

Old data, new tricks: Comprehensive computational analysis of 10 years of multi-center EuroFlow AML diagnosis data

Sarah Bonte

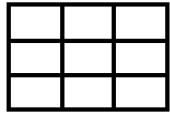


01.
Data

The dataset



Patients with **acute myeloid leukemia**
EuroFlow ALOT + AML diagnosis panel
2011 - 2023
8 centers
5379 FCS files from 799 patients

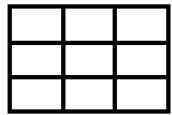


WHO 2008 classification

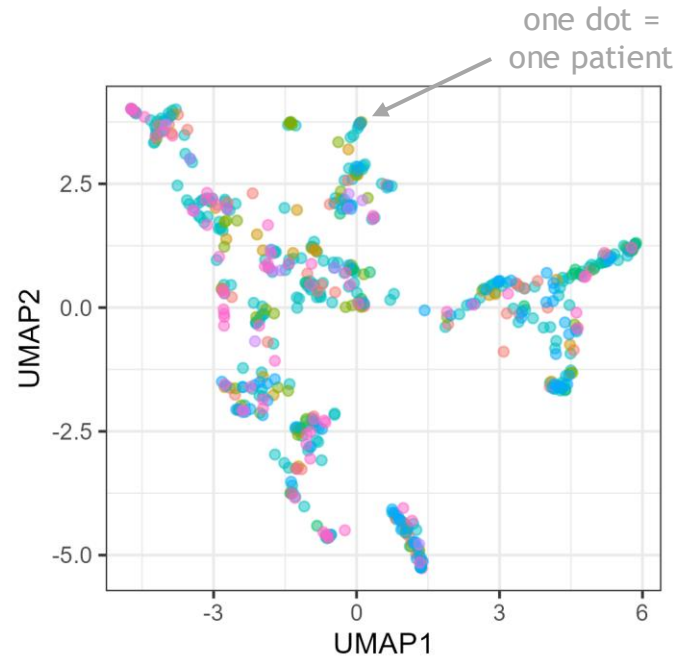


02. Results

EuroFlow standardization mitigates center-specific batch effects

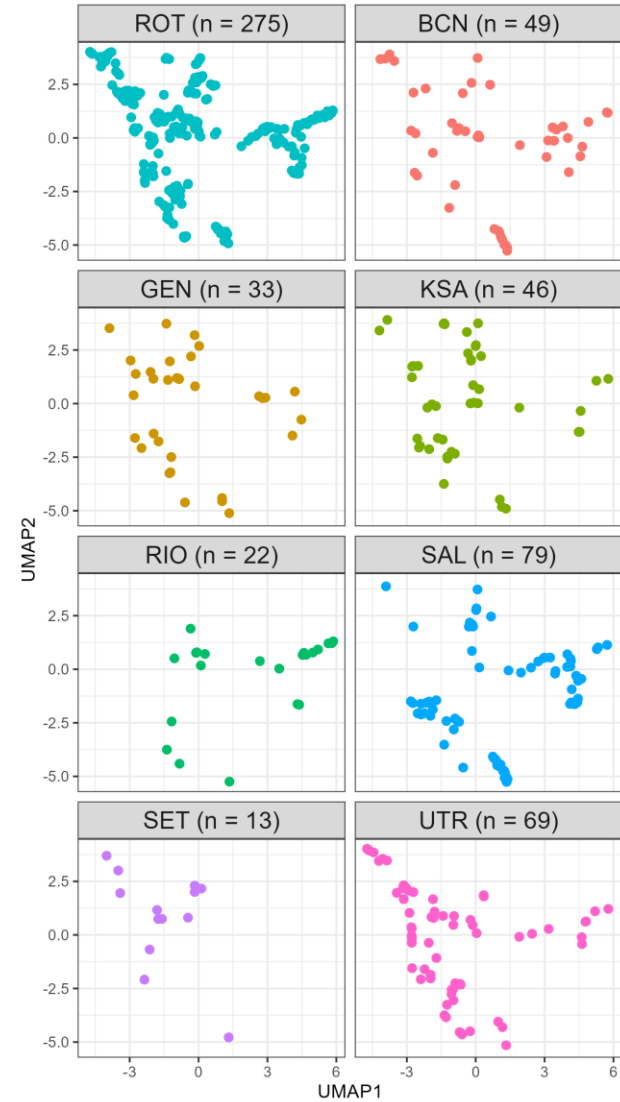
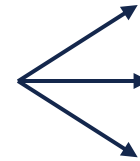


