# KSBi-BIML 2021

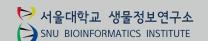
Bioinformatics & Machine Learning (BIML)
Workshop for Life Scientists

생물정보학 & 머쉰러닝 워크샵(온라인)

Mutational Signatures in Cancer Genomes

주영석







# Bioinformatics & Machine Learning for Life Scientists BIML-2021

안녕하십니까?

한국생명정보학회의 동계 워크샵인 BIML-2021을 2월 15부터 2월 19일까지 개최합니다. 생명정보학 분야의 융합이론 보급과 실무역량 강화를 위해 도입한 전문 교육 프로그램인 BIML 워크샵은 2015년에 시작하였으며 올해로 7차를 맞이하게 되었습니다. 유례가 없는 코로나 대유행으로 인해 올해의 BIML 워크숍은 온라인으로 준비했습니다. 생생한 현장 강의에서만 느낄 수 있는 강의자와 수강생 사이의 상호교감을 가질수 없다는 단점이 있지만, 온라인 강의의 여러 장점을 살려서 최근 생명정보학에서 주목받고 있는 거의 모든 분야를 망라한 강의를 준비했습니다. 또한 온라인 강의의한계를 극복하기 위해서 실시간 Q&A 세션 또한 마련했습니다.

BIML 워크샵은 전통적으로 크게 생명정보학과 AI, 두 개의 분야로 구성되어오고 있으며 올해 역시 유사한 방식을 채택했습니다. AI 분야는 Probabilistic Modeling, Dimensionality Reduction, SVM 등과 같은 전통적인 Machine Learning부터 Deep Learning을 이용한 신약개발 및 유전체 연구까지 다양한 내용을 다루고 있습니다. 생명정보학 분야로는, Proteomics, Chemoinformatics, Single Cell Genomics, Cancer Genomics, Network Biology, 3D Epigenomics, RNA Biology, Microbiome 등 거의 모든 분야가 포함되어 있습니다. 연사들은 각 분야 최고의 전문가들이라 자부합니다.

이번 BIML-2021을 준비하기까지 너무나 많은 수고를 해주신 BIML-2021 운영위원회의 김태민 교수님, 류성호 교수님, 남진우 교수님, 백대현 교수님께 커다란 감사를 드립니다. 또한 재정적 도움을 주신, 김선 교수님 (Al-based Drug Discovery), 류성호 교수님, 남진우 교수님께 감사를 표시하고 싶습니다. 마지막으로 부족한 시간에도 불구하고 강의 부탁을 흔쾌히 허락하시고 훌륭한 강의자료를 만드는데 노력하셨을 뿐만아니라 실시간 온라인 Q&A 세션까지 참여해 수고해 주시는 모든 연사분들께 깊이감사드립니다.

2021년 2월

한국생명정보학회장 김동섭

#### 강의개요

### 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.

\* 강의: 주영석 교수 (KAIST 의과학대학원)

#### **Curriculum Vitae**

#### Speaker Name: Young Seok Ju, MD, PhD.



#### ▶ Personal Info

Name Young Seok Ju

Title Associate Professor

Affiliation KAIST

#### **▶** Contact Information

Address 291 Daehak-Ro, Yuseong-Gu, Daejeon, 34141

Email ysju@kaist.ac.kr

Phone Number 042-350-4237

Research interest: Somatic muttions in human cells

#### **Educational Experience**

2007 M.D. in Medicine, Seoul National University College of Medicine, Korea

2010 Ph.D. in Genomic Medicine, Seoul National University College of Medicine, Korea

#### **Professional Experience**

2013-2015 Postdoctoral Fellow, Wellcome Sanger Institute, Cambridge UK

2015- Assistant/Associate Professor, KAIST, Daejeon, Korea

#### **Selected Publications (5 maximum)**

- 1. 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 epub ahead of print.
- 2. Yuan Y\*, **Ju YS**\*, Kim Y\*, Li J, Wang J, ..., & Liang H. Comprehensive molecular characterization of mitochondrial genomes in human cancers. *Nat Genet*. 2020 Mar;52(3):342-352. PMID:32024997.
- 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 Mar;145(3):1017-1021. PMID:31805313
- 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 Jun 13;177(7):1842-1857. PMID:31155235.
- 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



본 강의 자료는 한국생명정보학회가 주관하는 KSBi-BIML 2021 워크샵 온라인 수업을 목적으로 제작된것으로 해당 목적 이외의 다른 용도로 사용할 수 없음을 분명하게 알립니다. 수업 목적으 로 배포 및 전송 받은 경우에도 이를 다른 사람과 공유하거나 복 제, 배포, 전송할 수 없습니다.

만약 이러한 사항을 위반할 경우 발생하는 모든 법적 책임은 전적으로 불법 행위자 본인에게 있음을 경고합니다.

### 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

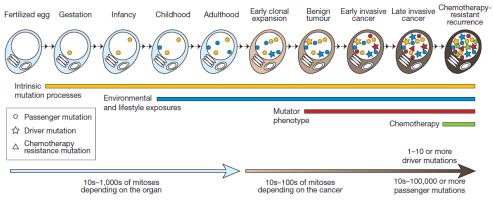
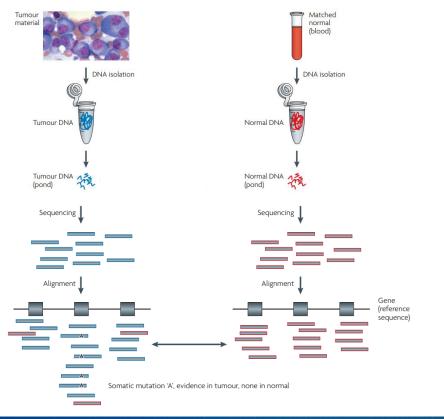


Figure 1 | The lineage of mitotic cell divisions from the fertilized egg to a single cell within a cancer showing the timing of the somatic mutations acquired by the cancer cell and the processes that contribute to them. Mutations may be acquired while the cell lineage is phenotypically normal, reflecting 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.

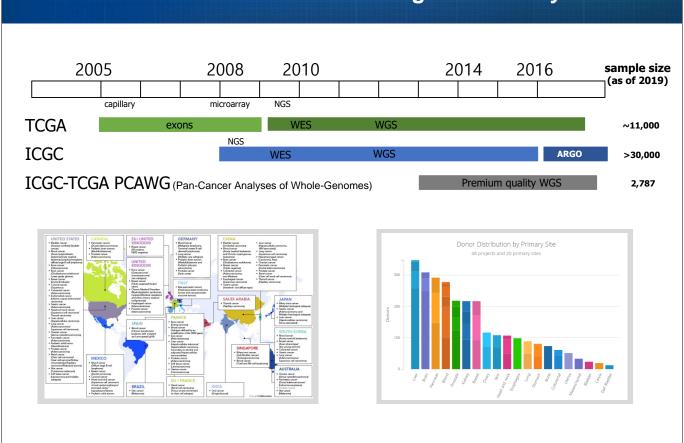
• 대부분의 산발성 암 (sporadic cancer) 의 원인은 체세포 돌연변이이다

# 돌연변이의 검출을 위한 전략

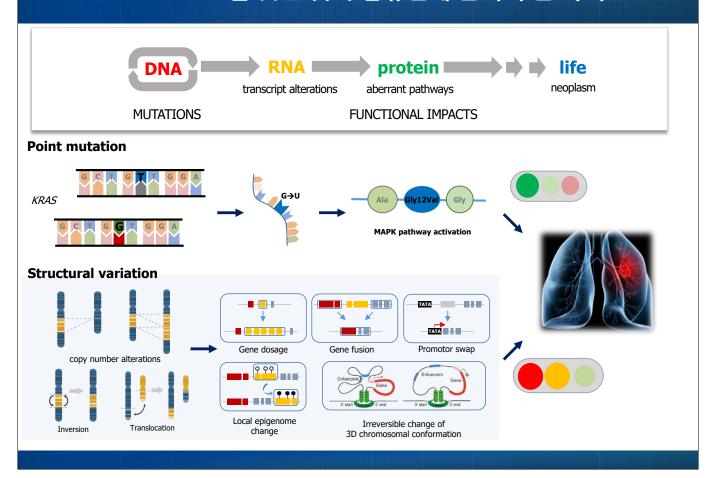


Meyerson M et al., Nat Rev Cancer (2010)

### International consortia for cancer genome analyses

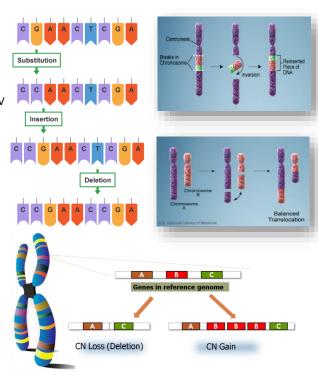


### Driver mutation을 찾는 것이 암유전체 분석의 한 목적

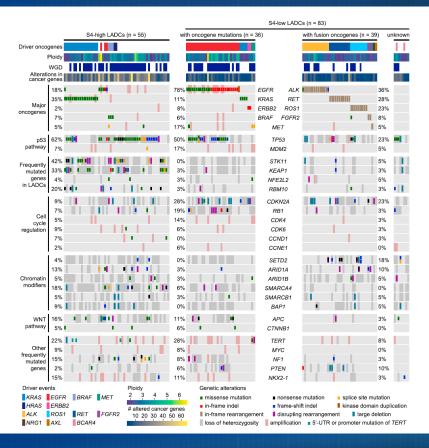


### 어떤 돌연변이가 있는가?

- 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)

