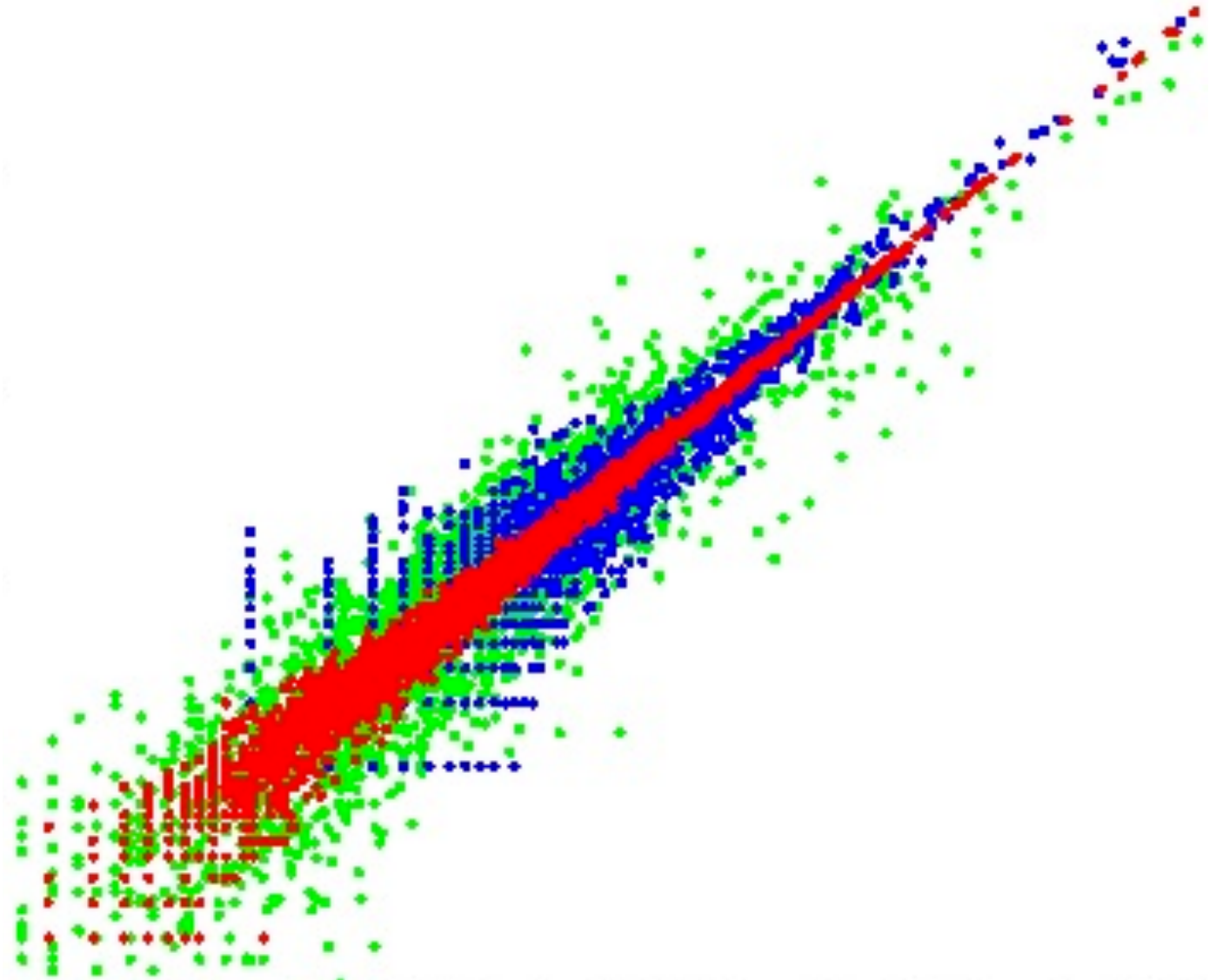


# Experimental Design

---



**The datasets might be larger, but basic science principles still apply...**

- Controls
- Replication
- Good experimental design

# The Importance of Controls in NGS Experiments

PNAS

## Hybrid DNA virus in Chinese patients with seronegative hepatitis discovered by deep sequencing

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

<sup>a</sup>Hematology Branch and <sup>c</sup>Systems Biology Center, National Heart, Lung, and Blood Institute, Bethesda, MD 20892; <sup>b</sup>Institute of Infectious Disease, Southwest Hospital, Third Military Medical University, Chongqing 400038, China; <sup>d</sup>Department of Embryology, Carnegie Institution for Science, Baltimore, MD 21218

