

# Improved NK cell markers from single-cell RNA-seq of PBMC populations with the new ROC-driven combiROC R package



UNIVERSITÀ  
DEGLI STUDI  
DI MILANO

I. Ferrari<sup>1,2</sup>, S. Mazzara<sup>3</sup>, A. Gobbini<sup>1</sup>, N. Di Marzo<sup>4</sup>, M. Crosti<sup>1</sup>, S. Abrignani<sup>1,5</sup>, R. Grifantini<sup>1</sup>  
<sup>4</sup>, M. Bombaci<sup>1\*</sup>, R.L. Rossi<sup>1\*</sup>

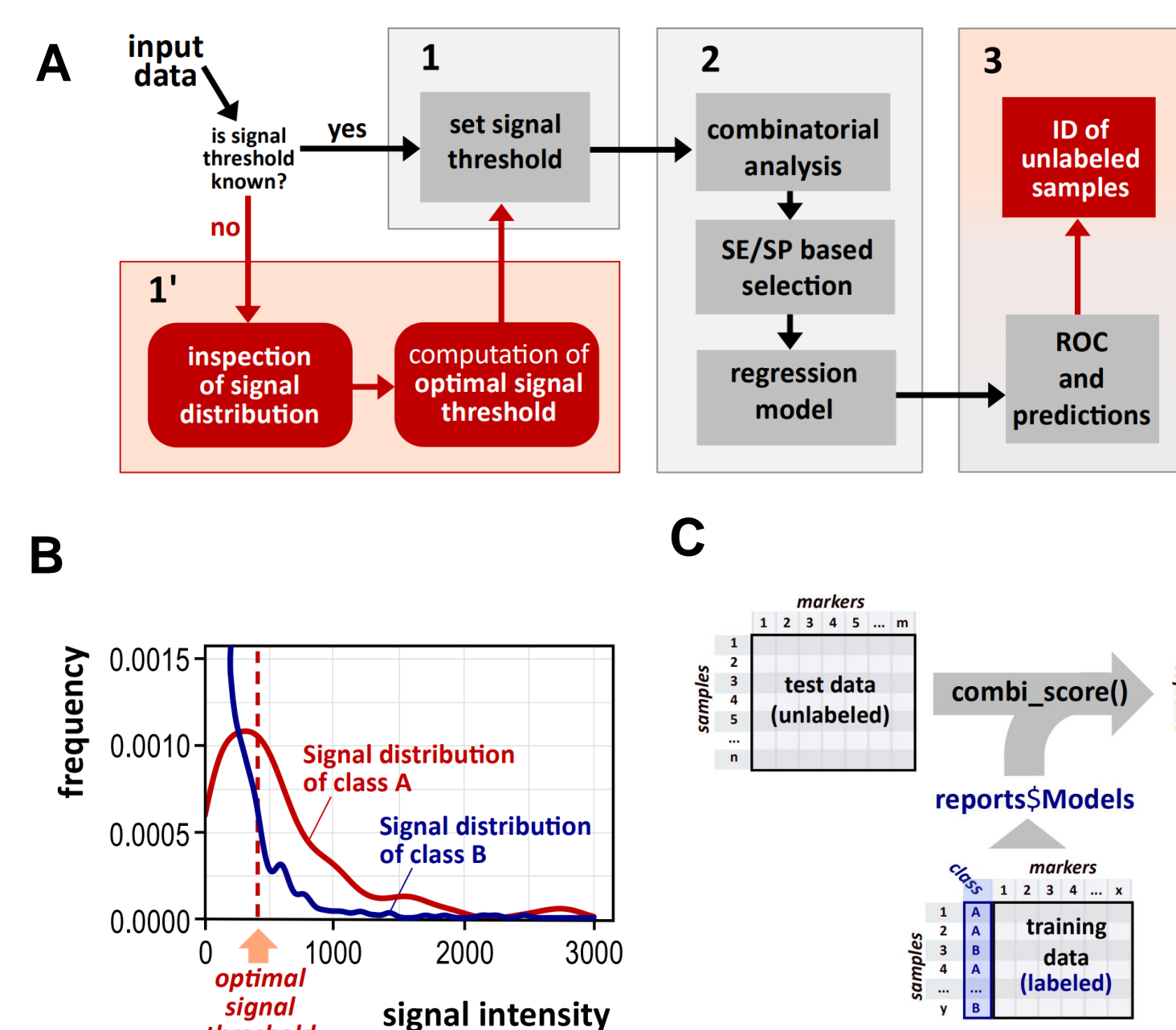
<sup>1</sup> National Institute of Molecular Genetics, Milan, Italy. <sup>2</sup> Department of Biosciences, University of Milan, Milan, Italy. <sup>3</sup> Department of computing Sciences and Bocconi Institute for Data Science and Analytics (BIDSA), Bocconi University, Milan, Italy. <sup>4</sup> CheckAb Srl, Milan, Italy. <sup>5</sup> Department of Pathophysiology and Transplantation, University of Milan, Milan, Italy.  
\* Corresponding authors