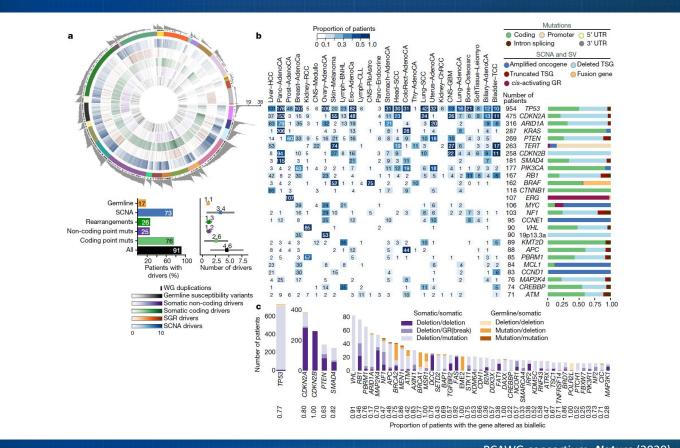


### Cancer genome에서 driver mutation의 분포

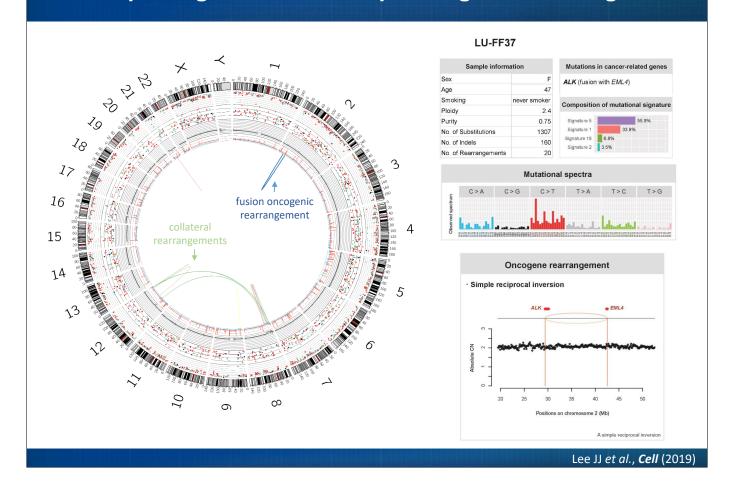


Lee JJ et al., **Cell** (2019)

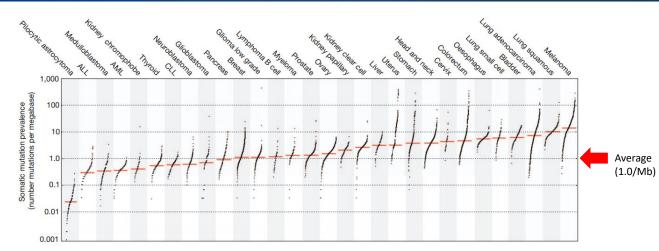
### **Driver mutations in pan-cancer genomes**



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



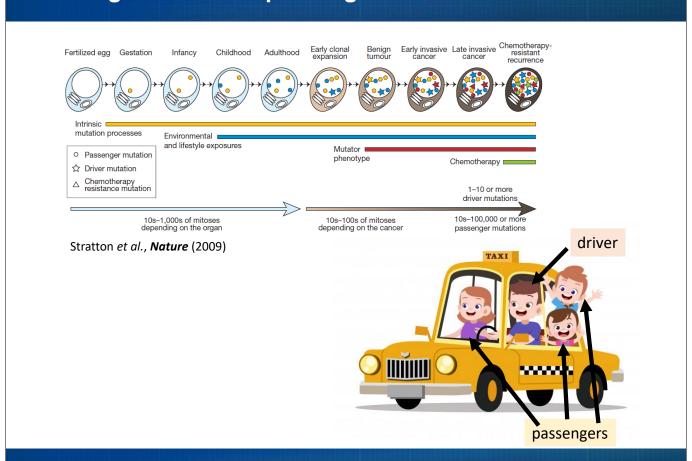
### 암유전체의 돌연변이들

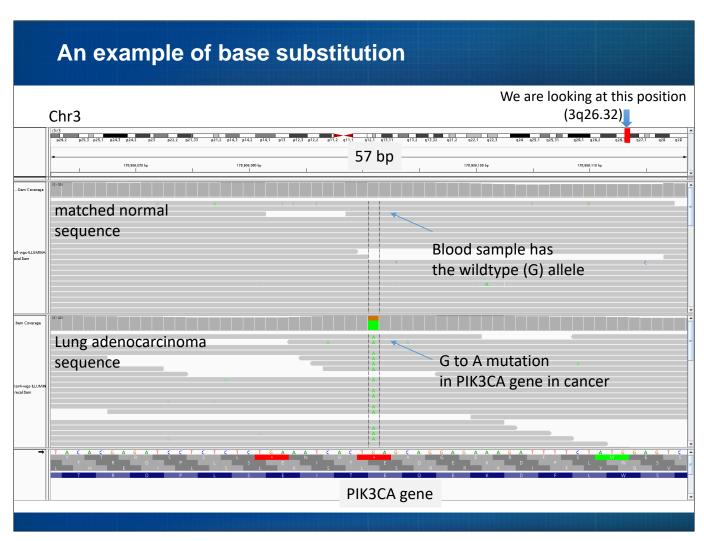


Alexandrov L et al., Nature (2013)

- 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?





### Mutational signature 개념을 접하다

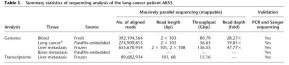


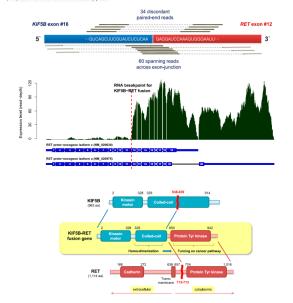
Genome Research (2012.3)

#### A transforming KIF5B and RET gene fusion in lung adenocarcinoma revealed from whole-genome and transcriptome sequencing

Young Seok Ju, <sup>1,2</sup> Won-Chul Lee, <sup>1,3</sup> Jong-Yeon Shin, <sup>1,4</sup> Seungbok Lee, <sup>1,3</sup> Thomas Bleazard, <sup>1</sup> Jae-Kyung Won, <sup>2</sup> Young Tae Kim, <sup>6,7</sup> Jong-Il Kim, <sup>1,3,4,8</sup> Jin-Hyoung Kang, <sup>2</sup> and Jeong-Sun Seol. <sup>1,2,3,4,8,10</sup> (Fammik Medicin Institute (KM), Medical Research Center, Seoul National University, Cental State Office of Seoul 110-799, Korea; <sup>4</sup> Mooragen Inc., Seoul 153-78, Korea; <sup>5</sup> Department of Bornekiod Science, Seoul National University Creations Seoul 110-799, Korea; <sup>6</sup> Mooragen Inc., Seoul 153-78, Korea; <sup>6</sup> Vender Seoul 13-78, Korea; <sup>6</sup> Vender Seoul 13-79, Korea; <sup>6</sup> Vender Seoul 110-799, Korea; <sup>6</sup> Vender Seoul 110-

t. Article, supplemental material, and pub-nome.org/cgi/doi/10.1101/gr.133645.111.





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

### Mutational signature 개념을 접하다 (2)

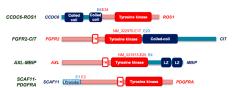
#### The transcriptional landscape and mutational profile of lung adenocarcinoma

JOILS-II NITH), JIL-HYOUNG RAING, "and TOUNG LAB KIM"

"Genomic Medicine Institute (CMI), Medical Research Center, Soul National University, Soul 110-799, Korea, "Department of Bochemistry, Seoul National University College of Medicine, Seoul 110-799, Korea, "Department of Bomedical Sciences, Seoul National University College of Medicine, Seoul 110-799, Korea, "Department of Bomedical Sciences, Seoul National University College of Medicine, Seoul National University College of Medicine, Seoul National University College of Medicine, Seoul 110-799, Korea, "Polysiston of Medical Coology, Research Institute of Medical Science, The Catholic University of Korea, Seoul 137-040, Korea, "Department of Theoric and Cardiovascular Surgery, Seoul National University Hospital, Seoul 110-799, Korea, "Department of Theoric and Cardiovascular Surgery, Seoul National University Hospital, Seoul 110-799, Korea, "Department of Theoric and Cardiovascular Surgery, Seoul National University Hospital, Seoul 110-799, Korea, "Department of Theoric and Cardiovascular Surgery, Seoul 137-040, Korea

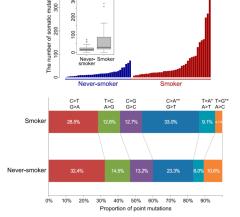
Genome Research 2109







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



Seo JS et al., Genome Res (2012b)

### Mutational signature 개념을 접하다 (3)







Elizabeth P Murchison



#### **ARTICLE**

oi:10.1038/eature12477

# Signatures of mutational processes in human cancer

A list of authors and their affiliations appears at the end of the paper

All cancers are caused by somatic mutations; however, understanding of the biological processes generaling these mutations is limited. The etailogue of somatic mutation from a cancer genome bears the signatures of the mutational processes that have been operative. Here we can layed 4,938,362 mutations from 7,042 cancers and extracted more than 20 distinct mutational signatures. Some are present in many cancer types, notably a signature artitrated to the APOBEC family of yiddine deaminases, whereas others are confined to a single cancer class. Certain signatures are associated with age of the patient at cancer diagnosis, known mutagenic exposures or defects in DAM anintennae, but many ared cryptic origin. In addition to these genome-wide mutational signatures, hypermutation localized to small genomic regions, ktategie, is found in many cancer types. The results reveal the diversity of mutational processes underlying

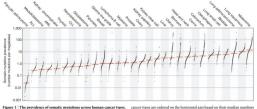
Somatic mutations found in cancer genomes' may be the consequence of the intrinsic slight infidelity of the DNA replication muchinery ecogenous or endogenous mutagen exposures, enzymatic modification of DNA, or delective DNA replication. In some cancer types, a sub-stantial proportion of somatic mutations are known to be generated by exposures, for example, beloes modified in lung cancers and tenance, for example, defective DNA mismatch repair in some colorectal cancers. However, our meterstanding of the mutations processes that cause somatic mutations in most cancer classes i remarkably limited.

remarkably limited.

