

miRNAs variants accuracy with sequencing platforms

Lorena Pantano, P.h. D

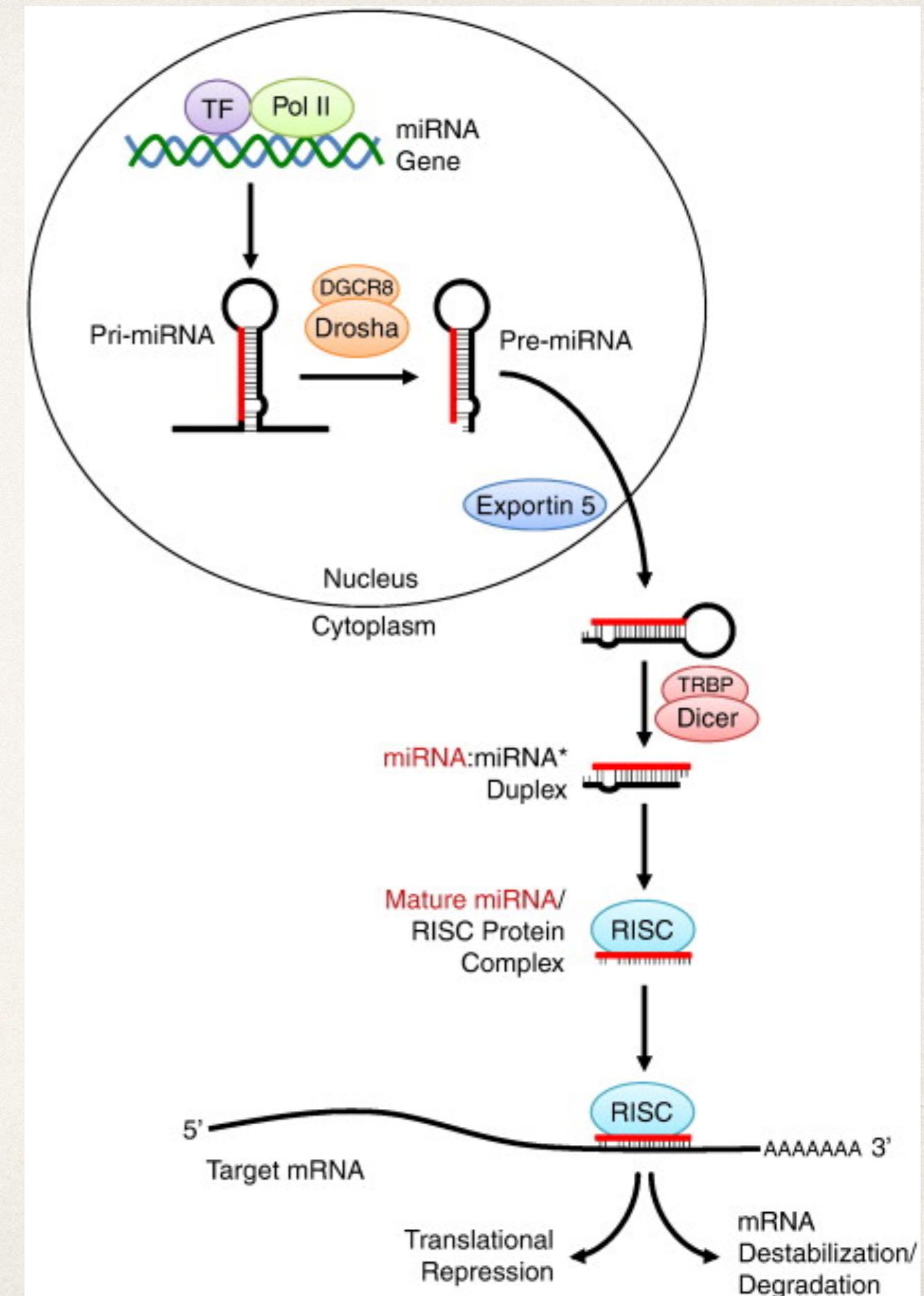
<https://lpantano.github.io>

03-18-2019

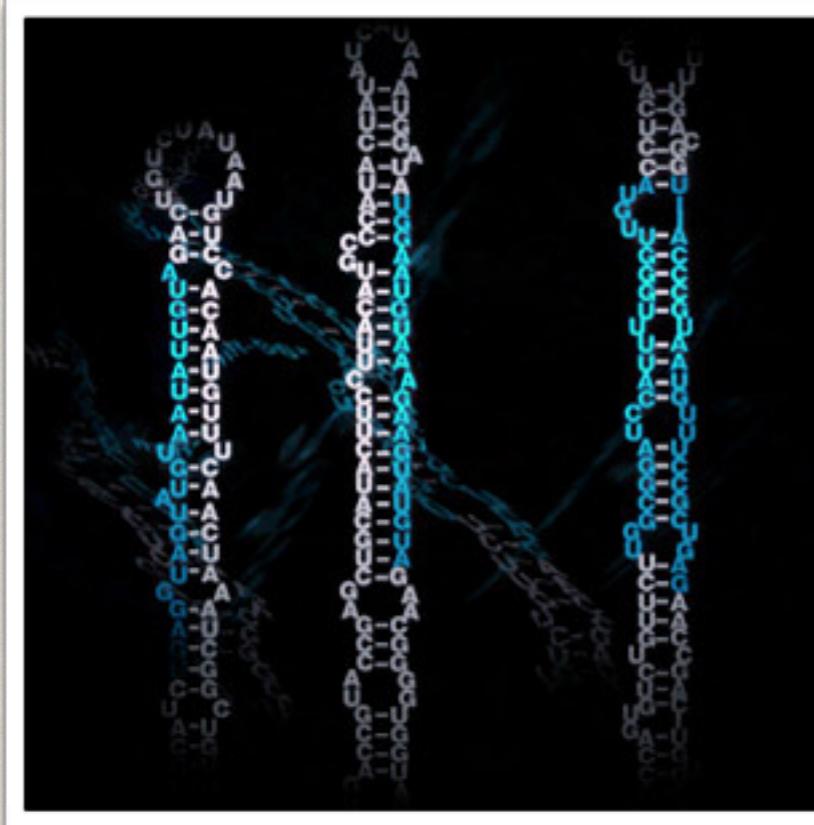
@lopantano

miRNAs

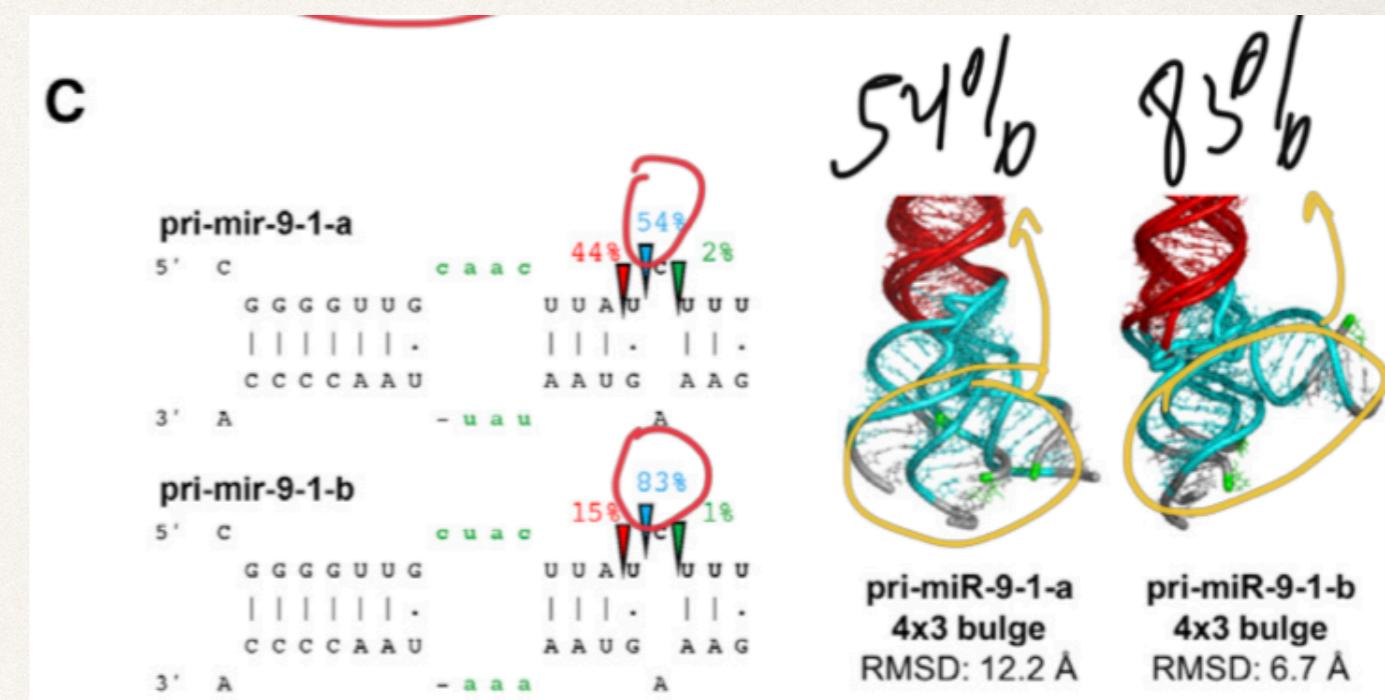
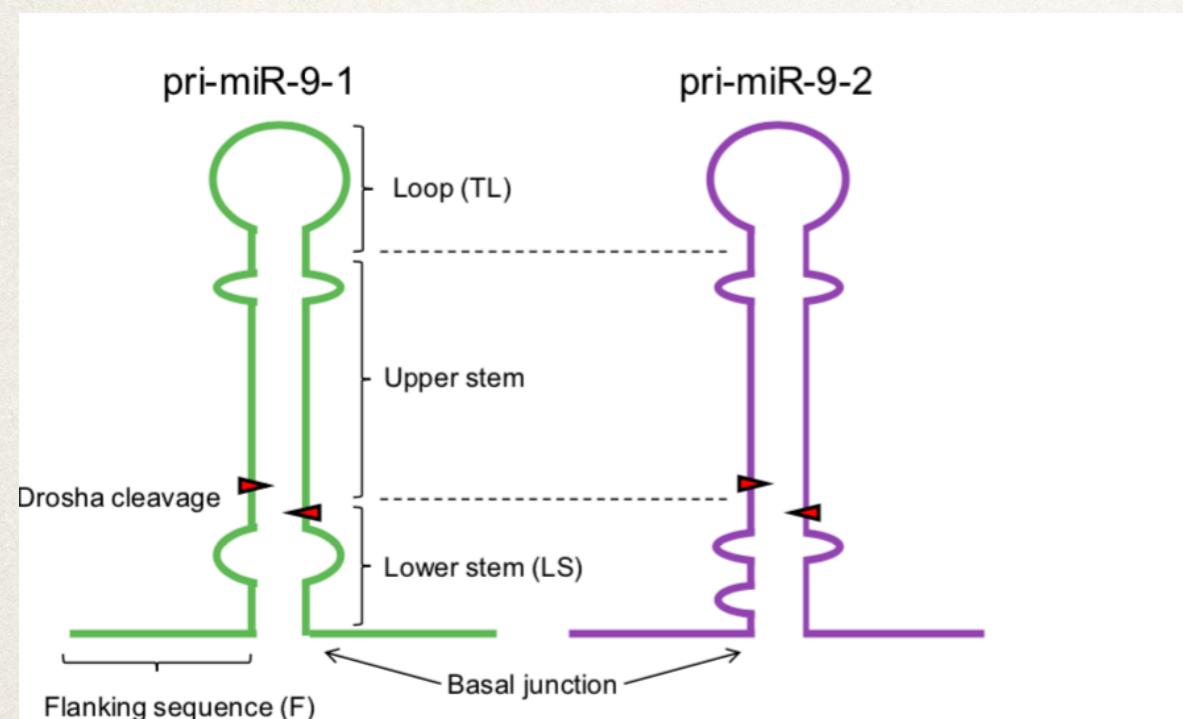
Gene regulation by imperfect complementary between seed region in miRNA and 3'UTR in the targeted RNA molecule.