## Abstract

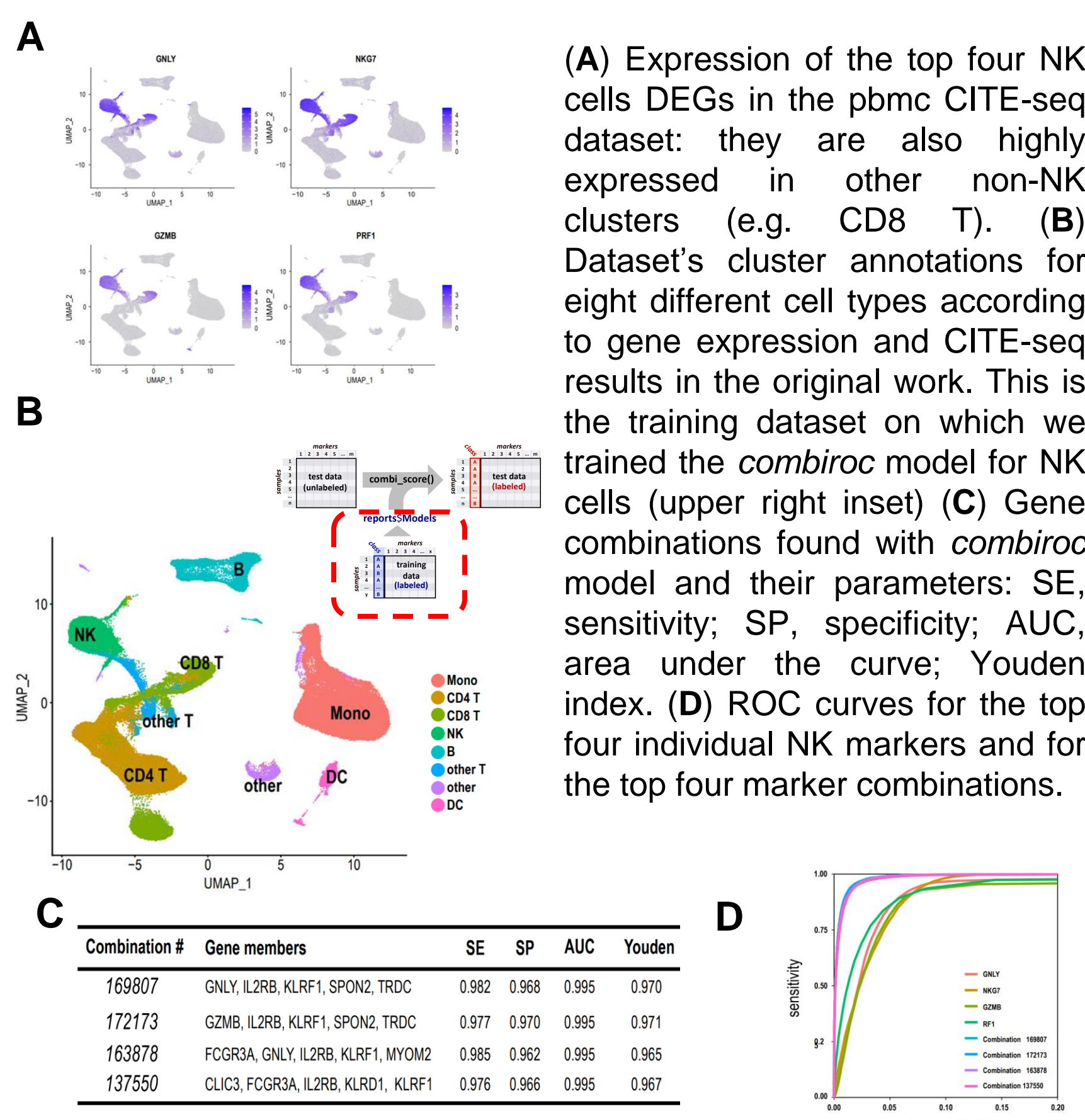
In this study, we introduce the **combiROC R package**, a tool for refining signatures in high-throughput omics data. Leveraging a ROC-driven combinatorial selection approach, this new package was designed to facilitate the identification of potent sub-signatures from single-cell RNA-seq experiments, enabling more efficient cell annotation using a reduced set of markers. By applying *combiROC* to Peripheral Blood Mononuclear Cells (PBMC) datasets, we identified non-canonical marker combinations for Natural Killer (NK) cells that aligned with the Human Protein Atlas (HPA). We further validated these combinations through cytometry staining and functional assays. **The single-cell workflow presented in this work significantly impacts marker signature research in general and transcriptomic gene signatures in particular.** It demonstrates that the top differentially expressed genes are not necessarily the most specific ones and that smaller signatures can be more powerful, regardless of the differential expression ranking of individual markers. This principle of "less is more" has the potential to re-evaluate existing gene signatures and bring forth new markers that may have gone unnoticed so far. <https://ingmbioinfo.github.io/combiROC/> (also on CRAN).

