

# KSBI-BIML 2026

Bioinformatics & Machine Learning(BIML)  
Workshop for Life Scientists

생명정보학 & 머신러닝 워크샵(온라인)



## Introduction to cancer-immune analysis

김상우 \_ 연세대학교



**KSBI**  
KOREAN SOCIETY FOR  
BIOINFORMATICS

| 한국생명정보학회



본 강의 자료는 한국생명정보학회가 주관하는 BIML 2026 워크샵을 목적으로  
제작된 것으로 해당 목적 이외의 다른 용도로 사용할 수 없음을 분명하게 알립니다.

이를 다른 사람과 공유하거나 복제, 배포, 전송할 수 없으며 만약 이러한 사항을 위반할 경우  
발생하는 **모든 법적 책임은 행위자 본인에게 있음**을 알립니다.

# KSBI-BIML 2026

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

한국생명정보학회가 주최하는 BIML-2026 동계 Bioinformatics & Machine Learning 교육 워크숍에 여러분을 초대합니다.

BIML 워크숍은 생명정보학 연구자들이 최신 AI바이오 분야의 인공지능 기반 분석 기술과 바이오 데이터 분석 기법을 이론과 실습을 통해 체계적으로 배울 수 있는 전문 교육 프로그램입니다. 2015년에 시작된 BIML 워크숍은 올해로 12년 차를 맞이하며, 국내 생명정보학 분야의 최초이자 최고 수준의 교육 프로그램으로 자리 잡았습니다. 이번 워크숍은 크게 인공지능바이오(AI바이오) 분야와 디지털바이오 분야, 두 분야로 구성됩니다.

AI바이오 분야에서는 생명정보 분석에 폭넓게 응용되고 있는 다양한 인공지능 기반 자료 모델링 기법을 다룰 예정입니다. 특히, 인공지능 심층학습을 활용한 단백질 구조 예측, 유전체 분석, 신약 개발에 대한 이론 및 실습 강의를 진행됩니다.

또한 디지털바이오 분야에서는 단일세포오믹스, 공간오믹스, 멀티오믹스, 메타오믹스에 대한 강의도 마련되어 있어, 연구자들의 분석 역량 강화에 실질적인 도움을 줄 것으로 기대됩니다.

또한 2024년부터 추가된 의료정보 자료 분석을 다루는 강의를 올해도 지속해서 운영하고자 합니다. 이는 최근 의료정보 자료 분석에 관한 연구 수요 증가를 반영한 것으로, 관련 연구를 수행하는 의과학자 및 의료정보 연구자들에게 유용한 지침을 제공할 것입니다.

또한, 올해도 생명정보학 기술의 다양화에 발맞춰 온라인 강좌를 대폭 확대했습니다. 올해는 무료 강좌 10개를 포함한 총 40개 이상의 강좌가 개설되며, 연구 주제에 맞는 강좌 추천과 강연료 할인 혜택도 제공합니다.

BIML-2026는 국내 주요 연구 중심 대학의 전임 교수 및 각 분야 최고 전문가들의 강의로 구성되어 있으며, 기초 이론부터 최신 연구 동향까지 아우르는 심도 있는 교육의 장이 될 것으로 확신합니다.

여러분의 많은 관심과 참여를 기대합니다!

2026년 2월

한국생명정보학회장 류 성 호

# Introduction to Cancer Immune Analysis

암은 인간의 면역과 밀접한 관계를 가진다. 암이 처음 생겨나는 과정에서 다양한 면역을 이겨내고 무력화시키기도 하고, 암을 치료하는 과정에서도 면역이 적극적으로 활용되기도 한다. 암이 가지는 신항원 (neoantigen) 은 암 면역치료의 핵심 타겟이 되는 한편, 암 주변의 미세환경 (microenvironment) 에 따라 그 효과가 달라지기도 한다. 이렇듯, 암의 예방과 치료에 대한 핵심전략으로 떠오르는 면역과의 상관성을 분석하는 것은 암 유전체학의 매우 중요한 부분이다.

본 강의에서는 WES, RNA-seq, Single-cell 및 Spatial Transcriptomics 를 기반으로 한 암 면역성과 미세환경을 분석하는 방법에 대한 전반적인 이론과 실습을 수행한다. 이를 통해 면역치료의 타겟, 바이오마커 발굴, 종양의 면역학적 특성을 이해할 수 있다.

강의는 다음의 내용을 포함한다:

- 암 면역성과 면역치료 전략 (이론)
- DNA-seq을 이용한 종양 내 신항원 예측 분석 (이론 및 실습)
- RNA-seq 을 이용한 종양미세환경 분석 (이론 및 실습)
- Single-cell 및 spatial transcriptomics를 이용한 종양미세환경 분석 (이론 및 실습)

\* 교육생준비물: 노트북 (메모리 8GB 이상, 디스크 여유공간 30GB 이상)

\* 강의 난이도: 중급

\* 강의: 김상우 교수 (연세대학교 의과대학) / 홍지윤 조교

# Curriculum Vitae

Speaker Name: Sangwoo Kim, Ph.D.



## ► Personal Info

Name Sangwoo Kim  
Title Associate Professor  
Affiliation Yonsei University College of Medicine

## ► Contact Information

Address 50 Yonsei-ro, Seodaemun-gu, Seoul 03722, Korea  
Email swkim@yuhs.ac

---

## Research Interest

Translational Genomics, Variant analysis, Cancer Genomics, Bioinformatics

## Educational Experience

2002 B.S. in Computer Science, KAIST, Korea  
2004 M.S. in Bioinformatics, KAIST, Korea  
2010 Ph.D. in Bioinformatics, KAIST, Korea

## Professional Experience

2010-2013 Post-doc Research Fellow, UC San Diego, USA  
2014-2020 Assistant Professor, Yonsei University College of Medicine, Korea  
2021-current Associate Professor, Yonsei University College of Medicine, Korea

## Selected Publications (5 maximum)

1. Yoo-Jin Ha, Seungseok Kang, Jisoo Kim, Jun Han Kim, Se-Young Jo, and **Sangwoo Kim\***, Comprehensive benchmarking and guidelines of mosaic variant calling strategies, *Nature Methods* 2023
2. Bhumsuk Keam, Min Hee Hong, Seong Hoon Shin, Seong Gu Heo, Ji Eun Kim, Hee Kyung Ahn, Yun-Gyoo Lee, Keon-Uk Park, Tak Yun, Keun-Wook Lee, Sung-Bae Kim, Sang-Cheol Lee, Min Kyung Kim, Sang Hee Cho, So Yeon Oh, Sang-Gon Park, Shinwon Hwang, Byung-Ho Nam, **Sangwoo Kim\***, Hye Ryun Kim\*, Hwan-Jung Yun\*, *Journal of Clinical Oncology* 2023
3. Tae-Min Kim, In Seok Yang, Byung-Joon Seung, Sejoon Lee, Dohyun Kim, Yoo-Jin Ha, Mi-kyoung Seo, Ka-Kyung Kim, Hyun Seok Kim, Jae-Ho Cheong, Jung-Hyang Sur, Hojung Nam, and **Sangwoo Kim\***, Cross-species Oncogenic Signatures of Breast Cancer in Canine Mammary Tumors, *Nature Communications* 2020
4. Se-Young Jot†, Eunyong Kim†, and **Sangwoo Kim\***, Impact of mouse contamination in genomic profiling of patient-derived models and best practice for robust analysis, *Genome Biology* 2019
5. Sora Kim†, Han Sang Kim†, Eunyong Kim, Min Goo Lee, Eui-Cheol Shin, Soonmyung Paik, and **Sangwoo Kim\***, Neopepsee: accurate genome-level prediction of neoantigens by harnessing sequence and amino acid immunogenicity information, *Annals of Oncology* 2018

# Introduction to Cancer Immune Analysis

2024 BIML  
연세대학교 김상우

1

## 강의 개론

### Introduction to Cancer Immune Analysis

암은 인간의 면역과 밀접한 관계를 가진다. 암이 처음 생겨나는 과정에서 다양한 면역을 이겨내고 무력화시키기도 하고, 암을 치료하는 과정에서도 면역이 적극적으로 활용되기도 한다. 암이 가지는 신항원 (neoantigen) 은 암 면역치료의 핵심 타겟이 되는 한편, 암 주변의 미세환경 (microenvironment) 에 따라 그 효과가 달라지기도 한다. 이렇듯, 암의 예방과 치료에 대한 핵심전략으로 떠오르는 면역과의 상관성을 분석하는 것은 암 유전체학의 매우 중요한 부분이다.

본 강의에서는 WES, RNA-seq, Single-cell 및 Spatial Transcriptomics 를 기반으로 한 암 면역성과 미세환경을 분석하는 방법에 대한 전반적인 이론과 실습을 수행한다. 이를 통해 면역치료의 타겟, 바이오마커 발굴, 종양의 면역학적 특성을 이해할 수 있다.

강의는 다음의 내용을 포함한다.

- 암 면역성과 면역치료 전략 (이론)
- DNA-seq을 이용한 종양 내 신항원 예측 분석 (이론 및 실습)
- RNA-seq 을 이용한 종양미세환경 분석 (이론 및 실습)
- Single-cell 및 spatial transcriptomics를 이용한 종양미세환경 분석 (이론 및 실습)

2

# Introduction to Cancer Immune and Immunotherapy

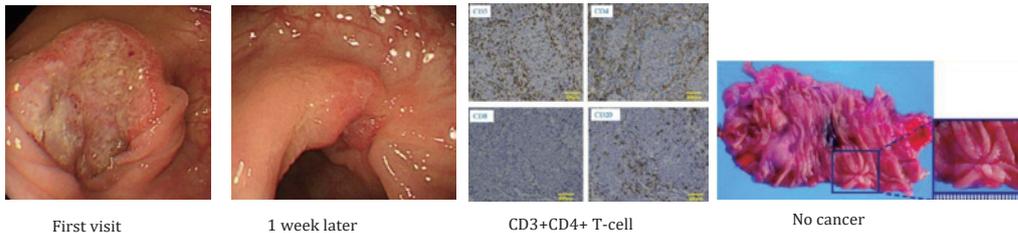
3

## *Cancer Immunotherapy:*

- Exploit **host's immune system** to treat cancer
- Generate or augment an immune response against cancer

# Immune and cancer

- Immunosuppressed patients have a higher risk for cancer
- Spontaneous regression occurs one in every 60,000 to 100,000 cancer cases



First visit

1 week later

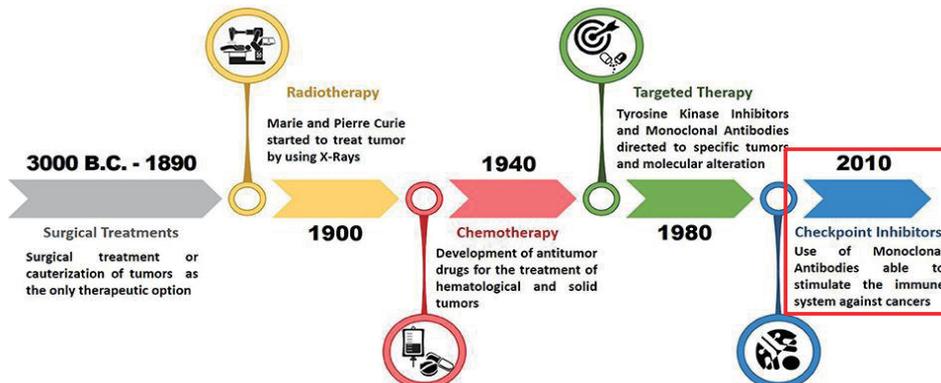
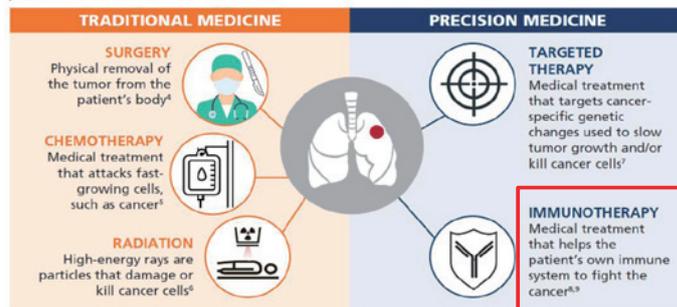
CD3+CD4+ T-cell

No cancer

Chida et al, Surg Case Rep 2017

## Cancer Immunotherapy as a new hope

Surgery, chemotherapy, and radiation have been the backbone of cancer treatment for decades, but recent advances are allowing doctors to further individualize their patients' treatment with precision medicine.<sup>2,3</sup>



# The history of immunotherapy

New York Times - July 29, 1908

## ERYSIPELAS GERMS AS CURE FOR CANCER

Dr. Coley's Remedy of Mixed  
Toxins Makes One Disease  
Cast Out the Other.

MANY CASES CURED HERE

Physician Has Used the Cure for 15  
Years and Treated 430 Cases—  
Probably 150 Sure Cures.

Following news from St. Lou's that  
two men have been cured of cancer in  
the City Hospital there by the use of a  
fluid discovered by Dr. William B.  
Coley of New York. It came out yester-  
day that nearly 100 cases of that sup-  
posedly incurable disease have been cured  
in this city during the last few years, all  
through the use of the fluid discovered  
by Dr. Coley.



erysipelas

### CONTRIBUTION TO THE KNOWLEDGE OF SARCOMA.<sup>1</sup>

By WILLIAM B. COLEY, M.D.,  
OF NEW YORK.

- I. A CASE OF PERIOSTEAL ROUND-CELLED SARCOMA OF THE METACARPAL BONE; AMPUTATION OF THE FOREARM; GENERAL DISSEMINATION IN FOUR WEEKS; DEATH SIX WEEKS LATER.
- II. THE GENERAL COURSE AND PROGNOSIS OF SARCOMA, BASED UPON AN ANALYSIS OF NINETY UNPUBLISHED CASES.
- III. THE TREATMENT OF SARCOMA BY INOCULATION WITH ERYSIPELAS, WITH A REPORT OF THREE RECENT (ORIGINAL) CASES.

**I.** THE patient a young lady, *æt.* 18, had been in perfect health from earliest childhood. The family history was likewise good with the exception of a remote tubercular tendency, and the fact that an ancestor, three generations before, had died of "cancer" of the lip, presumably epithelioma.

In the early part of July, 1890, she received a slight blow upon the back of the right hand. The hand became a little swollen and somewhat painful the first night. The next few days the pain became a trifle less and the swelling subsided, but did not entirely disappear. About a week later the swelling again began to increase very slowly, and the pain became more severe. She consulted a physician at the time of the injury, but there being no evidence of anything more than an ordinary bruise the usual local applications were applied.

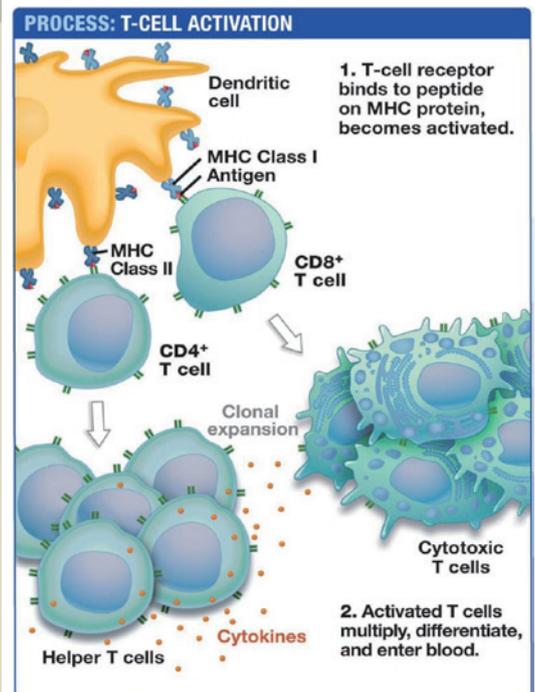
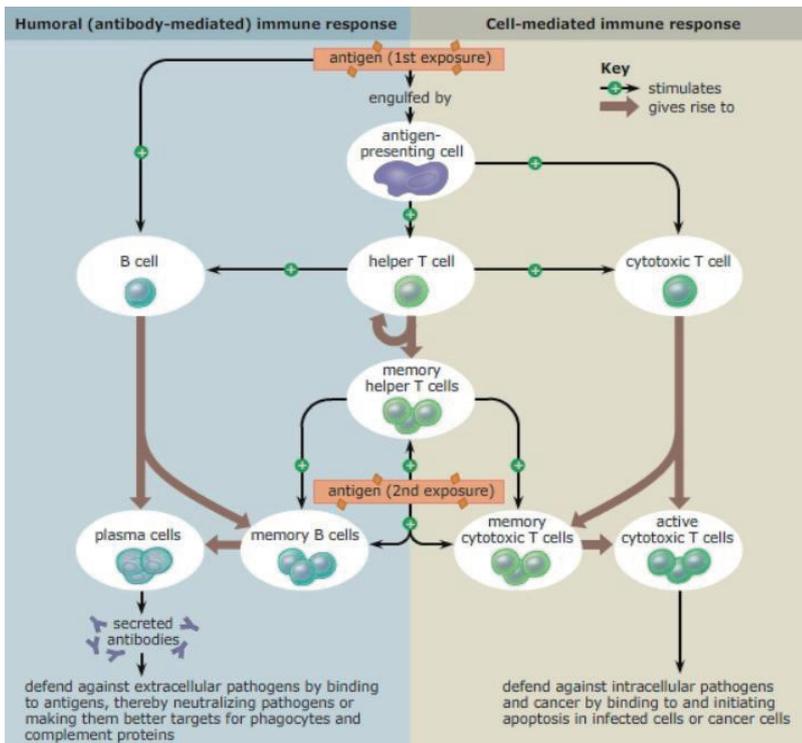
August 12. The pain and swelling continuing, she again sought

<sup>1</sup>Read before the Surgical Section of the New York Academy of Medicine, April 27, 1891. (With a report of three cases treated since).

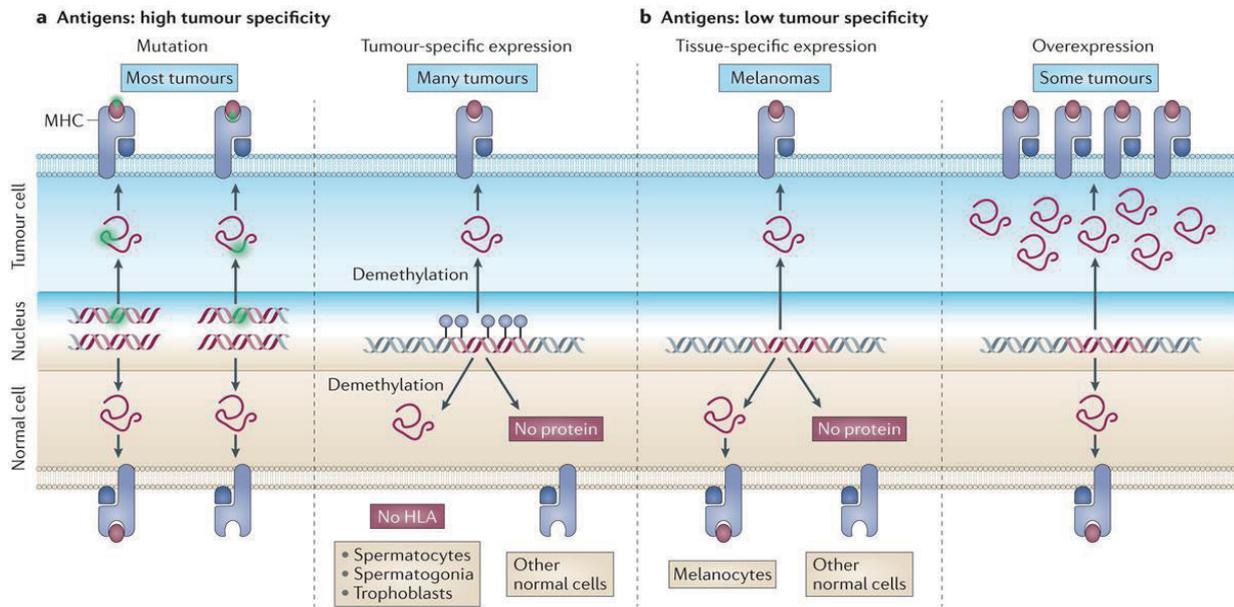
(199)

Coley, *Annals of Surgery*, 1981

# Adaptive Immunity / T-cell activation



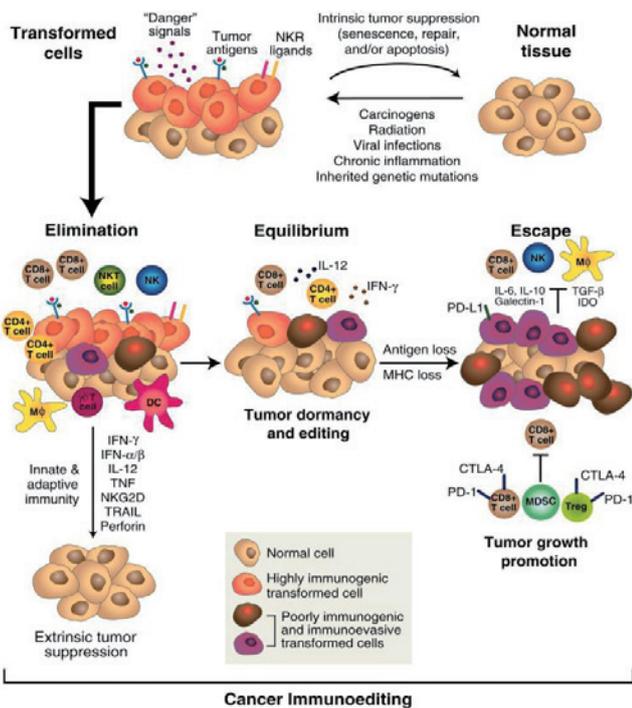
# Tumor Antigens



Nature Reviews | Cancer

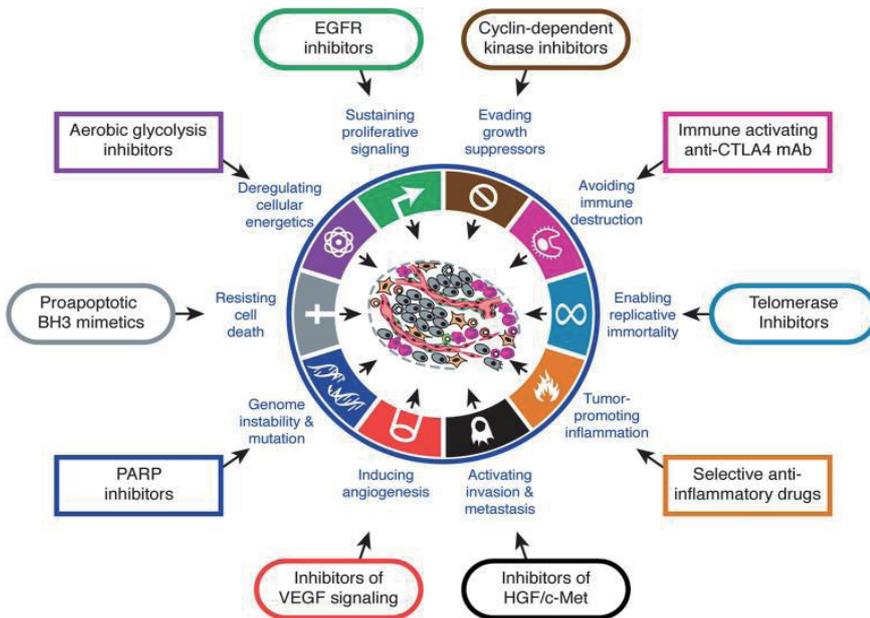
**TAA** (Tumor Associated Antigen): presented in tumor cells + (some normal cells)  
**TSA** (Tumor Specific Antigen): presented only in tumor cells

# Immunoediting of cancer



- **Elimination (immunosurveillance):**
  - Initial damage (possible destruction) of tumor cells by innate immune system
  - Tumor antigen presentation and attacked by CD4+, CD8+ T-cells
- **Equilibrium:**
  - Survived tumor cells do not progress and remain dormant
- **Escape:**
  - Cancer cells grow and metastasize due to the loss of control by the immune system

# Immune evasion



Hannahan and Weinberg, Hallmarks of cancer: The Next Generation, Cell 2011

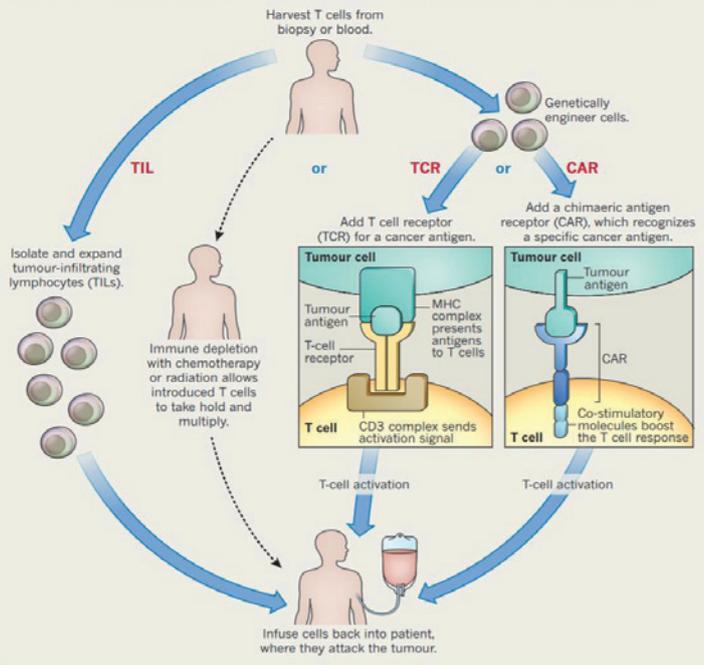
- Paralyze CTLs and NK cells by secreting TGF- $\beta$  or immunosuppressive factors
- Recruitment of regulatory T-cells (Tregs) and myeloid-derived suppressor cells (MDSCs)
- Loss of MHC class I expression

## CURRENT APPROACHES

# 1. Adoptive Cell Transfer

## CELLULAR ATTACK

Adoptive cell transfer (ACT) attacks cancer using either tumour-infiltrating lymphocytes (TILs) or genetically engineered T cells. Engineered cells are given either a new T-cell receptor (TCR) or an antibody-like molecule called a chimeric antigen receptor (CAR); both activate the T cell when they encounter a particular cancer antigen.