isomiRs



GGG**A**TGAGGTAGGTTGTATAGTT**T**TAGG
TGAGGTAGGTTGTATAGTT
ATGAGGTAGGTTGTATAGTT**T**
TGAGGTAGGTTGTATAGT**T**
TGAGGTAGGTTGTATAGTT**AA**
TGAGGTAGGTTGTATAGT**T****AA**



Biomarkers for the early-detection and monitoring of Huntington's Disease

David W. Salzman*, Joli Bregu[^], Nathan S. Ray*, and Richard H. Myers*[^]

*sRNAlytics Inc. AstraZeneca BioHub Incubator, 35 Gatehouse Drive, Waltham MA, 02451

[^]Boston University Medical School, Department of Neurology, 72 East Concord Street, Boston, MA, 02118

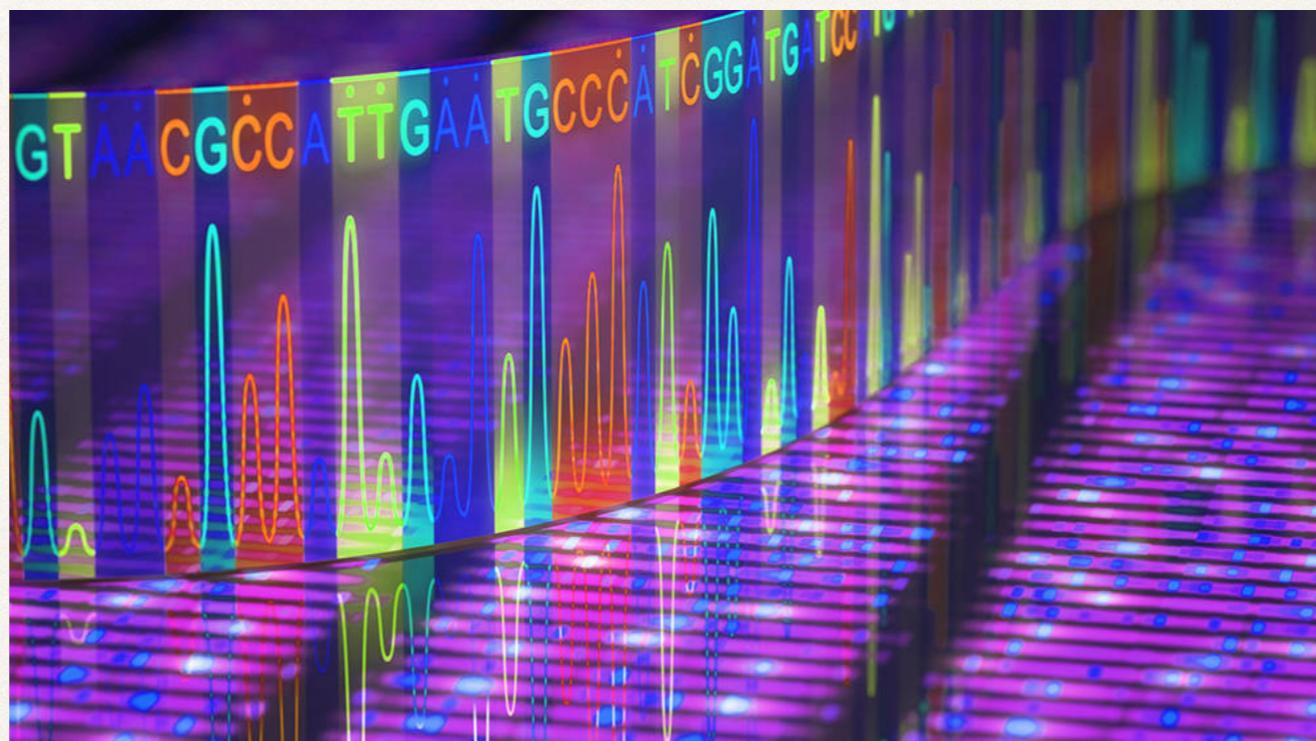
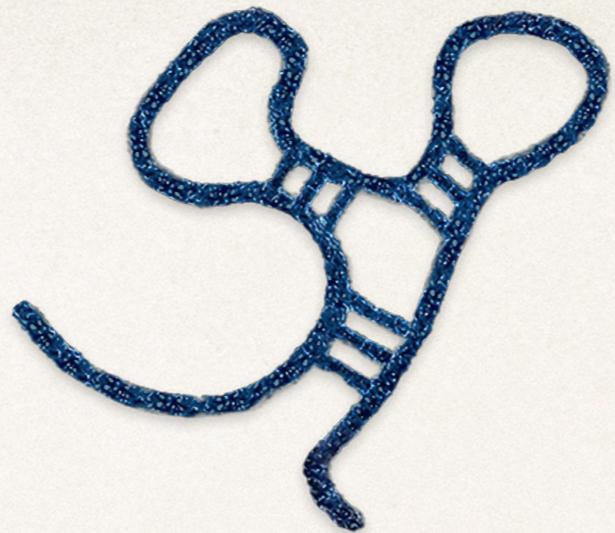
Huntington's Disease (HD) is a genetic disease caused by a CAG trinucleotide repeat in Exon1 of the huntingtin gene. Neurodegeneration results in the loss of cognitive and motor functions, and is caused by aggregation of mutant huntingtin protein in striatal neurons. Volumetric changes in the striatum can be detected decades before the manifestation of clinical phenotypes, indicating that therapeutic intervention would need to occur long before symptomatic presentation. In clinical practice and research settings, the Unified Huntington's Disease Rating Scale (UHDRS) is utilized to evaluate a patients overall physical and neurological health. UHDRS is also the most widely used outcome measure for establishing drug efficacy. However, ...

Nucleic Acids Res. 2017 Apr 7;45(6):2973-2985. doi: 10.1093/nar/gkx082.

Knowledge about the presence or absence of miRNA isoforms (isomiRs) can successfully discriminate amongst 32 TCGA cancer types.

Telonis AG¹, Magee R¹, Loher P¹, Chervoneva I², Londin E¹, Rigoutsos I¹.

 **Author information**

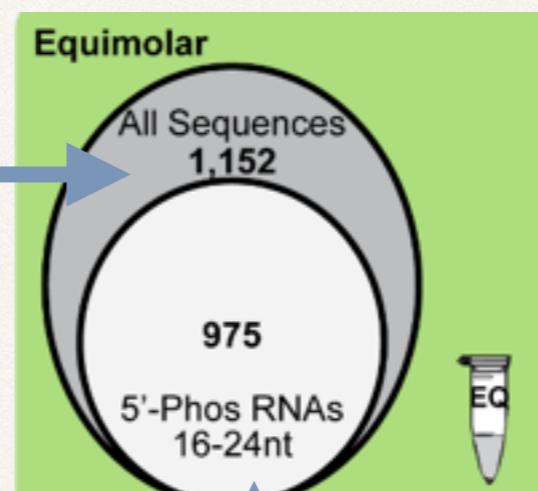


miRNAs

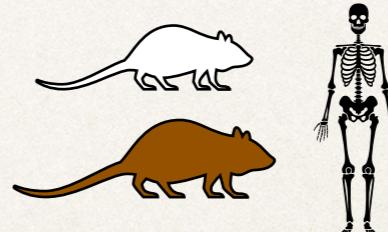
simplify idea of synthetic

Giraldez et al.

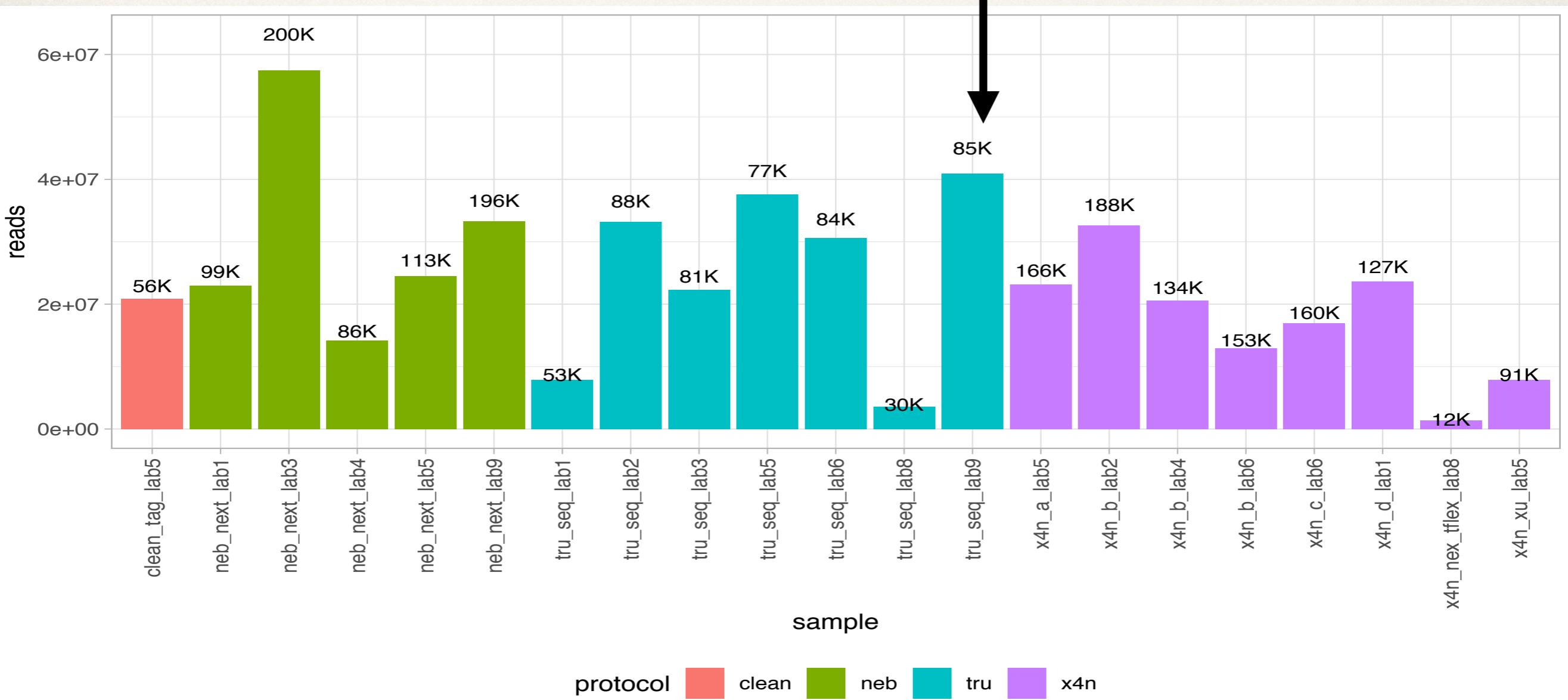
Tewari custom



Tewari synthetic



40 million reads and 85K different sequences.