Different mutational processes often generate different combinations of mutation types, termed 'signatures'. Until recently, mutational sig-

of frequently mutated calcher gibes, notatoy 12:53 (et. 4), Although of mortal control frequently and the state through the control frequently and the state through the control frequently and the control freque

Recent advances in sequencing technology have overcome past limitations of scale<sup>1</sup>. Thousands of somatic mutations can now be identified in a single cancers sample, offering the possibility of deciphrish mutational signatures even when several mutational processes are

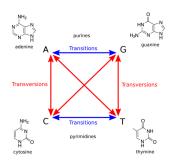


Every dot represents a sample whereas the red horizontal lines are the median numbers of mutations in the respective cancer types. The vertical axis (log scaled) shows the number of mutations per megabase whereas the different nomatic mutations. We thank G. Getz and colleagues for the design of this are 3s. ALL, acute hymphoblastic leukaemia; AML, acute myeloid leukaemia; L, chronic lymphocytic leukaemia.

22 AUGUST 2013 | VOL 500

Alexandrov L et al., Nature (2013)

### mutational origin: Mutation은 랜덤하게 생기는 것이 아니다



#### 6 classes of base substitutions

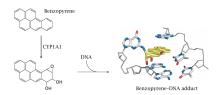
C>A (G>T), C>G (G>C), C>T (G>A)
T>A (A>T), T>C (A>G), T>G (A>C)



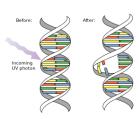
#### 돌연변이는 "랜덤" 이 아니라 DNA damage x DNA repair 과정

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

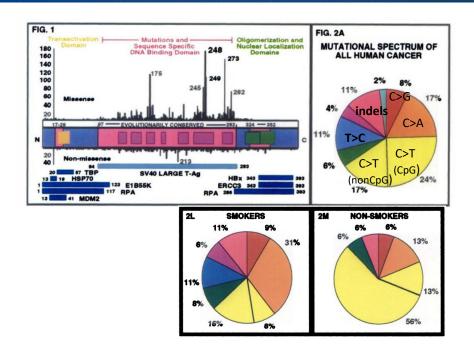
**Tobacco smoking** C>A substitutions



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



### Classical observation



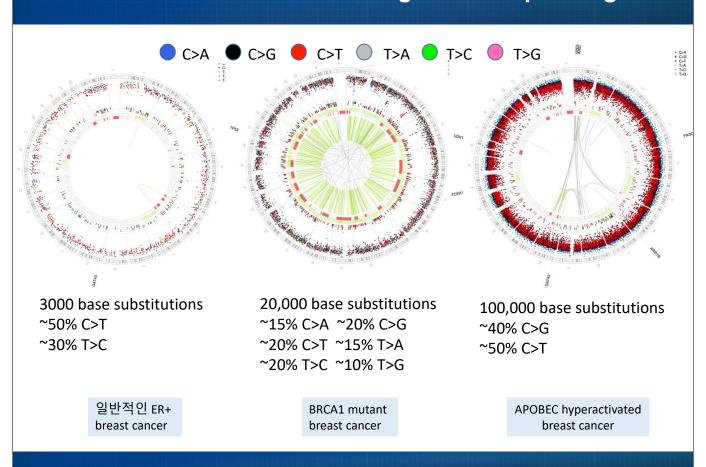
- 암종마다, 그리고 발암물질의 노출에 따라서 TP53 유전자에 생기는 돌연변이 패턴이 상이하다

Greenblatt et al., Cancer Research (1994)

### Mutational signature의 예시

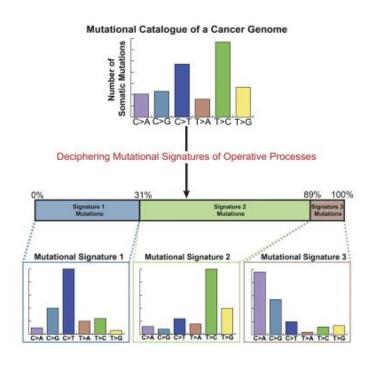
- 폐암의 전장 유전체 분석에서 20,000개의 base substitution 발견
  - 이 가운데 80%가 C>A mutations. 주된 돌연변이 발생기전은?
  - (흡연에 노출)
- 흑색종의 전장 유전체 분석에서 20,000 개의 base substitution 발견
  - 이 가운데 90%가 C>T mutations 이고 수백개의 CC>TT 도 같이 발견 주된 돌연변이 발생기전은?
  - (UV에 노출)
- 실제로는 하나의 암 유전체에서 발생하는 돌연변이들이 위와 같은 단일 돌연변이 발생 기전이 아니라, 여러 돌연변이 기전의 '조합' 으로 만들어지는 일이 훨씬 흔하다

# 실제 3개의 breast cancer whole-genome sequencing 결과



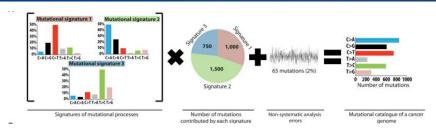
### Tumor의 돌연변이 스펙트럼은 이론적으로 n개의 서로 다른 Mutational process로 설명된다

하지만 우리는 n이 얼마인지도, 각각의 process의 spectrum도 알고있지 못한다



Alexandrov L et al., Cell reports (2013)

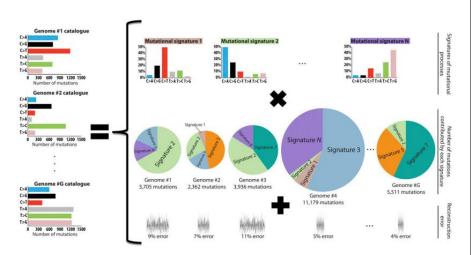
# Understanding mutational processes from mutational spectrum: a blind-source separation problem



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)

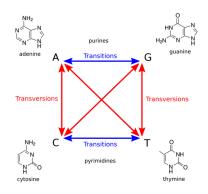


Alexandrov L et al., Cell reports (2013)

### Single base substitutions (SBS) into 96 subclasses

### • C>A (G>T)

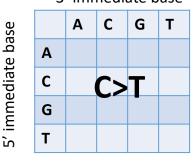
- C>G (G>C)
- C>T (G>A)
- T>A (A>T)
- T>C (A>G)
- T>G (A>C)



