# KSBi-BIML 2024

**Bioinformatics & Machine Learning(BIML)**
**Workshop for Life and Medical Scientists**

생명정보학 & 머신러닝 워크샵 (온라인)

# Mutational signatures
# in cancer genomes

주영석 _ KAIST

# KSBi-BIML 2024

## Bioinformatics & Machine Learning(BIML)
## Workshop for Life and Medical Scientists

안녕하십니까?

한국생명정보학회가 개최하는 동계 교육 워크샵인 BIML-2024에 여러분을 초대합니다. 생명정보학 분야의 연구자들에게 최신 동향의 데이터 분석기술을 이론과 실습을 겸비해 전달하고자 도입한 전문 교육 프로그램인 BIML 워크샵은 2015년에 시작하여 올해로 벌써 10년 차를 맞이하게 되었습니다. BIML 워크샵은 국내 생명정보학 분야의 최초이자 최고 수준의 교육프로그램으로 크게 인공지능과 생명정보분석 두 개의 분야로 구성되어 있습니다. 올해 인공지능 분야에서는 최근 생명정보 분석에서도 응용이 확대되고 있는 다양한 인공지능 기반 자료모델링 기법들에 대한 현장 강의가 진행될 예정이며, 관련하여 심층학습을 이용한 단백질구조예측, 유전체분석, 신약개발에 대한 이론과 실습 강의가 함께 제공될 예정입니다. 또한 단일세포오믹스, 공간오믹스, 메타오믹스, 그리고 롱리드염기서열 자료 분석에 대한 현장 강의는 많은 연구자의 연구 수월성 확보에 큰 도움을 줄 것으로 기대하고 있습니다.

올해 BIML의 가장 큰 변화는 최근 연구 수요가 급증하고 있는 의료정보자료 분석에 대한 현장 강의를 추가하였다는 것입니다. 특히 의료정보자료 분석을 많이 수행하시는 의과학자 및 의료정보 연구자들께서 본 강좌를 통해 많은 도움을 받으실 수 있기를 기대하고 있습니다. 또한 다양한 생명정보학 분야에 대한 온라인 강좌 프로그램도 점차 증가하고 있는 생명정보 분석기술의 다양화에 발맞추기 위해 작년과 비교해 5강좌 이상을 신규로 추가했습니다. 올해는 무료 강좌 5개를 포함하여 35개 이상의 온라인 강좌가 개설되어 제공되며, 연구 주제에 따른 연관된 강좌 추천 및 강연료 할인 프로그램도 제공되며, 온라인을 통한 Q&A 세션도 마련될 예정입니다. BIML-2024는 국내 주요 연구 중심 대학의 전임 교원이자 각 분야 최고 전문가들의 강의로 구성되었기에 해당 분야의 기초부터 최신 연구 동향까지 포함하는 수준 높은 내용의 강의가 될 것이라 확신합니다.

BIML-2024을 준비하기까지 너무나 많은 수고를 해주신 운영위원회의 정성원, 우현구, 백대현, 김태민, 김준일, 김상우, 장혜식, 박종은 교수님과 KOBIC 이병욱 박사님께 커다란 감사를 드립니다. 마지막으로 부족한 시간에도 불구하고 강의 부탁을 흔쾌히 허락하시고 훌륭한 현장 강의와 온라인 강의를 준비하시는데 노고를 아끼지 않으신 모든 강사분들께 깊은 감사를 드립니다.

2024년 2월

**한국생명정보학회장 이 인 석**

# Mutational signatures in cancer genomes

Cancer genome sequencing을 이용하면 우리는 무엇을 배울 수 있을까? 1차적으로는 최적화 약제를 선별하기 위한 cancer driver mutation을 찾기 위한 목표로 쓰인다. 하지만 Cancer genome에서 나오는 수 많은 돌연변이의 pattern, 즉 mutational signature를 체계적으로 분석하면 정상세포에서 암 세포로 돌변하는 과정중에서 돌연변이들을 만들어낸 기전을 이해할 수 있다.

본 강의에서는 암 세포에서 발견한 돌연변이로부터 mutational signature를 빠르게 추출하고 분석하는 방법을 설명한다. Mutational signature의 개념, signature를 calling하는 알고리즘 및 툴을 소개하며, 이를 실제 암 유전체 데이터에 적용하여 효율적이고 효과적인 분석을 할 수 있는 핵심 역량을 갖추는 것을 목표로 한다.


강의는 다음의 내용을 포함한다:

- Mutational signature 의 개념
- Mutational signature 의 calling algorithm 및 tools


* 강의 난이도: 초급


* 강의: 주영석 교수 (KAIST 의과학대학원)

## Curriculum Vitae

## Speaker Name: Young Seok Ju, M.D. Ph.D.

▶ **Personal Info**

| | |
|---|---|
| Name | Young Seok Ju |
| Title | Associate Professor |
| Affiliation | Grad School of Medical Science and Engineering, KAIST |

▶ **Contact Information**

| | |
|---|---|
| Address | 291 Daehak-ro Yuseong-gu, Daejeon 34141 |
| Email | ysju@kaist.ac.kr |
| Phone Number | 042-350-4237 |

---

**Research Interest**

Somatic mutation, somatic mosaicism, bioinformatics, mutational process

**Educational Experience**

| | |
|---|---|
| 2007 | M.D. in Medicine, Seoul Nat'l Univ College of Medicine, Seoul, Korea |
| 2010 | Ph.D. in Genomic Medicine, Seoul Nat'l Univ College of Medicine, Seoul, Korea |

**Professional Experience**

| | |
|---|---|
| 2013-2015 | Post-doc, Wellcome Sanger Institute, Daejeon, Korea |
| 2015- | Associate/Assistant Professor, KAIST |

**Selected Publications (5 maximum)**

1. Park S, Mali NM, Kim R, Choi JW, Lee J*,…,Oh J#, , **Ju YS#**. Clonal dynamics in early human embryogenesis inferred from somatic mutation. *Nature* 2021

2. Youk J*, Kim T*, Evans KV*, Jeong Y-I*, Hur Y*, Hong SP*, …, Kim YT#, Koh GY#, Choi B-S#, **Ju YS#**, Lee JH#. Three-dimensional human alveolar stem cell culture models reveal infection response to SARS-CoV-2. *Cell Stem Cell*. 2020

3. Lee JS, An Y, Yoon CJ, Kim JY, Kim KH, … , Lee EY# & **Ju YS#**. Germline gain-of-function mutation of STAT1 rescued by somatic mosaicism in immune dysregulation-polyendocrinopathy-enteropathy-X-linked-like disorder. *J Allergy Clin Immunol*. 2020

4. Lee JJ-K, Park S, Park H, Kim S, Lee J, … , **Ju YS#** & Kim YT#. Tracing oncogene rearrangements in the mutational history of lung adenocarcinoma. *Cell*. 2019

