

KSBi-BIML 2024

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

생명정보학 & 머신러닝 워크샵 (온라인)



Introduction to cancer-immune analysis

김상우 _ 연세대학교



KSBI
KOREAN SOCIETY FOR
BIOINFORMATICS

| 한국생명정보학회



본 강의 자료는 한국생명정보학회가 주관하는 BIML 2024 워크샵 온라인 수업을 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 전적으로 불법 행위자 본인에게 있음을 경고**합니다.

KSBI-BIML 2024

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

안녕하십니까?

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

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

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

2024년 2월

한국생명정보학회장 이 인 석

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 Transcriptomics를 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-seq을 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
Phone Number 010-3407-9861

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, Eunyoung 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†, Eunyoung 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

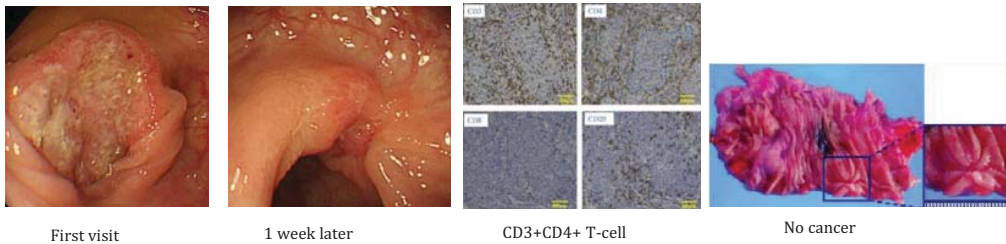
3

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

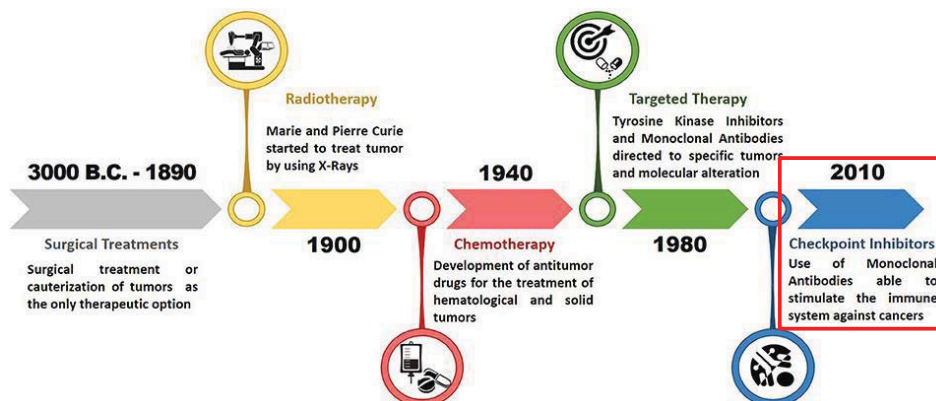
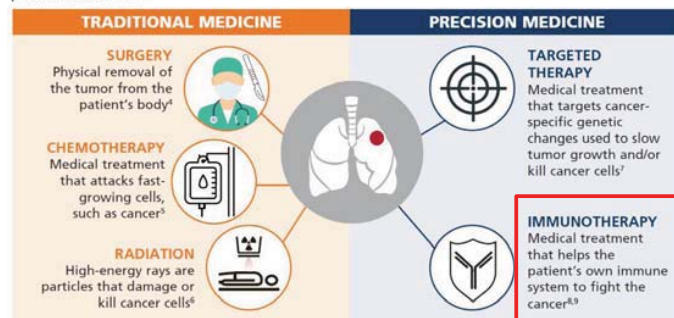
CD3+CD4+ T-cell

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.^{2,3}



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.¹

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.

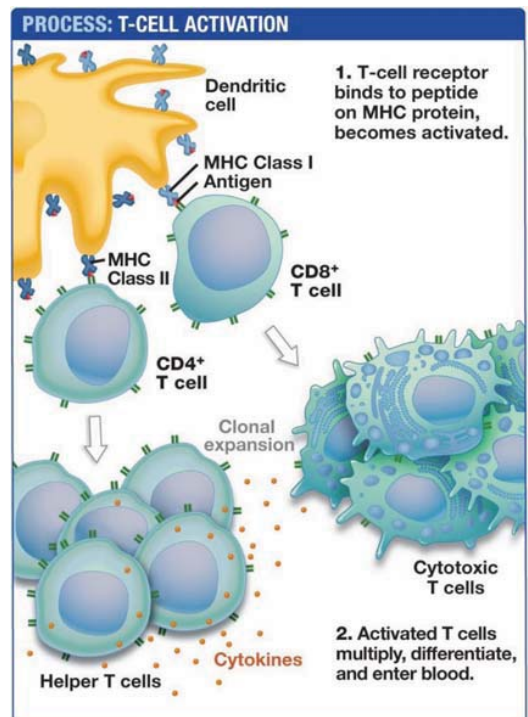
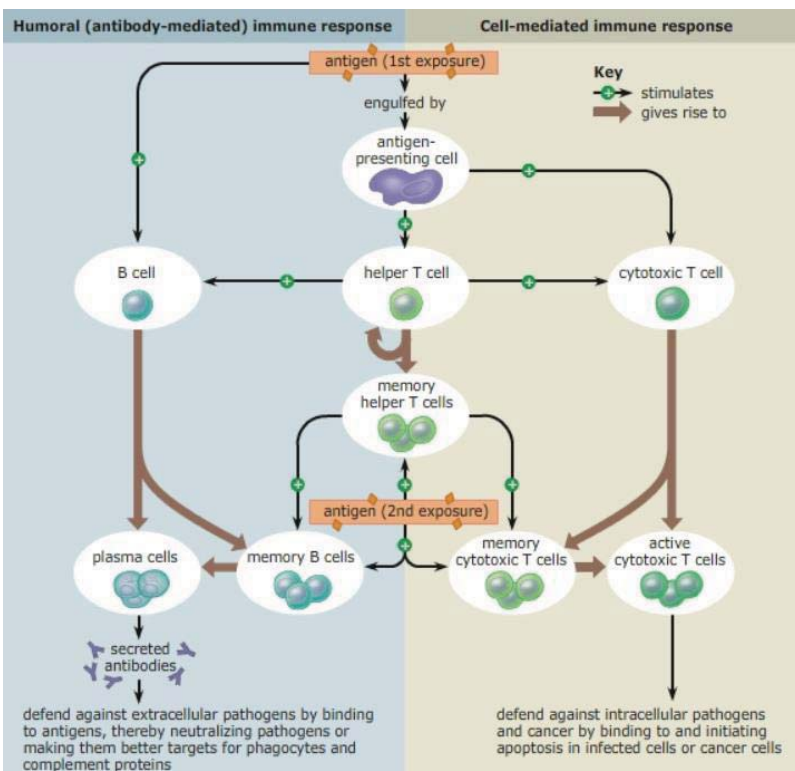
August 12. The pain and swelling continuing, she again sought

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

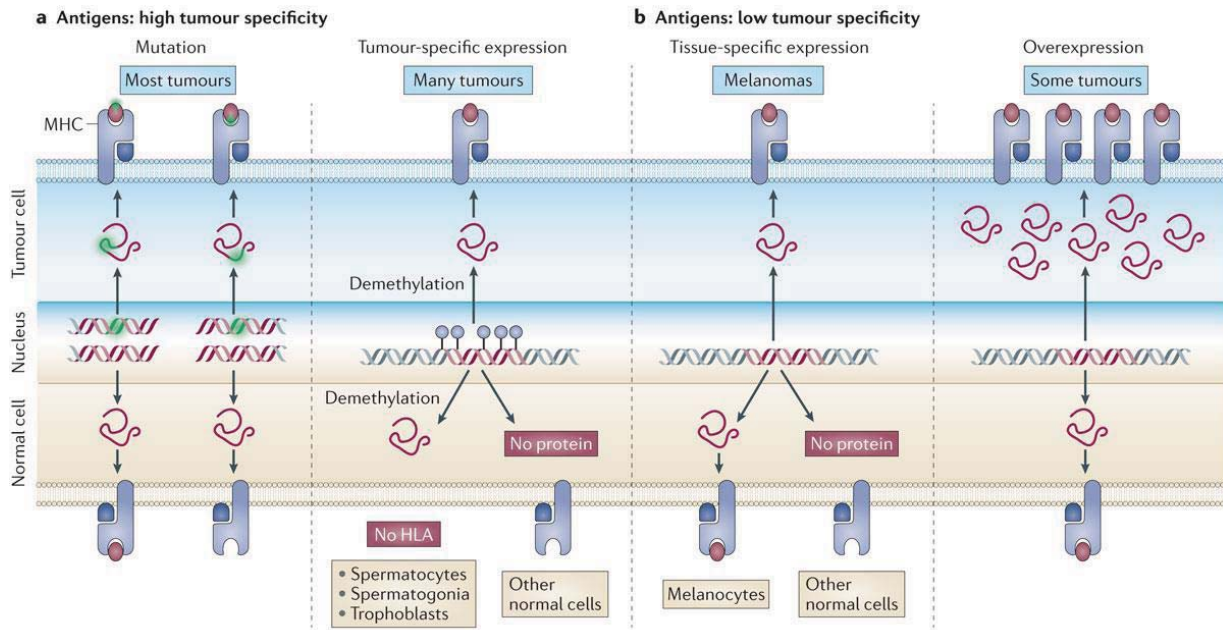
(199)

Coley, Annals of Surgery, 1981

Adaptive Immunity / T-cell activation



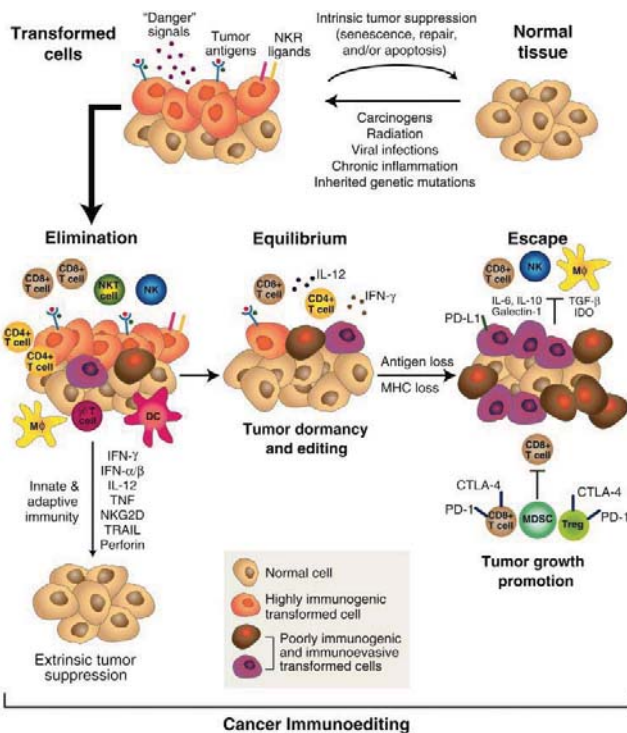
Tumor Antigens



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

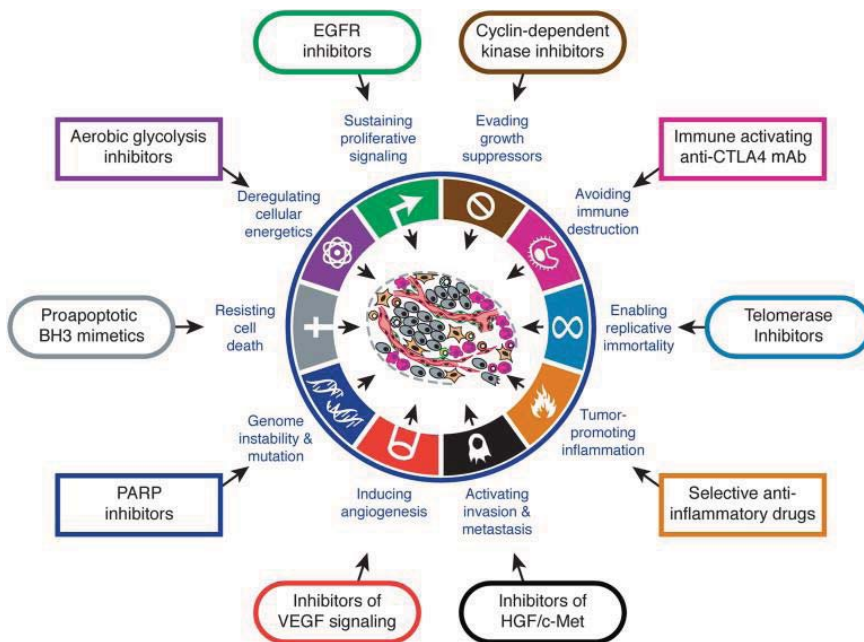
Nature Reviews | Cancer

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

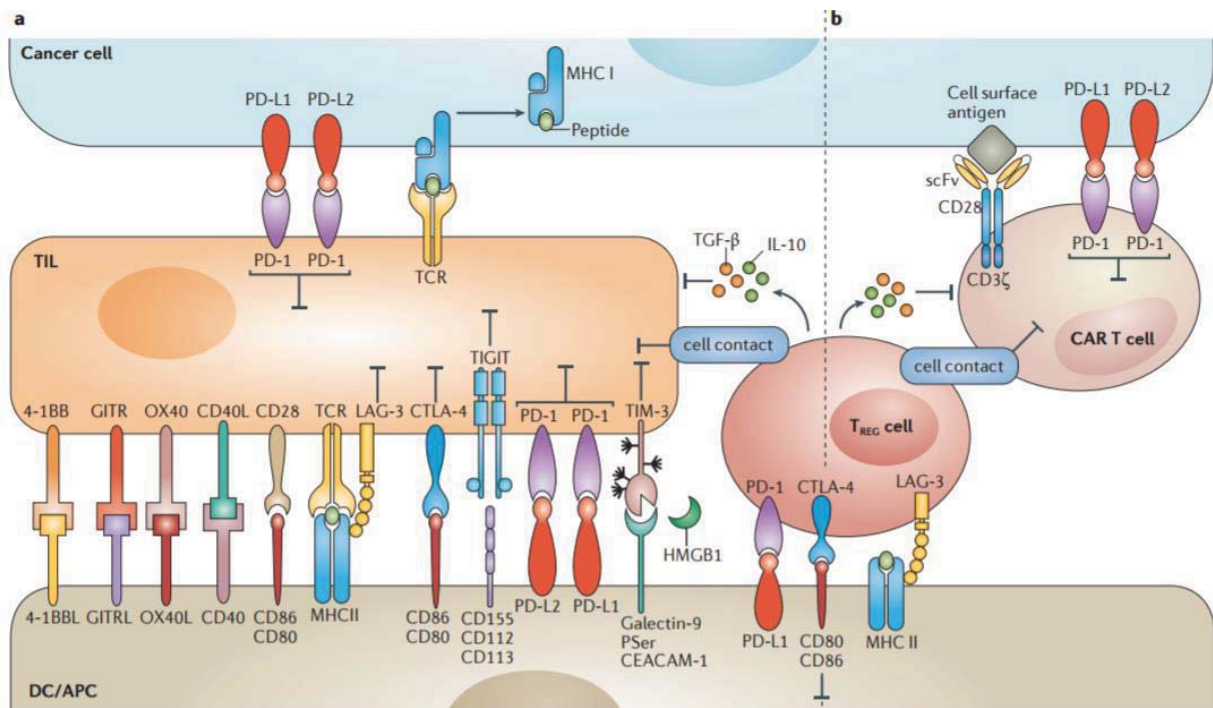


- Paralyze CTLs and NK cells by secreting TGF- β or immunosuppressive factors
- Recruitment of regulatory T-cell (Tregs) and myeloid-derived suppressor cells (MDSCs)
- Loss of MHC class I expression

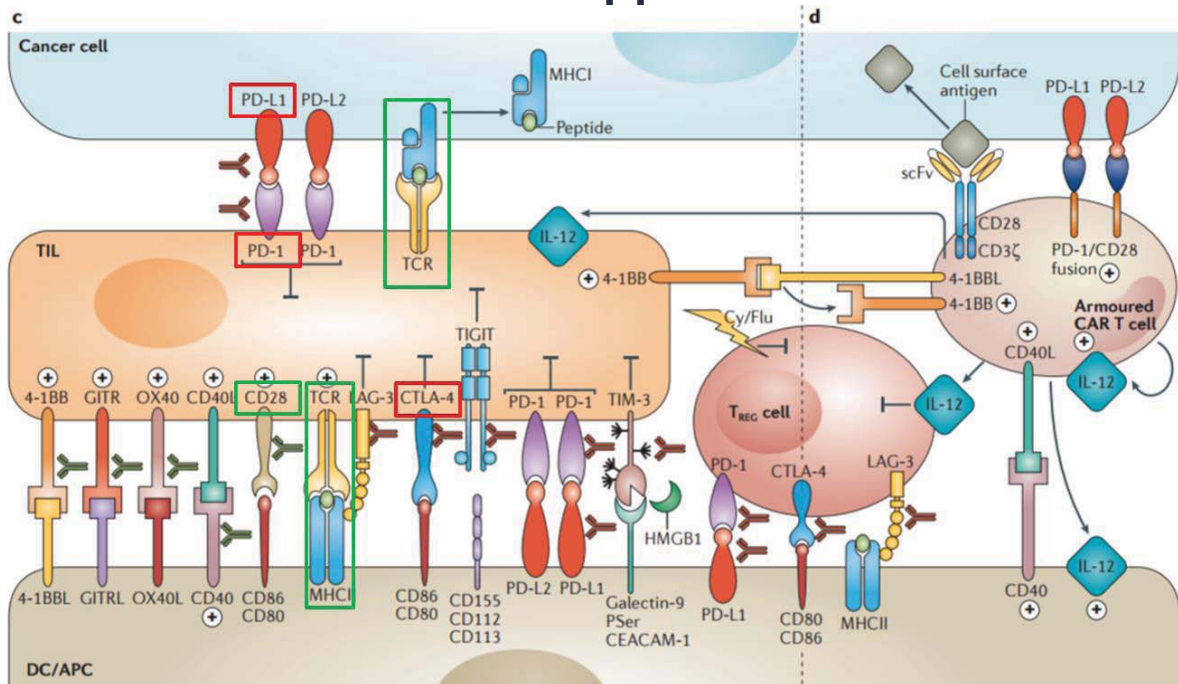
Hannahan and Weinberg, Hallmarks of cancer: The Next Generation, Cell 2011

