# UMI-Gen: a new UMI-based read simulator
## for variant calling evaluation in paired-end sequencing

*Vincent SATER*
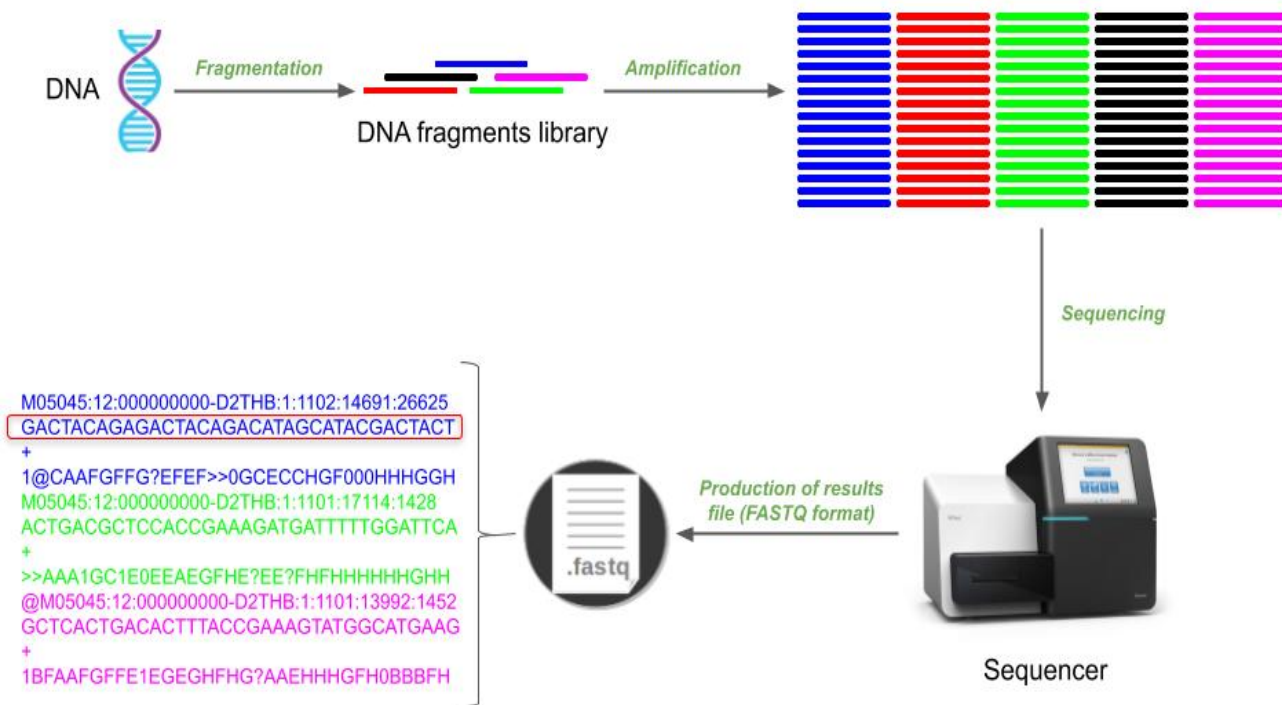
Thierry LECROQ, Pierre-Julien VIAILLY, Elise PRIEUR-GASTON, Philippe RUMINY, Caroline BÉRARD and Fabrice JARDIN.

