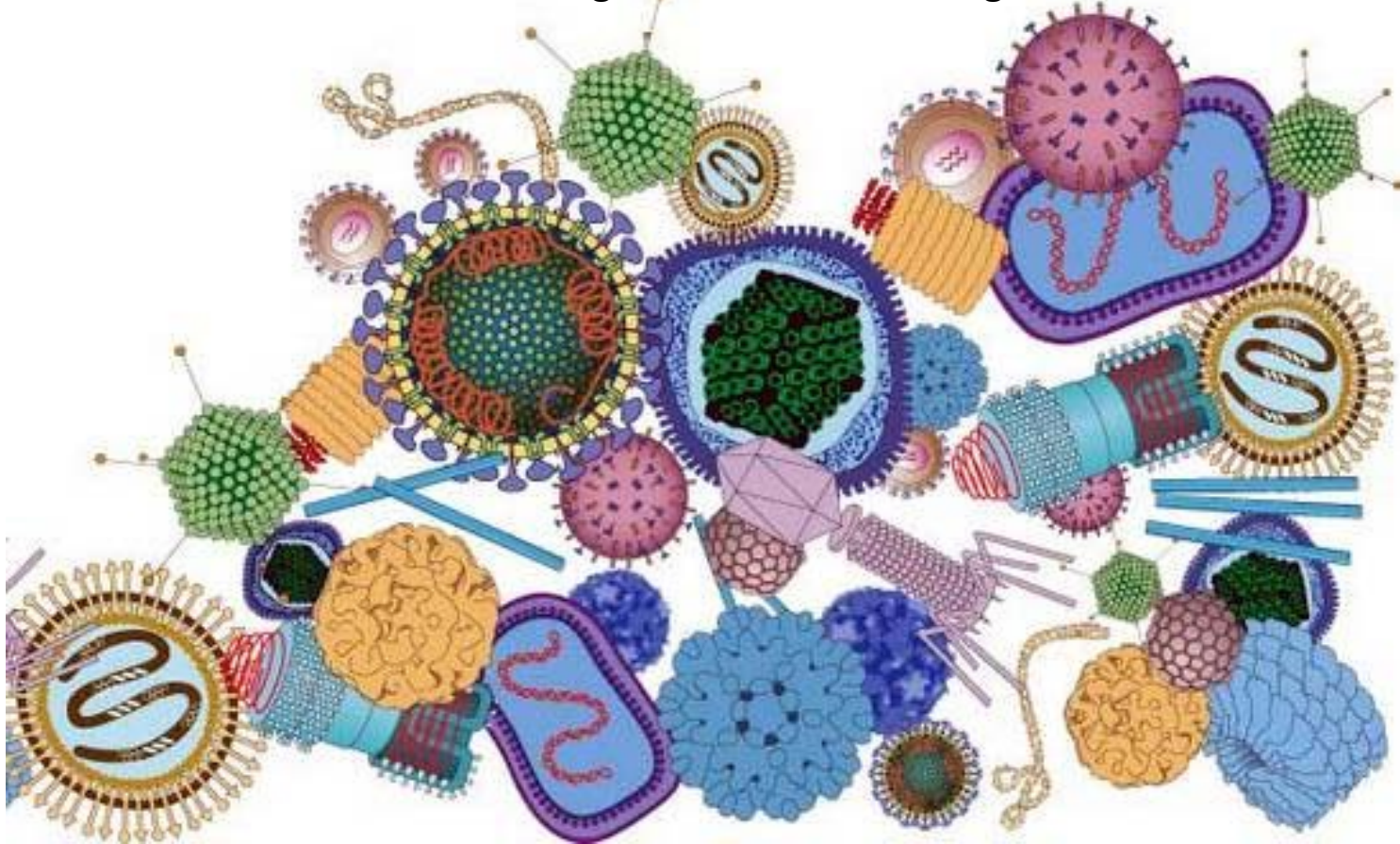


Identification of highly divergent viral sequences.

NIST workshop on standards for NGS detection of viral adventitious agents
in biologics/biomanufacturing

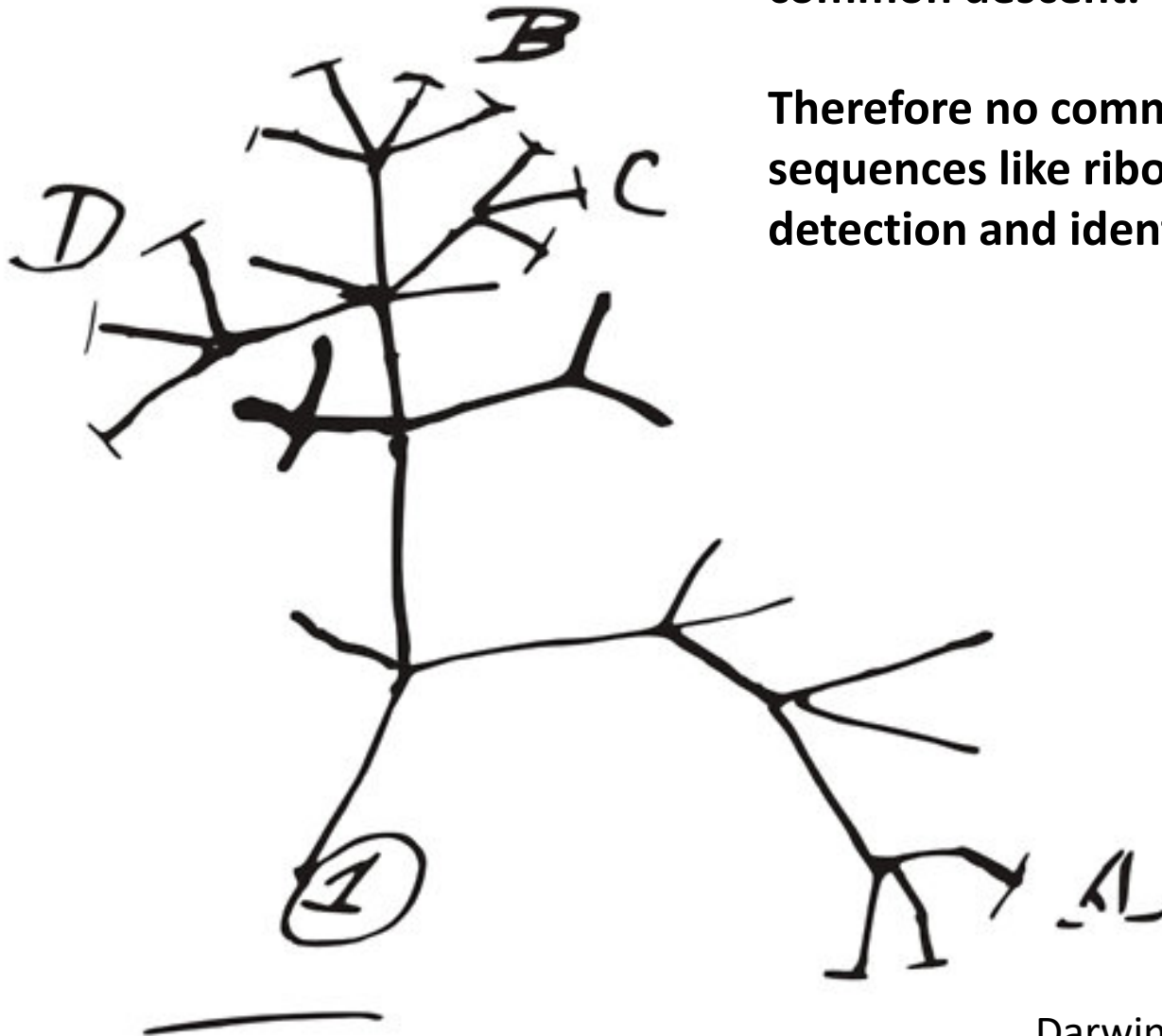


Eric Delwart
Vitalant Research Institute & UCSF Lab Medicine.

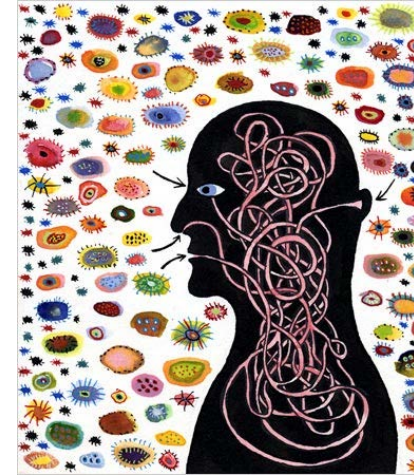
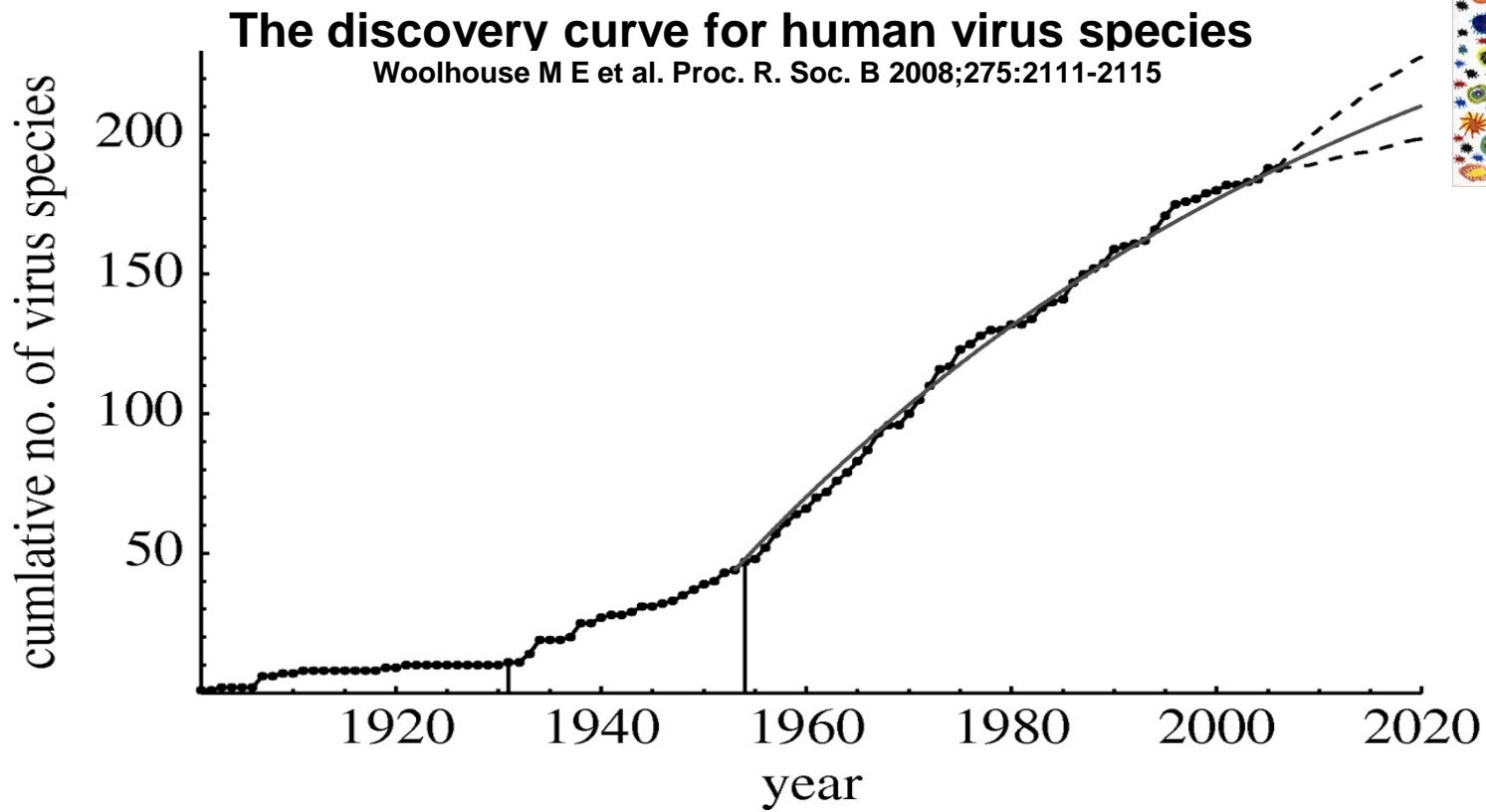
I think