CURRENT APPROACHES

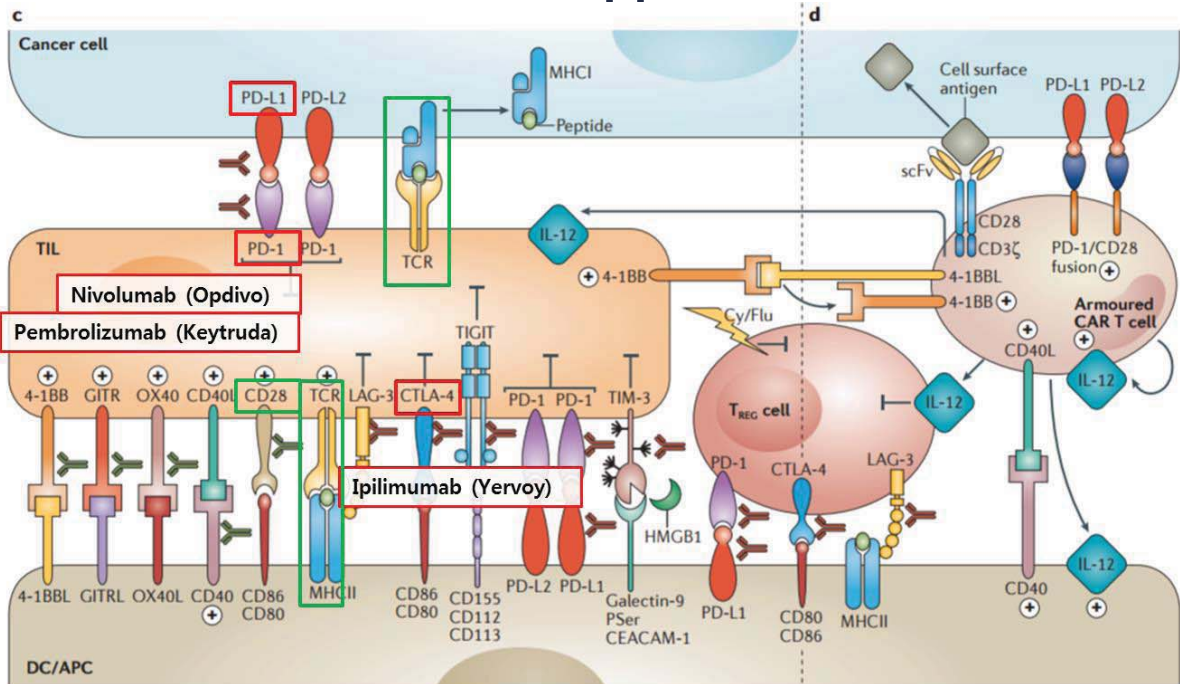
2. Checkpoint inhibitors



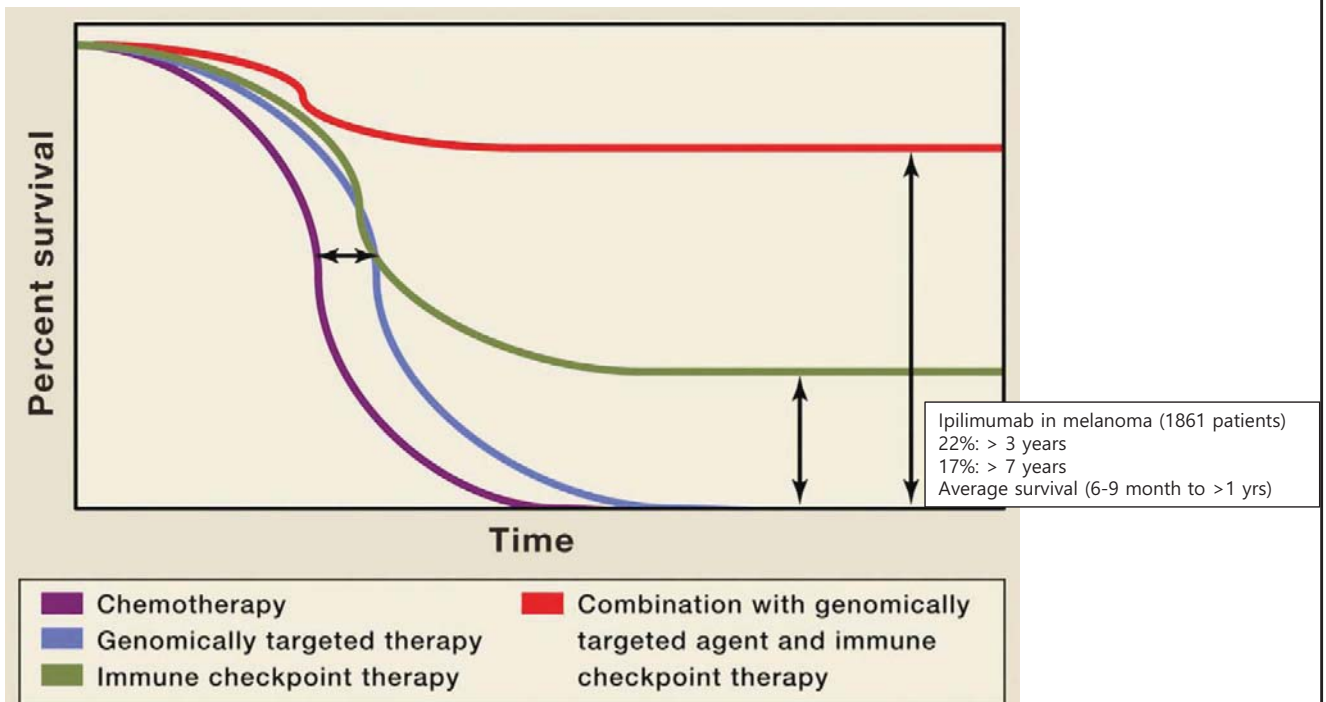
Immunomodulatory mAbs to overcome immunosuppression



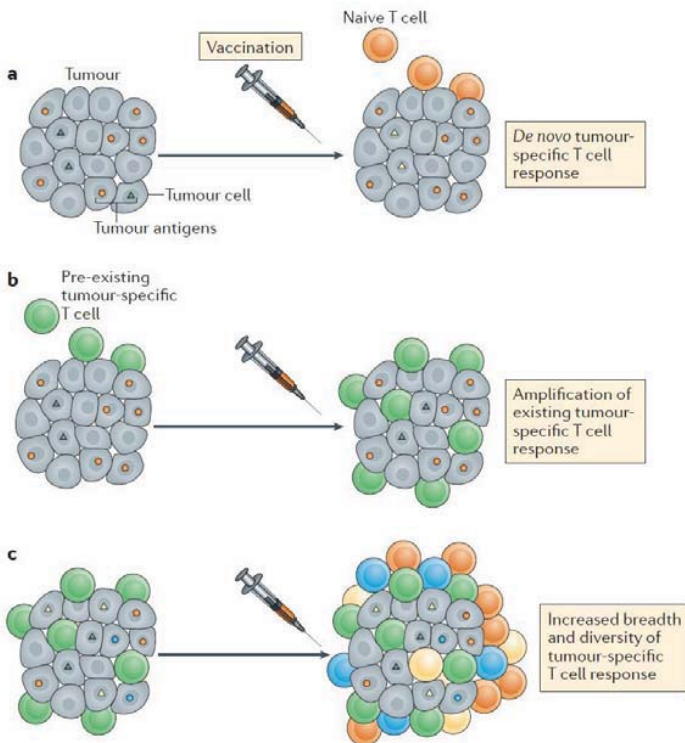
Immunomodulatory mABs to overcome immunosuppression



The benefits from cancer immunotherapy



3. Cancer Vaccine

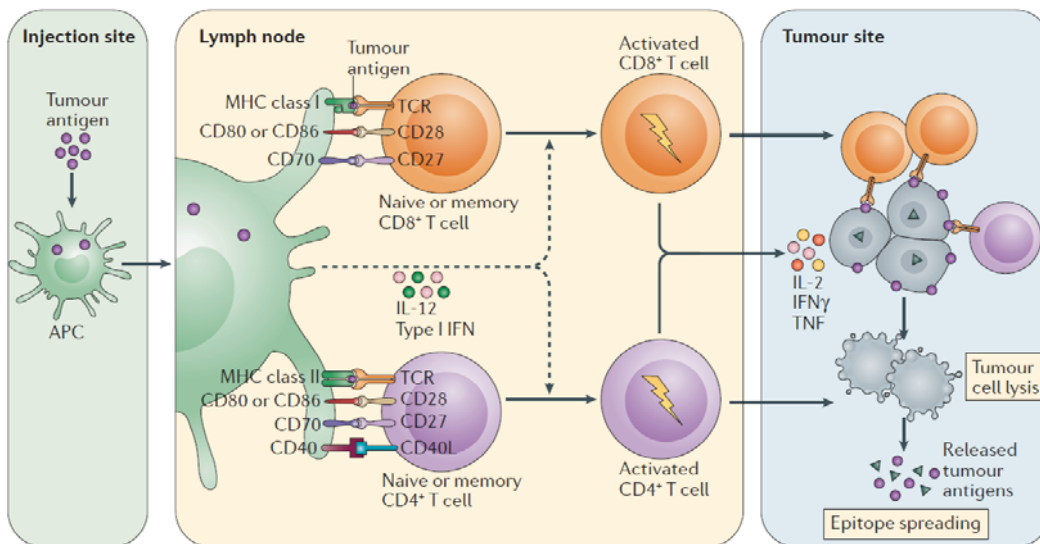


Cancer vaccines:

- 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

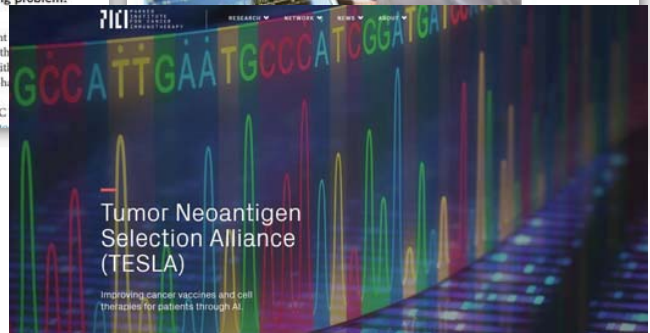
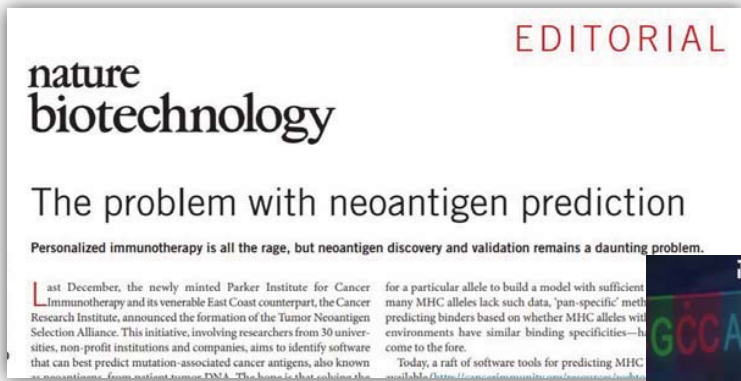
How cancer vaccine works



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 γ , 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

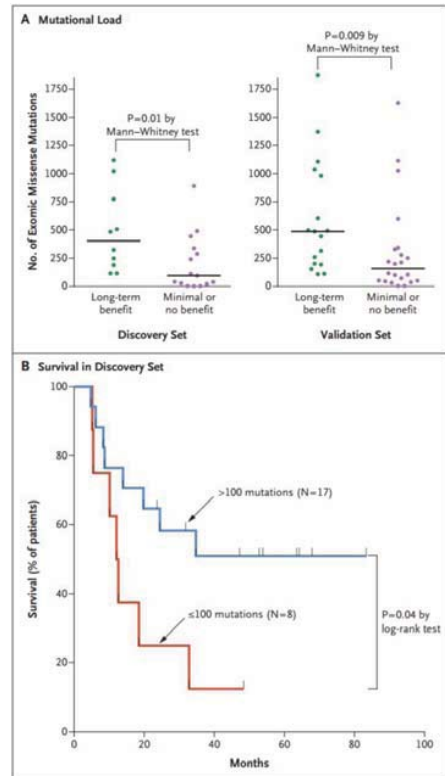
ORIGINAL ARTICLE

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

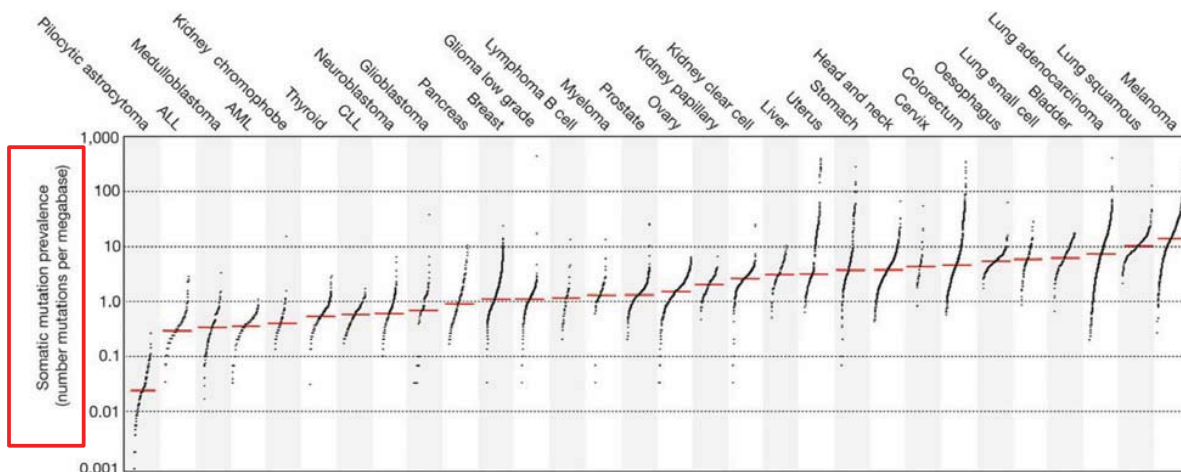
Alexandra Snyder, M.D., Vladimir Makarov, M.D., Taha Merghoub, Ph.D., Jianda Yuan, M.D., Ph.D., Jesse M. Zaretsky, B.S., Alexis Desrichard, Ph.D., Logan A. Walsh, Ph.D., Michael A. Postow, M.D., Phillip Wong, Ph.D., Teresa S. Ho, B.S., Travis J. Hollmann, M.D., Ph.D., Cameron Bruggeman, M.A., Kasthuri Kannan, Ph.D., Yanyun Li, M.D., Ph.D., Ceyhan Elipenahli, B.S., Caillan 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

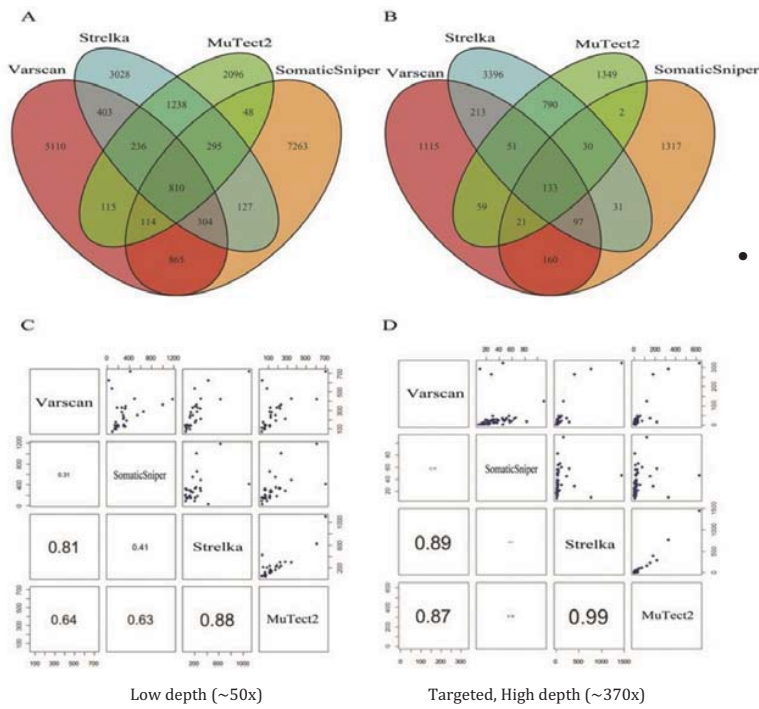


Tumor mutation burden



• Tumor Mutation Burden (TMB) =
$$\frac{\#total_somatic_mutation}{total_targeted_genome_size(Mb)}$$

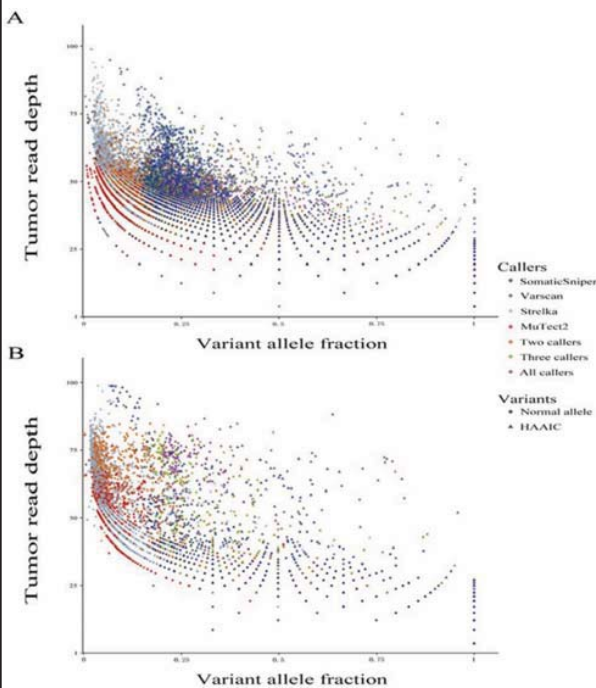
Inconsistence of somatic mutation calls



- 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

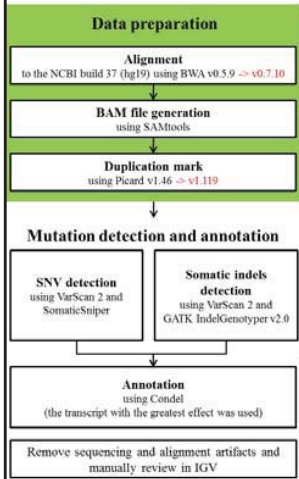
Fixing pipeline

mut/MB
(SNV)

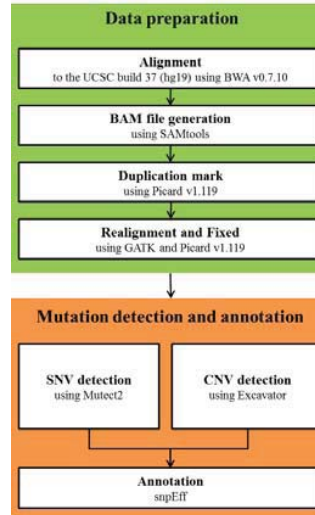
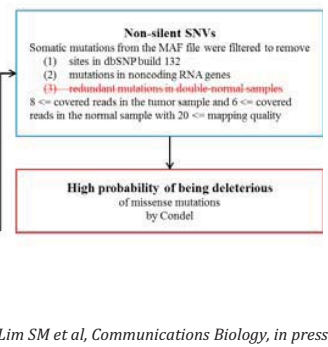
5.533993
5.178398
3.056166
1.459616
1.471475
1.453428
1.706356

mut/MB

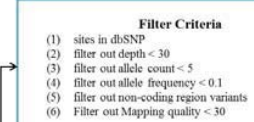
2.991871
2.4023
1.641857
1.27743
1.820113
1.108711
1.003712



TCGA Flow chart



In-house Flow chart



Potential pitfalls (use with care)

VIEWPOINT

Tumor Mutation Burden—From Hopes to Doubts

Alfredo Addeo, MD
Oncology Department, Geneva University Hospital, Geneva, Switzerland.

Giuseppe L. Barina, 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.

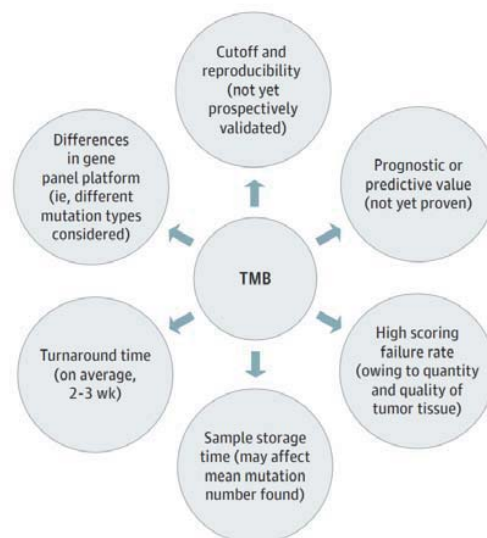
Over the past few years, the development of immune checkpoint inhibitors has altered the treatment paradigm in non-small cell lung cancer (NSCLC). Enrichment strategies have identified programmed death-ligand 1 (PD-L1) staining by immunohistochemistry to be a predictive biomarker in treatment-naïve patients with refractory NSCLC. In particular, Keynote-024 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 repertoires. 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,⁷ 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

team⁸ recently calculated TMB scores by whole-exome sequencing in a subset of patients from the CheckMate-026 study,⁹ 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.