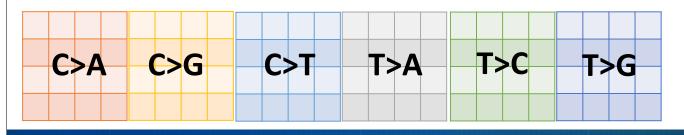
sequence context

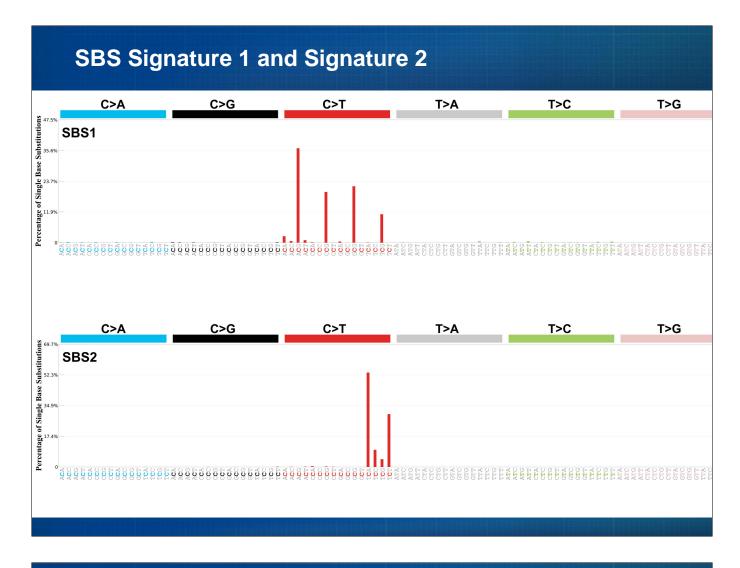
5'B - Wt > Var- 3'B

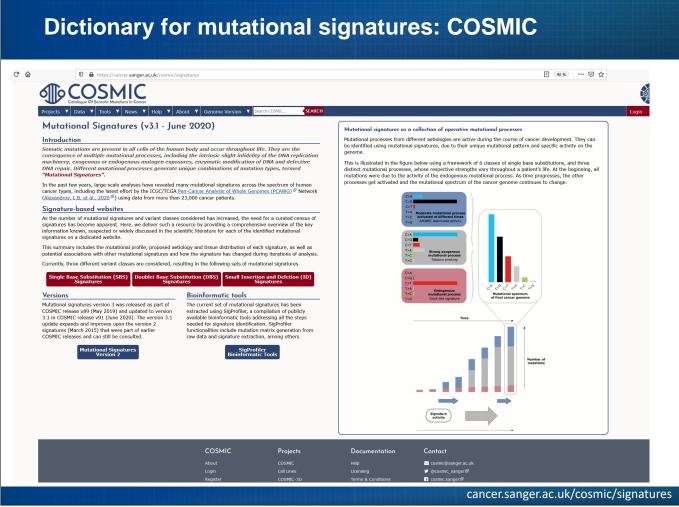
3' immediate base



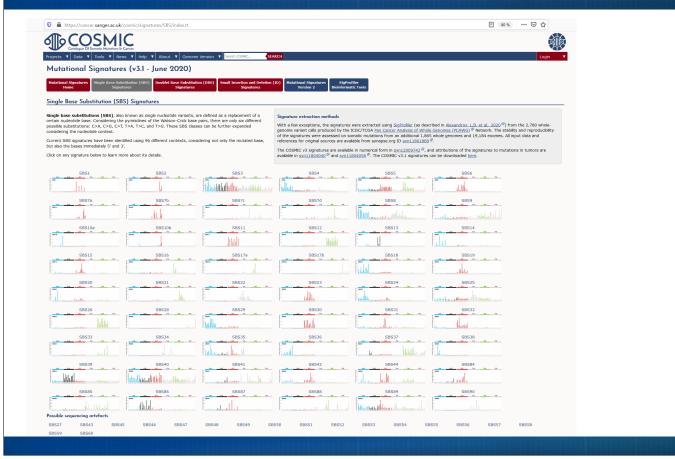
4 x 6 types x 4 = 96 subtypes

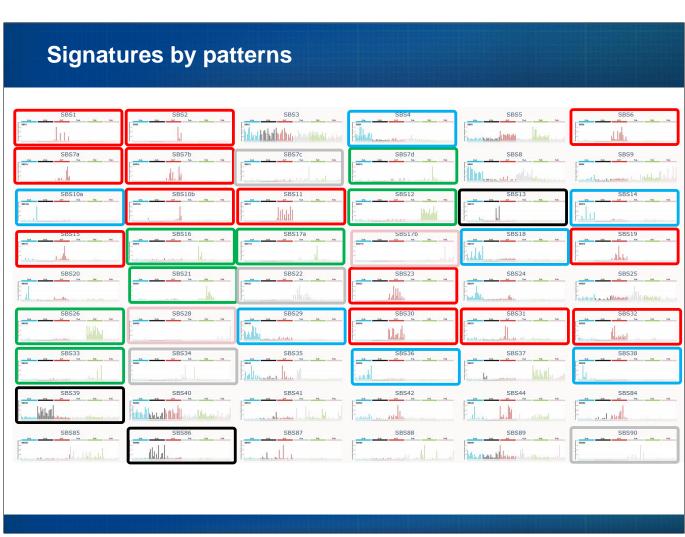


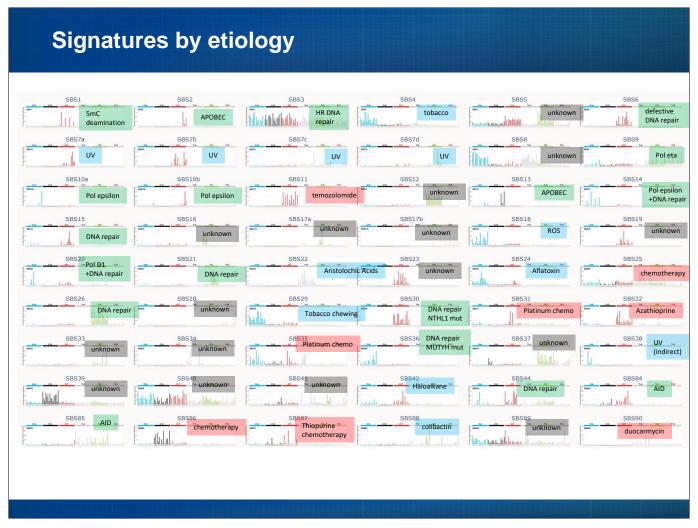


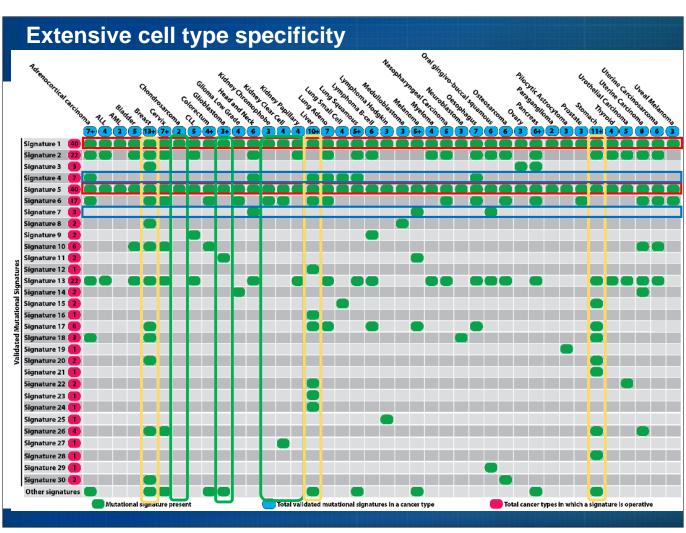


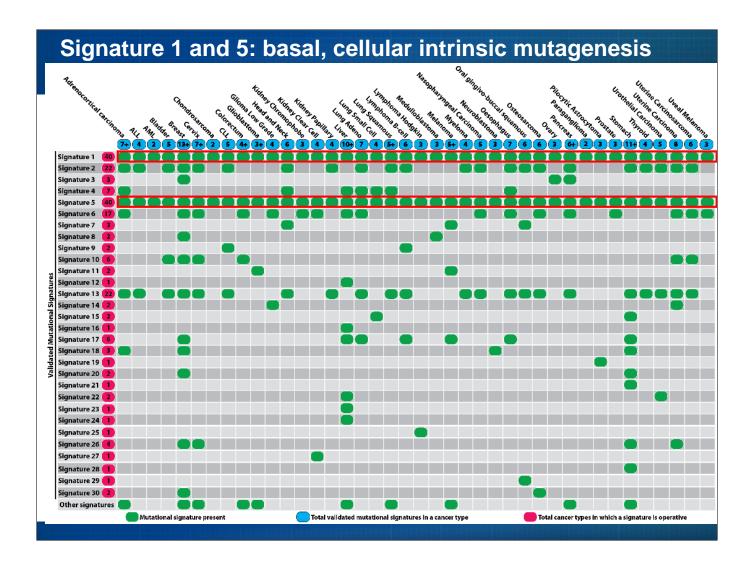
### 49 +5 biologic signatures in SBS mutations (v3)



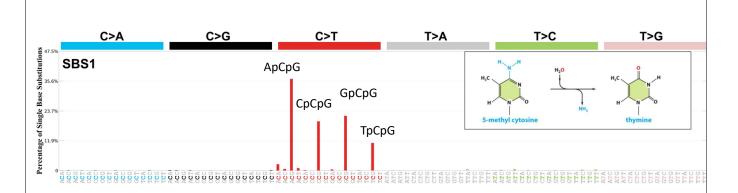












#### **Cancer types:**

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.

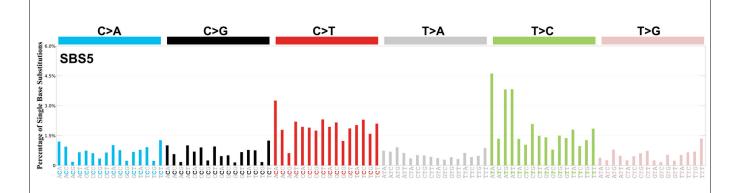
#### **Additional mutational features:**

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.

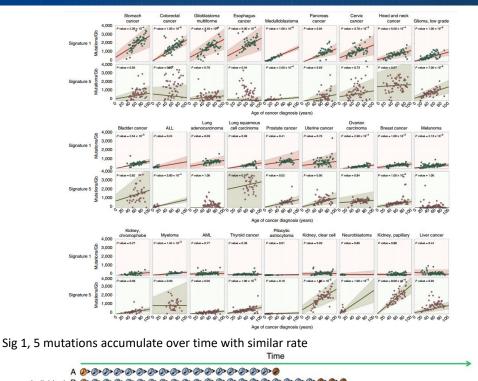
#### Proposed etiology:

The aetiology of Signature 5 is unknown.

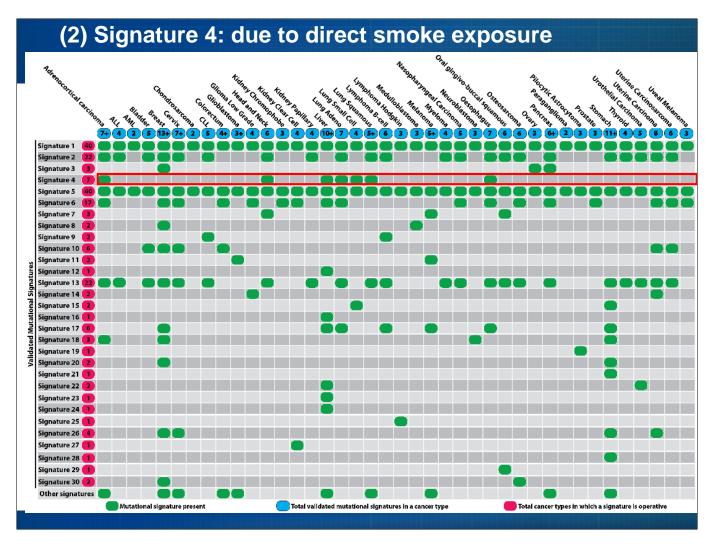
#### **Additional mutational features:**

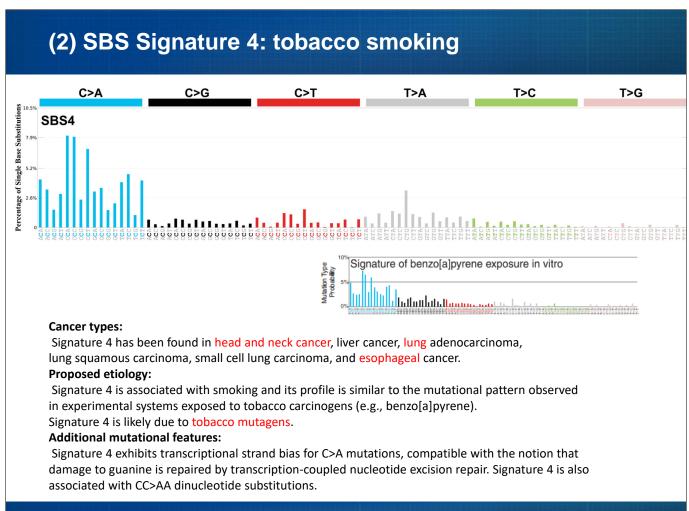
Signature 5 exhibits transcriptional strand bias for T>C substitutions at ApTpN context.

