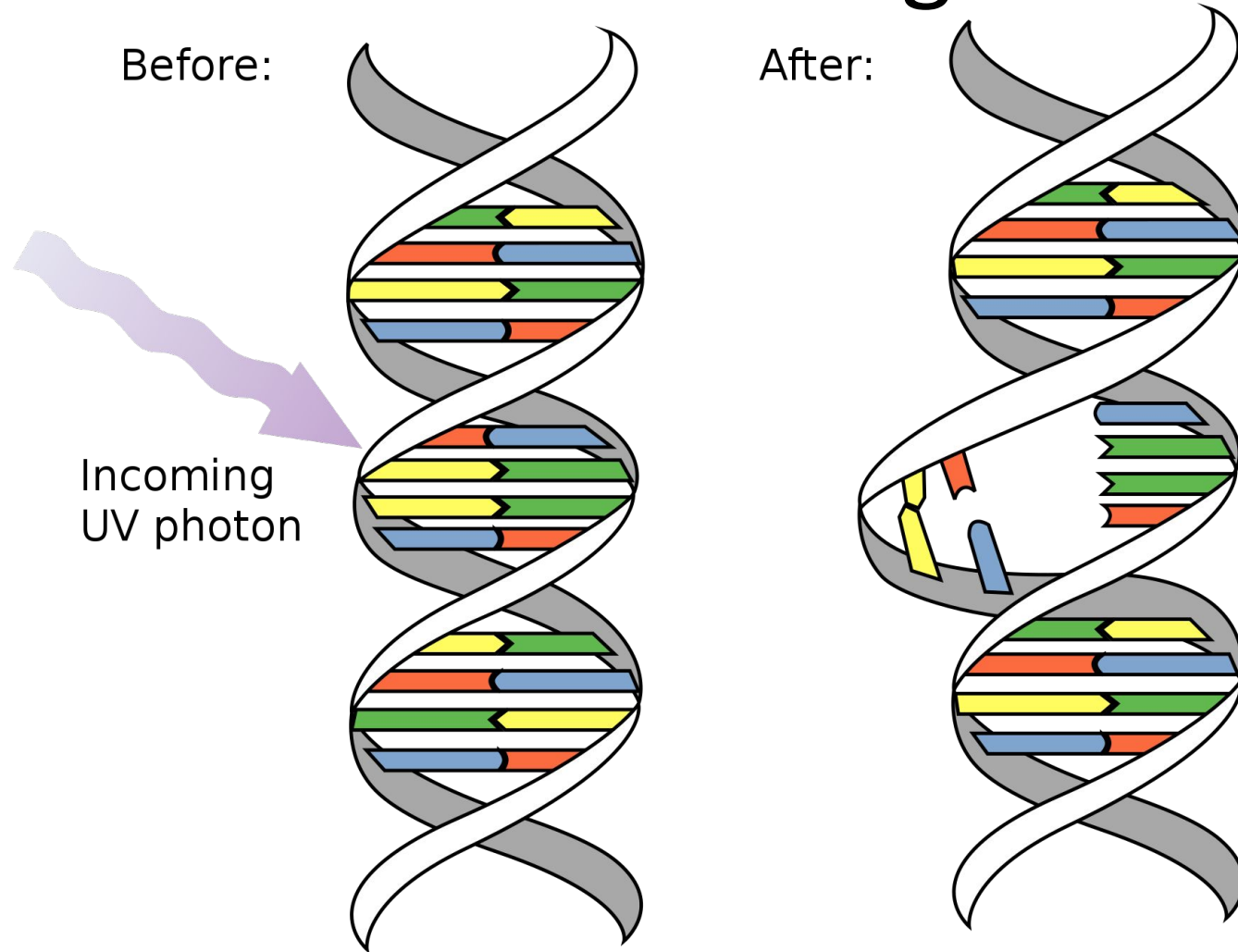


Benchmarking of Mutational Signature Assignment Tools

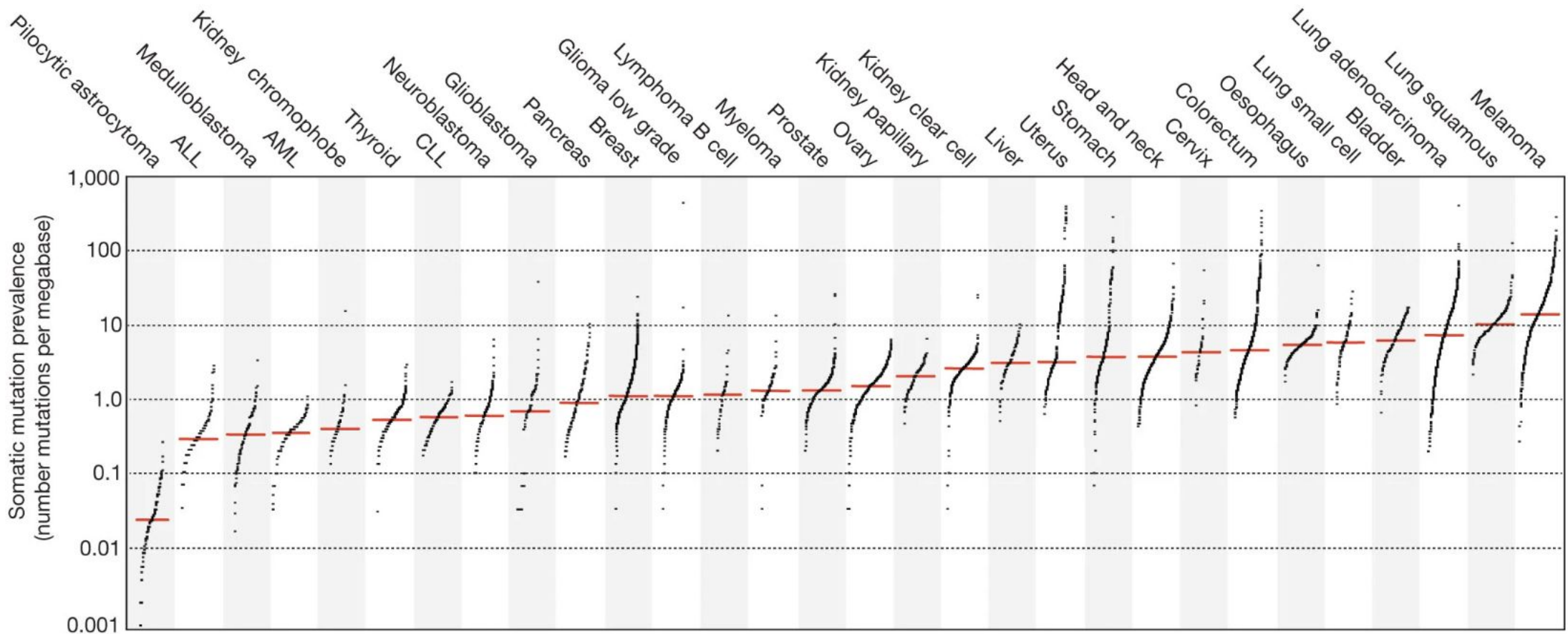
Xi (Sam) Wang, Dr. Marcos Díaz-Gay, Dr. Raviteja Vangara,
Dr. Ludmil B. Alexandrov

UC San Diego

Somatic Mutations: changes in DNA