Edited\* by Harvey Alter, National Institutes of Health, Bethesda, MD, and approved March 19, 2013 (received for review March 4, 2013)



## The Perils of Pathogen Discovery: Origin of a Novel Parvovirus-Like Hybrid Genome Traced to Nucleic Acid Extraction Spin Columns

Samia N. Naccache,<sup>a,b</sup> Alexander L. Greninger,<sup>a,b</sup> Deanna Lee,<sup>a,b</sup> Lark L. Coffey,<sup>c</sup> Tung Phan,<sup>c</sup> Annie Rein-Weston,<sup>a,b</sup> Andrew Aronsohn,<sup>d</sup> John Hackett, Jr.,<sup>e</sup> Eric L. Delwart,<sup>a,c</sup> Charles Y. Chiu<sup>a,b,f</sup>

Department of Laboratory Medicine, University of California, San Francisco, California, USA<sup>a</sup>; UCSF-Abbott Viral Diagnostics and Discovery Center, San Francisco, California, USA<sup>b</sup>; Blood Systems Research Institute, San Francisco, California, USA<sup>c</sup>; Center for Liver Disease, University of Chicago Medical Center, Chicago, Illinois, USA<sup>d</sup>; Abbott Diagnostics, Abbott Park, Illinois, USA<sup>e</sup>; Department of Medicine, Division of Infectious Diseases, University of California, San Francisco, California, USA<sup>f</sup>

# Reagents as a Source of Contamination

TABLE 1 PCR screening of commonly used viral nucleic acid extraction kits for parvovirus-like hybrid virus (PHV-1)<sup>a</sup>

Kit	Spin column	PCR result for:							
		Replicase, nt763-1010 (248 nt)		Bridge, nt1554-2044 (491 nt)		Capsid, nt1922-2044 (121 nt)		Capsid + NCR, nt3288-3448 (161 nt)	
		C	F	C	F	C	F	C	F
RNeasy MinElute cleanup kit	RNeasy MinElute column	+	+	-	+	+	+	+	+
RNeasy minikit	RNeasy minicolumn	+	+	+	+	+	+	+	+
QIAamp UltraSens virus kit	QIAamp minicolumn	+	+	-	+	+	+	+	+
QIAamp viral RNA minikit	QIAamp minicolumn	-	+	-	-	+	+	+	+
QIAamp DSP virus kit	QIAamp MinElute column	-	+	-	-	-	+	-	+
PureLink viral RNA/DNA minikit	PureLink viral column	-	-	-	-	-	-	-	-
TRIZOL LS kit	NA	-	-	-	-	-	-	-	-
EZ1 viral minikit v2.0	NA	-	-	-	-	-	-	-	-
Water, nuclease-free (Qiagen, Fisher Scientific, and Epicentre)	NA	-	-	-	-	-	-	-	-

<sup>a</sup> NCR, noncoding region; C, column elution; F, full extraction; nt, nucleotide; NA, not applicable.