**SeqBIM 2020**
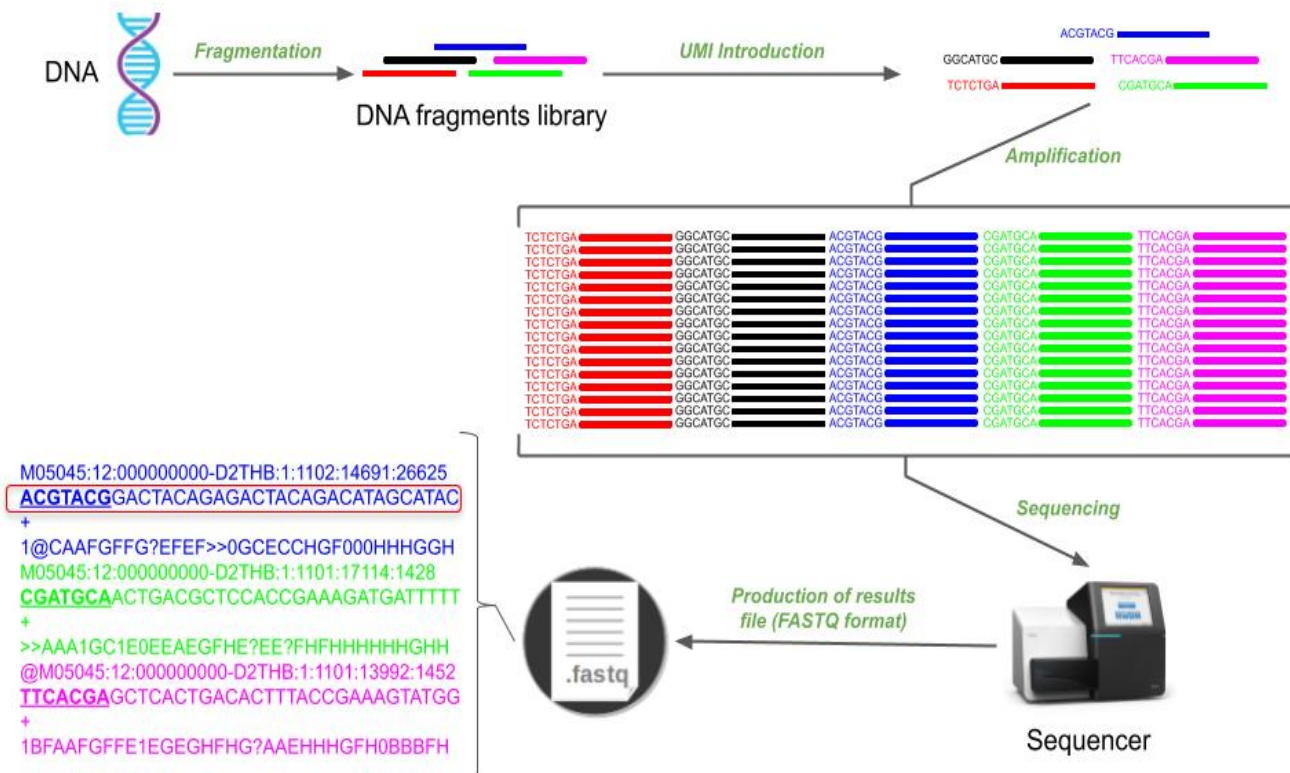
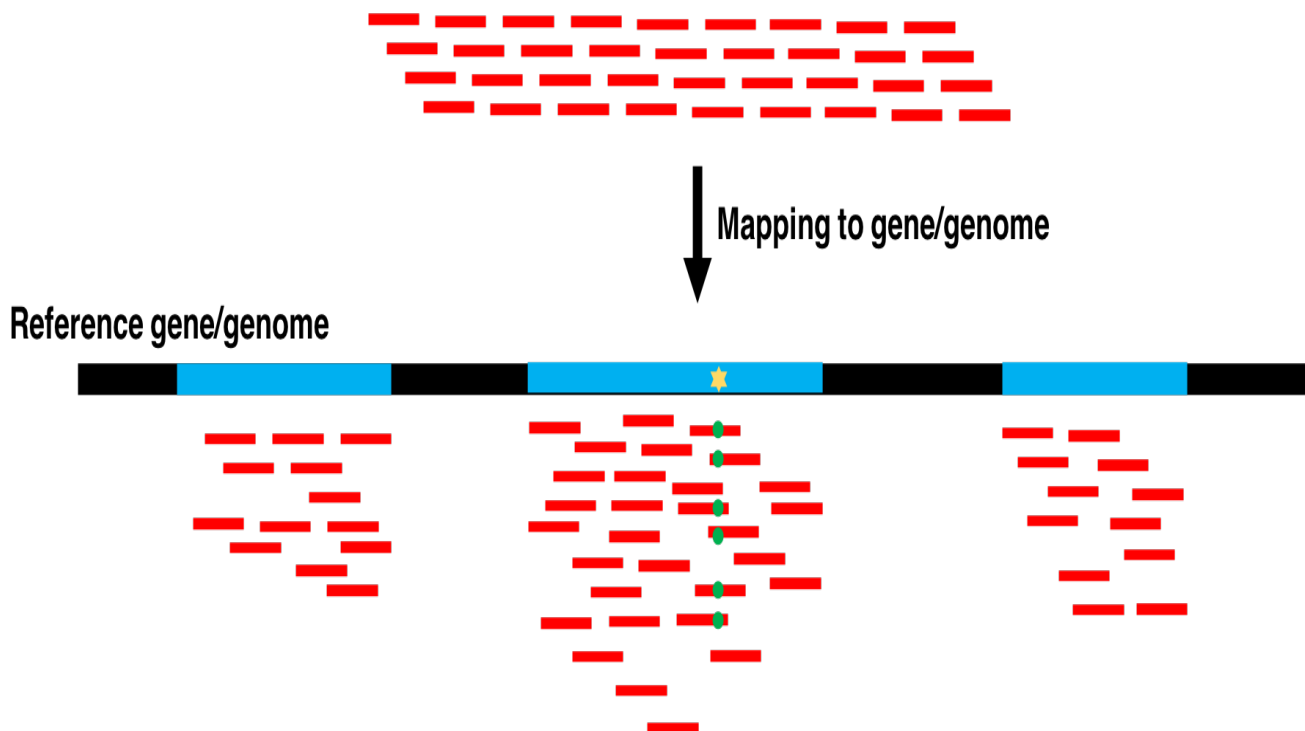23 November 2020

# Next Generation Sequencing (NGS)

# What are UMIs ?

- Unique Molecular Identifier

- Short arbitrary nucleotide sequences

- Increased usage in Next Generation Sequencing since 2015

# UMIs in NGS

# Aligning sequences to the reference genome

# Variant Calling

# The theory

# The theory

# The theory

# The theory

| Introduction | UMI-Gen | Results & Sample Validation | Application | Conclusion |
|---|---|---|---|---|
| ○○○○○○●○○ | ○○○ | ○○○○○○ | ○○○ | ○ |

Introduction

# The theory

# The theory

# The problem

Evaluating the performance of 3 variant callers VC1, VC2 and VC3 on the sample S1

| | VC1 | VC2 | VC3 |
|---|---|---|---|
| S1 | 142 | 99 | 210 |

# The problem

Evaluating the performance of 3 variant callers VC1, VC2 and VC3 on the sample S1

|      | VC1 | VC2 | VC3 |
|------|-----|-----|-----|
| **S1** | 142 | 99  | 210 |

No information about the variants present in S1 !!!

# The problem

Evaluating the performance of 3 variant callers VC1, VC2 and VC3 on the sample S1

|  | VC1 | VC2 | VC3 |
|---|---|---|---|
| **S1** | 142 | 99 | 210 |

False positives     ???

# The problem

Evaluating the performance of 3 variant callers VC1, VC2 and VC3 on the sample S1

|  | VC1 | VC2 | VC3 |
|---|---|---|---|
| S1 | 142 | 99 | 210 |

False positives ???

False negatives ???

# The problem

Evaluating the performance of 3 variant callers VC1, VC2 and VC3 on the sample S1

|     | VC1 | VC2 | VC3 |
| --- | --- | --- | --- |
| **S1** | 142 | 99  | 210 |

False positives   ???

False negatives  ???

True sensitivity   ???

# The problem

Evaluating the performance of 3 variant callers VC1, VC2 and VC3 on the sample S1

|  | VC1 | VC2 | VC3 |
|---|---|---|---|
| **S1** | 142 | 99 | 210 |

False positives　???