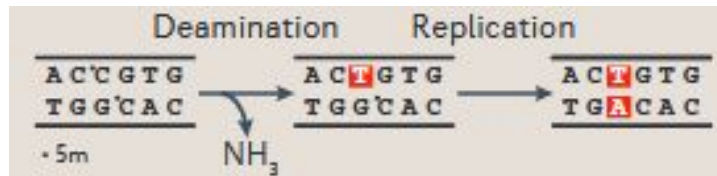
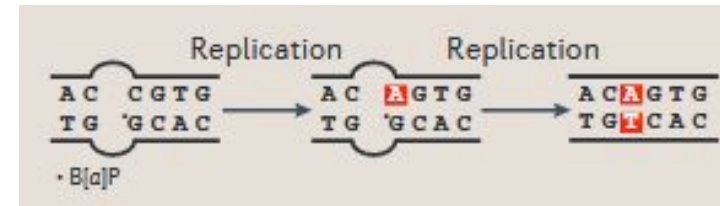


Somatic mutations are known to be generated by exposures



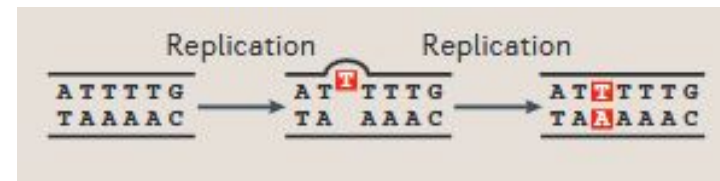
Patterns of mutations are linked to sources of DNA damage