## 1. CombiROC automates signal threshold setting and creates models to annotate unlabelled samples

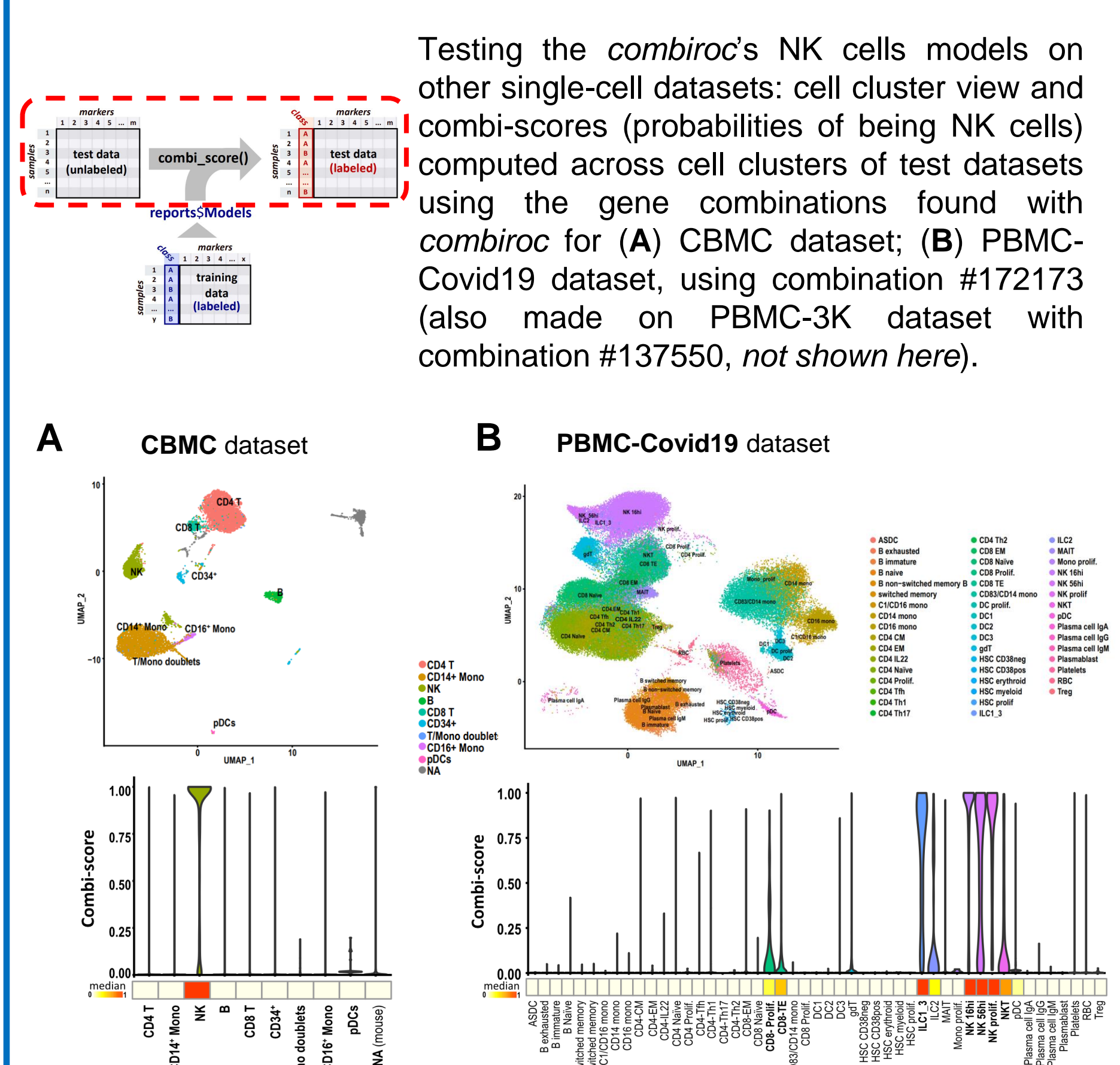


(A) *combiROC's* workflow. Red boxes are new features introduced by the *combiROC* package: alternative phase 1' for computation of signal threshold and part of phase 3 with labeling of unknown samples. (B) Display of the calculated optimal signal threshold (dashed line) on overlapping signal intensity distributions. (C) Samples of unlabeled data (left) can be associated with a class (right, red annotations), using regression models generated from labeled training data (blue annotations). "A" and "B" refer to the binary annotation of classes in the datasets.