Courtney Humpreies, Nature 504, S13-15, 2013

- **TILs** (tumor-infiltrating lymphocytes) - metastatic melanoma
  - tissue surrounding tumor may contain immune cells and antitumor activity
  - culture TILs and re-infuse
  - deplete endogenous immune cells
- **TCR** (T-cell receptor)
  - give cells new receptor
  - viral vector in patient's T-cell
  - T-cell receptor must be genetically match to the patient's immune type
- **CAR** (chimeric antigen receptor)
  - artificial, antibody-like protein
  - antibody (binding to cancer antigen)
  - cell activating receptor
  - stimulatory molecule

## Adverse effects and personalization

Table 1

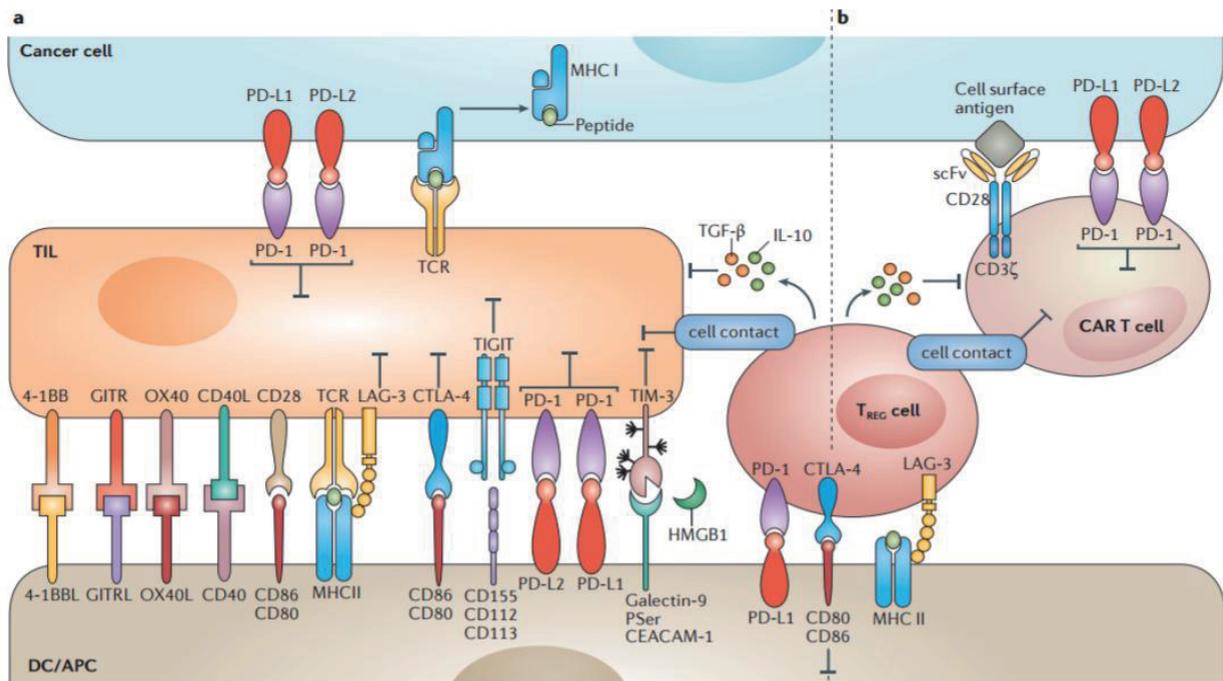
Select examples of adverse events resulting from clinical application of immunotherapies targeting public antigens

Antigen	Immunotherapy	Adverse event	Cause	Ref.
MART-1/MelanA	TCR	Fatal neural and cardiac toxicity	High levels of inflammatory cytokines alone or in combination with semi-acute heart failure and epileptic seizure	[30]
		Uveitis, Hearing loss, Loss of pigmentation	On-target activity of TCR-engineered T cells targeting normal cells expressing the cognate epitope	[24*]
NY-ESO-1	TCR + DC vaccination	Acute respiratory distress	High levels of inflammatory cytokines	[31]
	TCR (Affinity enhanced)	Skin rash with lymphocytosis, diarrheal syndrome	Autologous GVHD-like syndrome possibly due to loss of self-tolerance	[32]
MAGE-A3	TCR (Affinity enhanced)	Fatal cardiogenic shock	Cross-reactivity with an unrelated epitope from the Titin protein presented on cardiac tissue	[28]
	TCR (Affinity enhanced)	Mental status changes, comas, necrotizing leukoencephalopathy with extensive white matter defects	Reactivity to similar MAGE-A12-derived epitope presented on neural cells	[33]

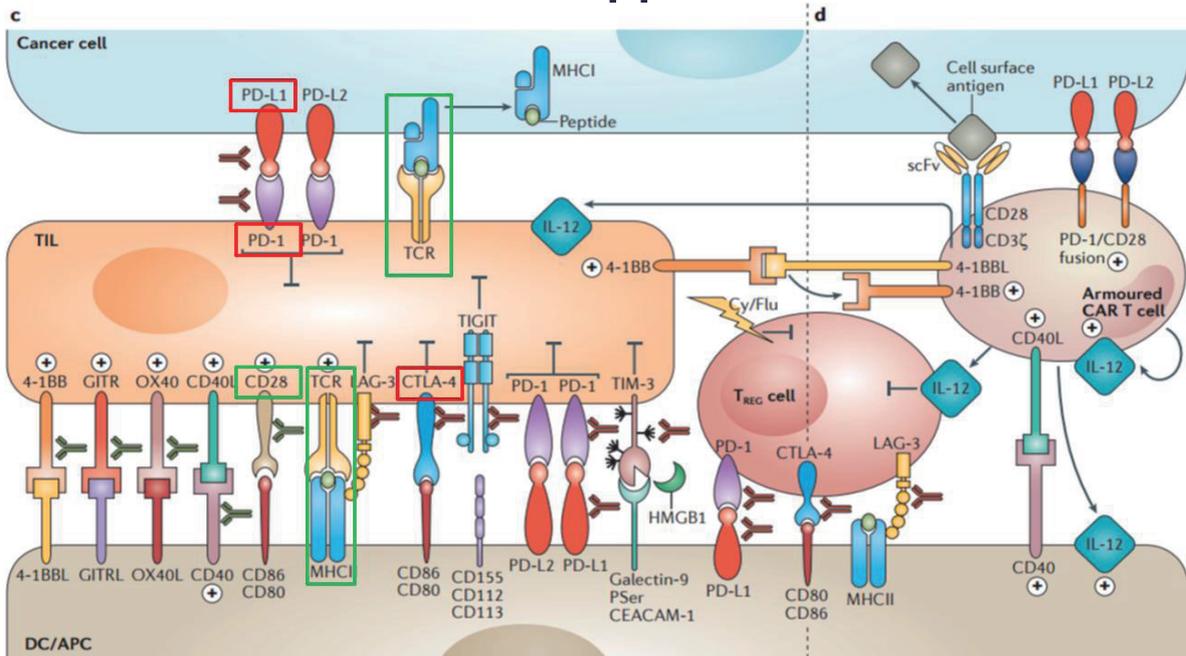
- Adverse effects in ACT
  - cytokine storm
- Need to target "tumor-specific" antigen
  - **Neoantigen?**

Courtney Humpreies, Nature 504, S13-15, 2013

## 2. Checkpoint inhibitors

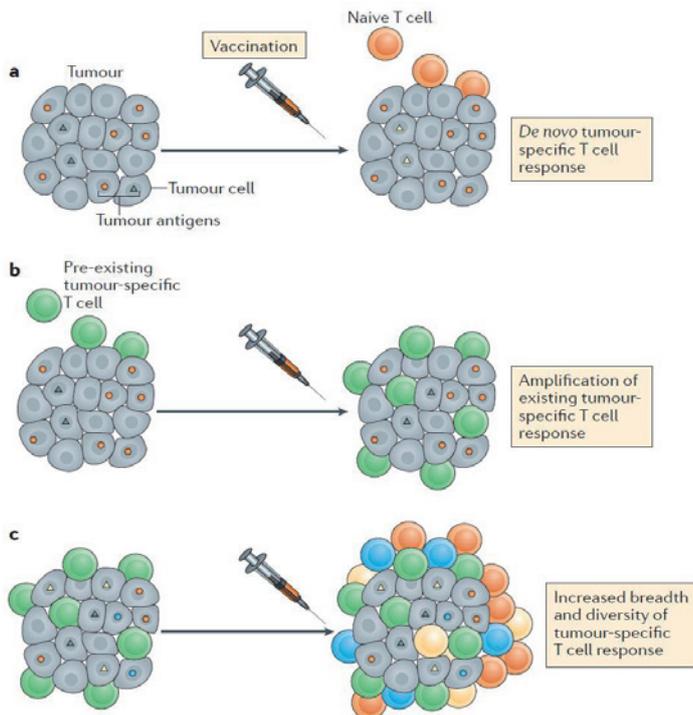


## Immunomodulatory mAbs to overcome immunosuppression





### 3. Cancer Vaccine

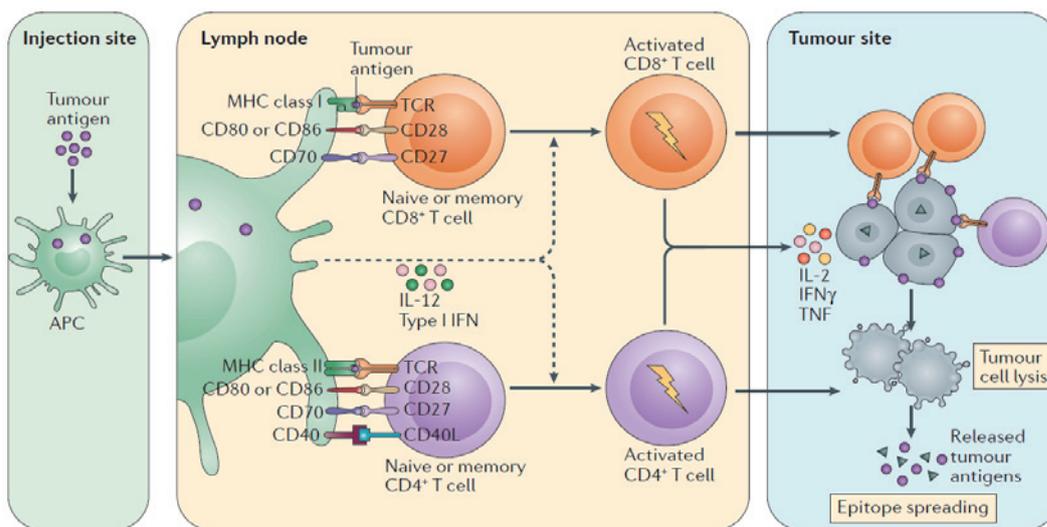


#### Cancer vaccines:

- Injection of tumor antigens
- generate new antigen-specific T-cell response
- amplification of existing T-cell response
- increase breadth and diversity of T-cell response

Hu et al, Nat. Rev. Immunol 2018

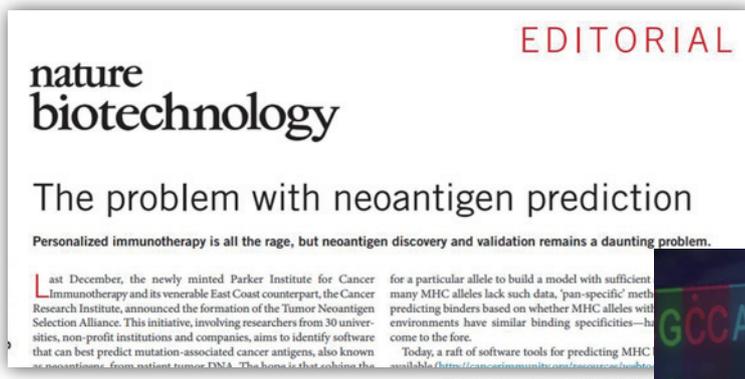
### How cancer vaccine works



Hu et al, Nat. Rev. Immunol 2018

- Antigen injection (or DC vaccine):
- Migration of APC to present antigens to T-cells (signal 1)
- Co-stimulatory signals (signal 2)
- Migration of T-cells to tumor site
- Kill tumor cells (cytotoxicity, IFN $\gamma$ , TNF..)

# Neoantigen prediction is a key challenge



- Neoantigen prediction for markers of checkpoint inhibitor
- Neoantigen prediction for finding tumor-specific (non-self) antigens for ACT

## TUMOR MUTATION BURDEN (TMB)

# Who can benefit from checkpoint inhibitor?

The NEW ENGLAND JOURNAL of MEDICINE

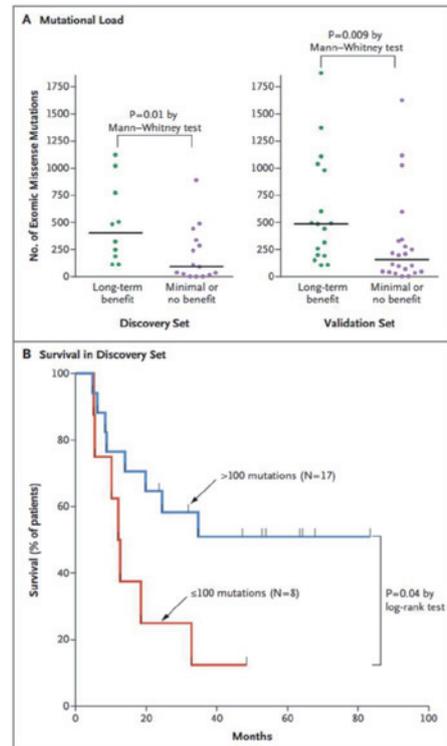
ORIGINAL ARTICLE

## Genetic Basis for Clinical Response to CTLA-4 Blockade in Melanoma

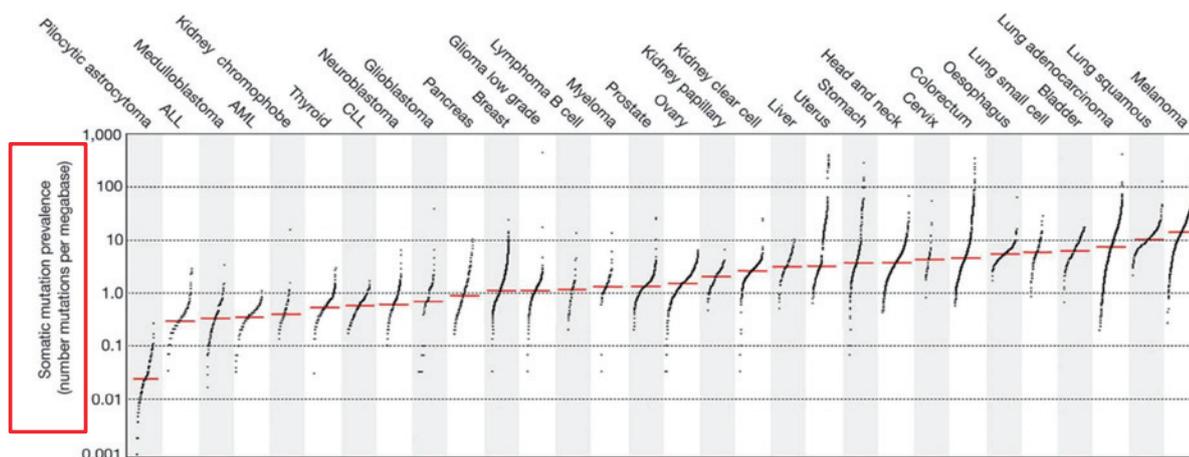
Alexandra Snyder, M.D., Vladimir Makarov, M.D., Taha Merghoub, Ph.D., Jianda Yuan, M.D., Ph.D., Jesse M. Zaretsky, B.S., Alexis Desrichard, Ph.D., Logan A. Walsh, Ph.D., Michael A. Postow, M.D., Phillip Wong, Ph.D., Teresa S. Ho, B.S., Travis J. Hollmann, M.D., Ph.D., Cameron Bruggeman, M.A., Kasthuri Kannan, Ph.D., Yanyun Li, M.D., Ph.D., Ceyhan Elipenahli, B.S., Cailian Liu, M.D., Christopher T. Harbison, Ph.D., Lisu Wang, M.D., Antoni Ribas, M.D., Ph.D., Jedd D. Wolchok, M.D., Ph.D., and Timothy A. Chan, M.D., Ph.D.

64 melanoma patients (25 discovery set, 39 validation set) treated with Ipilimumab .

Patients with high mutation burden: good survival, long-term benefit

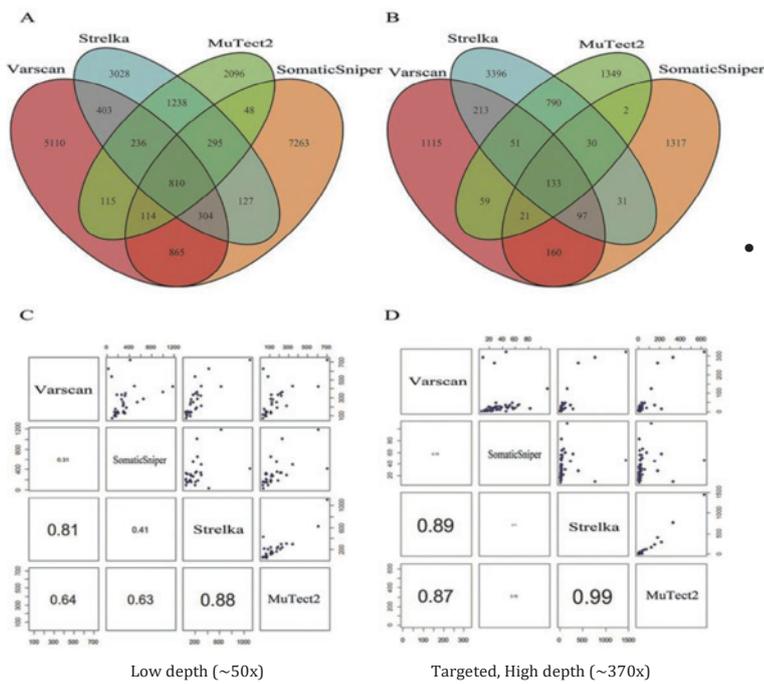


# Tumor mutation burden



• Tumor Mutation Burden (TMB) =  $\frac{\#total\_somatic\_mutation}{total\_targeted\_genome\_size(Mb)}$

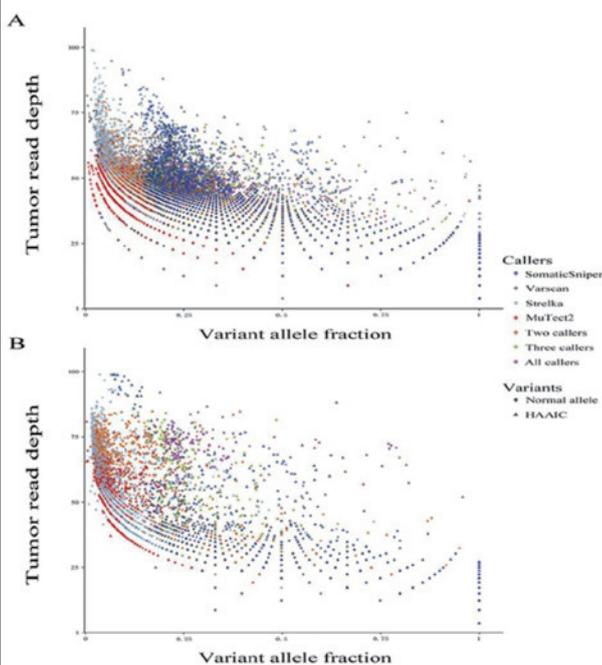
# Inconsistency of somatic mutation calls



- The number of somatic mutations are largely dependent on the variant caller used

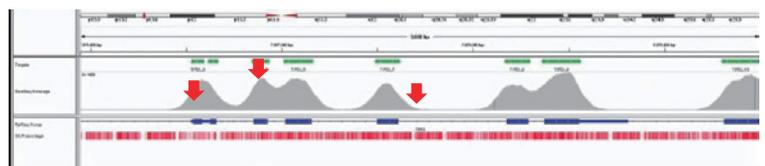
Cai et al, Sci Rep. 2016

# Tumor mutation burden



- The number of somatic mutations are largely dependent on the read depth

- And the read depth is simply not uniform



Cai et al, Sci Rep. 2016

# Fixing pipeline

mut/MB  
(SNV)

5.533993

5.178398

3.056166

1.459616

1.471475

1.453428

1.706356

mut/MB

2.991871

2.4023

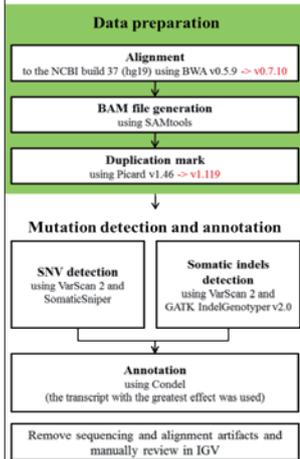
1.641857

1.27743

1.820113

1.108711

1.003712



## TCGA Flow chart

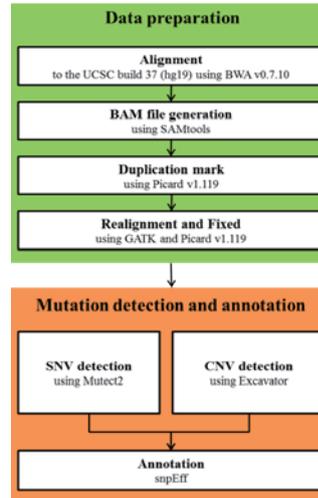
**Non-silent SNVs**  
Somatic mutations from the MAF file were filtered to remove

- (1) sites in dbSNP build 132
- (2) mutations in noncoding RNA genes
- (3) —redundant mutations in double-normal samples

8 <- covered reads in the tumor sample and 6 <- covered reads in the normal sample with 20 <- mapping quality

**High probability of being deleterious** of missense mutations by Condel

Lim SM et al, Communications Biology, in press



## In-house Flow chart

**Filter Criteria**

- (1) sites in dbSNP
- (2) filter out depth < 30
- (3) filter out allele count < 5
- (4) filter out allele frequency < 0.1
- (5) filter out non-coding region variants
- (6) Filter out Mapping quality < 30

# Potential pitfalls (use with care)

**VIEWPOINT**

### Tumor Mutation Burden—From Hopes to Doubts

**Alfredo Addo, MD**  
Oncology Department, Geneva University Hospital, Geneva, Switzerland.

**Giuseppe L. Banna, MD**  
Division of Medical Oncology, Cannizzaro Hospital, Catania, Italy.

**Glen J. Weiss, MD, MBA**  
Department of Medicine, Beth Israel Deaconess Medical Center, Harvard Medical School, Boston, Massachusetts.