False negatives　???

True sensitivity　???

True specificity　???

# The problem

Evaluating the performance of 3 variant callers VC1, VC2 and VC3 on the sample S1

|  | VC1 | VC2 | VC3 |
|---|---|---|---|
| S1 | 142 | 99 | 210 |

False positives ???

False negatives ???

True sensitivity ???

True specificity ???

→ Poor & biased comparison

# The solution

- Using a read simulator that mimic real tumor samples

- Real variants in the produced samples must be known

- No UMI-based read simulators available → we developed UMI-Gen

| Introduction | UMI-Gen | Results & Sample Validation | Application | Conclusion |
|---|---|---|---|---|
| ○○○○○○○○○ | ●○○ | ○○○○○○ | ○○○ | ○ |

UMI-Gen

# UMI-Gen's workflow



Caption:
- artifact  x variant  ▪ UMI tag 1  ▢ UMI tag 2  ⊞ UMI tag 3  ▨ UMI tag 4  ▨ UMI tag 5

| Introduction | UMI-Gen | Results & Sample Validation | Application | Conclusion |
|---|---|---|---|---|
| ○○○○○○○○ | ○●○ | ○○○○○○ | ○○○ | ○ |

UMI-Gen

# Background noise estimation

# Adding the variants

Introduction
○○○○○○○○○

UMI-Gen
○○●

Results & Sample Validation
○○○○○○

Application
○○○

Conclusion
○

UMI-Gen

# Adding the variants

| Introduction | UMI-Gen | Results & Sample Validation | Application | Conclusion |
|---|---|---|---|---|
| ○○○○○○○○○ | ○○○ | ●○○○○○ | ○○○ | ○ |

Results & Sample Validation

# Sample Production



list of **13** know mutations — CSV

+ FASTA reference
+ background noise

→ sample **1** (.bam)

list of **15** know mutations — CSV

→ sample **2** (.bam)

| sample | depth | inserted mutations frequencies | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0.90 | 0.80 | 0.70 | 0.60 | 0.5 | 0.4 | 0.3 | 0.2 | 0.1 | 0.05 | 0.01 | 0.005 | 0.001 |
| sample 1 | 1000 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 3 | 1 | 1 | **X** | **X** |
| sample 2 | 10 000 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 3 | 1 | 1 | 1 | 1 |

# Noise reproduction from control samples

|           | A | C  | G  | T   | Total |
|-----------|---|----|----|-----|-------|
| Control 1 | 0 | 11 | 10 | 874 | 895   |
| Control 2 | 0 | 1  | 7  | 843 | 851   |
| Control 3 | 0 | 2  | 2  | 860 | 864   |
| Control 4 | 0 | 6  | 9  | 965 | 980   |
| Control 5 | 1 | 2  | 4  | 867 | 874   |
| Control 6 | 3 | 2  | 2  | 880 | 887   |

# Noise reproduction from control samples

| | A | C | G | T |
|---|---|---|---|---|
| Control 1 | 0 | 0.01229 | 0.01117 | 0.97654 |
| Control 2 | 0 | 0.00117 | 0.00823 | 0.9906 |
| Control 3 | 0 | 0.00232 | 0.00231 | 0.99537 |
| Control 4 | 0 | 0.00611 | 0.00918 | 0.98469 |
| Control 5 | 0.00113 | 0.00228 | 0.00458 | 0.99199 |
| Control 6 | 0.00338 | 0.00226 | 0.00225 | 0.99211 |

Introduction
00000000

UMI-Gen
000

Results & Sample Validation
○●○○○○○

Application
000

Conclusion
○

Results & Sample Validation

# Noise reproduction from control samples

|           | A       | C       | G       | T       |
|-----------|---------|---------|---------|---------|
| Control 1 | 0       | 0.01229 | 0.01117 | 0.97654 |
| Control 2 | 0       | 0.00117 | 0.00823 | 0.9906  |
| Control 3 | 0       | 0.00232 | 0.00231 | 0.99537 |
| Control 4 | 0       | 0.00611 | 0.00918 | 0.98469 |
| Control 5 | 0.00113 | 0.00228 | 0.00458 | 0.99199 |
| Control 6 | 0.00338 | 0.00226 | 0.00225 | 0.99211 |
| **Average** | **0.00075** | **0.00441** | **0.00629** | **0.98855** |

