

宏基因组二代测序数据的病毒鉴定 自动化分析流程

李洋

liyang@ivdc.chinacdc.cn

中心实验室

卫生部医学病毒学重点实验室

中国疾病预防控制中心病毒病预防控制所

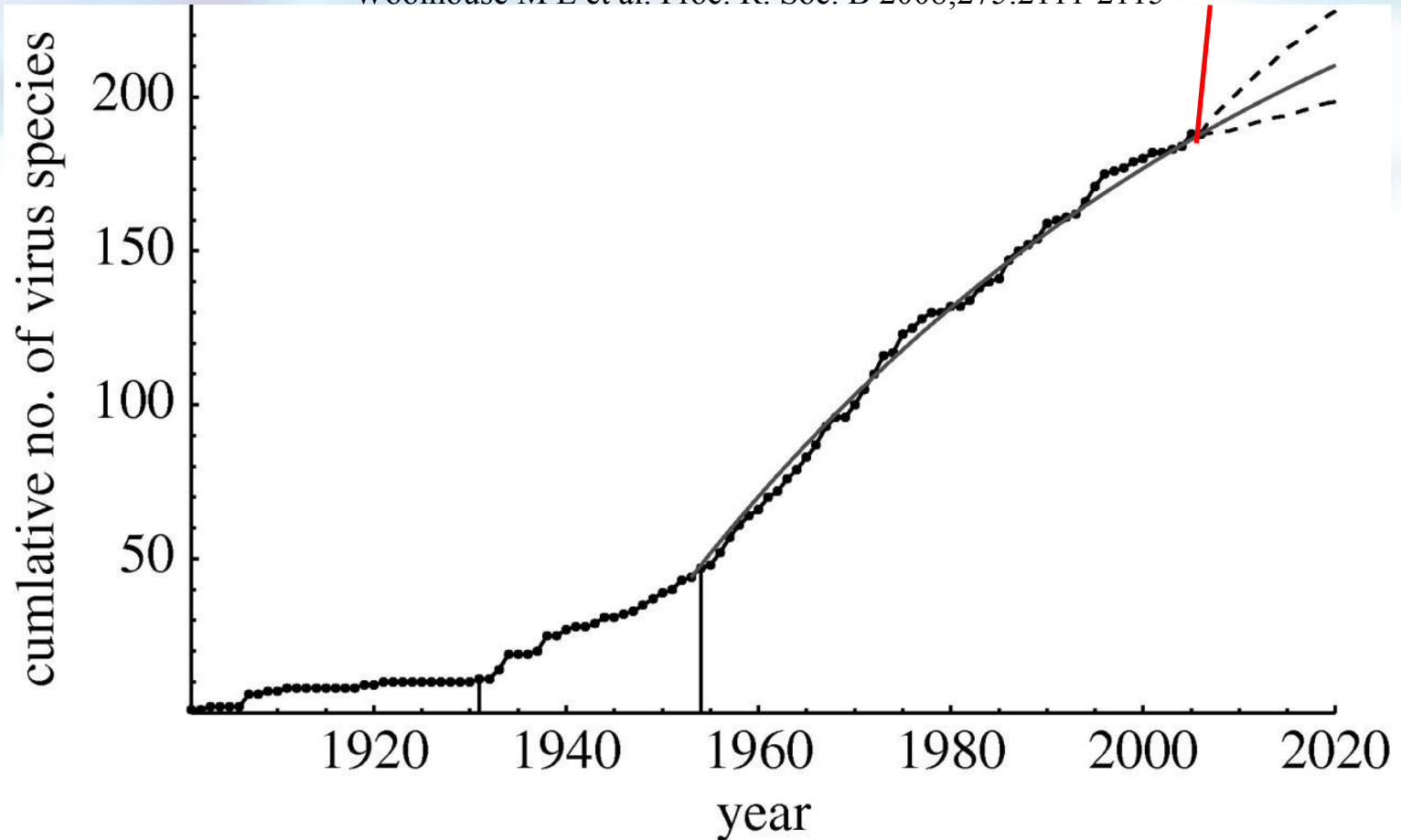


病毒与下一代测序 (NGS)

病毒鉴定自动化分析流程
(Virus Identification Pipeline)

可感染人病毒的发现时间曲线

Woolhouse M E et al. Proc. R. Soc. B 2008;275:2111-2115



随着越来越多的病毒被发现了，其中可能致病的病毒种类也随之增加。

检测技术概述

- **形态学（电镜）**
- **组织学（细胞培养）**
- **免疫学（ELISA, Western Blot, IF, NT, HI）**
- **分子生物学（PCR, RT-PCR, qPCR, mPCR, Microarray, Sequencing……）**