## 2. CombiROC unlocks hidden markers sub-signatures

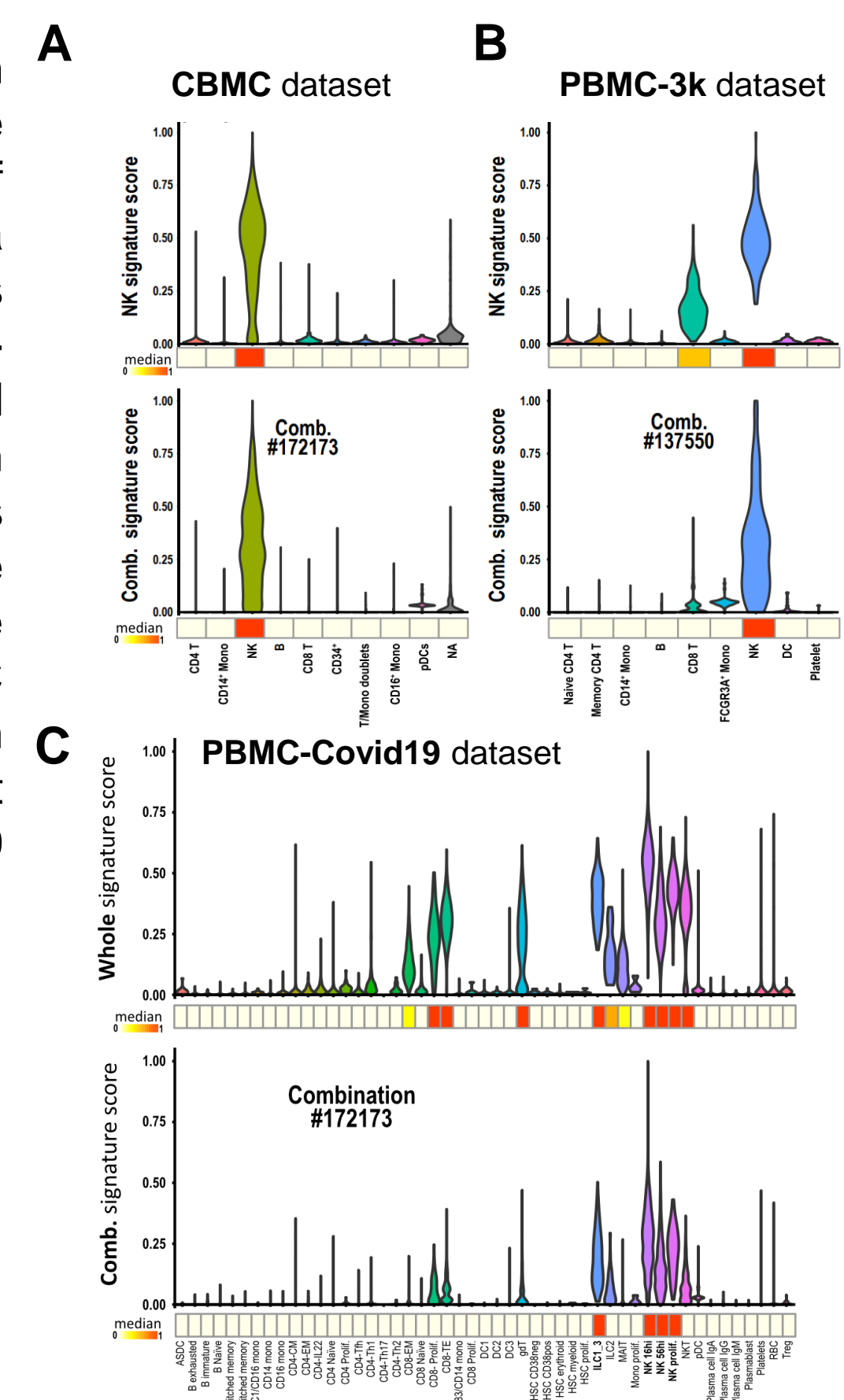


## 3. Sub-signatures accurately identifies NK cells in unlabelled single-cell datasets

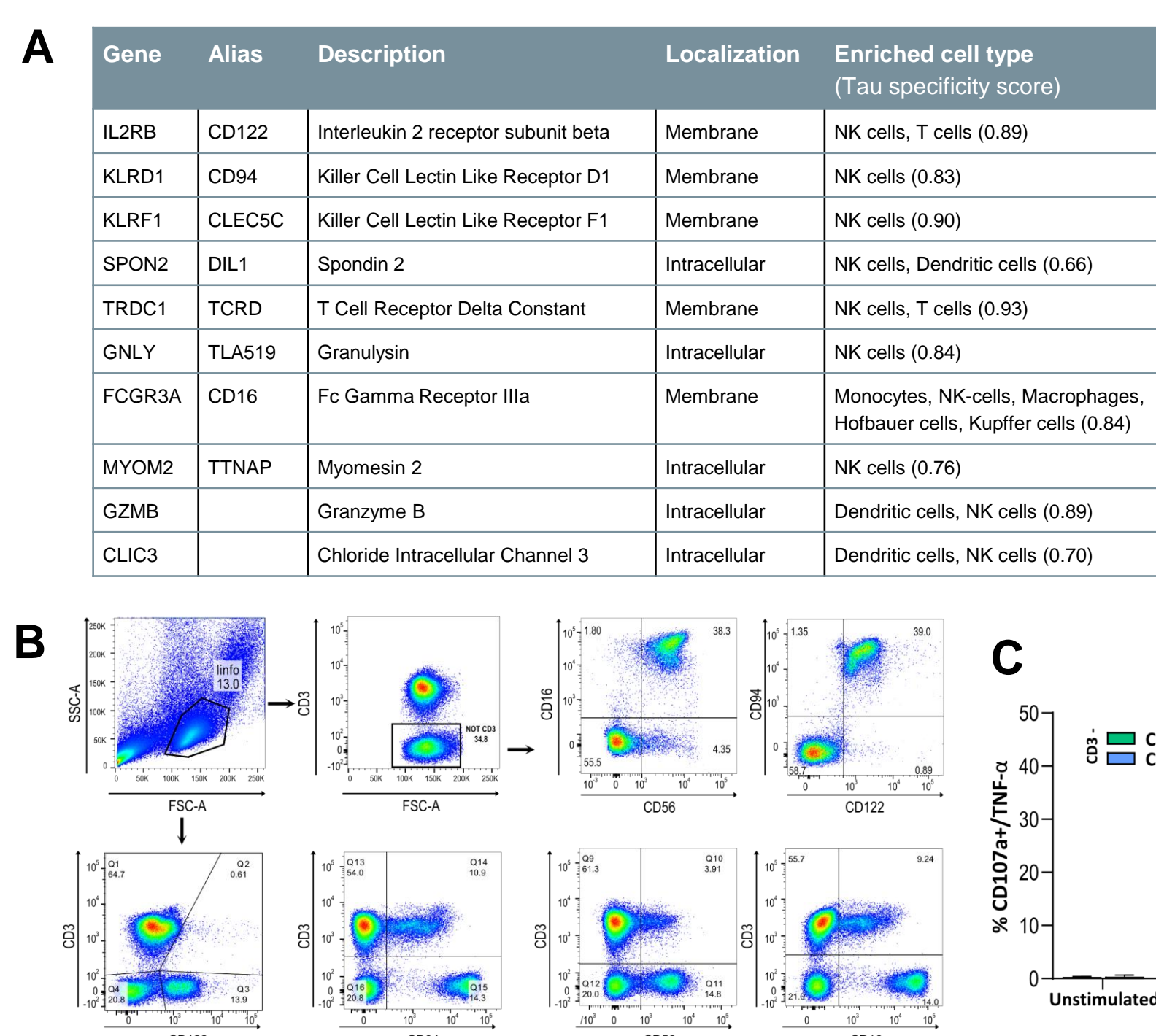


## 4. Revolutionary 5-gene cocktail better discriminates NK cells than 30-genes parent signature

Gene signature scores can be calculated to determine the descriptive power of any gene signature (Della Chiara et al., 2021): scores derived from *combiROC's* 5-genes combinations proved to be just as effective in identifying NK cells as those obtained from the entire 30-genes signature as observed in CBMC dataset (A), and even better in PBMC-3k dataset (B) and in PBMC-Covid19 dataset (C).