Unlike viruses all cells have some homologous genes from common descent.

Therefore no common sequences like ribosomes for detection and identification



Darwin 1837



Some “new” human viruses **since 2005**

- Parvoviruses:

- B19 erythrovirus (3 genotypes)
- **PARV4** (3 genotypes)
- **Bocavirus-1/2** (2 species)
- **Bufavirus-1/2** (2 species)
- **Tusavirus** (1 species)
- **Cutavirus** (1 species)

- Picornaviruses:

- Enterovirus (7 species >260 serotypes)
- Hepatovirus (1 species)
- Parechovirus (1 species)
- Kobuvirus (1 species)
- **Cosaviruses** (5 species >33 genotypes)
- **Salivirus** (1 species)
- **Cardioviruses** (1 species >8 genotypes)

- Astroviruses:

- Human mamastrovirus 1 (8 genotypes)
- **Human astroviruses MLB-1/2/3** (3 genotypes)
- **Human astroviruses HMO-1/2/3 VA4** (4 genotypes)
- **Human mamastrovirus 20** (1 genotype)

- Polyomaviruses:

- BKV/JC -KIPyV -WUPyV
- MCPyV -HPyV6 -HPyV7
- TSPyV -HMWPyV-HPyV9
-**HPyV10** -HPyV11

Sources of viral contaminations in Biologics

- History of contamination starting with live attenuated Yellow Fever Virus vaccine contaminated with HBV from “stabilizer” human serum.
- Frequent use of animal serum in cell cultures.
- Other animal product: porcine trypsin.
- Use of contaminated cell lines.

QC of Biologics for viruses

- Avoidance of laboratory contamination requires stringent QC similar to ancient DNA labs (segregation of amplified nucleic acids like PCR amplicons, libraries, and plasmids, unidirectional samples flow).
- NGS data management to detect possible source of contamination (virus present in recent library).
- Awareness of bleed-over between samples sequenced on same Illumina chip.
- Problem bigger for ID diagnostics handling more (and more diverse) samples.

Baoyan Xu^{a,b,1}, Ning Zhi^{a,1,2}, Gangqing Hu^{c,1}, Zhihong Wan^a, Xiaobin Zheng^d, Xiaohong Liu^a, Susan Wong^a, Sachiko Kajigaya^a, Keji Zhao^{c,3}, Qing Mao^{b,2}, and Neal S. Young^{a,3}

F	C
+	-
+	+
+	-
+	-
+	-
-	-
-	-
-	-
-	-

Experimental animals asymptotomatically infected

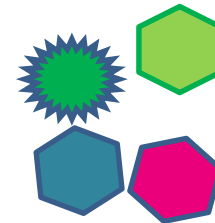
- Mice
- Rat
- Rhesus macaques
- African green monkey
- Zebrafish

Viral metagenomics



virus sized
fraction

→ DNase and RNase



Enriched
Virus
particles

RNA + DNA extraction

Random RT-PCR

5' GTCCATGCATGACTCGAGTCNNNNNNNNN3'

Illumina seq ←



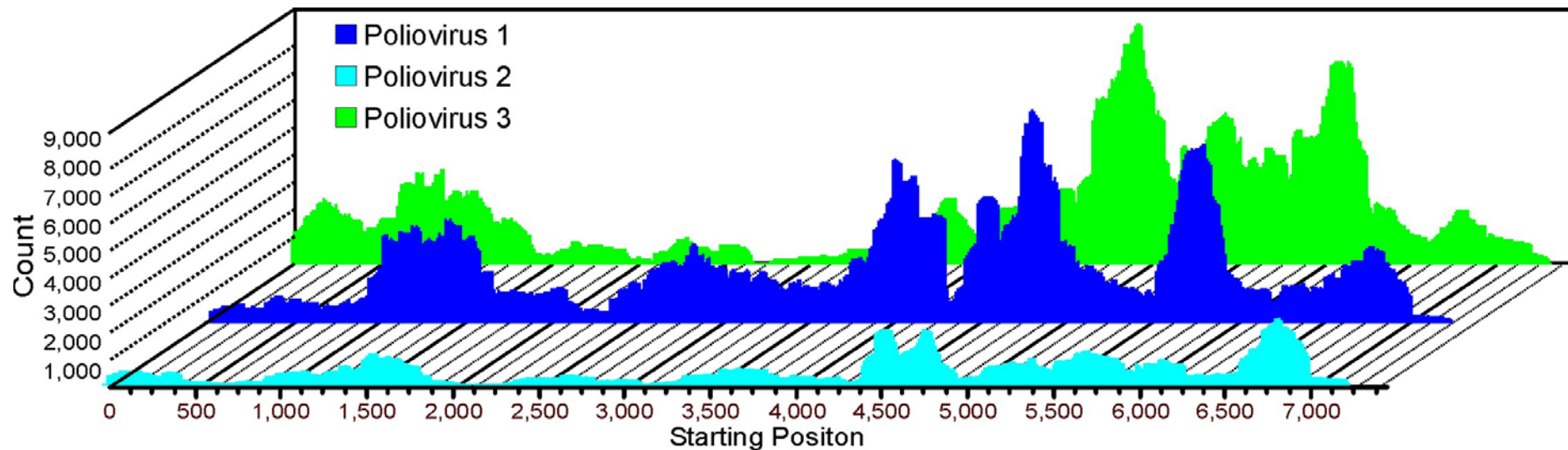
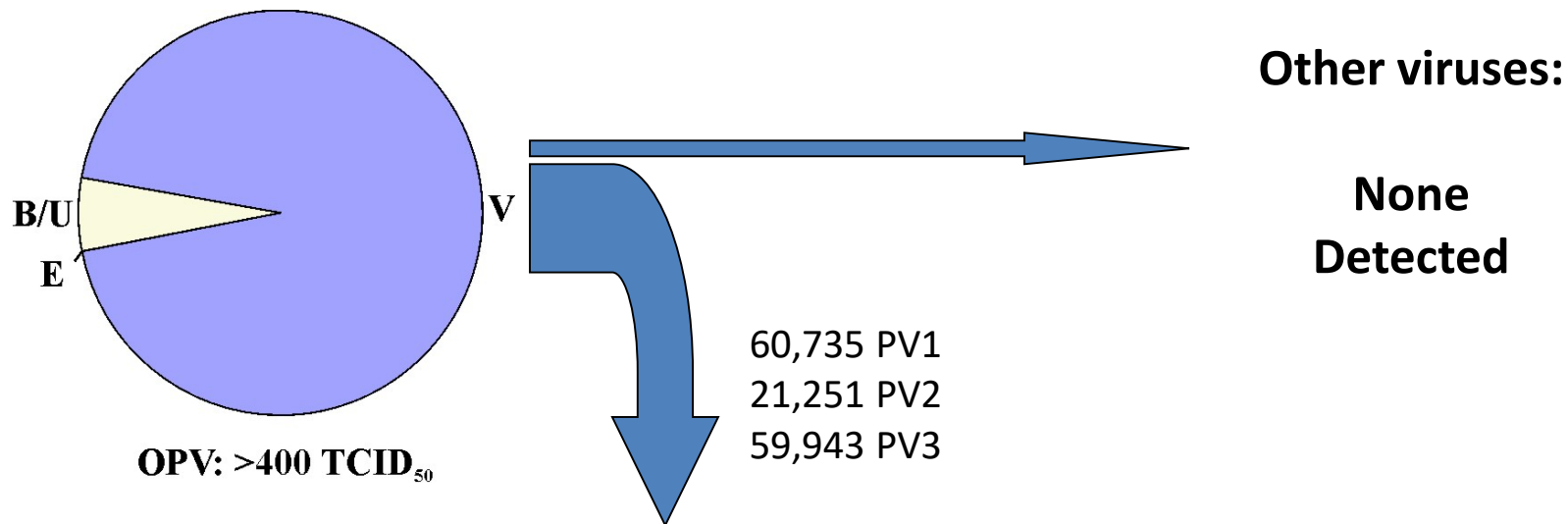
De novo assembly