为什么？

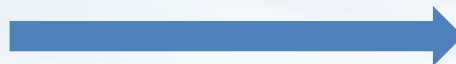
- 病毒变异
- 含量低，背景噪声大
- 新发疾病

接下来呢？

- 分离培养
- 高通量测序

临床样本

高通量测序

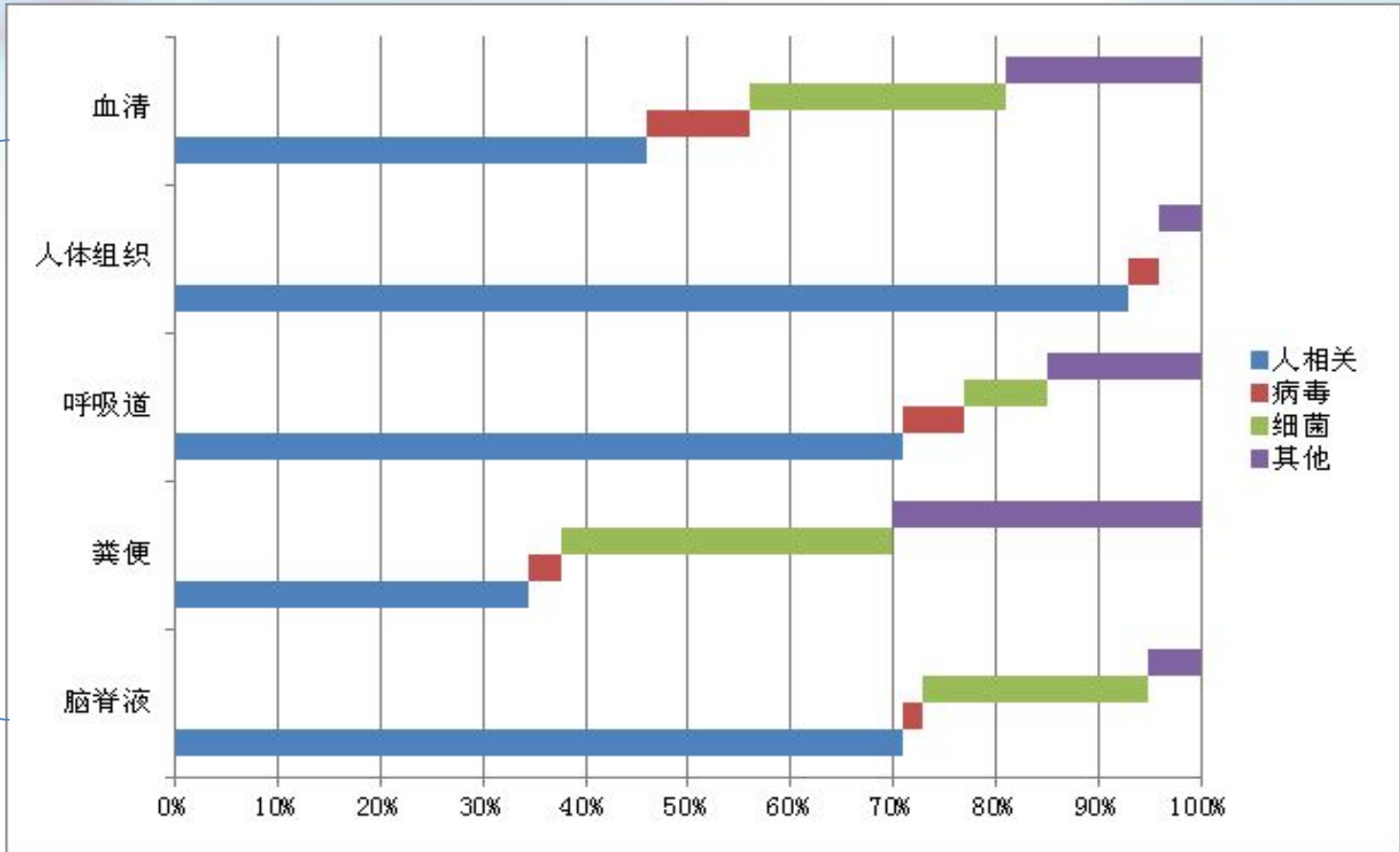


临床宏基因组

(Clinical Metagenomics)

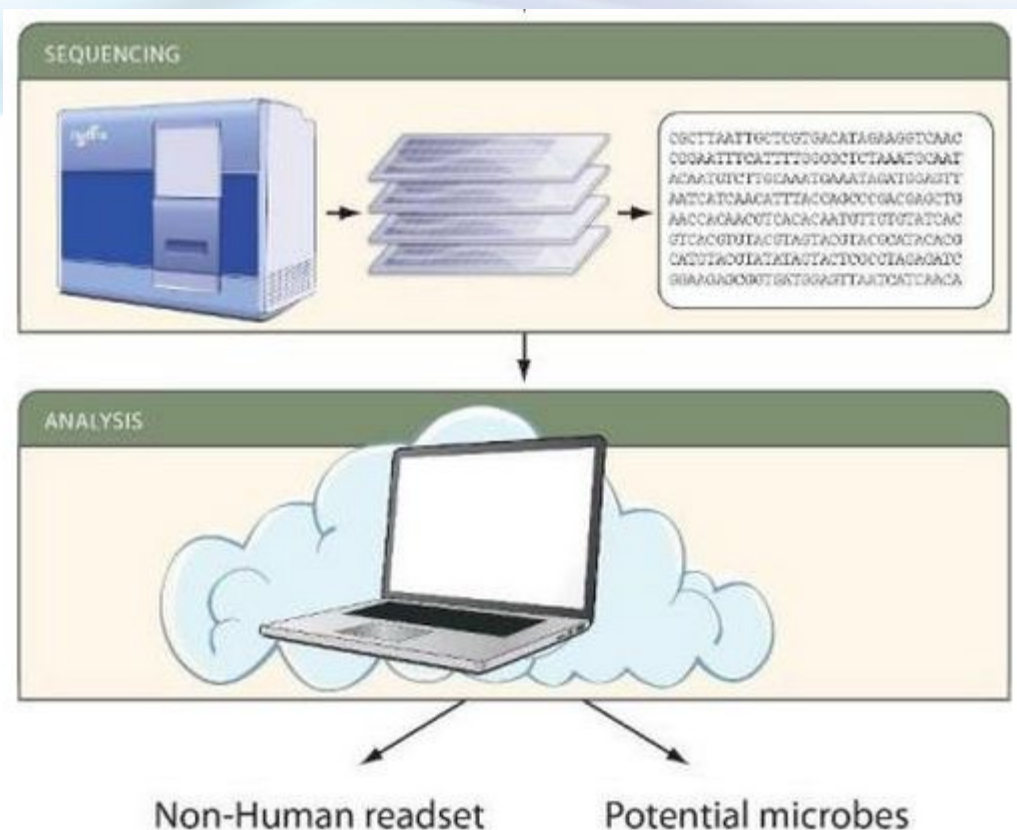
可以在不需要目标病原体的背景资料的情况下，
一次高通量测序试验就可以完成检测工作。

临床样本



病毒检测如大海捞针

主要难点



- 比对算法/分类算法
(BLAST)
- 参考数据库
(NT/NR)
- 结果信息?

Table 1. Summary of viruses identified in this study.

Patient code	Clinic virus ID	Virochip virus ID	Sequencing virus ID	Virus TaxID ^a	# virus reads	# initial reads	Fraction virus reads
187	DENV-2	DENV-2	Dengue virus 2	11060	4280	1.1E+06	3.9E−03
275	DENV-2	DENV-2	Dengue virus 2	11060	1511	1.6E+06	9.7E−04
282	DENV-2	DENV-2	Dengue virus 2	11060	699	1.6E+06	4.2E−04
266	DENV-2	DENV-2	Dengue virus 2	11060	135749	4.8E+06	2.8E−02
274	DENV-1	DENV-1	Dengue virus 1	11053	27	1.2E+06	2.3E−05
401 ^b	HAV	HAV	Hepatitis A virus	12092	2164	1.8E+05	1.2E−02
401 ^b	HAV	HAV	Hepatitis A virus	12092	4562	1.3E+06	3.5E−03
235	-	-	Human herpesvirus 6	10368	116	5.5E+06	2.1E−05
451	-	-	Human herpesvirus 6	10368	88	2.7E+06	3.2E−05
207	-	-	Human herpesvirus 6	10368	390	9.6E+06	4.1E−05
432	-	-	Human herpesvirus 6	10368	411	3.5E+06	1.2E−04
574	-	-	Human herpesvirus 6	10368	138	3.2E+06	4.4E−05
370	-	-	Human herpesvirus 6	10368	90	3.2E+06	2.9E−05
78	-	-	Human herpesvirus 6	10368	113	1.2E+06	9.8E−05

Yozwiak, N.L., Skewes-Cox, P., Stenglein, M.D., Balmaseda, A., Harris, E. and DeRisi, J.L. (2012) Virus identification in unknown tropical febrile illness cases using deep sequencing. *PLoS Negl Trop Dis*, **6**, e1485.