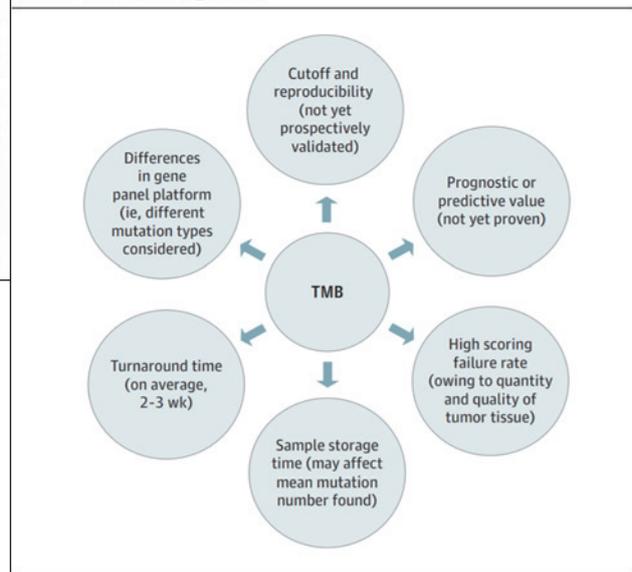
Over the past few years, the development of immune checkpoint inhibitors has altered the treatment paradigm in non-small cell lung cancer (NSCLC). Enrichment strategies have identified programmed death-ligand 1 (PD-L1) staining by immunohistochemistry to be a predictive biomarker in treatment-naïve patients with refractory NSCLC. In particular, Keynote-024<sup>1</sup> met its primary end points for overall survival (OS) and progression-free survival (PFS) in PD-L1 immunohistochemistry 50% or greater for pembrolizumab compared with platinum-based chemotherapy, validating PD-L1 immunohistochemistry as a biomarker for OS. Tumor mutation burden (TMB) has also emerged as a possible biomarker. The prevalence of somatic mutations among cancers ranges from 0.01 mutations/megabase (Mb) to more than 400 mutations/Mb. Some of these mutations lead to the translation of novel peptide epitopes or neoantigens that should enhance the immunogenicity of the tumor by eliciting T-cell responses. Initial studies of TMB were conducted by using whole-exome sequencing on tumor DNA and case-matched germline DNA.

In one study of advanced-stage NSCLC,<sup>2</sup> whole-exome sequencing was performed in 2 independent cohorts of patients with NSCLC (16 patients in one and 18 in the other) treated with pembrolizumab, and

team<sup>3</sup> recently calculated TMB scores by whole-exome sequencing in a subset of patients from the CheckMate-026 study,<sup>6</sup> a randomized phase 3 trial comparing nivolumab with platinum doublet chemotherapy as a first-line treatment in treatment-naïve patients with NSCLC with PD-L1 expression greater than 5%. Patients with a high TMB (defined as having ≥243 missense mutations) had a prolonged PFS (median PFS of 9.7 vs 5.8 months; hazard ratio [HR], 0.62; 95% CI, 0.38-1.00) and higher objective response rate (46.8% vs 28.3%) but a nonsignificant OS difference with nivolumab treatment vs chemotherapy.

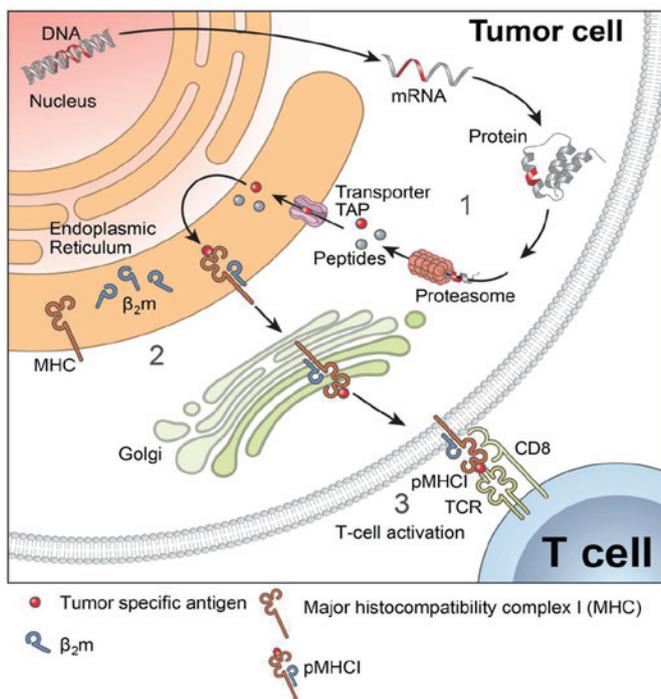
Guidelines from the European Society for Medical Oncology (ESMO) and ESMO Asia have already incorporated TMB as a possible biomarker in advanced NSCLC, recommending the combination of ipilimumab plus nivolumab as first-line treatment for patients with high TMB (>10 mutations/Mb). Supporting evidence stems from the CheckMate-227 trial, which reported results for first-line nivolumab plus ipilimumab vs platinum doublet chemotherapy.<sup>7</sup> That study showed an improved PFS in PD-L1-positive (HR, 0.62; 95% CI, 0.27-0.85) and -negative (HR, 0.48; 95% CI, 0.44-0.88) patients. At the time of publication, OS data did not meet the criteria for significance and await for release. The trial had

Figure. Pitfalls of Tumor Mutation Burden (TMB) for Clinical Application in Non-Small Cell Lung Cancer

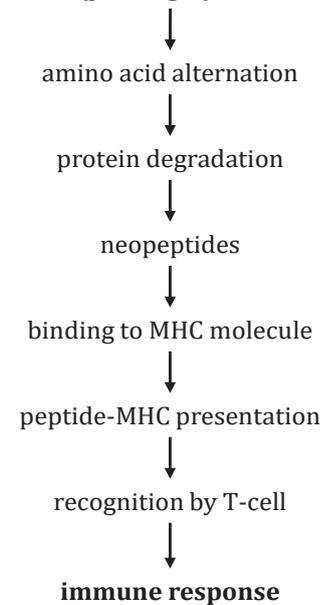


# HLA TYPING IN THE ANTIGEN PROCESSING

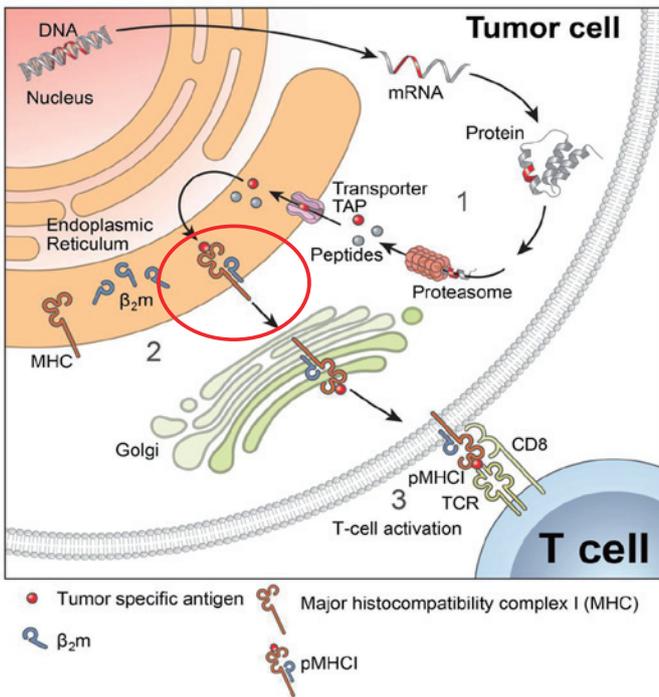
## Neoantigen processing



somatic (passenger) mutations



# Neoantigen processing



somatic (passenger) mutations

↓  
amino acid alternation

↓  
protein degradation

↓  
neopeptides

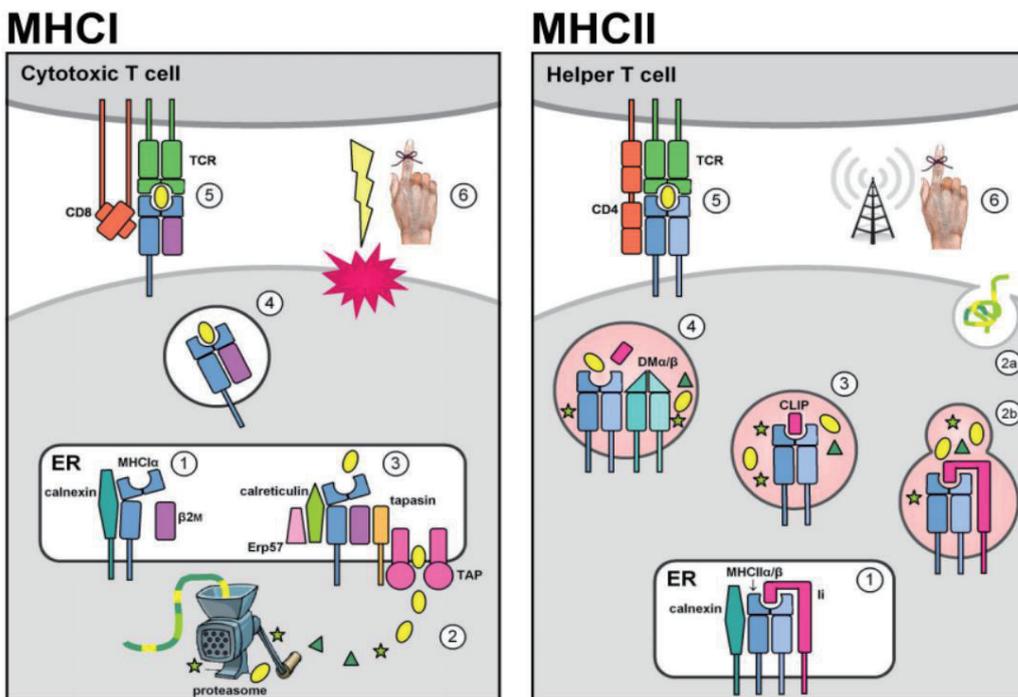
↓  
**binding to MHC molecule**

↓  
peptide-MHC presentation

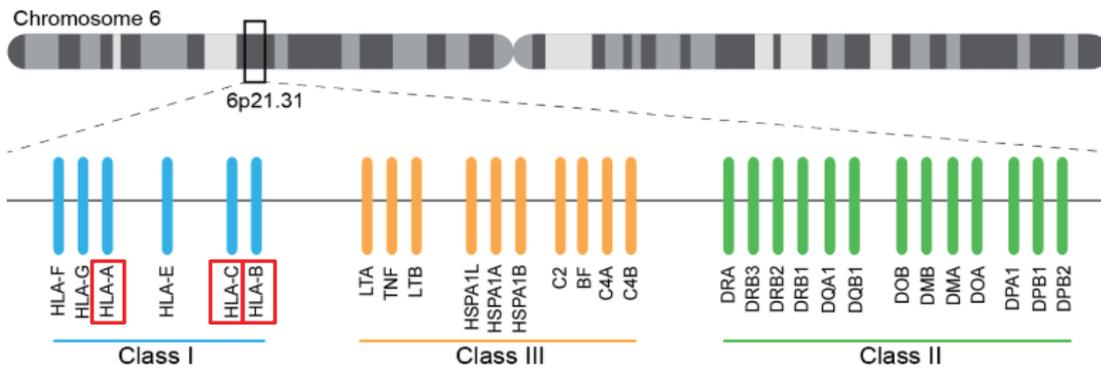
↓  
recognition by T-cell

↓  
**immune response**

# MHC (Major Histocompatibility Complex)



# HLA (Human Leukocyte Antigen)



AA Codon	3	10	15	20	AA Codon	30	35	40	45	50	55	60	65	70	75	80	85	90	
A*24:02:01:01	GC	TCC	CAC	TCC	ATG	AGG	TAT	TTC	TCC	ACA	TCC	GTG	TCC	CGG	CGC	GGC	GGG	GAG	CGC
A*24:156	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
A*24:191	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
AA Codon	25	30	35	40	AA Codon	10	15	20	25	30	35	40	45	50	55	60	65	70	
A*24:02:01:01	CGC	TTC	ATC	GCC	GTG	GGC	TAC	GTG	GAC	ACG	CAG	TTC	GTG	CGG	TTC	GAC	GAC	GAC	GCC
A*24:156	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
A*24:191	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
AA Codon	45	50	55	60	AA Codon	30	35	40	45	50	55	60	65	70	75	80	85	90	
A*24:02:01:01	CGC	AGC	CAG	AGG	ATG	GAG	CGG	CGC	CGC	CCG	TGG	ATA	GAG	CAG	GAG	GGG	CGG	CAG	TAT
A*24:156	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
A*24:191	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
AA Codon	65	70	75	80	AA Codon	55	60	65	70	75	80	85	90	95	100	105	110	115	
A*24:02:01:01	GAC	GAG	GAG	ACA	GGG	AAA	GTG	AAG	GCC	CAC	TCA	CAG	ACT	GAC	CGA	GAG	AAC	CTG	CGG
A*24:156	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
A*24:191	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
AA Codon	85	90	95	100	AA Codon	80	85	90	95	100	105	110	115	120	125	130	135	140	
A*24:02:01:01	CGC	CTC	CGC	TAC	TAC	AAC	CAG	AGC	GAG	GCC	G	AGA	CAC	AAC	TAC	GGG	GTT	GCT	GAG
A*24:156	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
A*24:191	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
AA Codon	105	110	115	120	AA Codon	100	105	110	115	120	125	130	135	140	145	150	155	160	
A*24:02:01:01	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
A*24:156	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
A*24:191	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

SBI 한국생명정보학회  
Korean Society for Bioinformatics

## HLA alleles are ethnic specific

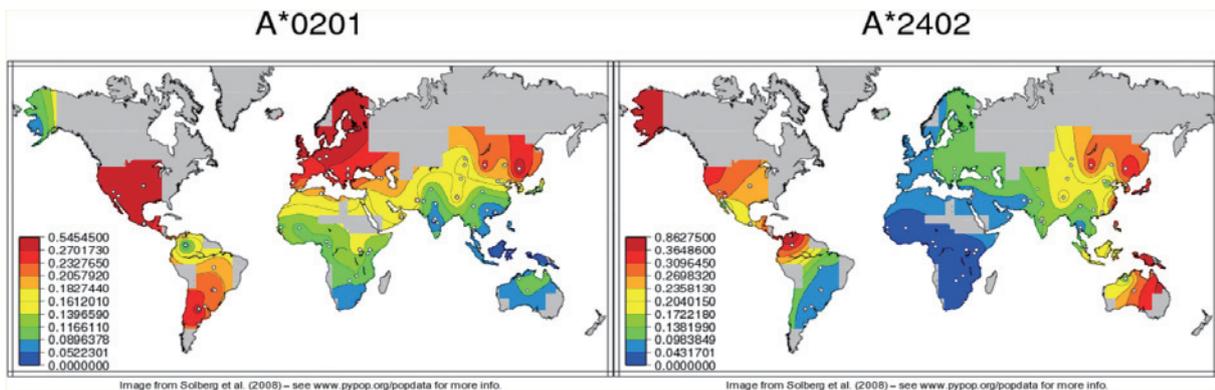
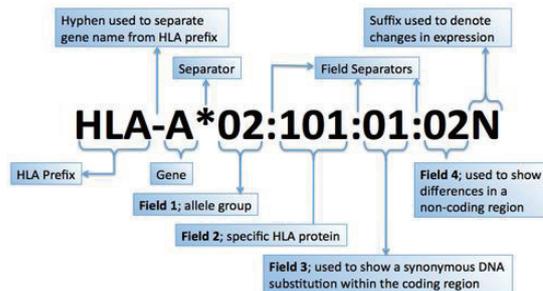


Image from Seberg et al. (2008) - see www.pyppop.org/bopdata for more info.

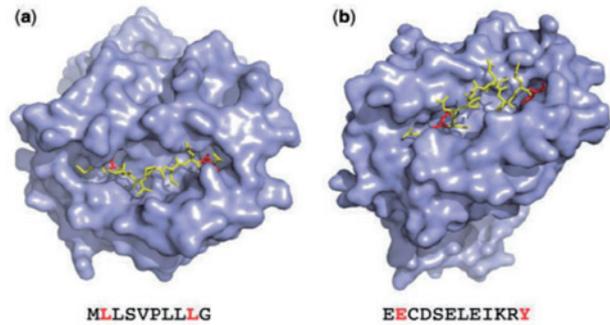
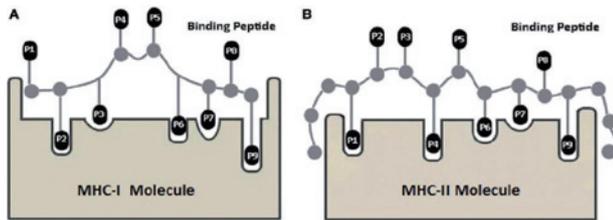
Image from Seberg et al. (2008) - see www.pyppop.org/bopdata for more info.



© SGE Marsh 04/10

SBI 한국생명정보학회  
Korean Society for Bioinformatics

# MHC-peptide binding

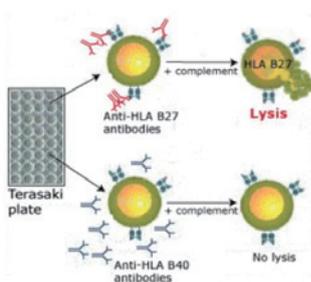


**Fig. 5.** 3D structures for two MHC class I molecules with bound peptides longer than 9 amino acids (PDB references 2CLR and 4JQX). (a) The 10mer peptide MLLSVPLLLG bound to HLA-A\*02:01 extends at the C terminus with a glycine (G) amino acid. The residues at the anchor positions P2 (L) and P9 (L) are highlighted. (b) The 12mer EECDSLEIKRY bound to HLA-B\*44:03 has anchors at its second (E) and last (Y) positions and bulges out from the middle of the MHC binding groove

But it is highly dependent on the HLA alleles  
 - That's why we need to know HLA allele (of the patient)

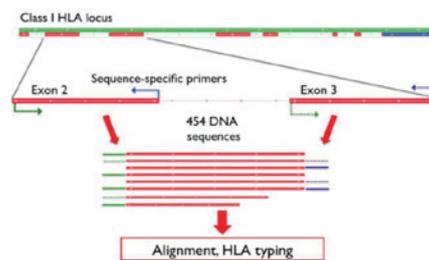
# HLA typing methods

## 1. Serology-based typing

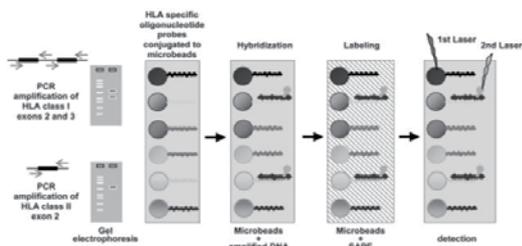


- Use of microcytotoxicity - complement mediated lysis
- Simple and low-cost
- Mostly used in HLA-A and HLA-B
- Can type allele groups and alleles only

## 2. Sanger sequencing



## 3. Sequence-specific Oligonucleotide Hybridization (SSO)



- Amplify targeted regions with biotin-labeled primers
- Hybridized sequences emit fluorescence

# NGS-based HLA typing

- **PROS**

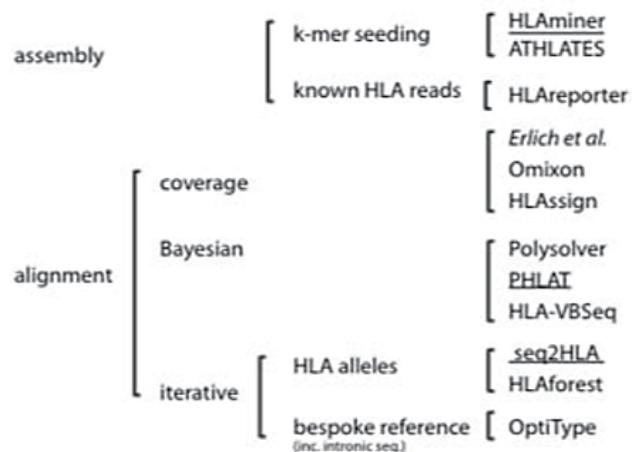
- Use of (already) produced NGS-data
- No extra-cost
- Fast

- **Threat**

- Short-read
- HLA genes are GC-rich: lower-sequencing coverage

# NGS-based HLA typing

## C Tool categories



Bauer et al, *Briefings in Bioinformatics*. 2018

# Assembly-based HLA typing

Warren et al. *Genome Medicine* 2012, 4:95  
<http://genomemedicine.com/content/4/1/95>



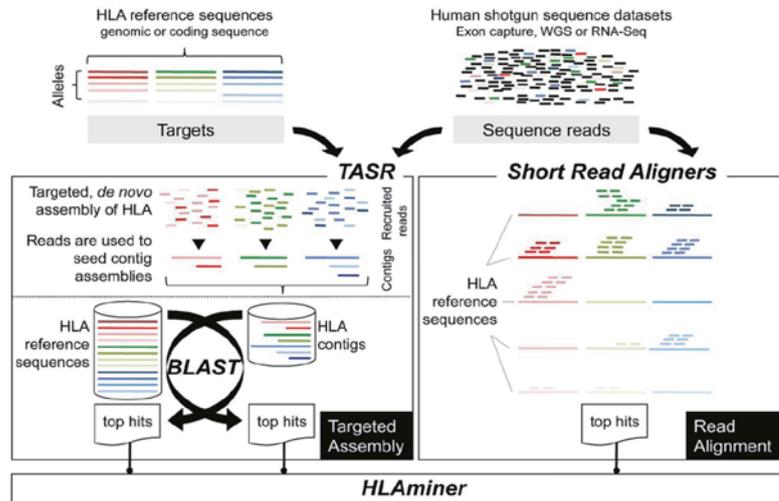
**METHOD** Open Access

## Derivation of HLA types from shotgun sequence datasets

René L Warren<sup>1</sup>, Gina Choe<sup>1</sup>, Douglas J Freeman<sup>1</sup>, Mauro Castellari<sup>1</sup>, Sarah Munro<sup>1</sup>, Richard Moore<sup>1</sup> and Robert A Holt<sup>1,2\*</sup>

**Abstract**  
 The human leukocyte antigen (HLA) is key to many aspects of human physiology and medicine. All current sequence-based HLA typing methodologies are targeted approaches requiring the amplification of specific HLA gene segments. Whole genome/exome and transcriptome shotgun sequencing can generate prodigious data but due to the complexity of HLA loci these data have not been immediately informative regarding HLA genotype. We describe HLAMiner, a computational method for identifying HLA alleles directly from shotgun sequence datasets (<http://www.biopython.org/platform/bioinformatics/software/hlaminer/>). This approach circumvents the additional time and cost of generating HLA-specific data and capitalizes on the increasing accessibility and affordability of massively parallel sequencing.