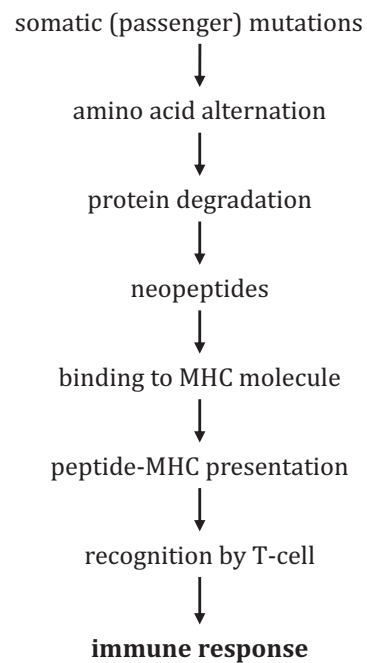
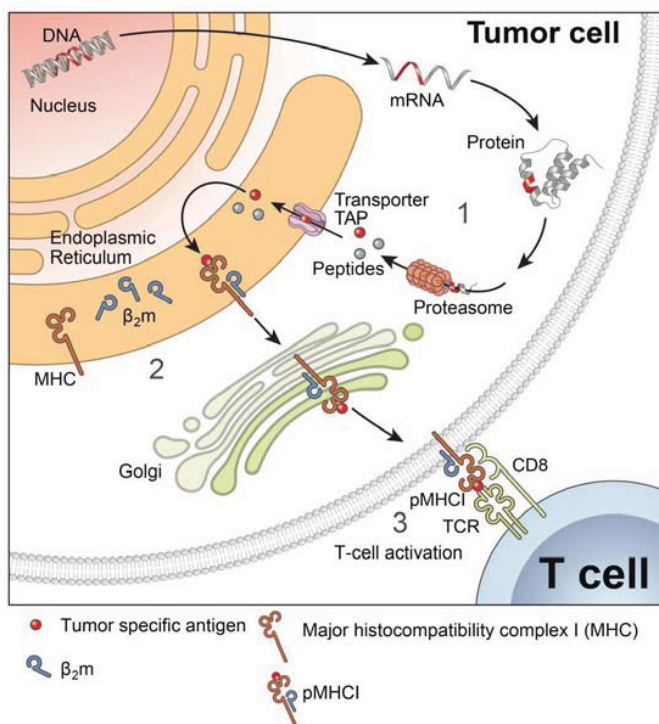
Guidelines from the European Society for Medical Oncology (ESMO) and ESMO Asia have already incorporated TMB as a possible biomarker in advanced NSCLC, recommending the combination of ipilimumab plus nivolumab as first-line treatment for patients with high TMB (>10 mutations/Mb). Supporting evidence stems from the CheckMate-227 trial, which reported results for first-line nivolumab plus ipilimumab vs platinum doublet chemotherapy.¹⁰ 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 trial's prespecified end point for analysis. The trial had

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

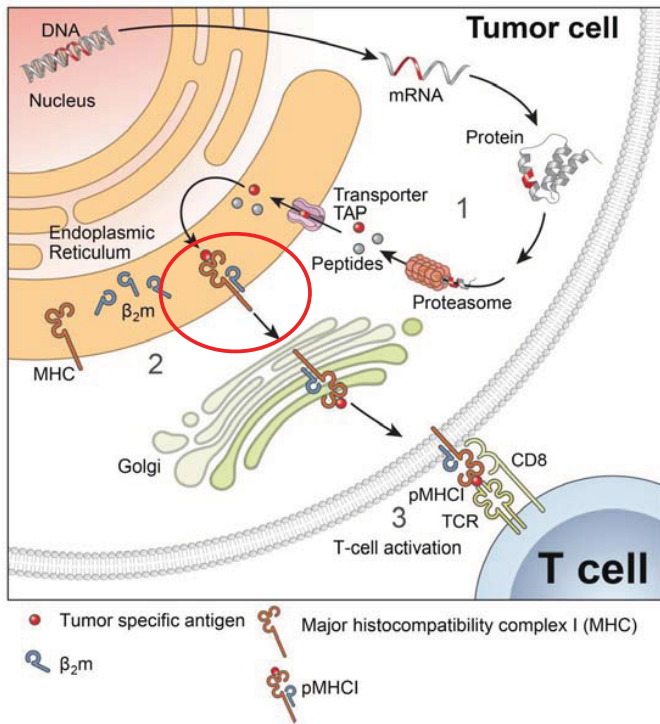


HLA TYPING IN THE ANTIGEN PROCESSING

Neoantigen processing



Neoantigen processing



somatic (passenger) mutations

↓
amino acid alternation

↓
protein degradation

↓
neopeptides

↓
binding to MHC molecule

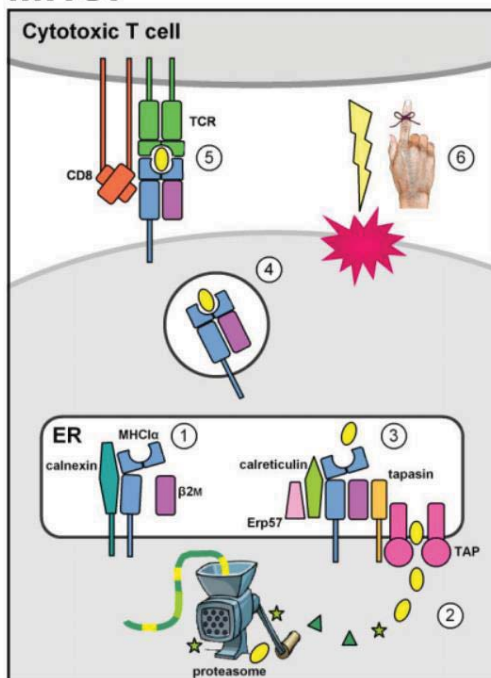
↓
peptide-MHC presentation

↓
recognition by T-cell

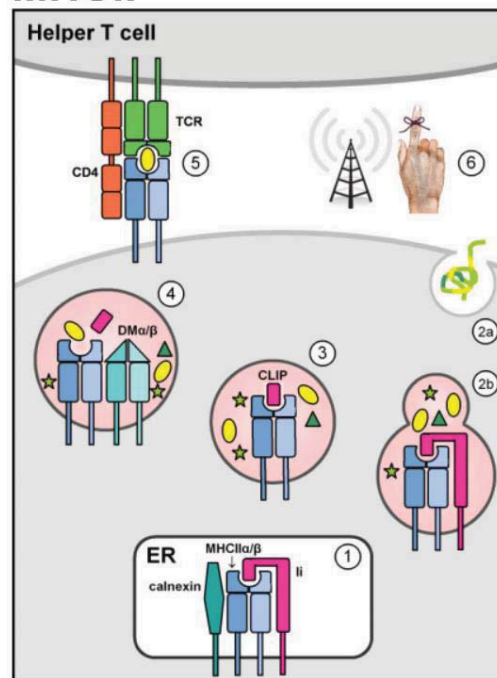
↓
immune response

MHC (Major Histocompatibility Complex)

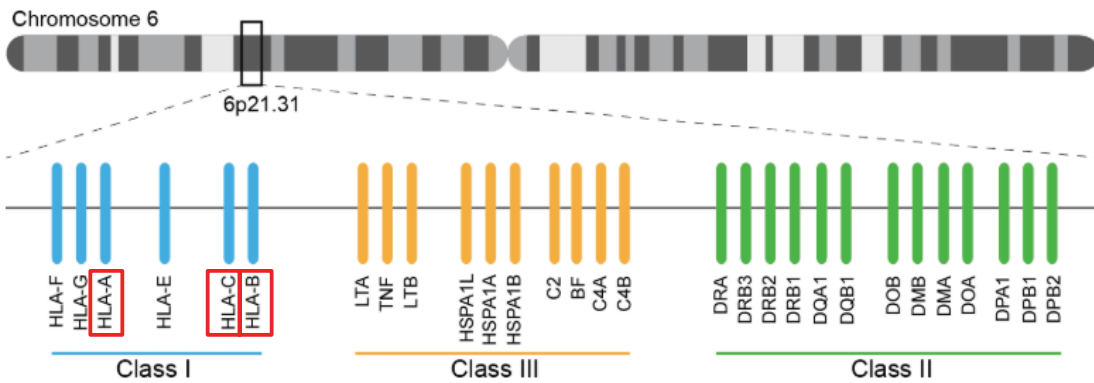
MHCI



MHCII



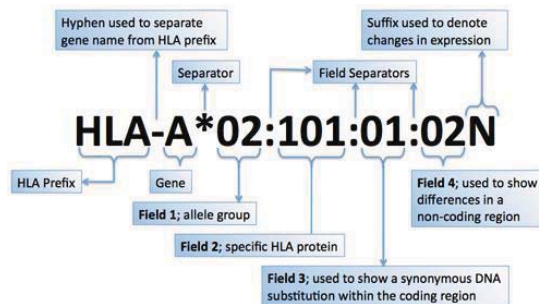
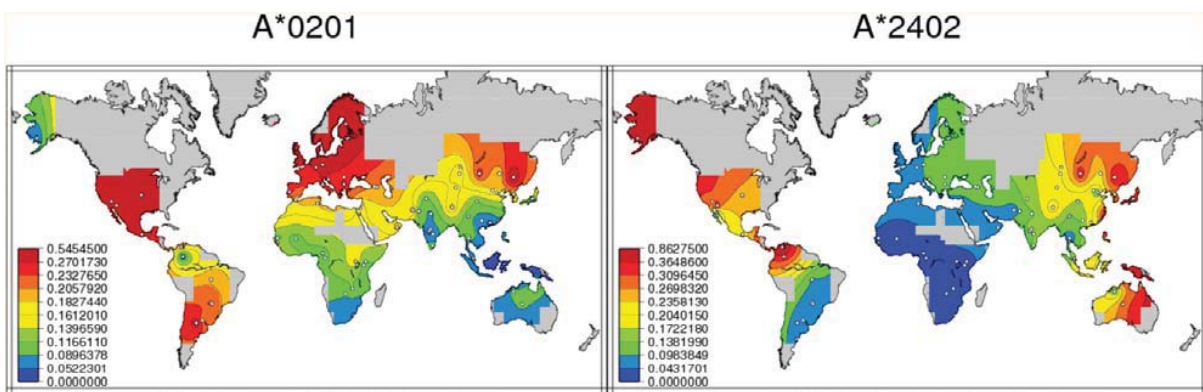
HLA (Human Leukocyte Antigen)



<p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p> <p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p> <p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p> <p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p> <p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p>	<p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p> <p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p> <p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p> <p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p> <p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p>	<p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p> <p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p> <p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p> <p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p> <p>AA Codon A*24:02:01:01 A*24:156 A*24:191</p>
---	---	---

SBI 한국생명정보학회
Korean Society for Bioinformatics

HLA alleles are ethnic specific



SBI 한국생명정보학회
Korean Society for Bioinformatics

MHC-peptide binding

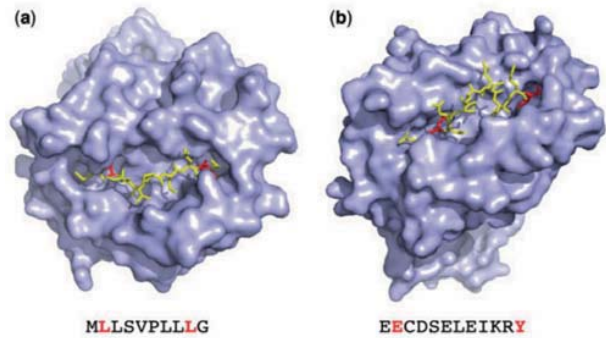
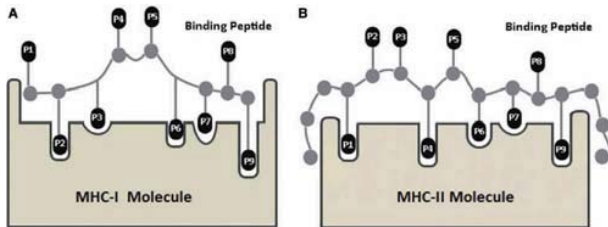
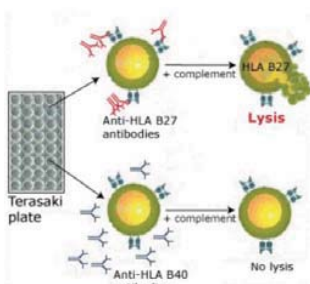


Fig. 5. 3D structures for two MHC class I molecules with bound peptides longer than 9 amino acids (PDB references 2CLR and 4JQX). (a) The 10mer peptide MLLSVPLLLG bound to HLA-A*02:01 extends at the C terminus with a glycine (G) amino acid. The residues at the anchor positions P2 (L) and P9 (L) are highlighted. (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)

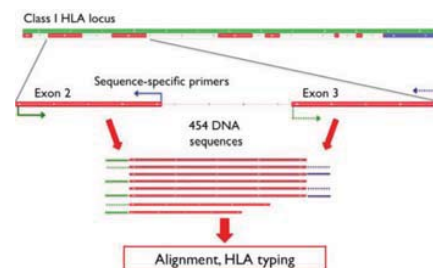
HLA typing methods

1. Serology-based typing

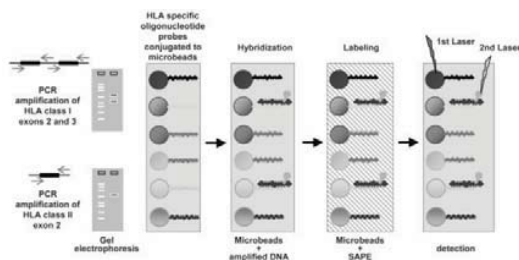


- Use of microcytotoxicity - complement mediated lysis
- Simple and low-cost
- Mostly used in HLA-A and HLA-B
- Can type allele groups and alleles only

2. Sanger sequencing



3. Sequence-specific Oligonucleotide Hybridization (SSO)

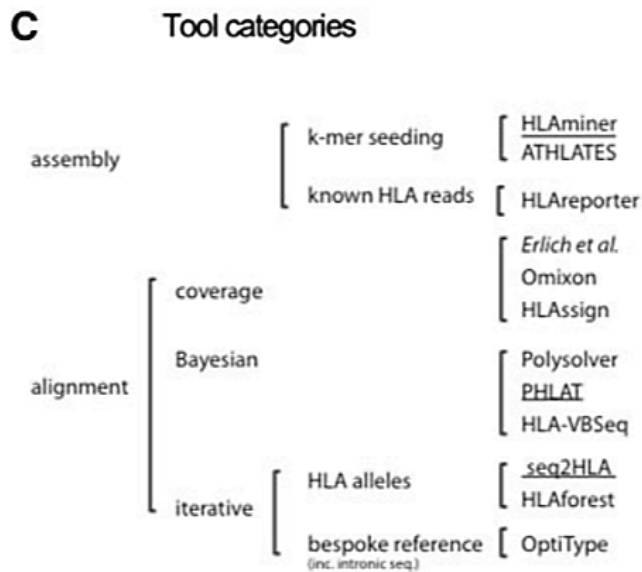


- Amplify targeted regions with biotin-labeled primers
- Hybridized sequences emit fluorescence

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



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

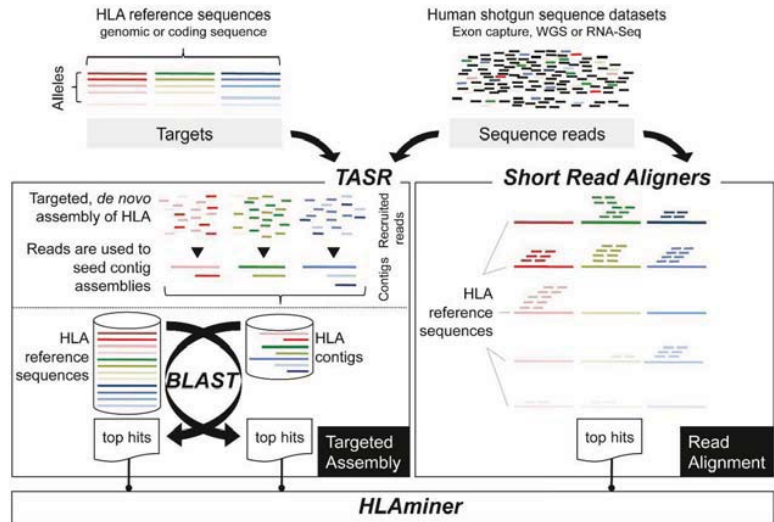
Derivation of HLA types from shotgun sequence datasets

René L. Warren¹, Gina Choe¹, Douglas J. Freeman¹, Mauro Castellani¹, Sarah Munro¹, Richard Moore¹ and Robert A. Holt^{1,2*}

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 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 HLAmimer, a computational method for identifying HLA alleles directly from shotgun sequence datasets (<http://www.biopython.org/platform/bioinformatics/software/HLAmimer>). 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.

HLAmimer



Alignment-based HLA typing

ANALYSIS

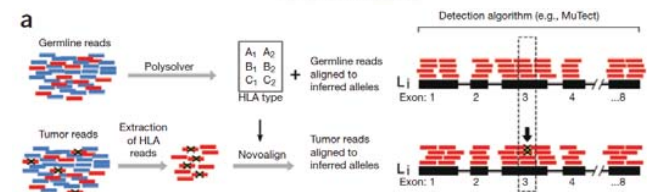
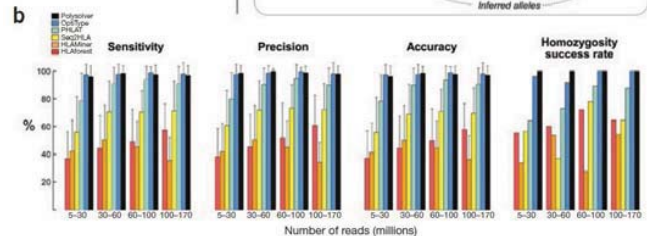
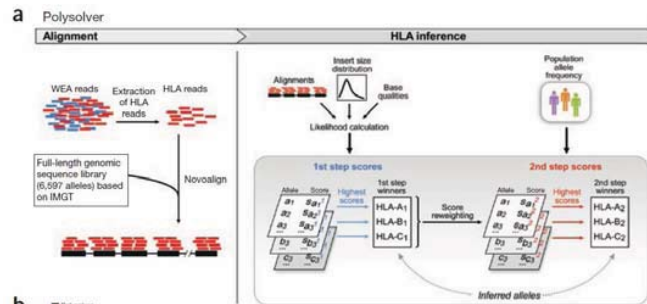
computational
BIOLOGY
nature
biotechnology

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

Sachet A Shukla¹⁻³, Michael S Rooney^{2,4}, Mohini Rajasagi^{1,5}, Grace Tiao², Philip M Dixon³, Michael S Lawrence², Jonathan Stevens⁶, William J Lane^{6,7}, Jamie L Dellagatta⁶, Scott Steedman², Carrie Soungner², Kristian Cibulski², Adam Kiezun², Nir Hacohen^{2,8,9}, Vladimir Brusic^{1,5}, Catherine J Wu^{1,2,5,8,11} & Gad Getz^{2,10,11}

Detection of somatic mutations in human leukocyte antigen (HLA) cancer genes *whole-exome sequencing* adenocarcinoma and diffuse large B-cell lymphoma¹⁻⁵. The HLA locus, located on chromosome 6, is among the most polymorphic

Polysolver



MHC BINDING PREDICTION

MHC-peptide binding

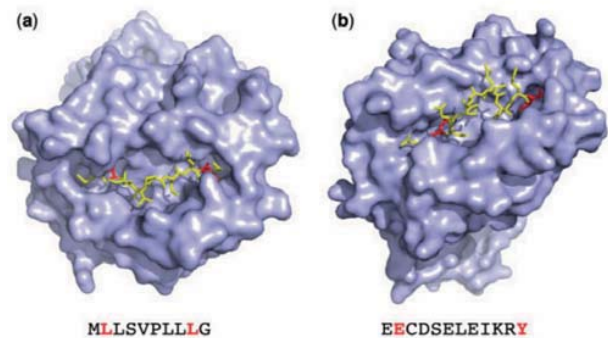
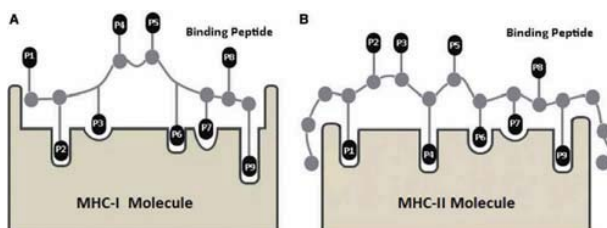
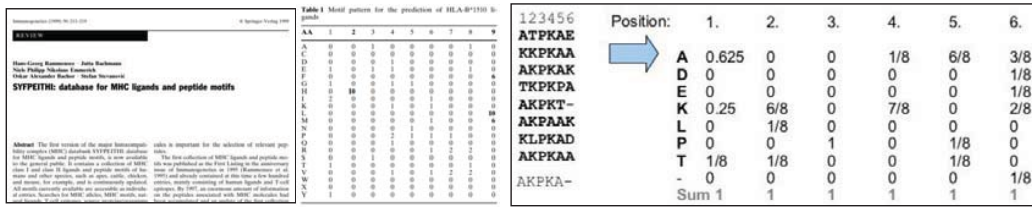


Fig. 5. 3D structures for two MHC class I molecules with bound peptides longer than 9 amino acids (PDB references 2CLR and 4JQX). (a) The 10mer peptide MLLSVPLLLG bound to HLA-A*02:01 extends at the C terminus with a glycine (G) amino acid. The residues at the anchor positions P2 (L) and P9 (L) are highlighted. (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

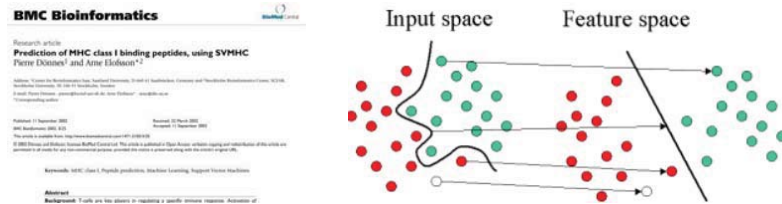
Can we predict if a given peptide will bind to MHC?