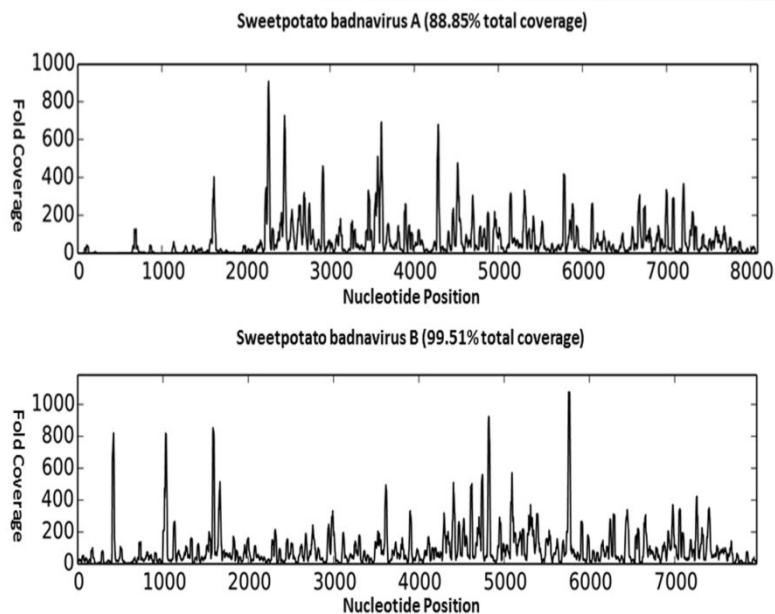


病毒与下一代测序 (NGS)

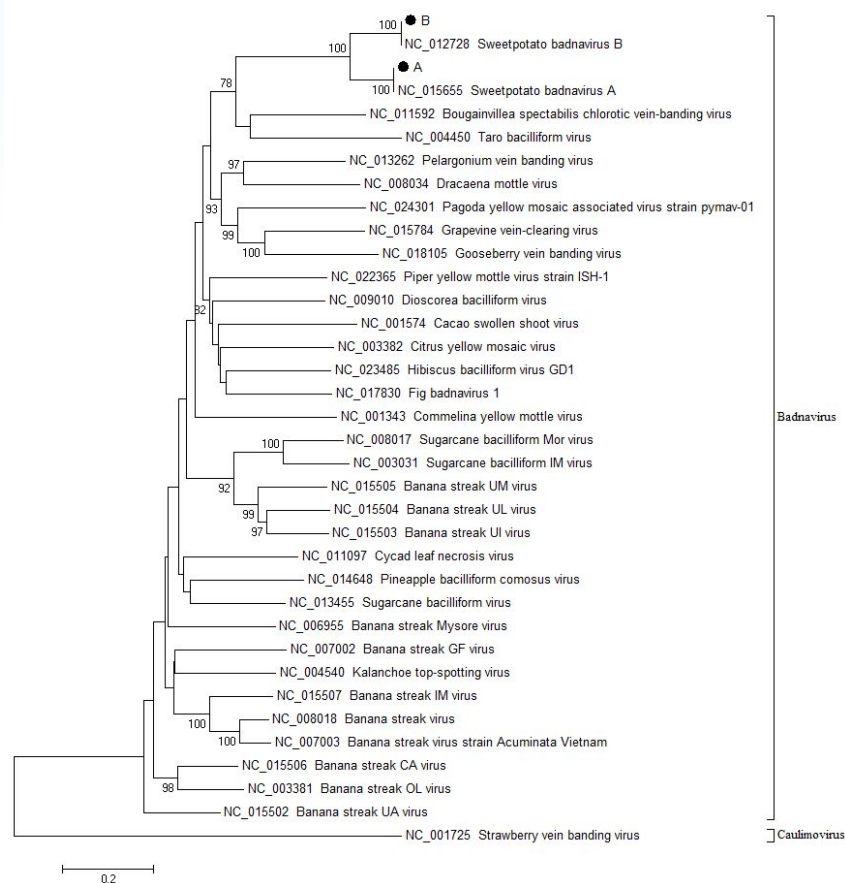
病毒鉴定自动化分析流程
(Virus Identification Pipeline)

什么样的分析结果

基因组覆盖度



聚类分析

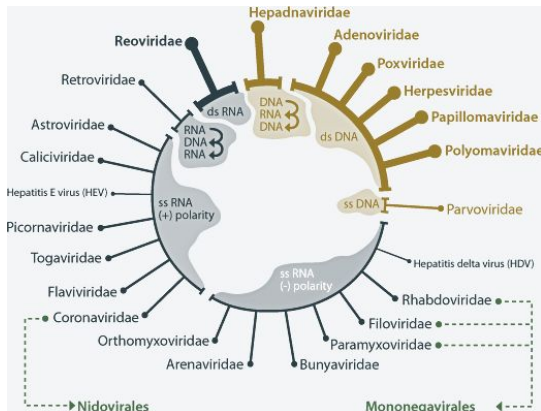


汇总报告

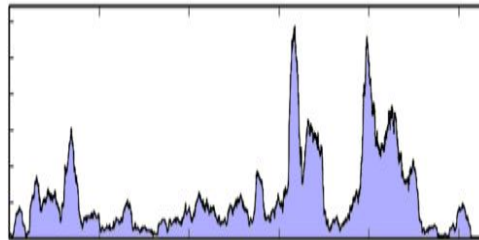


Virus Identification Pipeline (VIP)

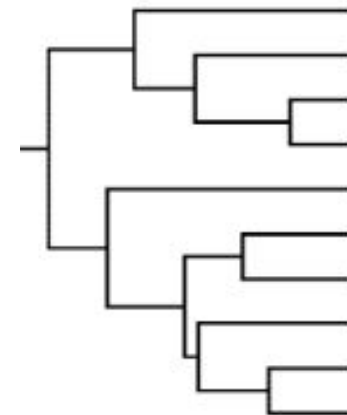
分类学



基因组覆盖度



聚类分析



高准确度分类检测

SRA	sample type	literature report	VIP_result	Coverage%	Reads
SRR453448	serum	Dengue virus 2	Dengue virus 2	99.14	64778
SRR453458	serum	Hepatitis A virus	Hepatitis A virus	57.15	2315
SRR1106123	serum	Hepatitis C virus,	Hepatitis C virus	100.00	108274
		Hepatitis G virus	Hepatitis G virus	100.00	751691
SRR1106548	plasma	HIV	HIV	96.61	29528
SRR1106553	nasal swab	H1N1	H1N1	98.25	12613
			Sapovirus	96.16	13019
			Rotavirus A	92.41	786
			Adeno-associated virus 2	80.74	1836
			Human parechovirus 6	76.36	4467
SRR1106550	stool	Sapovirus, Rotavirus A	Torque teno mini virus 5	42.09	270