## HLAMiner



# Alignment-based HLA typing

## ANALYSIS

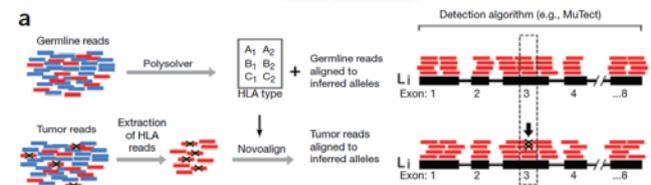
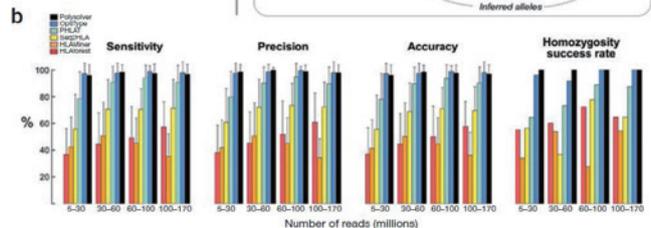
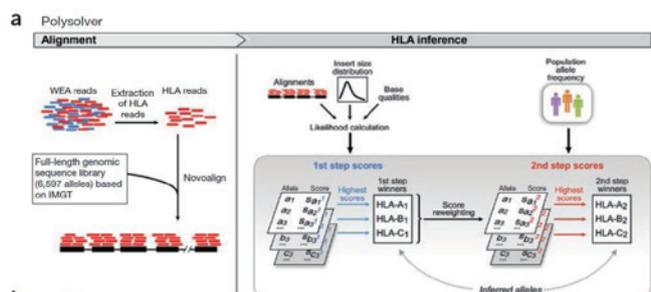
computational BIOLOGY  
 nature biotechnology

## Comprehensive analysis of cancer-associated somatic mutations in class I HLA genes

Sachet A Shukla<sup>1-3</sup>, Michael S Rooney<sup>2,4</sup>, Mohini Rajasagi<sup>1,5</sup>, Grace Tiao<sup>2</sup>, Philip M Dixon<sup>3</sup>, Michael S Lawrence<sup>2</sup>, Jonathan Stevens<sup>6</sup>, William J Lane<sup>6,7</sup>, Jamie L Dellagatta<sup>8</sup>, Scott Steelman<sup>1</sup>, Carrie Sougnez<sup>2</sup>, Kristian Cibulskis<sup>2</sup>, Adam Kiezun<sup>2</sup>, Nir Hacohen<sup>2,8,9</sup>, Vladimir Brusic<sup>1,3</sup>, Catherine J Wu<sup>1,2,5,11</sup> & Gad Getz<sup>2,10,11</sup>

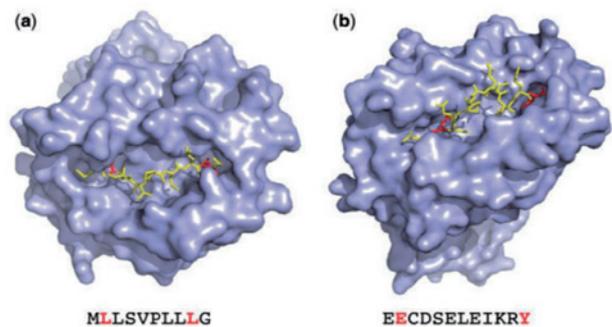
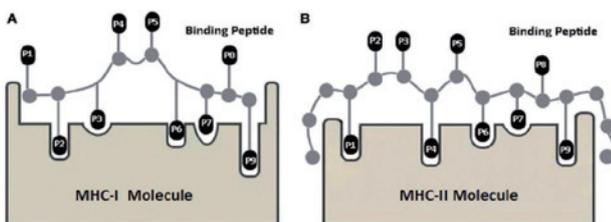
Detection of somatic mutations in human leukocyte antigen (HLA) genes using whole-exome sequencing in adenocarcinoma and diffuse large B-cell lymphoma<sup>1-5</sup>. The HLA locus, located on chromosome 6, is among the most polymorphic

## Polysolver



# MHC BINDING PREDICTION

## MHC-peptide binding

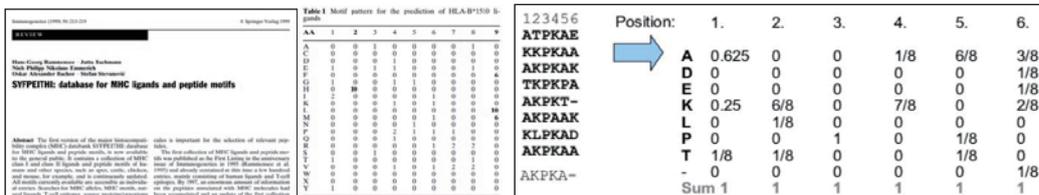


**Fig. 5.** 3D structures for two MHC class I molecules with bound peptides longer than 9 amino acids (PDB references 2CLR and 4JQX). (a) The 10mer peptide MLLSVPLLLG bound to HLA-A\*02:01 extends at the C terminus with a glycine (G) amino acid. The residues at the anchor positions P2 (L) and P9 (L) are highlighted. (b) The 12mer EECDSLEIKRY bound to HLA-B\*44:03 has anchors at its second (E) and last (Y) positions and bulges out from the middle of the MHC binding groove

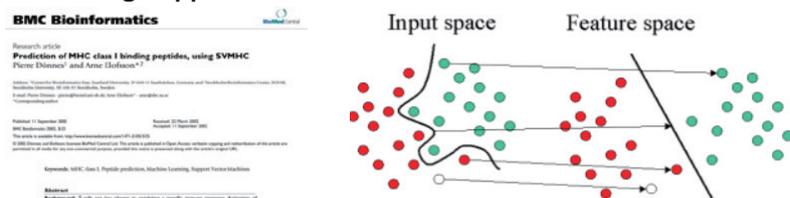
Can we predict if a given peptide will bind to MHC?

# Prediction algorithms

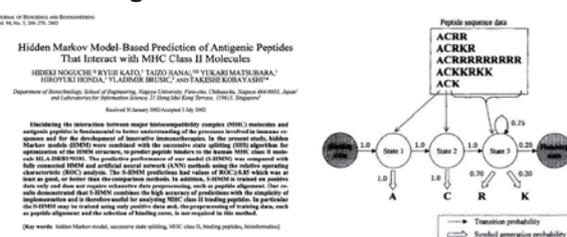
- SYFPEITHI: using PSSM



- SVMHC: using Support Vector Machine



- S-HMM: using Hidden Markov Model



SBI 한국생명정보학회  
 Korean Society for Bioinformatics

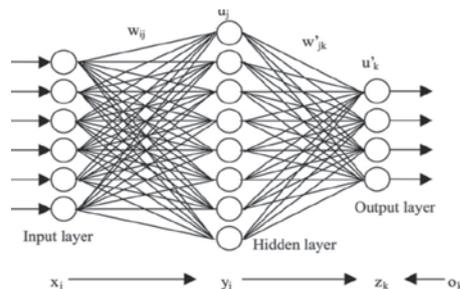
# ANN based algorithms

- NetMHC: Classification of MHC-I binding peptides using ANN

Reliable prediction of T-cell epitopes using neural networks with novel sequence representations

MORTEN NIELSEN,<sup>1</sup> CLAUD LUNDEGAARD,<sup>1</sup> PEDER WORNING,<sup>1</sup> SANNE LISE LAUTEMÖLLER,<sup>2</sup> KASPER LAMBERTH,<sup>2</sup> SØREN BUUS,<sup>2</sup> SØREN BRUNAK,<sup>1</sup> and OLE LUND<sup>1</sup>

**Abstract**  
 In this paper we describe an improved neural network method to predict T-cell class I epitopes. A novel input representation has been developed consisting of a combination of sparse encoding, Bloom encoding, and input derived from hidden Markov models. We demonstrate that the combination of several neural networks derived using different sequence-encoding schemes has a performance superior to neural networks derived using a single sequence-encoding scheme. The new method is shown to have a performance that is substantially higher than that of other methods. By use of manual information calculations we show that



- NetMHC-3.0

**BIinformatics APPLICATIONS NOTE**  
 Sequence analysis  
**Accurate approximation method for prediction of class I MHC affinities for peptides of length 8, 10 and 11 using prediction tools trained on 9mers**  
 Claus Lundegaard<sup>1</sup>, Ole Lund and Morten Nielsen  
 Center for Biological Sequence Analysis – CBS, Department of Systems Biology, The Technical University of Denmark – DTU, Kemitorvet Bldg. 205, 2800 Lyngby, Denmark.  
 Received on February 9, 2010; revised and accepted on April 4, 2010  
 Authors' addresses: clund@bioinformatics.sbi.dk, 2010

Approximation of 8, 10, 11 from 9 mer model

peptide	logscore	affinity(CW)
TEMLVENS	0.189	6489
TEMLVENS	0.351	6314
TEMLVENS	0.234	3979
TEMLVENS	0.353	3896
TEMLVENS	0.317	846
TEMLVENS	0.462	623
Genetic mean		2149

- NetMHC-4.0

(a) A I L D F T H L	(b) F Y G E R P L T R Y
X A I L D F T H L	F Y G E R P L T R Y
A X I L D F T H L	F Y G E R P L T R Y
A I X L D F T H L	F Y G E R P L T R Y
A I L X D F T H L	F Y G E R P L T R Y
A I L D X F T H L	F Y G E R P L T R Y
A I L D F X T H L	F Y G E R P L T R Y
A I L D F T X H L	F Y G E R P L T R Y
A I L D F T H X L	F Y G E R P L T R Y
A I L D F T H L X	F Y G E R P L T R Y

Gapped alignment to ANN : 9 to 8-11 mer

SBI 한국생명정보학회  
 Korean Society for Bioinformatics

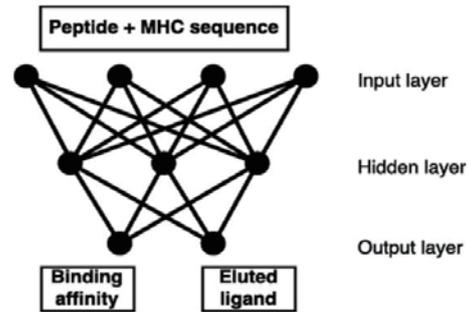
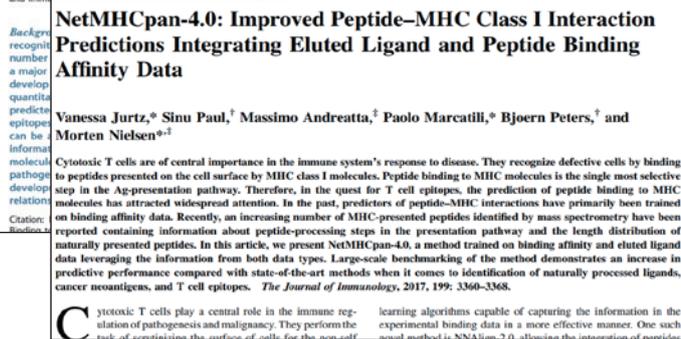
# Regarding all HLA-types at once

**NetMHCpan:** Prediction on all HLA-A/B alleles, simultaneously



Experimental data are biased to major HLA alleles  
 ▶ lack of training data in rare alleles  
 ▶ lack of accuracy

Build a classifier that work on HLA-peptide pair



# Too many methods. Need a consensus

**NetMHCcons:** Prediction on all HLA-A/B alleles, simultaneously



$$\text{NetMHCcons} = \begin{cases} \text{NetMHCpan} & \text{for } N_p < 50 \text{ and } N_b < 10 \\ \text{NetMHC} + \text{NetMHCpan} & \text{otherwise} \end{cases}$$

We demonstrate that a **simple combination of NetMHC and NetMHCpan gives the highest performance** when the allele in question is included in the training and is characterized by at least 50 data points with at least ten binders. Otherwise, NetMHCpan is the best predictor.

# Benchmarks and competitions

Journal of Immunological Methods 374 (2011) 28–34

Contents lists available at ScienceDirect

**Journal of Immunological Methods**

journal homepage: [www.elsevier.com/locate/jim](http://www.elsevier.com/locate/jim)

Research paper

**Prediction of epitopes using neural network based methods**

Claus Lundegaard<sup>a</sup>, Ole Lund, Morten Nielsen

*Center for Biological Sequence Analysis, DTU Systems Biology, Building 206, Technical University of Denmark, DK-2800 Lyngby, Denmark*

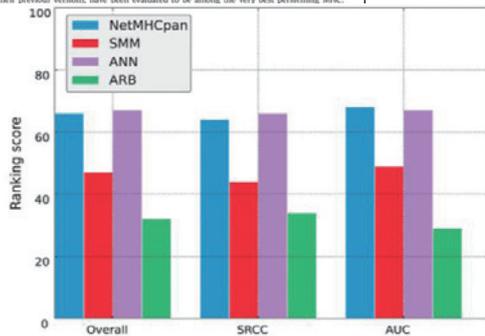
**ARTICLE INFO**

**ABSTRACT**

**Article history:**  
 Received 26 July 2010  
 Received in revised form 23 October 2010  
 Accepted 27 October 2010  
 Available online 31 October 2010

**Keywords:**  
 MHC  
 Binding  
 Prediction  
 Epitope  
 Discovery  
 T cell

In this paper, we describe the methodologies behind three different aspects of the NetMHC family for prediction of MHC class I binding, mainly to H2As. We have updated the prediction servers, NetMHC 3.2, NetMHCpan 2.2, and a new consensus method, NetMHCcons, which in their previous versions, have been evaluated to be among the very best performing MHC.



## 2nd Machine Learning Competition in Immunology 2012

Sponsors: InCoB 2012 and ICIW 2012

**Prediction task:**  
 Predict peptides naturally processed by MHC Class I pathway ("eluted peptides") for each target MHC molecule. For each target molecule, the competitors are asked to submit a set of predicted eluted peptides from the test set.

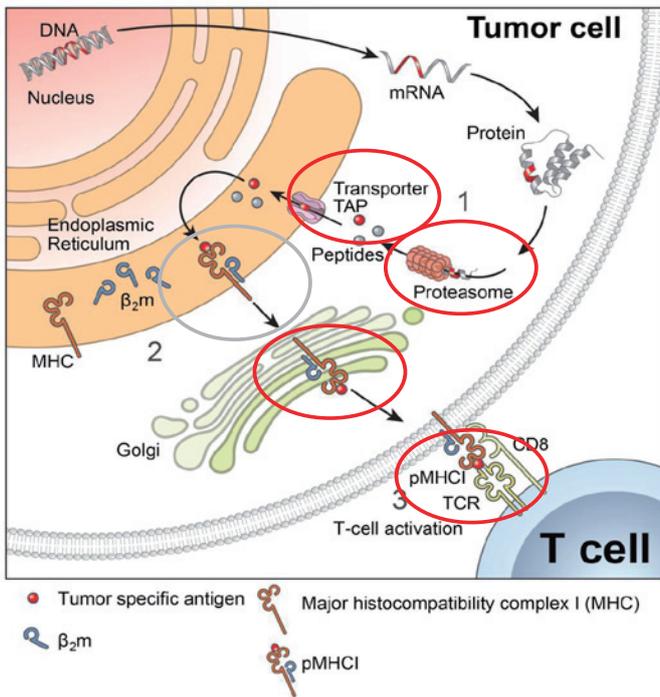


A total of 32 submissions were submitted for the competition. Of these, 24 submissions (Group 1) provided a set of thresholds (elution score based predictors) for each peptide and each MHC molecule. Another 8 submissions (Group 2) provided lists of peptides that were predicted as eluted from specific MHC molecules (eluted peptide list based predictors) for each of 8 studied MHC alleles. The NetMHC 3.2 server (1D-BENCH) results were used as a benchmark method.

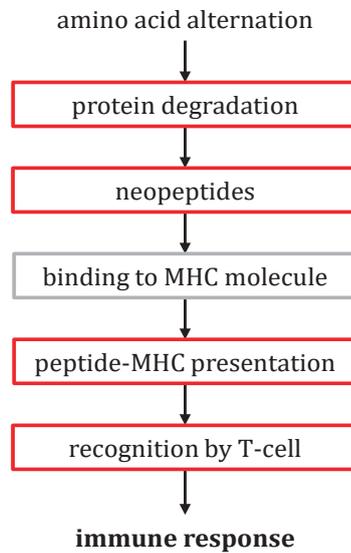
Winning Team	Predictor No.	Prediction Method	Winning Category
Lundegaard C, Lamberth K, Harndahl M, Buus S, Lund O, Nielsen M, Technical University of Denmark	1D, BENCH	NetMHC 3.2 (Reference)	Group 1: A*0201
Giguere S, Drouin A, Lacoste A, Laval University, Canada	2F	A Bayesian model averaging method over several SVMs using the GS kernel.	Group 1: B*0702, H-2D <sup>P</sup> , and H-2K <sup>D</sup>
Nielsen M, et al., Technical University of Denmark	9D	A combination of NetMHC, NetMHCpan and MHCkernel predictions.	Group 1: B*3501 and B*4403
Giguere S, Drouin A, Lacoste A, Laval University, Canada	2D	A SVM classifier and a novel string kernel (SS kernel).	Group 1: B*5301
Xiang Z, He Y, University of Michigan Medical School, Ann Arbor, MI, USA	20D	A position-specific scoring matrix (PSSM) with statistical P-value as the cutoff.	Group 1: B*5701
Yu Ting Wei, Department of Probability and Statistics, School of Mathematical Sciences, Peking University; Wen Jun Shen and Hau-San Wong, Department of Computer Science, City University of Hong Kong	14A	ConsMHC: a consensus program incorporating the results of kernelRLSpan-I, NetMHC, NetMHCpan and PickPocket by SVM	Group 2

# ANTIGEN PROCESSING STEPS

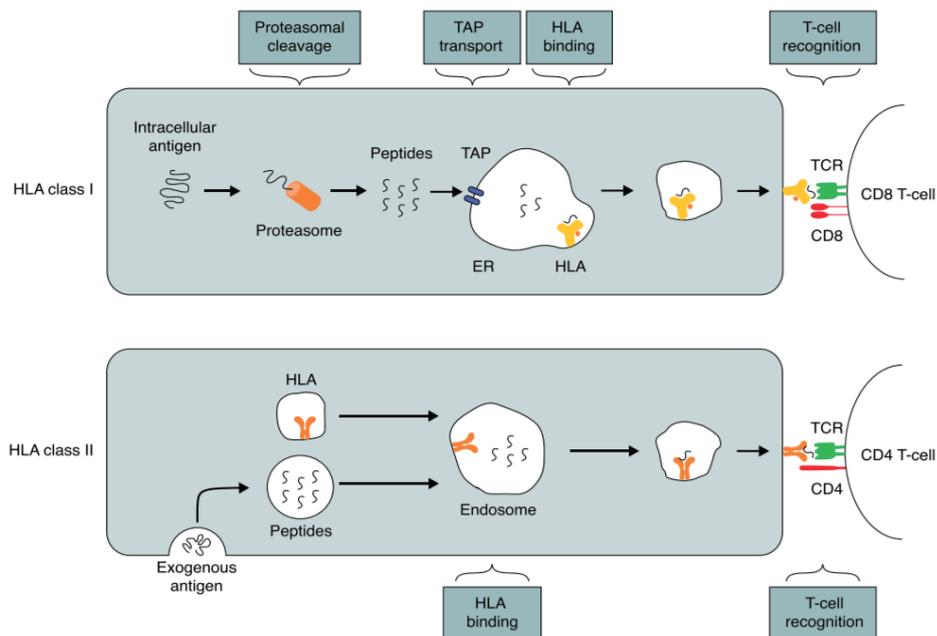
# Neoantigen processing revisited



somatic (passenger) mutations

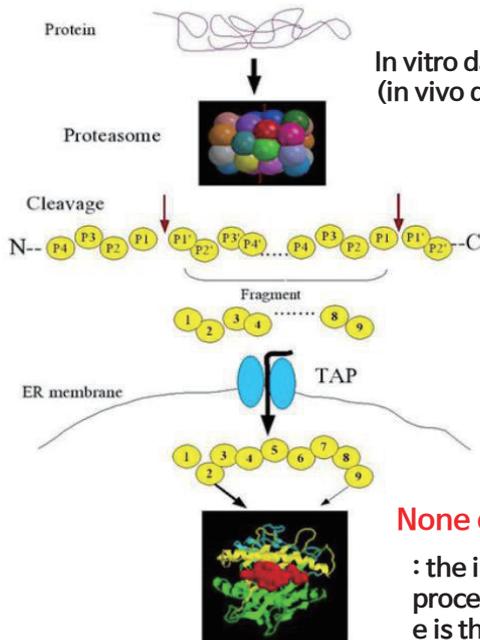


# Antigen Processing Pathways for MHC class I/II



Backert and Kohlbacher, *Genome Medicine*, 2015

# Proteasomal cleavage



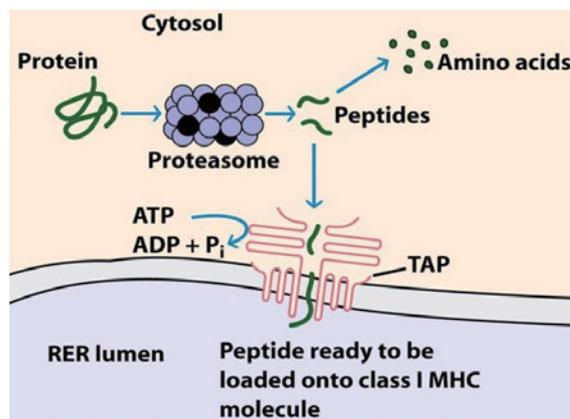
In vitro data created with purified proteasomes in the laboratory  
(in vivo data are harder to collect)

**C**-terminus: commonly determined by proteasomal cleavage  
**N**-terminus: can undergo further trimming by proteases located in the cytosol or ER

None of the predictors achieved an MCC above 0.3

: the in vitro data do not capture the full complexity of proteasomal processing in vivo. The value of predictions of proteasomal cleavage is thus rather limited

# TAP transport prediction



- Primarily owing to the scarcity of data, there are few published methods on TAP transport prediction.
- No unbiased blind benchmarks for TAP transport methods have been published so far, and a comparative assessment of the various methods is thus currently difficult

# Considering MHC-binding stability, not affinity

European Journal of  
Immunology

## Peptide-MHC class I stability is a better predictor than peptide affinity of CTL immunogenicity

Mikkel Harndahl<sup>1</sup>, Michael Rasmussen<sup>1</sup>, Gustav Roder<sup>1</sup>, Ida Dalgaard Pedersen<sup>1</sup>, Mikael Sørensen<sup>2</sup>, Morten Nielsen<sup>2</sup> and Søren Buus<sup>1</sup>

<sup>1</sup> Laboratory of Experimental Immunology, Faculty of Health Sciences, University of Copenhagen, Denmark

<sup>2</sup> Center for Biological Sequence Analysis, Department of Systems Biology, Technical University of Denmark, Denmark

Efficient presentation of peptide-MHC class I (pMHC-I) complexes to immune T cells should benefit from a stable peptide-MHC-I interaction. However, it has been difficult to distinguish stability from other requirements for MHC-I binding, for example, affinity. We have recently established a high-throughput assay for pMHC-I stability. Here, we have generated a large database containing stability measurements of pMHC-I complexes, and re-examined a previously reported unbiased analysis of the relative contributions of antigen processing and presentation in defining cytotoxic T lymphocyte (CTL) immunogenicity [Assarsson et al., J. Immunol. 2007. 178: 7890-7901]. Using an affinity-balanced approach, we demonstrated that immunogenic peptides tend to be more stably bound to MHC-I molecules compared with nonimmunogenic peptides. We also developed a bioinformatics method to predict pMHC-I stability, which suggested that 30% of the nonimmunogenic binders hitherto classified as "holes in the T-cell repertoire" can be explained as being unstably bound to MHC-I. Finally, we suggest that nonoptimal anchor

### Binding (kinetic) stability

We also developed a bioinformatics method to predict pMHC-I stability, which suggested that 30% of the nonimmunogenic binders hitherto classified as "holes in the T-cell repertoire" can be explained as being unstably bound to MHC-I.