# Noise reproduction from control samples

|  | A | C | G | T |
|---|---|---|---|---|
| Control 1 | 0 | 0.01229 | 0.01117 | 0.97654 |
| Control 2 | 0 | 0.00117 | 0.00823 | 0.9906 |
| Control 3 | 0 | 0.00232 | 0.00231 | 0.99537 |
| Control 4 | 0 | 0.00611 | 0.00918 | 0.98469 |
| Control 5 | 0.00113 | 0.00228 | 0.00458 | 0.99199 |
| Control 6 | 0.00338 | 0.00226 | 0.00225 | 0.99211 |
| **Average** | **0.00075** | **0.00441** | **0.00629** | **0.98855** |
| Theoretical S1 | 1 | 4 | 6 | 989 |

# Noise reproduction from control samples

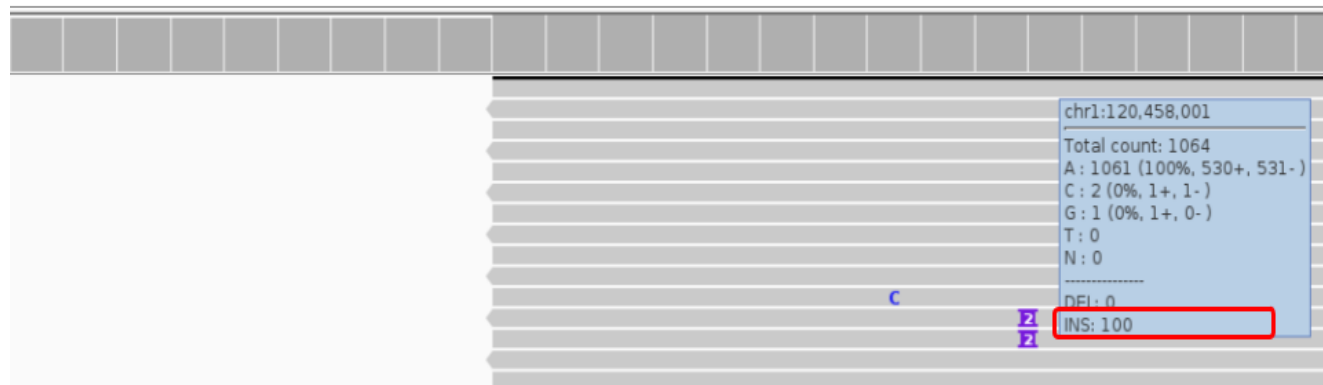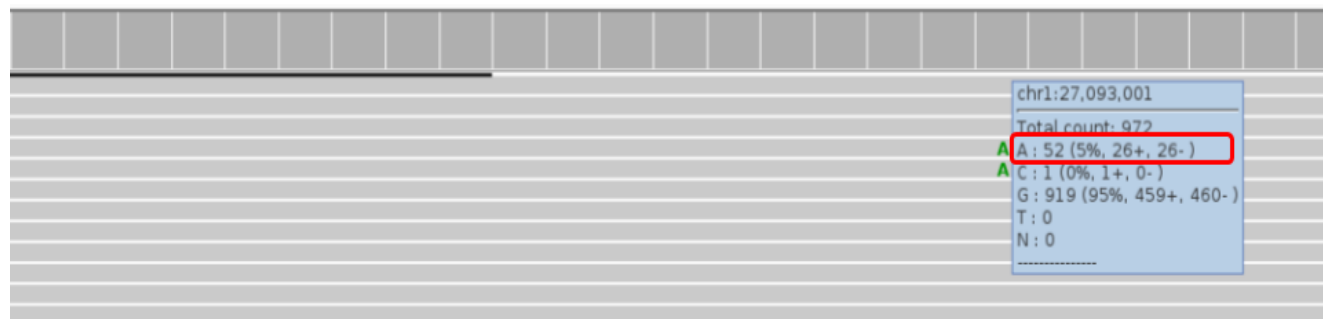|           | A       | C       | G       | T       |
|-----------|---------|---------|---------|---------|
| Control 1 | 0       | 0.01229 | 0.01117 | 0.97654 |
| Control 2 | 0       | 0.00117 | 0.00823 | 0.9906  |
| Control 3 | 0       | 0.00232 | 0.00231 | 0.99537 |
| Control 4 | 0       | 0.00611 | 0.00918 | 0.98469 |
| Control 5 | 0.00113 | 0.00228 | 0.00458 | 0.99199 |
| Control 6 | 0.00338 | 0.00226 | 0.00225 | 0.99211 |
| **Average** | **0.00075** | **0.00441** | **0.00629** | **0.98855** |
| Theoretical S1 | 1 | 4 | 6 | 989 |
| Theoretical S2 | 8 | 44 | 63 | 9885 |