# Reagents as a Source of Contamination

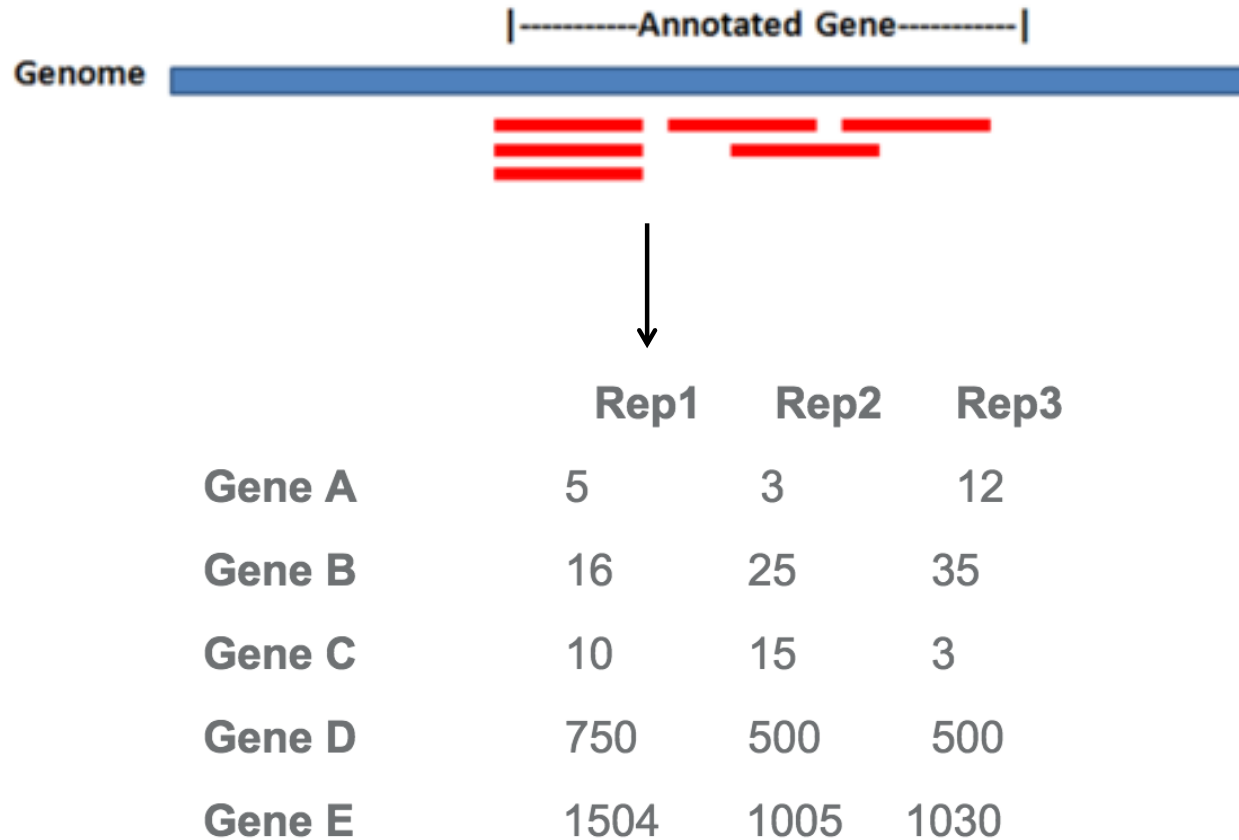
Contamination also prevalent in 16S microbiome studies