Prediction algorithms

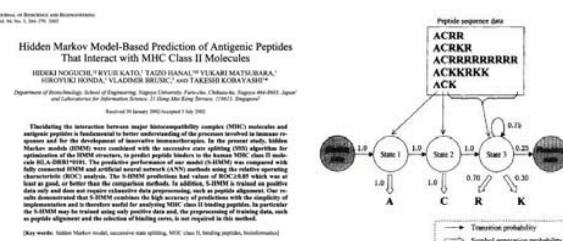
- SYFPEITHI: using PSSM



- SVMHC: using Support Vector Machine



- S-HMM: using Hidden Markov Model



SBI 한국생명정보학회
Korean Society for Bioinformatics

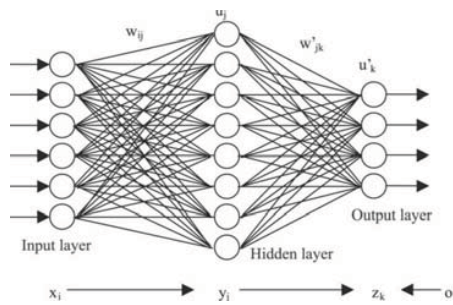
ANN based algorithms

NetMHC: Classification of MHC-I binding peptides using ANN

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

MORTEN NIELSEN,¹ CLAUS LUNDEGAARD,¹ PEDER WORNING,¹ SANNIE LISE LAUEMOLLER,² KASPER LAMBERTH,² SOREN BUUS,² SOREN BRINKA,³ AND OLE LIND¹
¹Center for Biological Sequence Analysis, BioCentrum-DTU, Technical University of Denmark, DK-2800 Lyngby, Denmark
²Department of Experimental Immunology, Institute of Medical Microbiology and Immunology, University of Copenhagen, Blegdamsvej 3C, DK-2200 Copenhagen, Denmark
(Received November 14, 2002; Accepted February 19, 2003)

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 mutual information calculations we show that



NetMHC-3.0

BIOINFORMATICS APPLICATIONS NOTE
doi:10.1093/bioinformatics/btl038

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 Lundegaard¹, Ole Lund¹ and Morten Nielsen¹
¹Center for Biological Sequence Analysis - CBS, Department of Systems Biology, The Technical University of Denmark - DTU, Kemitorvet Bldg. 206, 2800 Lyngby, Denmark
Received on February 9, 2006; revised and accepted on April 4, 2006
Advance Access publication April 16, 2006
Associate Editor: Stuart B. Pevsley

Approximation of 8, 10, 11 from 9 mer model

FIG. 3.

A. EIGHTMER → { EIGHTMER, EIGHTMER, EIGHTMER, EIGHTMER }

B. TENLEVEN → { TENLEVEN, TENLEVEN, TENLEVEN, TENLEVEN }

C.

peptide	log(odds ratio)	log(odds ratio)
TENLEVEN	0.189	6489
TENLEVEN	0.325	6318
TENLEVEN	0.234	3979
TENLEVEN	0.355	1896
TENLEVEN	0.377	846
TENLEVEN	0.481	653

GENETIC CODE 2148 aa

NetMHC-4.0

Sequence analysis
Gapped sequence alignment using artificial neural networks: application to the MHC class I system
Massimo Andreatta¹ and Morten Nielsen^{1,2*}

(a) A I L D F T H L

Sequence	Score
X A I L D F T H L	0.043
A X I L D F T H L	0.013
A T X L D F T H L	0.562
A I L X D F T H L	0.743
A I L D X F T H L	0.425
A I L D F X T H L	0.523
A I L D F T X H L	0.505
A I L D F T H X L	0.366
A I L D F T H L X	0.013

(b) F Y G E R L T R Y

Sequence	Score
F Y G E R L P T R Y	0.103
F Y G E R L T R Y	0.012
F Y G R L P T R Y	0.378
F Y G E R L T R Y	0.466
F Y G E P L T R Y	0.462
F Y G E R L T R Y	0.712
F Y G E R P T R Y	0.609
F Y G E R P L R Y	0.598
F Y G E R P L T Y	0.309
F Y G E R P L T R Y	0.111

Gapped alignment to ANN : 9 to 8~11 mer

SBI 한국생명정보학회
Korean Society for Bioinformatics

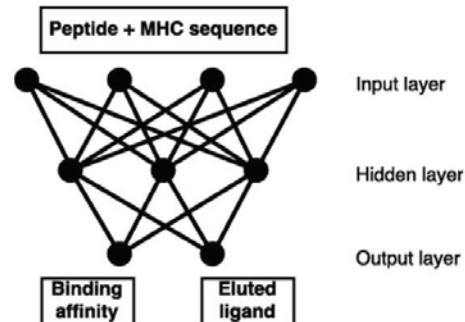
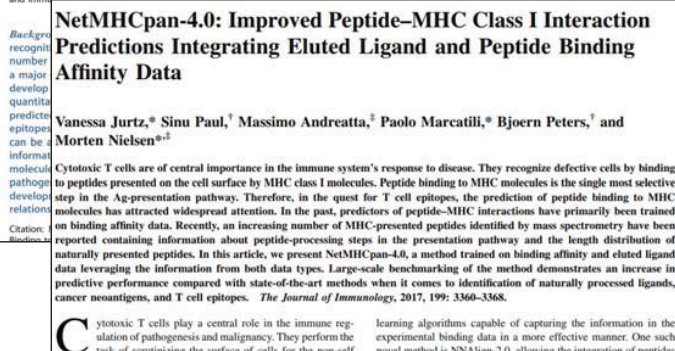
Regarding all HLA-types at once

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



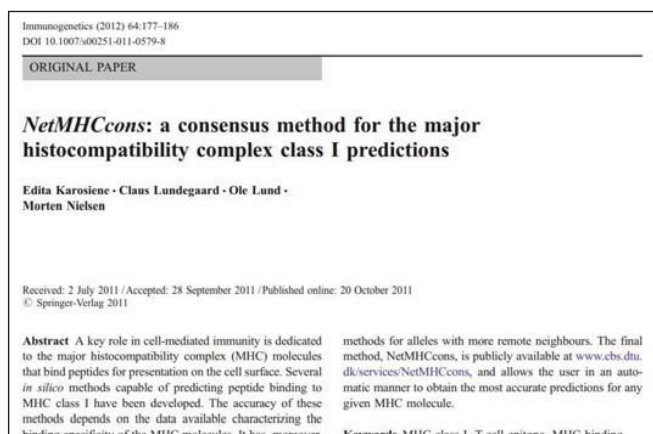
Experimental data are biased to major HLA alleles
 ▶ lack of training data in rare alleles
 ▶ lack of accuracy

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) 26–34

Contents lists available at ScienceDirect

Journal of Immunological Methods

ELSEVIER

journal homepage: www.elsevier.com/locate/jim

Research paper

Prediction of epitopes using neural network based methods

Claus Lundegaard^a, Ole Lund, Morten Nielsen

Center for Biological Sequence Analysis, DTU Systems Biology, Building 208, Technical University of Denmark, DK-2800 Lyngby, Denmark

ARTICLE INFO

Article history:
Received 30 July 2010
Received in revised form 23 October 2010
Accepted 27 October 2010
Available online 31 October 2010

ABSTRACT

In this paper, we describe the methodologies behind three different aspects of the NetMHC family for prediction of MHC Class I binding, mainly to HLA. 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

Keywords:
MHC
Binding
Prediction
Epitope
Discovery
T cell

Metric	NetMHCpan	SMM	ANN	ARB
Overall	~68	~48	~68	~35
SRCC	~65	~45	~68	~35
AUC	~68	~50	~68	~30

SBI 한국생명정보학회
Korean Society for Bioinformatics

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 a target molecule, the competitors are asked to submit a set of predicted eluted peptides from the test set.

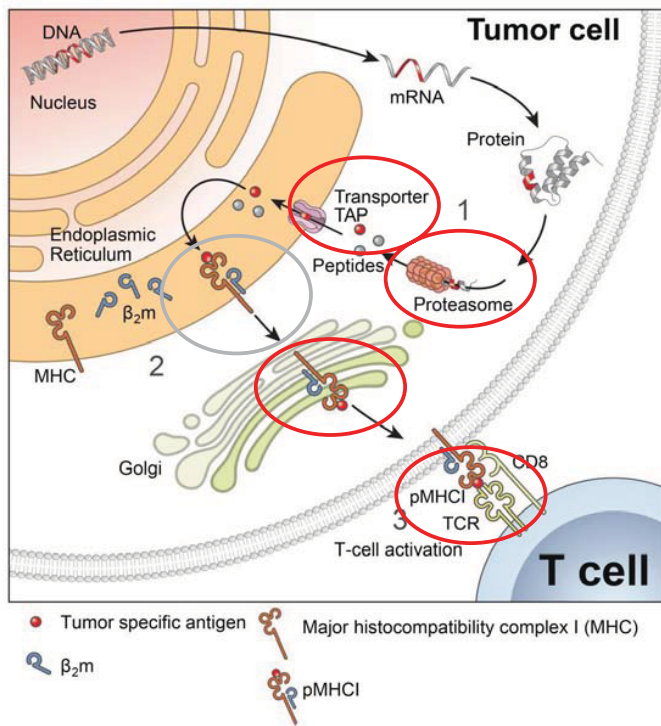


A total of 32 submissions were submitted for the competition. Of these, 24 submissions (Group 1) provided a set of thresholds (elution score based predictors) for each peptide and each MHC molecule. Another 8 submissions (Group 2) provided lists of peptides that were predicted as eluted from specific 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, Hårdahl M, Baus 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 ^b , and H-2K ^b
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 (GS 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

amino acid alternation

protein degradation

neopeptides

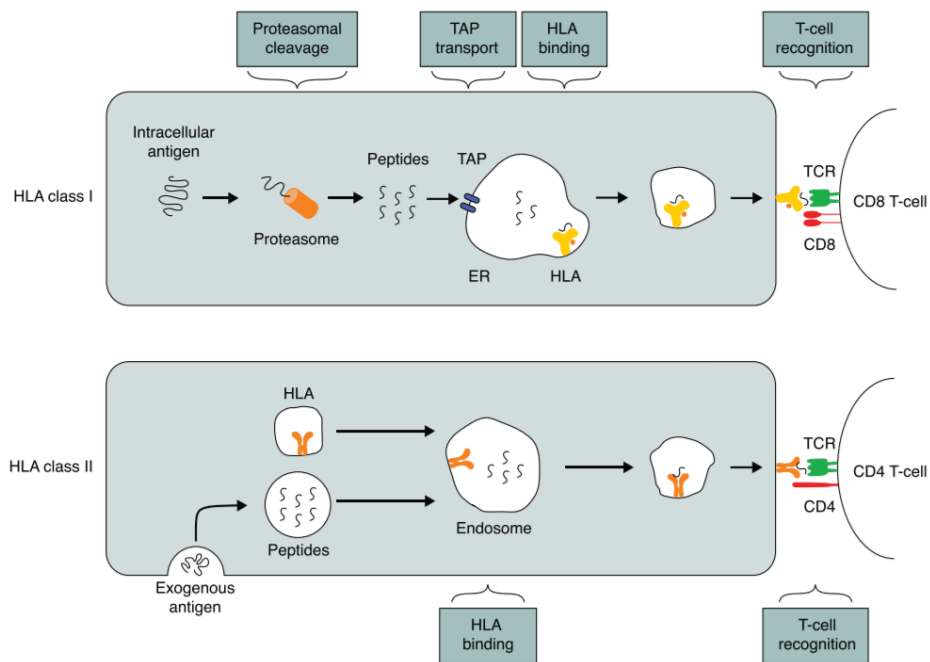
binding to MHC molecule

peptide-MHC presentation

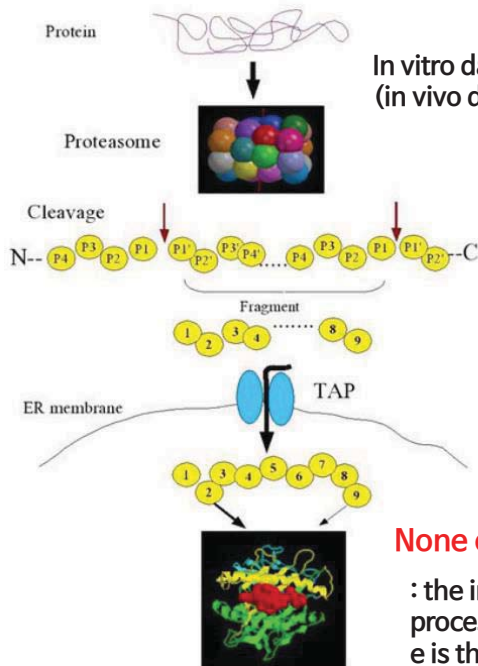
recognition by T-cell

immune response

Antigen Processing Pathways for MHC class I/II



Proteasomal cleavage



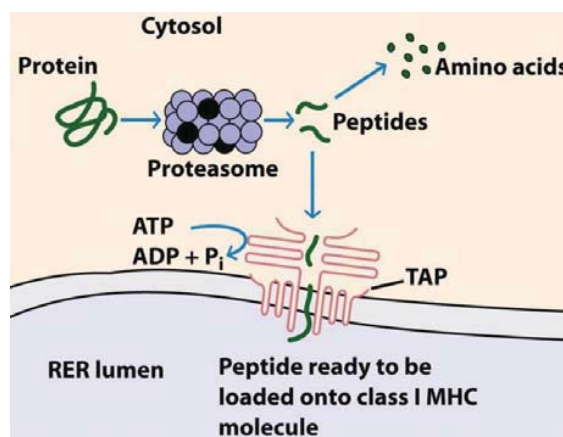
In vitro data created with purified proteasomes in the laboratory (in vivo data are harder to collect)

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 Harndahl¹, Michael Rasmussen¹, Gustav Roder¹, Ida Dalgaard Pedersen¹, Mikael Sørensen², Morten Nielsen² and Søren Buus¹

¹ Laboratory of Experimental Immunology, Faculty of Health Sciences, University of Copenhagen, Denmark

² 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.



Prediction on the stability

NetMHCstab: predicting stability of pMHC-I complexes

Immunology
The Journal of Cells, Molecules, Systems and Technology
IMMUNOLOGY ORIGINAL ARTICLE

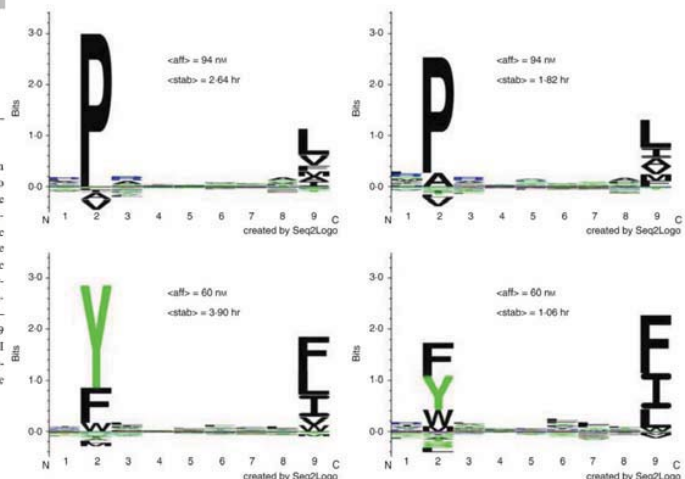
NetMHCstab – predicting stability of peptide-MHC-I complexes; impacts for cytotoxic T lymphocyte epitope discovery

Kasper W. Jørgensen,^{1,*} Michael Rasmussen,^{2,*} Søren Buus² and Morten Nielsen^{1,3}

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

Summary

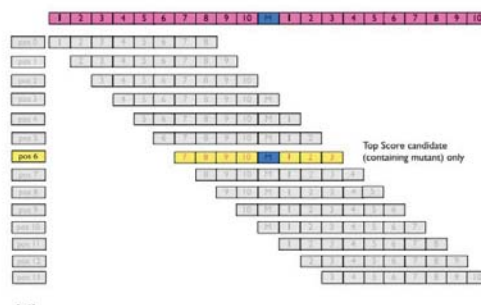
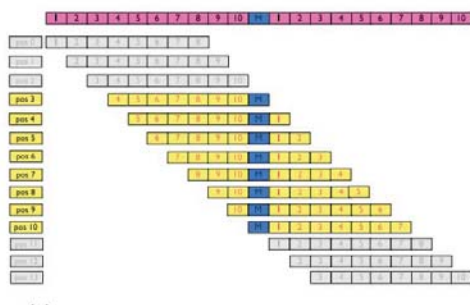
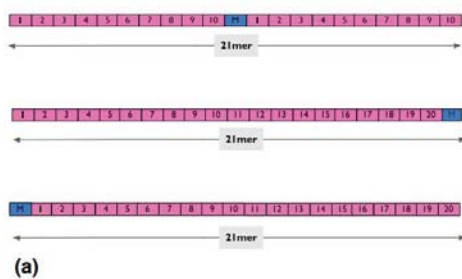
Major histocompatibility complex class I (MHC-I) molecules play an essential role in the cellular immune response, presenting peptides to cytotoxic T lymphocytes (CTLs) allowing the immune system to scrutinize ongoing intracellular production of proteins. In the early 1990s, immunogenicity and stability 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



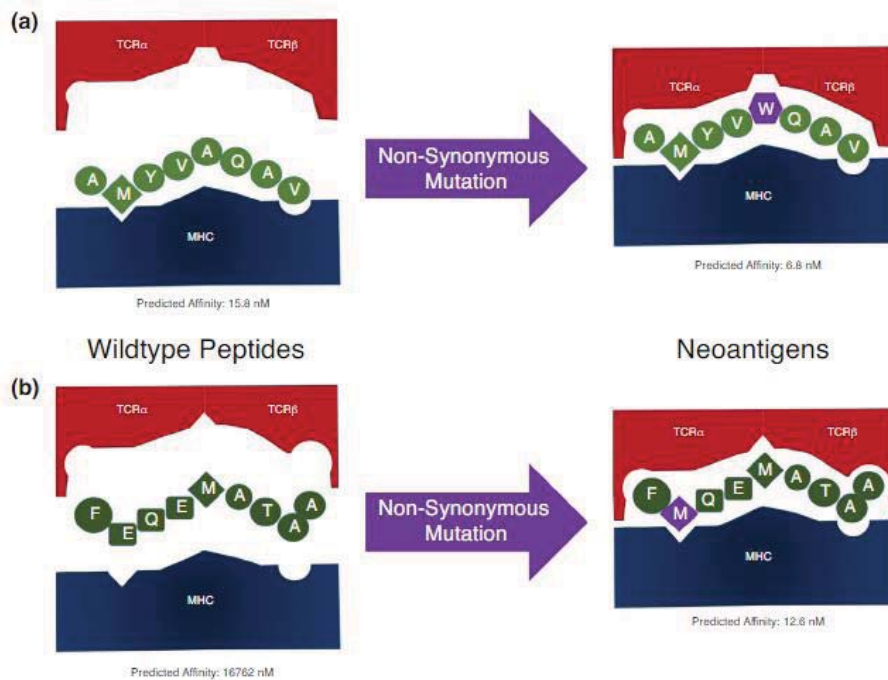
stable

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

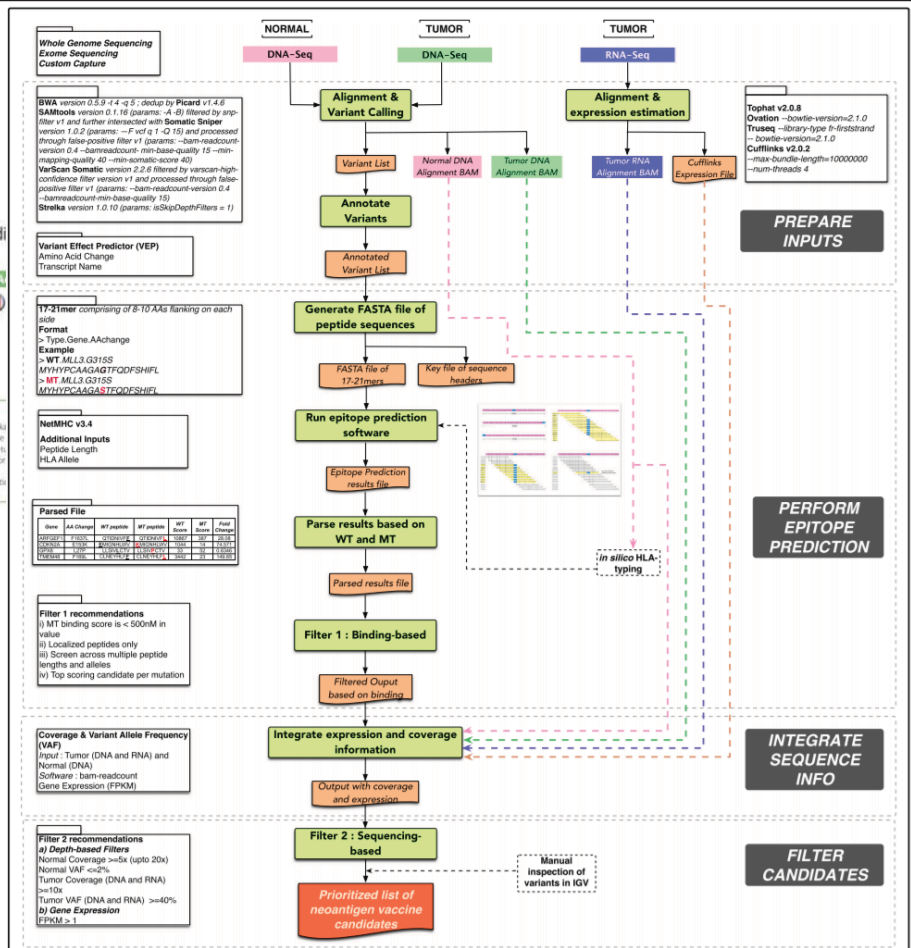
Genome Medi

METHOD Open Access

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

Jazreet Hundal¹, Beatriz M. Camero², Allegra A. Presti¹, Gerald P. Linette^{1,2,3,4}, Ebaner R. Mardis^{1,2,3,4,5} and Malachi Griffith^{1,2,3,4}

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 Antigen (pVAC-Seq) that integrates tumor mutational 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

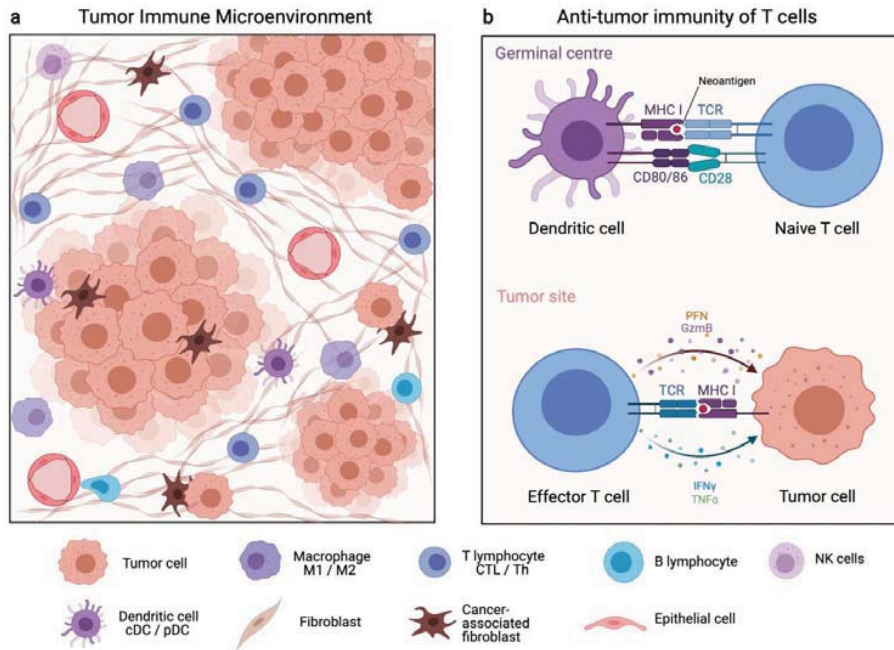
- Is MHC-I binding affinity the only applicable feature?
- Is IC_{50} under 50nM (or 500nM) an acceptable cut-off?