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

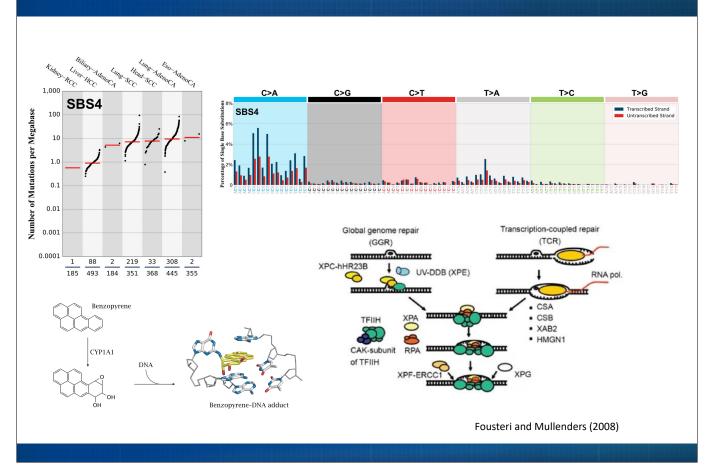


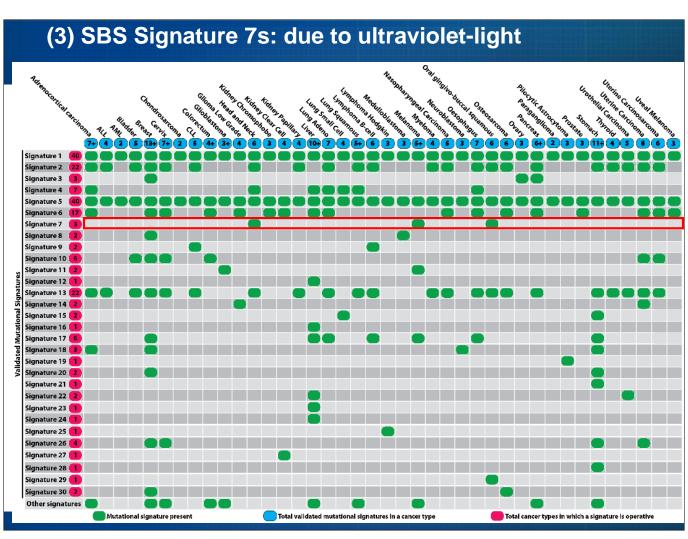
Alexandrov L et al., Nature Genet (2015)

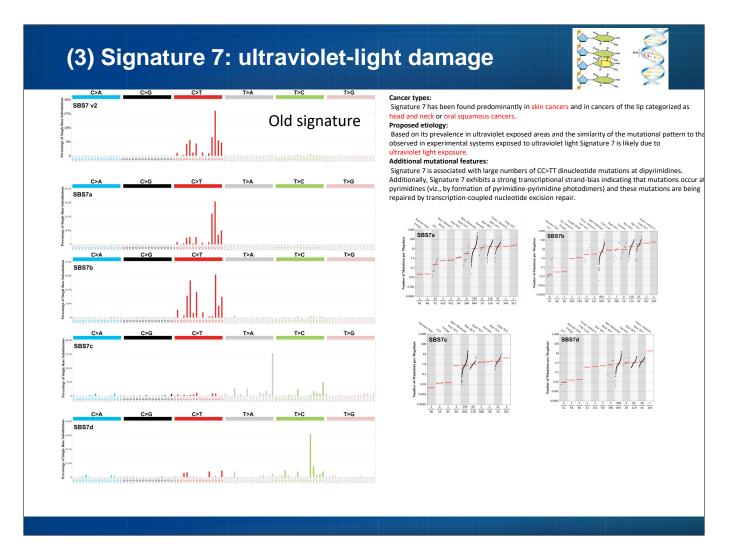


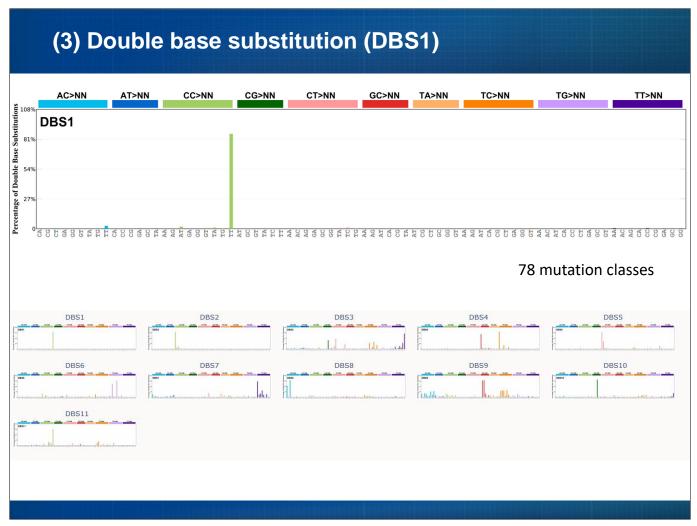


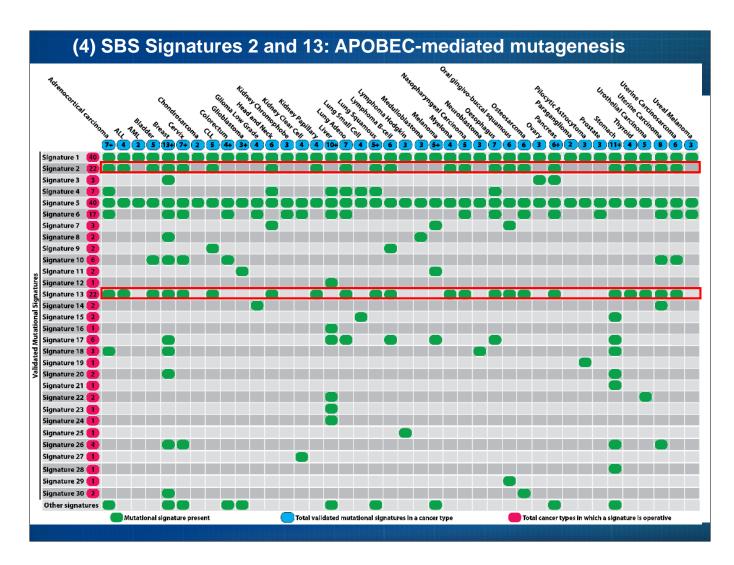
### (2) SBS Signature 4: mutational burden and strand bias

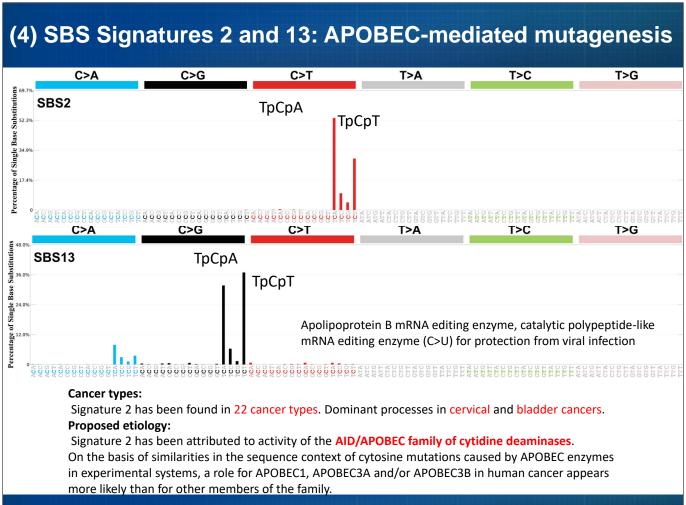




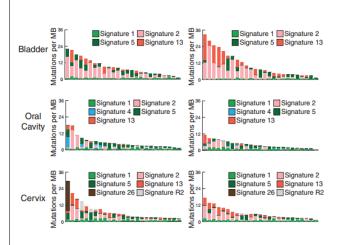








### **APOBEC-mediated mutations**



EGFR T30M (detected in clinical setting)