特异性和灵敏度结果

SRR1106553	BLAST	VIP		H1N1 (JF915184 - JF915191)
		Tru	Neg	
	Tru	51,509.0	89.0	
	Neg	15.0	0.0	

SRR1106548	BLAST	VIP		HIV (AF063223.1)
		Tru	Neg	
	Tru	30,108.0	51.0	
	Neg	452.0	0.0	

SRR1170797	BLAST	VIP		BVDV (JN380086.1)
		Tru	Neg	
	Tru	35,554.0	1,087.0	
	Neg	65.0	0.0	

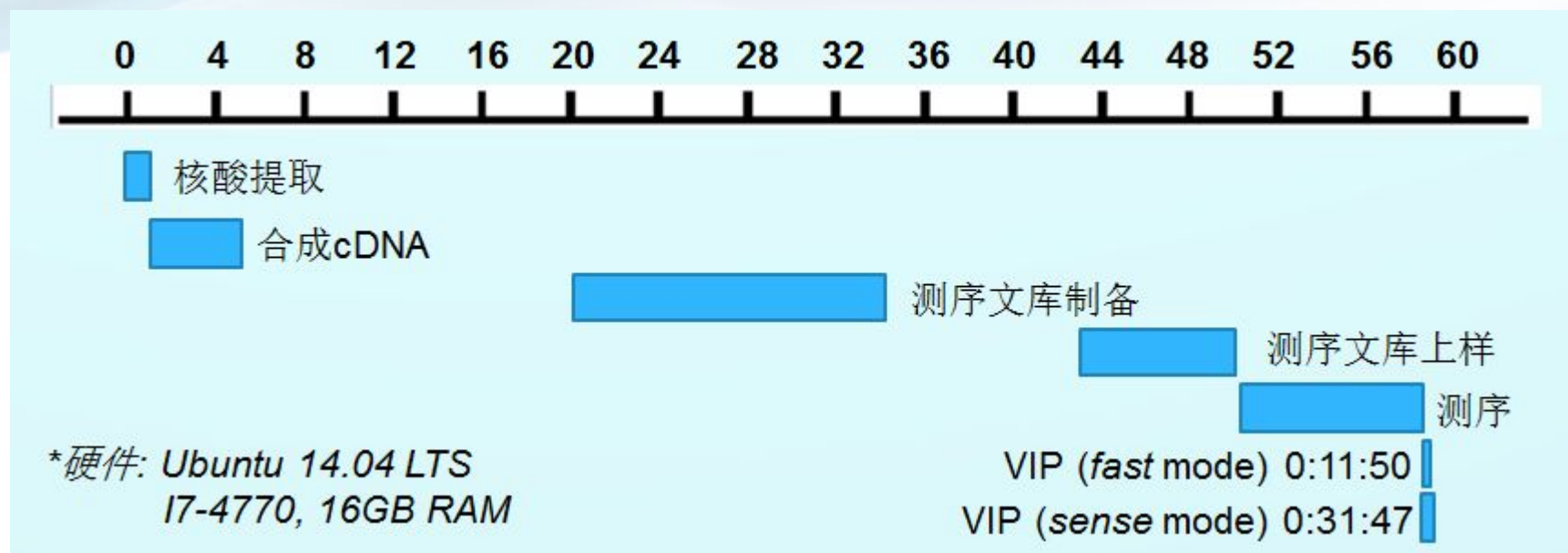
特异性和灵敏度结果

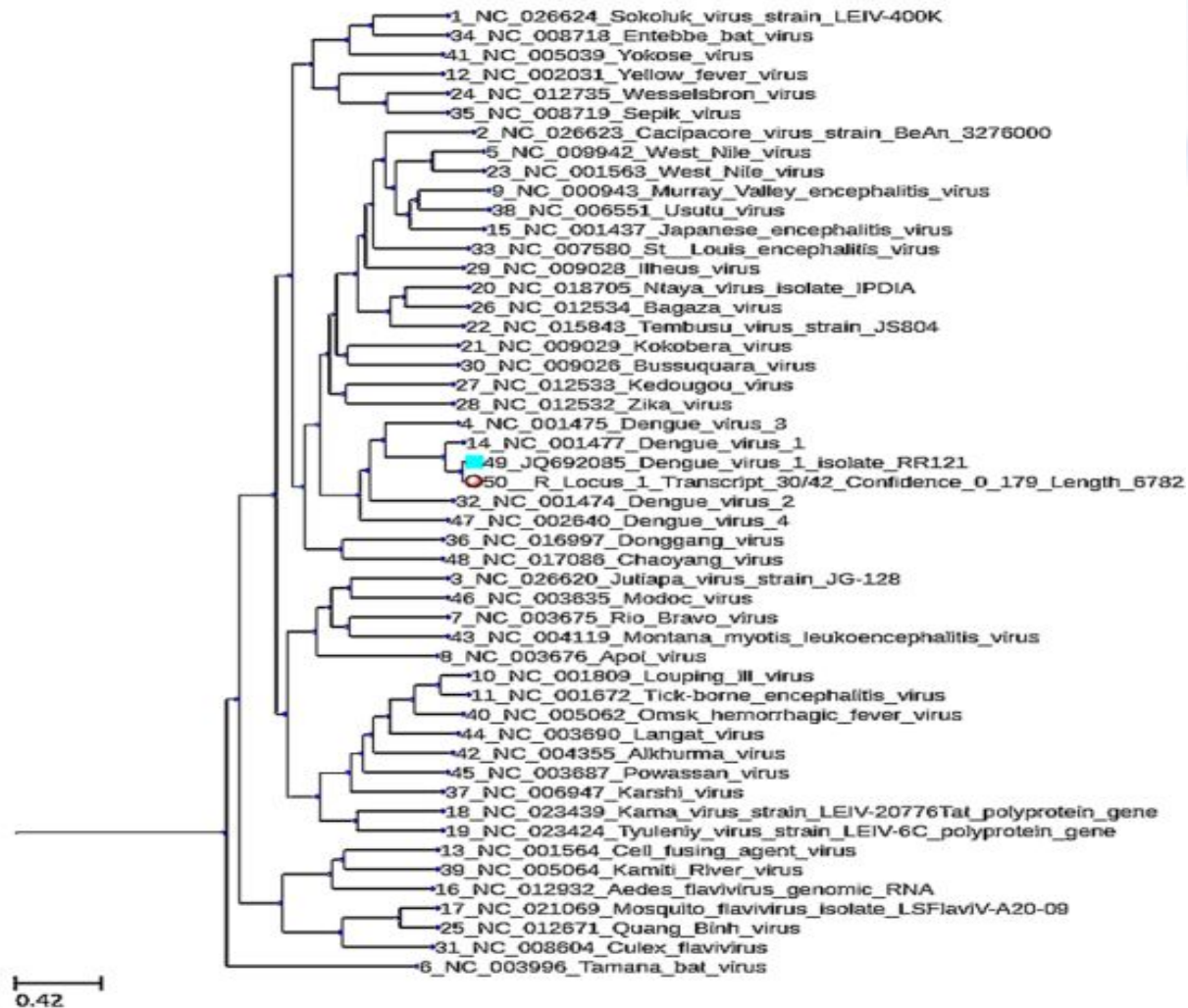
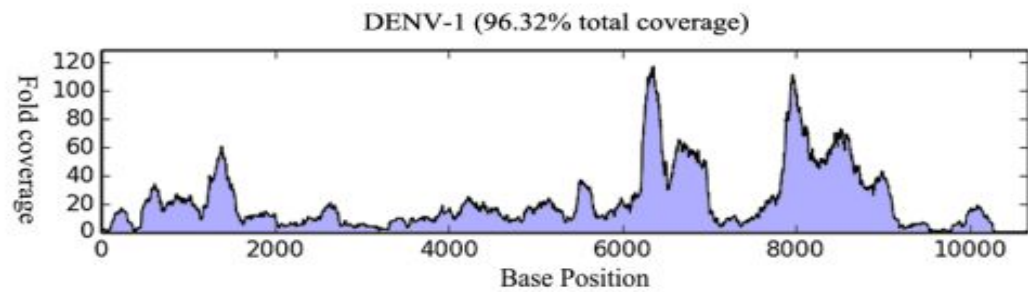
$$\text{TPR} = \text{TP}/(\text{TP} + \text{FN})$$

