

A test for the comparison of gene-set transcriptomic profiles of vaccines

Arthur Hughes¹² Denis Agniel⁴⁵ Rodolphe Thiébaut¹²³ Boris Hejblum¹²

¹University of Bordeaux, BPH INSERM U1219, INRIA SISTM ²Vaccine Research Institute ³CHU de Bordeaux

⁴Harvard Medical School, Boston, USA ⁵RAND Corporation, Santa Monica, USA



Introduction

Comparing gene expression profiles between vaccines has high potential

- Understanding vaccine mechanisms
- Identifying biomarkers

But there are many challenges with high-dimensional data...

- ⚠ Interpretation of results
- ⚠ Sensitivity to investigator choices
- ⚠ Low signal-noise ratio



Gather genes in *gene sets* defined by biological function

- ✓ Reduce dimension
- ✓ Ease interpretation
- ✓ Maximise signal

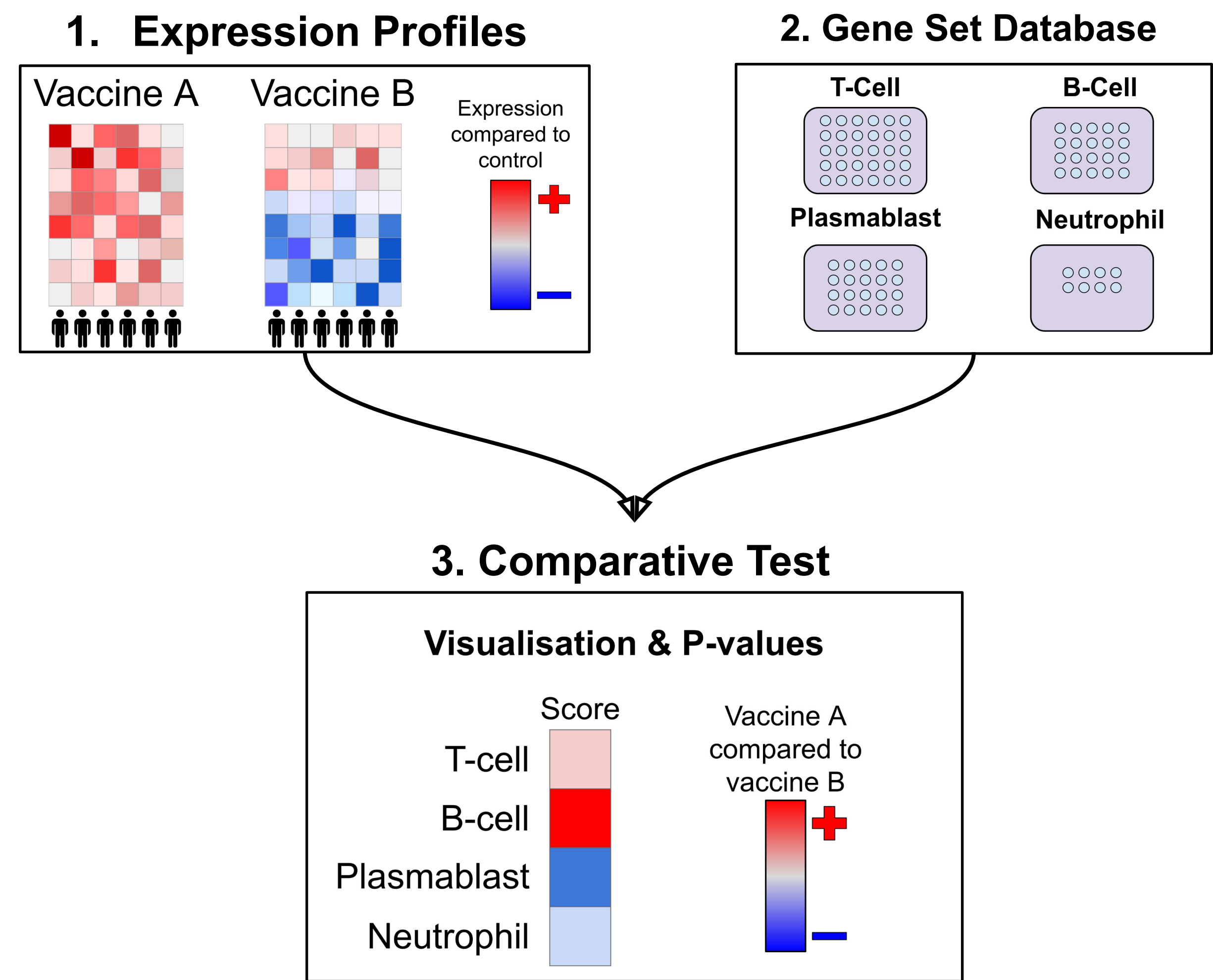


Figure 1. Illustration of the gene set comparative test.

