

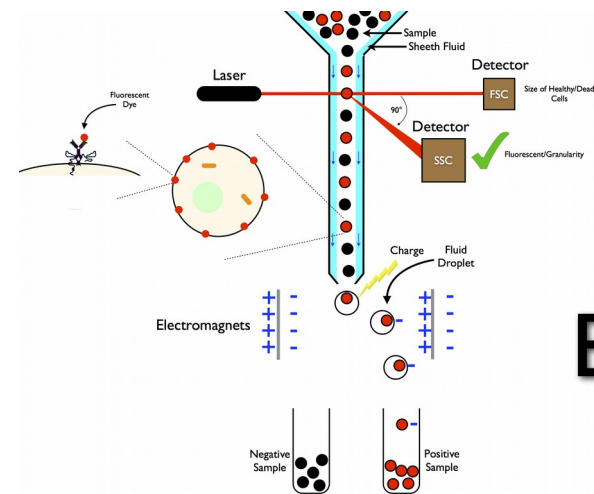
A decorative graphic at the top of the slide consists of numerous thin, overlapping wavy lines in various colors including red, orange, yellow, green, blue, and purple, creating a sense of motion and data flow.

# Klastrovanie buniek podľa cytometrických dát

Marek Behún

# Popis úlohy

Input: flow cytometry dáta (100k+, 13 až 50 dimenzií)  
(napríklad bunky z kostnej drene od zdravého darcu)  
Chceme: rozoznať populácie rovnakých buniek



**B**

# Problém

Pri veľkom počte buniek štandardné klastrovacie algoritmy nefungujú dobre:

šum

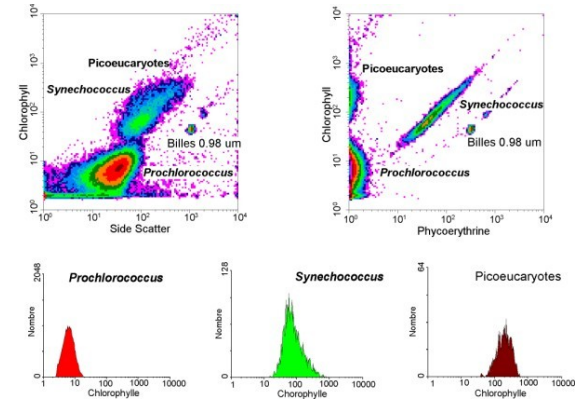
evolučná kontinuita

...

Lepšie algoritmy sú zase pomalé (rádovo 100 tisíce buniek).

# Hodnotenie presnosti

- časti datasetov sú olabelované tzv. “manuálnym gatingom”
  - ide o vizuálne pozorovanie scatterplotov, histogramov, ...
  - limitácie: subjektivita, chyba operátora, zlá reproducibilita, ťažkosti detekovať nové typy buniek
- používame F1 skóre
- hľadáme priradenie clusterov s najlepším súčtom F1 (cez všetky permutácie)



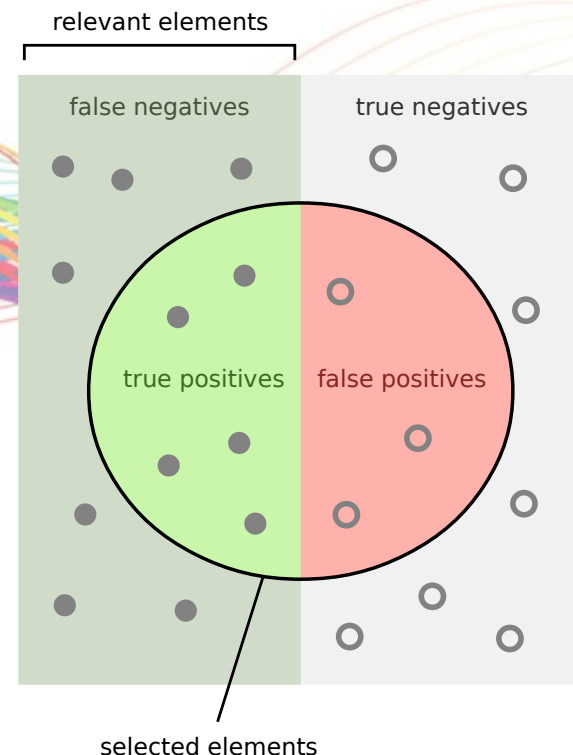


# F1 skóre

- harmonický priemer *precision* a *recall*

$$F_1 = \left( \frac{\textit{precision}^{-1} + \textit{recall}^{-1}}{2} \right)^{-1}$$