$$\text{FPR} = \text{FP}/(\text{FP} + \text{TN})$$

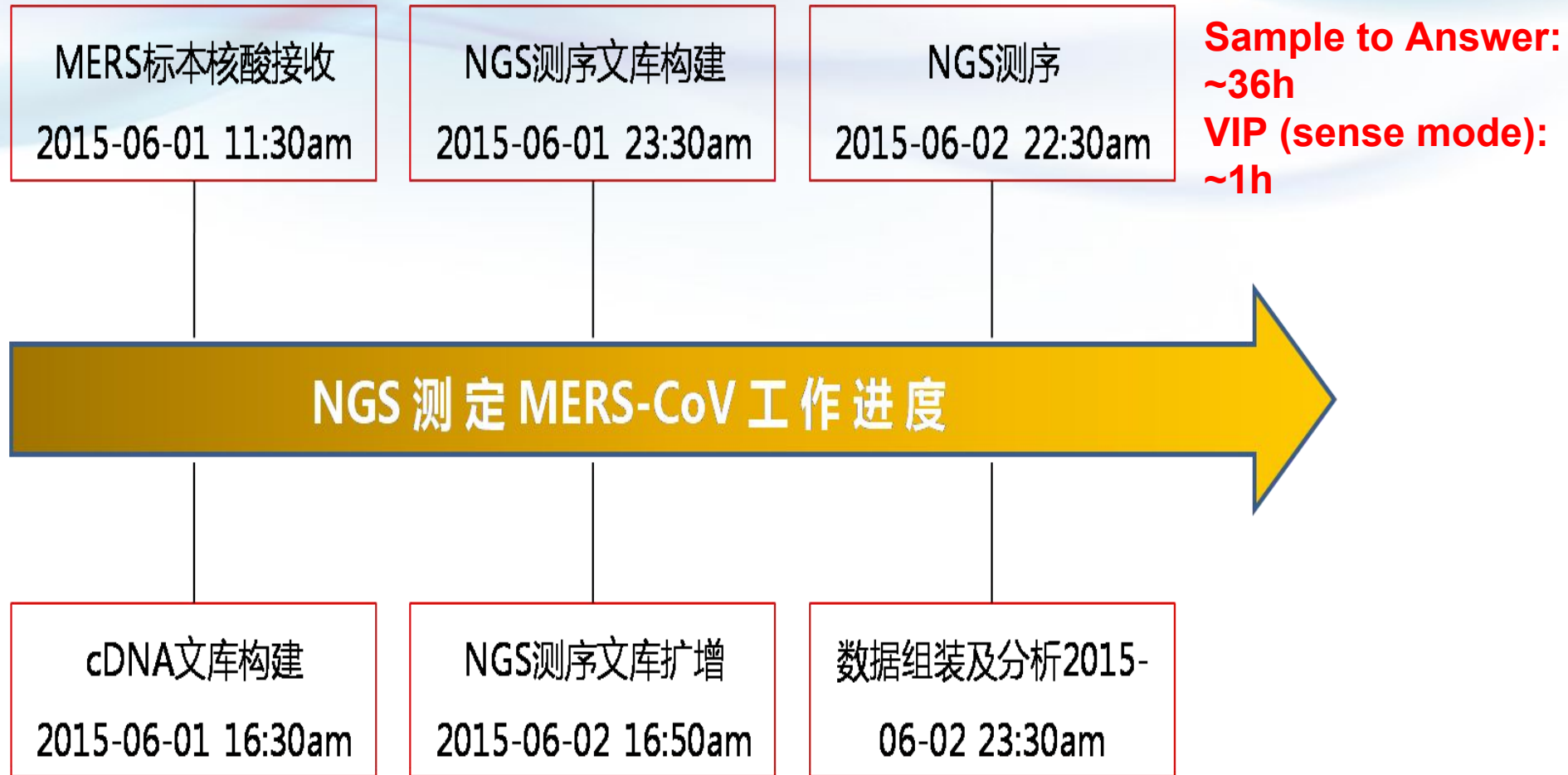
SRA	Viruses	% True Positive Rate	% False Positive Rate
SRR1170797	BVDV	97.03	100.00
SRR1106548	HIV	99.83	100.00
SRR1106553	H1N1	99.82	100.00