Discovery of new features

- Can we find a new feature for immunogenicity prediction?

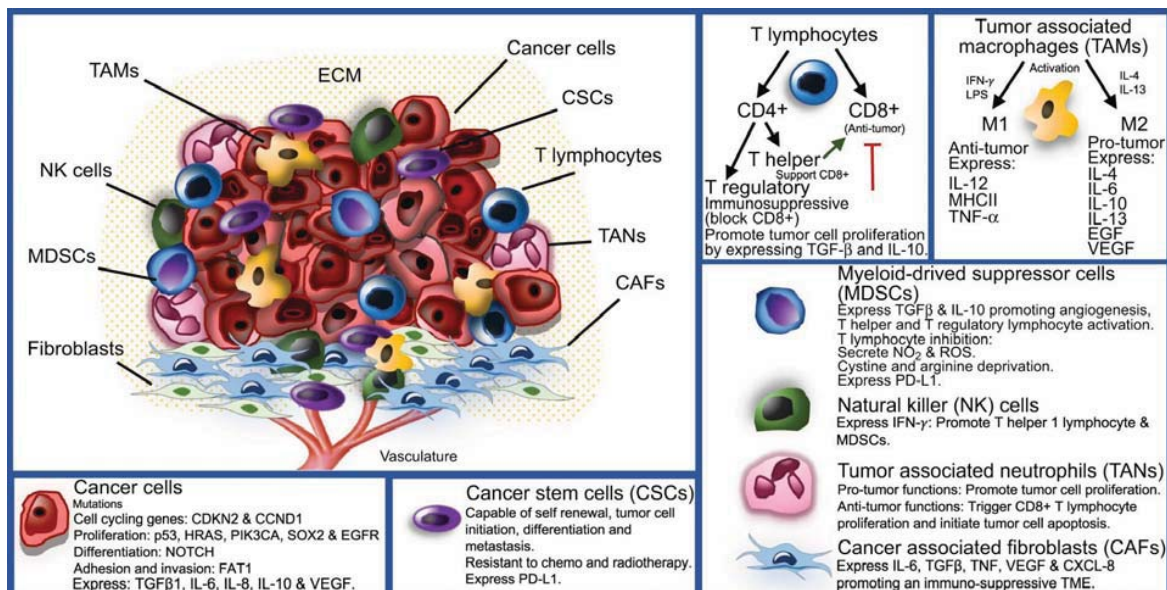
IDENTIFYING TUMOR IMMUNE MICROENVIRONMENT

Tumor Immune Microenvironment



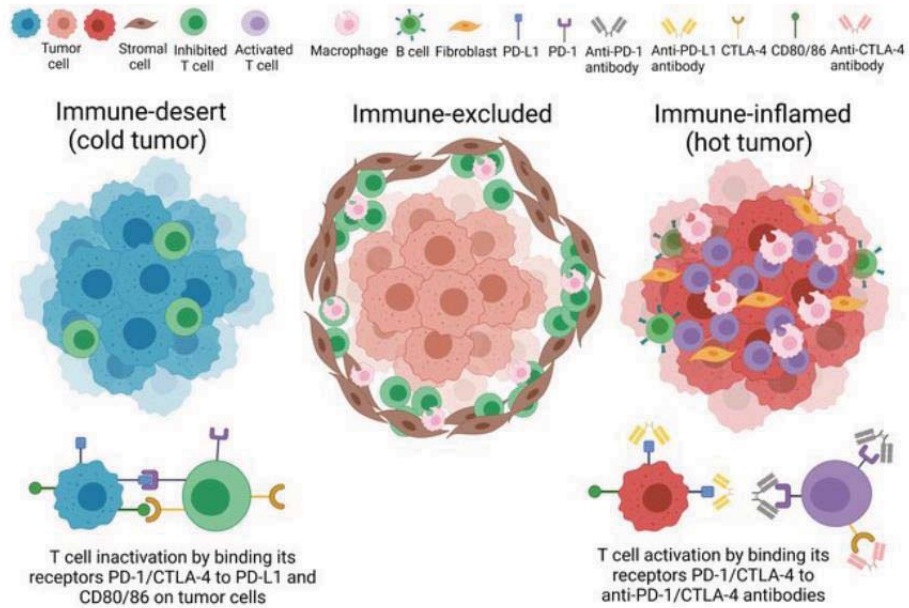
- A complex, organ-like structure (tumor cells, immune cells, fibroblasts, vascular endothelial cells, and other stromal cells)
- Immune cells + secreted factors (cytokines, chemokines, growth factors)

Tumor Immune Microenvironment



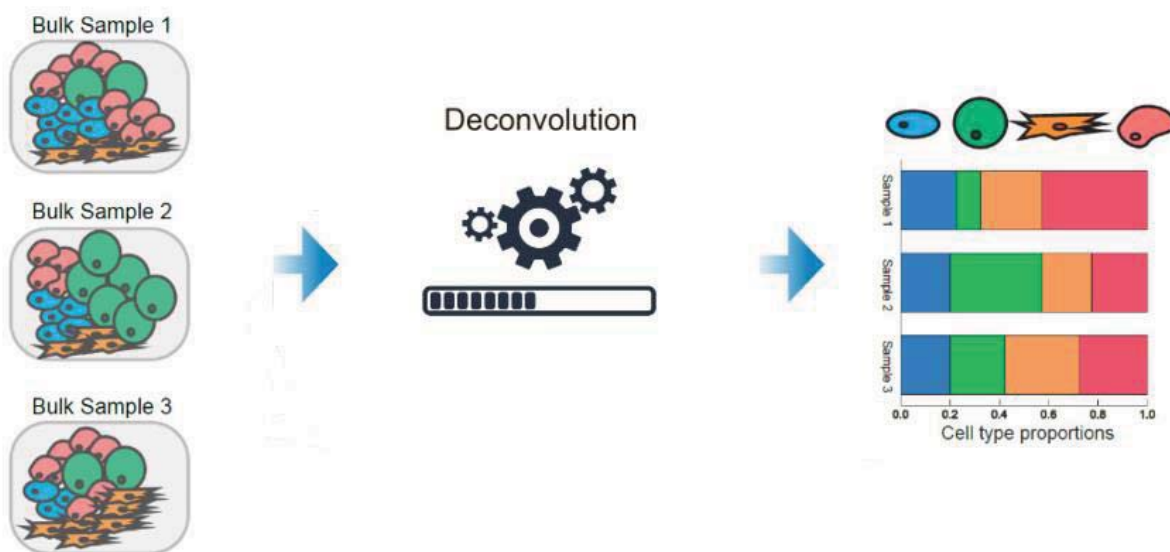
- TME components often inhibit or promote anti-tumor immunity
- But their roles are not definitive, and can be context-specific

Tumor Immune Phenotype

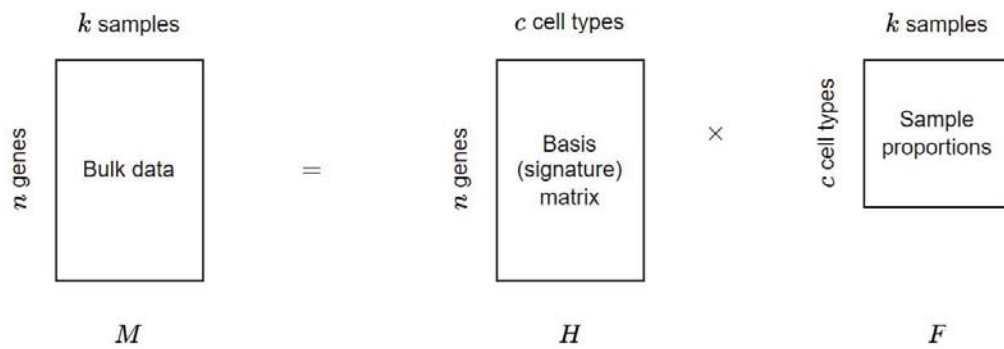


- Immune inflamed: immune cells infiltrated the tumor
- Immune excluded: immune cells are restricted to the stroma
- immune desert: T cells are not recruited

Identifying TME by Cell-type decomposition

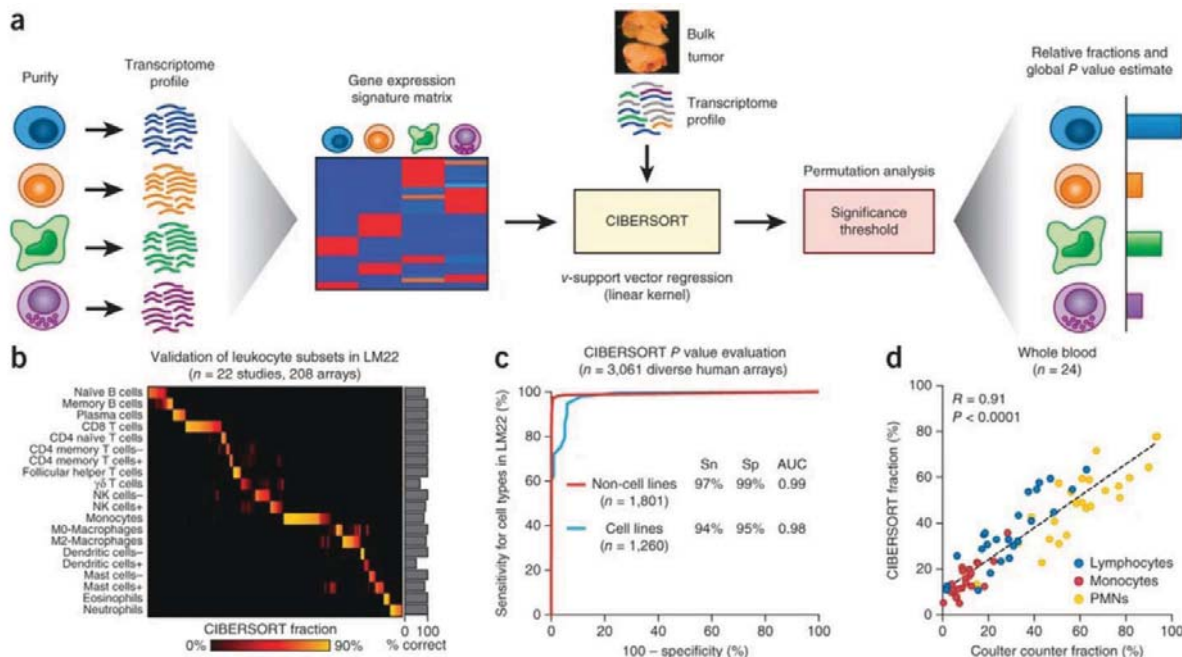


Identifying TME by Cell-type decomposition



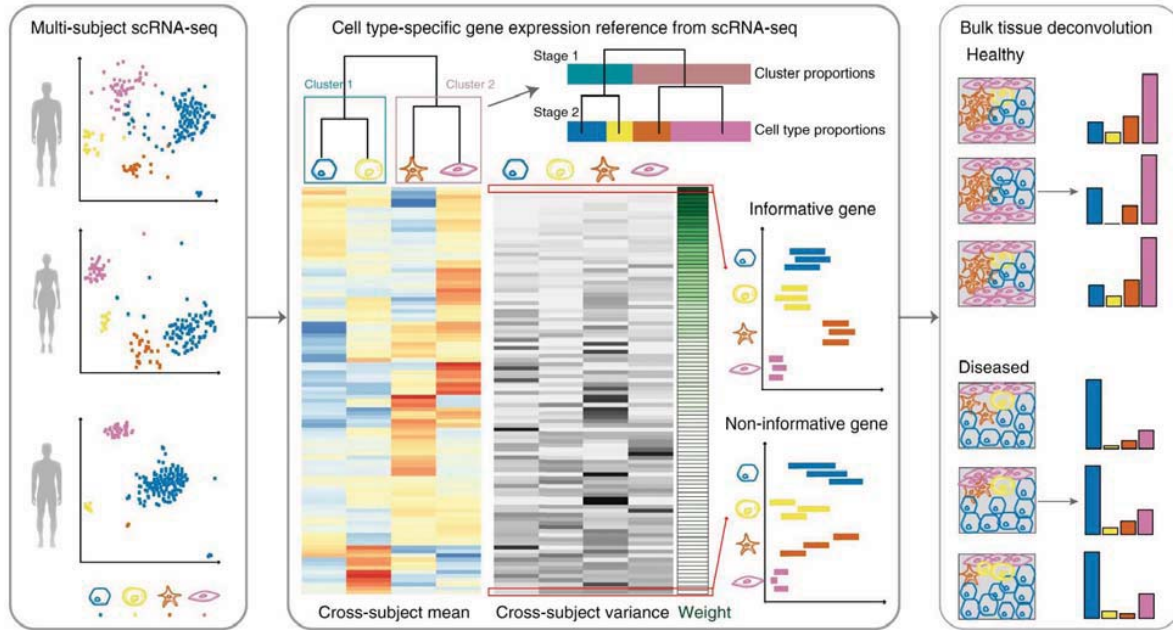
Problem	Given	Estimate	Requires
Estimate cell type proportions from bulk profile and signature matrix	M, H	F	$n > c$
Generate signature matrix from bulk profile and known cell type proportions	M, F	H	$k > c$
Estimate bulk profile from signature matrix and cell type proportions	H, F	M	none

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	DeconRNASeq	DCQ	NNLS	dtangle	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	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	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	Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended
Sturm, G. et al., 2019	Recommended	Not Recommended	Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended	Not Recommended

Recommended (Red), Not Recommended (Black), Not Evaluated (White)

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](#) & [Marius List](#)

Protocol | [First Online: 09 March 2020](#)

5035 Accesses | 88 Citations | 1 Altmetric

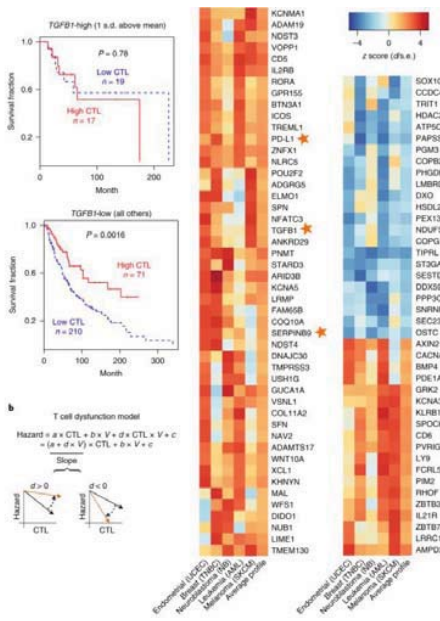
Part of the [Methods in Molecular Biology](#) book series (MIMB volume 2120)

Abstract

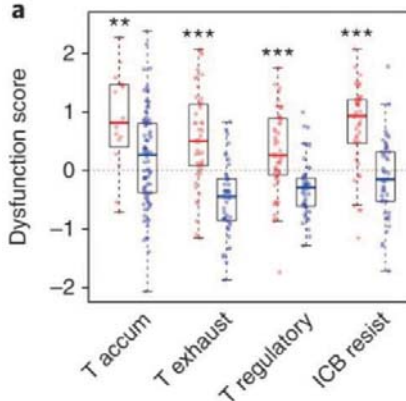
Since the performance of in silico approaches for estimating immune-cell fractions from bulk RNA-seq data can vary, it is often advisable to compare results of several methods. Given numerous dependencies and differences in input and output format of the various computational methods, comparative analyses can become quite complex. This motivated us to develop *immunedconv*, 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 *immunedconv* 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 *immunedconv* results can be summarized graphically.

- Each tool has its own pros and cons, and do not agree each other
- Immunedconv 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)

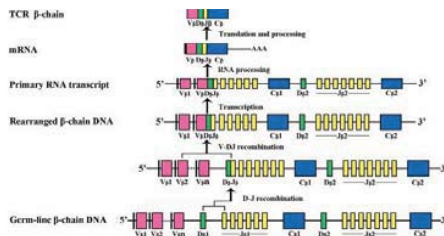
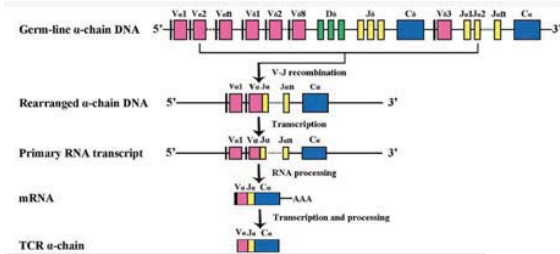
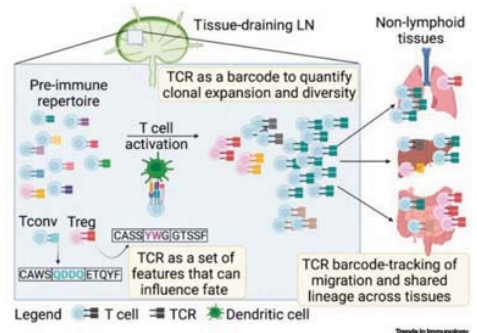
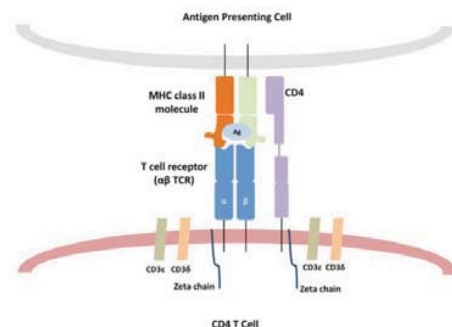


TGFβ is interacting with CTL infiltration because:
 - Survival of High vs. Low CTL-infiltrated patients are discriminated only when TGFβ is highly expressed



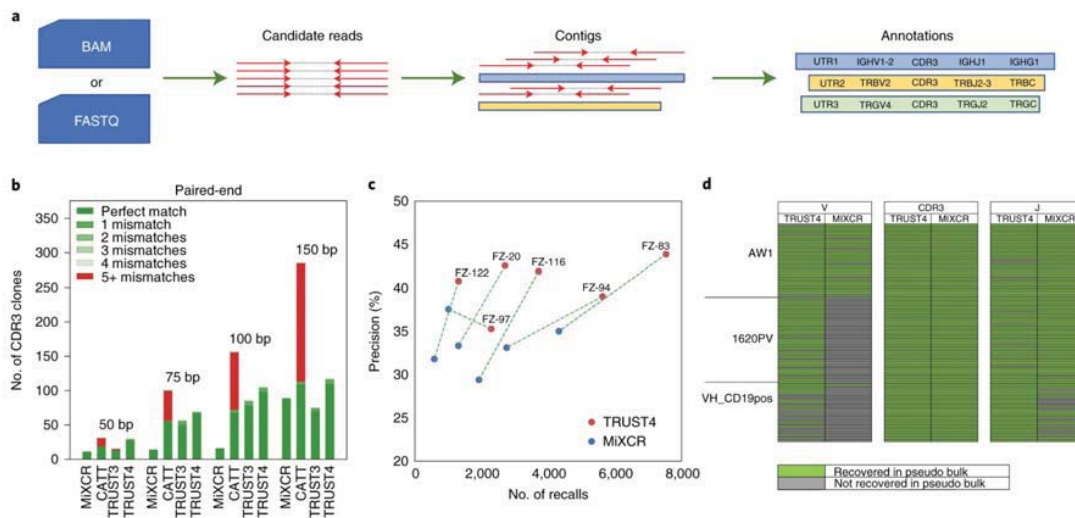
- Predicting T-cell dysfunction model: high infiltration but dysfunctional T-cells or excluded T-cells
- Extract T-cell dysfunction genes from interaction test in treatment naïve data
- Calculate T-cell dysfunction score