Environmental exposures
Tobacco smoking or chewing



Normal cellular activities
Spontaneous deamination of
methylated cytosines

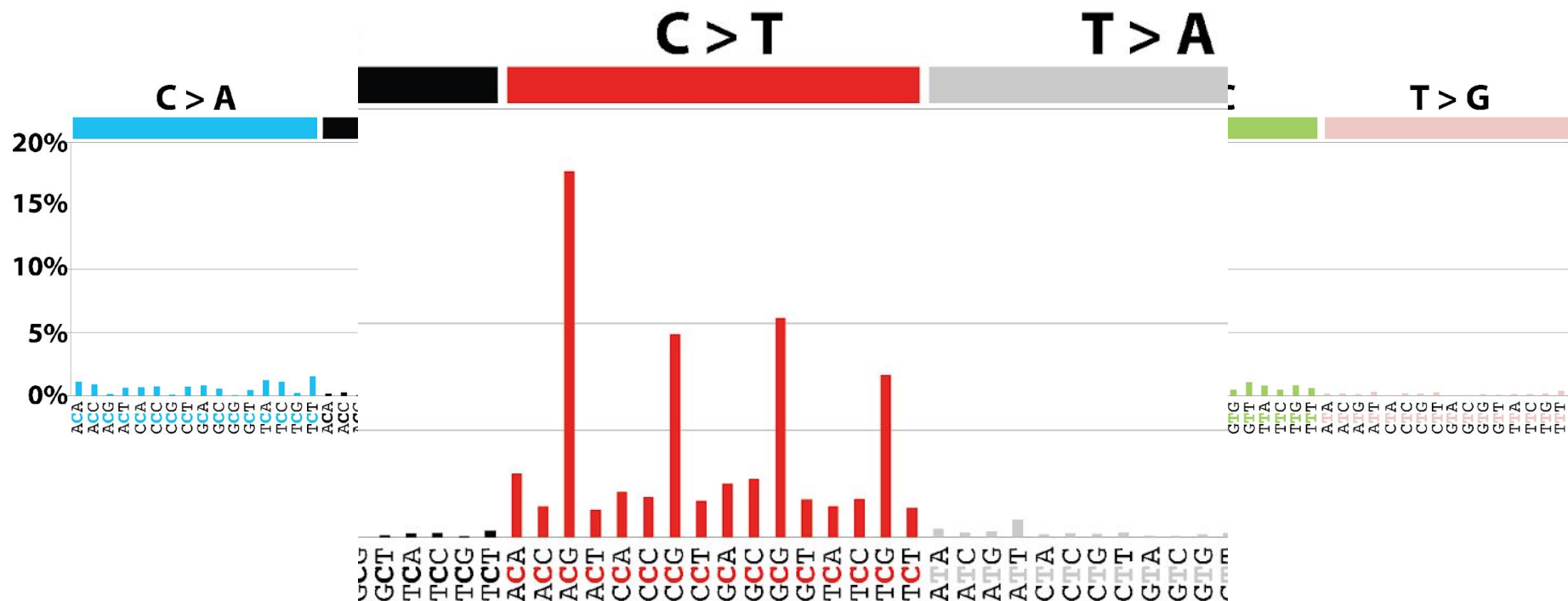
Failure in DNA replication or repair
Aberrant mismatch repair pathway



Mutational Signature: defined by base substitutions and context

Six classes of single-base mutations
Reported by pyrimidine

Adding 5' and 3' adjacent bases
96 possibilities considering context



Signature Assignment

- The MATRIX:



$$X = W \times H$$

Mutation Matrix (given)

Signatures Matrix
(standard cosmic_v3)

Activities/Exposures Matrix
(WHAT WE WANT!)

Mutation Type	Sample 1	Sample 2	...
A[C>A]A			
A[C>A]C			
...			
T[C>G]T			

=

Mutation Type	SBS1	SBS2	...
A[C>A]A			
A[C>A]C			
...			
T[C>G]T			

X

Signatures	Sample 1	Sample 2	...
SBS1			
SBS2			
...			
SBS60			

In other words:

Mutational signature assignment is the process of finding the most contributing factors associated with an individual's potential or existing cancer causes.