5. Lee JK., Lee J, Kim S, Kim S, Youk J, …, Kim TM# & **Ju YS#**. Clonal history and genetic predictors of transformation into small cell carcinomas from lung adenocarcinomas. *Journal of Clinical Oncology* 2017 Sep 10;35(26):3065-3074. PMID:28498782

# KSBi-BIML

Mutational signatures in cancer genomes

주영석 (KAIST)   ysju@kaist.ac.kr

---

## 분석의 목적: 왜 암 유전체를 분석하는가?

- 목적에 따라 다양한 접근법을 이용할 수 있음

  - 임상 의사: 환자 암에서 clinically actionable target을 발굴, 진료에 응용
    (EGFR activating mutation 발굴)

  - Genomics, Bioinformatics 에 관심이 있는 학부생, 대학원생, 박사 후 연구원 등
    새로운 돌연변이 발굴, technology/bioinformatics 개발, 논문 출판

  - 회사나 연구소의 전문 연구원
    Pipeline 구축 등

# 암 유전체 분석의 시작: 돌연변이의 검출



Figure 1 | The lineage of mitotic cell divisions from the fertilized egg to a single cell within a cancer showing the timing of the somatic mutations acquired by the cancer cell and the processes that contribute to them. Mutations may be acquired while the cell lineage is phenotypically normal, reflecting both the intrinsic mutations acquired during normal cell division and the effects of exogenous mutagens. During the development of the cancer other processes, for example DNA repair defects, may contribute to the mutational burden. Passenger mutations do not have any effect on the cancer cell, but driver mutations will cause a clonal expansion. Relapse after chemotherapy can be associated with resistance mutations that often predate the initiation of treatment.

- 대부분의 산발성 암 (sporadic cancer) 의 원인은 체세포 돌연변이이다

Stratton *et al.*, *Nature* (2009)

# 돌연변이의 검출을 위한 전략



Meyerson M *et al.*, *Nat Rev Cancer* (2010)

International consortia for cancer genome analyses



Driver mutation을 찾는 것이 암유전체 분석의 한 목적

# 어떤 돌연변이가 있는가?

- 크기에 의한 분류
  - Small (point-mutation):
    - base substitution (SNV, SNP), short-indel
  - Large:
    - Copy number variation, genome rearrangements, SV

- 유전체 위치에 의한 분류
  - Coding mutation (in the protein coding region)
    - Non-sense/frameshift (truncating, stop-gain)
    - Missense (non-synonymous)
    - Silent (synonymous)
  - UTR, intronic, splicing-junction
  - intergenic (between two genes)



# Cancer genome에서 driver mutation의 분포

# Driver mutations in pan-cancer genomes

# An example of genome-wide sequencing of a cancer genome

# 암유전체의 돌연변이들



Alexandrov L *et al., **Nature*** (2013)

- WGS (3,000 Mb) ➜ 3,000 (1,000 – 100,000 substitutions)
- WES (~50 Mb) ➜ 50 (10 – 1,000 substitutions)
- Targeted-gene seq. (covers ~1 Mb) ➜ 10 (1 – 100 substitutions)

---

# Cancer genomics에서 passenger mutation은 쓸모가 없는가?



Stratton *et al., **Nature*** (2009)

driver

passengers

# An example of base substitution



We are looking at this position
(3q26.32)

Chr3

57 bp

matched normal sequence

Blood sample has the wildtype (G) allele

Lung adenocarcinoma sequence

G to A mutation in PIK3CA gene in cancer

PIK3CA gene

# Mutational signature 개념을 접하다



Genome Research (2012.3)

Ju YS *et al.*, Genome Res (2012a)

# Mutational signature 개념을 접하다 (2)



The transcriptional landscape and mutational profile of lung adenocarcinoma — Seo JS et al., Genome Research (2012)

Dr. Myles Axton
(former Chief Editor
@ Nature Genetics)

# Mutational signature 개념을 접하다 (3)



Ludmil B Alexandrov

Elizabeth P Murchison

## ARTICLE

### Signatures of mutational processes in human cancer

# mutational origin: Mutation은 랜덤하게 생기는 것이 아니다



돌연변이는 "랜덤" 이 아니라 DNA damage x DNA repair 과정

**6 classes of base substitutions**

C>A (G>T), C>G (G>C), C>T (G>A)

T>A (A>T), T>C (A>G), T>G (A>C)

**Spontaneous cytosine deamination**
C>T substitutions
(mostly at CpG context)

**Tobacco smoking**
C>A substitutions

**Ultraviolet (UV) light**
C>T substitutions
(CC>TT)



---

# Classical observation



- 암종마다, 그리고 발암물질의 노출에 따라서 TP53 유전자에 생기는 돌연변이 패턴이 상이하다

Greenblatt *et al.*, *Cancer Research* (1994)

## Mutational signature의 예시

- 폐암의 전장 유전체 분석에서 20,000개의 base substitution 발견
  - 이 가운데 80%가 C>A mutations. 주된 돌연변이 발생기전은?
  - (흡연에 노출)

- 흑색종의 전장 유전체 분석에서 20,000 개의 base substitution 발견
  - 이 가운데 90%가 C>T mutations 이고 수백개의 CC>TT 도 같이 발견
    주된 돌연변이 발생기전은?
  - (UV에 노출)

- 실제로는 하나의 암 유전체에서 발생하는 돌연변이들이
  위와 같은 단일 돌연변이 발생 기전이 아니라,
  여러 돌연변이 기전의 '조합' 으로 만들어지는 일이 훨씬 흔하다

## 실제 3개의 breast cancer whole-genome sequencing 결과



● C>A  ● C>G  ● C>T  ○ T>A  ● T>C  ● T>G

3000 base substitutions
~50% C>T
~30% T>C

20,000 base substitutions
~15% C>A  ~20% C>G
~20% C>T  ~15% T>A
~20% T>C  ~10% T>G

100,000 base substitutions
~40% C>G
~50% C>T

일반적인 ER+
breast cancer

BRCA1 mutant
breast cancer

APOBEC hyperactivated
breast cancer

## Tumor의 돌연변이 스펙트럼은 이론적으로 n개의 서로 다른 Mutational process로 설명된다

하지만 우리는 n이 얼마인지도, 각각의 process의 spectrum도 알고있지 못한다



Alexandrov L *et al.*, ***Cell reports*** (2013)

## Understanding mutational processes from mutational spectrum: a blind-source separation problem



Somatic mutations explored in a sample can be explained by linear sum of different exposures

With genome big-dataset

& using NMF
(or other equivalent algorithms)



Alexandrov L *et al.*, ***Cell reports*** (2013)