Misper clone

PLICAC ESSEK + F726K

ROS F1035

MISCA

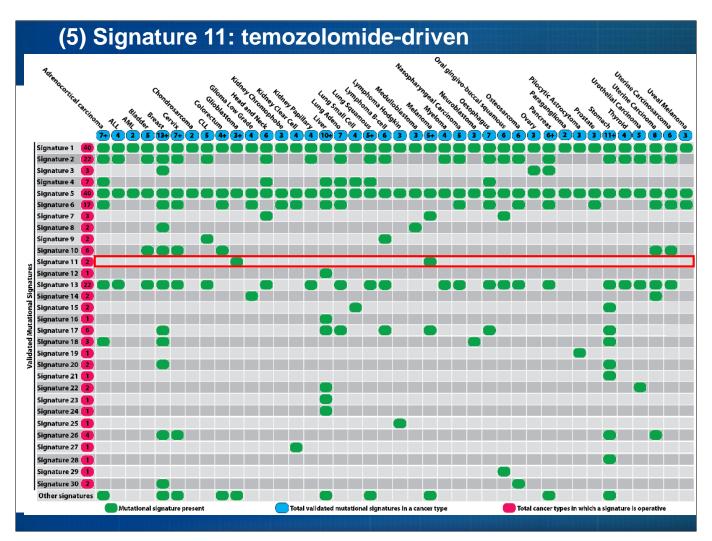
MISC

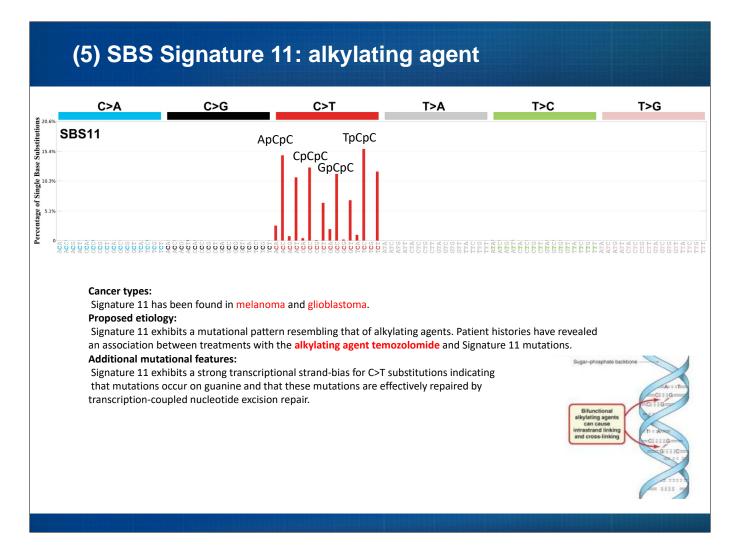
Alexandrov L et al., Science (2016)

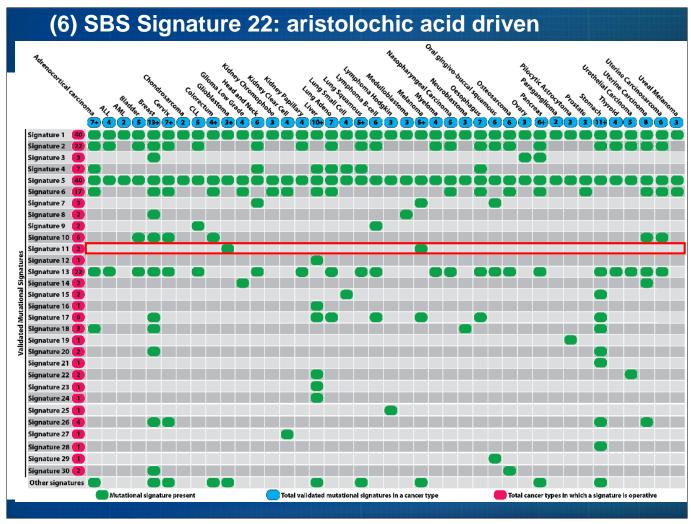
Lee J et al., J Clin Oncol (2017)

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

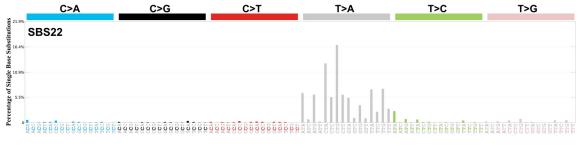
Activated in the late branch in lung cancers. (Episodically activating?)

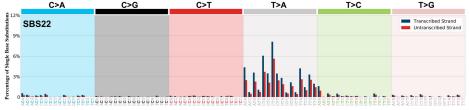


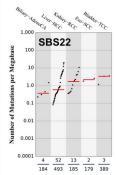




### (6) SBS Signature 22: aristolochic acids







#### 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:

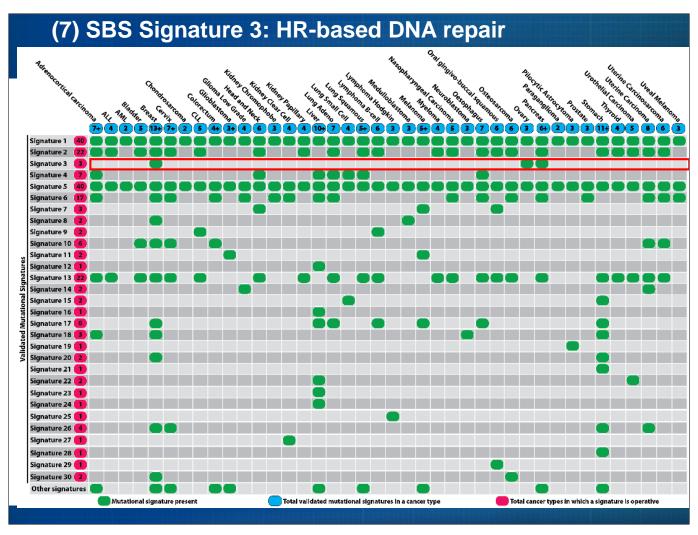
Signature 22 has a very high mutational burden in urothelial carcinoma; however, its mutational burden is much lower in liver cancers.

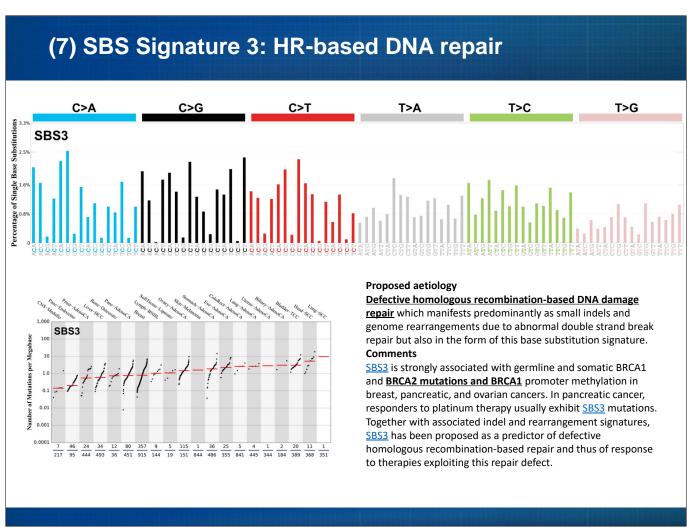
### (6) SBS signature 22: aristolochic acids



Aristolochia clematitis (쥐방울덩굴, 동북마두령, 관목통)

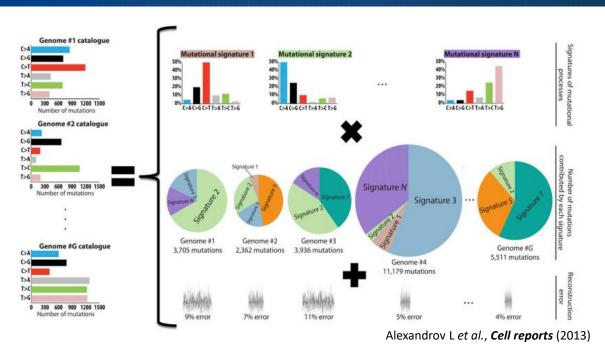
https://www.accessdata.fda.gov/cms ia/importalert 141.html





# 어떻게 mutational signature를 구할 것인가?

# 몇 개의 sample, 몇 개의 돌연변이가 필요할까?



- (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

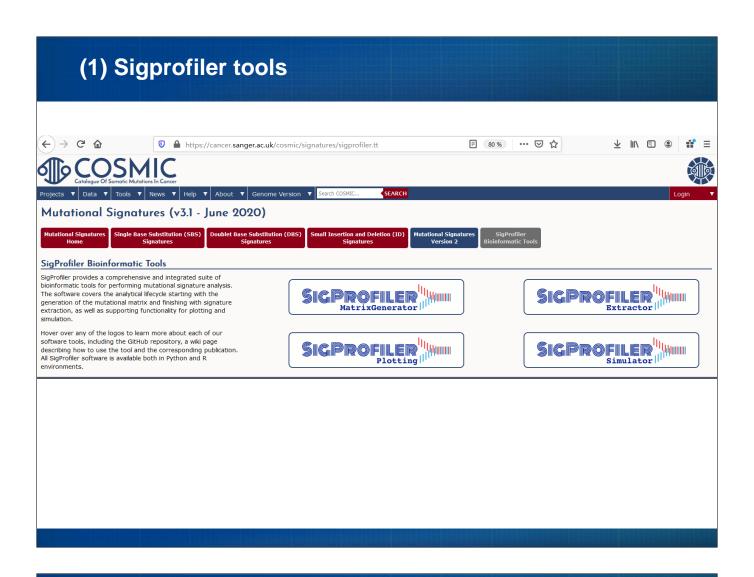
deconsructSigs (Genome Biology 2016)