An accurate assignment means accurate discovery of causes of cancer!

Problem

There are no systematic benchmarks for all published mutational signature assignment tools.



So, what do we do in this project: Benchmark them all!

Current published signature assignment tools:

- Our own tools at the Alexandrov Lab: SigProfilerSingleSample, SigProfilerExtractor (decomposition module)
- 15 other tools using following algorithms:
 - NNLS (non-negative least squares)
 - Multiple linear regression
 - Quadratic programming
 - Cone projection
 - Bootstrapping

How do we benchmark?

- Git! Automation scripts! Benchmarking performance scripts!
- Qualitative & Quantitative Performance Analysis

The screenshot shows a GitHub repository page for 'marcos-diazg Test GitHub Sam'. The repository has 3 branches and 0 tags. The commit history shows 69 commits, with the latest commit being 7e04a2b on Jun 22. The repository contains 14 folders, each with a commit message and a timestamp. The right sidebar shows the repository's description, a README link, and sections for Releases, Packages, and Languages.

Folder Name	Commit Message	Timestamp
00_Benchmark_input_data/Referen...	Add ref signatures files for SPSS manual edit	6 months ago
01_SigProfilerSingleSample	Test GitHub Sam	4 months ago
02_decompTumor2Sig	Add install and run scripts for dt2s	7 months ago
03_deconstructSigs	Reformat comment lines dS and MP	7 months ago
04_MSA	Add run_2 script for formatting results of MSA	7 months ago
05_MutaGene	Update run_2 script for MutaGene	7 months ago
06_Mutalisk	Update run_1 script for Mutalisk	7 months ago
07_MutationalCone	Update install script mutationalcone	7 months ago
08_MutationalPatterns	Update comments on MP run script	6 months ago
09_Palimpsest	Update run script for Palimpsest to use manually provided ref signatur...	6 months ago
10_QPsig	Add install and run scripts for QPsig	7 months ago
11_sigfit	Update run script for sigfit	7 months ago
12_sigLASSO	Add install and run script for sigLASSO	7 months ago
13_SignatureToolsLib	Update run script SignatureToolsLib	7 months ago
14_SignatureEstimation	Add run script for SignatureEstimation	7 months ago

About
Benchmarking scripts for mutational signature refitting analysis tools
[Readme](#)

Releases
No releases published
[Create a new release](#)

Packages
No packages published
[Publish your first package](#)

Languages
R 86.0% Shell 9.8% Python 4.2%

Current Progress?

Qualitative & Quantitative

Qualitative Analysis Method

- For the assignment results of each tool, we want:
 - True positive, true negative, false positive, false negative to get:
 - Average precision, sensitivity, and specificity across all scenarios per noise level
- A yes or no questions:

	SP.Syn.Panc-	SP.Syn.Panc-	SP.Syn.Panc-	SP.Syn.Panc-	SP.Syn.Panc-
SBS1	554.298921	1378.31487	827.105158	1116.7868	1576.23622
SBS2	204.743328	0	256.835317	61.8438701	0
SBS3	0	0	0	0	0
SBS5	4410.98461	2156.46037	7293.59929	2590.98303	1340.81445
SBS13	329.166888	133.140279	70.0235047	225.297873	108.256413
SBS17a	0	0	0	0	0
SBS17b	0	0	0	464.135025	0
SBS18	1894.12041	0	0	0	706.566553
SBS28	0	0	0	0	0
SBS30	0	0	475.873322	0	0
SBS40	0	2160.76295	2475.98966	0	0