# Single base substitutions (SBS) into 96 subclasses

- C>A (G>T)
- C>G (G>C)
- C>T (G>A)
- T>A (A>T)
- T>C (A>G)
- T>G (A>C)

sequence context

5'B - Wt > Var- 3'B

3' immediate base

|   | A | C | G | T |
|---|---|---|---|---|
| A |   |   |   |   |
| C |   | C>T |   |   |
| G |   |   |   |   |
| T |   |   |   |   |

5' immediate base

4 x    6 types    x 4

= 96 subtypes

C>A    C>G    C>T    T>A    T>C    T>G

# SBS Signature 1 and Signature 2

# Dictionary for mutational signatures: COSMIC

# 49 +5 biologic signatures in SBS mutations (v3)

Signatures by patterns



Signatures by etiology

Extensive cell type specificity


Signature 1 and 5: basal, cellular intrinsic mutagenesis

# (1) SBS Signature 1: 5mC deamination



**Cancer types:**
Signature 1 has been found in all cancer types and in most cancer samples.
**Proposed etiology:**
Signature 1 is the result of an endogenous mutational process initiated by spontaneous deamination of 5-methylcytosine.
**Additional mutational features:**
Signature 1 is associated with small numbers of small insertions and deletions in most tissue types.
**Comments:**
The number of Signature 1 mutations correlates with age of cancer diagnosis.

# (1) SBS Signature 5: unknown mechanism



**Cancer types:**
Signature 5 has been found in all cancer types and most cancer samples.
**Proposed etiology:**
The aetiology of Signature 5 is unknown.
**Additional mutational features:**
Signature 5 exhibits transcriptional strand bias for T>C substitutions at ApTpN context.

# (1) SBS Signatures 1, 5; clock-like property



Sig 1, 5 mutations accumulate over time with similar rate

# (2) Signature 4: due to direct smoke exposure

# (2) SBS Signature 4: tobacco smoking



Signature of benzo[a]pyrene exposure in vitro

**Cancer types:**
Signature 4 has been found in head and neck cancer, liver cancer, lung adenocarcinoma, lung squamous carcinoma, small cell lung carcinoma, and esophageal cancer.

**Proposed etiology:**
Signature 4 is associated with smoking and its profile is similar to the mutational pattern observed in experimental systems exposed to tobacco carcinogens (e.g., benzo[a]pyrene).
Signature 4 is likely due to tobacco mutagens.

**Additional mutational features:**
Signature 4 exhibits transcriptional strand bias for C>A mutations, compatible with the notion that damage to guanine is repaired by transcription-coupled nucleotide excision repair. Signature 4 is also associated with CC>AA dinucleotide substitutions.

# (2) SBS Signature 4: mutational burden and strand bias



Fousteri and Mullenders (2008)

# (3) SBS Signature 7s: due to ultraviolet-light



Mutational signature present | Total validated mutational signatures in a cancer type | Total cancer types in which a signature is operative

# (3) Signature 7: ultraviolet-light damage



Old signature

**Cancer types:**
Signature 7 has been found predominantly in skin cancers and in cancers of the lip categorized as head and neck or oral squamous cancers.

**Proposed etiology:**
Based on its prevalence in ultraviolet exposed areas and the similarity of the mutational pattern to that observed in experimental systems exposed to ultraviolet light Signature 7 is likely due to ultraviolet light exposure.

**Additional mutational features:**
Signature 7 is associated with large numbers of CC>TT dinucleotide mutations at dipyrimidines. Additionally, Signature 7 exhibits a strong transcriptional strand-bias indicating that mutations occur at pyrimidines (viz., by formation of pyrimidine-pyrimidine photodimers) and these mutations are being repaired by transcription-coupled nucleotide excision repair.

# (3) Double base substitution (DBS1)



78 mutation classes



# (4) SBS Signatures 2 and 13: APOBEC-mediated mutagenesis

# (4) SBS Signatures 2 and 13: APOBEC-mediated mutagenesis



**Cancer types:**
Signature 2 has been found in 22 cancer types. Dominant processes in cervical and bladder cancers.
**Proposed etiology:**
Signature 2 has been attributed to activity of the **AID/APOBEC family of cytidine deaminases**.
On the basis of similarities in the sequence context of cytosine mutations caused by APOBEC enzymes in experimental systems, a role for APOBEC1, APOBEC3A and/or APOBEC3B in human cancer appears more likely than for other members of the family.

# APOBEC-mediated mutations



Alexandrov L *et al., Science* (2016)

Lee J *et al., J Clin Oncol* (2017)

Activated in many cancer types
including cervical, bladder, breast and lung cancers.

Activated in the late branch in lung cancers.
(Episodically activating?)

# (5) Signature 11: temozolomide-driven



# (5) SBS Signature 11: alkylating agent



**Cancer types:**
Signature 11 has been found in melanoma and glioblastoma.
**Proposed etiology:**
Signature 11 exhibits a mutational pattern resembling that of alkylating agents. Patient histories have revealed an association between treatments with the **alkylating agent temozolomide** and Signature 11 mutations.
**Additional mutational features:**
Signature 11 exhibits a strong transcriptional strand-bias for C>T substitutions indicating that mutations occur on guanine and that these mutations are effectively repaired by transcription-coupled nucleotide excision repair.

# (6) SBS Signature 22: aristolochic acid driven



Mutational signature present · Total validated mutational signatures in a cancer type · Total cancer types in which a signature is operative