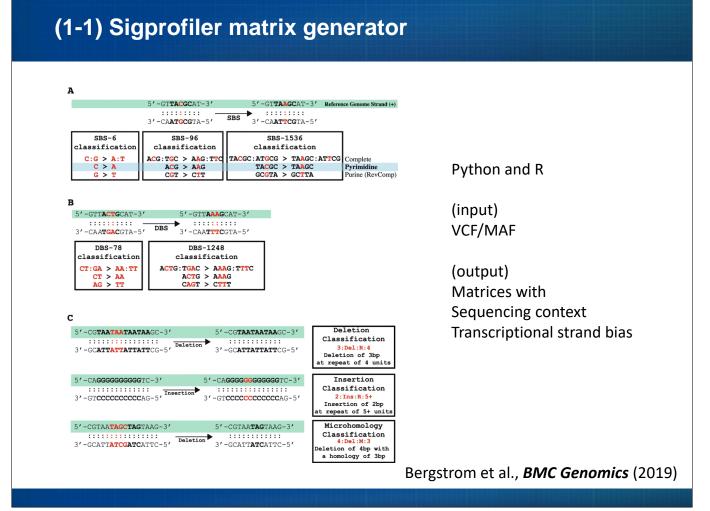
SignatureEstimation (*Bioinformatics* 2018) YAPSA (R Package v 1.16.0)

#### Web interfaces

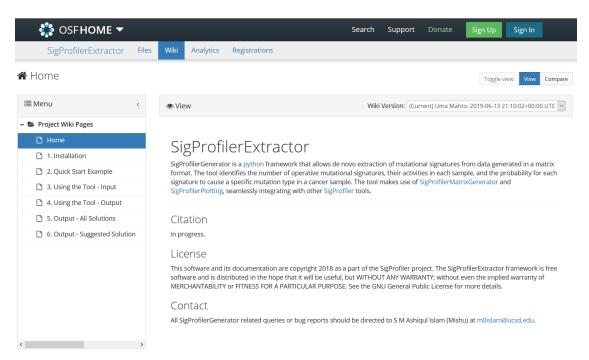
MutaGene (*NAR* 2017) mSignatureDB (*NAR* 2018) MuSiCa (*BMC Bioinformatics* 2018) Mutalisk (*NAR* 2018)

#### (1) Sigprofilers **业 Ⅲ ۩ ◎ #** ≡ https://cancer.sanger.ac.uk/cosmic/signatures E 80% · · · ☑ ☆ എ் COSMIC Projects ▼ Data ▼ Tools ▼ News ▼ Help ▼ About ▼ Genome Version ▼ Search CO Mutational Signatures (v3.1 - June 2020) Mutational signatures as a collection of operative mutational processes Mutational processes from different aetiologies are active during the course of cancer development. They can be identified using mutational signatures, due to their unique mutational pattern and specific activity on the genome. Introduction Somatic mutations are present in all cells of the human body and occur throughout life. They are the consequence of multiple mutational processes, including the intrinsic slight infidelity of the DNA replication machinery, exogenous or endogenous mutagen exposures, enzymatic modification of DNA and defective DNA repair. Different mutational processes generate unique combinations of mutation types, termed "Mutational Signatures". This is illustrated in the figure below using a framework of 6 classes of single base substitutions, and three distinct mutational processes, whose respective strengths vary throughout a patient's life. At the beginning, all mutations were due to the activity of the endogenous mutational process. As time progresses, the other processes get activated and the mutational spectrum of the cancer genome In the past few years, large-scale analyses have revealed many mutational signatures across the spectrum of human cancer types, including the latest effort by the ICGC/TCGA <u>Pan-Cancer Analysis</u> of <u>Whole Genomes (PCAWG)</u> <sup>©</sup> Network (<u>Alexandrov, L.B. et al., 2020</u> <sup>©</sup>) using data from more than 23,000 cancer patients. continues to change Signature-based websites As the number of mutational signatures and variant classes considered has increased, the need for a curated census of signatures an availant dasses ovinated and as assessed, the need to a curated census of signatures has become apparent. Here, we deliver such a resource by providing a comprehensive overview of the key information known, suspected or widely discussed in the scientific literature for each of the identified mutational signatures on a dedicated website. This summary includes the mutational profile, proposed aetiology and tissue distribution of each signature, as well as potential associations with other mutational signatures and how the signature has changed during iterations of analysis. Currently, three different variant classes are considered, resulting in the following sets of mutational signatures. Versions **Bioinformatic tools** Mutational signatures version 3 was released as part of COSMIC release v89 (May 2019) and The current set of mutational signatures has been extracted using SigProfiler, a compilation updated to version 3.1 in COSMIC release v91 (June 2020). The version 3.1 update expands and improves upon the version 2 signatures (March 2015) that were part of earlier COSMIC releases and can still be consulted. of publicly available bioinformatic tools addressing all the steps needed for signature identification. SigProfiler functionalities include mutation matrix generation from raw data and signature extraction, among others.





### (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)

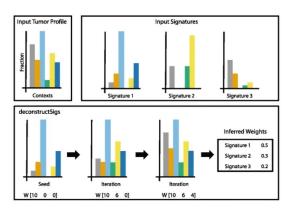
mSignatureDB (NAR 2018)

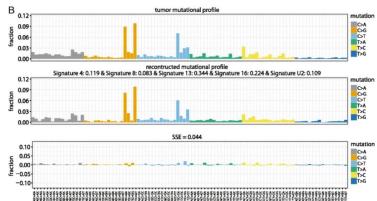
MuSiCa (BMC Bioinformatics 2018)

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)

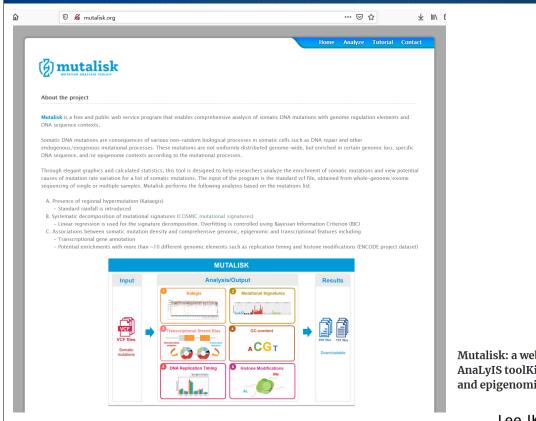
mSignatureDB (NAR 2018)

MuSiCa (BMC Bioinformatics 2018)

Mutalisk (NAR 2018)

# (3) Web interfaces: Mutalisk

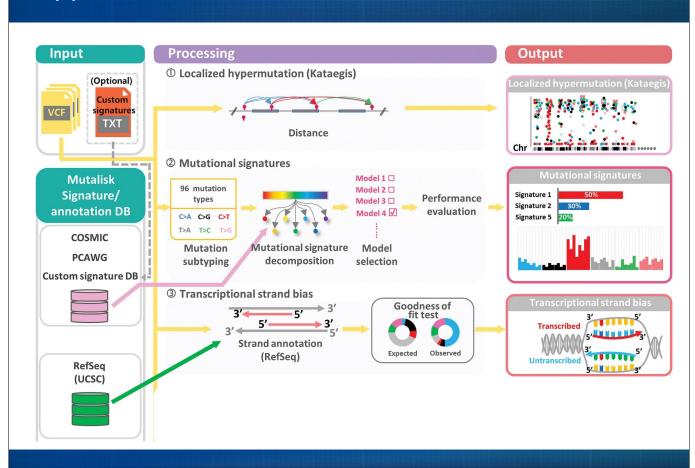
#### http://mutalisk.org



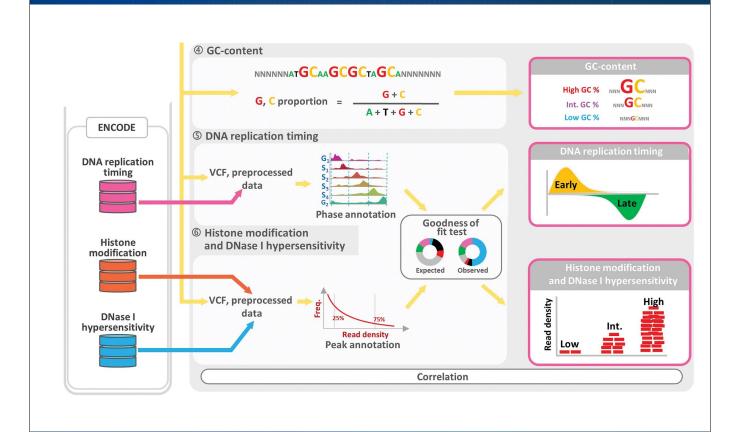
Mutalisk: a web-based somatic MUTation AnaLyIS toolKit for genomic, transcription and epigenomic signatures 3