Ground Truth

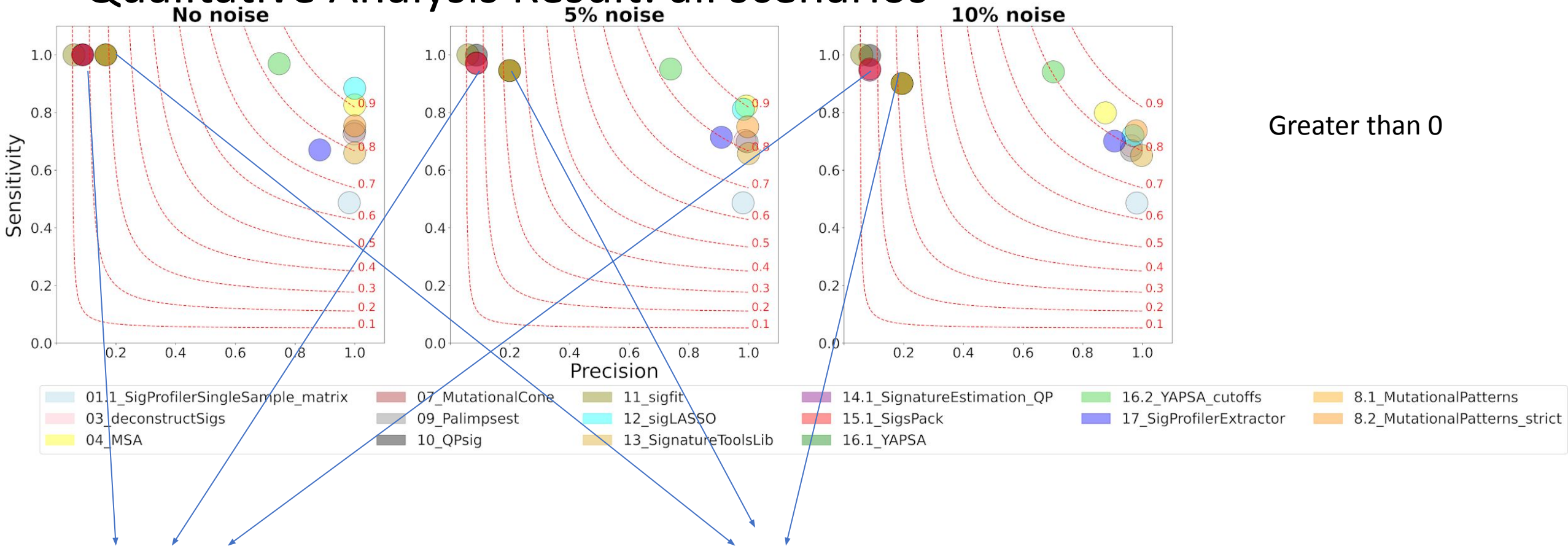
VS.

	SP.Syn.Panc-	SP.Syn.Panc-	SP.Syn.Panc-	SP.Syn.Panc-	SP.Syn.Panc-
SBS1	554.504682	1377.42657	824.560442	1116.68054	1569.43852
SBS2	0	0	0	0	0
SBS3	0	0	0	0	0
SBS5	4399.98309	2156.51546	7144.99565	2580.38002	1189.15106
SBS13	0	0	0	0	0
SBS17a	0	0	0	0	0
SBS17b	0	0	0	463.263627	0
SBS18	1893.90681	0	0	0	709.020081
SBS28	0	0	0	0	0
SBS30	0	0	0	0	0
SBS40	0	2124.91202	2455.60728	0	0
Unknown	545.605418	170.145944	976.83663	297.675817	259.390335

Assignment from
03_deconstructSigs

* Sample shown the scenario 2 without noise ground truth activities vs. activity results from tool 03_deconstructSigs