# (6) SBS Signature 22: aristolochic acids



**Cancer types:**
Signature 22 has been found in urothelial (renal pelvis) carcinoma and liver cancers.

**Proposed aetiology:**
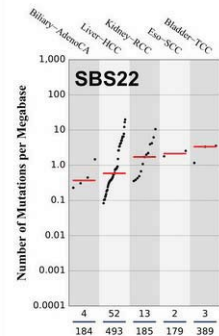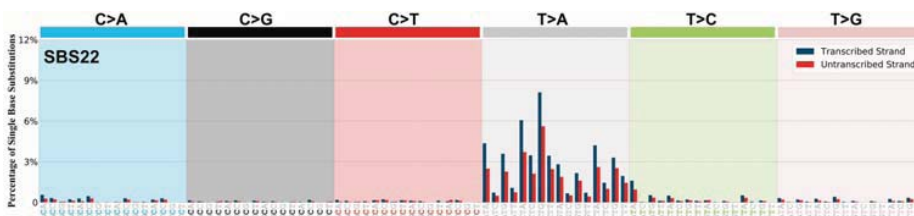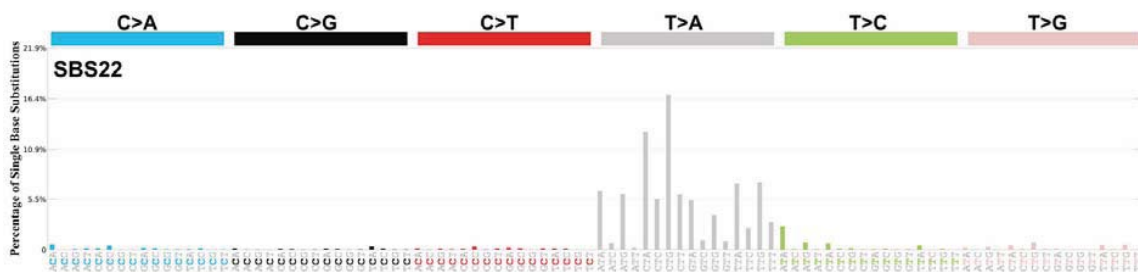Signature 22 has been found in cancer samples with known exposures to aristolochic acid.
Additionally, the pattern of mutations exhibited by the signature is consistent with the one
previous observed in experimental systems exposed to aristolochic acid.
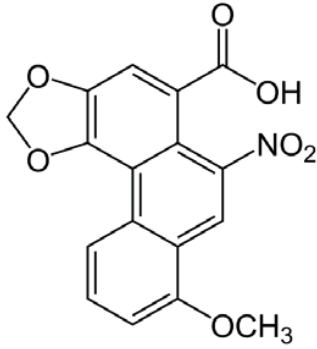
**Additional mutational features:**
Signature 22 exhibits a very strong transcriptional strand bias for T>A mutations indicating adenine damage
that is being repaired by transcription-coupled nucleotide excision repair.
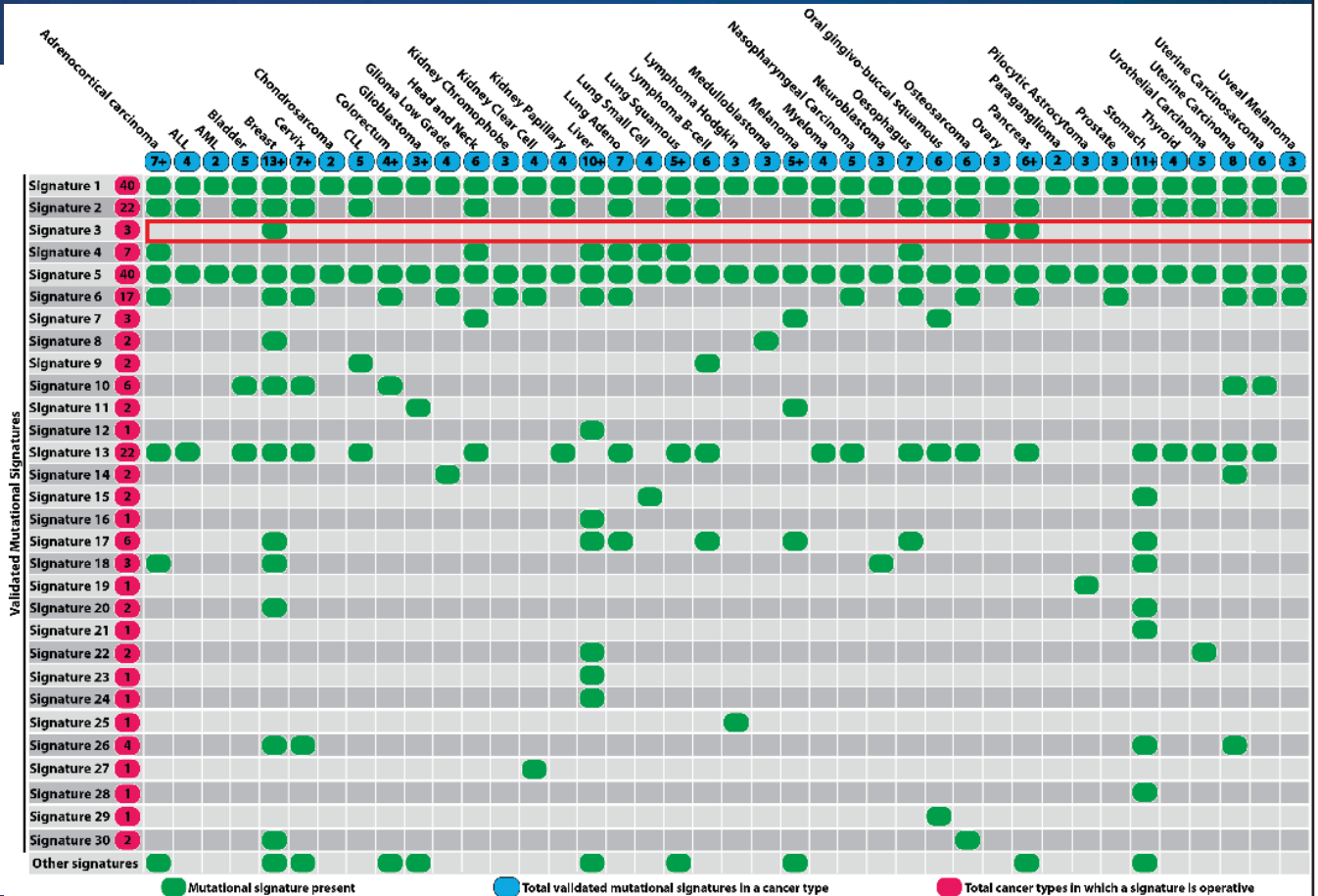
**Comments:**
Signature 22 has a very high mutational burden in urothelial carcinoma; however, its mutational burden is much lower in liver cancers.
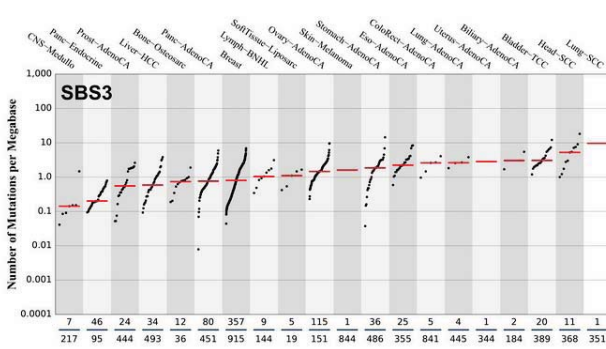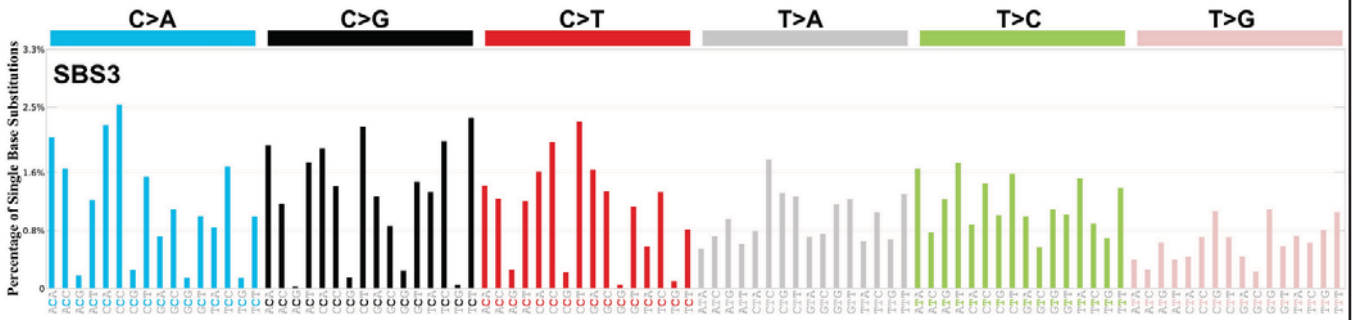
## (6) SBS signature 22: aristolochic acids