“十二五”考核





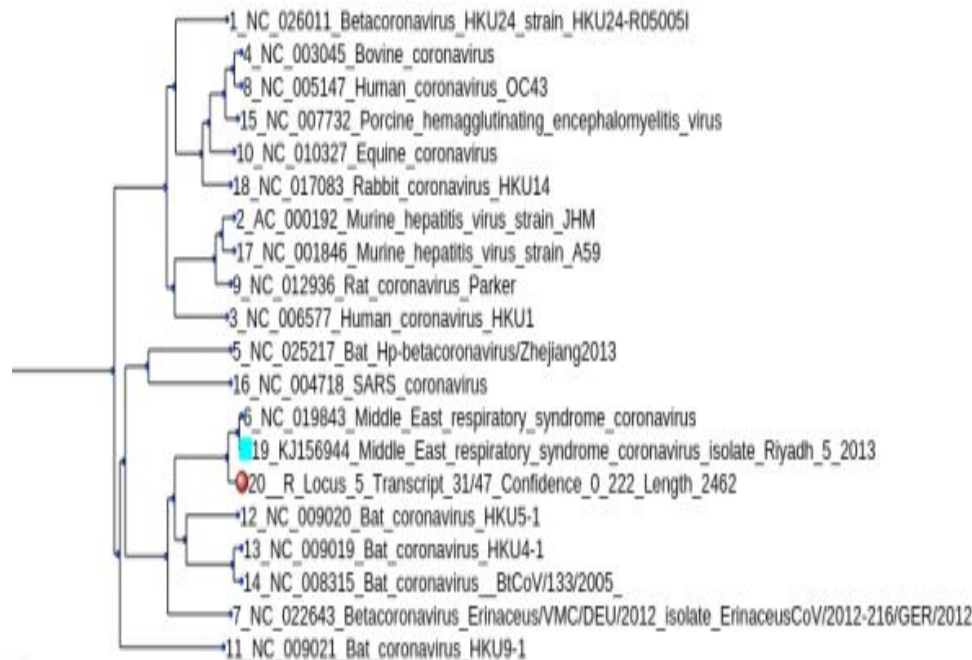
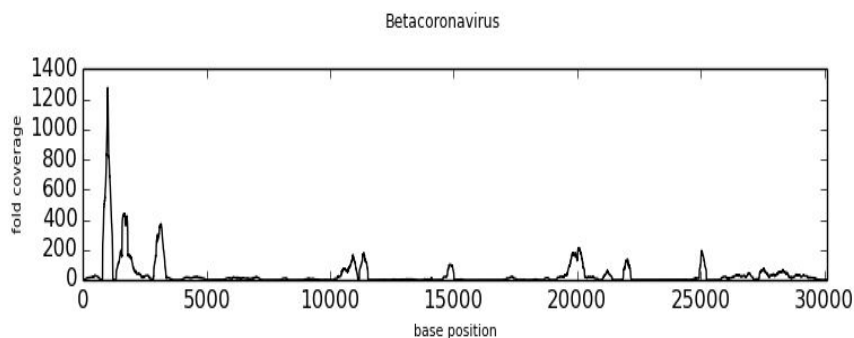
MERS应急



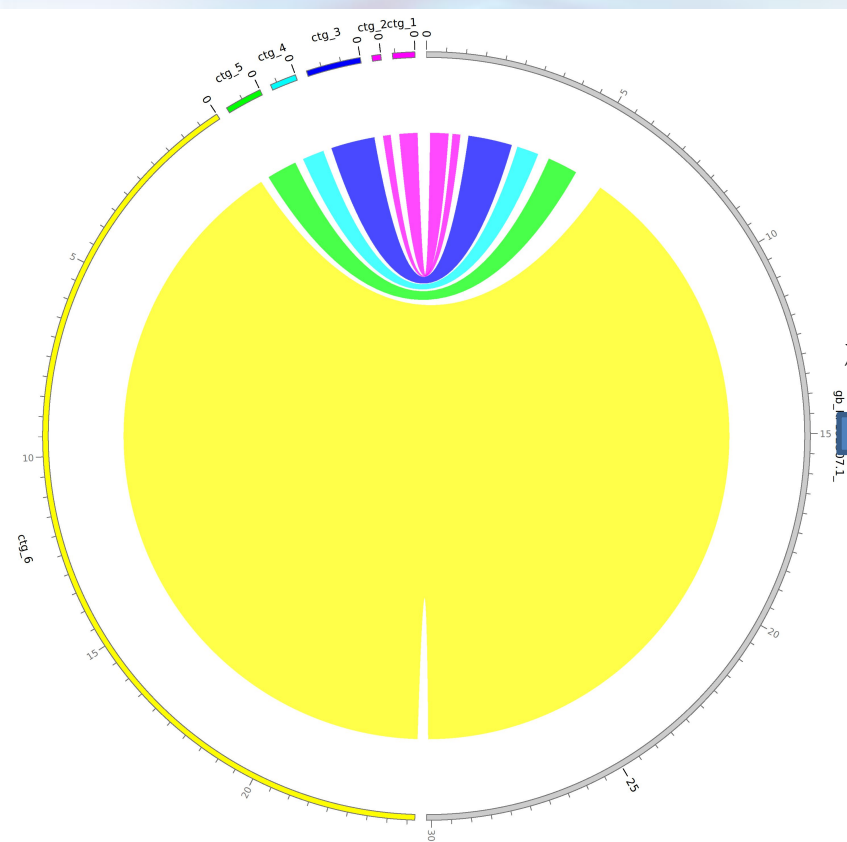
NGS Timeline

MERS应急

Species	Genus	GI	%Cove rage	Read s_hit	Reads_ num	Average depth of coverage
Middle East respiratory syndrome coronavirus, complete s genome	Betacoronaviru s	51126266 2	67.95	5,980	6,340	3,506.33
基因组覆盖度(67.95%)			进化分析			

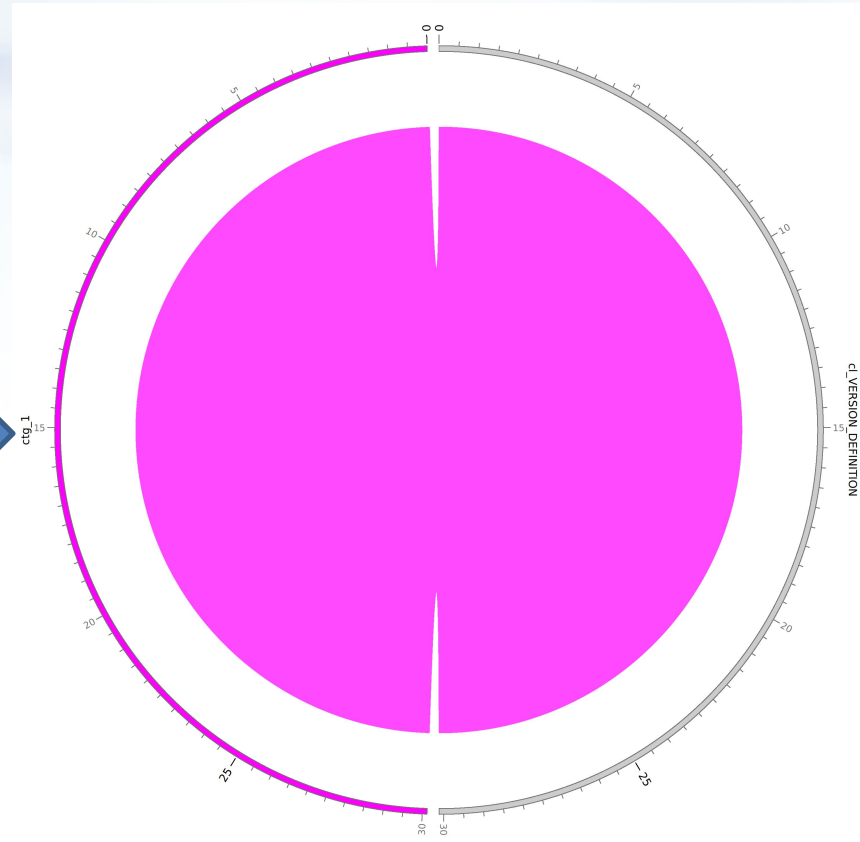


MERS GD01全基因组绘制



Draft genome
(95.8% coverage)