mirTOP

miRNA transcriptome open project

<http://mirtop.github.io>

Repositories 5 People 20 Teams 2 Projects 0 Settings

Find a repository... Type: All Language: All Customize pinned repositories [New](#)

incubator

Where all ideas and discussions happen to lead to new repositories

● R ★ 3 ⚡ 4 Updated a minute ago



miRTOP.github.io

project for small RNA standard annotations

● HTML ★ 2 ⚡ 1 MIT 1 issue needs help Updated 4 days ago



mirtop

command lines tool to annotate miRNAs with a standard mirna/isomir naming

formatter mirna gff isomirs smallrna-seq

● Python ★ 6 ⚡ 13 MIT 17 issues need help Updated 5 days ago



simulator

first ideas and brainstorming for small RNA simulator

● C++ Updated 25 days ago

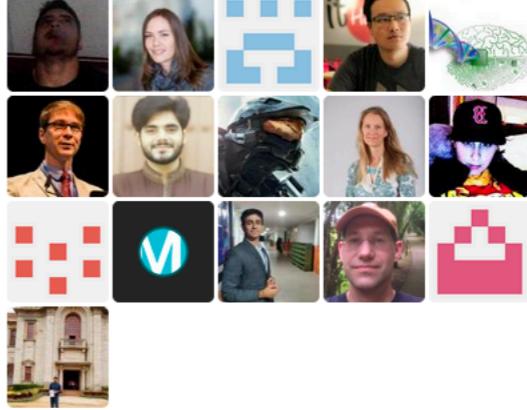


Top languages

● C++ ● Python ● R ● HTML

People

20 >



Invite someone



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH



Harvard Medical School
Research Computing



JOHNS HOPKINS
SCHOOL *of* MEDICINE

SciLifeLab



National Institutes of Health

O UNIVERSITY OF
OREGON



UNIVERSIDAD
DE GRANADA

U **HEALTH**
UNIVERSITY OF UTAH

BROAD
INSTITUTE



BC
CAN
CER

Provincial Health Services Authority

CRG[®]
Centre
for Genomic
Regulation



Jefferson

HOME OF SIDNEY KIMMEL MEDICAL COLLEGE

BRIGHAM HEALTH



**BRIGHAM AND
WOMEN'S HOSPITAL**



Main projects:

- a format
 - toolkit
-



Format derived from GFF3: mirGFF3

a

```
## mirGFF3. VERSION 1.1
## source-ontology: miRBasev21 doi:10.25504/fairsharing.hmgte8
## COLDATA: sample1
```

e

b

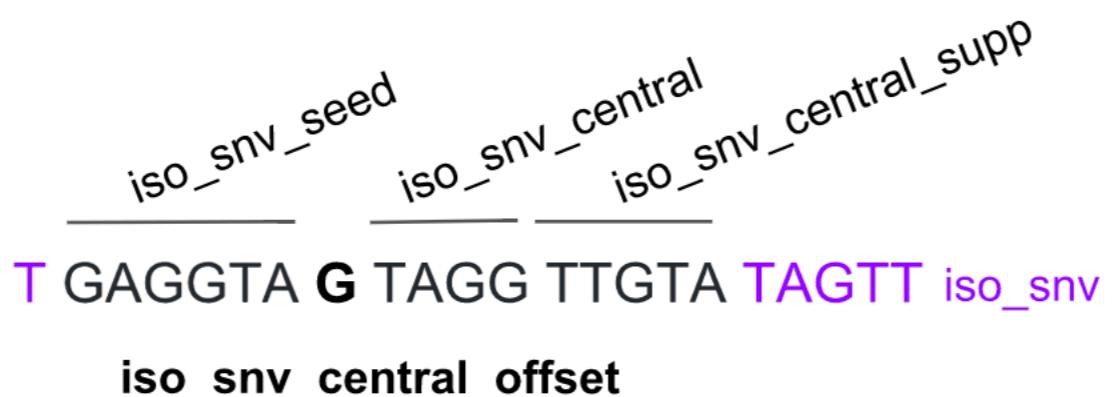
Read=GATGAGGTAGTAGGTTGTATAGTT -> UID=**iso-24-5URPV39QFE**

Read=ATGAGGTAGTAGGTTGTATAGTT -> UID=**iso-23-I0S31NSL0E**

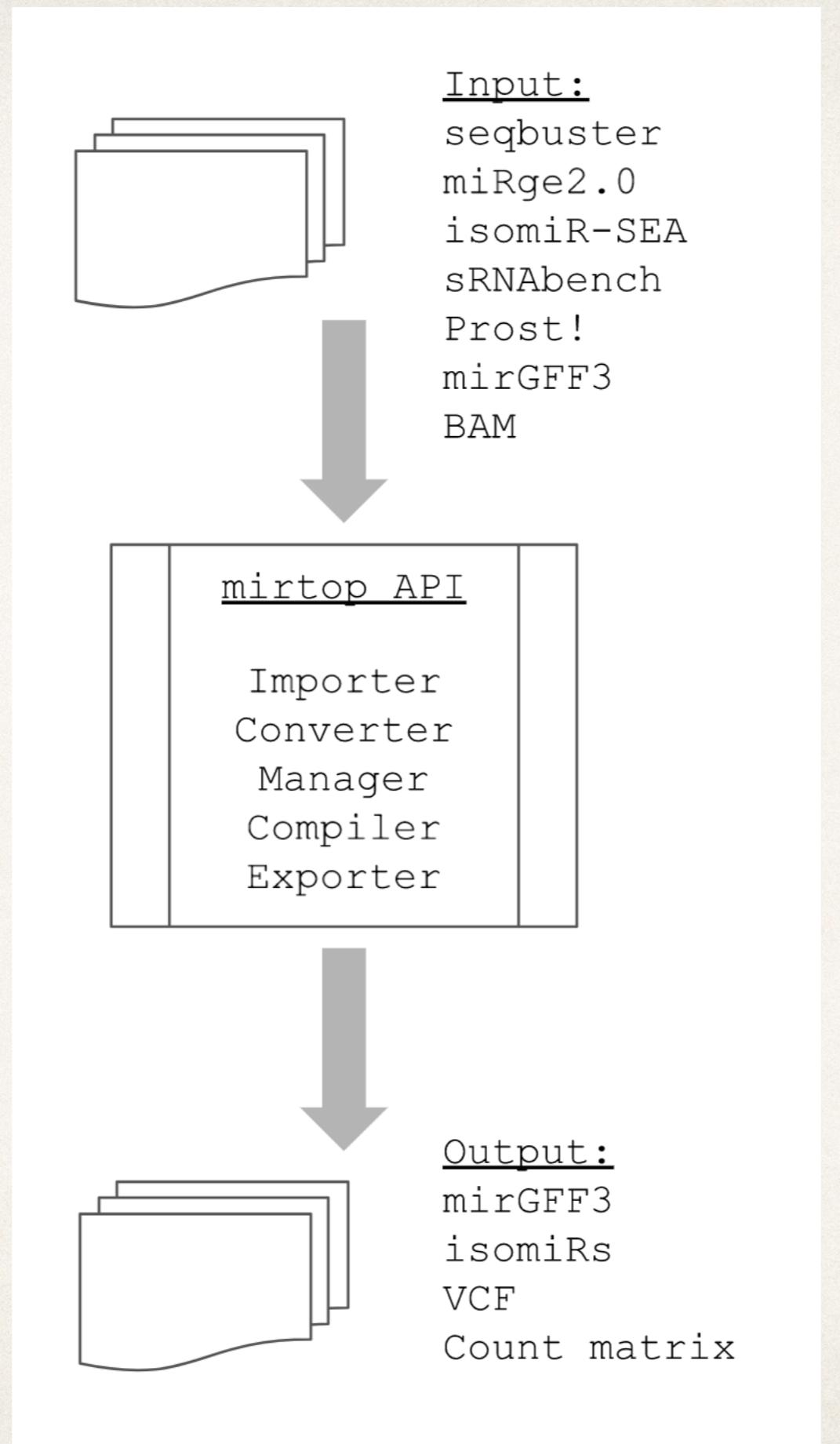
c

GGG A TGAGGTAGTAGGTTGTATAGTT T TAGG	Precursor
TGAGGTAGTAGGTTGTATAGTT	User defined reference
A TGAGGTAGTAGGTTGTATAGTT T	iso_5p:-1, iso_3p:+1
T GAGGTAGTAGGTTGTATAGT T	iso_5p:+1, iso_3p:-1
TGAGGTAGTAGGTTGTATAGTT AA	iso_add:2
T GAGGTAGTAGGTTGTATAGT T AA	iso_5p:+1, iso_3p:-1, iso_add:2