TCR repertoire



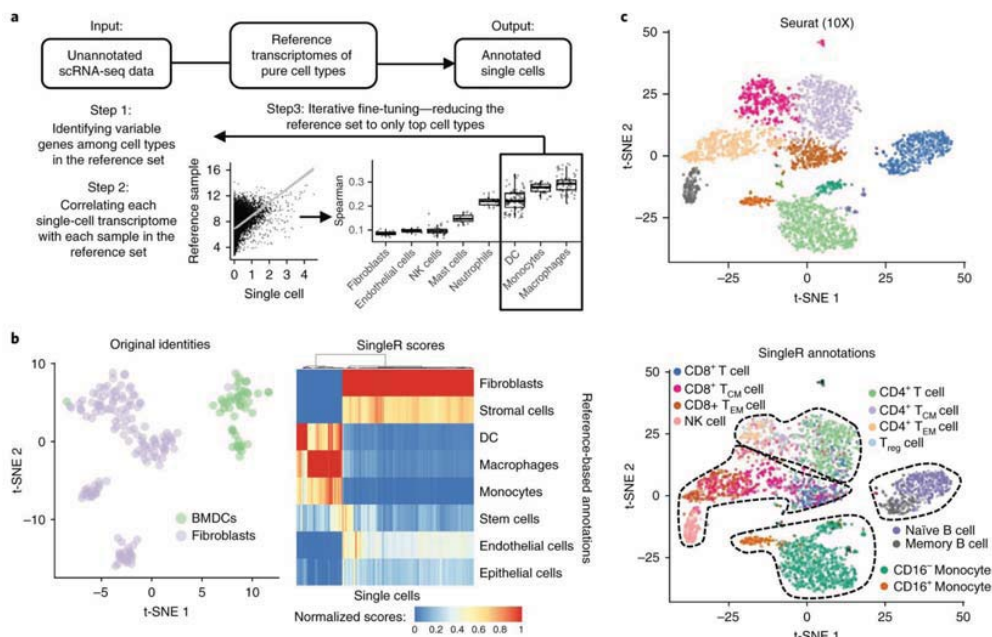
- T-cell diversity and clonality is the overall resultant response to the complex T-cell immune environment
- Diversity is inversely related to clonality
- High clonality is generally a marker for good response

TCR repertoire reconstruction



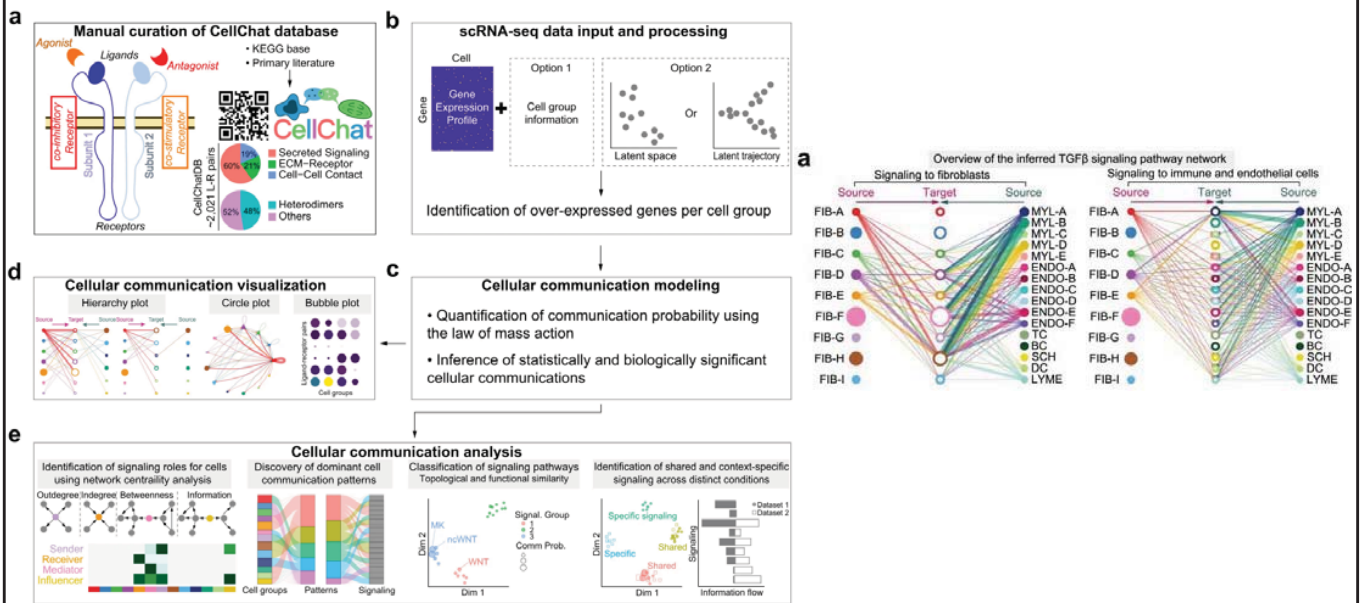
- Generally, TCR or BCR sequencing is employed for repertoire reconstruction
- Conventional bulk RNA-seq can be also used using specialized tools, such as TRUST
- TRUST 1) extract TCR/BCR candidate reads, 2) assemble to form contigs, 3) identify somatic hypermutations, 4) reconstruct repertoire

Use of single- and spatial transcriptomics



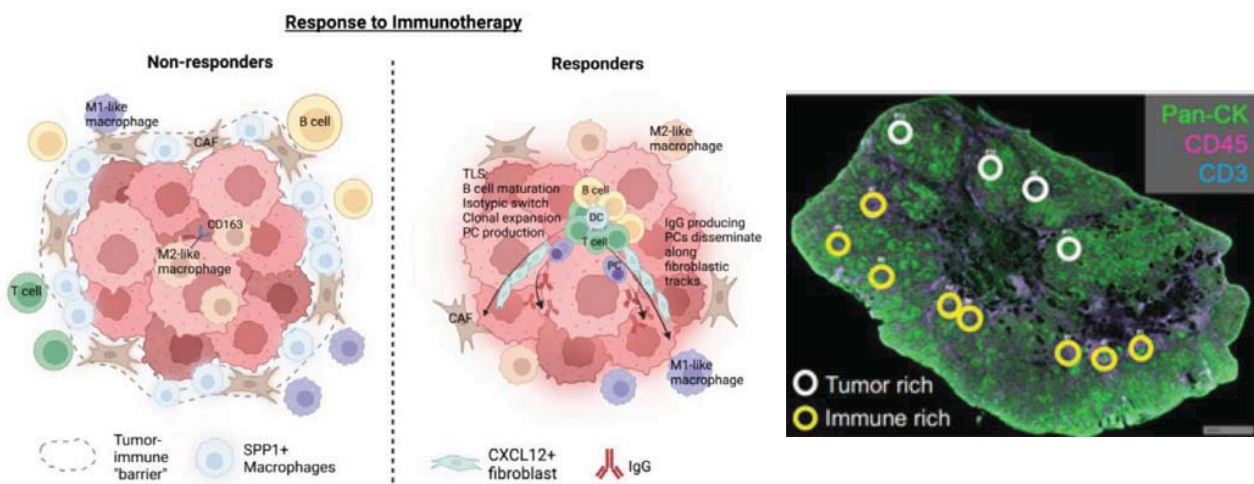
- In single cell sequencing, a complex decomposition is not necessary once the single cells are well clustered.
- Clusters should be annotated using reference gene expression

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



- Using spatial transcriptomics, 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

Thank you

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



2024 BIML 실습에 필요한 파일 다운로드 링크

<https://onedrive.live.com/?authkey=%21AF8kWqyBkTS6jqQ&id=B6775C185E600E18%211225&cid=B6775C185E600E18>

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/jyhgong906/BIML_2024/Bulk_WES/Data

실습 스크립트: /home/jyhgong906/BIML_2024/Bulk_WES/Script

환경 변수 설정

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

# PATH #
HLA_PATH=/home/jyhgong906/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/jyhgong906/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>

VEP (Variant Effect Predictor)

The screenshot shows the Ensembl Variant Effect Predictor (VEP) website. At the top, there's a navigation bar with 'BLAST/BLAT | VEP | Tools | Biomart | Downloads | Help & Docs | Blog'. Below that, a search bar is present. The main content area is titled 'Ensembl Variant Effect Predictor (VEP)' and includes a brief description: 'VEP determines the effect of your variants (SNPs, insertions, deletions, CNVs or structural variants) on genes, transcripts, and protein sequence, as well as regulatory regions.' It lists key features like 'Genes and Transcripts affected by the variants', 'Location of the variants', and 'Consequence of your variants'. There are also links for 'What's new in release 111?' and 'VEP interfaces' which includes 'Web interface', 'Command line tool', and 'REST API'.

- 변이의 종류에 따른 영향(예: missense, nonsense, frameshift 등)을 포함한 상세한 annotation을 제공한다.
- 유전체 데이터에서 발견된 변이의 기능적 영향을 분석하고, 이 정보를 바탕으로 생물학적 해석을 가능하게 함.
- VEP를 통한 변이 annotation은 변이가 단백질에 미치는 영향을 이해함으로써, potent neoantigen을 예측할 수 있음.

VEP annotation

```

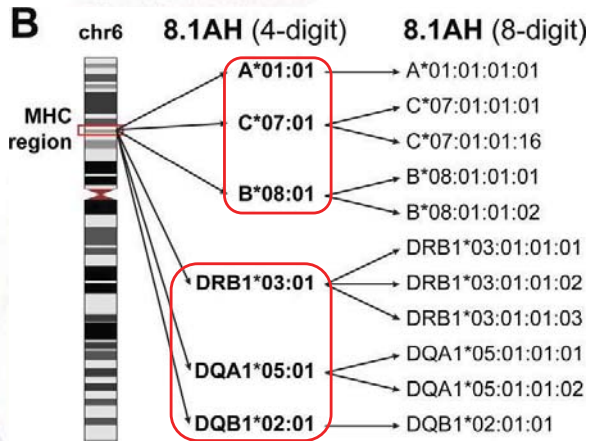
$ {vep_run} \
-i ${HLA_PATH}/${patient_id}/${vcf_format} \
-o ${HLA_PATH}/${patient_id}/${ann_format} \
--vcf \
--symbol \
--terms SO \
--tsl \
--hgvs \
--fasta ${ref} \
--force_overwrite \
--assembly GRCh38 \
--plugin Wildtype \
--plugin Frameshift \
--offline \
--cache \
--dir_cache /data/public/VEP/104 \
--dir_plugins /data/public/VEP/104/Plugins \
--pick \
--transcript_version \
    
```

```

chr9 5988439 0 0 C PASS AS_FilterStatus:SITE;AS_SB_TABLE:106,183|33
27:DP:274;ECNT:1;GERM:93;MBD:29;37:MFRL:223,231;MMD:60,60;MPOS:34;NALOD:-4,274e-01;NL00D:-4,86;POP
AF:5,60;TL00:170,22;CSQ:C[m]issense_variant[Moderate][KIAA2026][ENSG00000183354][Transcript][ENST0000029
9933.0][protein_coding][2/8][ENST00000399933.0:c.700C>G][ENSP0000032815.3;p.Arg234Gly][1395][700][234][R
][G/Ga][Gaa][1][1][HGNC][HGNC:23378][5][HSPVGP:GAMEPAGEERPPPAEGEEDDEEVAAMAAATSGPAPHRSASSLEADODEEEMAM
TGGCCGCELELTYE000RLLIGEELOEKHRLTAPFLPGLGVATAEVEAEQPSGRGGRAPFQPPGKLLQWEKFA56QVGGITFVADFLR
WETTCYRLHGVHMTSAGKQLEMLKOKLALSRHLREKTTAVTSRSRYVLEDEKGTACTSTGRTSPRLSAGTSQVSTFTWVLRQEEJAME
EKLRIGERQEKAEASQKLEEMKRLAQAPMTCMEIMEIPAIGHFLCAOILNLPVEIVFYELERLCLLPQOMAFSLHMTSLSPRRHPTLHR
PITLRYRTEAALRQKVOQMYTAVGQTEPNQCAEKLLCPQFKVLEWNPLEKPFHELPYQKWLKRLGDFVYEQEVDPAIDPCACPKMLDNDID
LGYDVENAVYHPQFGADVRYKORPQAEFPFPDPTIKORVPRIKLEKLCQVYVSYNGEHRCSRDSLSSFKKCEVSPKNSGQPKAKMLDNDID
SVENQVSYVEIRIRRRPCEIKKTCCKENLEKPRSPGVTGQEPSPQELRLENOYKGEARLXKPELNTKSCDIHNGSHSDIENLQK
VARDLLEGLSFRKRLKLLPFAKRRKXKXKLEKNEKLELRWREKLELSAKSYRPELQKLLIKKRRKRRKRSKRSKRAKTKRKYLVKS
TFPEPFLICTMLDELRELTKTENELKDLNSRKSQKRYHRQAKWELHSTLIRLLELNPWKMLKAFQRNRSRSLKXQDVRPQDPDHFREL
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPKSISKSDTEPLDLEKDFSDDMKLEIDFPMARSKLLKELPKSKDLPKTLTKLKRQKQTYVDQDSKELSPRKKALK
STNEITVENLESVDIDCFSEKHPTEPSPFASLDSVPSVTLQKQTPQIALLAKNIGKVTLNQPLSPITGRNALAVEKPLSPPEASIPKALTC
HINTKGPLQMVYKPCQQLPLDLONSWIKIQVPHVQDKTEKIMQVLLPKNIFVIOHKEGKAVKVEPQDKQTEHQSSFPDITLHNSIASVF
WISQYVETLQPTPTFRKTLTSLISSARPPQISPVPSVSNLTPSWKTSQSEAGKAVNSVANSFSPASQSTFSTTQPLSSITLHCTMCTGQSPF
HFAQDITADSEAKKELKTCYRDSLWRTGKQDTPWPKNSPWTSSVSPFAPQLEWTFSSPSTSSSPLSPFATITLSSLSQSP
PFCVTDPSAPFSMKMLKLLGHTGSSGLGHVDKTSHPSPKSSLCSSITLPSSTSSVSISSAANGONNMLHPTK000VYDTSYSP
VTRSEATAATNDVTSQPKVLMVSAQSLSSGNGATNMPALSTQVSAQKLVINAPVSPSTLTVAESLQKTLPPHAKVAKVTEPOPTO
LPSVVGTPKINSSPASVSDIKVIGLNLGQAVITVSGVPAIPSNLNLQVTPKEDSKSKYLLPLLSGNSVPSVSNVSNLQNTVSNVSVAR
AVWLVYTGANLGLSPVYASAGAAPVPLVSDNSTRIMPILSNRCSLSSGLTVAITWGTGLASVLLSTTPWPKRLCALOITVPTVVA
LPTTATTSKATLITMPSANVGRATRSVLSKRSRBLDQPSHPTSTVPTMTWPKOTELSSLSSTSPKTKTTSNPSASLPMOJALWTRPYSAPAP
WTFDEGDDIQAQSPKPK
```

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.

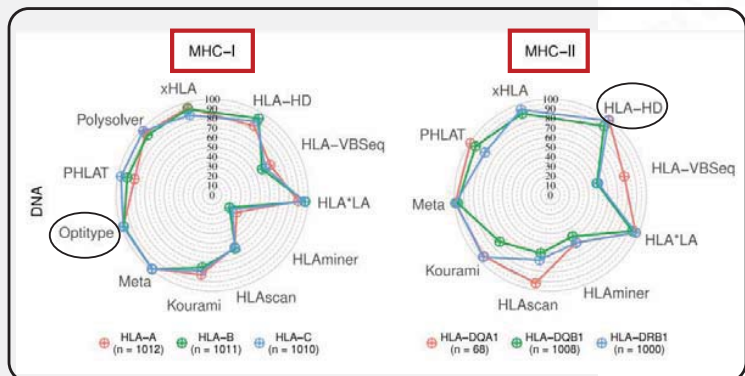
HLA typing (MHC class I, II)

```
# HLA typing - MHC class I (OptiType) #
python2 ${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 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}
```

```
[jyhong906@master ACC_T_01]$ pwd
/home/jyhong906/BIML_2024/Bulk_WES/Data/OptiType/ACC_T_01
[jyhong906@master ACC_T_01]$ ll
total 1132
-rw-r--r-- 1 jyhong906 jyhong906 1154428 Feb  8 17:26 ACC_T_01_coverage_plot.pdf
-rw-r--r-- 1 jyhong906 jyhong906  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	2732.2449999999956

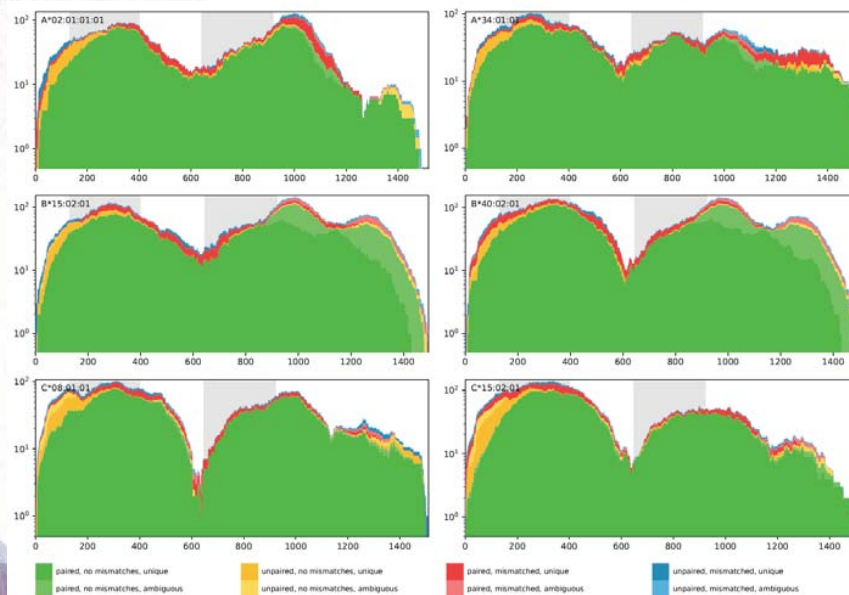
MHC class I – ACC_T_01.MHC.I.list.txt

```
A-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)

```
[jyhong96@master result]$ pwd
/home/jyhong96/BIMB_2024/Box_MES/Data/HLA-HD/ACC_T_01/result
[jyhong96@master result]$ ll
total 1676
-rw-r--r-- 1 jyhong96 jyhong96 4711 Feb 9 17:42 ACC_T_01_A.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 117230 Feb 9 17:42 ACC_T_01_A.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 965 Feb 9 17:42 ACC_T_01_B.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 85393 Feb 9 17:42 ACC_T_01_B.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 156 Feb 9 17:42 ACC_T_01_C.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 85495 Feb 9 17:42 ACC_T_01_C.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 268 Feb 9 17:42 ACC_T_01_DMA.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 59235 Feb 9 17:42 ACC_T_01_DMA.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 899 Feb 9 17:42 ACC_T_01_DMB.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 54183 Feb 9 17:42 ACC_T_01_DMB.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 376 Feb 9 17:43 ACC_T_01_DPA1.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 64783 Feb 9 17:43 ACC_T_01_DPA1.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 619 Feb 9 17:41 ACC_T_01_DPB.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 49139 Feb 9 17:41 ACC_T_01_DPB.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 1949 Feb 9 17:42 ACC_T_01_DPA2.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 45288 Feb 9 17:42 ACC_T_01_DPA2.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 249 Feb 9 17:43 ACC_T_01_DPA2.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 3220 Feb 9 17:43 ACC_T_01_DPA2.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 2924 Feb 9 17:41 ACC_T_01_DPB1.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 58668 Feb 9 17:41 ACC_T_01_DPB1.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 748 Feb 9 17:42 ACC_T_01_DPA1.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 75414 Feb 9 17:42 ACC_T_01_DPA1.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 1568 Feb 9 17:42 ACC_T_01_DPB1.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 47342 Feb 9 17:42 ACC_T_01_DPB1.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 697 Feb 9 17:42 ACC_T_01_DRA.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 109544 Feb 9 17:42 ACC_T_01_DRA.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 433 Feb 9 17:40 ACC_T_01_DRB1.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 33356 Feb 9 17:40 ACC_T_01_DRB1.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 182 Feb 9 17:43 ACC_T_01_DRB2.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 1159 Feb 9 17:43 ACC_T_01_DRB2.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 181 Feb 9 17:43 ACC_T_01_DRB3.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 12652 Feb 9 17:43 ACC_T_01_DRB3.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 0 Feb 9 17:41 ACC_T_01_DRB4.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 185 Feb 9 17:43 ACC_T_01_DRB5.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 17365 Feb 9 17:43 ACC_T_01_DRB5.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 73 Feb 9 17:43 ACC_T_01_DRB6.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 2616 Feb 9 17:43 ACC_T_01_DRB6.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 14 Feb 9 17:43 ACC_T_01_DRB7.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 0 Feb 9 17:41 ACC_T_01_DRB7.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 14 Feb 9 17:41 ACC_T_01_DRB8.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 0 Feb 9 17:41 ACC_T_01_DRB8.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 80 Feb 9 17:43 ACC_T_01_DRB9.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 7838 Feb 9 17:43 ACC_T_01_DRB9.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 1080 Feb 9 17:43 ACC_T_01_E.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 89401 Feb 9 17:43 ACC_T_01_E.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 494 Feb 9 17:42 ACC_T_01_F.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 550 Feb 9 17:44 ACC_T_01_final.result.txt
-rw-r--r-- 1 jyhong96 jyhong96 56222 Feb 9 17:42 ACC_T_01_F.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 685 Feb 9 17:43 ACC_T_01_G.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 75325 Feb 9 17:43 ACC_T_01_G.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 219 Feb 9 17:43 ACC_T_01_H.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 100146 Feb 9 17:43 ACC_T_01_H.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 80 Feb 9 17:43 ACC_T_01_I.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 87786 Feb 9 17:44 ACC_T_01_J.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 92 Feb 9 17:42 ACC_T_01_K.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 31139 Feb 9 17:42 ACC_T_01_L.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 95 Feb 9 17:43 ACC_T_01_L.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 40231 Feb 9 17:43 ACC_T_01_L.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 222 Feb 9 17:44 ACC_T_01_T.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 11897 Feb 9 17:44 ACC_T_01_T.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 118 Feb 9 17:44 ACC_T_01_V.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 27658 Feb 9 17:44 ACC_T_01_V.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 193 Feb 9 17:42 ACC_T_01_W.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 5123 Feb 9 17:42 ACC_T_01_W.read.txt
-rw-r--r-- 1 jyhong96 jyhong96 343 Feb 9 17:44 ACC_T_01_Y.est.txt
-rw-r--r-- 1 jyhong96 jyhong96 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)
```