BLAST

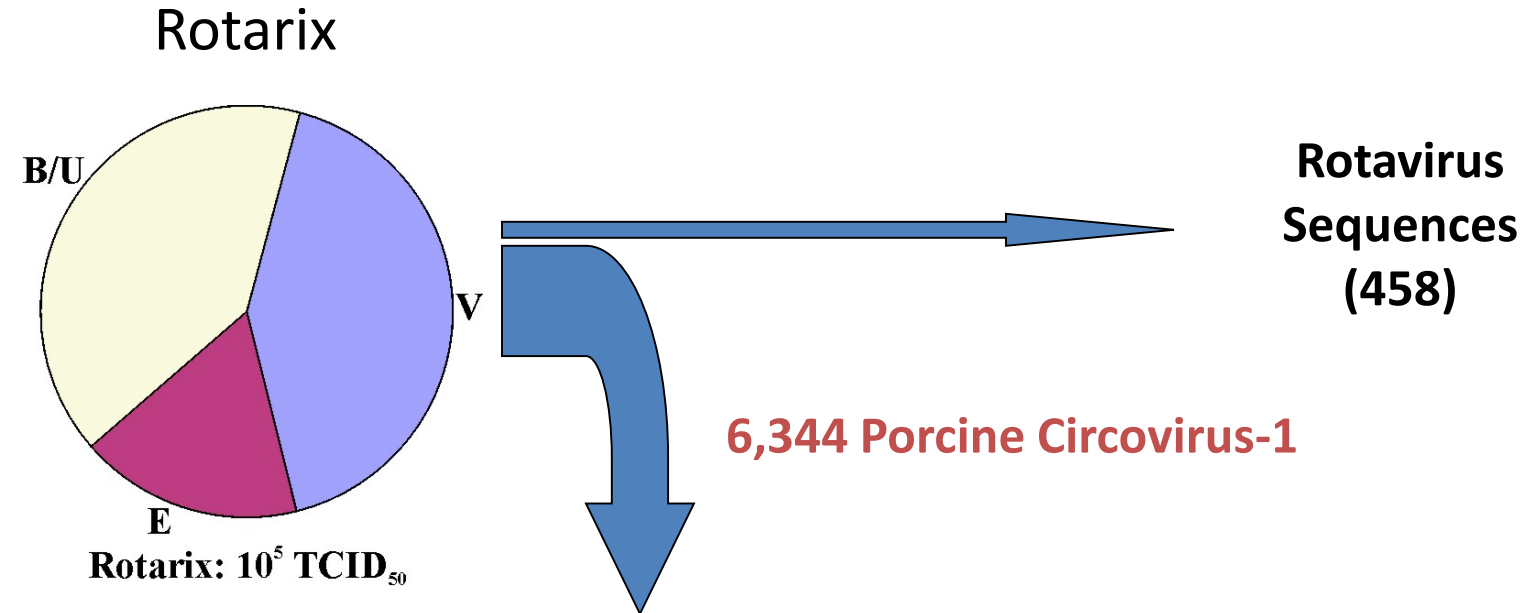
against all viruses



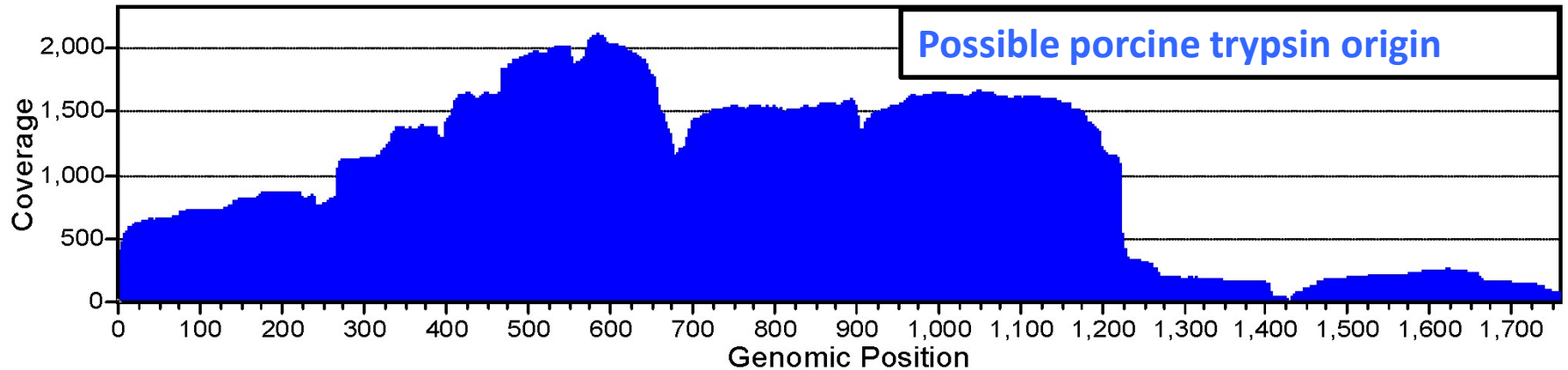
Oral poliovirus 1,2,3



Metagenomics is useful tool to detect adventitious virus contaminations



Rotarix (GlaxoSmithKline) Contamination with Porcine Circovirus



Method is non-specific

Basic

Local

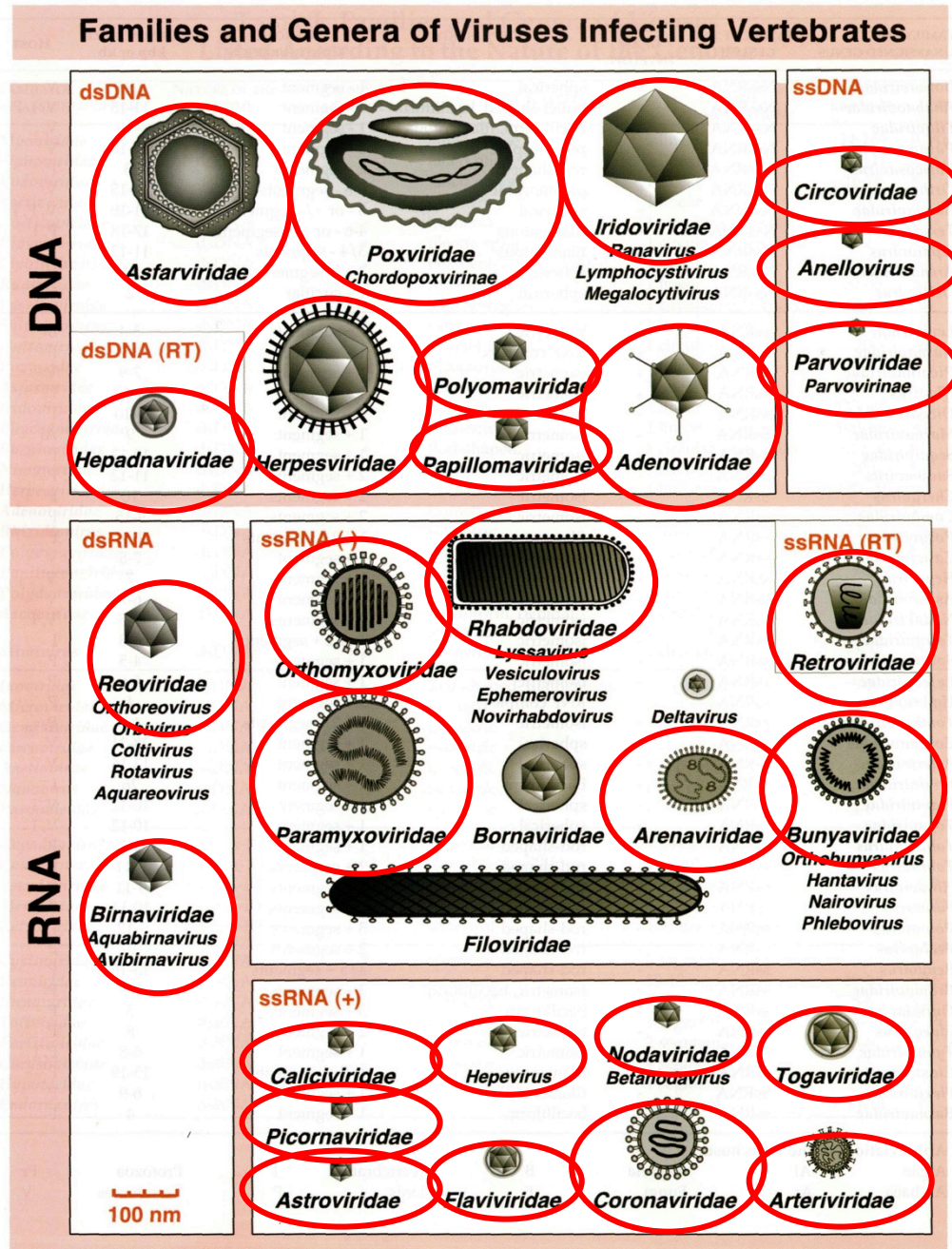
Alignment

Search

Tool

Viruses detected:

BLASTx E score $<10^{-5}$



Library prep options

- Enrich VLP or extract total NA?
- Total NA versus RNA or DNA?
- Random amplification, primers targeting many viruses, combination?
- Random RT-PCR or Rolling Circle Amplification?
- Libraries finished by ligation or transposon?