d



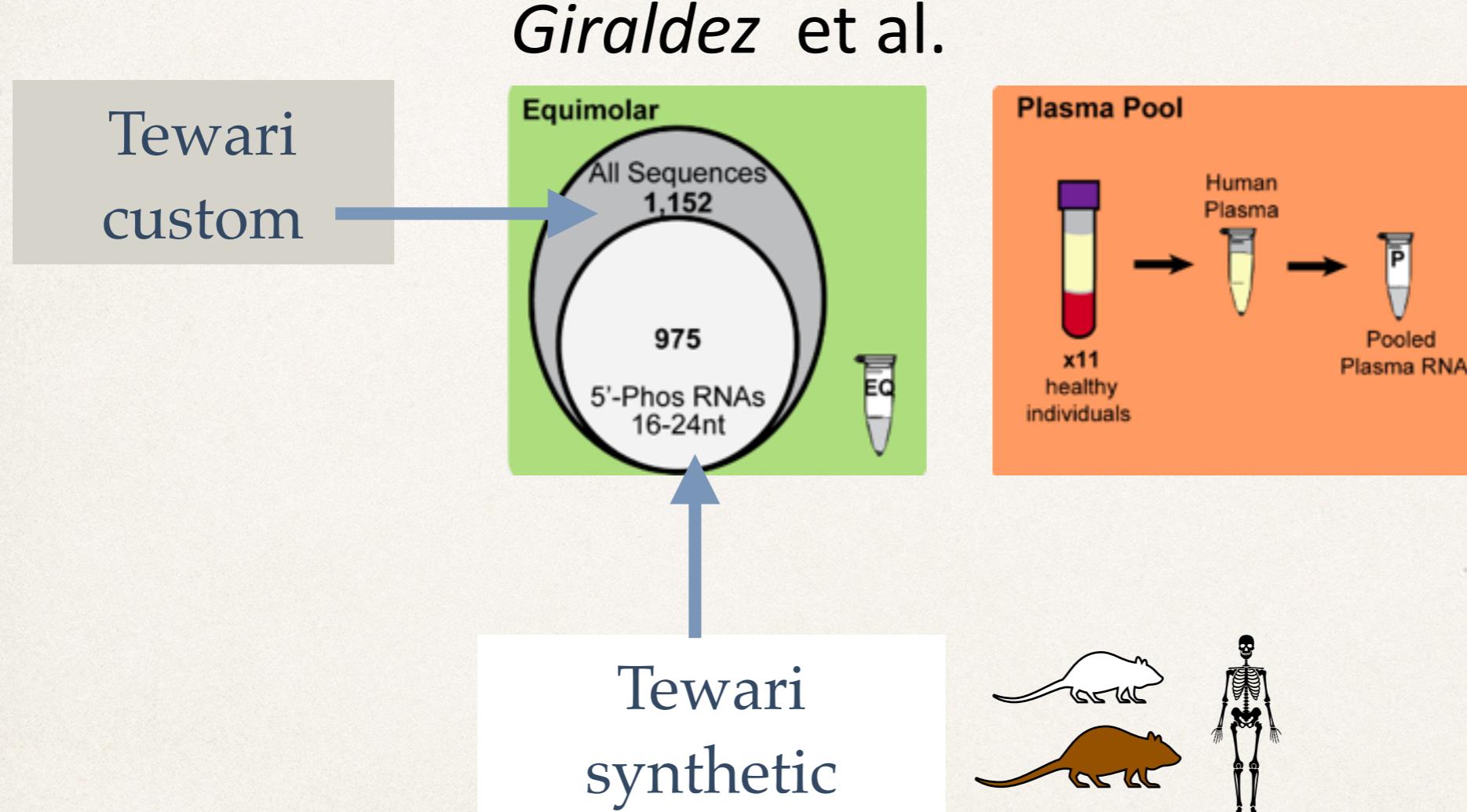
Python API: mirtop



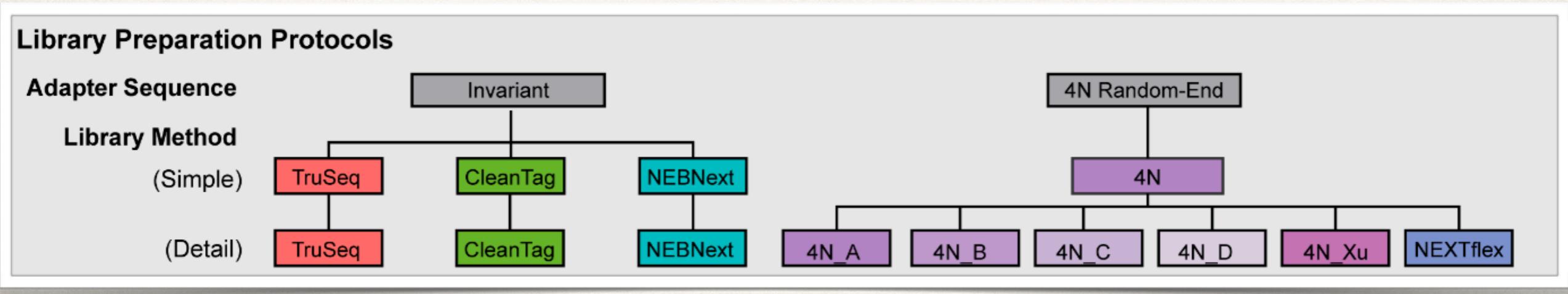
“How well NGS detect isomiRs?”

– *Scientists*

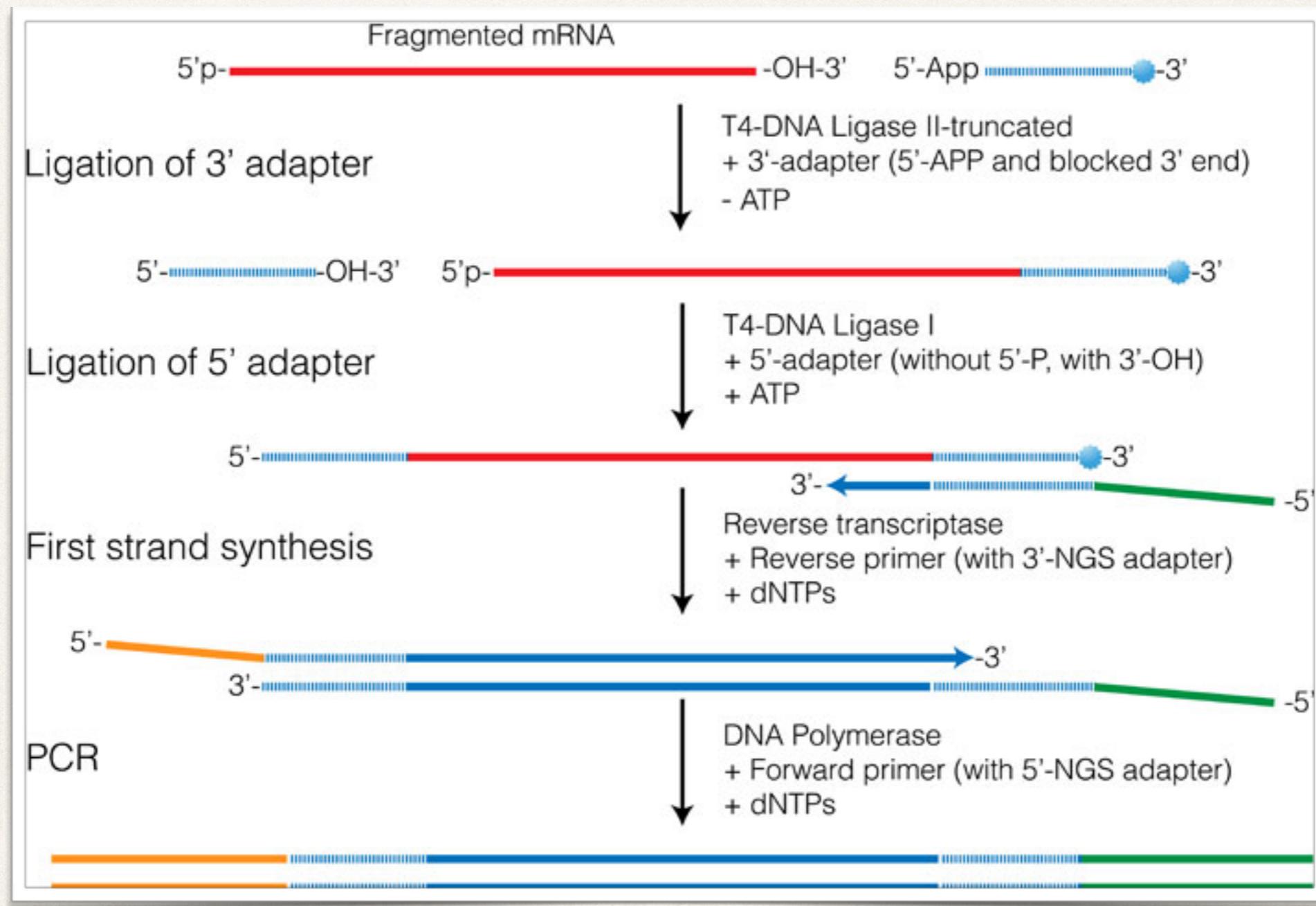
Samples suitable for benchmarking



Different protocols



Protocols



Variations

	3 adapter	5 adapter	PEG
TrueSeq			
NEBNext			Yes
NextFlex	4N	4N	Yes
CleanTag	Methylphosphonate	2'0me	
CATS	poly-T	TSO with rX	
SMARTer	poly-T	TSO	