Aristolochia clematitis
(쥐방울덩굴, 동북마두령, 관목통)

https://www.accessdata.fda.gov/cms_ia/importalert_141.html

## (7) SBS Signature 3: HR-based DNA repair



Mutational signature present — Total validated mutational signatures in a cancer type — Total cancer types in which a signature is operative
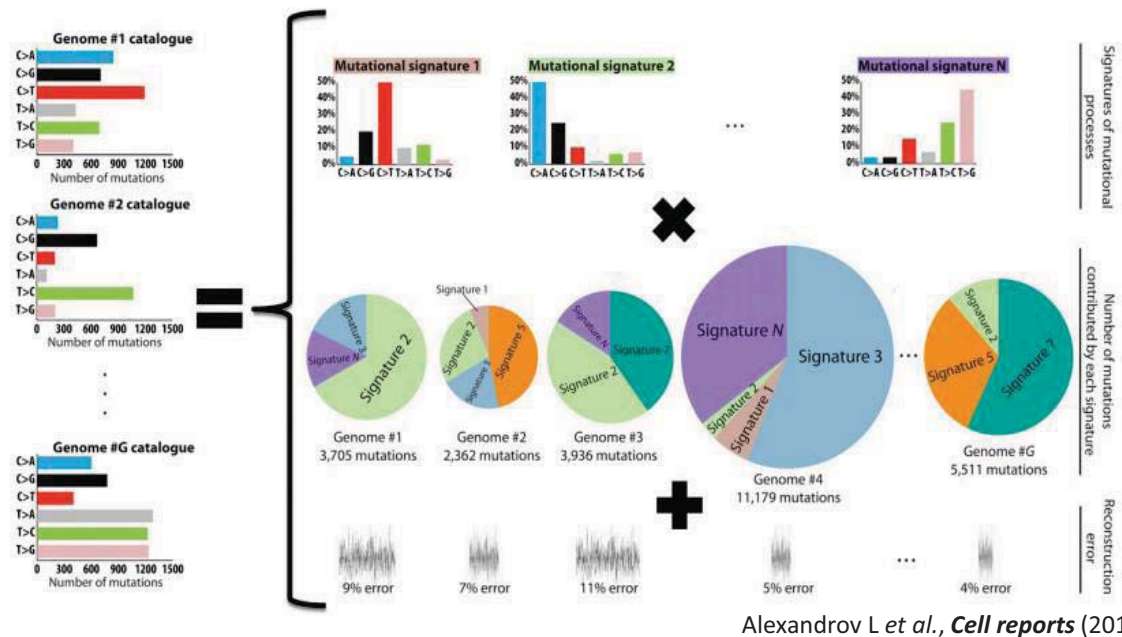
**Proposed aetiology**
**Defective homologous recombination-based DNA damage repair** which manifests predominantly as small indels and genome rearrangements due to abnormal double strand break repair but also in the form of this base substitution signature.
**Comments**
SBS3 is strongly associated with germline and somatic BRCA1 and **BRCA2 mutations and BRCA1** promoter methylation in breast, pancreatic, and ovarian cancers. In pancreatic cancer, responders to platinum therapy usually exhibit SBS3 mutations. Together with associated indel and rearrangement signatures, SBS3 has been proposed as a predictor of defective homologous recombination-based repair and thus of response to therapies exploiting this repair defect.

---

# 어떻게 mutational signature를 구할 것인가?

## 몇 개의 sample, 몇 개의 돌연변이가 필요할까?

Alexandrov L *et al.*, ***Cell reports*** (2013)

(1) For inferring *de novo* mutational signature: many whole-genome sequences
(2) For fitting known signatures: whole-genome, (exome?)

---

## Tools for extracting mutational signatures

### Inferring *de novo* signatures

Alexandrov, MatLab (***Nature*** 2013)
EMu (***Genome Biology*** 2013)
Maftools (***Genome Res*** 2018)
MutationalPatterns (***Genome Med*** 2018)
MutSpec (***BMC Bioinformatics*** 2016)
SigFit (***BioRxiv*** 2020)
SigMiner (***medRxiv*** 2020)
SignatureAnalyzer (***Nature Commun*** 2015)
SignatureToolsLib (***Nat Cancer*** 2020)
SigneR (***Bioinformatics*** 2017)
SomaticSignatures (***Bioinformatics*** 2015)
SigProfiler (COSMIC)

### Fitting known signatures

deconsructSigs (***Genome Biology*** 2016)
SignatureEstimation (***Bioinformatics*** 2018)
YAPSA (R Package v 1.16.0)

### Web interfaces

MutaGene (***NAR*** 2017)
mSignatureDB (***NAR*** 2018)
MuSiCa (***BMC Bioinformatics*** 2018)
Mutalisk (***NAR*** 2018)

# (1) Sigprofilers



# (1) Sigprofiler tools

# (1-1) Sigprofiler matrix generator



Python and R

(input)
VCF/MAF

(output)
Matrices with
Sequencing context
Transcriptional strand bias

Bergstrom et al., *BMC Genomics* (2019)

# (1-2) Sigprofiler Extractor



Somatic mutation matrix ➔ NMF ➔ model selection (# of signatures and stability)
➔ Detection of de novo mutational signatures ➔ comparison with known signatures