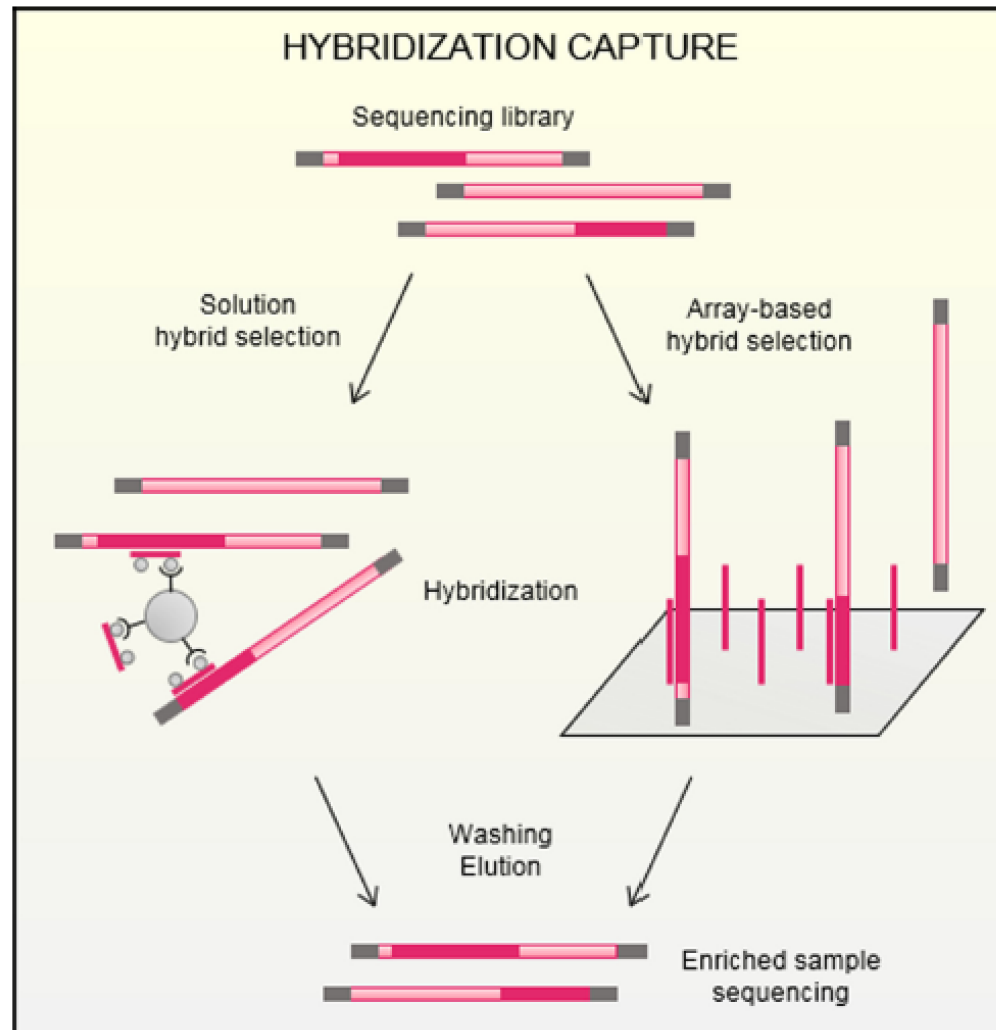
Sensitivity approaching real time PCR using 25 viruses pool (NIBSC)

Virus	Ct	N1	N230	N231	N12	N2	N3	N227	N221	N225	N226	N32	N4	N42	N5	N232	N233	N6
RVA	24.49	5588	3617	6000	1535	4401	7110	5277	5459	2896	1938	1375	9004	1360	18357	24380	24988	1302
HHV5	28.95	156	100	282	62	9	536	457	322	447	298	57	376	40	514	148	112	6
HHV3	29.02	1533	176	43	338	2227	836	826	678	330	219	303	1799	292	1600	350	564	814
HPeV3	29.35	11893	5877	7608	4452	16548	19423	18843	10190	3507	2494	3891	13545	3666	17973	19938	24802	6163
AdV2	29.71	1	0	1	1	249	94	284	301	260	113	18	0	11	16	6	6	0
AstV	30.53	0	0	0	6	69	219	0	0	14	7	1	0	8	24	36	70	0
HHV1	30.59	0	0	0	0	0	0	0	0	11	4	2	0	1	3	0	0	0
CVB4	30.72	1	0	1	6	0	64	12	12	24	4	7	5	9	60	111	189	0
HRV A39	31.16	0	0	0	0	0	0	2	0	6	2	0	0	0	11	32	39	0
HHV4	31.27	0	0	0	3	0	8	0	5	34	31	8	2	11	26	12	6	295
PIV4	31.83	0	0	0	11	1	385	38	53	24	18	18	1	19	18	37	47	0
HMPV A	31.86	0	0	0	0	0	46	0	5	26	35	5	0	0	15	44	40	0
IFV A H1N1	32.02	0	0	0	0	0	7	0	2	2	3	4	0	0	0	0	0	0
HHV2	32.48	0	0	0	0	0	0	0	2	7	9	1	0	0	0	1	0	0
SaV C12	33.37	0	0	0	1	11	0	95	45	14	2	0	0	0	28	84	70	0
PIV2	33.87	3	0	2	35	1	882	102	102	253	164	42	0	30	2	4	4	0
RSV A2	34.33	0	0	0	0	0	0	0	0	4	12	0	0	0	0	5	2	0
PIV1	34.43	4	0	0	57	5	22	32	39	44	26	46	5	59	78	161	207	1480
CoV 229E	36.48	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
NV GI	ND	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
NV GII	ND	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
IFV B	ND	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0
IFV A H3N2	ND	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0
PIV3	ND	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0
AdV41	ND	0	0	0	0	0	0	0	0	4	4	0	0	0	0	0	0	0
Virus detected		8	4	7	12	10	13	11	14	20	20	15	8	12	15	17	15	6

Target capture

uses synthetic RNA baits on magnetic beads to pull down targeted viral sequences increasing sensitivity.

Restricted to already known viruses or close relatives.



Data analysis options

- Use curated viral genome/proteome databases. Update frequently.
- Limit to known human/mammalian/vertebrate/eukaryotic viruses?
- Include virus of unknown tropism (eg many CRESS-DNA viruses, many members of *Picornavirales* order)?
- Nucleotide similarity for fast analyses and close similarities only.
- Protein similarity for more divergent viruses.

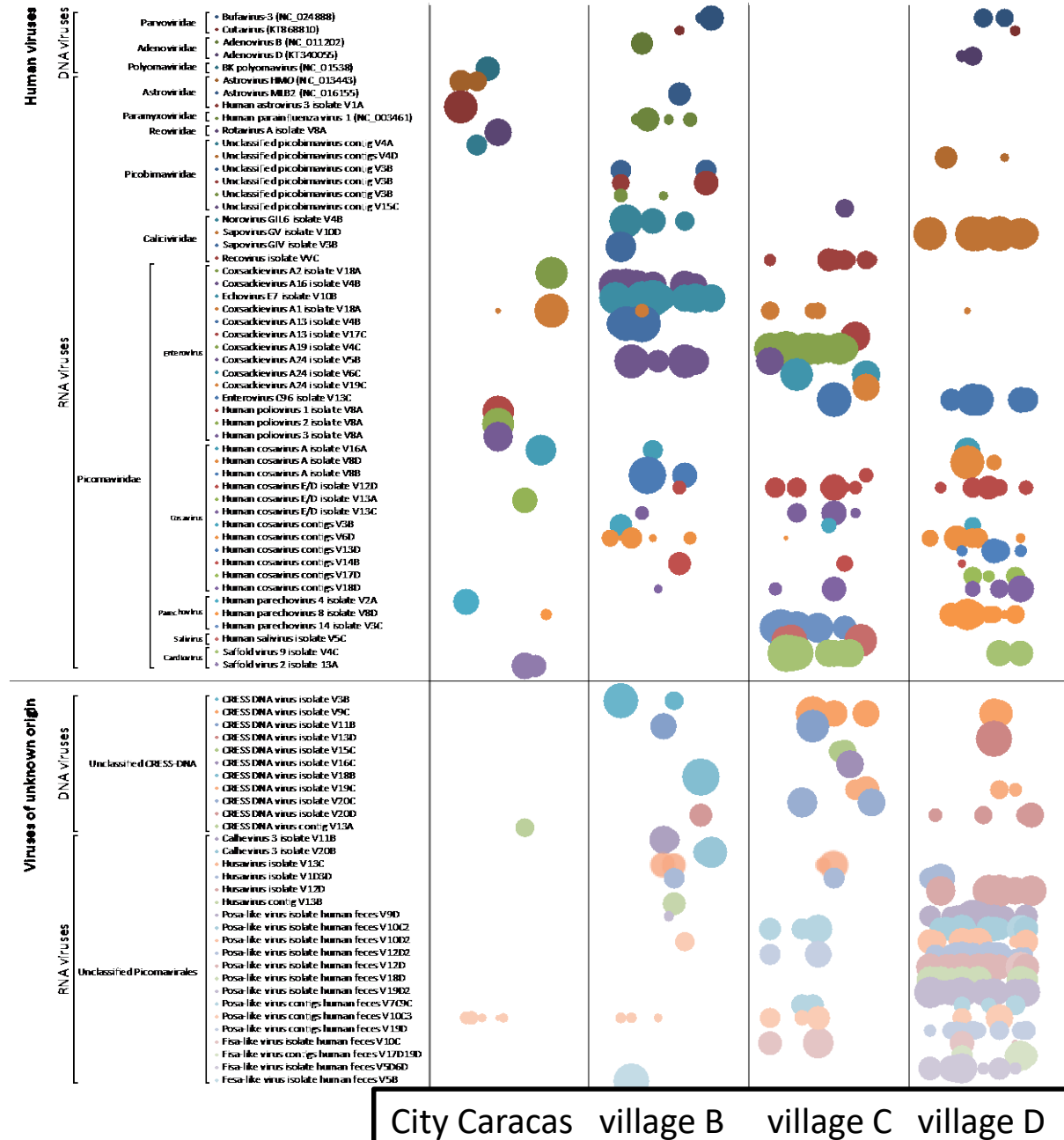
Distribution of eukaryotic viruses in feces from 80 healthy children