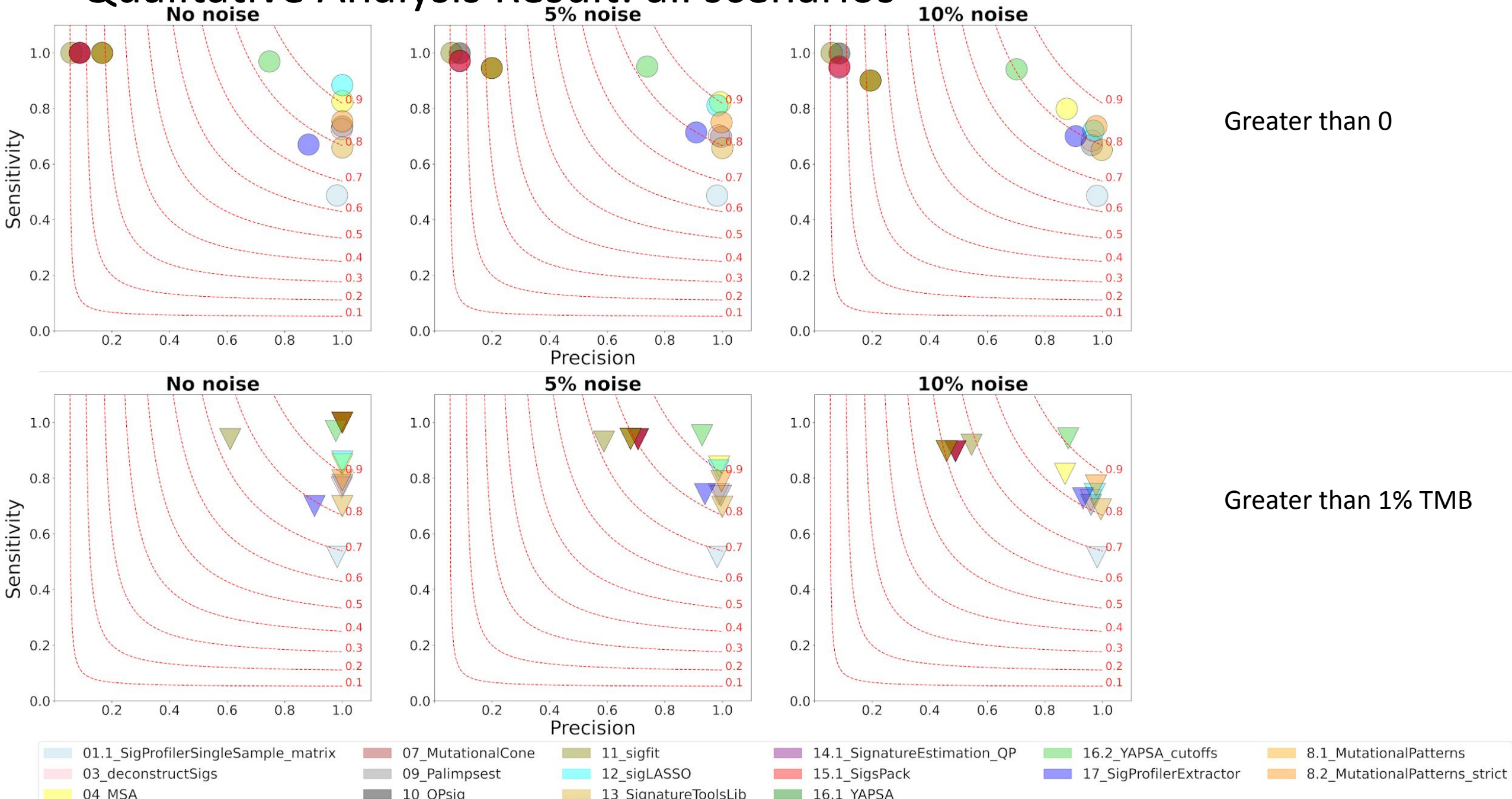
Qualitative Analysis Result: all scenarios



Cluster 1: QPSig, SigsPack,
SignatureEstimation_QP

Cluster 2: MutationalPatterns,
YAPSA (normal), MutationalCone

Qualitative Analysis Result: all scenarios



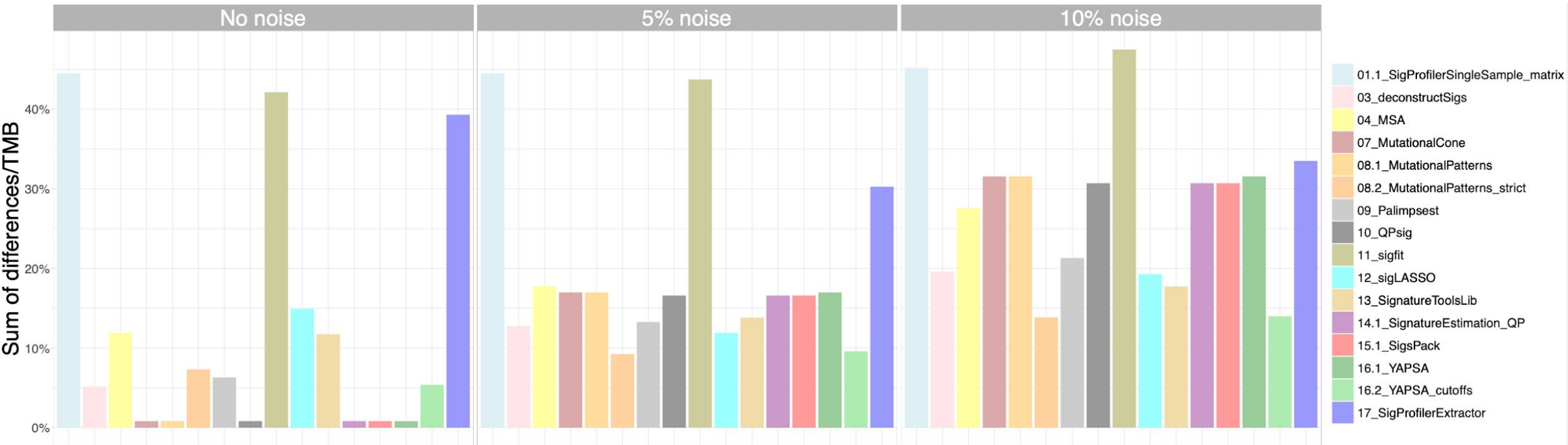
A step forward: Qualitative Analysis Method