Working Linear Mixed Effects Model

Test statistic derived from a working linear-mixed effects model:

$$\mathbf{y}_i^G = \alpha_0 + \mathbf{X}_i \alpha + \Phi_i \beta + \Phi_i \xi_i + \epsilon_i$$

- $\mathbf{y}_i^G = ((\mathbf{y}_i^1)^T, \dots, (\mathbf{y}_i^p)^T)^T$ - expression of p genes in set G for individuals $i = 1, \dots, n$
- \mathbf{X}_i - matrix of baseline covariates to control for
- Φ_i - K time-dependent variables whose association with \mathbf{y}_i^G is to be tested

- $\Rightarrow \beta$ - fixed effects of testing variables
- $\Rightarrow \xi_i \sim \mathcal{N}(\mathbf{0}, \Sigma_\xi)$ - individual-level random effects of testing variables
- $\Rightarrow \epsilon_i \sim \mathcal{N}(\mathbf{0}, \Sigma_i)$ - random error terms

Null Hypothesis

$$H_0 : \beta = \mathbf{0}, \xi_i = \mathbf{0}$$

Variance-Component Score Test

The derived **variance-component score test** statistic [1] is $Q = \mathbf{q}^T \mathbf{q}$ with

$$\mathbf{q}^T = n^{-1/2} \sum_{i=1}^n (\mathbf{y}_i^G - (\alpha_0 + \mathbf{X}_i \alpha))^T \Sigma_i^{-1} \Phi_i$$

- Central limit theorem $\Rightarrow Q \underset{+}{\sim} \sum_{k=1}^{pK} a_k \chi_1^2$
where a_k is the k th eigenvalue of $\text{cov}(\mathbf{q})$

- ✓ Only requires estimation of model under H_0 (i.e. linear model)!
- ✓ Type 1 error control relies only on the central limit theorem!

Gene-Set Correlation

- ⚠ Dependence between genes in same set



Take into account in residual variance Σ_i

- Fit linear model for each gene $g \in G$:
$$y_{ij}^g = \alpha_0^g + \mathbf{x}_i^T \alpha^g + \phi_{ij}^T \beta^g + e_{ij}^g$$
- OLS estimates residuals $r_{ij}^g = y_{ij}^g - \hat{y}_{ij}^g \Rightarrow \mathbf{r}^g = (r_{11}^g, \dots, r_{1t}^g, \dots, r_{n1}^g, \dots, r_{nt}^g)^T$
- Entry $(g1, g2)$ of $\Sigma_i \Rightarrow$ **covariance between gene-wise residuals** $\forall g1, g2 \in G$

$$[\Sigma_i]_{g1, g2} = \text{Cov}(\mathbf{r}^{g1}, \mathbf{r}^{g2})$$

Application - Yellow Fever vs Flu Vaccines

- Compare early gene expression signatures of yellow fever and influenza vaccines*
- Gene sets: BloodGen3 Modules [2] - immunology focused sets

	Score	P-Value
M8.3 - Type 1 Interferon		1.8×10^{-8}
M10.1 - Interferon		2.0×10^{-7}
M15.127 - Interferon		2.3×10^{-7}
M15.86 - Interferon		2.4×10^{-7}
M13.17 - Interferon		8.8×10^{-5}
M16.30 - Complement		1.6×10^{-4}
M15.64 - Interferon		7.7×10^{-4}
M16.40 - Monocytes		1.1×10^{-3}
M15.36 - Protein Modification		2.2×10^{-3}
M14.65 - Monocytes		2.2×10^{-3}

Score: YF
Compared
to Flu

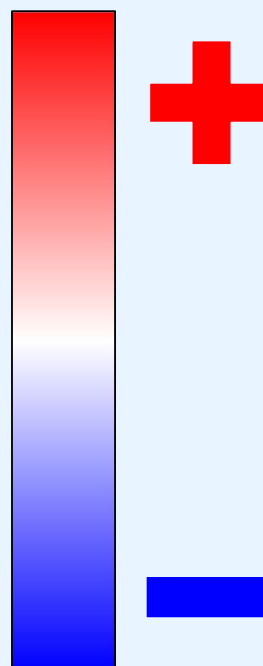


Figure 2. Scores and p-values for the top 10 most significant sets.

- \Rightarrow 300/382 gene sets significant**
- \Rightarrow Interferon modules strongly up-regulated early in YF compared to flu vaccine

*Public data from Immune Signatures Data Resource [3] - gene expression prior to 7 days post-vaccination of 1090 samples from 709 individuals vaccinated with either YF-17D or inactivated seasonal influenza vaccines. Analysis controls for age and sex, but not correlation structures.
**P-values corrected for multiple testing with Benjamini-Hochberg procedure

Summary

- Gene sets resolve problems with transcriptomic data
- Variance-component score test is a flexible, powerful method
- Biologically interpretable differences found between two vaccines
- Future work : explore properties of test under correlation structure estimation

References

- Marine Gauthier, Denis Agniel, Rodolphe Thiébaut, and Boris P Hejblum. Dearseq: A variance component score test for rna-seq differential analysis that effectively controls the false discovery rate. *NAR Genomics and Bioinformatics*, 2(4), 2020.
- Matthew C. Altman and et al. Development of a fixed module repertoire for the analysis and interpretation of blood transcriptome data. *Nature Communications*, 12(1), July 2021.
- Joann Diray-Arce and et al. The immune signatures data resource, a compendium of systems vaccinology datasets. *Scientific Data*, 9(1), 2022.

Arthur Hughes is supported by the Digital Public Health Graduate school, funded by the PIA 3 (Investments for the Future - Project reference: 17-EURE-0019)