Source of Variation

Synthesis

Shipment

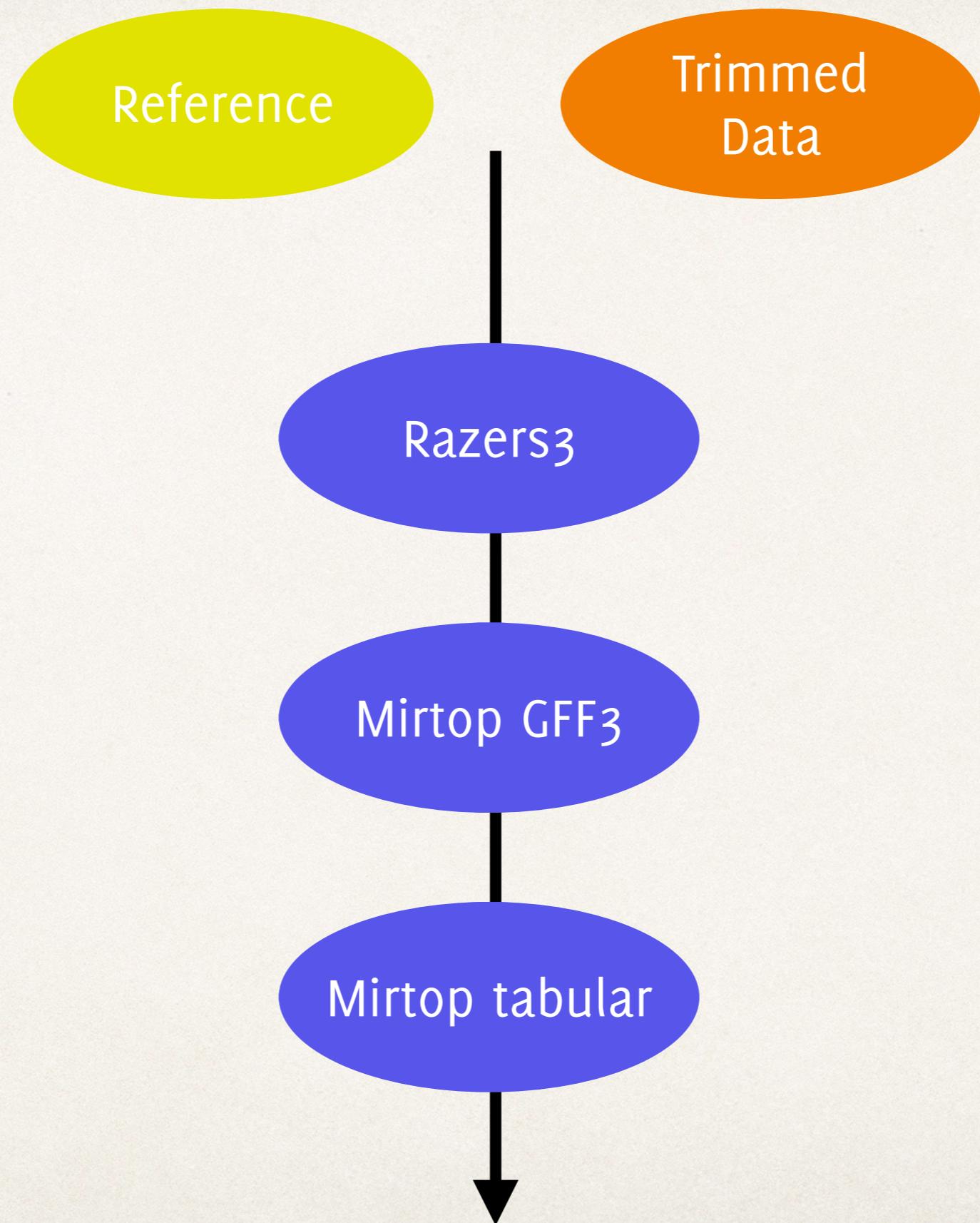
Library
preparation

Sequencing

Analysis

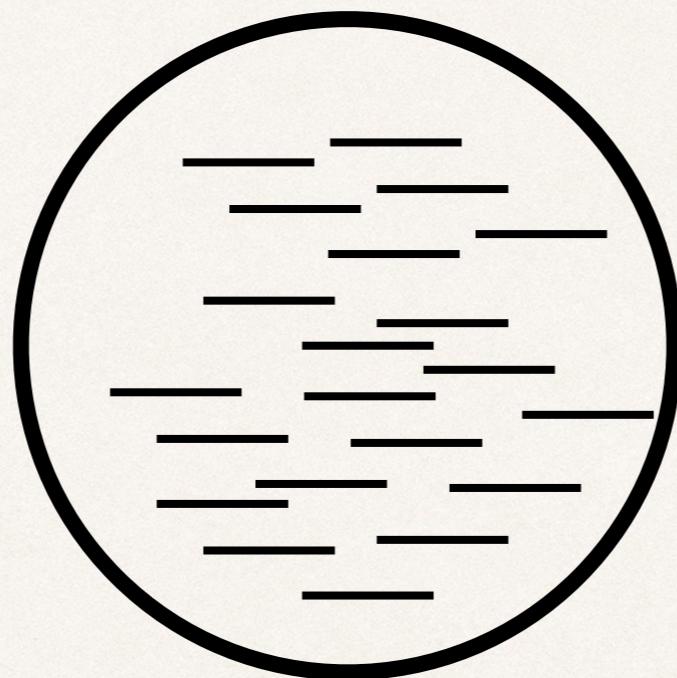
Pipeline razers3 + mirtop

snakemake



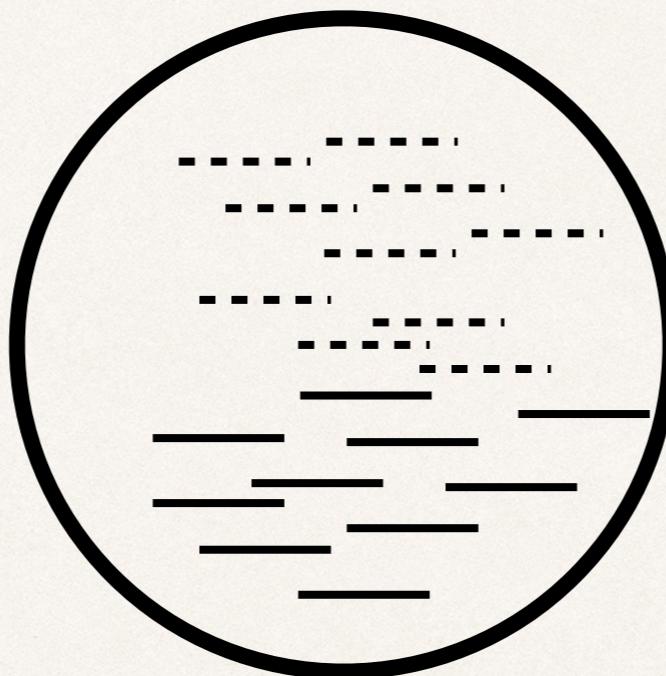
Methods and Metrics

Sample



Library size = 20

Sample

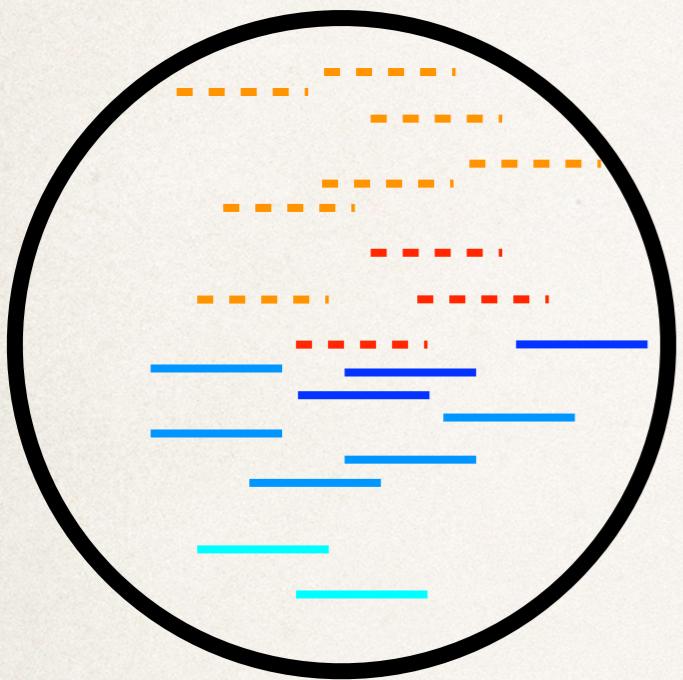


miRNA 1 = 10

miRNA 2 = 10

Library size = 20

Sample

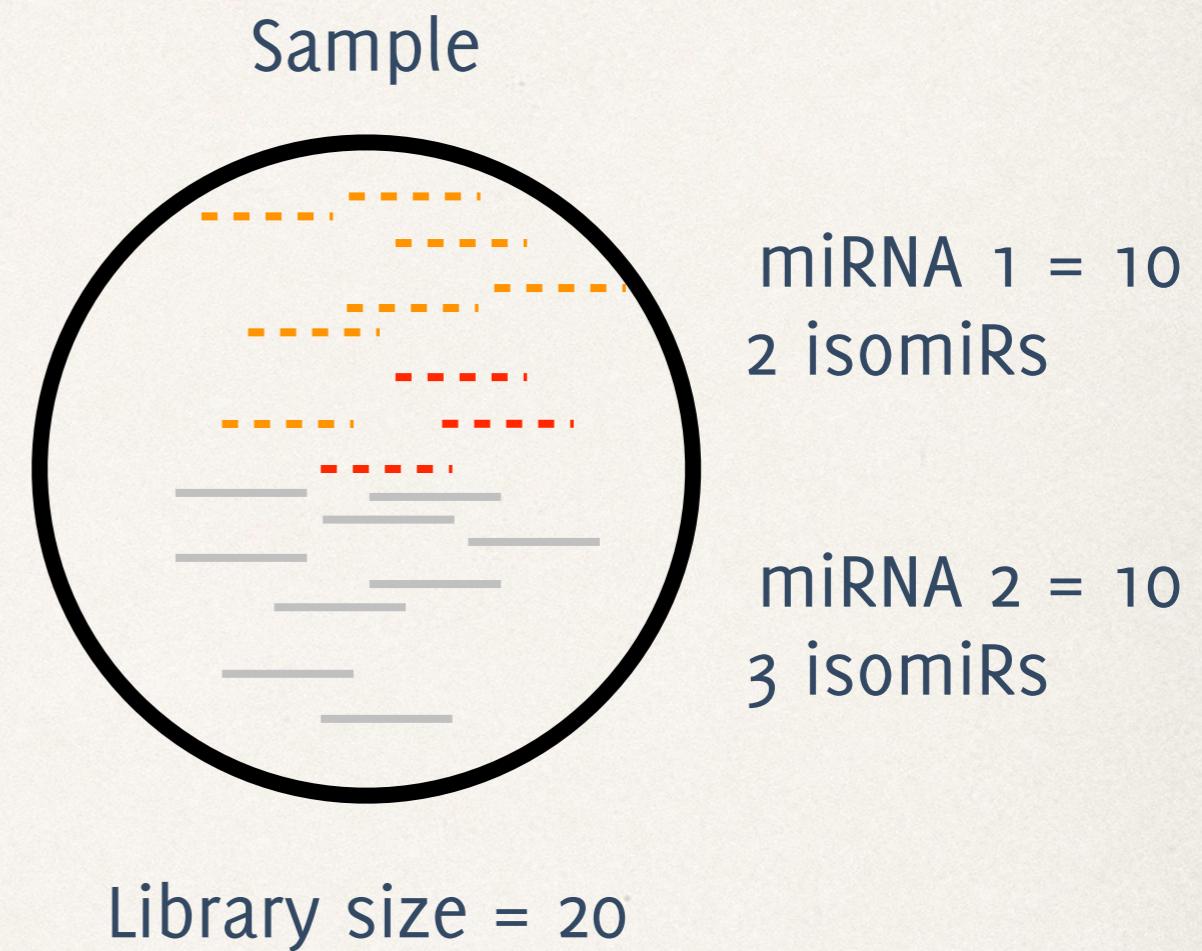
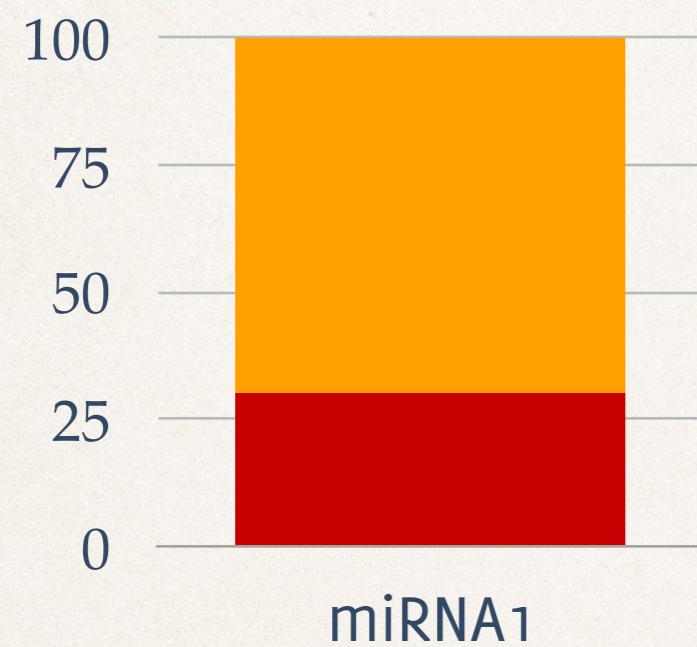


Library size = 20

miRNA 1 = X10
isomiR 1.1 = X7
(match perfectly spike-in)
isomiR 1.2 = X3
(NOT match perfectly spike-in)

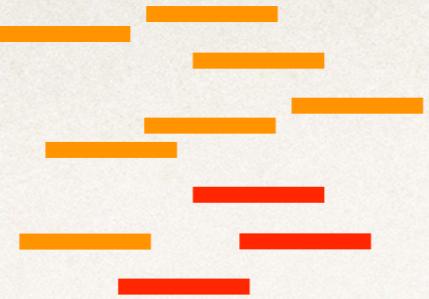
miRNA 2 = X10
isomiR 2.1 = X3
(match perfectly spike-in)
isomiR 2.2 = X5
(NOT match perfectly spike-in)
isomiR 2.2 = X2
(NOT match perfectly spike-in)

IMPORTANCE = isomiR 1.1 is $7/10=70\%$ of the miRNA1 reads



Data analysis - Filters

All sequences in a miRNA



NO

Human?

YES

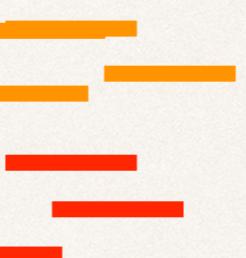


Is the spike-in
detected?



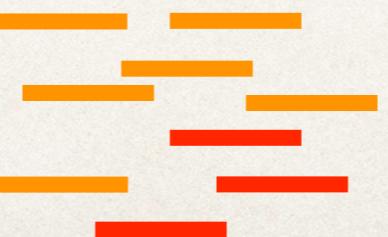
NO

YES



Any sequence in
group mapped to
other miRNA?