# Noise in the produced samples

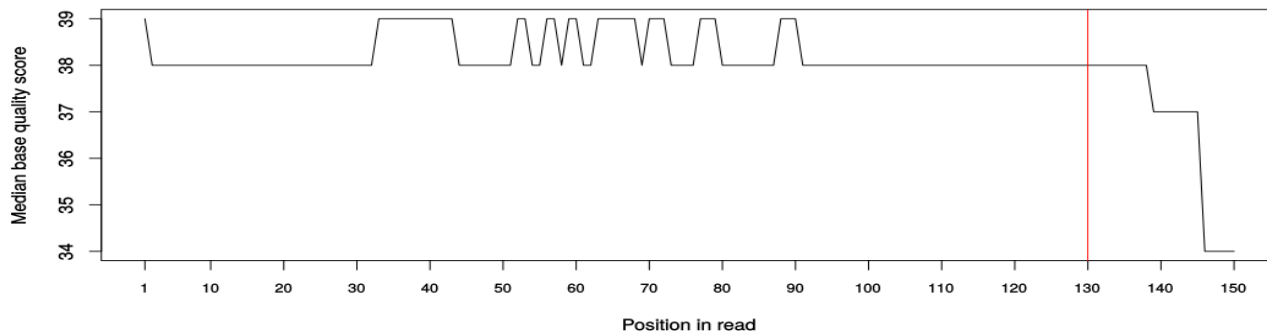# Real variants in the produced samples

# Real variants in the produced samples

Introduction
○○○○○○○○○

UMI-Gen
○○○

Results & Sample Validation
○○○○○●○

Application
○○○

Conclusion
○

Results & Sample Validation

# Quality score reproduction

Introduction
○○○○○○○○○

UMI-Gen
○○○

Results & Sample Validation
○○○○○●○

Application
○○○

Conclusion
○

Results & Sample Validation

# Quality score reproduction

Introduction
00000000
UMI-Gen
000
Results & Sample Validation
000000●
Application
000
Conclusion
○

Results & Sample Validation

# %GC comparison

Introduction
00000000

UMI-Gen
000

Results & Sample Validation
000000

Application
●00

Conclusion
○

Application

# Evaluation on 4 variant callers

- 2 raw-read-based variant callers: SiNVICT & OutLyzer

- 2 UMI-based variant callers: UMI-VarCal & DeepSNVMiner

- Sample 1 has a depth of 1000x with 13 variants

- Sample 2 has a depth of 10,000x with 15 variants

Introduction
00000000

UMI-Gen
000

Results & Sample Validation
000000

Application
0●0

Conclusion
0

Application

# Evaluation on 4 variant callers

sample 1

| Variant Caller | True positives | False Positives | False Negatives | Sensitivity (%) | Specificity (%) |
|---|---|---|---|---|---|
| SiNVICT | 8 | 233 | 5 | 61.5 | 99.7 |
| OutLyzer | | | | | |
| DeepSNV Miner | | | | | |
| UMI-VarCal | | | | | |

sample 2

| Variant Caller | True positives | False Positives | False Negatives | Sensitivity (%) | Specificity (%) |
|---|---|---|---|---|---|
| SiNVICT | 8 | 455 | 7 | 53.4 | 99.4 |
| OutLyzer | | | | | |
| DeepSNV Miner | | | | | |
| UMI-VarCal | | | | | |