FlowSOM-derived
cell population %

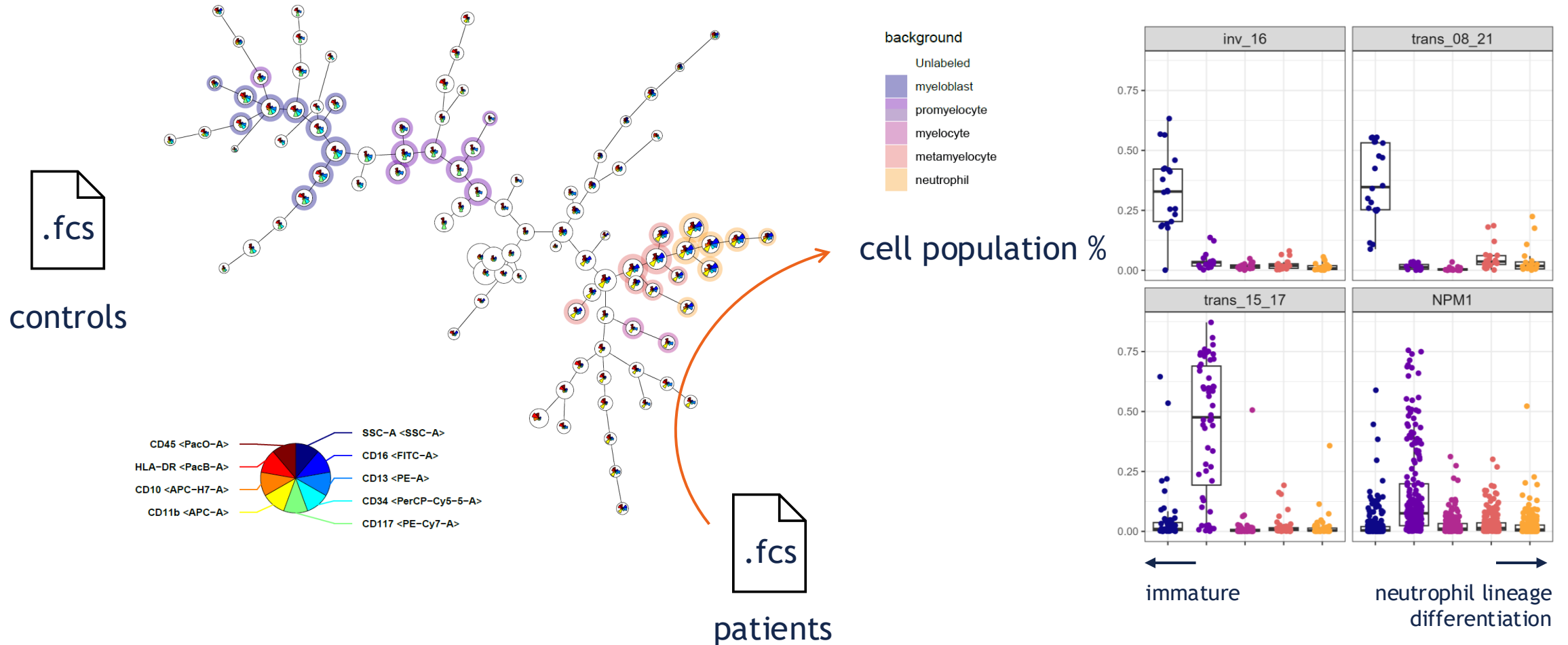


Center

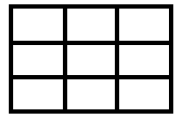
- BCN
- GEN
- KSA
- RIO
- ROT
- SAL
- SET
- UTR



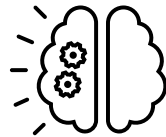
FlowSOM analysis locates maturation arrest at various stages for distinct WHO classes



A random forest machine learning model can predict WHO genetic classes from immunophenotypic data