NO

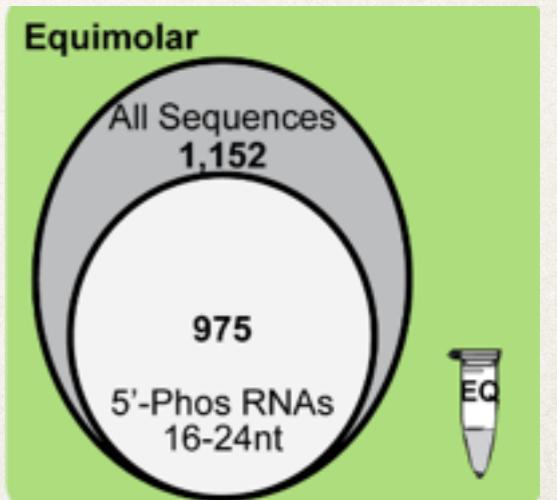


YES

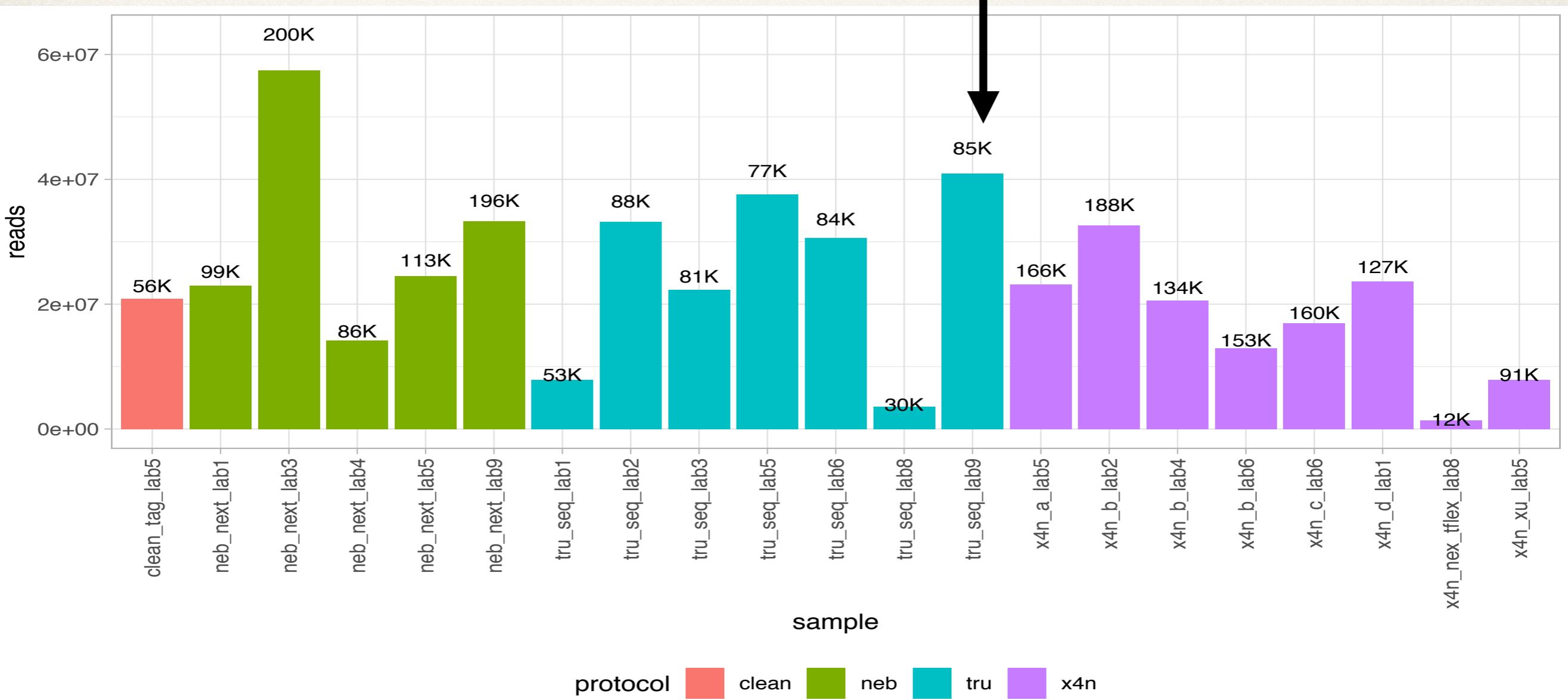


synthetic: library size

Giraldez et al.

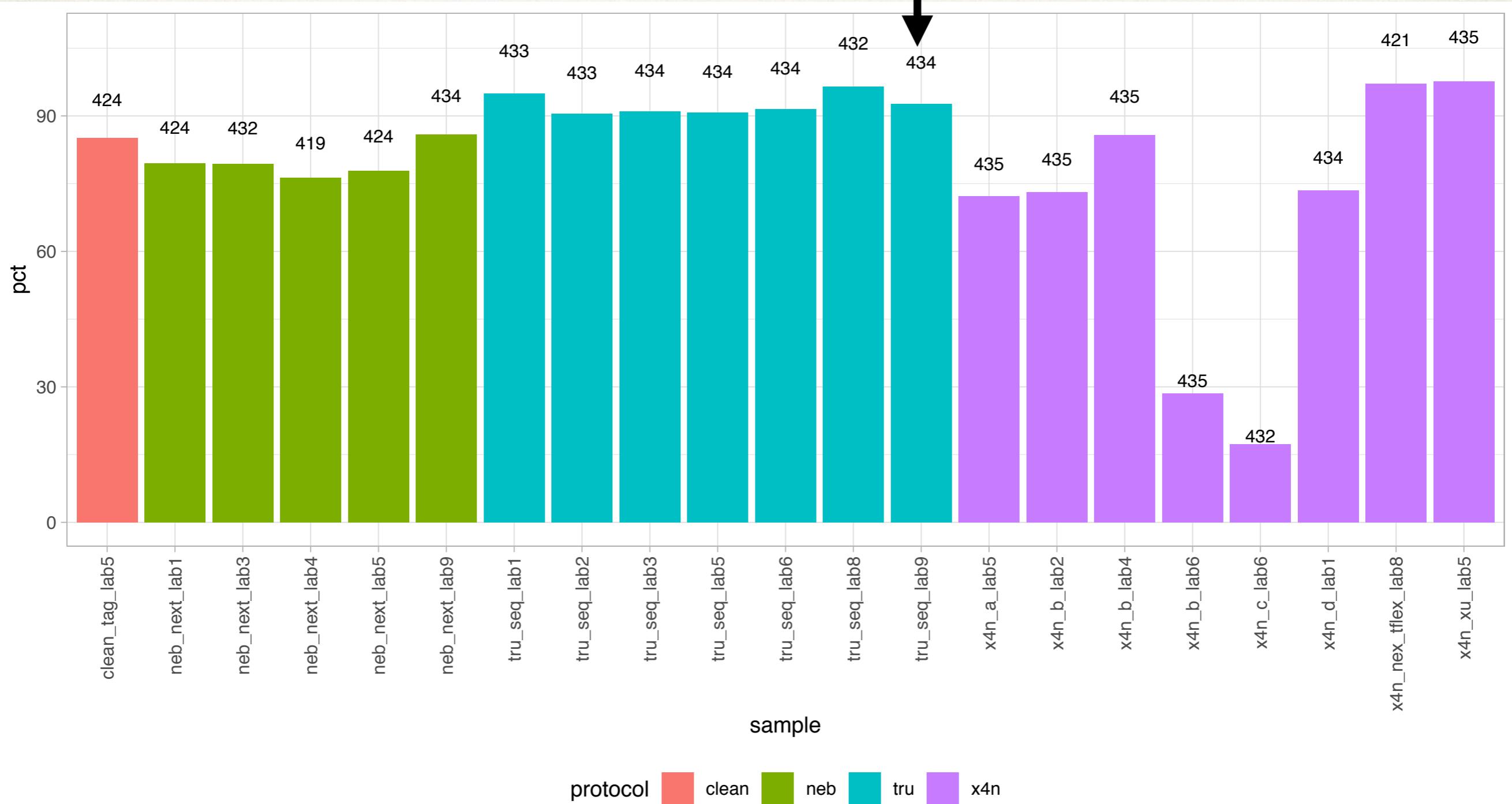


40 million reads and 85K different sequences.



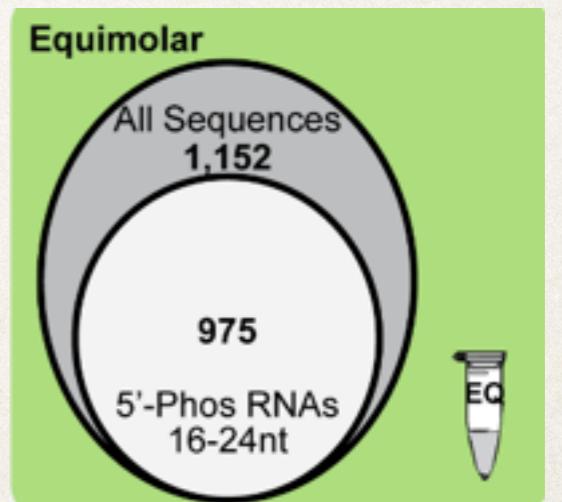
synthetic: spike-ins are the top sequence?

90% (434 being total detected miRNAs) of miRNAs
with top sequence to be the expected one.