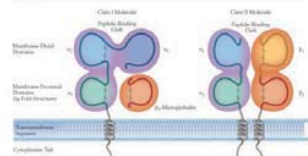
MHC class II - ACC_T_01_final.result.txt

```
D
R
C
DRB1
DQA1
DQB1
DPA1
DPB1
DMA
DMB
DQA
DQB
DRA
DRB2
DRB3
DRB4
DRB5
DRB7
DRB8
DRB9
DPA2
E
F
G
H
J
K
L
T
V
W
Y
HLA-A*02:01:01 HLA-A*34:01:01
HLA-B*48:02:01 HLA-B*15:02:01
HLA-C*15:02:01 HLA-C*08:01:01
HLA-DRB1*15:02:01 HLA-DRB1*12:02:01
HLA-DQA1*01:02:01 HLA-DQA1*06:01:01
HLA-DQB1*03:01:01 HLA-DQB1*05:02:01
HLA-DPA1*01:03:01 HLA-DPA1*02:02:02
HLA-DPB1*02:01:02 HLA-DPB1*01:01:01
HLA-DMA*01:02:01 -
HLA-DMB*01:01:01 -
HLA-DQA*01:01:04 HLA-DQA*01:01:01
HLA-DQB*01:01:01 HLA-DQB*01:01:01
HLA-DRA*01:02:02 HLA-DRA*01:01:01
HLA-DRB2*01:01 -
HLA-DRB3*03:01:03 -
DRB4 Not typed Not typed
HLA-DRB5*01:01:01 -
HLA-DRB6*02:01 -
DRB7 Not typed Not typed
DRB8 Not typed Not typed
HLA-DRB9*01:02:01 -
HLA-DPA2*01:01:02 HLA-DPA2*02:01:02
HLA-E*01:03:01 HLA-E*01:03:02
HLA-F*03:01:01 -
HLA-G*01:01:01 HLA-G*01:01:03
HLA-H*01:01:01 HLA-H*02:27
J HLA-J*01:01:01 HLA-J*01:01:01
K HLA-K*01:02 -
L HLA-L*01:02 -
T HLA-T*03:01:01 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,DQB1*03:01,DQA1*06:01,DQB1*05:02
```

Differences between MHC Class I and MHC Class II



pVACseq (neoantigen prediction)

```
# pVACseq - MHC class I (NetMHCpan) #
allele='cat $(opttype_PATH)/$(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) \
--iedb-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) \
--iedb-install-directory $(IEDB_MHCII)
```

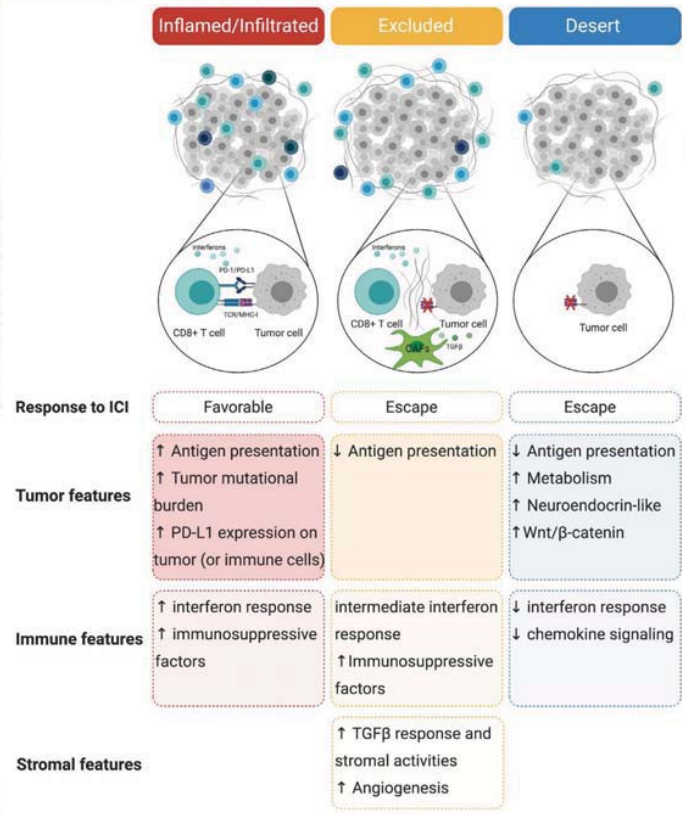
MHC class I - ACC_T_01.filtered.tsv

Chromosome	Start	Stop	Reference	Variant	Transcript	Transcript Support Level	Transcript Length	Biotype	Ensembl	Gene ID	Variant Type	Mutation	Protein Position							
chr6	7884402	7884403	C	T	ENST00000379757.9	1	432	protein_coding	ENSG00000292664	missense	E/K	378	TXND5C	ENST00000379757.9:c.11326G>A	ENSP00000369081.4:p.G1					
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	
ts	False	0	False	0	False	0	False	0	False	0	False	0	False	0	False	0	False	0	False	0

MHC class II - ACC_T_01.filtered.tsv

Chromosome	Start	Stop	Reference	Variant	Transcript	Transcript Support Level	Transcript Length	Biotype	Ensembl	Gene ID	Variant Type	Mutation	Protein Position							
chr6	7884402	7884403	C	T	ENST00000379757.9	1	432	protein_coding	ENSG00000292664	missense	E/K	378	TXND5C	ENST00000379757.9:c.11326G>A	ENSP00000369081.4:p.G1					
u378lys	HLA-DQA1*01:02/DQA1*03:01	16	3	14	KKEFPLAGVKIAKV	KKEFPLAGVKIAEV	NetMHCIIpan	375.12	373.19	0.995	NetMHCIIpan	2.7	2.6	136	0.365	NA	NA	NA	NA	
ts	False	0	False	0	False	0	False	0	False	0	False	0	False	0	False	0	False	0	False	0

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)
```



	ACC_T_01	ACC_T_02	ACC_T_03	ACC_T_04	ACC_T_05	ACC_T_06
A1BG	8.818338e+01	1.198802e+01	1.740194e+00	1.715061e+00	3.228799e-01	4.781870e-01
A1BG-AS1	3.672197e+00	3.738398e+00	4.397913e+00	6.015022e+00	3.514066e+00	3.898722e+00
A2M	1.355442e+01	4.722011e+02	3.209852e+03	6.858230e+02	1.632307e+03	1.699864e+03
A2M-AS1	1.064123e+01	9.269908e+00	1.445803e+00	2.594986e+00	4.641643e+00	4.817001e+00
ADHL1	2.559659e+00	6.410235e-02	2.060730e+01	4.204261e+00	1.992093e+01	2.370282e+00
ADHL1-2	1.029029e+00	8.979255e-01	8.810581e-01	4.436115e-01	1.992157e+00	5.781857e-01
ADGALT	1.379352e+00	8.890345e-01	7.178419e+01	4.843958e+01	2.676064e+00	1.468356e+00
AAAS	2.211448e+01	3.813641e+01	1.042619e+01	1.312209e+01	3.434816e+01	2.119581e+01
AACS	4.047587e+00	6.734105e+00	1.282836e+01	4.972145e+00	7.884083e+00	7.189079e+00
AACSP1	1.709853e+02	1.732525e+01	1.780513e+01	7.902435e+00	1.694886e+02	1.049168e+02
AADAT	6.649355e+01	1.523508e+02	4.424247e+00	8.832889e+00	3.314472e+01	5.990574e+01
AGAB	4.843727e+01	5.716138e+01	2.498052e+01	3.258767e+01	7.531622e+01	6.415446e+01
AAK1	6.028251e+01	2.060464e+01	1.052317e+01	2.832064e+01	5.113273e+01	6.736277e+01
AAOC	2.037598e+01	1.420794e+01	1.032013e+01	1.083977e+01	1.921214e+01	9.252845e+00
AAMP	7.387827e+01	5.991051e+01	8.170519e+01	1.143184e+02	7.019441e+01	7.876269e+01
AAO2	4.404381e+01	4.147078e+01	1.561102e+01	5.613361e+01	4.643148e+01	4.299489e+01
AAO3	1.004734e+01	1.253398e+01	2.893168e+01	6.459862e+01	1.642113e+01	8.409363e+01
AAOS2	2.682950e+01	3.728108e+01	1.591450e+01	2.954066e+01	2.653228e+01	4.278427e+01
AAOSD1	3.296782e+00	2.541961e+00	2.502514e+00	1.855305e+00	3.253168e+00	3.301350e+00
AAOSD	6.604800e+01	9.405322e+01	1.468431e+01	2.318409e+01	6.964260e+01	7.562925e+01
AAOSPP1	7.508998e+01	7.245487e+01	3.527344e+01	3.131916e+01	6.760113e+01	4.198889e+01
AAOS	1.847466e+02	1.802919e+02	1.599966e+01	1.155745e+01	1.139355e+02	5.821509e+01
AAICR	1.252096e+00	2.073114e+00	1.774436e+00	6.720870e+00	4.776974e+00	9.712573e-01
AATF	3.050248e+01	3.270385e+01	2.870022e+01	2.317741e+01	2.952149e+01	2.831465e+01
AATK	9.388976e-02	2.154731e-01	4.928081e-01	1.026395e+00	4.098307e-01	1.612430e-01
AAT	1.995402e+02	3.782676e+01	3.670583e+00	1.678356e+00	7.722701e+00	5.595163e+00
AICA1	2.699329e+02	1.176262e+02	2.447518e+01	3.352699e+01	1.362989e+02	7.617583e+01
AICA10	2.123992e+01	1.581991e+01	3.396801e-01	2.780478e+00	4.086010e+00	6.709590e+00
AICA11P	1.501424e+01	2.589371e-01	3.363030e+00	6.891310e+00	1.909869e+01	2.353918e+01
AICA12	5.03625e-01	5.971615e-01	6.773390e+00	2.474813e+00	1.571988e+00	7.180188e-01
AICA13	5.218595e+01	3.363434e+02	6.713081e+01	1.339675e+02	3.32375e+01	1.342150e+02
AICA17P	1.281842e-01	1.470647e-01	1.447764e-01	3.07719e-01	3.186664e-01	4.200203e-01
AICA2	4.991191e+00	4.200612e+00	6.044670e+00	5.405588e+00	5.909786e+00	8.635339e+00
AICA3	2.544006e+00	5.417196e+00	9.721827e+00	4.674464e+01	7.827281e+01	2.354581e+01
AICA4	9.600918e-01	2.955527e-01	2.808188e-01	4.589043e+00	5.256945e-01	5.538210e-01

Load library & normalization

```
#####
# Load packages #
#####
library(GeoTcgaData) # normalized TPM
library(immudeconv) # Cell type decomposition
library(ComplexHeatmap) # Visualization
library(ggplot2) # Visualization
library(ggpol) # 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/GeoTcgaData
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 = 1)
P_idx <- SGC_groups$Condition == "PRIMARY"; M_idx <- SGC_groups$Condition == "METASTASIS"
```

Install.packages()

Immune cell deconvolution

```

tool_df_list <- c()
for (idx in seq(length(deconvolution_methods))){
  tool <- deconvolution_methods[idx]

  if (tool %in% "cibersort_abs") {
    set_cibersort_binary('/data/project/BIML_2024/Bulk_RNA/Script/CIBERSORT.R'
  )
    set_cibersort_mat('/data/project/BIML_2024/Bulk_RNA/Script/LM22.txt')
    print(tool)
    assign(tool, as.data.frame(deconvolute(normalized_TPM, tool)))
    assign(paste0(tool, '_df'), data.frame(get(tool),
      row.names = get(tool)$cell_type)[
    ,-1])
    # saveRDS(get(tool), file = paste0("/data/project/BIML_2024/Bulk_RNA/Data
    /ImmuneDeconv/",tool, ".rds"))
  }

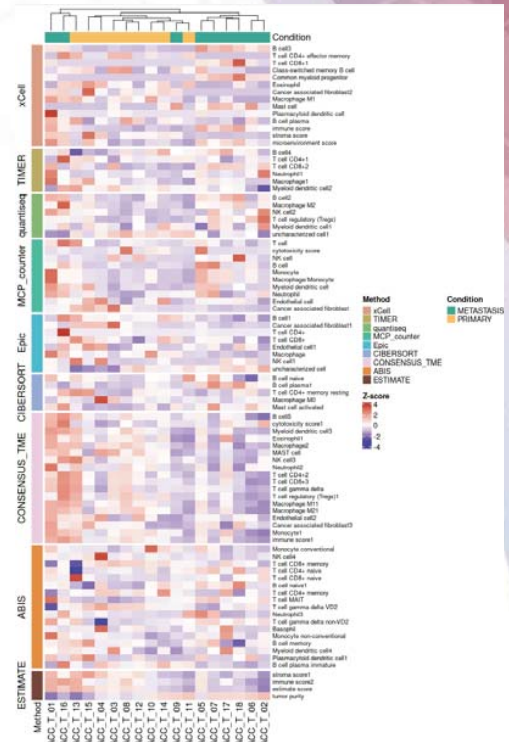
  else if (tool %in% c("timer", "consensus_tme")) {
    print(tool)
    assign(tool, deconvolute(normalized_TPM, tool,
      indications=c(rep("hnscc", ncol(normalized_TPM
    ))))
    assign(paste0(tool, '_df'), data.frame(get(tool),
      row.names = get(tool)$cell_type)[
    ,-1])
    # saveRDS(get(tool), file = paste0("/data/project/BIML_2024/Bulk_RNA/Data
    /ImmuneDeconv/",tool, ".rds"))
  }

  else if (!tool %in% c("cibersort")) {
    print(tool)
    assign(tool, as.data.frame(deconvolute(normalized_TPM, tool)))
    assign(paste0(tool, '_df'), data.frame(get(tool),
      row.names = get(tool)$cell_type)[
    ,-1])
    # saveRDS(get(tool), file = paste0("/data/project/BIML_2024/Bulk_RNA/Data
    /ImmuneDeconv/",tool, ".rds"))
  }

  tool_df_list <- append(tool_df_list, paste0(tool, '_df'))
}
# read_rds("/data/project/BIML_2024/Bulk_RNA/Data/ImmuneDeconv/normalized_TPM.rds")

for (tool_df_idx in setdiff(tool_df_list, "cibersort_df")) {
  tmp_tool_df <- get(tool_df_idx)
  tmp_tool_df <- tmp_tool_df[which(rowSums(is.na(tmp_tool_df)) != ncol(tmp_tool_df)),]
  tmp_tool_df <- tmp_tool_df[which(rowSums(tmp_tool_df == 0) < (ncol(tmp_tool_df) * 0
  .2)),]
  assign(tool_df_idx, tmp_tool_df)
}

```



Immune cell deconvolution – cell type

```

cell_type_list <- c()
for (tool in c("mcp_counter_df", "epic_df", "quantiseq_df", "xcell_df", "cibersort_abs_df",
"timer_df", "consensus_tme_df", "abis_df", "estimate_df")) {
  print(row.names(get(tool)))
  cell_type_list <- c(cell_type_list, row.names(get(tool)))
}; unique_cell_type <- unique(cell_type_list)

comb_cell_type_list <- c()
comb_p_list <- c()
for (cell_type in unique_cell_type) {
  print(cell_type)
  comb_cell_type_list <- c(comb_cell_type_list, cell_type)

  p_list <- list()
  for (tool in c("mcp_counter_df", "epic_df", "quantiseq_df", "xcell_df",
"cibersort_abs_df", "timer_df", "consensus_tme_df", "abis_df", "estimate_df")) {
    test_mat <- as.matrix(get(tool))

    if (sum(row.names(test_mat) == cell_type) >= 1) {
      test <- wilcox.test(test_mat[cell_type, p_idx],
test_mat[cell_type, M_idx])
      p_list <- c(p_list, test$p.value)
    }
  }; comb_p <- combineParallelPValues(p_list, method = "fisher"); print(comb_p$p.value)
}; comb_p_list <- c(comb_p_list, comb_p.p.value)
}; cor_idx <- mixedorder(comb_p_list, decreasing = F)
comb_cell_type_list <- comb_cell_type_list[cor_idx]; comb_p_list <- comb_p_list[cor_idx]
cell_type_df <- data.frame(cell_type = seq(length(comb_cell_type_list)),
p_value = seq(length(comb_cell_type_list)),
freq = length(comb_cell_type_list))
cell_type_df$cell_type <- comb_cell_type_list; cell_type_df$p_value <- comb_p_list; row.names
(cell_type_df) <- cell_type_df$cell_type