Human DNA viruses

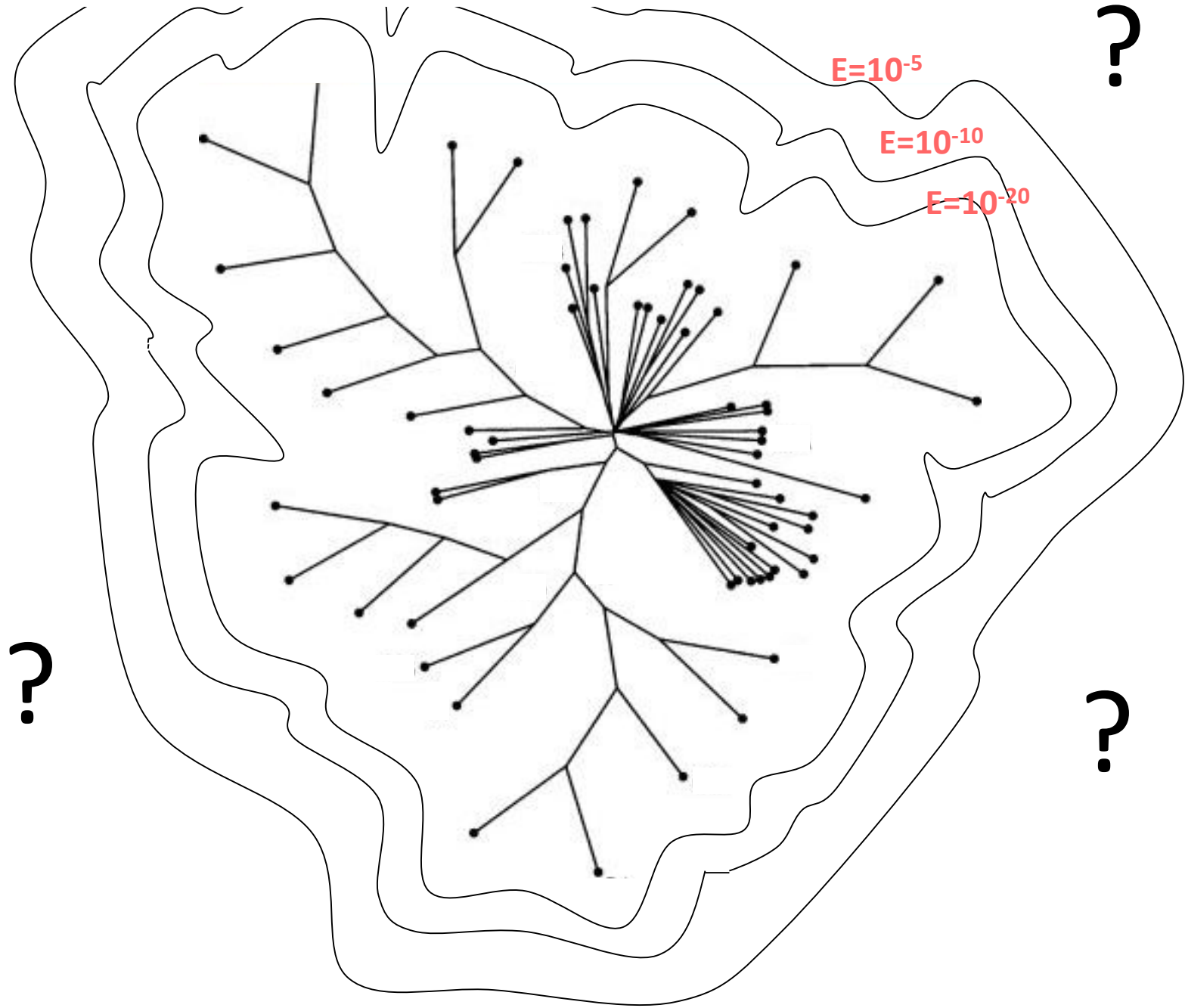
Human RNA viruses

DNA viruses unknown tropism

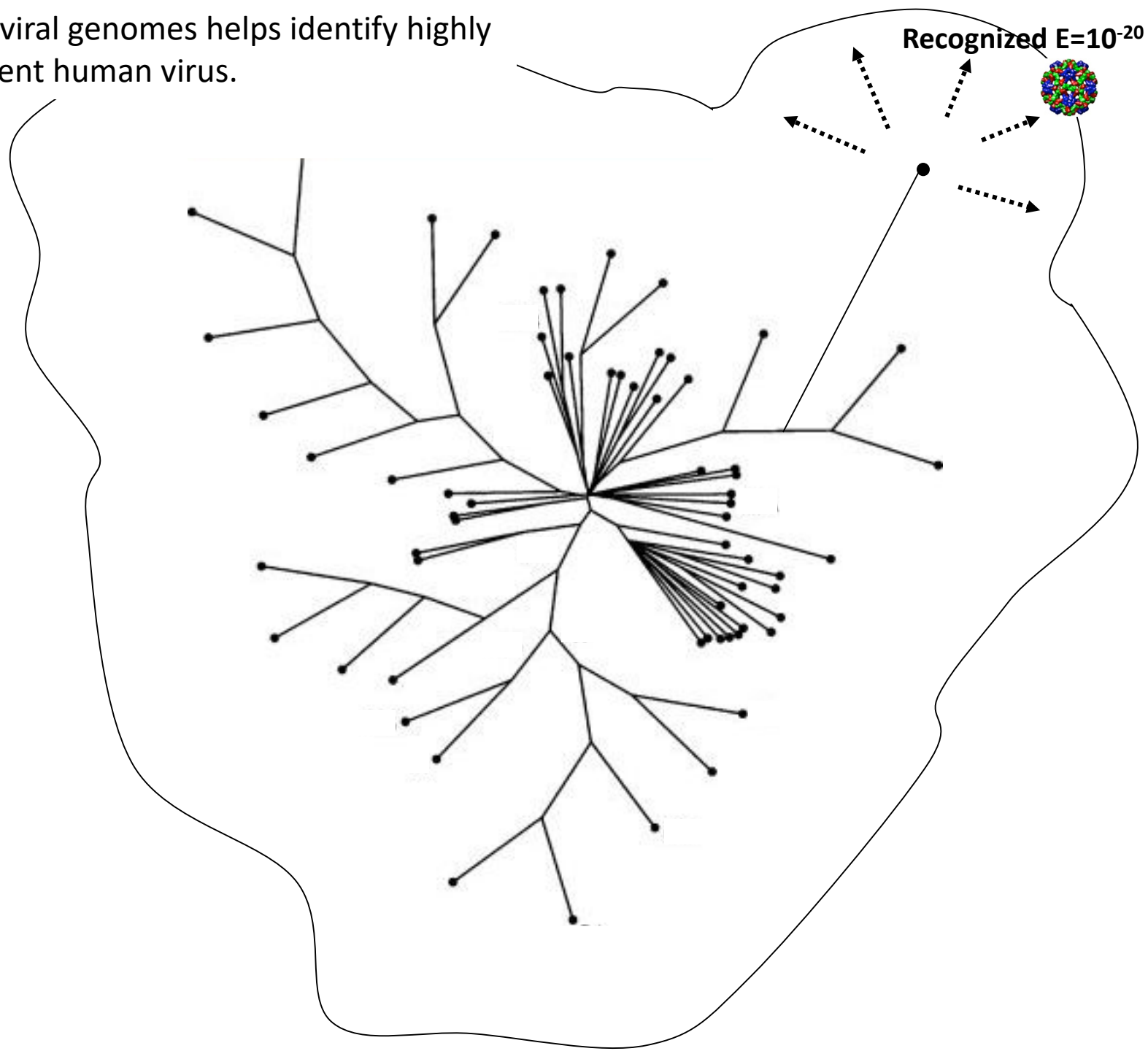
RNA viruses unknown tropism



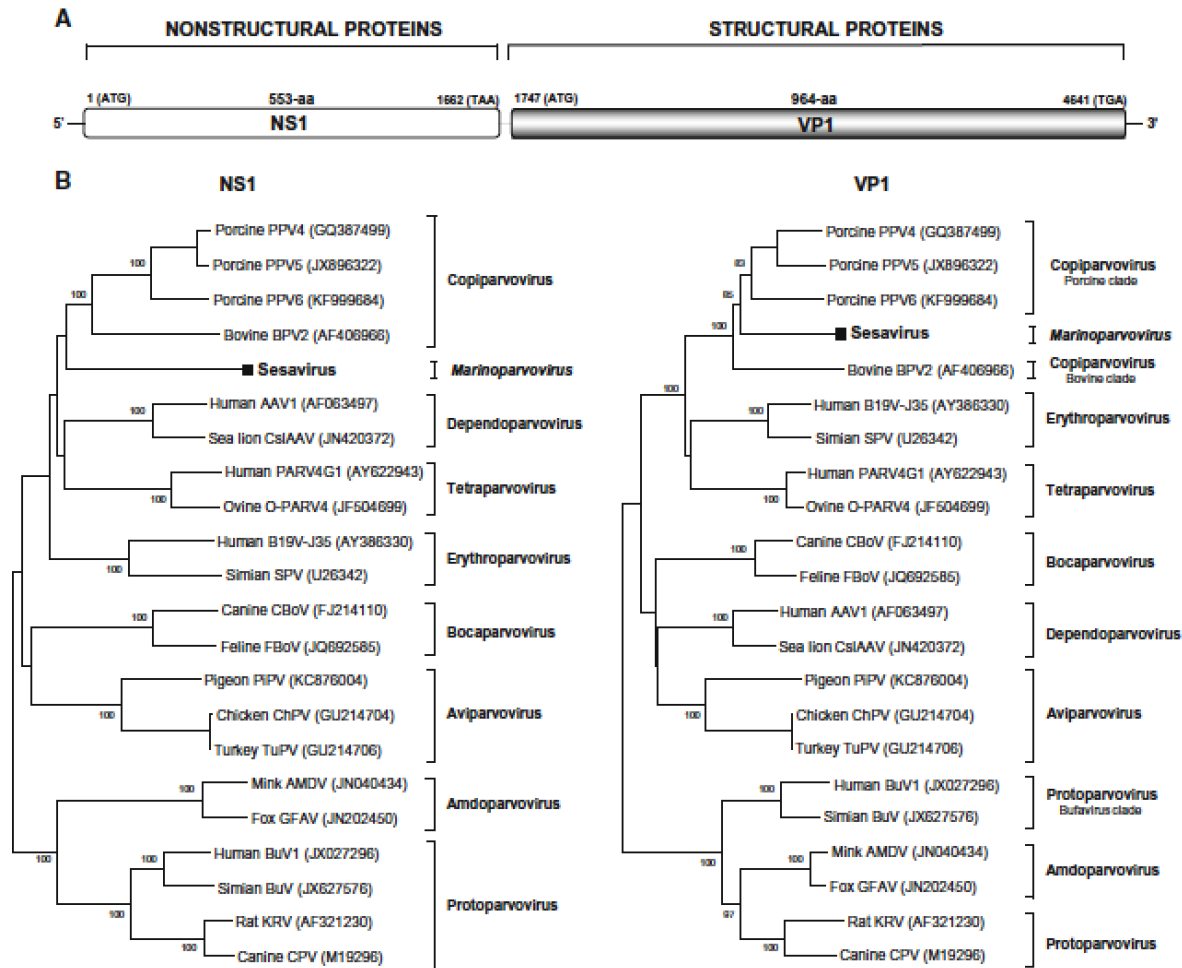
Highly divergent viruses are not recognized by similarity search.