# Tools for extracting mutational signatures

**Inferring *de novo* signatures**

Alexandrov, MatLab (*Nature* 2013)
EMu (*Genome Biology* 2013)
Maftools (*Genome Res* 2018)
MutationalPatterns (*Genome Med* 2018)
MutSpec (*BMC Bioinformatics* 2016)
SigFit (*BioRxiv* 2020)
SigMiner (*medRxiv* 2020)
SignatureAnalyzer (*Nature Commun* 2015)
SignatureToolsLib (*Nat Cancer* 2020)
SigneR (*Bioinformatics* 2017)
SomaticSignatures (*Bioinformatics* 2015)
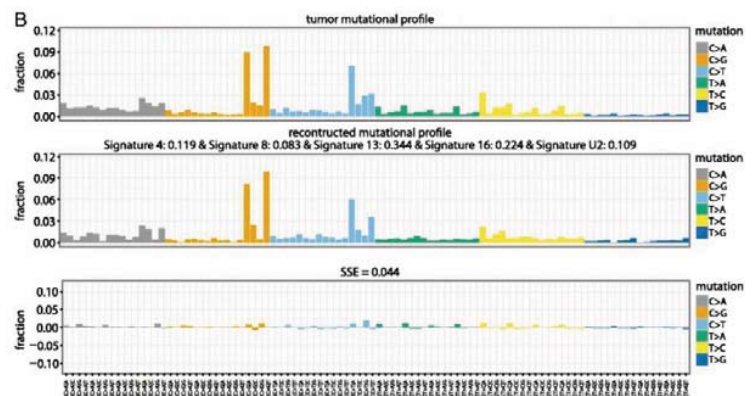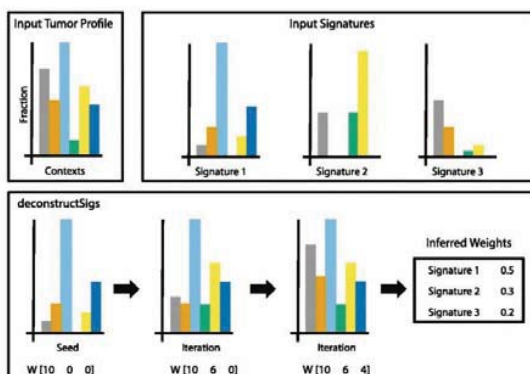SigProfiler (COSMIC)

**Fitting known signatures**

deconsructSigs (*Genome Biology* 2016)
SignatureEstimation (*Bioinformatics* 2018)
YAPSA (R Package v 1.16.0)

**Web interfaces**

MutaGene (*NAR* 2017)
mSignatureDB (*NAR* 2018)
MuSiCa (*BMC Bioinformatics* 2018)
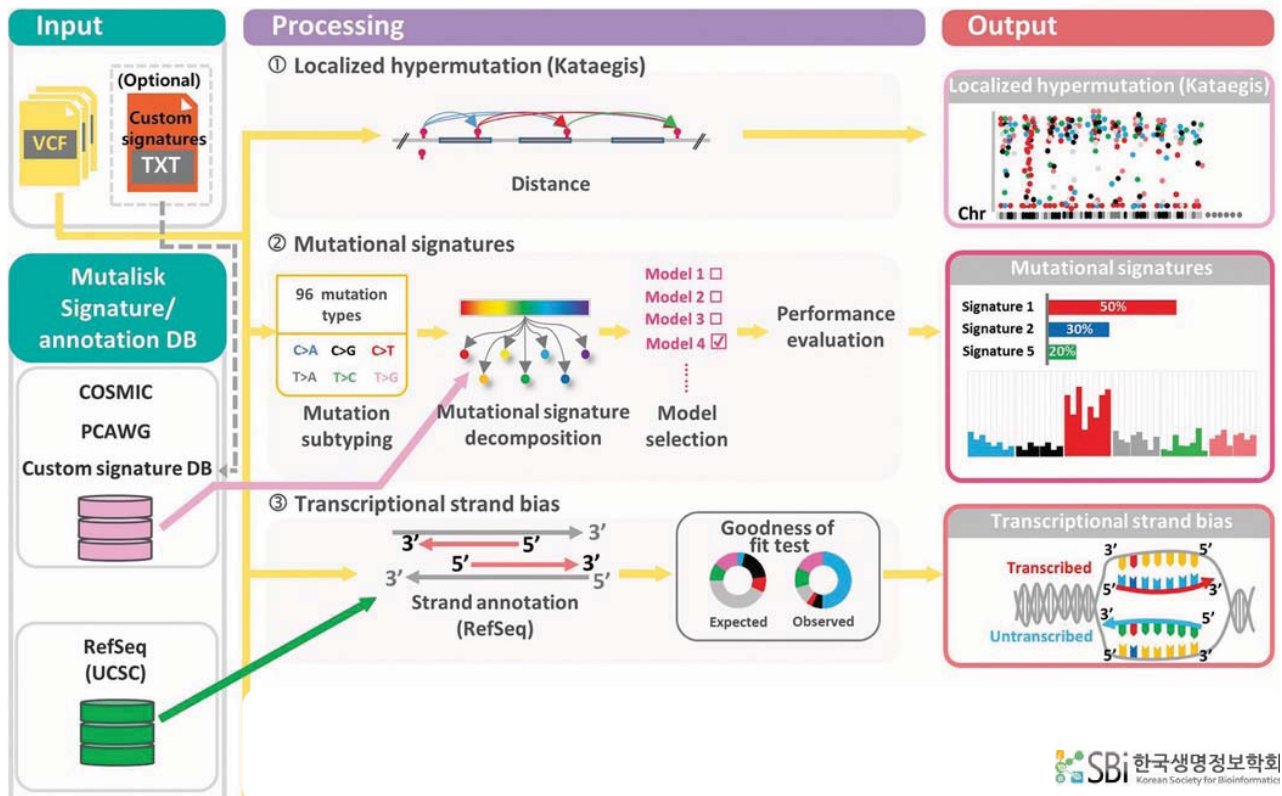Mutalisk (*NAR* 2018)

---

# (2) deconstructSigs

R based package. Mutation matrix as an input (sample, chr, pos, ref, alt)



Rosenthal et al., *Genome Biology* (2016)

# Tools for extracting mutational signatures

**Inferring *de novo* signatures**

Alexandrov, MatLab (***Nature*** 2013)
EMu (***Genome Biology*** 2013)
Maftools (***Genome Res*** 2018)
MutationalPatterns (***Genome Med*** 2018)
MutSpec (***BMC Bioinformatics*** 2016)
SigFit (***BioRxiv*** 2020)
SigMiner (***medRxiv*** 2020)
SignatureAnalyzer (***Nature Commun*** 2015)
SignatureToolsLib (***Nat Cancer*** 2020)
SigneR (***Bioinformatics*** 2017)
SomaticSignatures (***Bioinformatics*** 2015)
SigProfiler (COSMIC)

**Fitting known signatures**

deconsructSigs (***Genome Biology*** 2016)
SignatureEstimation (***Bioinformatics*** 2018)
YAPSA (R Package v 1.16.0)

**Web interfaces**

MutaGene (***NAR*** 2017)
mSignatureDB (***NAR*** 2018)
MuSiCa (***BMC Bioinformatics*** 2018)
Mutalisk (***NAR*** 2018)