**Table 1 List of contaminant genera detected in sequenced negative 'blank' controls**

Phylum	List of constituent contaminant genera
Proteobacteria	<p>Alpha-proteobacteria:</p> <p><i>Afipia</i>, <i>Aquabacterium</i><sup>e</sup>, <i>Asticcacaulis</i>, <i>Aurantimonas</i>, <i>Beijerinckia</i>, <i>Bosea</i>, <i>Bradyrhizobium</i><sup>d</sup>, <i>Brevundimonas</i><sup>c</sup>, <i>Caulobacter</i>, <i>Craurococcus</i>, <i>Devosia</i>, <i>Hoeflea</i><sup>e</sup>, <i>Mesorhizobium</i>, <i>Methylobacterium</i><sup>c</sup>, <i>Novosphingobium</i>, <i>Ochrobactrum</i>, <i>Paracoccus</i>, <i>Pedomicrobium</i>, <i>Phyllobacterium</i><sup>e</sup>, <i>Rhizobium</i><sup>c,d</sup>, <i>Roseomonas</i>, <i>Sphingobium</i>, <i>Sphingomonas</i><sup>c,d,e</sup>, <i>Sphingopyxis</i></p> <p>Beta-proteobacteria:</p> <p><i>Acidovorax</i><sup>c,e</sup>, <i>Azoarcus</i><sup>e</sup>, <i>Azospira</i>, <i>Burkholderia</i><sup>d</sup>, <i>Comamonas</i><sup>c</sup>, <i>Cupriavidus</i><sup>c</sup>, <i>Curvibacter</i>, <i>Delftia</i><sup>e</sup>, <i>Duganella</i><sup>a</sup>, <i>Herbaspirillum</i><sup>a,c</sup>, <i>Janthinobacterium</i><sup>e</sup>, <i>Kingella</i>, <i>Leptothrix</i><sup>a</sup>, <i>Limnobacter</i><sup>e</sup>, <i>Massilia</i><sup>c</sup>, <i>Methylophilus</i>, <i>Methyloversatilis</i><sup>e</sup>, <i>Oxalobacter</i>, <i>Pelomonas</i>, <i>Polaromonas</i><sup>e</sup>, <i>Ralstonia</i><sup>b,c,d,e</sup>, <i>Schlegelella</i>, <i>Sulfuritalea</i>, <i>Undibacterium</i><sup>e</sup>, <i>Variovorax</i></p> <p>Gamma-proteobacteria:</p> <p><i>Acinetobacter</i><sup>a,d,c</sup>, <i>Enhydrobacter</i>, <i>Enterobacter</i>, <i>Escherichia</i><sup>a,c,d,e</sup>, <i>Nevskia</i><sup>e</sup>, <i>Pseudomonas</i><sup>b,d,e</sup>, <i>Pseudoxanthomonas</i>, <i>Psychrobacter</i>, <i>Stenotrophomonas</i><sup>a,b,c,d,e</sup>, <i>Xanthomonas</i><sup>b</sup></p>
Actinobacteria	<i>Aeromicrobium</i> , <i>Arthrobacter</i> , <i>Beutenbergia</i> , <i>Brevibacterium</i> , <i>Corynebacterium</i> , <i>Curtobacterium</i> , <i>Dietzia</i> , <i>Geodermatophilus</i> , <i>Janibacter</i> , <i>Kocuria</i> , <i>Microbacterium</i> , <i>Micrococcus</i> , <i>Microlunatus</i> , <i>Patulibacter</i> , <i>Propionibacterium</i> <sup>e</sup> , <i>Rhodococcus</i> , <i>Tsukamurella</i>
Firmicutes	<i>Abiotrophia</i> , <i>Bacillus</i> <sup>b</sup> , <i>Brevibacillus</i> , <i>Brochothrix</i> , <i>Facklamia</i> , <i>Paenibacillus</i> , <i>Streptococcus</i>
Bacteroidetes	<i>Chryseobacterium</i> , <i>Dyadobacter</i> , <i>Flavobacterium</i> <sup>d</sup> , <i>Hydrothalea</i> , <i>Niastella</i> , <i>Olivibacter</i> , <i>Pedobacter</i> , <i>Wautersiella</i>
Deinococcus-Thermus	<i>Deinococcus</i>
Acidobacteria	Predominantly unclassified Acidobacteria Gp2 organisms

The listed genera were all detected in sequenced negative controls that were processed alongside human-derived samples in our laboratories (WTSI, ICL and UB) over a period of four years. A variety of DNA extraction and PCR kits were used over this period, although DNA was primarily extracted using the FastDNA SPIN