Novel viral genomes helps identify highly divergent human virus.



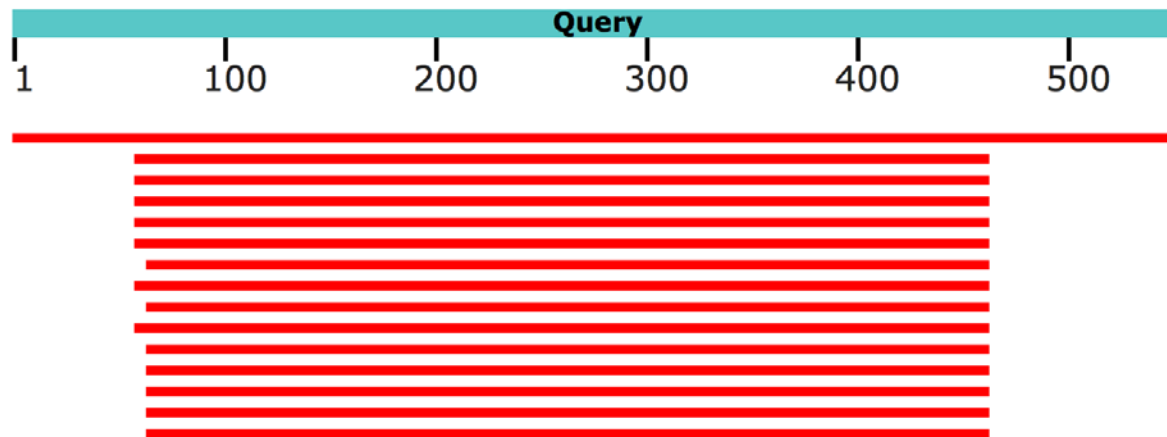
Identifying a new parvovirus genus



Protein BLASTx of novel NS

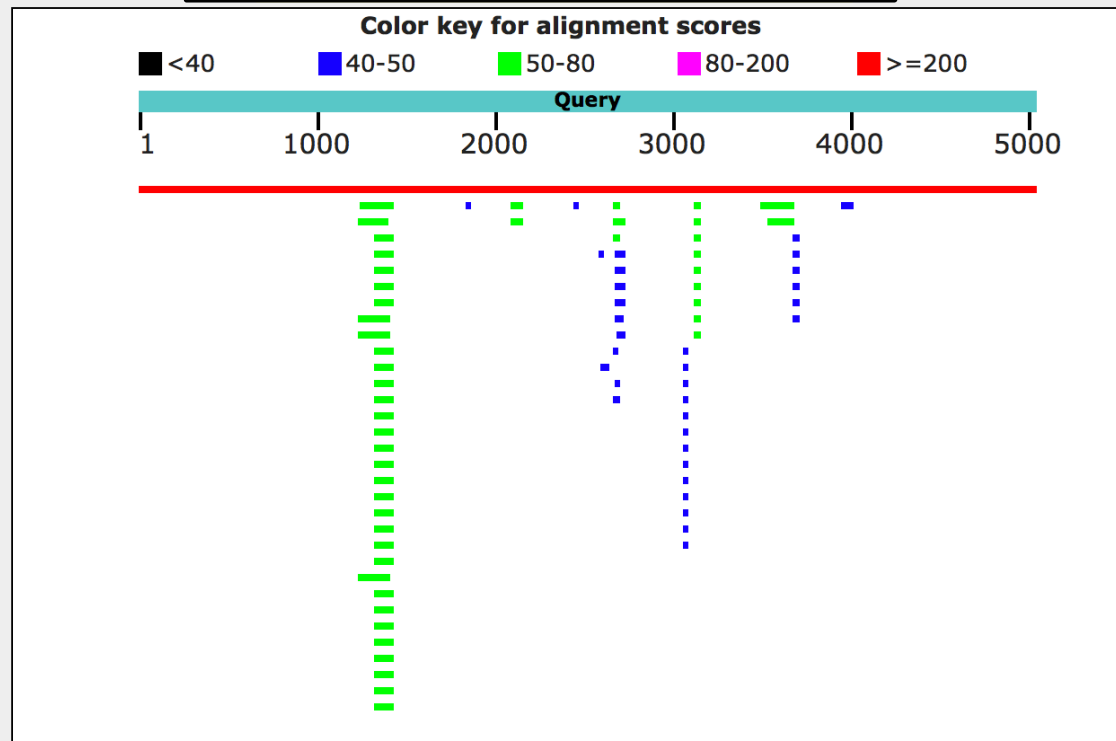
Color key for alignment scores

<40
 40-50
 50-80
 80-200
 >=200



	Max Score	Total Score	Query Cover	E value	Per. Ident	Accession
NS1 [Sesavirus CSL10538]	1162	1162	100%	0.0	100.00%	YP_009116876.1
replicase [Porcine parvovirus 5]	241	241	72%	3e-68	33.02%	AFU93187.1
nonstructural protein 1 [Porcine parvovirus 5]	239	239	72%	2e-67	32.79%	ANM72084.1
replicase [Porcine parvovirus 5]	239	239	72%	2e-67	32.95%	YP_008888533.1
replicase [Porcine parvovirus 5]	237	237	72%	2e-67	32.56%	QBA84740.1
replicase [Porcine parvovirus 5]	239	239	72%	2e-67	32.79%	QBA84795.1
replicase [Porcine parvovirus 5]	236	236	71%	3e-67	33.18%	QBA84791.1

DNA BLASTn of novel genome



	Max Score	Total Score	Query Cover	E value	Per. Ident	Accession
Sesavirus CSL10538 NS1 and VP1 genes, complete cds	9106	9106	100%	0.0	100.00%	KM035804.1
Goat protoparvovirus strain C4 nonstructural protein 1 gene, partial cds	70.7	70.7	3%	1e-06	68.72%	MH835426.1
Bosavirus sp. strain 738 putative structural protein gene, partial cds	69.8	69.8	3%	1e-06	68.53%	MH248790.1
Goat protoparvovirus strain C3 nonstructural protein 1 gene, partial cds	66.2	66.2	3%	1e-05	68.18%	MH835425.1
Chicken parvovirus strain GX-CH-PV-12, complete genome	59.9	59.9	2%	0.002	72.38%	KX133419.1
Galliform aveparvovirus 1 isolate IPV, complete genome	59.9	59.9	2%	0.002	72.38%	KU569162.1
Chicken parvovirus isolate ADL120686, complete genome	59.9	59.9	2%	0.002	72.38%	KJ486491.1

Deposited after reference

A divergent calicivirus with borderline E score 3×10^{-4} (23% protein identity)

all tag blast logout

Plasma pools 5

Plasma pools 6

Cancer Tissue

Rat Stool

Danish respiratory 1

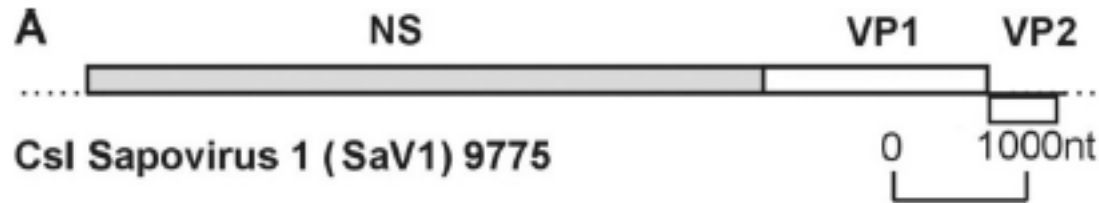
Danish respiratory 2

Pig

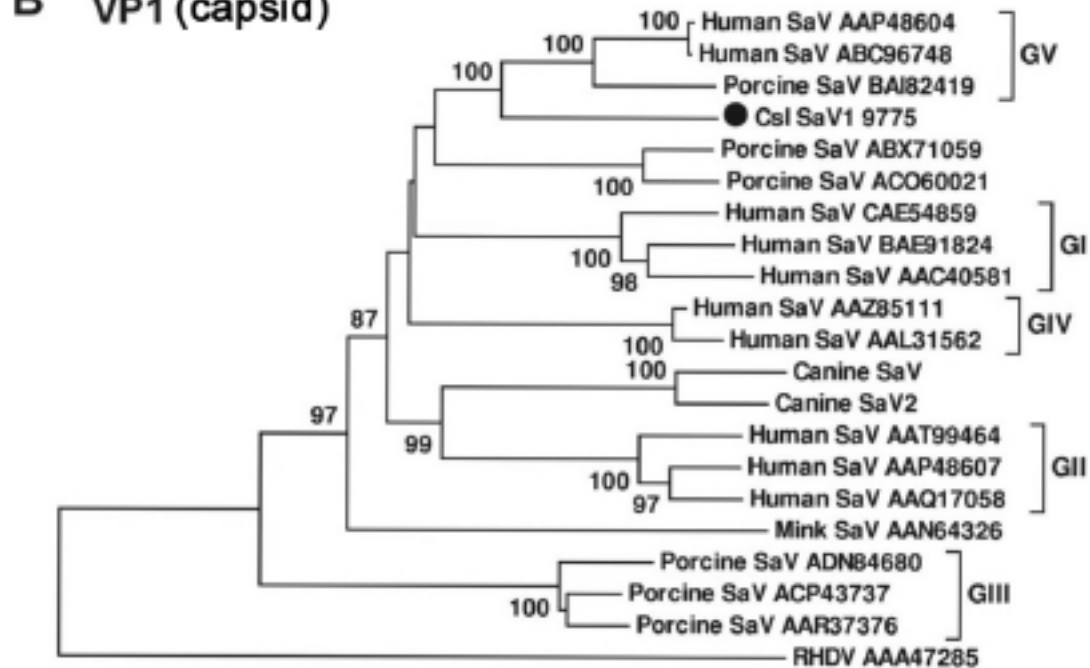
RatVole

Sea lions

vhit	
1125	1927
1136	956
1137	673
1140	890
1141	574
1148	393
1153	1034
1157	1518
1162	1056
1166	767
1169	865
1170	873
1174	1081
1181	822
1182	1316
1185	1234
1187	1309
1194	1185
1199	960