SBI 한국생명정보학회  
Korean Society for Bioinformatics

# Prediction on the stability

## NetMHCstab: predicting stability of pMHC-I complexes

Immunology  
The Journal of cells, molecules, systems and technologies  
IMMUNOLOGY ORIGINAL ARTICLE

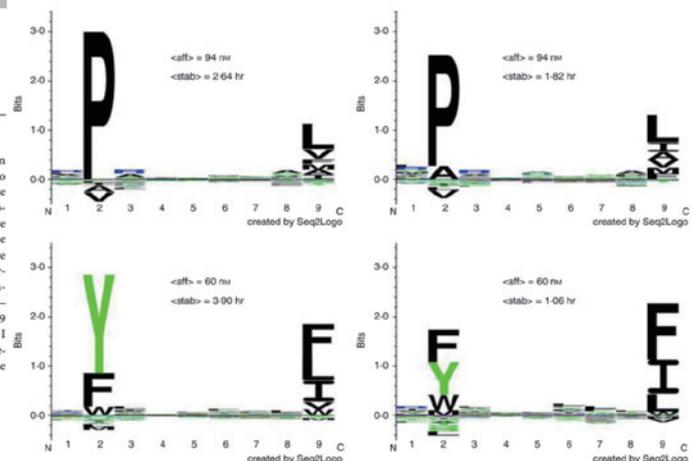
### NetMHCstab – predicting stability of peptide-MHC-I complexes; impacts for cytotoxic T lymphocyte epitope discovery

Kasper W. Jørgensen,<sup>1,2</sup> Michael Rasmussen,<sup>2,3</sup> Søren Buus<sup>2</sup> and Morten Nielsen<sup>1,3</sup>

<sup>1</sup>Department of Systems Biology, Centre for Biological Sequence Analysis, Technical University of Denmark, Lyngby, <sup>2</sup>Laboratory of Experimental Immunology, University of Copenhagen, Copenhagen N, Denmark, and <sup>3</sup>Instituto de Investigaciones Biológicas, Universidad Nacional de San Martín, San Martín, Buenos Aires, Argentina

#### Summary

Major histocompatibility complex class I (MHC-I) molecules play an essential role in the cellular immune response, presenting peptides to cytotoxic T lymphocytes (CTLs) allowing the immune system to scrutinize ongoing intracellular production of proteins. In the early 1990s, immunogenicity and stability of the peptide-MHC-I (pMHC-I) complex were shown to be correlated. At that time, measuring stability was cumbersome and time consuming and only small data sets were analysed. Here, we investigate this fairly unexplored area on a large scale compared with earlier studies. A recent small-scale study demonstrated that pMHC-I complex stability was a better correlate of CTL immunogenicity than peptide-MHC-I affinity. We here extended this study and analysed a total of 5509 distinct peptide stability measurements covering 10 different HLA class I molecules. Artificial neural networks were used to construct stability predictors capable of predicting the half-life of the pMHC-I complex. These



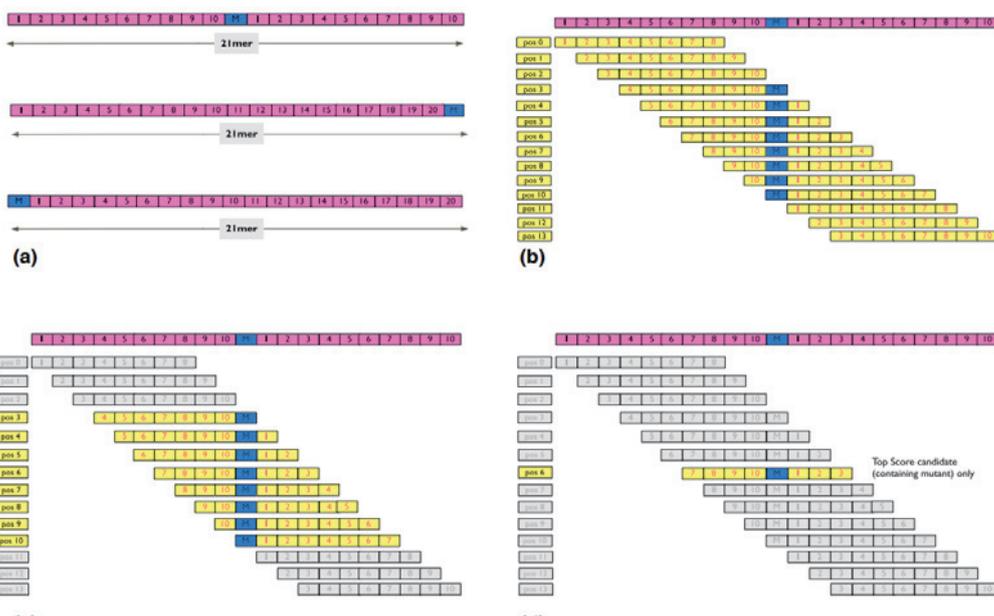
SBI 한국생명정보학회  
Korean Society for Bioinformatics

stable

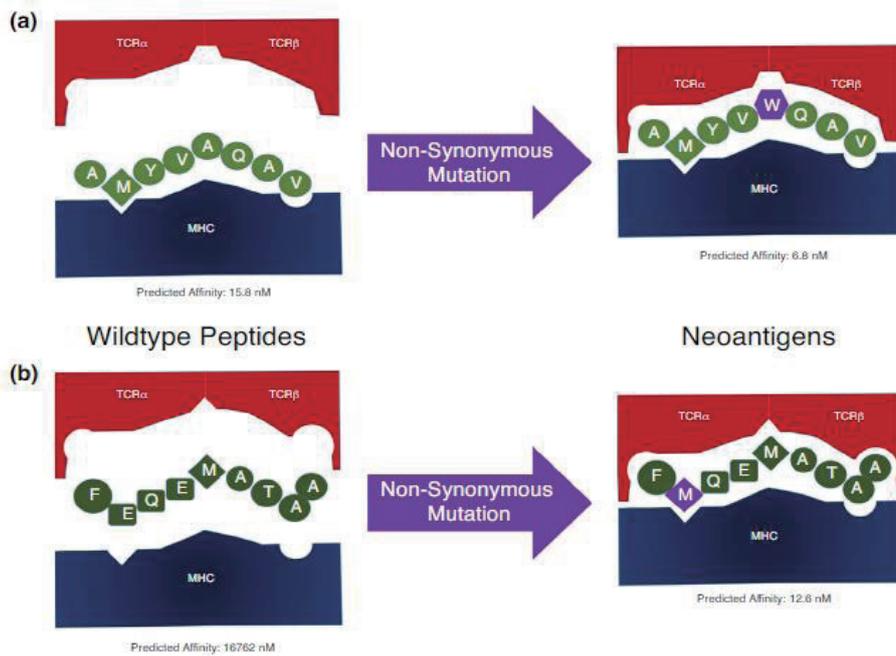


# NEOANTIGEN ANALYSIS & INTEGRATED PIPELINES

## Somatic mutation derived neopeptide



# And Neoantigens



Oiseth et al, *J Cancer Metastasis and Treatment*, 2017

## Overall Pipeline

Hundal et al. *Genome Medicine* (2016) 8:11  
DOI 10.1186/s13073-016-0264-5

Genome Medi

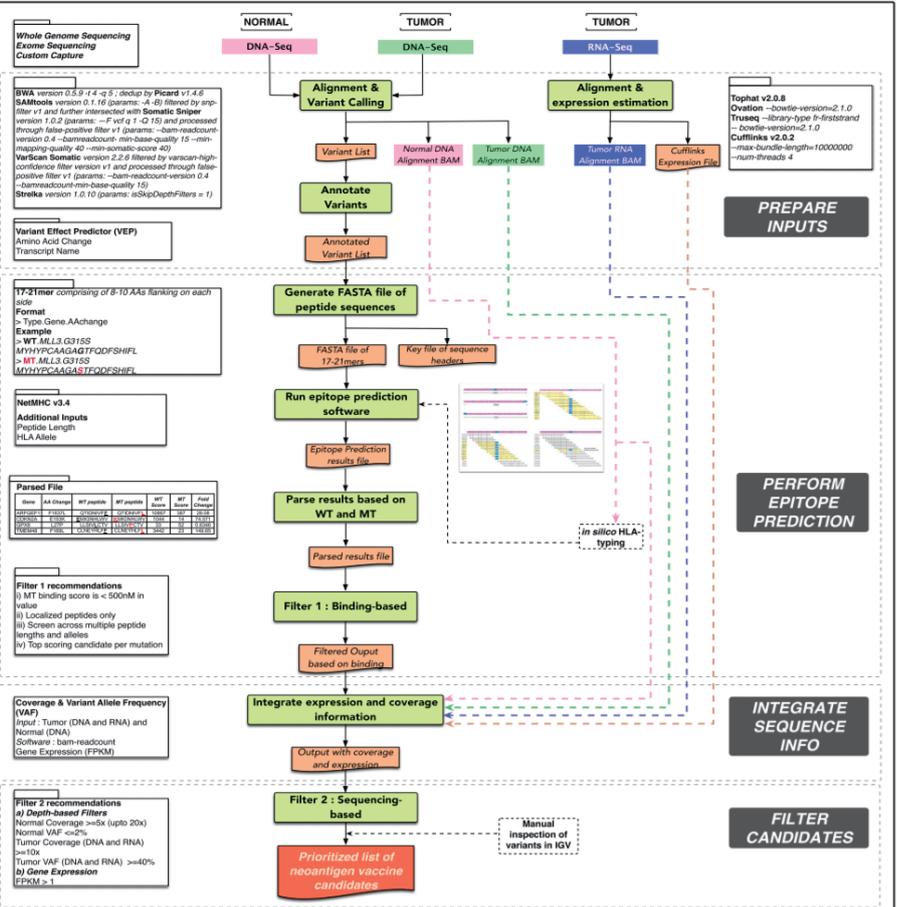
METHOD Open A

### pVAC-Seq: A genome-guided *in silico* approach to identifying tumor neoantigens

Jonen Hundal<sup>1</sup>, Beatriz M. Carrero<sup>2</sup>, Allegra A. Petri<sup>1</sup>, Gerald P. Linette<sup>2</sup>, Oki L. Griffin<sup>1,2,3,4</sup>, Elaine R. Mardis<sup>1,5,6,7</sup> and Malachi Griffith<sup>1,4,5</sup>

#### Abstract

Cancer immunotherapy has gained significant momentum from recent clinical successes of checkpoint blockade inhibition. Massively parallel sequence analysis suggests a connection between mutational load and response to this class of therapy. Methods to identify which tumor-specific mutant peptides (neoantigens) can elicit anti-tumor T cell immunity are needed to improve predictions of checkpoint therapy response and to identify targets for vaccines and adoptive T cell therapies. Here, we present a flexible, streamlined computational workflow for identification of personalized Variant Antigens by Cancer Sequencing (pVAC-Seq) that integrates tumor mutant and expression data (DNA- and RNA-Seq). pVAC-Seq is available at <https://github.com/griffithlab/pVAC-Seq>.



# Things need to be resolved for practical application

## Genome-level application

- Bulk/batched prediction of genome-level antigens
- Should be able to process all steps from NGS sequencing to final call
- Automated report with rich annotation and candidate suggestion

## Use of more information

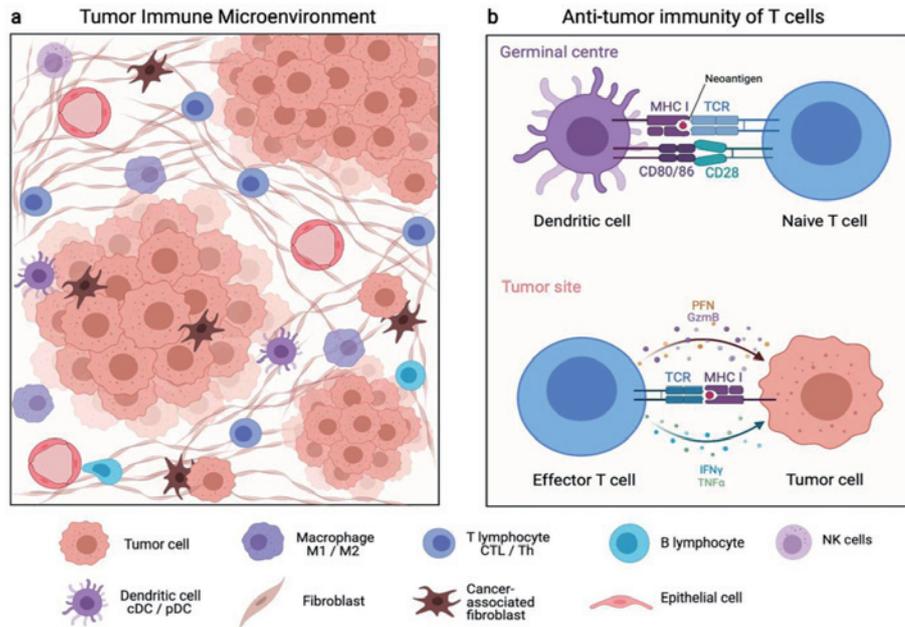
- Is MHC-I binding affinity the only applicable feature?
- Is  $IC_{50}$  under 50nM (or 500nM) an acceptable cut-off?

## Discovery of new features

- Can we find a new feature for immunogenicity prediction?

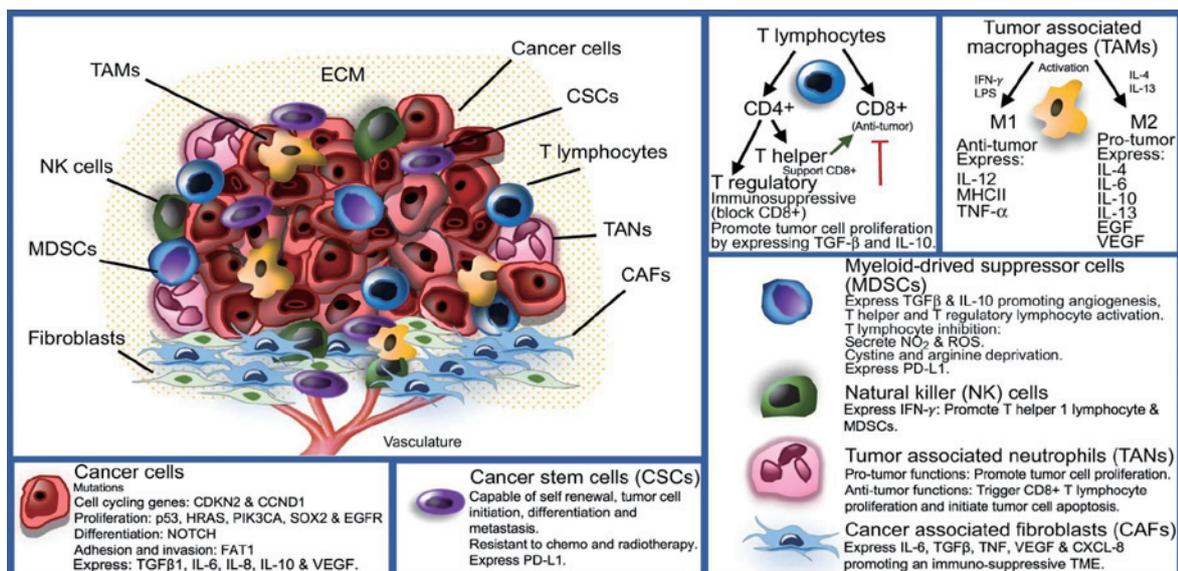
# IDENTIFYING TUMOR IMMUNE MICROENVIRONMENT

# Tumor Immune Microenvironment



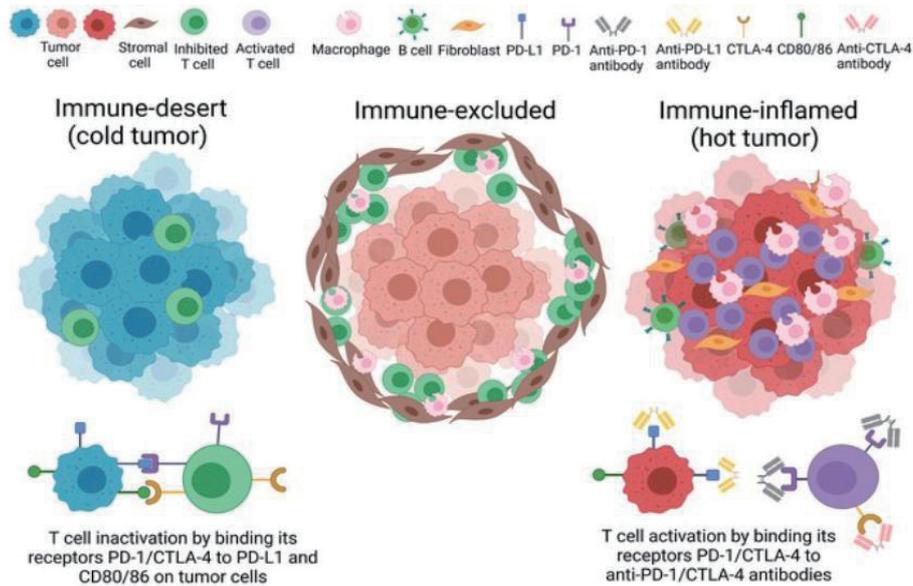
- A complex, organ-like structure (tumor cells, immune cells, fibroblasts, vascular endothelial cells, and other stromal cells)
- Immune cells + secreted factors (cytokines, chemokines, growth factors)

# Tumor Immune Microenvironment



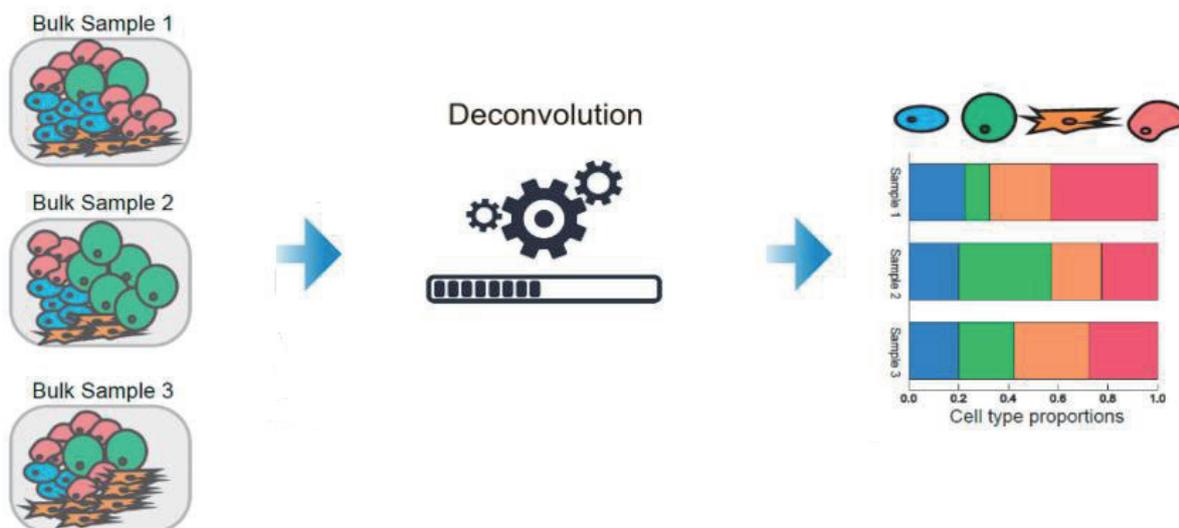
- TME components often inhibit or promote anti-tumor immunity
- But their roles are not definitive, and can be context-specific

# Tumor Immune Phenotype

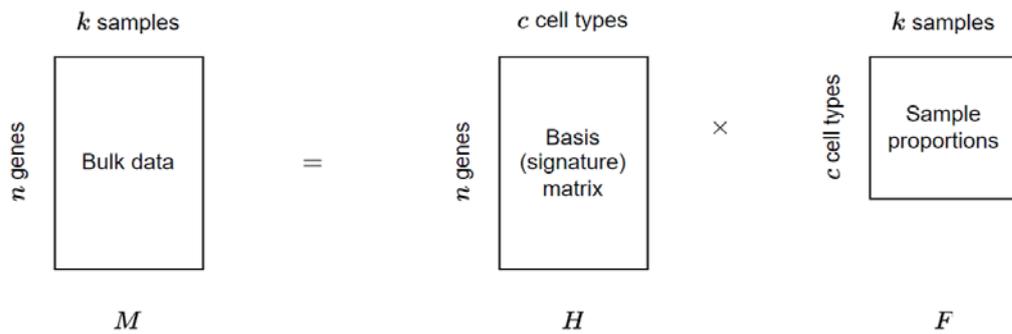


- Immune inflamed: immune cells infiltrated the tumor
- Immune excluded: immune cells are restricted to the stroma
- immune desert: T cells are not recruited

# Identifying TME by Cell-type decomposition

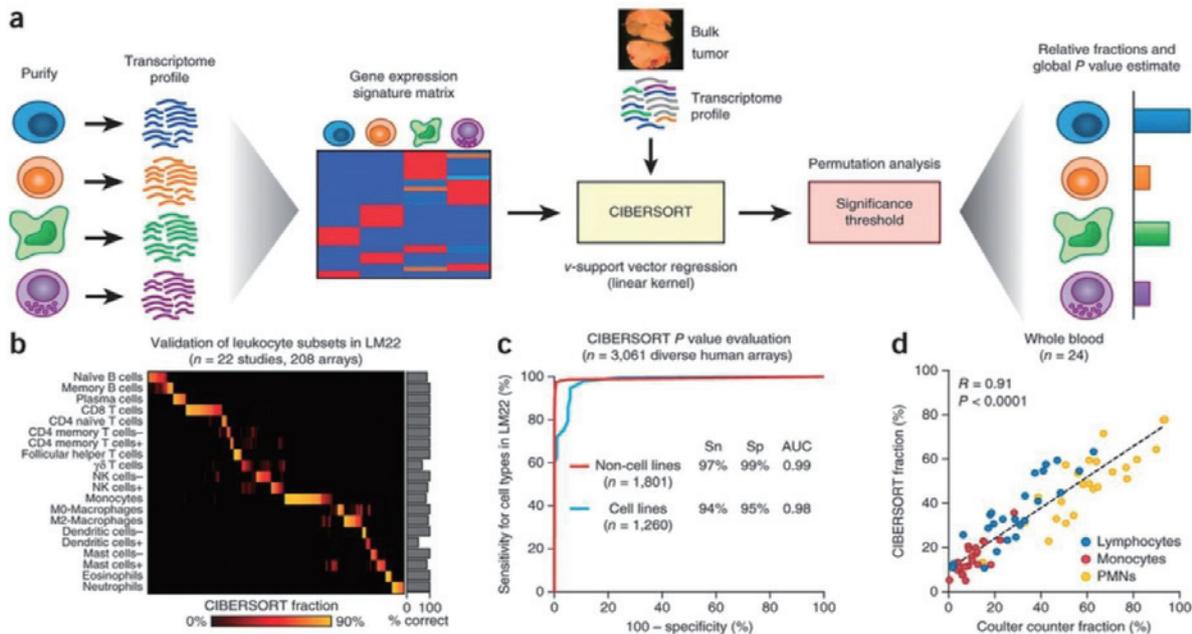


# Identifying TME by Cell-type decomposition



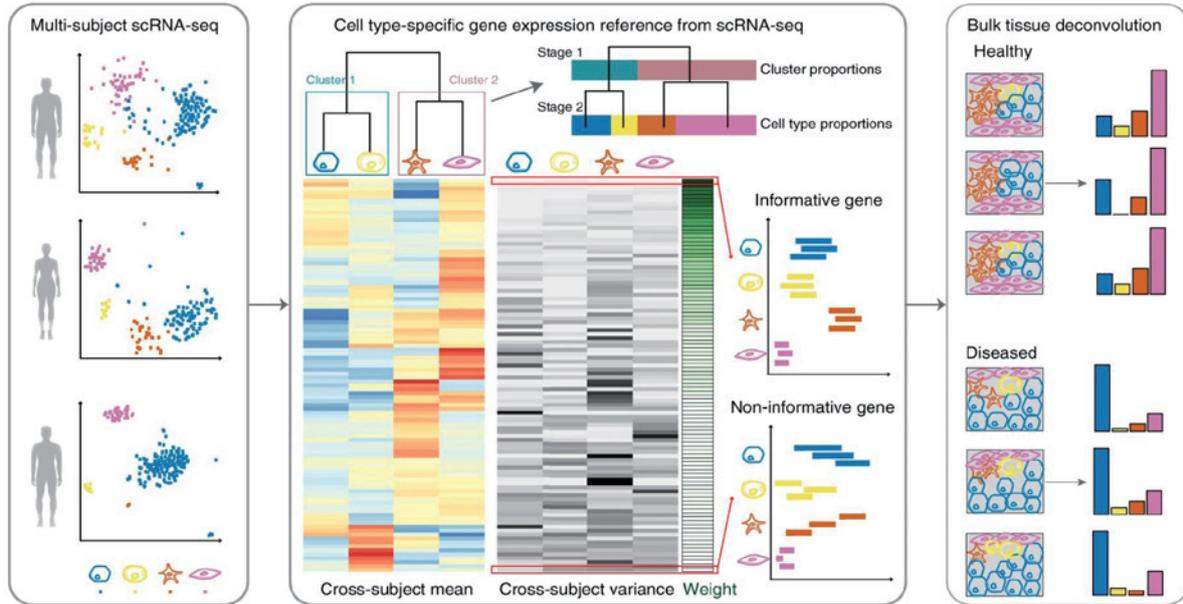
Problem	Given	Estimate	Requires
Estimate cell type proportions from bulk profile and signature matrix	$M, H$	$F$	$n > c$
Generate signature matrix from bulk profile and known cell type proportions	$M, F$	$H$	$k > c$
Estimate bulk profile from signature matrix and cell type proportions	$H, F$	$M$	none

## CIBERSORT



- Given a validated leukocyte gene signature matrix (LM22), deconvolute a  $n$  input bulk gene expression profile to generate cell-type fractions
- Support vector regression

# MuSiC



- Utilize scRNA-seq from multiple subjects, identifying reference cell type-specific gene expression
- Extract genes that are informative (low cross-subject variance)

# ImmuneDeconvR

	CIBERSORT	EPIC	MuSiC	DSA	TIMER	DeconvSeq	DCQ	NNS	dangle	Xcell	LinSeed	MCP-counter
Sutton, G. J. et al., 2022	Recommended	Not Recommended	Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended
Nadel, B.B. et al., 2021	Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended
Jin, H. et al., 2021	Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended
Cobos, S. et al., 2021	Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended
Sturm, G. et al., 2019	Not Recommended											

■ Recommended   
 ■ Not Recommended   
  Not Evaluated

Home > Bioinformatics for Cancer Immunotherapy > Protocol

## ImmuneDeconv: An R Package for Unified Access to Computational Methods for Estimating Immune Cell Fractions from Bulk RNA-Sequencing Data

Gregor Sturm, Francesca Finotello & Markus List

Protocol | First Online: 03 March 2020

5035 Accesses | 88 Citations | 1 Altmetric

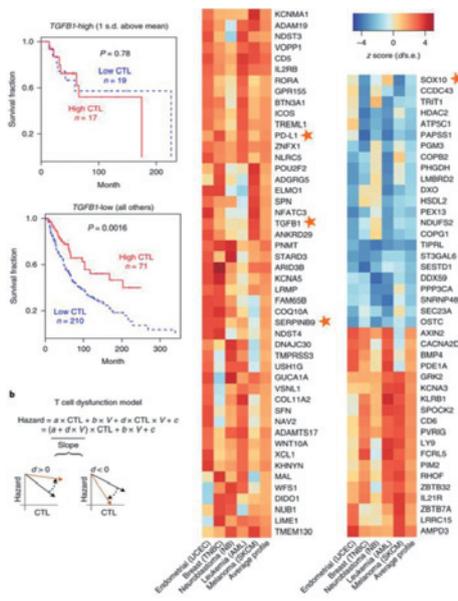
Part of the *Methods in Molecular Biology* book series (MIMB, volume 2120)

### Abstract

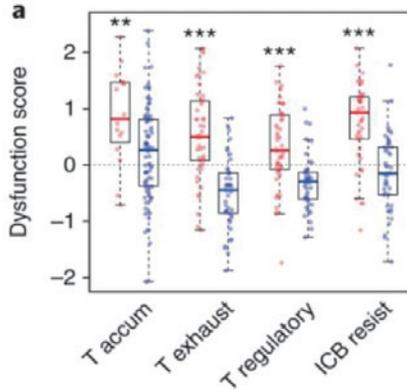
Since the performance of in silico approaches for estimating immune-cell fractions from bulk RNA-seq data can vary, it is often advisable to compare results of several methods. Given numerous dependencies and differences in input and output format of the various computational methods, comparative analyses can become quite complex. This motivated us to develop *immuneDeconv*, an R package providing uniform and user-friendly access to seven state-of-the-art computational methods for deconvolution of cell-type fractions from bulk RNA-seq data. Here, we show how *immuneDeconv* can be installed and applied to a typical dataset. First, we give an example for obtaining cell-type fractions using *quanTiseq*. Second, we show how dimensionless scores produced by MCP-counter can be used for cross-sample comparisons. For each of these examples, we provide R code illustrating how *immuneDeconv* results can be summarized graphically.