---

# (3) Web interfaces: Mutalisk
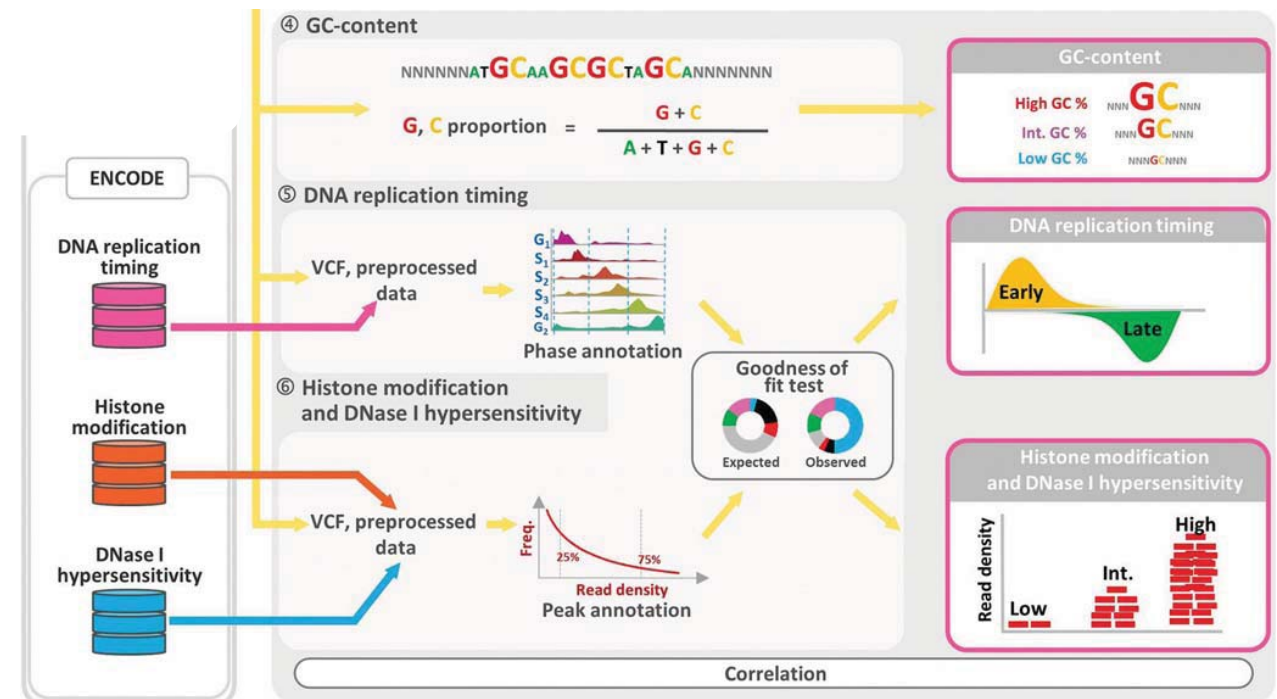
http://mutalisk.org



Mutalisk: a web-based somatic MUTation
AnaLyIS toolKit for genomic, transcription
and epigenomic signatures

Lee JK et al., NAR 2018

# (3) Workflow in Mutalisk



# (3) Workflow in Mutalisk (2)

# (3) Input for Mutalisk



# (3) Output in Mutalisk (1)

# (3) Output in Mutalisk (2)



# Genome QC with mutational signatures

# Amplification/sequencing artifacts make unique signatures



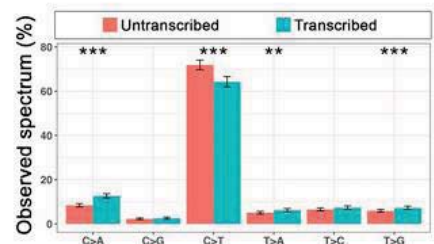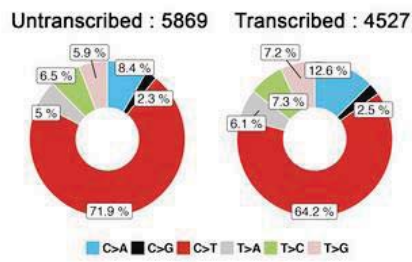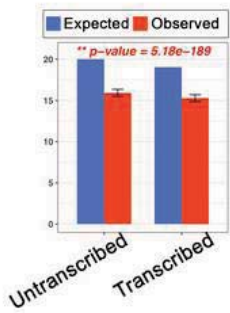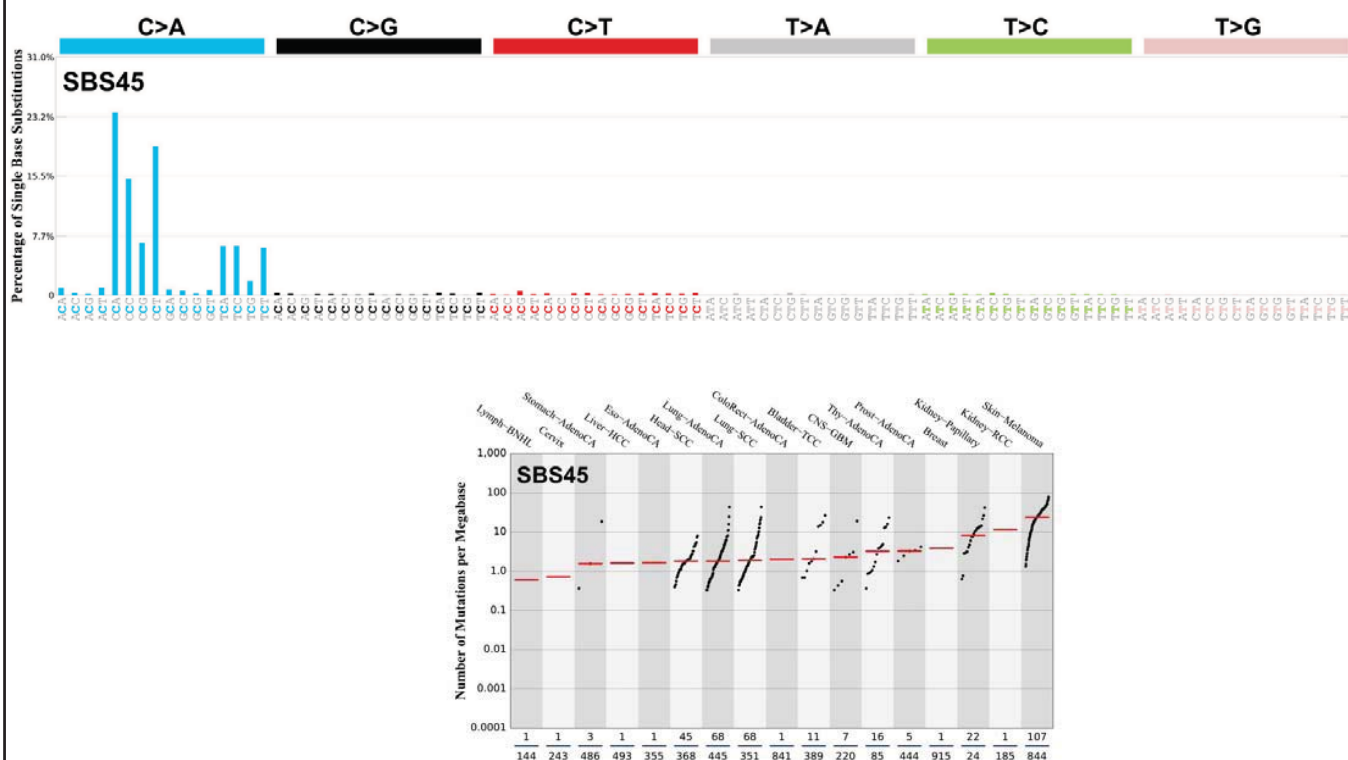# SBS45, a signature of 8-oxoG artifact