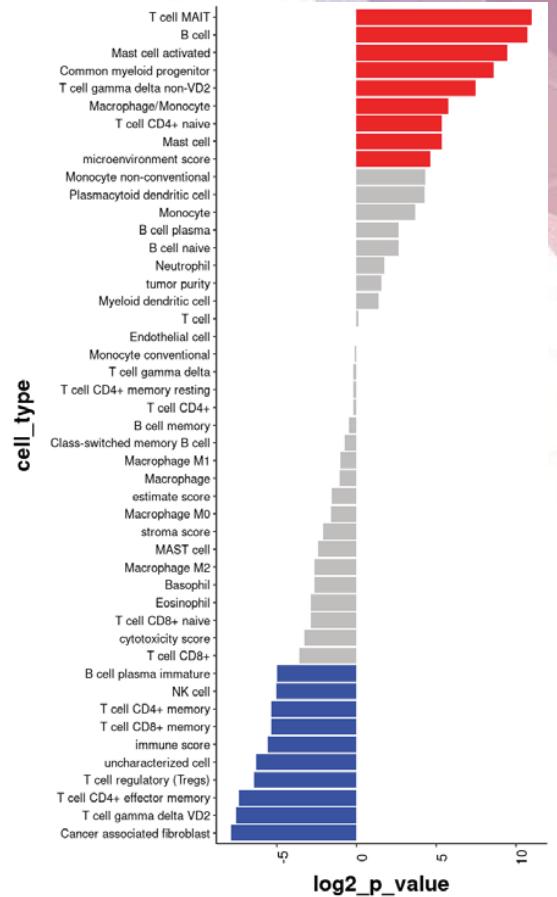
freq_list <- c()
for (cell_type in comb_cell_type_list) {
  print(cell_type)

  cell_type_line <- c()
  for (tool in c("mcp_counter_df", "epic_df", "quantiseq_df", "xcell_df",
"cibersort_abs_df",
"timer_df", "consensus_tme_df", "abis_df", "estimate_df")) {
    test_mat <- as.matrix(get(tool))

    if (sum(row.names(test_mat) == cell_type) >= 1) {
      cell_type_line <- rbind(test_mat[cell_type,], cell_type_line)
    }
  }; mean_value <- apply(cell_type_line, 2, mean)

  if (mean(mean_value[p_idx]) > mean(mean_value[M_idx])) {
    cell_type_df[cell_type,]$freq <- "primary"
  } else {
    cell_type_df[cell_type,]$freq <- "metastasis"
  }
}; cell_type_df$log2_p_value <- -log2(cell_type_df$p_value)
cell_type_df[cell_type_df$freq == "primary",]$log2_p_value <- cell_type_df[cell_type_df$freq
== "primary",]$log2_p_value * -1
cell_type_df <- cell_type_df %>% mutate(color = ifelse(abs(log2_p_value) > -log2(0.05) & freq
== "metastasis", "red",
ifelse(abs(log2_p_value) > -log2(0.05)
& freq == "primary", "blue",
ifelse(abs(log2_p_value) < -log2
(0.05), "grey", "NA"))))
cell_type_df$cell_type <- factor(cell_type_df$cell_type, levels = cell_type_df[order
(cell_type_df$log2_p_value, decreasing = F),]$cell_type)

```



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

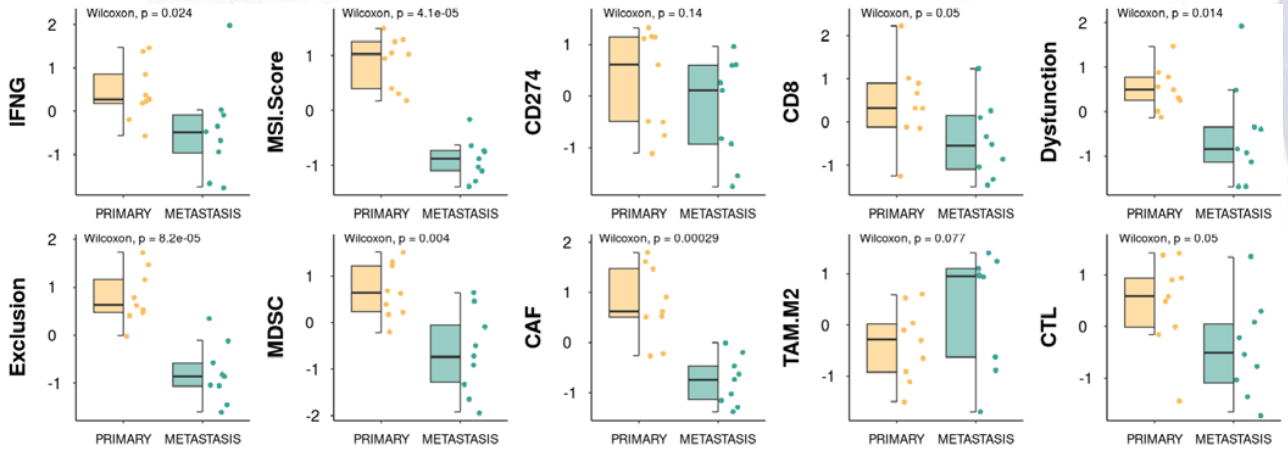
Signatures of T cell dysfunction and exclusion predict cancer immunotherapy response

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

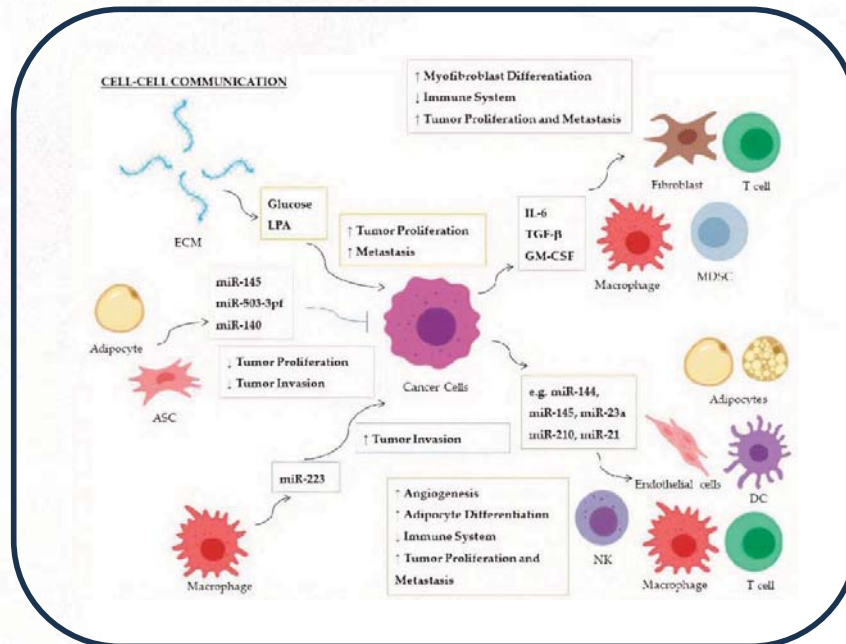
59k Accesses | 2269 Citations | 157 Altmetric Metrics

	IFNG	PSI_Score	CD274	CD8	DysFunction	Exclusion	MDSC	CAF	TAM_M2	CTL
ACC_T_16	2.7719452	0.30270286	0.5687956	2.2232168	1.9386887	-0.45261846	-0.051164299	0.002992495	-0.036737680	2.20361677
ACC_T_03	0.50849971	0.88954761	1.2365711	-2.2484679	-0.09121424	1.49101306	0.089843357	-0.139496973	-0.0063551964	-2.38583075
ACC_T_13	2.07460238	0.78585513	1.0657304	1.2062644	1.4815029	0.71345090	0.009393304	0.127709458	-0.019163032	1.53371222
ACC_T_15	0.3877009	0.80772155	1.0323718	0.5037738	0.90075486	1.25315294	0.02534704	0.15470957	0.0125360288	0.59997088
ACC_T_08	1.93107510	0.95170196	1.0657304	1.8198446	0.60117765	0.57082723	0.039209776	0.04777594	-0.0201004187	2.25467003
ACC_T_11	-0.79437342	0.54803795	-0.7130554	-0.2614655	0.53027842	-0.43202697	0.071956051	-0.019512709	-0.0023782193	-0.25541541
ACC_T_04	0.34080156	0.52457338	-0.4547424	-0.5755183	0.00180791	0.30928623	0.035830870	0.079250652	-0.02256481511	-0.03773558
ACC_T_09	-1.36521189	0.27038910	-0.8619188	0.5310158	0.52125515	0.32742377	0.039151032	-0.016429675	0.0255731042	-0.42955880
ACC_T_14	0.25552226	0.81373064	-0.4547424	1.6208195	0.29709455	0.49583761	0.041898668	0.046300336	-0.0241610421	1.45824411
ACC_T_12	-1.20027231	0.67802523	0.607156	0.0100306	0.0527366	1.01367462	0.07715012	0.055399598	-0.011648665	2.20437956
ACC_T_10	-0.27233313	0.61581230	-1.0270287	-0.2138809	0.32945539	-0.03174874	0.013298078	-0.053384041	0.0091615238	-0.80380597
ACC_T_06	-0.68389175	0.1019429	-0.7707709	-2.6987001	-1.60770152	-0.05589544	0.030360855	-0.036787790	0.0201927166	-2.23030923
ACC_T_17	-2.45351201	0.22641100	-1.0309392	-2.4062188	-1.01629434	-0.63583409	-0.027377810	-0.093482711	0.0267252800	-2.7131673
ACC_T_02	-2.56012913	0.15919592	-1.4323044	-1.9733693	-1.0725935	-0.60406541	-0.000513432	-0.114392118	0.0064246500	-1.79229383
ACC_T_18	-0.48025992	0.06943787	0.5578338	-1.4733693	-0.86554040	-0.83671020	-0.070262950	-0.000010792	0.0210833504	1.37529126
ACC_T_07	-0.59015469	0.20026171	0.1063724	-0.9822203	-0.78011890	-0.85666243	-0.039151173	-0.106680569	0.0302452300	-0.8343846
ACC_T_01	-0.11522170	0.44700400	0.0999215	-0.5071828	-0.29822962	-1.18600129	-0.180706009	-0.052000200	-0.0137160043	0.42510174
ACC_T_05	0.03725959	0.10103876	0.2438759	0.2783981	-0.34984742	-1.31150590	-0.009348783	-0.083941825	-0.0191899342	0.07702116



TME and cell to cell interaction

Cell-cell communication within the tumor microenvironment

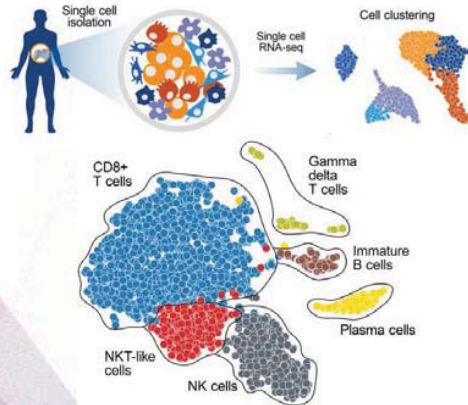


TME and cell to cell interaction

Cell type annotation

Single R

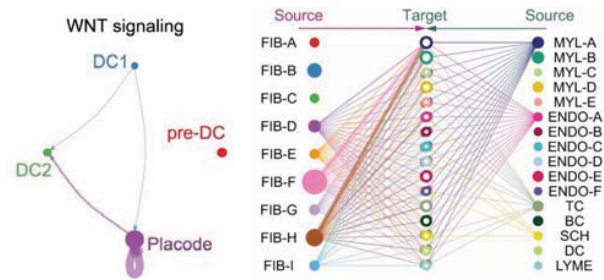
- computational method for unbiased cell type recognition of scRNA-seq
- SingleR's annotations combined with Seurat, a processing and analysis package designed for scRNA-seq



Cell to cell interaction

CellChat

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



Ianevski, 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

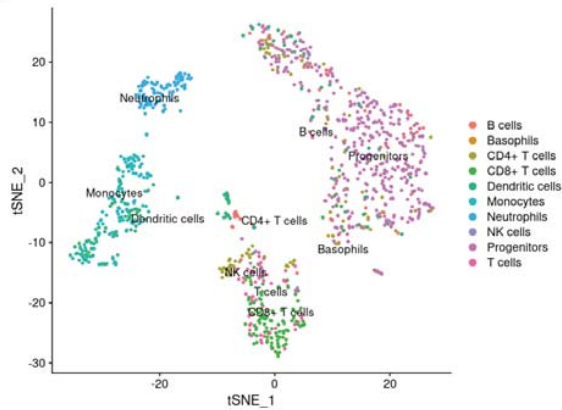
Cell type annotation

```

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")

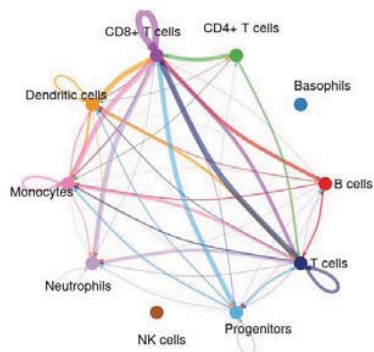
```

Visualization

```

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

```

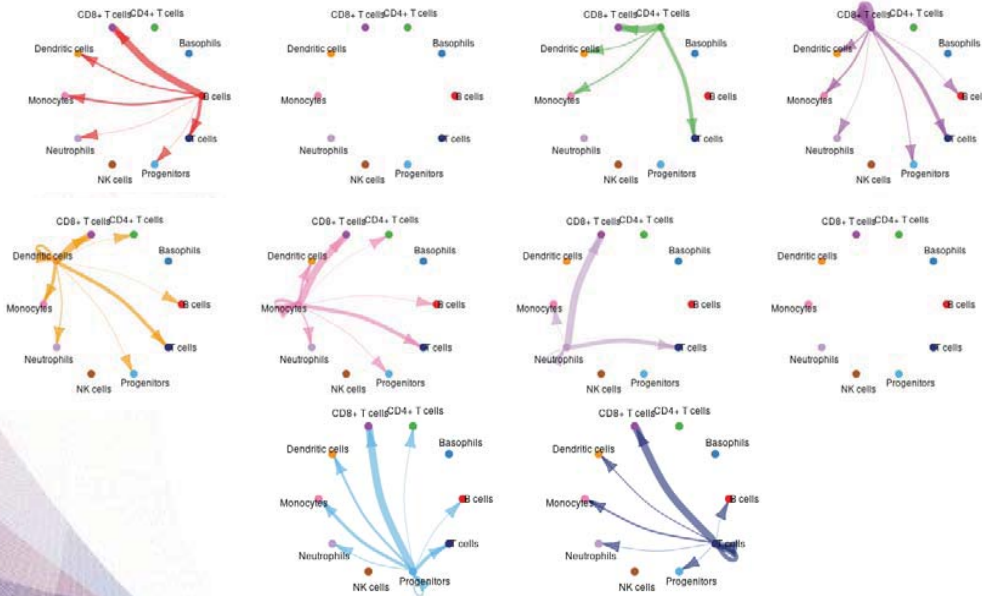


Visualization

```

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()

```

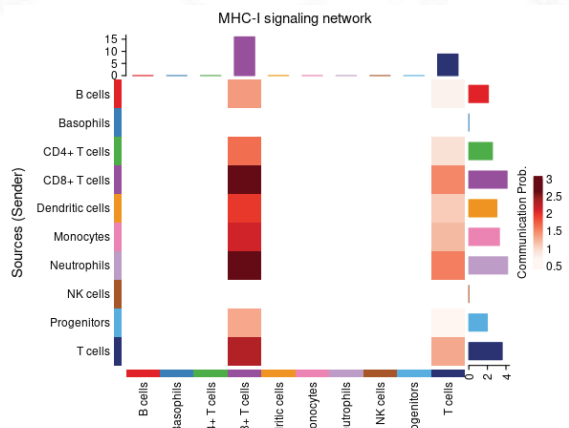
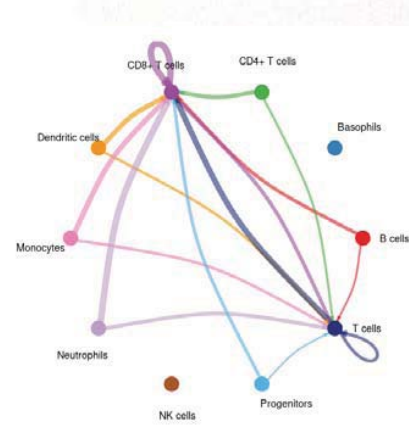


Visualization

```

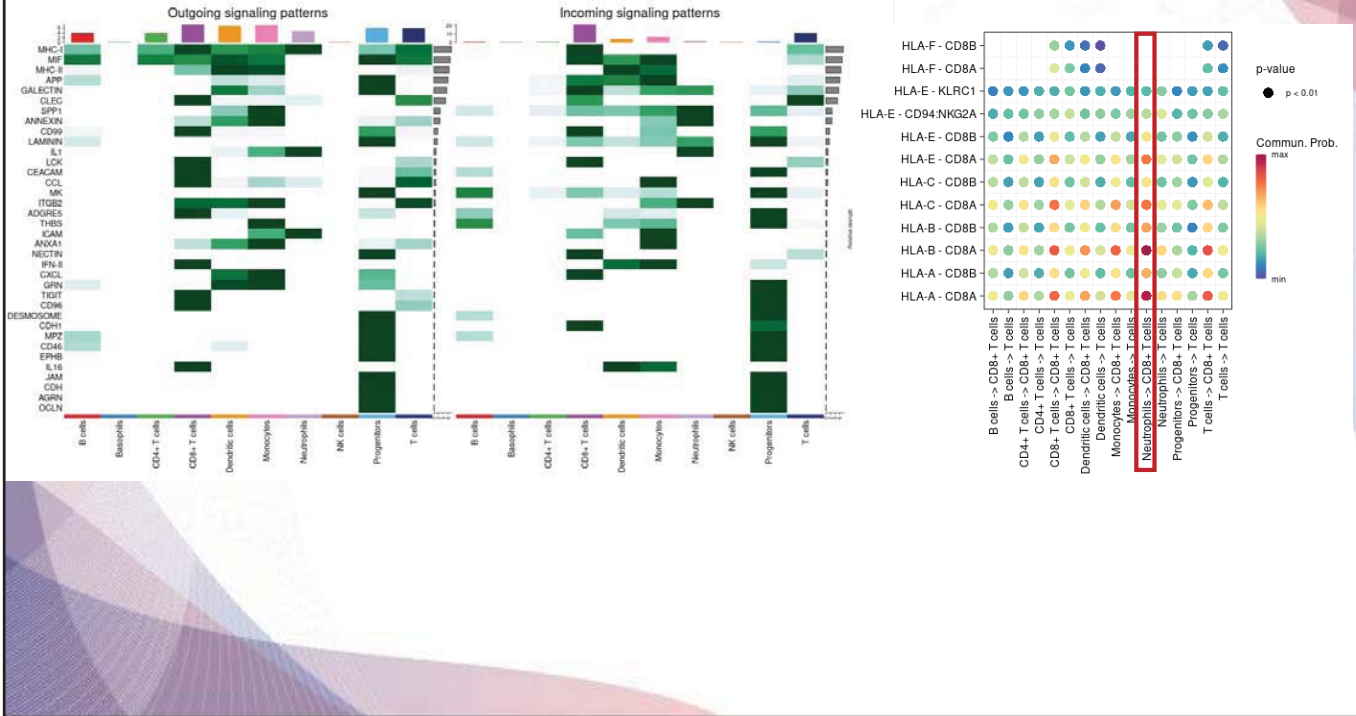
pathways.show <- c("MHC-I")
netVisual_aggregate(cellchat, signaling = pathways.show, layout = "circle")
netVisual_heatmap(cellchat, signaling = pathways.show, color.heatmap = "Reds") #특정 생물학적 경로 내 상호작용

```

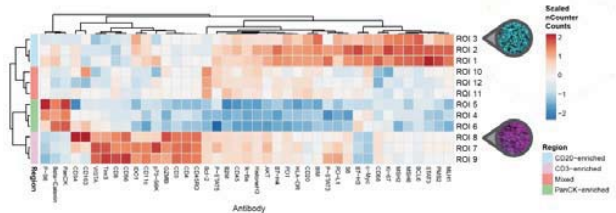
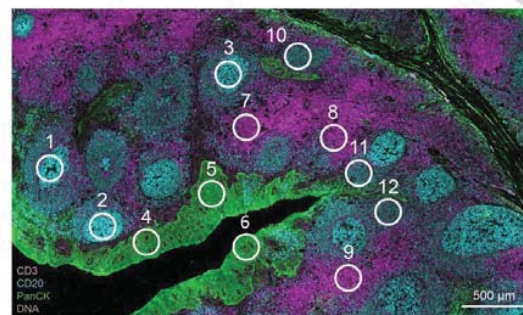
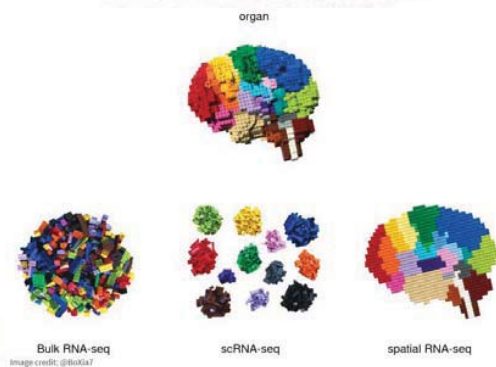


Visualization

```
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

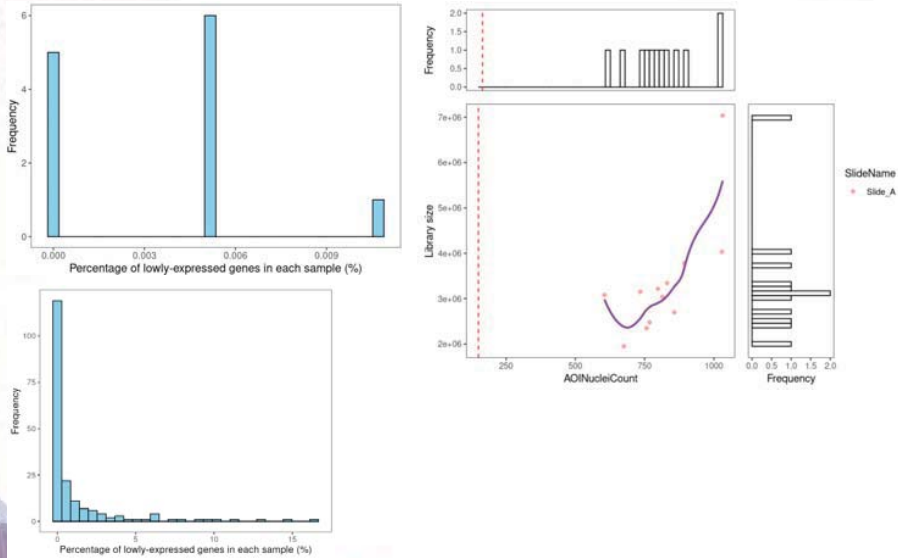
```
source("/data/project/BIML_2024/GeoMX/Script/GeoMX_function.R")  
  
#####  
# Load library #  
#####  
library(tidyverse)  
library(standR)  
library(SpatialExperiment)  
library(edgeR)  
library(limma)  
library(msigdb)  
library(GSEABase)  
library(SpatialDecon)  
library(speckle)  
  
# Visualization #  
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)
```