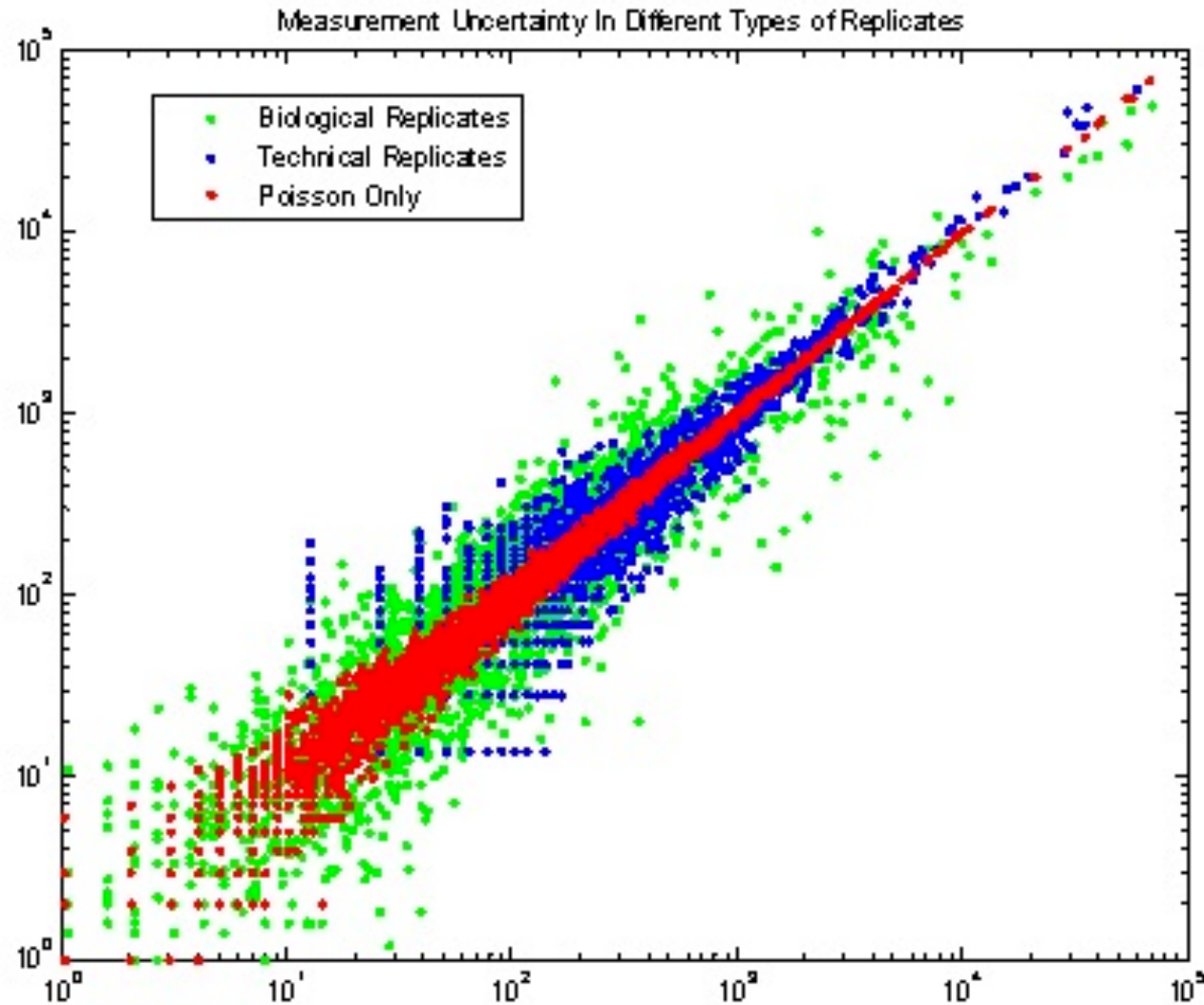
# Sources of Variance in NGS Experiments



1. Poisson sampling variance
2. Technical variation introduced during library construction and sequence
3. Biological variation between samples

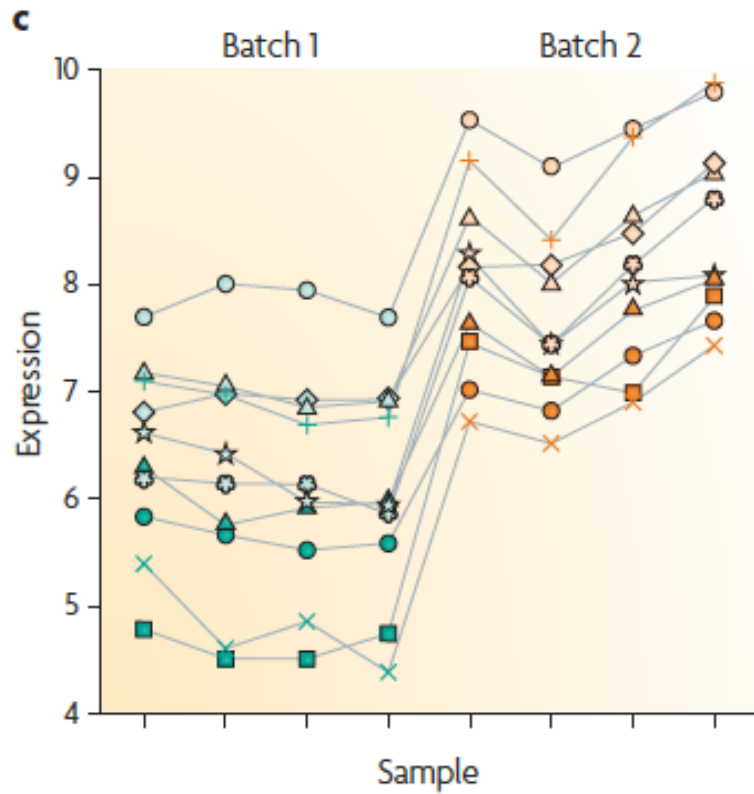
Variation in Read Counts among Replicates