- Each tool has its own pros and cons, and do not agree each other
- ImmuneDeconv provides a unified access to immune decomposition tools, so users can see different results and finally find a consensus

# Predicting of T-cell evasion mechanisms (TIDE)

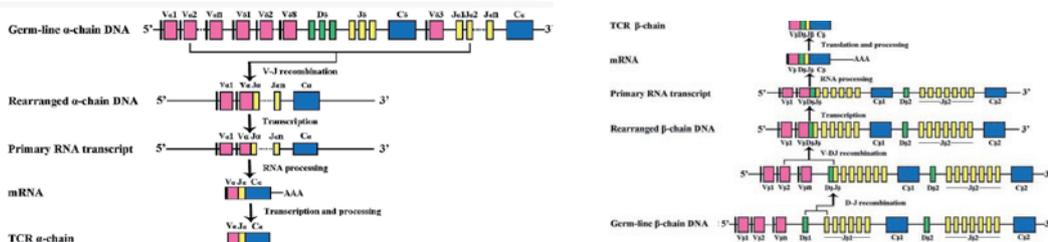
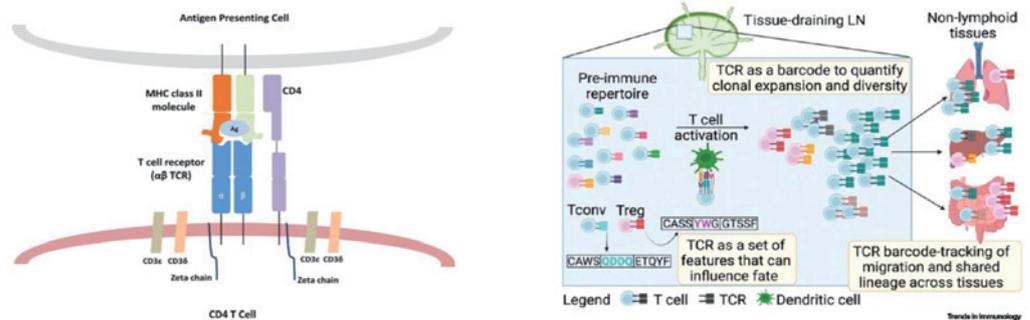


TGFB is interacting with CTL infiltration because:  
 - Survival of High vs. Low CTL-infiltrated patients are discriminated only when TGFB is highly expressed



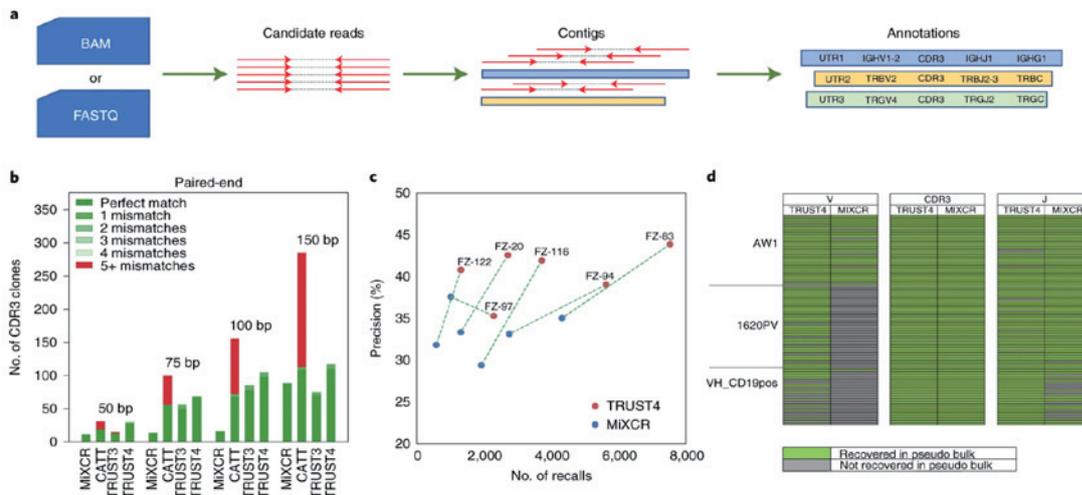
- Predicting T-cell dysfunction model: high infiltration but dysfunctional T-cells or excluded T-cells
- Extract T-cell dysfunction genes from interaction test in treatment naïve data
- Calculate T-cell dysfunction score

# TCR repertoire



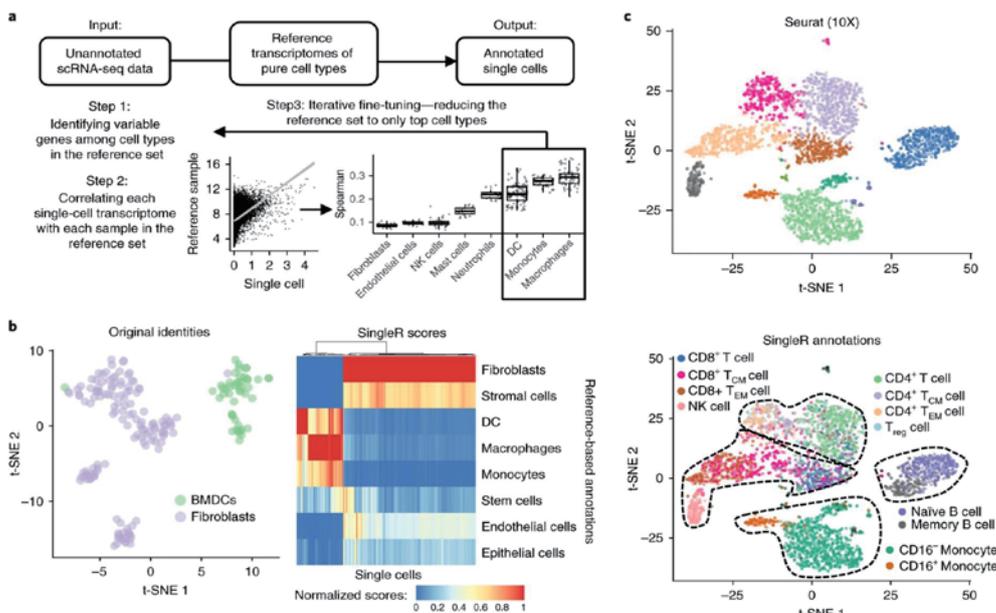
- T-cell diversity and clonality is the overall resultant response to the complex T-cell immune environment
- Diversity is inversely related to clonality
- High clonality is generally a marker for good response

# TCR repertoire reconstruction



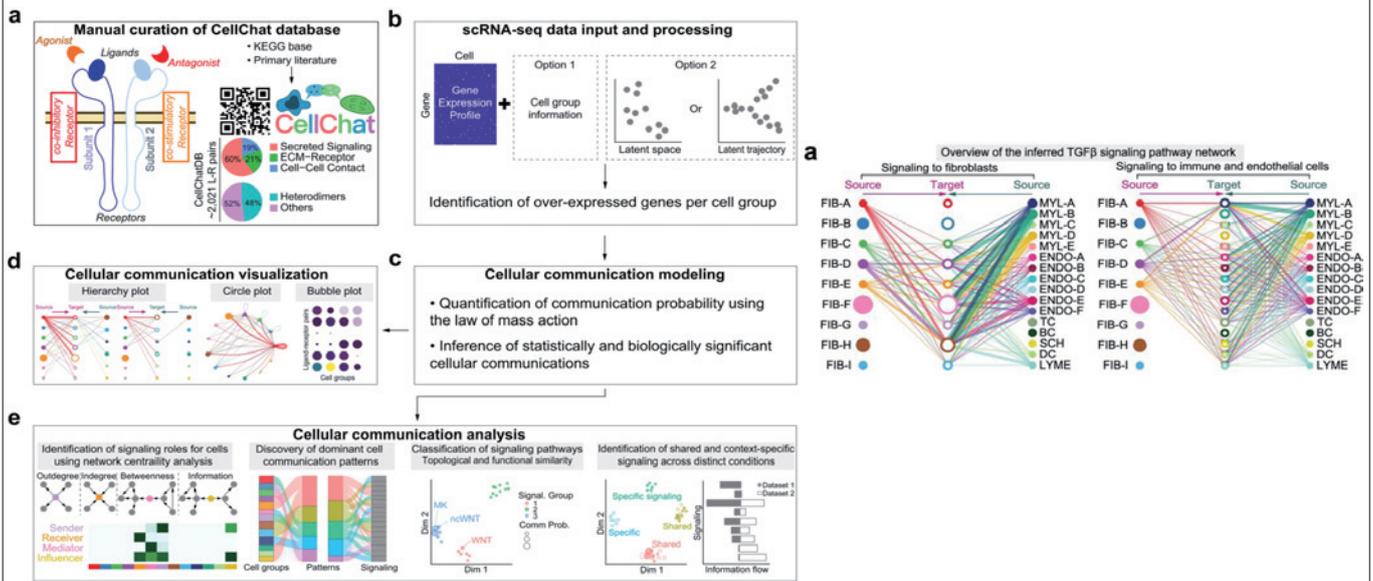
- Generally, TCR or BCR sequencing is employed for repertoire reconstruction
- Conventional bulk RNA-seq can be also used using specialized tools, such as TRUST
- TRUST 1) extract TCR/BCR candidate reads, 2) assemble to form contigs, 3) identify somatic hypermutations, 4) reconstruct repertoire

# Use of single- and spatial transcriptomics



- In single cell sequencing, a complex decomposition is not necessary once the single cells are well clustered.
- Clusters should be annotated using reference gene expression

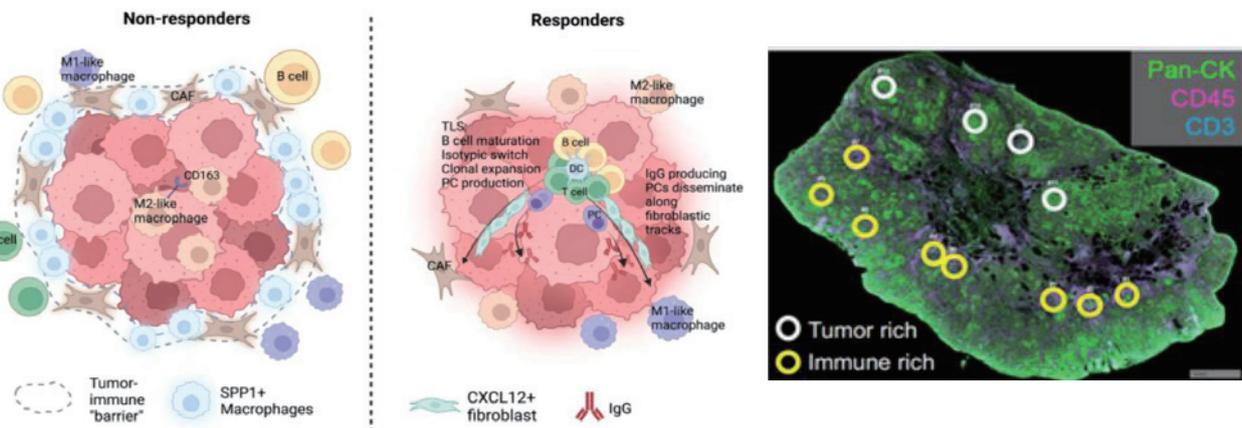
# Use of single- and spatial transcriptomics



- Cell type-level gene expression with predefined ligand-receptor interactions, cellular communications can be inferred, wherein which cell type influenced others through effector molecules

# Use of single- and spatial transcriptomics

## Response to Immunotherapy



- Using spatial transcriptomic, we can profile gene expression at the selected region of interest (ROI).
- Not only the abundance, but also the localization of immune cells direct the tumor immune microenvironment
- Similar bulk cell sequencing analysis techniques can be also applied to the spatial transcriptomics data

# Conclusion

- 다양한 cancer immunotherapy 의 발전으로 자신의 면역 시스템을 이용한 치료가 각광받고 있음
- 더 큰 효과와 적은 부작용을 위하여 환자, 종양 특이적 antigen 발굴이 필요함
- HLA type, MHC binding, Antigen processing 등 다양한 step 단계를 예측할 수 있는 computational algorithm 이 존재하며, 발전하고 있음
- Bulk, single, spatial transcriptomics 를 이용하여, 종양 주변의 면역환경인 Tumor immune microenvironment를 알아내고, 종양의 면역치료에 대한 환경에 따라 최적의 치료를 할 수 있음
- 결과적으로, NGS 에 기반하여 면역항암치료의 반응을 예측하고, 환자 특이적 치료를 할 수 있는 분석을 진행할 수 있음

# Thank you

Your success is our success. We've prescription for your business.  
We are professional communication group.



# KSBi-BIML 2024

## Introduction to cancer-immune analysis

### 실습용 도구 및 환경 안내



#### CLI (Command-Line Interface)



#### 특징

- 운영체제
- 무료 오픈소스
- 높은 통용성
- 높은 안정성
- 서버 환경으로 자주 사용됨

#### 특징

- 명령어 입력 방식 (아이콘 사용 X)
- 보다 가벼움
- 보다 안정적
- 자동화 용이

## 실습용 도구 및 환경 안내

R



### 특징

- 생물정보학 분석 필수 프로그램
- 무료 오픈소스

R studio



### 특징

- R 사용 보조
- 변수 관리, 명령어 입력 및 기록, figure 생성 등을 위한 통합 환경 제공
- 무료 오픈소스로 사용 가능

Bioconductor



### 특징

- 생물정보학 분석용 패키지 모음
- 무료 오픈소스 프로그램 사용 보조

## 실습 진행 순서

1. DNA-seq을 이용한 neoantigen prediction
2. Bulk RNA-seq을 이용한 tumor immune microenvironment 분석
3. Single cell RNA-seq을 이용한 cell-to-cell interaction prediction
4. Spatial RNA-seq을 이용한 TME 분석

## 실습용 데이터 안내

### DNA-seq을 이용한 neoantigen prediction

#### Prerequisites

##### Raw bam file (GRCh38)

- ACC\_T\_01.recaled.bam
- ACC\_T\_01.recaled.bai

##### Processed vcf file – Mutect2

- ACC\_T\_01.PASS.somatic.vcf

#### Processed data

##### Processed fastq

- ACC\_T\_01.chr6\_1.fastq
- ACC\_T\_01.chr6\_2.fastq

##### HLA typing (MHC class I) – OptiType

- ACC\_T\_01.MHC.I.processed.tsv
- ACC\_T\_01.MHC.I.list.txt

##### HLA typing (MHC class II) – HLA-HD

- ACC\_T\_01.MHC.II.processed.tsv
- ACC\_T\_01.MHC.II.list.txt

##### pVACseq (NetMHCpan, NetMHCIIpan)

- ACC\_T\_01.filtered.tsv (MHC Class I)
- ACC\_T\_01.filtered.tsv (MHC Class II)

실습 데이터: /home/jyhong906/BIML\_2024/Bulk\_WES/Data

실습 스크립트: /home/jyhong906/BIML\_2024/Bulk\_WES/Script

## 환경 변수 설정

```
#!/usr/bin/env bash # shebang
#$ -cwd # 현재 디렉토리 내 실행

# PATH #
HLA_PATH=/home/jyhong906/BIML_2024/Bulk_WES/Data # Input data, 결과 저장 디렉토리
optitype_PATH=${HLA_PATH}/OptiType # MHC class I 관련 HLA typing 결과 저장 디렉토리
hlahd_PATH=${HLA_PATH}/HLA-HD # MHC class II 관련 HLA typing 결과 저장 디렉토리

# MAKE FOLDER #
Path_list=(${HLA_PATH} ${optitype_PATH} ${hlahd_PATH})
for path in ${path_list[@]}; do
    mkdir -p $path # 상위 디렉토리 모두 생성
Done

# FILE #
Ref=/home/jyhong906/Project/Reference/Ref/hg38/genome.fa # Reference genome
IEDB_MHCI=/opt/Yonsei/IEDB-MHC_I # 사전 설치 필요
IEDB_MHCII=/opt/Yonsei/IEDB-MHC_II # 사전 설치 필요
hlahd_freq=/opt/Yonsei/HLA-HD/hlahd.1.7.0/freq_data
hlahd_split=/opt/Yonsei/HLA-HD/hlahd.1.7.0/HLA_gene.split.3.50.0.txt
hlahd_dict=/opt/Yonsei/HLA-HD/hlahd.1.7.0/dictionary

# EXECUTE #
vep_run=/opt/Yonsei/ensembl-vep/104.3/vep
optitype_run=/opt/Yonsei/OptiType/1.3.4/OptiTypePipeline.py
hlahd_run=hlahd.sh
pvacseq_run=/opt/Yonsei/python/3.8.1/bin/pvacseq

# SAMPLE #
patient_id=ACC_T_01

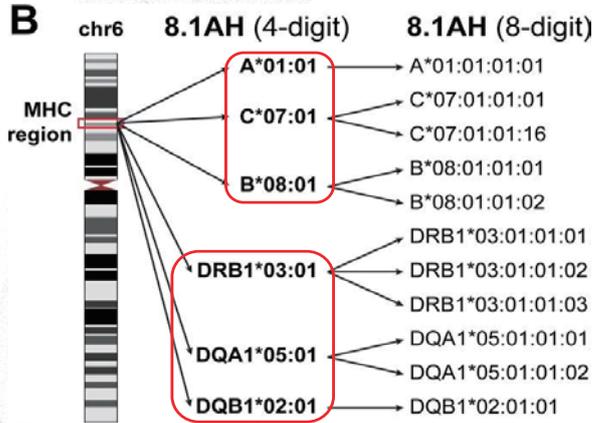
# FORMAT #
bam_format=.recaled.bam
chr6_bam_format=.sorted.chr6.bam
chr6_fastq1_format=.chr6_1.fastq
chr6_fastq2_format=.chr6_2.fastq
vcf_format=.PASS.somatic.vcf
ann_format=.vep.PASS.somatic.vcf
```

IEDB I, II installation  
<https://pvactools.readthedocs.io/en/latest/install.html#iedb-install>



## BAM to chr6 fastq

```
samtools view -h -b ${HLA_PATH}/${patient_id}${bam_format} chr6 > ${HLA_PATH}/${patient_id}${chr6_bam_format}
samtools fastq -1 ${HLA_PATH}/${patient_id}${chr6_fastq1_format} -2 ${HLA_PATH}/${patient_id}${chr6_fastq2_format} -F 4 ${HLA_PATH}/${patient_id}${chr6_bam_format}
```



Systematic genetic analysis of the MHC region reveals mechanistic underpinnings of HLA type associations with disease.

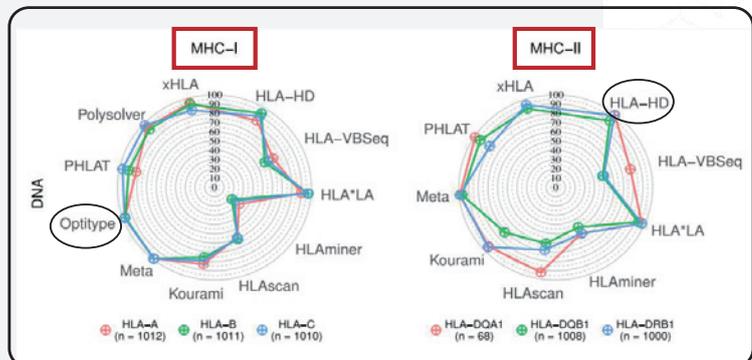
## HLA typing (MHC class I, II)

```
# HLA typing - MHC class I (OptiType) #
python2 ${optiype_run} \
-i ${HLA_PATH}/${patient_id}${chr6_fastq1_format} ${HLA_PATH}/${patient_id}${chr6_fastq2_format} \
-e 4 \
--dna \
-v \
-c /opt/Yonsei/OptiType/1.3.4/config.ini \
-o ${optiype_PATH}/${patient_id} \
--prefix ${patient_id}

# convert format #
python3 source_make_MHC_list.py MHC_I ${optiype_PATH} ${patient_id}

# HLA typing - MHC class II (HLA-HD) #
${hlahd_run} \
-t 10 \
-m 50 \
-f ${hlahd_freq} \
${HLA_PATH}/${patient_id}${chr6_fastq1_format} \
${HLA_PATH}/${patient_id}${chr6_fastq2_format} \
${hlahd_split} \
${hlahd_dict} \
${patient_id} \
${hlahd_PATH}

# convert format #
python3 source_make_MHC_list.py MHC_II ${hlahd_PATH} ${patient_id}
```



Benchmark of tools for in silico prediction of MHC class I and class II genotypes from NGS data

## HLA typing (MHC class I, II)

```
# HLA typing - MHC class I (OptiType) #
python2 ${optitype_run} \
-i ${HLA_PATH}/${patient_id}${chr6_fastq1_format} ${HLA_PATH}/${patient_id}${chr6_fastq2_format} \
-e 4 \
--dna \
-v \
-c /opt/Yonsei/OptiType/1.3.4/config.ini \
-o ${optitype_PATH}/${patient_id} \
--prefix ${patient_id}

# convert format #
python3 source_make_MHC_list.py MHC_I ${optitype_PATH} ${patient_id}
```

```
jjyhong906@master ACC_T_01$ pwd
/home/jyhong906/BIML_2024/Bulk_WES/Data/OptiType/ACC_T_01
jjyhong906@master ACC_T_01$ ll
total 1132
-rw-r--r-- 1 jjyhong906 jjyhong906 1154428 Feb  8 17:26 ACC_T_01.coverage_plot.pdf
-rw-r--r-- 1 jjyhong906 jjyhong906 337 Feb  8 17:26 ACC_T_01_result.tsv
```

HLA\_PATH=/home/jyhong906/BIML\_2024/Bulk\_WES/Data

MHC class I – ACC\_T\_01\_result.tsv

	A1	A2	B1	B2	C1	C2	Reads	Objective
0	A*02:01	A*34:01	B*40:02	B*15:02	C*15:02	C*08:01	2893.0	2762.8149999999955
1	A*02:01	A*34:01	B*15:02	B*40:06	C*15:02	C*08:01	2871.0	2741.8049999999953
2	A*02:01	A*34:05	B*40:02	B*15:02	C*15:02	C*08:01	2863.0	2734.1549999999955
3	A*34:01	A*02:16	B*40:02	B*15:02	C*15:02	C*08:01	2861.0	2732.244999999996

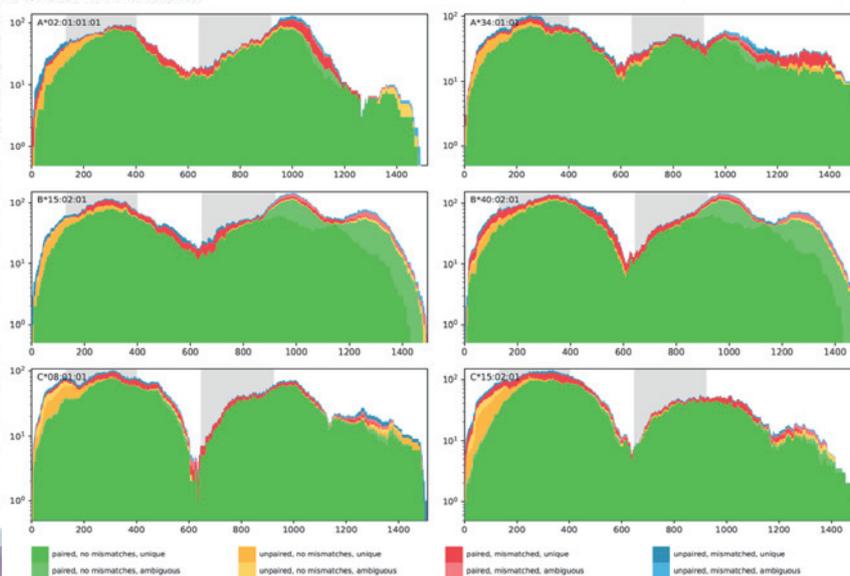
MHC class I – ACC\_T\_01.MHC.I.list.txt

```
HLA-A*02:01,HLA-A*34:01,HLA-B*40:02,HLA-B*15:02,HLA-C*15:02,HLA-C*08:01
```