B VP1 (capsid)



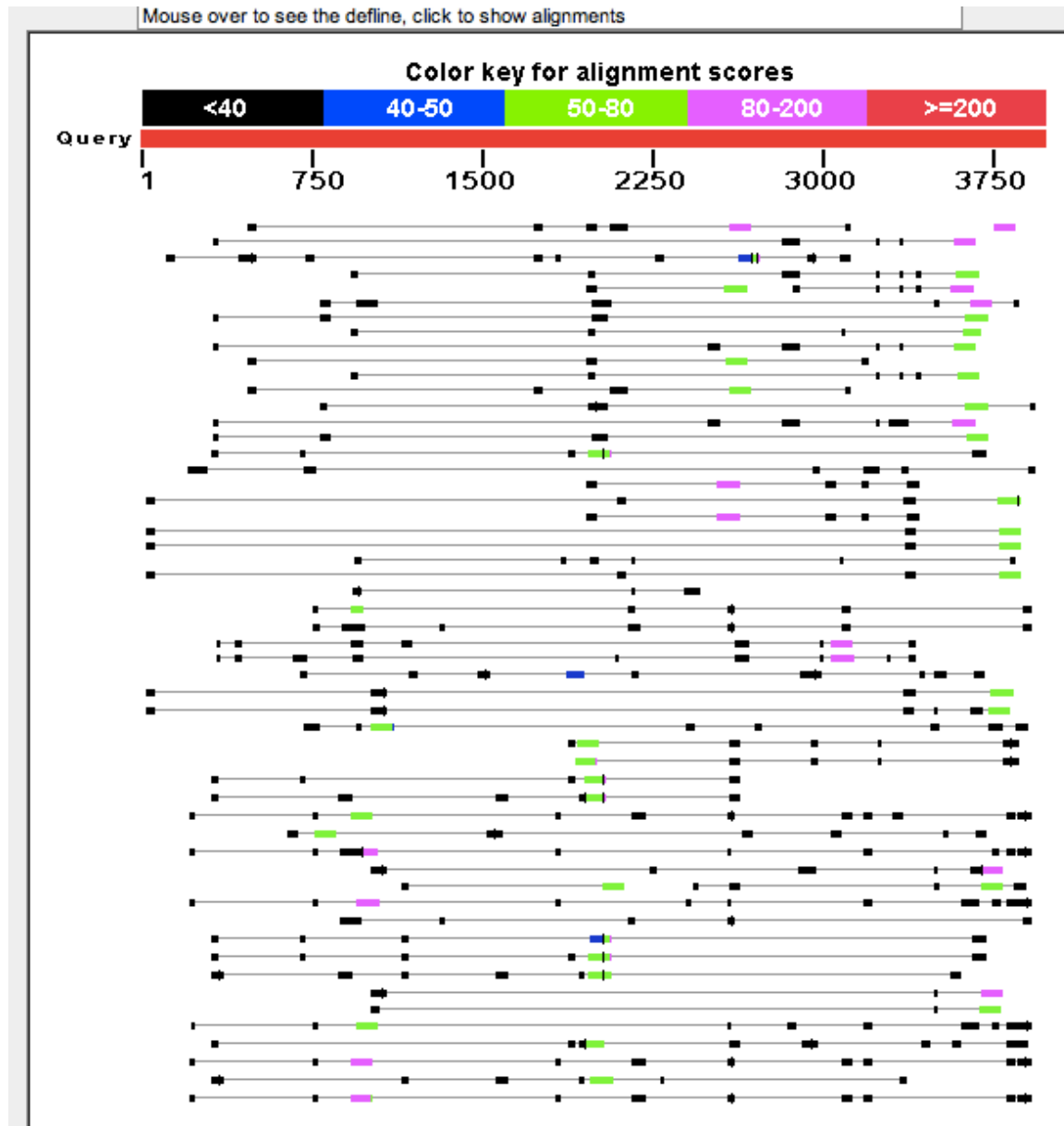
```

-VVPLKSDPDETRWHARTHNYTF 323
 V K+ PD +W ++T
TTVTAKTQPD-LKWGLDADHTA 1170
-----SEVFFQAGQPNVSRQ 503
 ++V P A++ +
PLPDVIATQVRPYQDPVVAINGE 1230
TAAGPSYGVVKQSVFNPDGAPAK 683
++ GP + VK+ DG P
SSNGPWFPGVKKDYTTEDGGPNP 1290
AGLTRPIFCFD 760
 G R ++ D
DGKRRLWGCD 1336
-----
licivirus]
    
```

CSL-9806_01301	1	170	56	
----------------	---	-----	----	--

MINING EXISTING DATABASE:

Re-tBLASTx with “new” viral genomes



[illegible]

Asfaviridae

[illegible]

Large fraction of reads are unidentifiable by BLAST

- This “**dark matter**” may contain “new” viral families not yet found in GenBank’s Virus RefSeq
- Weak protein similarity searches using sequence profiles (HHPRED)
- Most new viral families are phages.

Virus Taxonomy

International Union of Microbiological Societies
Virology Division

Andrew M. Q. King • Michael J. Adams
Eric B. Carstens • Elliot J. Lefkowitz



Viruses - taxonomy groups

- | | | |
|----------------------------|--------------------------------------|--------------------|
| Adenoviridae | Alloherpesviridae | Alphaflexiviridae |
| Alphatetraviridae | Alvernnaviridae | Amalgaviridae |
| Ampullaviridae | Anelloviridae | Arenaviridae |
| Arteriviridae | Ascoviridae | Asfarviridae |
| Astroviridae | Baculoviridae | Barnaviridae |
| Benyviridae | Benyvirus | Betaflexiviridae |
| Bicaudaviridae | Bidnaviridae | Birnaviridae |
| Bornaviridae | Bromoviridae | Bunyaviridae |
| Caliciviridae | Carmotetraviridae | Caudovirales |
| Caulimoviridae | Chrysoviridae | Cilevirus |
| Circoviridae | Closteroviridae | Coronaviridae |
| Corticoviridae | Cystoviridae | Deltavirus |
| Dicistroviridae | Dinodnavirus | Emaravirus |
| Endornaviridae | Filoviridae | Flaviviridae |
| Fusarividae | Fuselloviridae | Gammaxflexiviridae |
| Geminiviridae | Globuloviridae | Hepadnaviridae |
| Hepeviridae | Herpesvirales | Herpesviridae |
| Hypoviridae | Hytrosaviridae | Idaeovirus |
| Iflaviridae | Inoviridae | Iridoviridae |
| Leviviridae | Ligamenvirales | Lipothrixviridae |
| Luteoviridae | Malacoherpesviridae | Marnaviridae |
| Marseilleviridae | Megabirnaviridae | Mesoniviridae |
| Microviridae | Mimiviridae | Mononegavirales |
| Myoviridae | Nanoviridae | Narnaviridae |
| Nidovirales | Nimaviridae | Nodaviridae |
| Nudiviridae | Nyamiviridae | Ophioviridae |
| Orthomyxoviridae | Ourmiavirus | Papillomaviridae |
| Paramyxoviridae | Partitiviridae | Parvoviridae |
| Permutotetraviridae | Phycodnaviridae | Picobirnaviridae |
| Picornavirales | Picornavirales environmental samples | Picornaviridae |
| Plasmaviridae | Podoviridae | Polemovirus |
| Polydnaviridae | Polyomaviridae | Potyviridae |
| Poxviridae | Quadriviridae | Reoviridae |
| Retro-transcribing viruses | Retroviridae | Rhabdoviridae |
| Rhizidiovirus | Roniviridae | Rudiviridae |
| Salterprovirus | Satellites | Secoviridae |
| Siphoviridae | Sobemovirus | Sphaerolipoviridae |
| Tectiviridae | Tenuivirus | Togaviridae |
| Tombusviridae | Totiviridae | Turriviridae |
| Tymovirales | Tymoviridae | Umbrovirus |



"I'm searching for my keys."

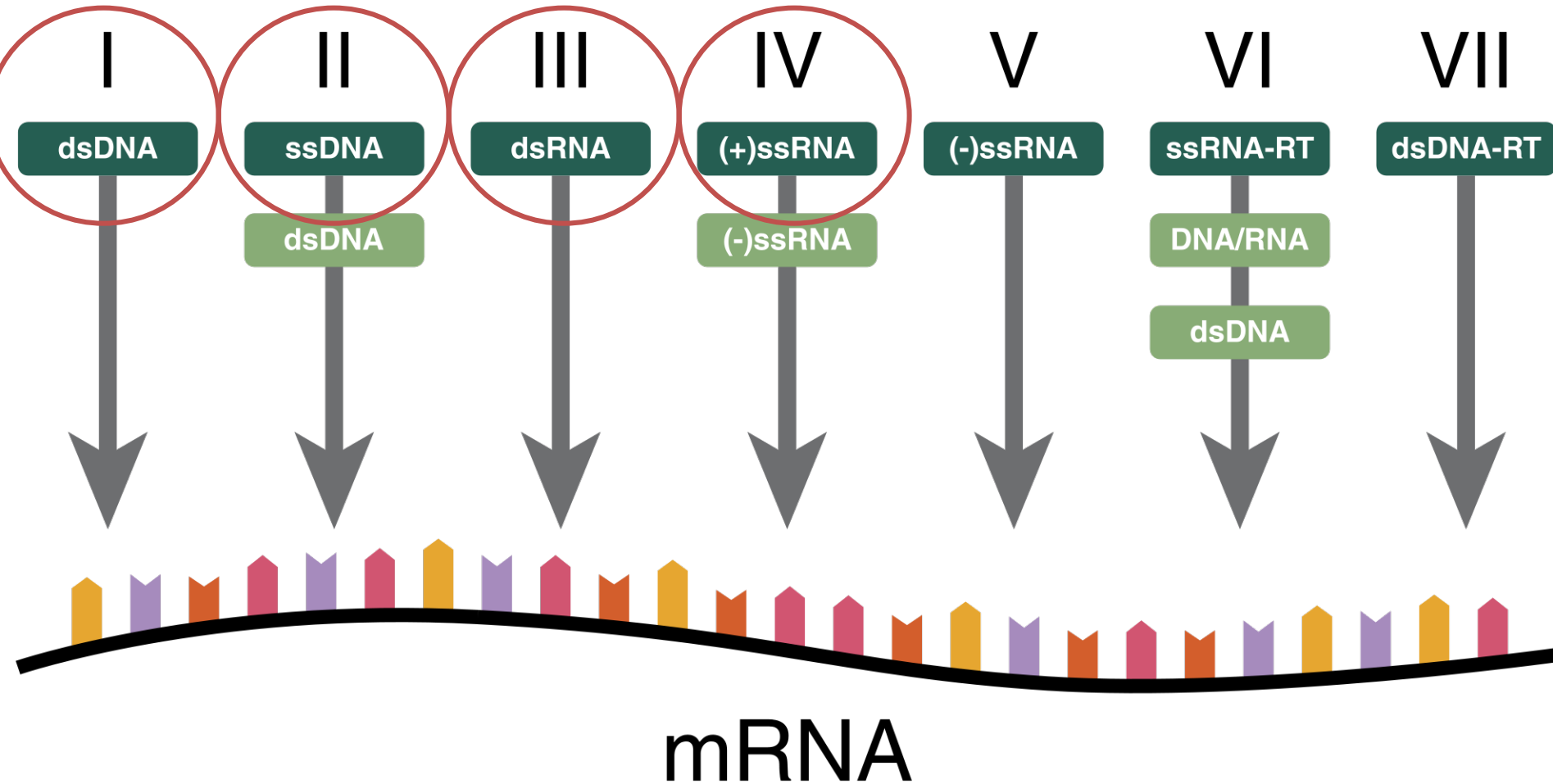
A wide diversity of virion and genome size