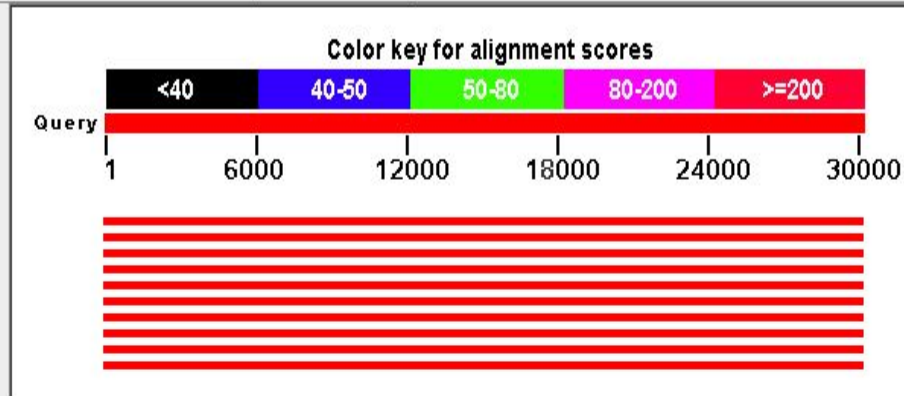
校准



Complete genome
(100.0% coverage)

Distribution of 10 Blast Hits on the Query Sequence

Mouse over to see the define, click to show alignments



Descriptions

Sequences producing significant alignments:

Select: [All](#) [None](#) Selected: 0

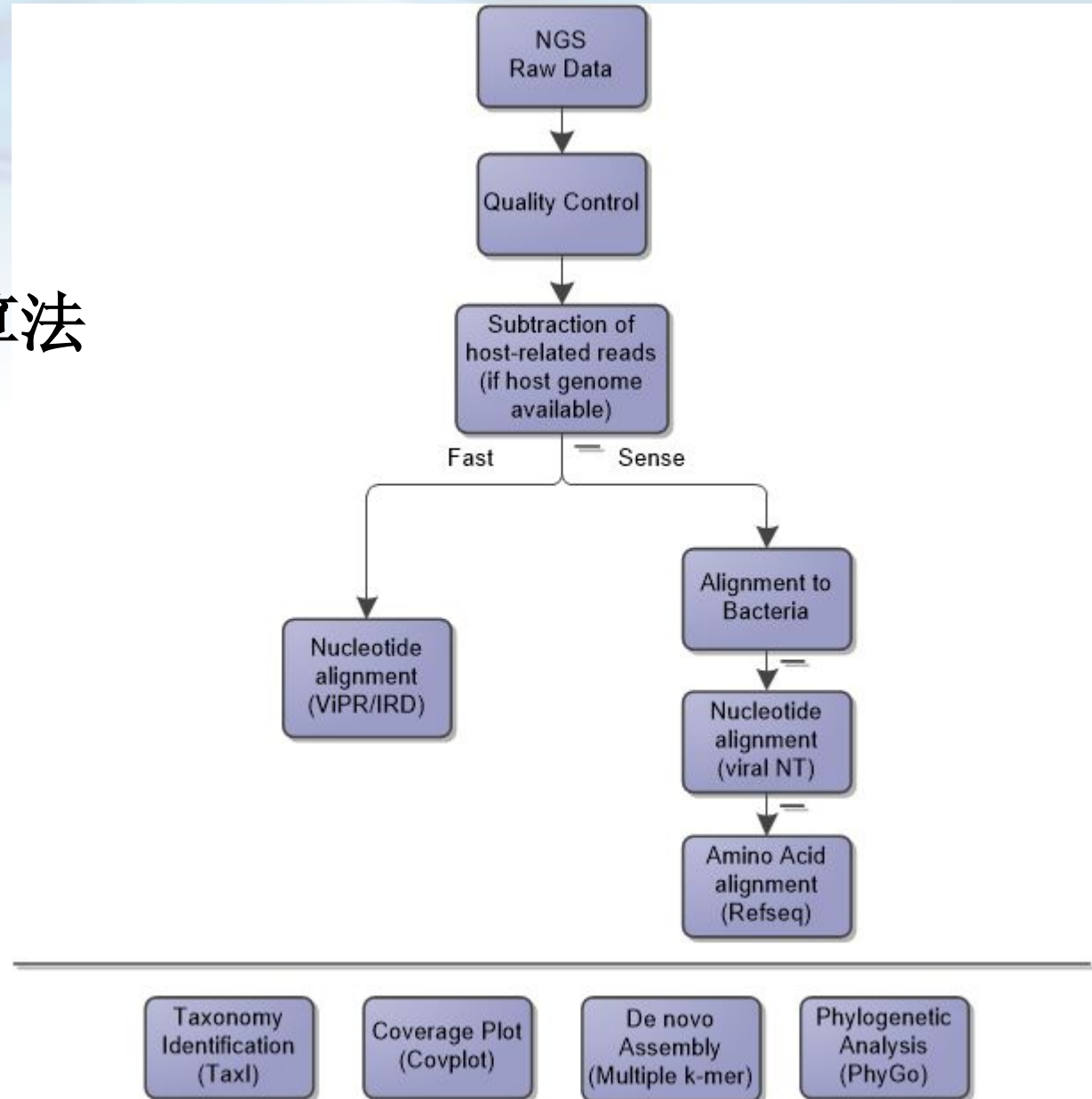
[Alignments](#) [Download](#) [GenBank](#) [Graphics](#) [Distance tree of results](#)