- Simple!
- Compare the assignment results of each tools vs. ground truth
- By calculating:
 - Sum of absolute differences by TMB
 - percentage mutation mistakes detected by the tools
 - Average cos_sim scores
 - We are applying this on Activities matrix this time!

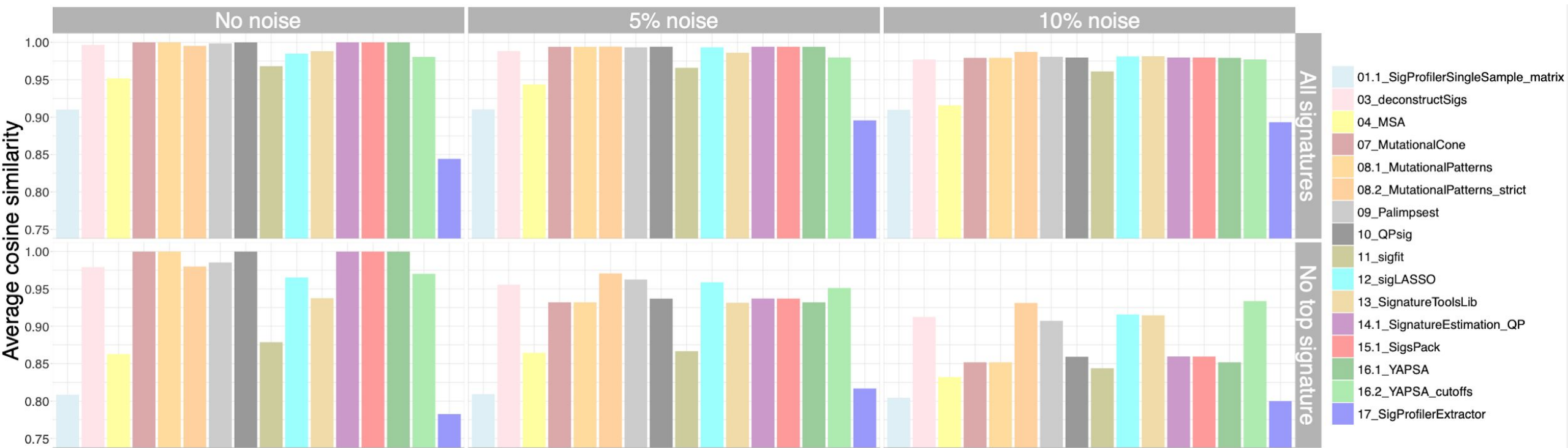
Sum of Absolute differences by TMB

	Ground Truth	Tool A	Abs. diff.	Abs. diff / TMB
SBS1	554.2989207	2000	1445.70108	19.6%
SBS2	204.7433281	200	4.74332809	0.1%
SBS3	0	0	0	0.0%
SBS5	4410.984607	4000	410.984607	5.6%
SBS13	329.1668878	300	29.1668878	0.4%
SBS17a	0	0	0	0.0%
SBS17b	0	0	0	0.0%
SBS18	1894.120409	1500	394.120409	5.3%
SBS28	0	0	0	0.0%
SBS30	0	0	0	0.0%
SBS40	0	1000	1000	13.5%
TMB	7393.314152			44.4%

Quantitative Analysis Result: sum of absolute differences by TMB



Qualitative Analysis Result: average cos_sim



Conclusion

- Consistent patterns are shown across all analysis
- All tools' accuracies decrease with noise in the samples
- Next step: need to do all the VCF tools and tools that require lots of time, finish run_time analysis, then possible finish benchmarking!

Future Impacts!

- Know where to improve
- Improve our own tool!
- Improve accuracy -> clinically applicable -> happy patients!

Thank you! Questions?

THE CHOICE IS YOURS