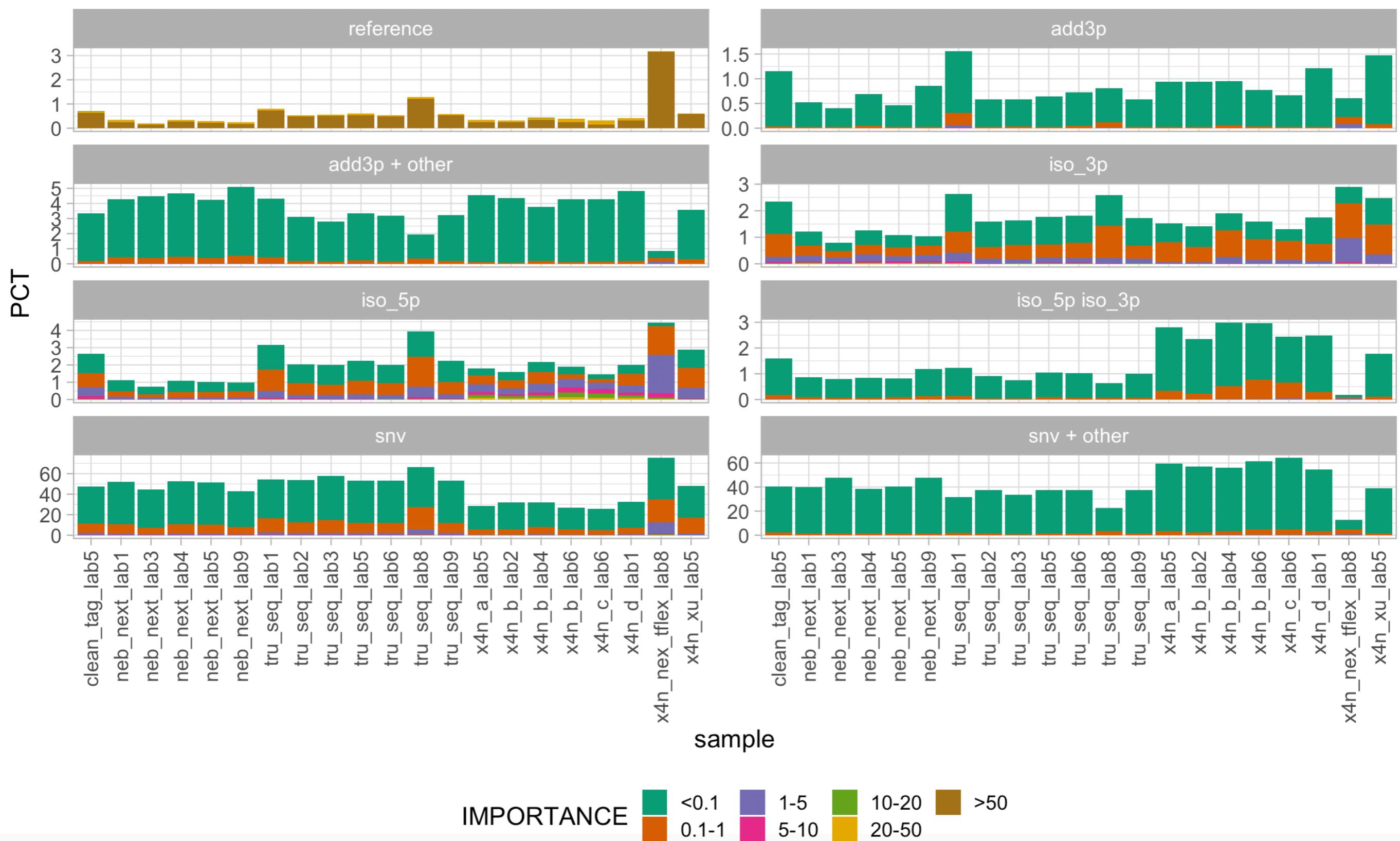


synthetic: Importance of the ‘isomiRs’ detected

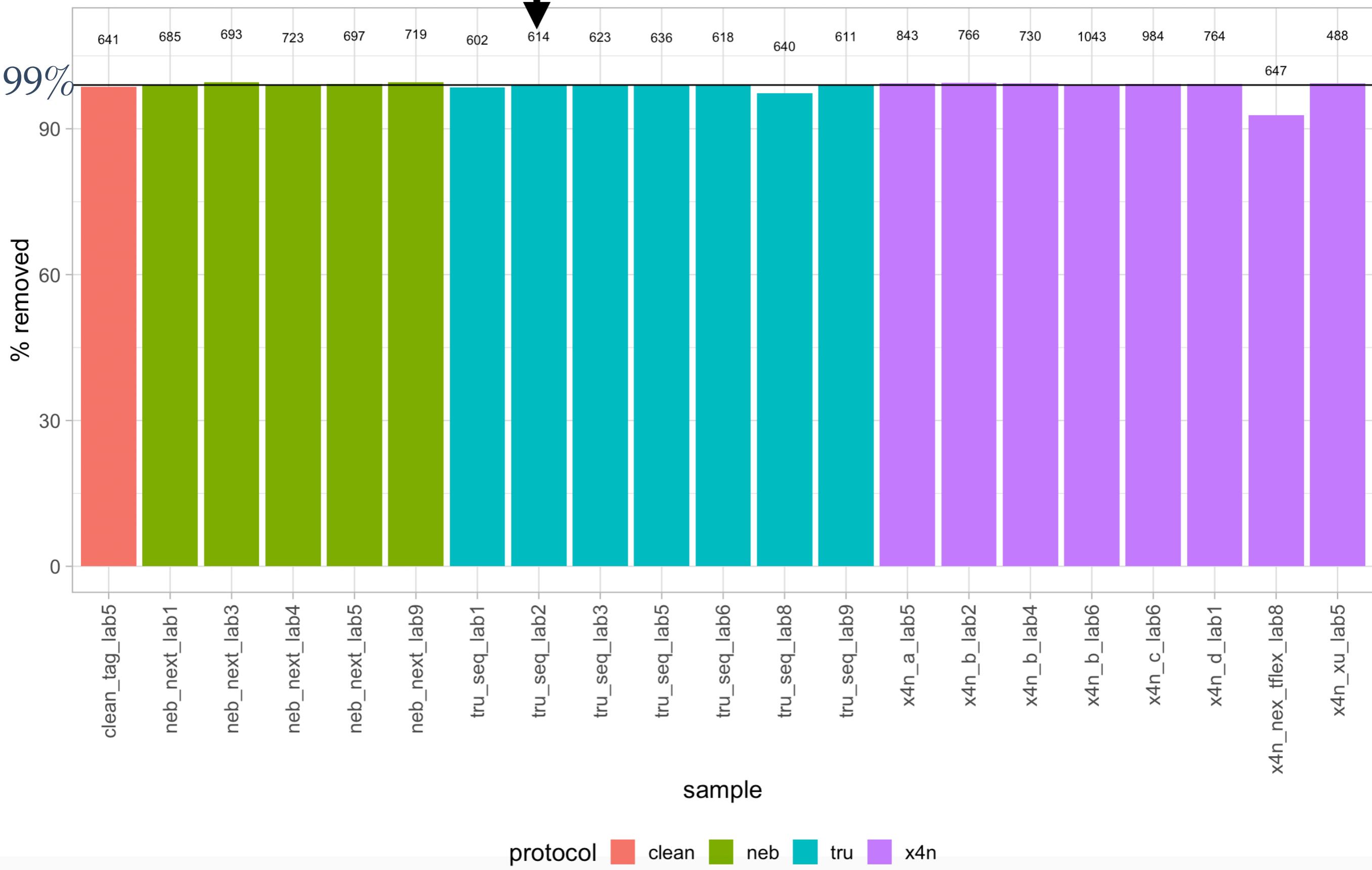
Giraldez et al.



isomiRs importance by type (miRNAs > 1000 counts)



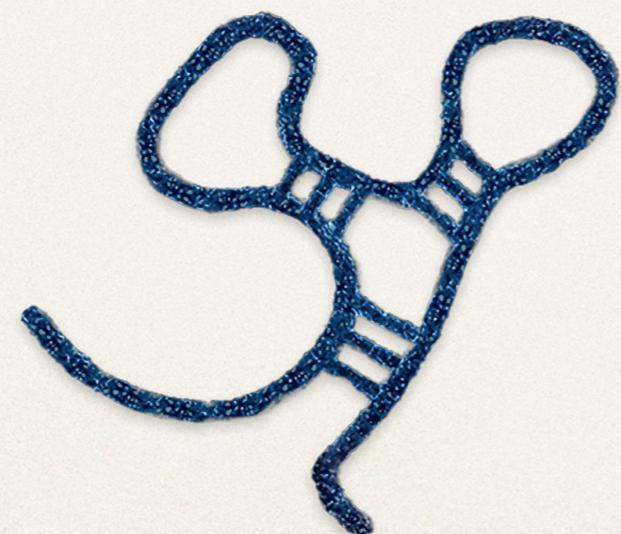
From 85K to 614 unique sequences.
99% of sequences are removed.



isomiRs importance by type (miRNAs > 1000 counts)