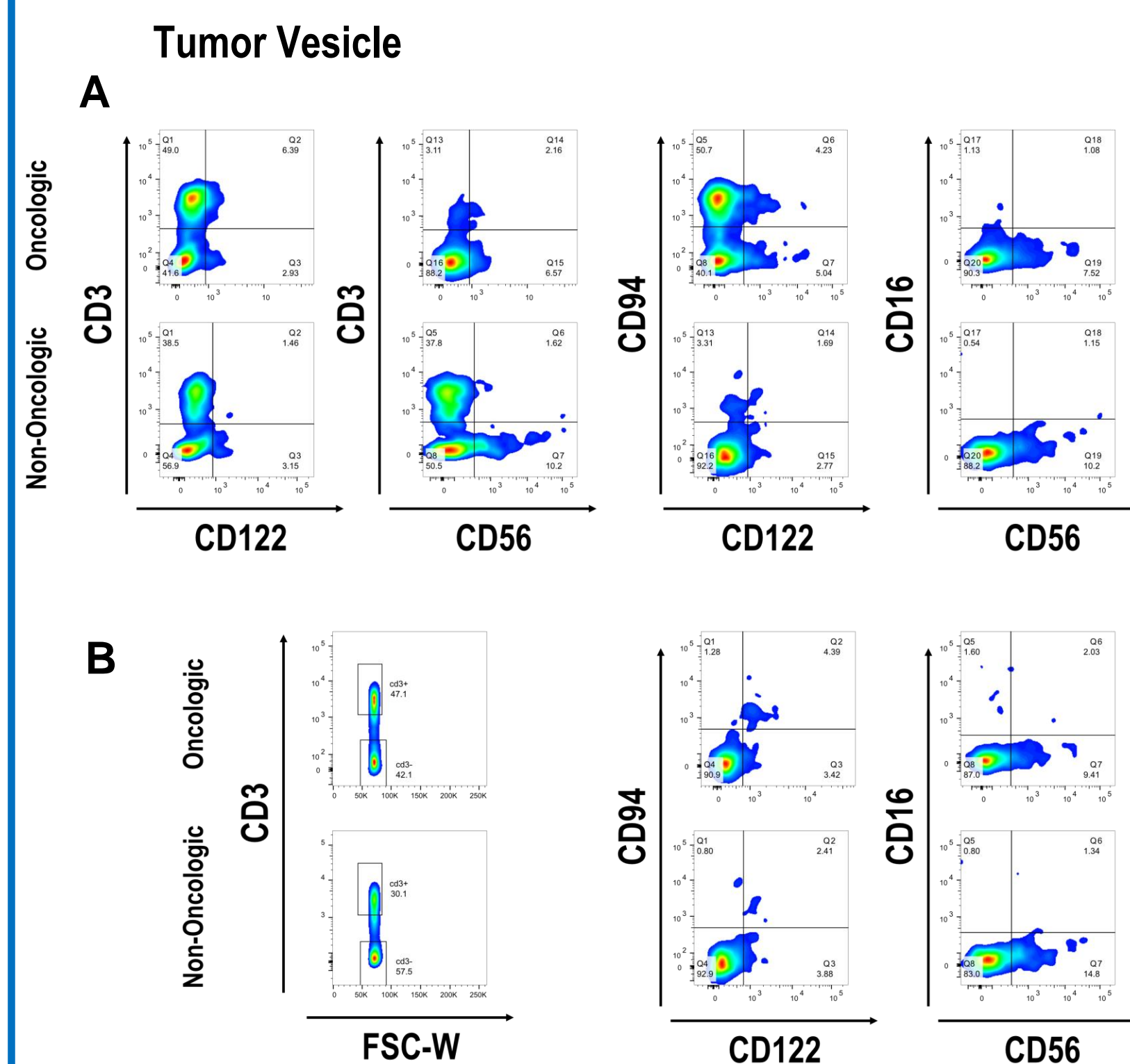


## 5. IL2RB (CD122) is specifically associated with highly functional NK cells



(A) The ten genes belonging to the four selected combination (from panel 2C) with cell type enrichment (Tau scores) as they are annotated in the "Blood and Immune" subset of Human Protein Atlas (HPA, <https://www.proteinatlas.org/>). (B) Gating strategy to discriminate among the different NK cells subsets. Live lymphocytes were gated based on side and forward scatter dot plot display during the acquisition process. The NK specific markers were then measured along with CD3 staining. (C) Fraction of CD3+/CD56+ and CD3+/CD122+ cells expressing CD107a after stimulation with PMA/ionomycin and with K562 cells. Percentages result from four independent experiments.

## 6. CD94/CD122 expression identifies NK cell population in tumours



(A) CD56 expression in the CRC tumor is significantly higher than CD122, although it remains comparable when contrasted with the non-tumor region, where CD122 exhibits notable differences in expression.

(B) In CD3-negative cells revealed that while the expression of the CD56/CD122 combination remains consistent between the tumor and non-tumor regions, the CD94/CD122 combination is uniquely present in the tumor region and absent in the non-tumor region.

## SUMMARY

# **CombiROC simplifies cell identification:** *combiROC* streamlines the single-cell RNA-seq data analysis, making it easier to pinpoint specific cell populations by reducing the number of marker genes to consider.  