## HLA typing (MHC class I, II)

```
# HLA typing - MHC class I (OptiType) #
python2 ${optitype_run} \
-i ${HLA_PATH}/${patient_id}${chr6_fastq1_format} ${HLA_PATH}/${patient_id}${chr6_fastq2_format} \
-e 4 \
--dna \
-v \
-c /opt/Yonsei/OptiType/1.3.4/config.ini \
-o ${optitype_PATH}/${patient_id} \
--prefix ${patient_id}

# convert format #
python3 source_make_MHC_list.py MHC_I ${optitype_PATH} ${patient_id}
```



# HLA typing (MHC class I, II)

```

[jyhong@master result]$ pwd
/home/jyhong996/HLA_2024/MS/Data/HLA-HD/ACC_T_01/result
[jyhong@master result]$ ll
total 1676
-rw-rw-r-- 1 jyhong996 jyhong996 4711 Feb 9 17:42 ACC_T_01.A.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 117230 Feb 9 17:42 ACC_T_01.A.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 965 Feb 9 17:42 ACC_T_01.B.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 89293 Feb 9 17:42 ACC_T_01.B.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 1568 Feb 9 17:42 ACC_T_01.C.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 85495 Feb 9 17:42 ACC_T_01.C.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 205 Feb 9 17:42 ACC_T_01.D.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 59235 Feb 9 17:42 ACC_T_01.DA1.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 899 Feb 9 17:42 ACC_T_01.DM8.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 54181 Feb 9 17:42 ACC_T_01.DM8.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 376 Feb 9 17:42 ACC_T_01.DPA2.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 64703 Feb 9 17:43 ACC_T_01.DPA2.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 610 Feb 9 17:41 ACC_T_01.DPB1.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 49198 Feb 9 17:41 ACC_T_01.DPB1.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 1949 Feb 9 17:42 ACC_T_01.DPA1.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 42920 Feb 9 17:42 ACC_T_01.DPA1.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 240 Feb 9 17:43 ACC_T_01.DPA2.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 3228 Feb 9 17:43 ACC_T_01.DPA2.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 2524 Feb 9 17:41 ACC_T_01.DPB1.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 56660 Feb 9 17:41 ACC_T_01.DPA1.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 740 Feb 9 17:42 ACC_T_01.DPA1.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 75414 Feb 9 17:42 ACC_T_01.DPA1.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 1568 Feb 9 17:42 ACC_T_01.DPB1.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 47342 Feb 9 17:42 ACC_T_01.DPB1.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 697 Feb 9 17:42 ACC_T_01.DPA1.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 10944 Feb 9 17:42 ACC_T_01.DPB1.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 431 Feb 9 17:40 ACC_T_01.DPB1.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 33356 Feb 9 17:40 ACC_T_01.DPB1.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 102 Feb 9 17:43 ACC_T_01.DPB2.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 181 Feb 9 17:43 ACC_T_01.DPB3.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 12652 Feb 9 17:43 ACC_T_01.DPB3.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 14 Feb 9 17:43 ACC_T_01.DPB4.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 9 Feb 9 17:41 ACC_T_01.DPB4.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 185 Feb 9 17:43 ACC_T_01.DPB5.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 17808 Feb 9 17:43 ACC_T_01.DPB5.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 73 Feb 9 17:43 ACC_T_01.DPB5.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 2610 Feb 9 17:43 ACC_T_01.DPB6.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 14 Feb 9 17:43 ACC_T_01.DPB7.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 9 Feb 9 17:43 ACC_T_01.DPB7.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 14 Feb 9 17:41 ACC_T_01.DPB8.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 88 Feb 9 17:43 ACC_T_01.DPB8.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 7838 Feb 9 17:43 ACC_T_01.DPB9.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 1888 Feb 9 17:43 ACC_T_01.E.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 89461 Feb 9 17:42 ACC_T_01.F.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 494 Feb 9 17:42 ACC_T_01.F.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 950 Feb 9 17:44 ACC_T_01.G.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 213 Feb 9 17:43 ACC_T_01.H.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 56222 Feb 9 17:43 ACC_T_01.G.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 685 Feb 9 17:43 ACC_T_01.G.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 75325 Feb 9 17:43 ACC_T_01.G.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 108146 Feb 9 17:43 ACC_T_01.H.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 352 Feb 9 17:44 ACC_T_01.J.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 87780 Feb 9 17:44 ACC_T_01.J.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 92 Feb 9 17:42 ACC_T_01.K.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 31139 Feb 9 17:42 ACC_T_01.K.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 85 Feb 9 17:43 ACC_T_01.L.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 49231 Feb 9 17:43 ACC_T_01.L.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 232 Feb 9 17:44 ACC_T_01.T.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 11807 Feb 9 17:44 ACC_T_01.T.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 116 Feb 9 17:44 ACC_T_01.V.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 27658 Feb 9 17:44 ACC_T_01.V.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 193 Feb 9 17:42 ACC_T_01.W.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 5122 Feb 9 17:42 ACC_T_01.W.read.txt
-rw-rw-r-- 1 jyhong996 jyhong996 343 Feb 9 17:44 ACC_T_01.Y.est.txt
-rw-rw-r-- 1 jyhong996 jyhong996 71265 Feb 9 17:44 ACC_T_01.Y.read.txt

```

```

# HLA typing - MHC class II (HLA-HD) #
$(hlahd_run) \
-t 10 \
-m 50 \
-f $(hlahd_freq) \
$(HLA_PATH)/$(patient_id)/$(chr6_fastq1_format) \
$(HLA_PATH)/$(patient_id)/$(chr6_fastq2_format) \
$(hlahd_split) \
$(hlahd_dict) \
$(patient_id) \
$(hlahd_PATH)

# convert format #
python3 source_make_MHC_list.py MHC_II $(hlahd_PATH) $(patient_id)

```

MHC class II - ACC\_T\_01\_final.result.txt

```

A HLA-A*02:01:01 HLA-A*34:01:01
B HLA-B*40:02:01 HLA-B*15:02:01
C HLA-C*15:02:01 HLA-C*08:01:01
DRB1 HLA-DRB1*15:02:01 HLA-DRB1*12:02:01
DQA1 HLA-DQA1*01:02:01 HLA-DQA1*06:01:01
DOB1 HLA-DOB1*03:01:01 HLA-DOB1*05:02:01
DPA1 HLA-DPA1*01:03:01 HLA-DPA1*02:02:02
DPB1 HLA-DPB1*02:01:02 HLA-DPB1*01:01:01
DMA HLA-DMA*01:02:01 -
DMB HLA-DMB*01:01:01 -
DOA HLA-DOA*01:01:04 HLA-DOA*01:01:01
DOB HLA-DOB*01:01:01 -
DRA HLA-DRA*01:02:02 HLA-DRA*01:01:01
DRB2 HLA-DRB2*01:01:01 -
DRB3 HLA-DRB3*03:01:03 -
DRB4 Not typed Not typed
DRB5 HLA-DRB5*01:01:01 -
DRB6 HLA-DRB6*02:01:01 -
DRB7 Not typed Not typed
DRB8 Not typed Not typed
DRB9 HLA-DRB9*01:02:01 -
DPA2 HLA-DPA2*01:01:02 HLA-DPA2*02:01:01
E HLA-E*01:03:01 HLA-E*01:03:02
F HLA-F*01:01:01 -
G HLA-G*01:01:01 HLA-G*01:01:03
H HLA-H*01:01:01 HLA-H*02:27
I HLA-I*01:01:01 HLA-I*01:01:01
K HLA-K*01:02 -
L HLA-L*01:02 -
T HLA-T*01:01:01 HLA-T*02:01:01
V HLA-V*01:01:01 -
W HLA-W*03:01:01 -
Y HLA-Y*01:01 HLA-Y*03:01

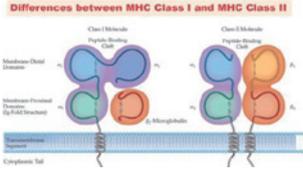
```

MHC class II - ACC\_T\_01.MHC.II.list.txt

```

DRB1*15:02, DRB1*12:02, DQA1*01:02, DOB1*03:01, DQA1*06:01, DQB1*05:02

```



# pVACseq (neoantigen prediction)

```

# pVACseq - MHC class I (NetMHCpan) #
allele='cat $(optype_PATH)/$(patient_id).MHC.II.list.txt'
$(pvacseq_run) run $(HLA_PATH)/$(patient_id)/$(ann_format) \
$(patient_id) \
$(allele) \
NetMHCpan \
-e1 8,9,10,11 \
--pass-only \
$(HLA_PATH)/$(patient_id) \
--edb-install-directory $(IEDB_MHCI)

# pVACseq - MHC class II (NetMHCIIpan) #
allele='cat $(hlahd_PATH)/$(patient_id).MHC.II.list.txt'
$(pvacseq_run) run $(HLA_PATH)/$(patient_id)/$(ann_format) \
$(patient_id) \
$(allele) \
NetMHCIIpan \
-e2 12,13,14,15,16,17,18 \
--pass-only \
$(HLA_PATH)/$(patient_id) \
--edb-install-directory $(IEDB_MHCI)

```

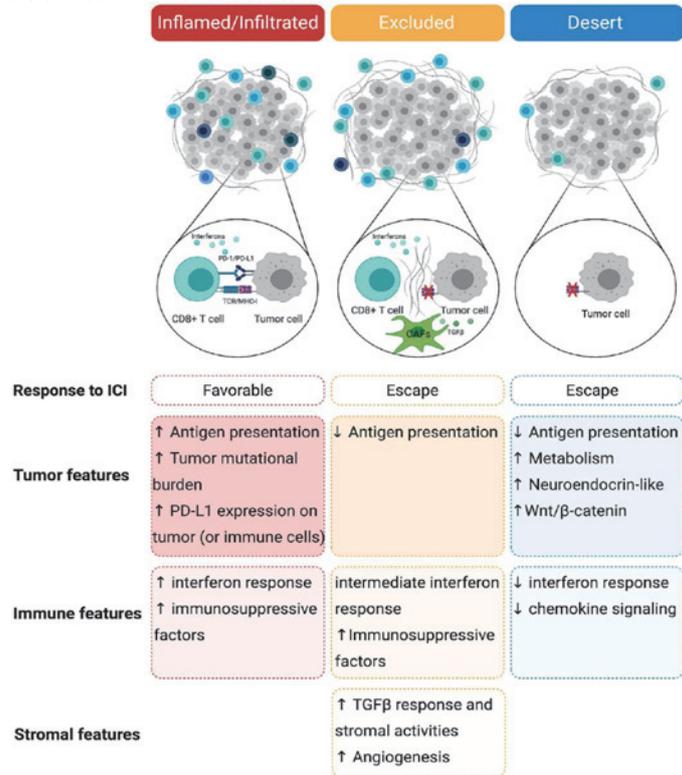
MHC class I - ACC\_T\_01.filtered.tsv

Gene Name	HGVSC	HGVSP	HLA Allele	Peptide Length	Sub-peptide Position	Mutation	MT Epitope Seq	MT Epitope Seq	Method	Best MT IC50 Score	Corr														
U378Lys	HLA-A*02:01	10	2	9	GLAGVKIAKV	GLAGVKIAEV	NetMHCpan	59.53	17.05	0.286	NetMHCpan	0.48	0.16	136	0.365	NA	NA	NA	NA	NA	NA				
chr9	128583017	C	3	6	ENST00000372739.7	1	2477	protein_coding	ENSG00000197694	missense	O/E	681	SPTAN1	ENST00000372739.7	c.2041C>G	ENSP00000369081.4	p.G1	0.57	0.6	96	0.321	NA	NA	NA	NA

MHC class II - ACC\_T\_01.filtered.tsv

Gene Name	HGVSC	HGVSP	HLA Allele	Peptide Length	Sub-peptide Position	Mutation	MT Epitope Seq	MT Epitope Seq	Method	Best MT IC50 Score	Corr										
u378Lys	HLA-DQA1*01:02/DRB1*03:01	16	3	14	KKEFPLAGVKIAKVD	KKEFPLAGVKIAEVD	NetMHCIIpan	375.12	373.19	0.995	NetMHCIIpan	2.7	2.6	136	0.365	NA	NA	NA	NA	NA	NA

# Tumor immune microenvironments (TIME)



## 실습용 데이터 안내

### Bulk RNA-seq을 이용한 tumor immune microenvironment 분석

#### Prerequisites

Gene quantification file – HTseq 등

- ~.htseq.count.txt

Raw fastq 파일

- ACC\_T\_01\_1.fastq.gz
- ACC\_T\_02\_1.fastq.gz

#### Processed data

Normalized expression matrix

- normalized\_TPM.rds

Cell type decomposition

- abis.rds
- cibersort\_abs.rds
- consensus\_tme.rds
- epic.rds
- estimate.rds
- mcp\_counter.rds
- quantiseq.rds
- timer.rds
- xcell.rds

Immune cell repertoire

- TRUST4\_dat.rds

Tumor immune dysfunction and exclusion

- TIDE\_dat.rds



실습 데이터: /home/jyhong906/BIML\_2024/Bulk\_RNA/Data

실습 스크립트: /home/jyhong906/BIML\_2024/Bulk\_RNA/Script

# TPM (Transcripts Per Million) normalization

```
#####
# Load expression data & TPM normalization #
#####
SGC_dir <- "/data/project/BIML_2024/Bulk_RNA/Sample"
SGC_files <- list.files(SGC_dir)
SGC_path <- paste0(SGC_dir, "/", SGC_files)
SGC_names <- gsub(".htseq.count.txt", "", SGC_files)

tmp_df_list <- c()
for (idx in seq(SGC_names)) {
  tmp_df <- read.table(
    SGC_path[idx],
    header = F,
    sep = "\t",
    stringsAsFactors = F,
    col.names = c("Symbol", SGC_names[idx])
  )
  tmp_df_list[[idx]] <- tmp_df
}

SGC_count_df <- Reduce(merge, tmp_df_list)[-c(1:5),]
rownames(SGC_count_df) <- SGC_count_df$Symbol; SGC_count_df <- SGC_count_df[,-1]

# Gene filter #
SGC_mat <- as.matrix(SGC_count_df[rowSums(SGC_count_df) >= 1]) >= ncol(SGC_count_df,))

load("/data/project/BIML_2024/Bulk_RNA/ImmuneDeconv/gene_cov.rda")
normalized_TPM <- countToTpm(SGC_mat,
                             keyType = "SYMBOL",
                             gene_cov = gene_cov)
```



# Load library & normalization

```
#####
# Load packages #
#####
library(GeoIcgaData) # normalized IPM
library(ImmuneDeconv) # Cell type decomposition
library(ComplexHeatmap) # Visualization
library(ggplot2) # Visualization
library(ggpubl) # Visualization
library(gridExtra) # Visualization
library(ggpubr) # statistics
library(circlize) # color
library(metapod) # combined p-value
library(gtools) # processing
library(dplyr) # processing

#####
# Load expression data & TPM normalization #
#####
SGC_dir <- "/data/project/BIML_2024/Bulk_RNA/Data/ImmuneDeconv"
SGC_files <- list.files(SGC_dir,pattern = "*.txt")
SGC_path <- paste0(SGC_dir, "/", SGC_files)
SGC_names <- gsub(".htseq.count.txt", "", SGC_files)

tmp_df_list <- c()
for (idx in seq(SGC_names)) {
  tmp_df <- read.table(
    SGC_path[idx],
    header = F,
    sep = "\t",
    stringsAsFactors = F,
    col.names = c("Symbol", SGC_names[idx])
  )
  tmp_df_list[[idx]] <- tmp_df
}

SGC_count_df <- Reduce(merge, tmp_df_list)[-c(1:5),]
rownames(SGC_count_df) <- SGC_count_df$Symbol; SGC_count_df <- SGC_count_df[,-1]

# Gene filter #
SGC_mat <- as.matrix(SGC_count_df[rowSums(SGC_count_df) >= 1]) >= ncol(SGC_count_df,))

load("/data/project/BIML_2024/Bulk_RNA/Data/ImmuneDeconv/gene_cov.rda") # https://github.com/YuLab-SMU/GeoIcgaData
normalized_TPM <- countToTpm(SGC_mat,
                             keyType = "SYMBOL",
                             gene_cov = gene_cov)

# saveRDS(normalized_TPM, file = "/data/project/BIML_2024/Bulk_RNA/Data/ImmuneDeconv/normalized_TPM.rds")
# read_rds("/data/project/BIML_2024/Bulk_RNA/Data/ImmuneDeconv/normalized_TPM.rds")

# Primary & metastasis 정보 확인 #
SGC_groups <- read.table("/data/project/BIML_2024/Bulk_RNA/Script/SGC_groups.txt",
                        sep = "\t",
                        header = T)
P_idx <- SGC_groups$Condition == "PRIMARY"; M_idx <- SGC_groups$Condition == "METASTASIS"
```

Install.packages()





# TIDE (Tumor immune dysfunction and exclusion)

nature medicine

Explore content About the journal Publish with us

nature > nature medicine > articles > article

Article | Published: 20 August 2018

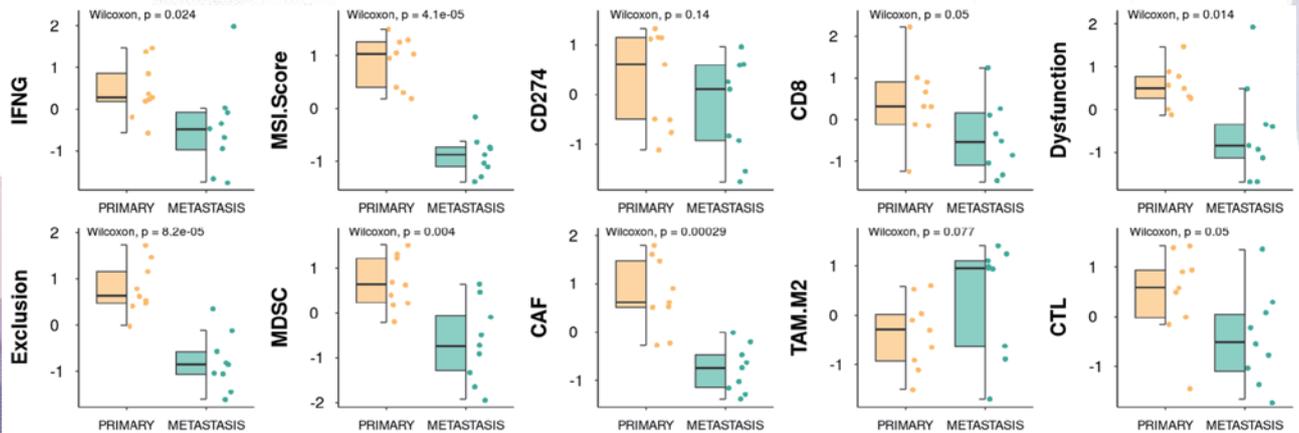
## Signatures of T cell dysfunction and exclusion predict cancer immunotherapy response

Peng Jiang, Shengqing Gu, Deng Pan, Jingxin Fu, Avinash Sahu, Xihao Hu, Ziyi Li, Nicole Traugh, Xia Bu, Bo Li, Jun Liu, Gordon J. Freeman, Myles A. Brown, Kai W. Wucherpfennig & X. Shirley Liu

Nature Medicine 24, 1550–1558 (2018) | Cite this article

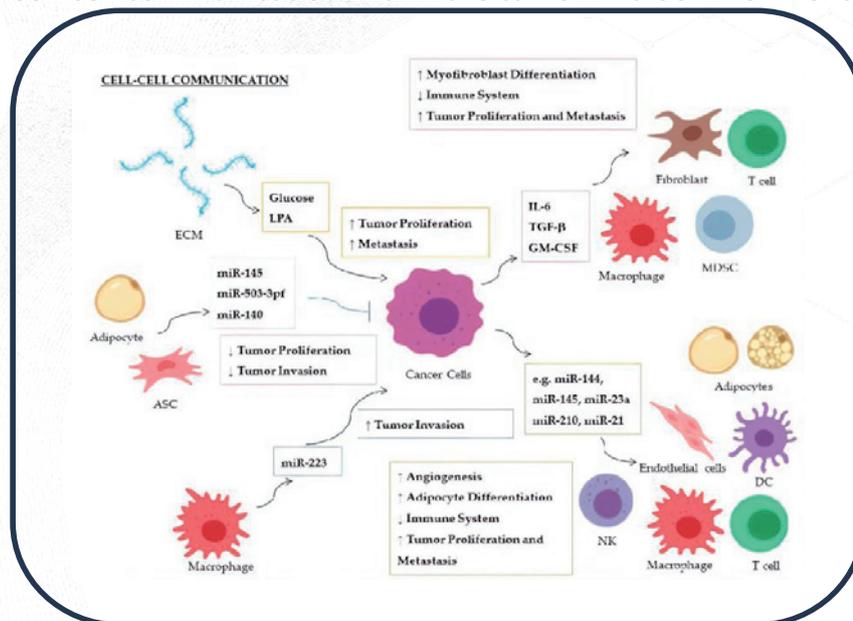
59k Accesses | 2269 Citations | 157 Altmetric | Metrics

	IFNG	MSI_Score	CD274	CD8_Dys_Function	Exclusion	MDSC	CAF	TAM_M2	CTL
ACC_T_16	2.7713852	0.38276256	0.5687156	2.212168	-1.93880887	-0.45261866	-0.051164829	0.002992495	-0.0169539769
ACC_T_03	0.50849971	0.88954761	1.2365711	-2.2484679	0.09121424	1.49101366	0.08643357	0.139496973	-0.0061515064
ACC_T_13	0.47460238	0.78839511	1.0657308	-2.062644	-1.48158199	0.71345098	-0.009193384	0.127708458	-0.0138163032
ACC_T_15	0.38777099	0.8877155	1.0322718	0.561730	0.90075486	1.25115264	0.02554704	0.1547399587	0.0125169288
ACC_T_08	1.93107510	0.95170156	1.0657308	1.8106446	0.68317765	0.57807223	0.039269776	0.047777594	0.0201084187
ACC_T_11	0.70437342	0.56803795	-0.7138054	-0.2614655	0.53027842	0.431916957	0.071556991	-0.01933786	-0.0037182193
ACC_T_04	0.34888156	0.58457518	-0.4547434	0.5755181	0.80198791	0.38928623	0.015838679	0.079258652	-0.0325641511
ACC_T_09	1.36521189	0.27838910	-0.8615188	0.510158	0.52122515	0.32742377	0.03151832	-0.016428675	0.0235731842
ACC_T_14	0.22552226	0.81374864	-0.4542434	1.6208195	0.25709455	0.409191761	0.041188668	0.0465380336	0.0241616021
ACC_T_12	1.28827211	0.87803523	0.5687156	4.0100306	0.05527366	1.01367462	0.077175812	0.055393958	0.0116444665
ACC_T_10	-0.27233313	0.61581230	-1.0279287	-0.2138869	0.32945239	0.03174874	0.013288879	-0.013384841	0.0002163238
ACC_T_06	0.44383125	0.18194626	-0.780789	-2.6097801	-1.68728152	-0.05389544	0.10180885	-0.03027798	0.031527166
ACC_T_17	2.45351281	0.22641188	1.6385932	-2.4862188	-1.61628434	-0.65583489	0.027577810	0.03402711	0.0267525209
ACC_T_02	2.36832913	0.15911562	-1.4312844	-1.9733693	-1.07253935	-0.68486541	-0.008513432	-0.114391218	0.028426059
ACC_T_18	0.48029992	0.88043171	0.578118	-1.4138691	-0.38258486	-0.81071026	-0.01929298	0.000101972	0.0218813988
ACC_T_07	0.59015469	0.26926171	0.1863724	-0.9822203	-0.78611899	-0.85664243	-0.031513173	-0.106685509	0.0302452309
ACC_T_01	-0.11512170	0.44708480	0.8999215	-0.5671828	-0.29827862	-1.18586179	-0.106106089	0.052049289	0.0137160843
ACC_T_05	0.41229290	0.31019376	0.2418296	0.2781981	-0.14084122	-1.11159958	-0.489148783	-0.081041825	-0.0192899162



## TME and cell to cell interaction

### Cell-cell communication within the tumor microenvironment



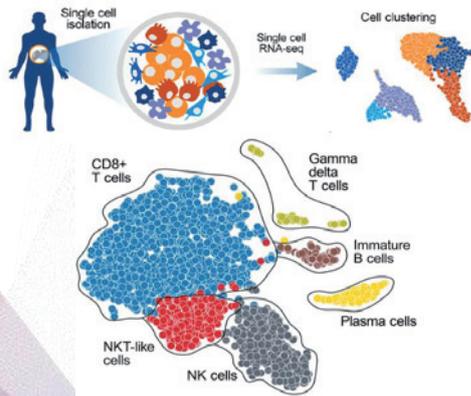
Conti I, Varano G, Simioni C, Laface I, Milani D, Rimondi E, Neri LM. miRNAs as Influencers of Cell-Cell Communication in Tumor Microenvironment. *Cells*. 2020 Jan 15;9(1):220. doi: 10.3390/cells9010220. PMID: 31952362; PMCID: PMC7016744.