Classification criteria	Nucleic acid																					
	RNA										DNA											
	Icosahedral					Helical					Icosahedral			Helical		Complex						
	Naked			Enveloped		Enveloped					Naked		Enveloped		Naked/Env. (cytoplasmic)	Enveloped	Enveloped (cytoplasmic)					
	ds 10-18 seg.	ds 2 seg.	(+) ss cont.	(+) ss cont.	(+) ss cont.	(+) ss cont.	(+) ss 2 copies	(+) ss cont.	(-) ss cont.	(-) ss cont.	(-) ss 3 seg.	(-) ss 8 seg.	(-) ss cont.	(-) ss 2 seg.	ss linear (+) or (-)	ds circular	ds linear	ds circle gapped	ds linear	ds linear	ds circular	ds linear (x linked)
Baltimore class	III	III	IV	IV	IV	IV	VI	IV	V	V	V	V	V	V	II	I	I	I	I	I	I	I
Family name	Reo	Birna	Calici	Picorna	Flavi	Toga	Retro	Corona	Filo	Rhabdo	Bunya	Orthomyxo	Paramyxo	Arena	Parvo	Papova	Adeno	Hepadna	Herpes	Irido	Baculo	Pox
Virion polymerase	(+)	(+)	(-)	(-)	(-)	(-)	(+)	(-)	(+)	(+)	(+)	(+)	(+)	(+)	(-)	(-)	(-)	(+)	(-)	(-)	(-)	(+)
Virion diameter (nm)	60-80	60	35-40	28-30	40-50	60-70	80-130	80-160	80 X 790-14,000	70-85 X 130-380	90-120	90-120	150-300	50-300	18-26	45-55	70-90	42	150-200	125-300	60 X 300	170-200 X 300-450
Genome size (total in kb)	22-27	7	8	7.2-8.4	10	12	3.5-9	16-21	12.7	13-16	13.5-21	13.6	16-20	10-14	5	5-8	36-38	3.2	120-200	150-350	100	130-280

Viruses can vary in genome size >1000 fold and in volume >10000 fold

VIRUSES CAN HAVE DIFFERENT TYPES GENOME
(Baltimore classification)

Class



Classic NGS QC standard

For viral particle purification, nucleic acid amplification, and sequencing.

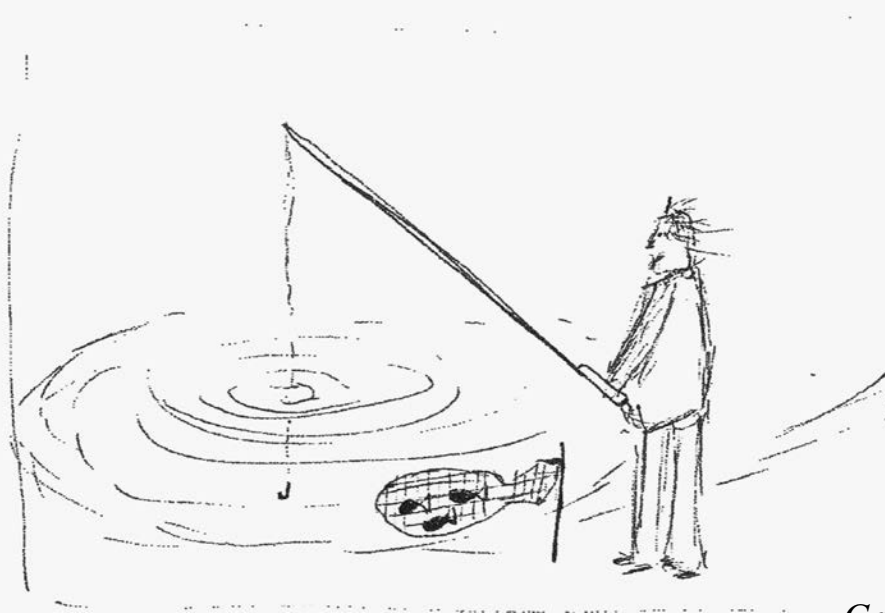
- Include small and large virions +/- lipid envelope. Include ds and ss DNA and RNA viruses.
- Quantify each virus by digital PCR.
- Pool should have approx. same genome equivalent of different viruses.
- Compare efficiency different NGS methods using reads per million (RPM) and % genome coverage.
- Set minimum standard LOD expected for each of the different spiked virus.
- Low concentration unlikely to reduce yield of possible contaminating virus

Synthetic NGS QC standard

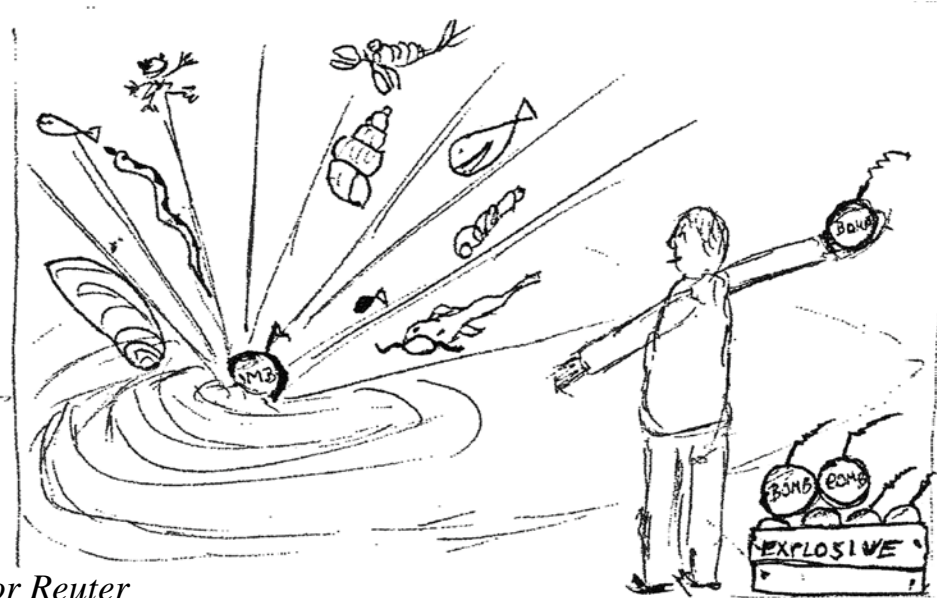
For nucleic acid amplification, and sequencing.

- Synthesize kb long nucleic acids of all known viral genome type: DNA and RNA, single and double stranded.
- Add to extracted nucleic acid to measure efficiency of library construction and sequencing of different type molecules.
- Spike biological sample to measure stability or incorporate into lipid particles of different size to mimic virions.
- Use any sequences. Keep typical viral GC/TA ratio.
- For the same nucleic acid structures (eg ds RNA) use different sequences at different concentrations to measure limit of sensitivity.
- Safer/cheaper/reproducible.

“conventional PCR”



“metagenomics”



Gabor Reuter



Vitalant Research Institute (aka Blood Systems) & UCSF

Xutao Deng-Bioinformatics

Eda Altan

Liz Fahsbender

Beatrix Kapusinszky (CA Dept Public Health)

Tung Phan (Indiana University, Indianapolis)

Juliana Siquera (Rio de Janeiro,Cancer institute)

Linlin Li (CA Dept Public Health)

Terry Ng (CDC)

Nikola Kondov (10X Genomics)

Lark Coffey (UC Davis)

Wen Zhang (Jiangsu University)

Tongling Shan (Shanghai Vet. Res. Institute)

Amit Kapoor (Ohio State U, Columbus)

Joe Victoria (Boehringer-Ingelheim)

Morris Saffold Jones (Mariposa Dept of Health)

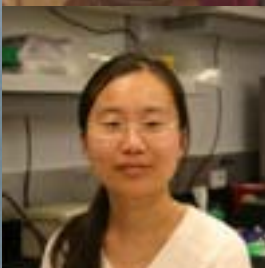
WHO polio eradication group in Pakistan

Sohail Zaidi

Oxford University

Peter Simmonds

Many other collaborators
for human, animal, and
environmental samples



UC Davis

Pat Pesavento

Alison Van Eenennaam

Ardesch Ardeschir

Denis Hartigan-O'Connor

UK NIBSC

Phillip Minor

Ed Mee

UCSF

Samia Naccache

Charles Chiu

Hungarian Natl Ref Lab

Gabor Reuter

Akos Boros

UCLA

Nelson Friemer

IDEXX

Christian Leutenegger

Funding:

NHLBI, NIAID

Vitalant Research Institute