# Evaluation on 4 variant callers

sample 1

| Variant Caller | True positives | False Positives | False Negatives | Sensitivity (%) | Specificity (%) |
|---|---|---|---|---|---|
| SiNVICT | 8 | 233 | 5 | 61.5 | 99.7 |
| OutLyzer | 11 | 98 | 2 | 84.6 | 99.9 |
| DeepSNV Miner | | | | | |
| UMI-VarCal | | | | | |

sample 2

| Variant Caller | True positives | False Positives | False Negatives | Sensitivity (%) | Specificity (%) |
|---|---|---|---|---|---|
| SiNVICT | 8 | 455 | 7 | 53.4 | 99.4 |
| OutLyzer | 12 | 330 | 3 | 80 | 99.6 |
| DeepSNV Miner | | | | | |
| UMI-VarCal | | | | | |

# Evaluation on 4 variant callers

sample 1

| Variant Caller | True positives | False Positives | False Negatives | Sensitivity (%) | Specificity (%) |
|---|---|---|---|---|---|
| SiNVICT | 8 | 233 | 5 | 61.5 | 99.7 |
| OutLyzer | 11 | 98 | 2 | 84.6 | 99.9 |
| DeepSNV Miner | 12 | 37 | 1 | 92.3 | 99.95 |
| UMI-VarCal | | | | | |

sample 2

| Variant Caller | True positives | False Positives | False Negatives | Sensitivity (%) | Specificity (%) |
|---|---|---|---|---|---|
| SiNVICT | 8 | 455 | 7 | 53.4 | 99.4 |
| OutLyzer | 12 | 330 | 3 | 80 | 99.6 |
| DeepSNV Miner | 14 | 2 | 1 | 93.4 | 99.99 |
| UMI-VarCal | | | | | |

# Evaluation on 4 variant callers



sample 1

| Variant Caller | True positives | False Positives | False Negatives | Sensitivity (%) | Specificity (%) |
|---|---|---|---|---|---|
| SiNVICT | 8 | 233 | 5 | 61.5 | 99.7 |
| OutLyzer | 11 | 98 | 2 | 84.6 | 99.9 |
| DeepSNV Miner | 12 | 37 | 1 | 92.3 | 99.95 |
| UMI-VarCal | 13 | 0 | 0 | 100 | 100 |



sample 2

| Variant Caller | True positives | False Positives | False Negatives | Sensitivity (%) | Specificity (%) |
|---|---|---|---|---|---|
| SiNVICT | 8 | 455 | 7 | 53.4 | 99.4 |
| OutLyzer | 12 | 330 | 3 | 80 | 99.6 |
| DeepSNV Miner | 14 | 2 | 1 | 93.4 | 99.99 |
| UMI-VarCal | 15 | 0 | 0 | 100 | 100 |

Introduction
○○○○○○○○○

UMI-Gen
○○○

Results & Sample Validation
○○○○○○

Application
○○●

Conclusion
○

Application

# Performance



Variation of UMI-Gen's performance with sample depth

Introduction
○○○○○○○○○

UMI-Gen
○○○

Results & Sample Validation
○○○○○○

Application
○○○

Conclusion
●

# Conclusion

- UMI-based read simulator that is capable of producing artificial samples in which real variants are known

- Most of the parameters are customizable which allows total control over the simulation

- Generates FASTQ, BAM and SAM files

- Published in the Computational and Structural Biotecnology Journal and available for Python 3 on GitLab: https://gitlab.com/vincent-sater/umigen

# Thank you for the attention