We assume real molecules will be detected across replicates



Sample 1

Sample 2

Sample 3

Sample 4

Synthetic miRNAs

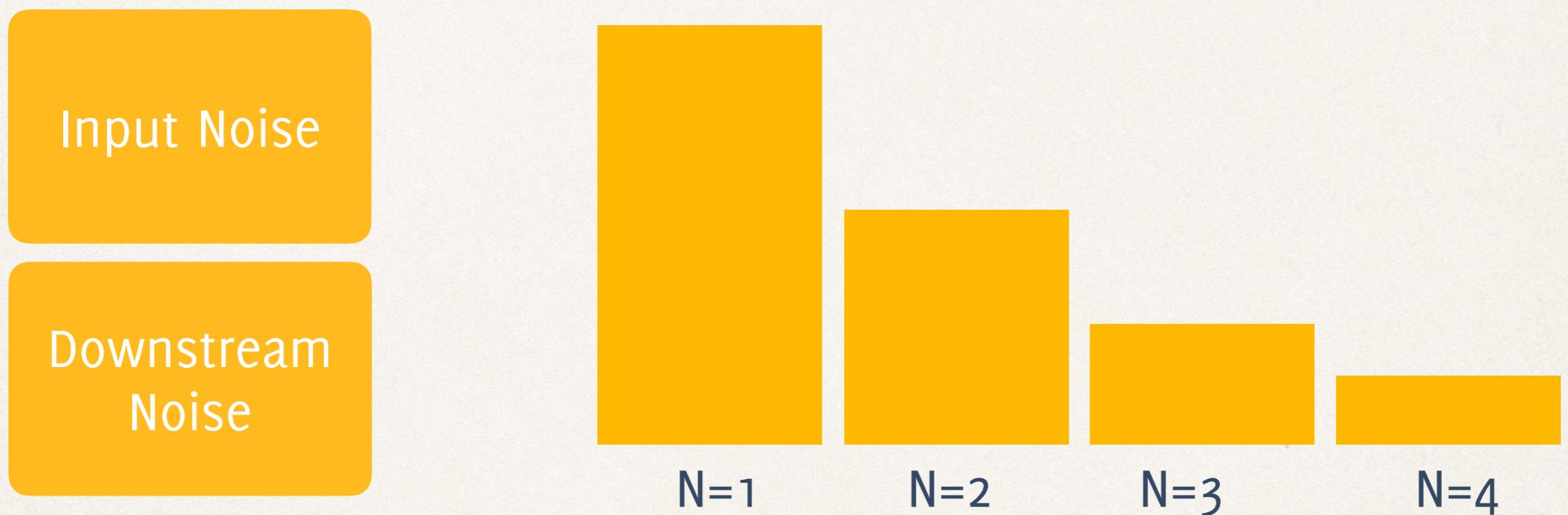
All

N=1

N=2

N=3

N=4



Narry Kim Data

Protocol 1

30 spikeins x 3 replicates

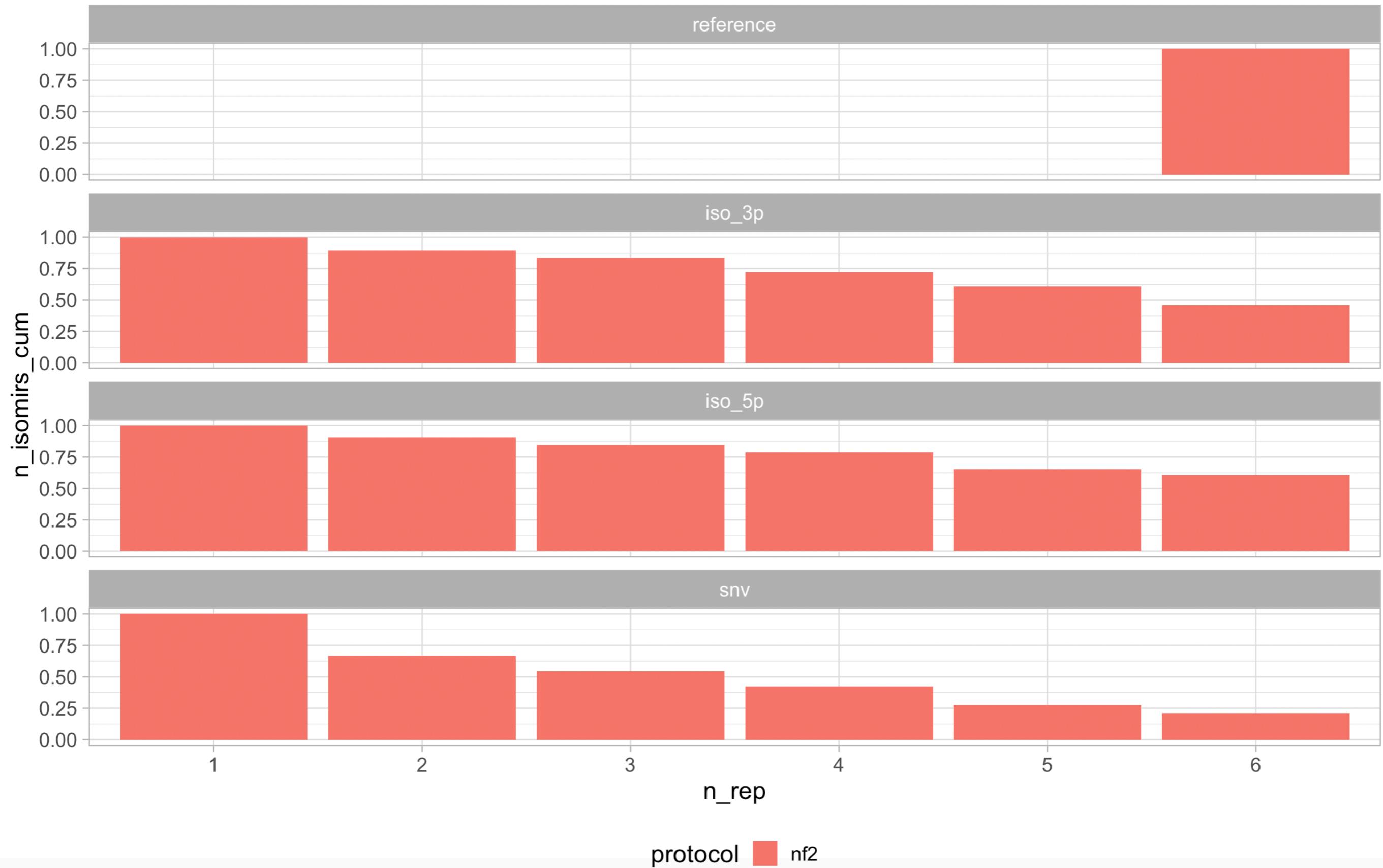
Protocol 2

30 spikeins x 4 replicates



Van Dijk Data

12 spikeins x 6 replicates



Human plasma

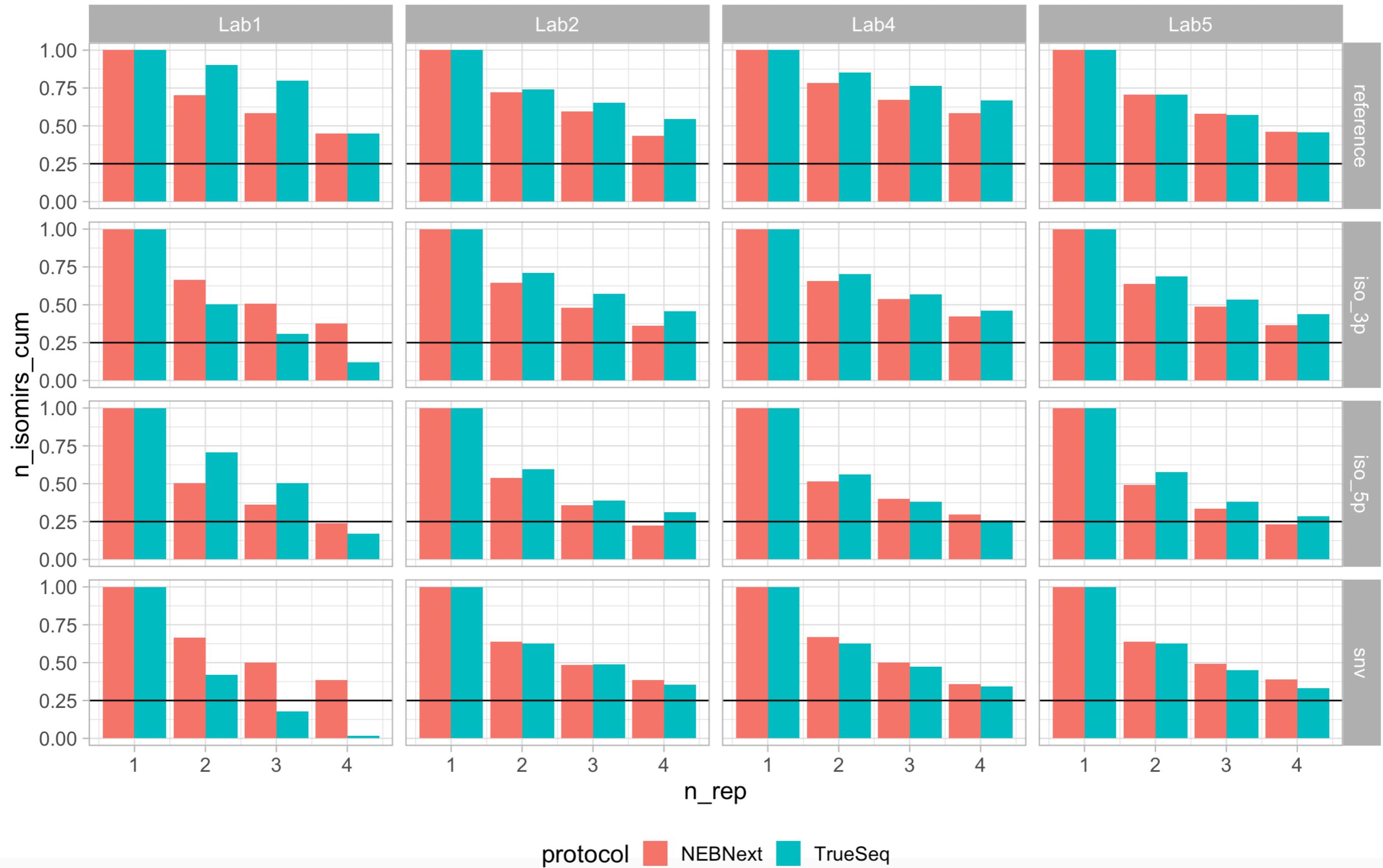
4 laboratories

Protocol 1

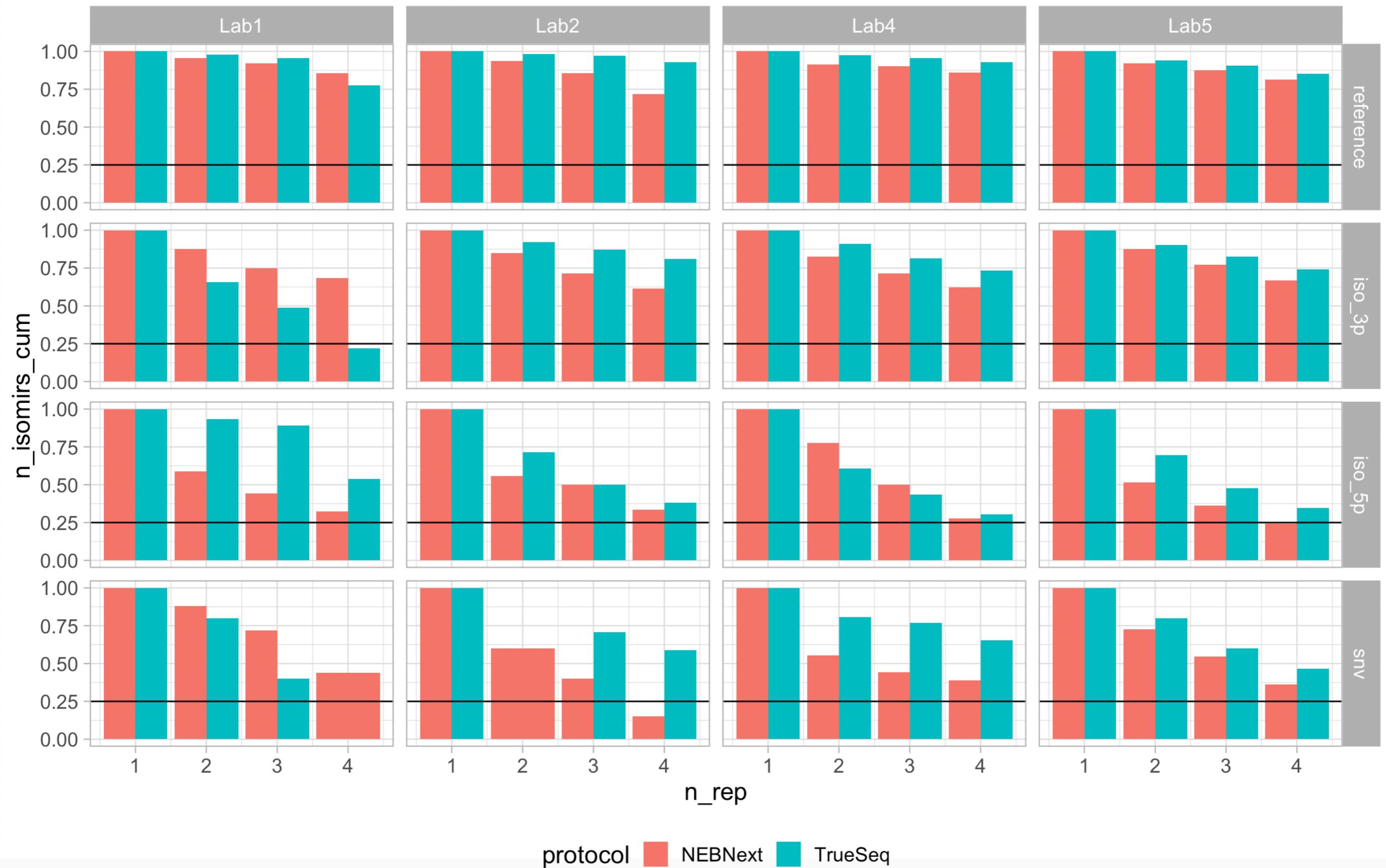
x4 replicates

Protocol 2

x4 replicates



Remove isomiRs > 10% & total > 100

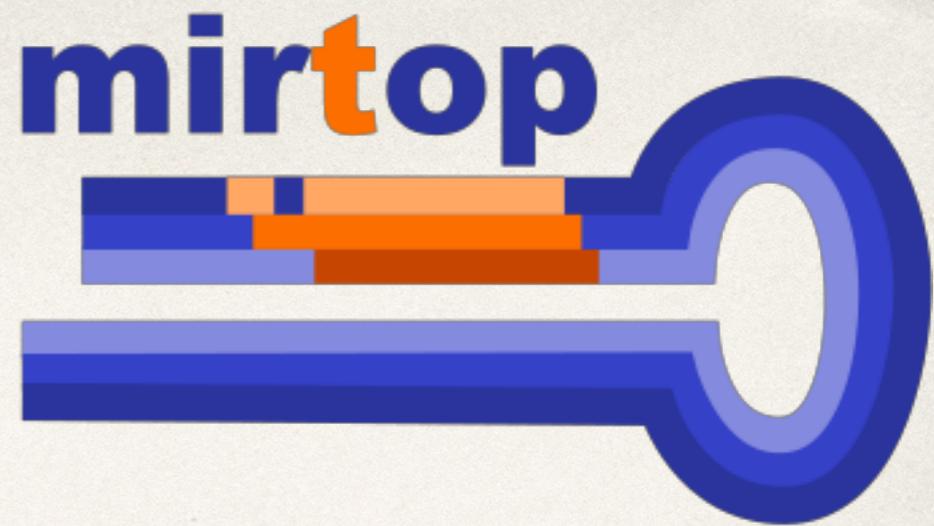


Summary

- ✓ Each miRNA generates a diversity of isomiRs:
 - ✓ 99% contribute to < 10% of the miRNA abundance
 - ✓ Independent of pipelines and data sets
- ✓ >90% of miRNAs affected
- ✓ Custom 4N protocols perform worst, BUT NEXTFLEX shows good performance

Guidelines

- ✓ NEXFlex/TrueSeq shows less unexpected sequences
- ✓ Removing isomiRs <10% of importance may clean your data
- ✓ Increase your replicates to 6 if possible
- ✓ Don't filter by sample individually



https://github.com/miRTop/incubator/tree/master/projects/tewari_equimolar

<https://github.com/miRTop/mirGFF3>

<https://github.com/miRTop/mirtop>



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH



Harvard Medical School
Research Computing

JOHNS HOPKINS
SCHOOL of MEDICINE

SciLifeLab

NIH
National Institutes of Health

UNIVERSITY OF
OREGON

Thomas Desvignes
Phillipe Loher
Karen Ellbeck
Bastian Fromm
Gianvito Urgese
Isidore Rigoutsos
Michael Hackenberg
Ioannis S. Vlachos
Marc K. Halushka

BROAD
INSTITUTE

Jeffery Ma, Jason Sydes, Yin Lu, Ernesto Aparicio-Puerta,
Shruthi Bandyadka, Victor Barrera, Peter Batzel,
Rafa Allis, Roderic Espin

CRG
Centre for Genomic
Regulation



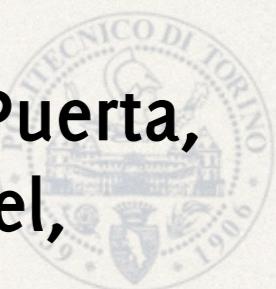
Jefferson

HOME OF SIDNEY KIMMEL MEDICAL COLLEGE



BRIGHAM HEALTH

BRIGHAM AND
WOMEN'S HOSPITAL



Kieran O'Neill, Eric Londin, Aristeidis G. Telonis,
Elisa Ficarra, John H. Postlethwait,



Provincial Health Services Authority

Thank you!