FlowSOM-derived
cell population %



random forest
machine learning model

AML_CEBPA	5	0	0	0	0	0	0
AML_inv_16	0	10	0	0	0	0	0
AML_MLL	0	0	12	0	0	1	10
AML_trans_08_21	0	0	0	13	0	0	0
AML_trans_15_17	0	0	0	1	34	4	3
AML_NPM1	1	0	4	0	1	139	26
other	2	4	5	1	7	10	91
	AML_CEBPA	AML_inv_16	AML_MLL	AML_trans_08_21	AML_trans_15_17	AML_NPM1	other

predicted label

"true" label






Sofie Van Gassen
Yvan Saey
Mattias Hofmans
Rosan Olsman
Vincent van der Velden

all EuroFlow collaborators



SarahM.Bonte@UGent.be

Old data, new tricks: comprehensive computational analysis of 10 years of multi-center EuroFlow acute myeloid leukemia diagnosis data

Sarah Bonte^{1,2}, Sofie Van Gassen^{1,2}, Rosan Olsman¹, Yvan Saey^{1,2}, Mattias Hofmans^{3,4}, Vincent van der Velden¹

1 Data Mining and Modelling for Biomedicine group, VIB Center for Inflammation Research, Ghent, Belgium; 2 Department of Mathematics, Computer Science and Statistics, Ghent University, Ghent, Belgium; 3 Department of Immunology, Erasmus MC, University Medical Center Rotterdam, Rotterdam, The Netherlands; 4 Department of Laboratory Medicine, Ghent University Hospital, Ghent, Belgium; 5 Department of Diagnostic Sciences, Ghent University, Ghent, Belgium

Acute myeloid leukemia (AML) is a malignancy with high inter- and intra-patient heterogeneity, both at the genotypic and phenotypic level, posing challenges in diagnosis, prognosis and treatment. The most recent edition of the WHO classification divides AML into AML with defining genetic aberrancies and AML defined by differentiation. Here, we investigated AML immunophenotypes and WHO classes, with a comprehensive computational analysis on an extensive AML dataset.

Data & data analysis

File selection

Reference FlowSOM

One FlowSOM model was created per tube of the EuroFlow AML panel, with non-lymphocyte cells from one center (large cohort, both pediatric and adult)

Compensation QC

Per center and per marker: Files were removed if $q5 < \text{median}(q5) - 2 * \text{MAD}(q5)$

Random forest model

Metacluster percentages (MCK) were used as input features. Within a 5-fold cross-validation, Synthetic Minority Oversampling Technique (SMOTE) was used to overcome WHO class imbalance.

First, one RF model using the complete training dataset was used to determine feature importance. Features with correlation higher than 0.4 were removed, resulting in 119 from a total of 358 features being used for training of the model.

RF model predictions result in a prediction score for each WHO class. The class with the highest score was assigned as the predicted WHO class.

Results

EuroFlow standardization mitigates center-specific batch effects

UMAP with FlowSOM-derived metacluster % (tube 4)

CytoBatchFlagR framework for batch effect detection (results for tube 4 of the panel)

The combination of EuroFlow standardization and robust FlowSOM metaclustering resulted in batch effect correction between centers not being required for further downstream analysis.

FlowSOM analysis locates maturation arrest at various stages for distinct WHO classes

FlowSOM based on 10 healthy controls

FlowSOM-based van Dongen et al metaclustering Leukemia (2012)

With FlowSOM we were able to identify neutrophil maturation stages (tube 1). Some genetic alterations lead to a maturation arrest at earlier stages. Patients from WHO classes like *inv(16)* and *t(8;21)* have more immature cells, whereas cells from other patient groups like *NPM1* and *t(15;17)* are more differentiated towards the neutrophil lineage. For the WHO diagnostic groups defined by differentiation, labels corresponded to the observed maturation arrest.

A random forest machine learning model can predict WHO genetic classes from immunophenotypic data

When performing 100 iterations of the RF model, 31 patients were wrongly predicted in each iteration. We investigated these in-depth and doublechecked their assigned WHO labels.

SHAP analysis

snippet from heatmap clustering all patients based on MCK over all tubes

The random forest model could accurately predict several genetic abnormalities, and these results were reproducible in an independent validation cohort (data not shown). In addition, some of the mispredictions hinted to actual mislabeling, or missing information, from the metadata.

When training machine learning models for clinical use, it is imperative to use correct metadata labels in the training dataset.