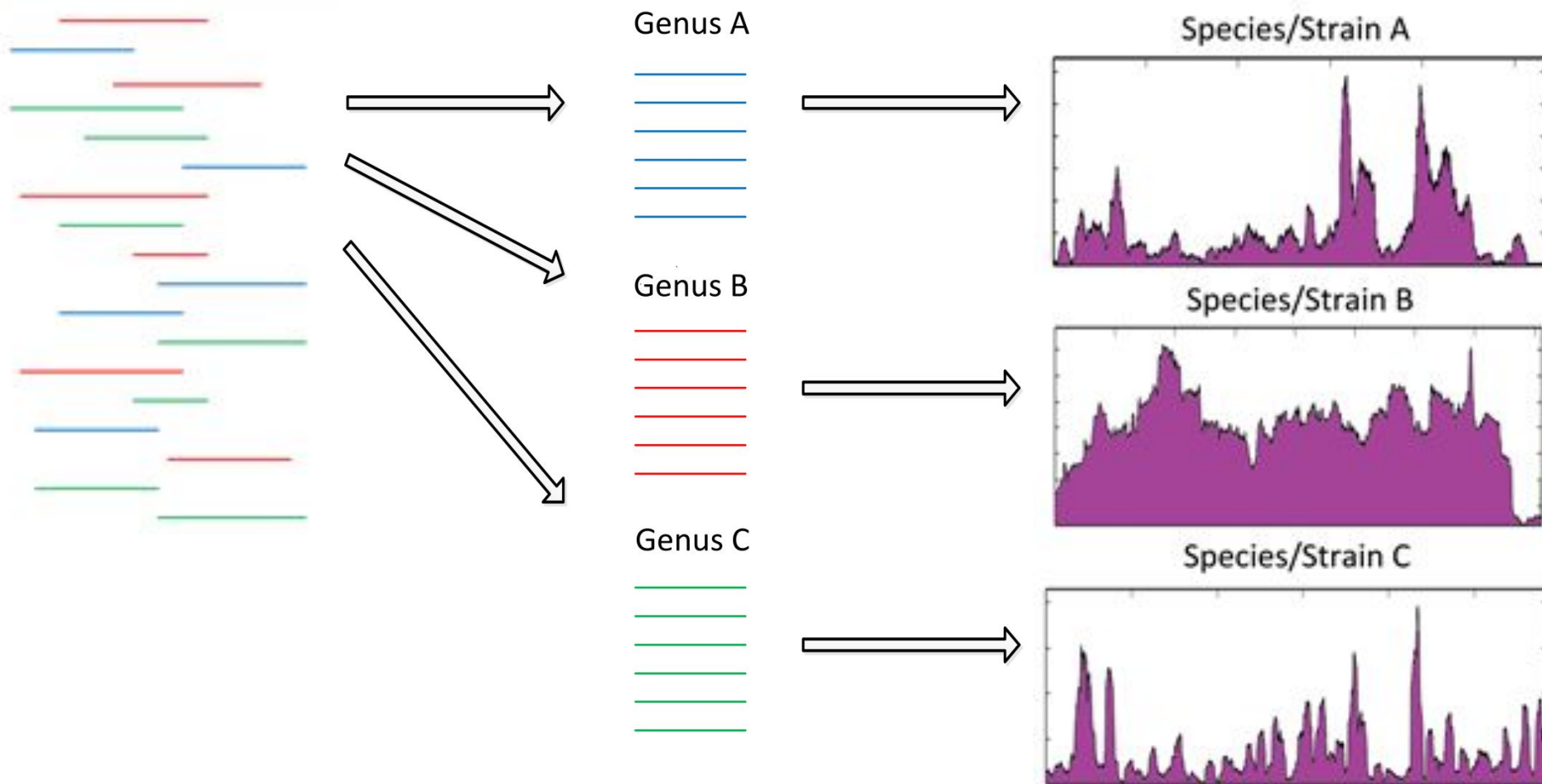
	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	Middle East respiratory syndrome coronavirus isolate Hafr-Al-Batin 1 2013, complete genome	55238	55238	99%	0.0	99%	KF600628.1
<input type="checkbox"/>	Middle East respiratory syndrome coronavirus isolate Camel/Qatar 2 2014, complete genome	55197	55197	99%	0.0	99%	KJ650098.1
<input type="checkbox"/>	Middle East respiratory syndrome coronavirus strain Florida/USA-2 Saudi Arabia 2014, complete genome	55184	55184	100%	0.0	99%	KJ829365.1
<input type="checkbox"/>	Middle East respiratory syndrome coronavirus isolate Qatar4, complete genome	55182	55182	99%	0.0	99%	KF961222.1
<input type="checkbox"/>	Middle East respiratory syndrome coronavirus isolate Florida/USA-2 Saudi Arabia 2014, complete genome	55173	55173	100%	0.0	99%	KP223131.1

VIP流程示意图

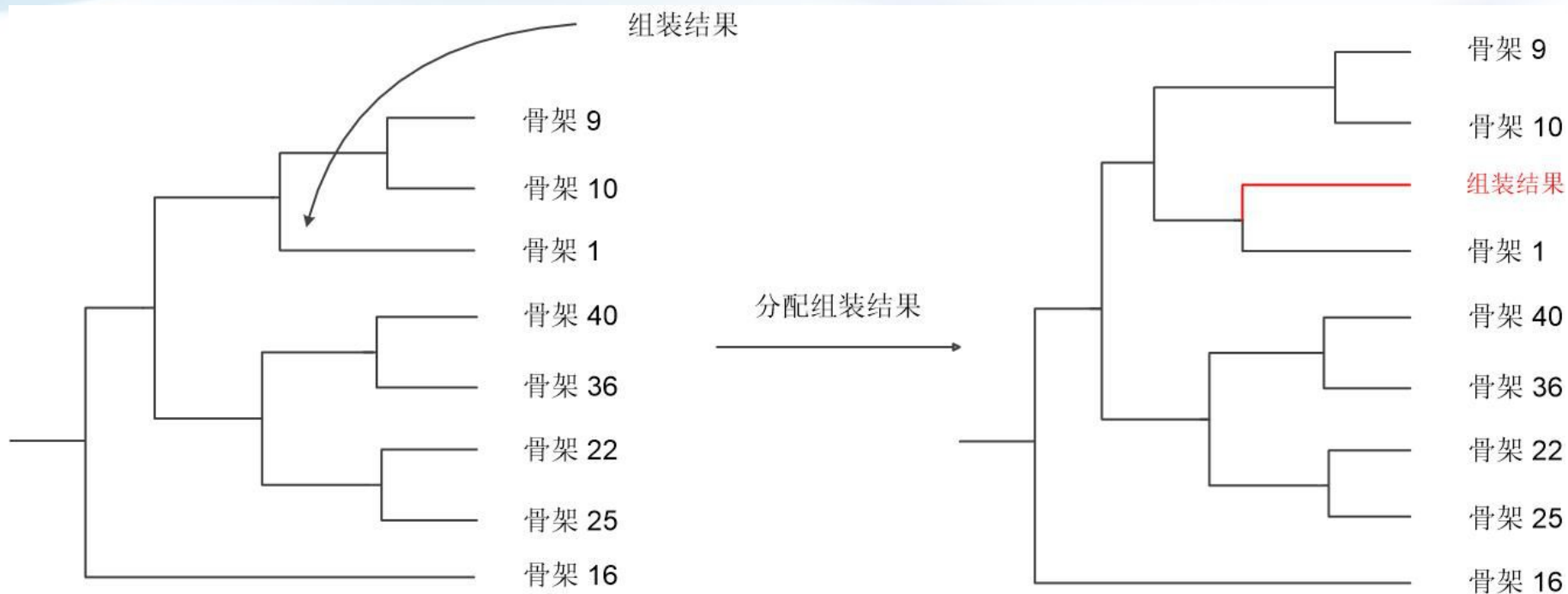
- 比对算法/分类算法
- 参考数据库
- 结果信息？



分类算法示意图



进化算法示意图



总结



中心实验室：坚持是一种美德





Thank
you