Lee JK et al., NAR 2018

### (3) Workflow in Mutalisk



### (3) Workflow in Mutalisk (2)



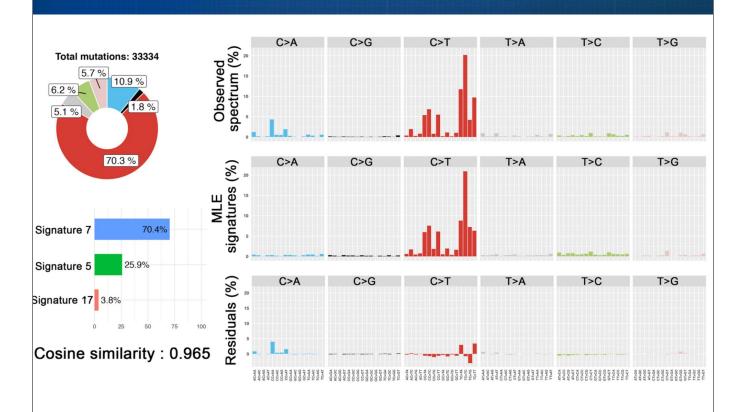
### (3) Input for Mutalisk



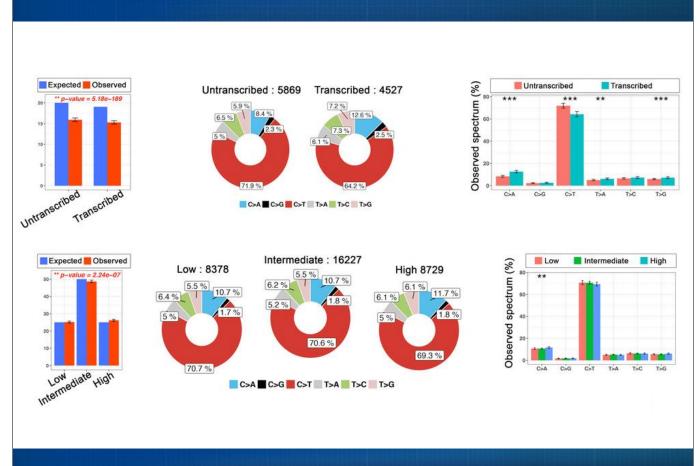
This site is optimized for Chrome. **DEMO** The following shows an example of how to run Mutalisk using the sample data. 1. Genome assembly 3. Mutational signatures 4. Genomic & epigenomic annotation GRCh37/hg19 [Homo sapiens (human)] 3-1. MLE method Linear Regression ☑ Localized hypermutation (kataegis) ☑ Transcriptional strand bias ☑ GC content 2. Input file [ ENCODE dataset reference cell ] The input file format of this tool is VCF file. 3-2. Cancer type User Selection You can select multiple files (max 300). GM12878 (Blood - Normal) 3-3. Select the mutational signatures. The total size of mutliple files should be less ☑ DNA replication timing ☐ Signature1  $\square \, Signature 2$ ☐ Signature3 ☑ DNasel hypersensitivity  $\square$  Signature4 ☐ Signature5 ☐ Signature6 ☑ Histone modification + Add Files  $\square$  Signature 7  $\square$  Signature 8 ☐ Signature9 ☐ Signature 10 ☐ Signature 11 ☐ Signature 12 No Files Selected ☐ Signature 13  $\square$  Signature 14 ☐ Signature 15 ☐ Signature 16 ☐ Signature 17 ☐ Signature 18 Reference to the genomic/epigenomic data:

\*\* The ENCODE Project & UCSC genome browser ☐ Signature 19 ☐ Signature 20 ☐ Signature21 ☐ Signature 22 ☐ Signature23 ☐ Signature24 ☐ Signature 26 ☐ Signature 25  $\square$  Signature 27 ☐ Signature 28 ☐ Signature 29 ☐ Signature30 Select All Deselect All Reference to the mutational signatures: \*\* Signatures of Mutational Processes in Human Cancer PCAWG - SigProfiler (provisional) Custom signatures

# (3) Output in Mutalisk (1)

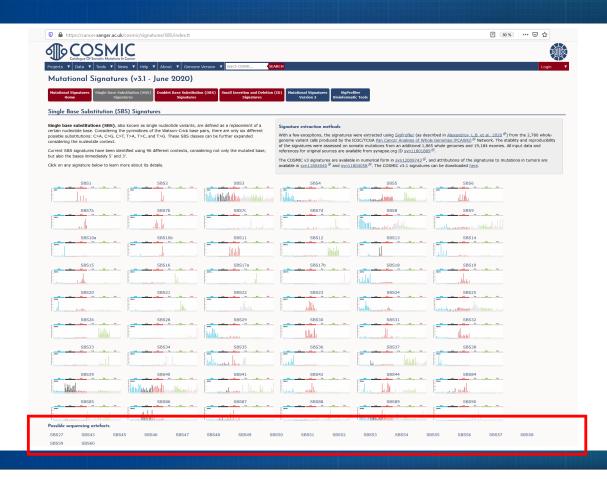


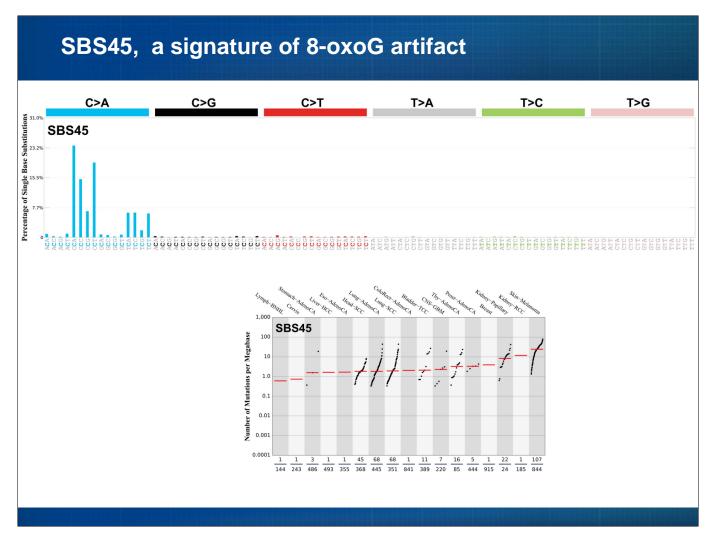
### (3) Output in Mutalisk (2)

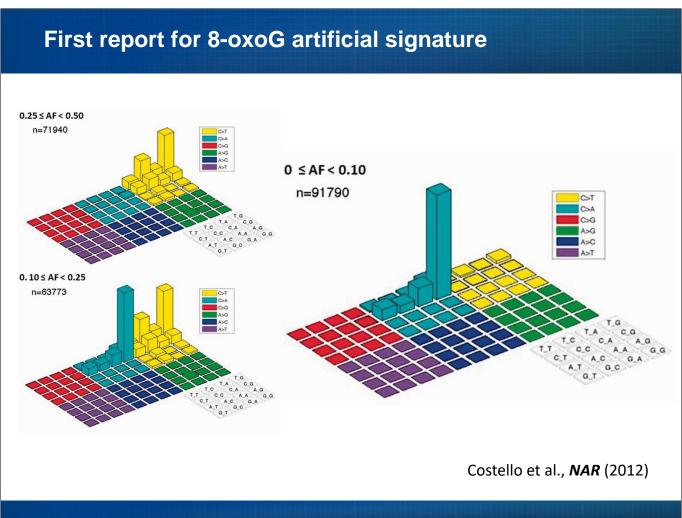


Genome QC with mutational signatures

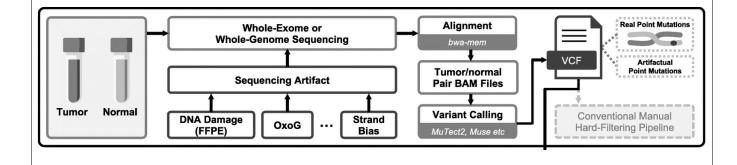
#### Amplification/sequencing artifacts make unique signatures



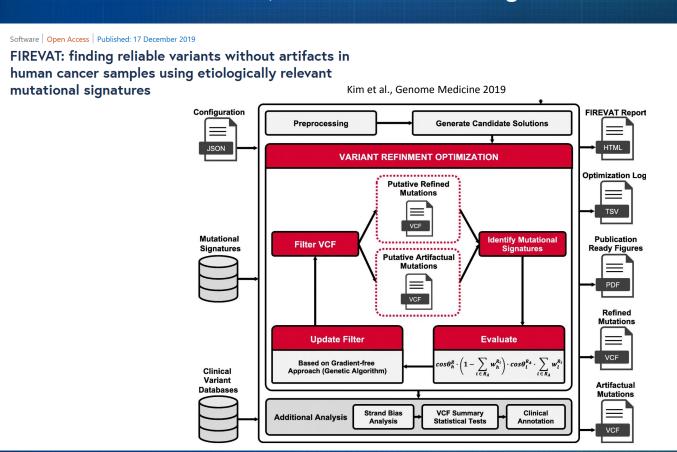




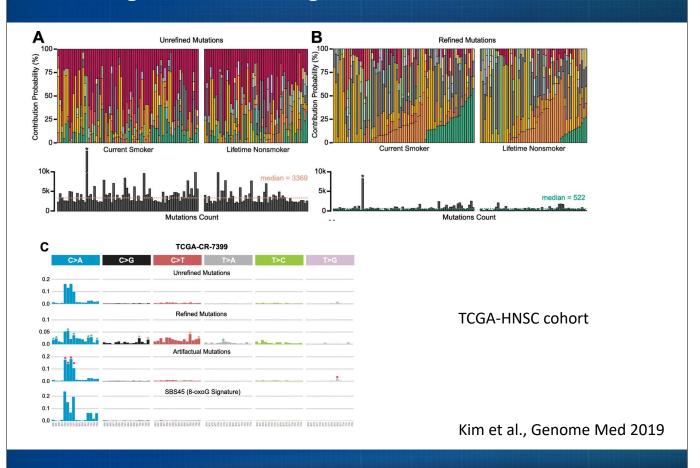
### A typical pipeline for cancer genome analyses



### Workflow in FIREVAT, a software for filtering artifacts



### **Filtering mutations using FIREVAT**



### 전망

- 돌연변이 signature 개수는 총 몇 개가 될까?
- 돌연변이 signature 각각의 원인을 규명할 수 있을까?
- Structural variation의 signature는 무엇이 있을까?

### **Summary**

- 돌연변이는 random 하게 생기지 않는다
- Mutational signature 개념을 이용하여 정확한 variant calling 을 할 수 있다
- Mutational signature 개념을 이용하여 돌연변이가 만들어 진 원인을 추적할 수 있다
- Mutational signature를 구하는 tool을 이해하고 사용할 수 있다.