# TME and cell to cell interaction

## Cell type annotation

### Single R

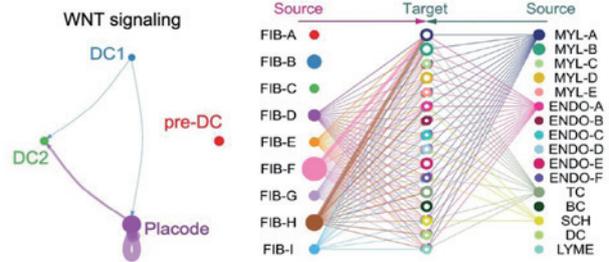
- computational method for unbiased cell type recognition of scRNA-seq
- SingleR's annotations combined with Seurat, a processing and analysis package designed for scRNA-seq



## Cell to cell interaction

### CellChat

- Infer cell-cell communication networks
- easy-to-use tool for extracting and visualizing



lanevski, A., Giri, A.K. & Aittokallio, T. Fully-automated and ultra-fast cell-type identification using specific marker combinations from single-cell transcriptomic data. *Nat Commun* **13**, 1246 (2022).

# 실습용 데이터 안내

## Single cell RNA-seq을 이용한 cell to cell interaction prediction

### Prerequisites

#### Raw single cell data

- Barcodes.tsv
- Features.tsv
- matrix.mtx

#### Human single cell reference

- monaco.ref.rda
- hpca.ref.rda
- dice.ref.rda

### Processed data

#### Seurat object

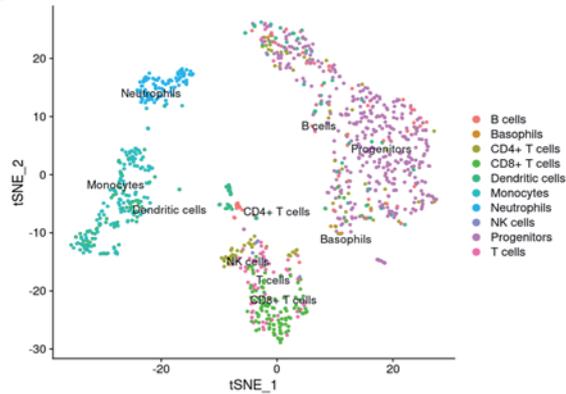
- CRC\_obj.rda
- CRC\_count.rda

## Cell type annotation

```

library('dplyr')
library('Seurat')
library('SingleR')
library('CellChat')
library('ADImpute')

# Cell type/state annotation #
load("/data/project/BIML_2024/scRNA/ref/monaco.ref.rda") # Reference single cell data. celldex::MonacoImmuneData() 로 다운 가능
load("/data/project/BIML_2024/scRNA/CRC_obj.rda") # Seurat object
load("/data/project/BIML_2024/scRNA/CRC_count.rda") # Single cell expression count file
monaco.main <- SingleR(method='single',sc_data=CRC_count, ref_data=monaco.ref@assays@data@listData$logcounts,types=monaco.ref$label.main)
CRC_obj@meta.data$monaco.main <- monaco.main$labels1
CRC_obj_monaco.main <- SetIdent(CRC_obj, value = "monaco.main")
DimPlot(CRC_obj_monaco.main, reduction = "tsne", label = TRUE, repel = TRUE, group.by = 'monaco.main')
    
```



## CellChat

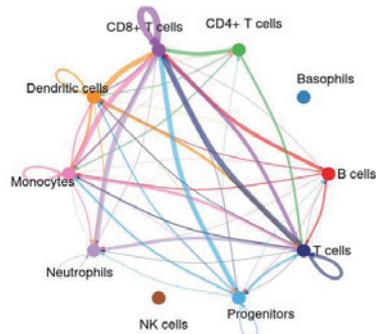
```

# CellChat object #
CellChatDB <- CellChatDB.human
cellchat <- createCellChat(object = CRC_obj_monaco.main, group.by = "monaco.main", assay = "RNA")
cellchat@DB <- CellChatDB
cellchat <- subsetData(cellchat)
cellchat <- identifyOverExpressedGenes(cellchat)
cellchat <- identifyOverExpressedInteractions(cellchat)
cellchat <- computeCommunProb(cellchat)
cellchat <- filterCommunication(cellchat, min.cells = 10)
cellchat <- computeCommunProbPathway(cellchat)
cellchat <- aggregateNet(cellchat)
cellchat <- netAnalysis_computeCentrality(cellchat, slot.name = "netP")
    
```

## Visualization

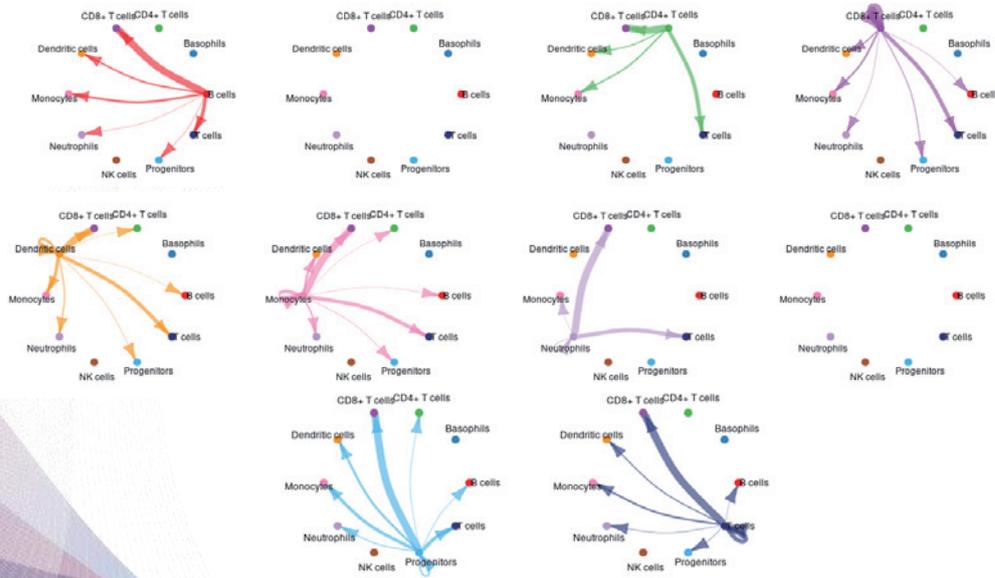
```

netVisual_circle(cellchat@net$weight, weight.scale = T, label.edge = F, title.name = "Interaction weights/strength") #전체 세포 상호작용
    
```



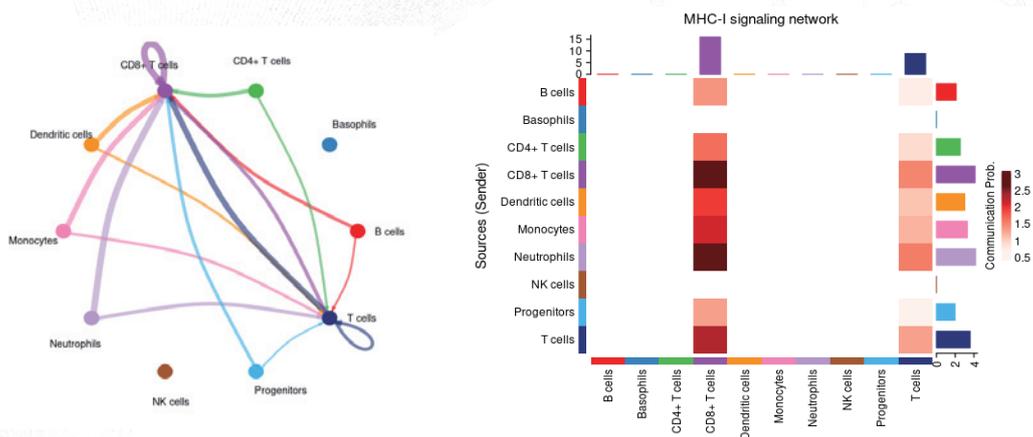
## Visualization

```
mat <- cellchat@net$weight
par(mfrow = c(3,4))
for (i in 1:nrow(mat)) {
  mat2 <- matrix(0, nrow = nrow(mat), ncol = ncol(mat), dimnames = dimnames(mat))
  mat2[i, ] <- mat[i, ]
  netVisual_circle(mat2, weight.scale = T, title.name = rownames(mat)[i])
}
dev.off()
```



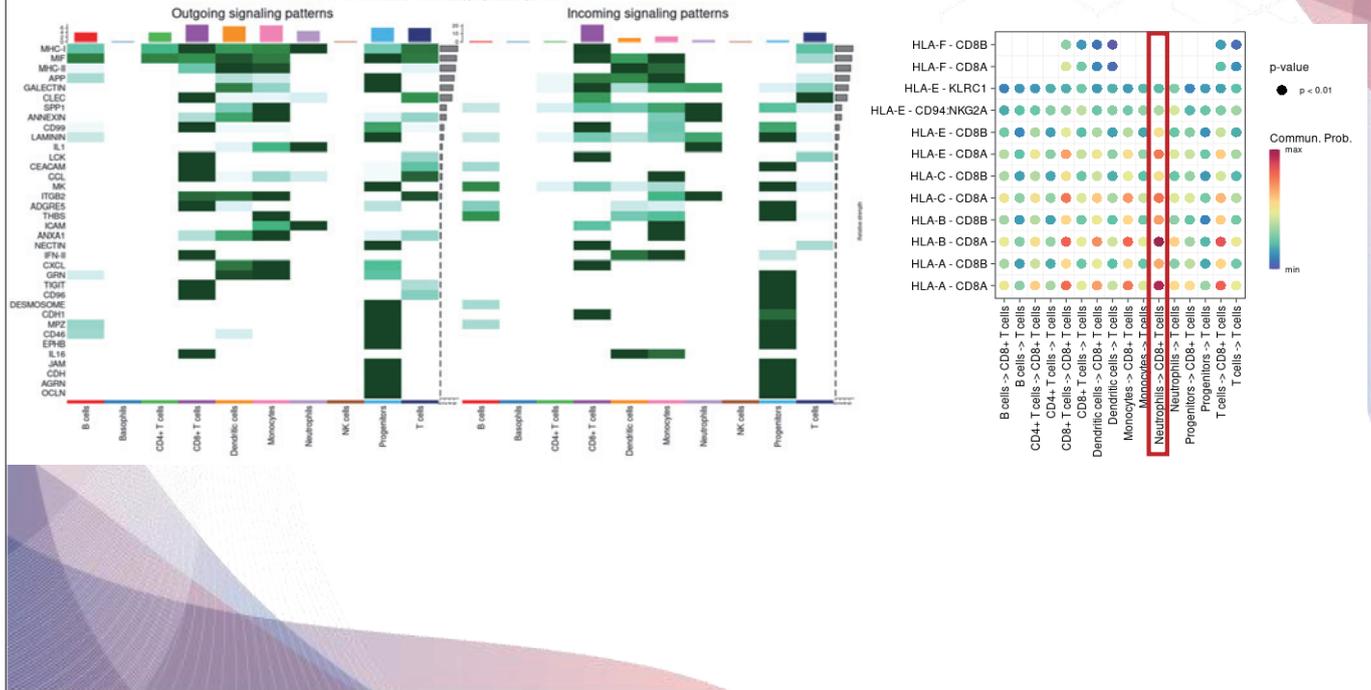
## Visualization

```
pathways.show <- c("MHC-I")
netVisual_aggregate(cellchat, signaling = pathways.show, layout = "circle")
netVisual_heatmap(cellchat, signaling = pathways.show, color.heatmap = "Reds") #특정 생물학적 경로 내 상호작용
```

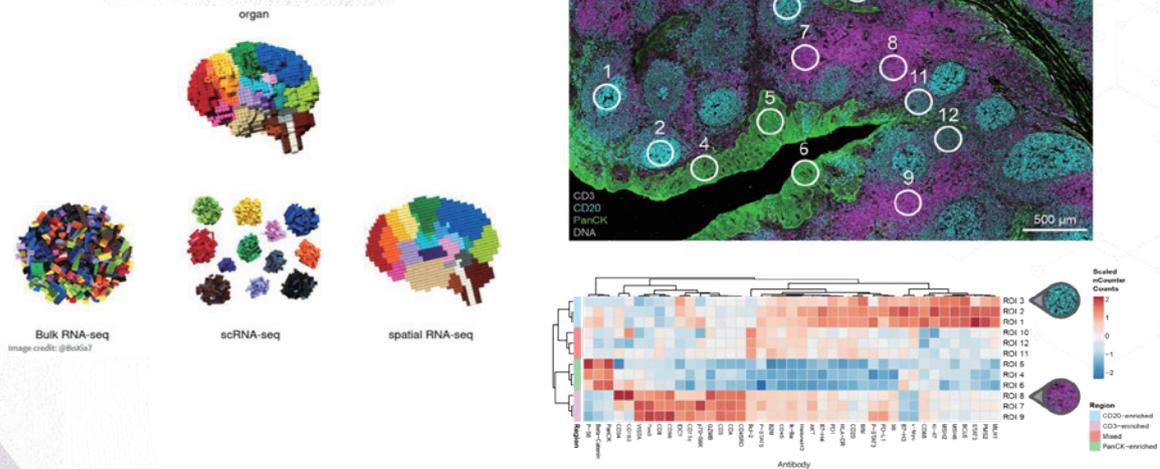


# Visualization

```
ht1 <- netAnalysis_signalingRole_heatmap(cellchat, pattern = "outgoing", width = 20, height = 20, font.size = 13, font.size.title = 20); ht1
ht2 <- netAnalysis_signalingRole_heatmap(cellchat, pattern = "incoming", width = 20, height = 20, font.size = 13, font.size.title = 20); ht2
ht1 + ht2
netVisual_bubble(cellchat, signaling=pathways.show, remove.isolate = FALSE) #Outgoing/Incoming signaling
```



# What is spatial transcriptomics?



# 실습용 데이터 안내

## spatial RNA-seq을 이용한 DEG, GSEA 분석

### Prerequisites

#### Processed GoeMX data

- count.rds
- anno.rds
- genemeta.txt
- msigdb\_hs.RData

실습 데이터: /home/jyhong906/BIML\_2024/GeoMX/Data

실습 스크립트: /home/jyhong906/BIML\_2024/GeoMX/Script

<https://cumulus.readthedocs.io/en/stable/geomxngs/index.html#convert-fastq-files-into-dcc-files-by-the-nanostring-geomx-digital-spatial-ngs-pipeline>

## Preparation

```
source("/data/project/BIML_2024/GeoMX/Script/GeoMX_function.R")

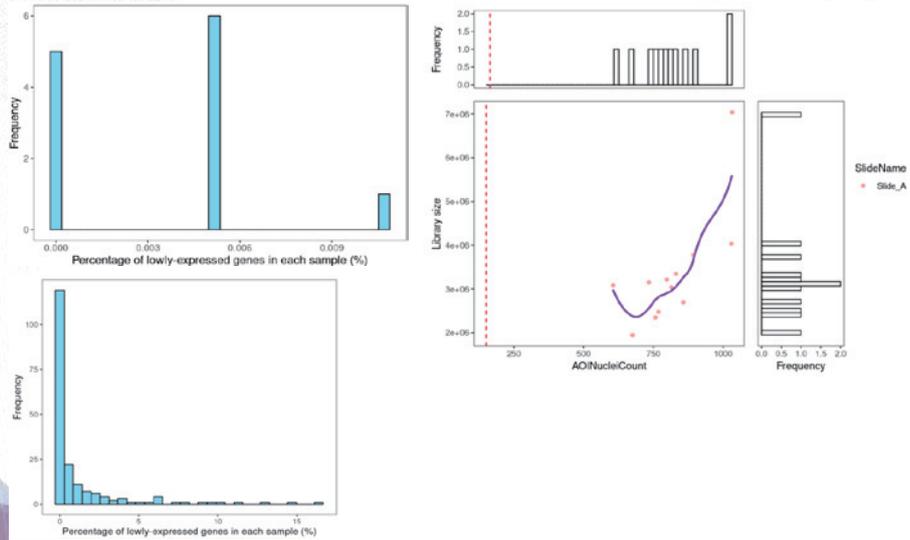
#####
# Load library #
#####
library(tidyverse)
library(standR)
library(SpatialExperiment)
library(edgeR)
library(limma)
library(msigdb)
library(GSEABase)
library(SpatialDecon)
library(speckle)

# Visualization #
library(ggplot2)
library(ggalluvial)
library(ggrepel)
library(DT)

#####
# Load data #
#####
countFile <- read_rds("/data/project/BIML_2024/GeoMX/Data/count.rds") %>% as.data.frame(); head(countFile)
sampleAnnoFile <- read_rds("/data/project/BIML_2024/GeoMX/Data/anno.rds") %>% as.data.frame(); head(sampleAnnoFile)
featureAnnoFile <- read_tsv("/data/project/BIML_2024/GeoMX/Data/genemeta.txt") %>% as.data.frame()
spe <- readGeoMX(countFile, sampleAnnoFile, featureAnnoFile)
```

## QC

```
#####  
# QC #  
#####  
# Gene level QC #  
spe <- addPerROIQC(spe, rm_genes = TRUE)  
plotGeneQC(spe, ordannots = "regions", col = regions, point_size = 2)  
  
# ROI level QC #  
plotROIQC(spe, x_threshold = 150, color = SlideName)  
qc <- colData(spe)$AOINucleiCount > 150; spe <- spe[, qc]  
  
# PCA #  
spe <- scater::runPCA(spe)  
pca_results <- reducedDim(spe, "PCA")  
plotPairPCA(spe, col = SlideName, precomputed = pca_results, n_dimension = 4)  
plotPairPCA(spe, col = class, precomputed = pca_results, n_dimension = 4)
```



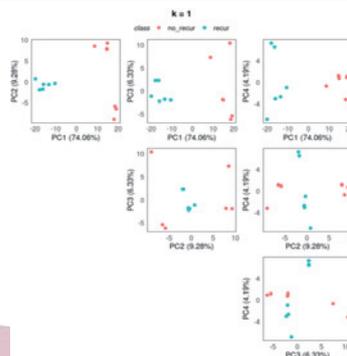
## Normalization

```
#####  
# Normalization #  
#####  
TMM, RPKM, TPM, CPM  
#####  
spe_tmm <- geomxNorm(spe, method = "TMM")  
plotRLEExpr(spe_tmm, assay = 2, color = SlideName) + ggtitle("TMM")
```

TMM: 각 샘플의 Library size를 이용하여 각 발현 수치를 보정하는 방법

## Batch correction

```
#####  
# Batch correction #  
#####  
spe <- findNCGs(spe, batch_name = "SlideName", top_n = 300)  
# for(i in seq(3)){  
#   spe_ruv <- geomxBatchCorrection(spe, factors = "class",  
#     NCGs = metadata(spe)$NCGs, k = i)  
#   #  
#   print(plotPairPCA(spe_ruv, assay = 2, n_dimension = 4, color = class, title = paste0("k = ", i)))  
# }  
# spe_ruv <- geomxBatchCorrection(spe, factors = "class",  
#   NCGs = metadata(spe)$NCGs, k = 1)
```



## DEGs

```
#####
# DEG #
#####
dge <- SE2DGEList(spe)
# design <- model.matrix(~0 + class + cov_W1 + cov_W2, data = colData(spe_ruv)); colnames(design) <- gsub("^class", "", colnames(design)); colnames(design) <- gsub(" ", "_", colnames(design))
# design <- model.matrix(~0 + class, data = colData(spe)); colnames(design) <- gsub("^class", "", colnames(design)); colnames(design) <- gsub(" ", "_", colnames(design))

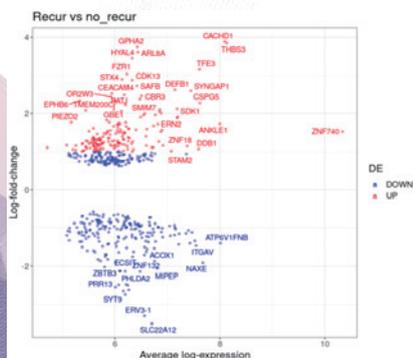
# IM vs CT
contr.matrix <- makeContrasts(
  BVT = recur - no_recur,
  levels = colnames(design)); keep <- filterByExpr(dge, design); table(keep)
dge_all <- dge[keep, ]
v <- voom(dge_all, design, plot = TRUE)

fit <- lmFit(v)
fit_contrast <- contrasts.fit(fit, contrasts = contr.matrix)
efit <- eBayes(fit_contrast, robust = TRUE)
results_efit <- decideTests(efit, p.value = 0.05)
de_results_BVT <- topTable(efit, coef = 1, sort.by = "P", n = Inf)
de_genes_topTable_BVT <- topTable(efit, coef = 1, sort.by = "P", n = Inf, p.value = 0.05)

de_results_BVT <- de_results_BVT %>%
  mutate(DE = ifelse(logFC > 0 & adj.P.Val < 0.05, "UP",
    ifelse(logFC < 0 & adj.P.Val < 0.05, "DOWN", "NOT DE"))); cor <- c("DOWN" = "#0000FF",
    "NOT DE" = "#808080",
    "UP" = "#FF0000")

DEG_vis(de_results_BVT,
  DEG_count = 500,
  title = "Recur vs no_recur",
  cor = cor)

updn_cols <- c(RColorBrewer::brewer.pal(6, 'Greens')[2], RColorBrewer::brewer.pal(6, 'Purples')[2])
DEG_table_vis(de_genes_topTable_BVT,
  DEG_count = 500,
  title = "Recur vs. no_recur (limma-voom)",
  cor = updn_cols)
```



logFC	AveExpr	P.Value	adj.P.Val	
SLC22A12	-3.508	6.705	3.699e-15	5.207e-11
ERV3-1	-3.304	6.563	8.363e-12	6.791e-9
SYT9	-2.747	6.196	3.078e-11	1.916e-8
BDKRB1	-2.656	6.160	3.521e-8	0.00004009
PIGT	-2.626	6.263	4.726e-11	2.596e-8
VPS50	-2.552	6.011	1.732e-9	3.762e-7
CD200R1	-2.501	6.246	1.774e-12	2.549e-9
PRR13	-2.490	6.074	1.958e-7	0.0001493

## GSEA (GeneSet Enrichment Analysis)

```
#####
# GSEA #
#####
# msigdb_hs <- getMsigdb(version = "7.2"); save(msigdb_hs, file = "/data/project/BIML_2024/GeoMX/msigdb_hs.RData")
load("/data/project/BIML_2024/GeoMX/Data/msigdb_hs.RData")
msigdb_hs <- appendKEGG(msigdb_hs)

sc <- listSubCollections(msigdb_hs)

gsc <- c(subsetCollection(msigdb_hs, c('h')),
  subsetCollection(msigdb_hs, 'c2', sc[grep("^CP:", sc)]),
  subsetCollection(msigdb_hs, 'c5', sc[grep("^GO:", sc)])) %>%
  GeneSetCollection()

fry_indices <- ids2indices(lapply(gsc, genes), rownames(v), remove.empty = FALSE)
names(fry_indices) <- sapply(gsc, setName)
gsc_category <- sapply(gsc, function(x) bcCategory(collectionType(x)))
gsc_category <- gsc_category[sapply(fry_indices, length) > 5]

gsc_subcategory <- sapply(gsc, function(x) bcSubCategory(collectionType(x)))
gsc_subcategory <- gsc_subcategory[sapply(fry_indices, length) > 5]

fry_indices <- fry_indices[sapply(fry_indices, length) > 5]
names(gsc_category) = names(gsc_subcategory) = names(fry_indices)

fry_indices_cat <- split(fry_indices, gsc_category[names(fry_indices)])
fry_res_out <- lapply(fry_indices_cat, function(x) {
  limma::fry(v, index = x, design = design, contrast = contr.matrix[,1], robust = TRUE)
})

fry_res_sig <- post_fry_format(fry_res_out, gsc_category, gsc_subcategory) %>%
  as.data.frame() %>%
  filter(FDR < 0.25) # 0.05 ~ 0.25

# GSEA - vis
GSEA_vis(df = fry_res_sig,
  DEG_type = "Up",
  cor <- "#E41E22",
  cnt = 10,
  title <- "up-deg")
GSEA_vis(df = fry_res_sig,
  DEG_type = "Down",
  cor <- "#377EB8",
  cnt = 10,
  title <- "down-deg")
```