# Sources of Variance in NGS Experiments

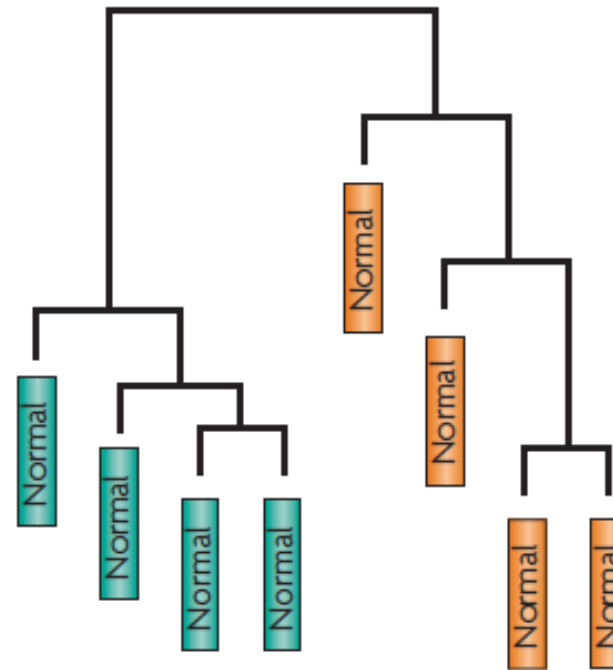


- Aim for a minimum of three biological replicates in “counting” experiments like differential-expression studies.
- Technical replicates do not address biological variance.
- Doing more replicates is often better than sequencing a small number of replicates more deeply.

# Batch Effects and Sources of Bias



**d**



- Reagent lots
- Technicians
- Library prep batches
- Sequencer runs/lanes



# Illumina's Recommendations for Reducing Index Hopping

**Table 1: Best Practices for Reducing Index Hopping**

Mitigation/Recommendation	Benefit/Outcome
Prepare dual indexed libraries with unique indexes <sup>a</sup>	Converts index hopped reads to undetermined
Sequence one 30x human genome per lane <sup>b</sup>	Avoids pooling and index hopping
Remove adapters (cleanup, spin columns, etc) <sup>c</sup>	Reduces levels of index hopping
Store prepared libraries at recommended temperature of $-20^{\circ}\text{C}$ <sup>c</sup>	Reduces levels of index hopping
Pool similar RNA-Seq samples together	Reduces contamination between high and low-expressors

Is this good scientific practice?

# Experimental Design Practice

## Study 1: Pathogen Discovery

You have observed a die-off of a species of frogs in a local lake and suspect that they may be experiencing an epidemic caused by a novel viral pathogen. You would like to use next generation sequencing to identify candidate viruses that may be responsible for this disease outbreak.

## Study 2: Differential Gene Expression

You are interested in how a bacterial infection alters gene expression in a species of shrimp that you study, and you have the ability to experimentally inoculate the shrimp and grow them in culture either with or without the bacterium.

### Describe the following features of your experimental design

- Sampling scheme, including plans for replication and/or controls
- Type(s) of nucleic acid to sample and any enrichment/depletion methods
- Sequencing platform and type of sequencing library
- Best practices to be used that will avoid batch effects, pseudoreplication, and artefacts