# **Smaller signatures for better insight:** smaller marker combinations, identified by *combiROC*, offer more precise insights than traditional differential expression rankings.  
# **Discovering overlooked markers:** the "less is more" approach reveals potentially hidden markers, shedding light on previously unnoticed cell characteristics.  
# **Translational potential:** highly performant *combiROC*-identified marker combinations have diagnostic and therapeutic applications, enhancing our understanding of cell populations.

## REFERENCES

- Mazzara et al. Sci. Rep. 2017. CombiROC: an interactive web tool for selecting accurate marker combinations of omics data.
- Bombaci & Rossi. Methods Mol. Biol. 2019. Computation and selection of optimal biomarker combinations by integrative ROC analysis using *combiROC*.
- Satija et al. Nat. Biotechnol. 2015. Spatial reconstruction of single-cell gene expression data.
- Stephenson et al. Nat. Med. 2021. Single-cell multi-omics analysis of the immune response in COVID-19.
- Hao, Y. et al. Cell. 2021. Integrated analysis of multimodal single-cell data.
- Della Chiara et al. Nat. Commun. 2021. Epigenomic landscape of human colorectal cancer unveils an aberrant core of pan-cancer enhancers orchestrated by YAP/TAZ.



## WHAT WE'RE WORKING ON

We are currently working on improvements:

- Algorithm improvement by integration of feature selection and parallelization
- Porting and release of a Python package
- Application of the algorithm to elusive immunological populations



ferrari@ingm.org  
bombaci@ingm.org  
rossi@ingm.org