# First report for 8-oxoG artificial signature



0.25 ≤ AF < 0.50
n=71940

0.10 ≤ AF < 0.25
n=63773

0 ≤ AF < 0.10
n=91790

Costello et al., *NAR* (2012)

# A typical pipeline for cancer genome analyses

# Workflow in FIREVAT, a software for filtering artifacts



Software | Open Access | Published: 17 December 2019

FIREVAT: finding reliable variants without artifacts in human cancer samples using etiologically relevant mutational signatures
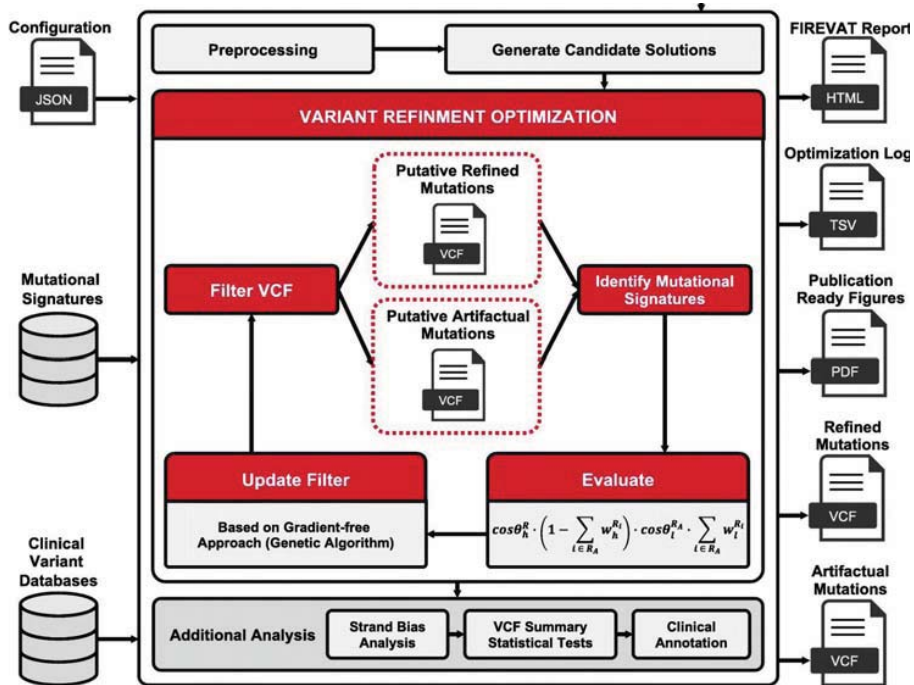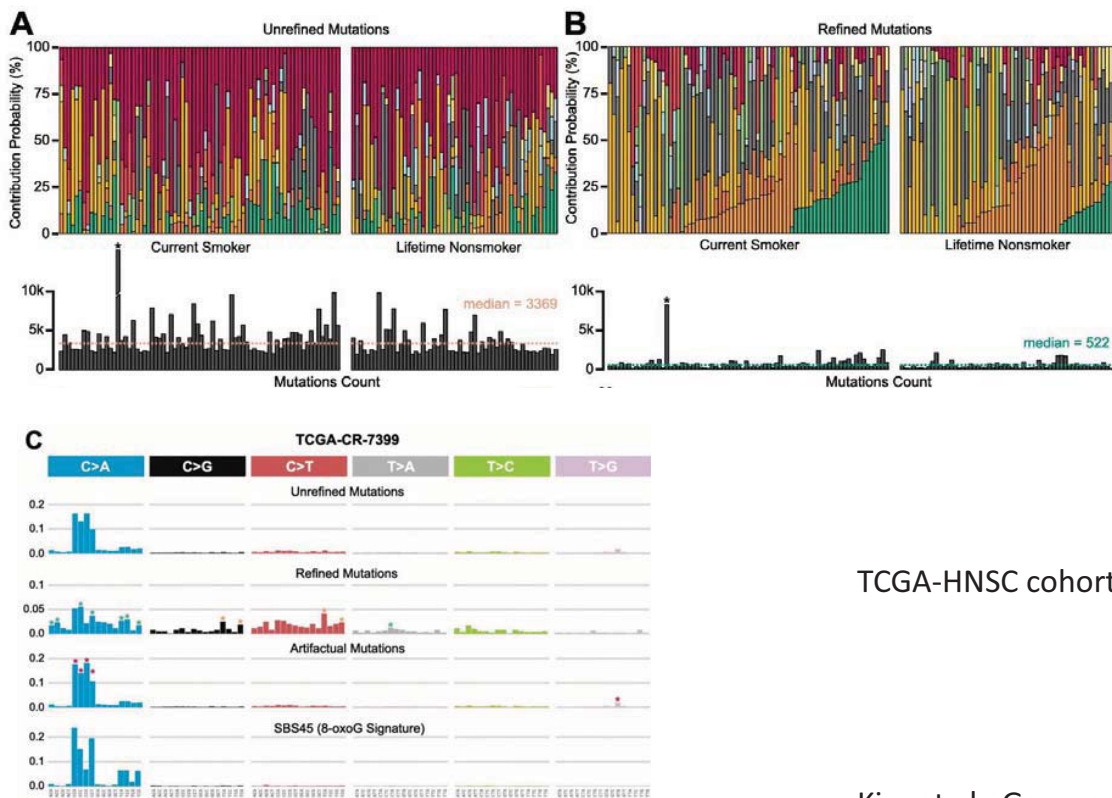
Kim et al., Genome Medicine 2019

# Filtering mutations using FIREVAT



TCGA-HNSC cohort

Kim et al., Genome Med 2019

## 전망

- 돌연변이 signature 개수는 총 몇 개가 될까?

- 돌연변이 signature 각각의 원인을 규명할 수 있을까?

- Structural variation의 signature는 무엇이 있을까?

SBi 한국생명정보학회
Korean Society for Bioinformatics

## Summary

- 돌연변이는 random 하게 생기지 않는다

- Mutational signature 개념을 이용하여 정확한 variant calling 을 할 수 있다

- Mutational signature 개념을 이용하여 돌연변이가 만들어진 원인을 추적할 수 있다

- Mutational signature를 구하는 tool을 이해하고 사용할 수 있다.

SBi 한국생명정보학회
Korean Society for Bioinformatics