- $F_1$  je v  $[0, 1]$
- $F_1 = 1$  znamená perfektný *precision* a *recall*
- $F_1 = 0$  najhorší



How many selected items are relevant?

Precision =



How many relevant items are selected?

Recall =



# Použité datasety

## Levine\_13dim

- 13 dimenzií
- 167 034 záznamov
- 49% manuálne gatovaných

## Samusik\_all

- 39 dimenzií
- 841 643 záznamov
- 61% manuálne gatovaných

V oboch prípadoch surové dáta prešli transformáciou  $\arcsinh(x/5)$

# KMeans: Samusik

**KMeans**



$F_1$  skóre: 0.547035

Čas: 97s

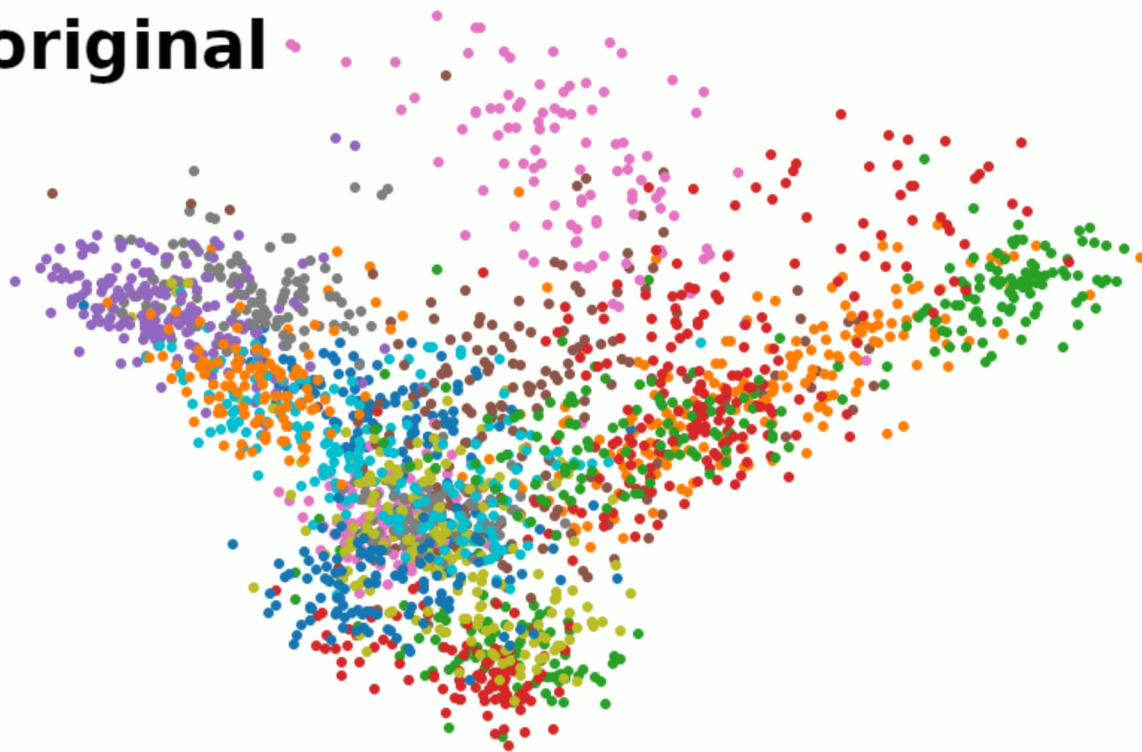






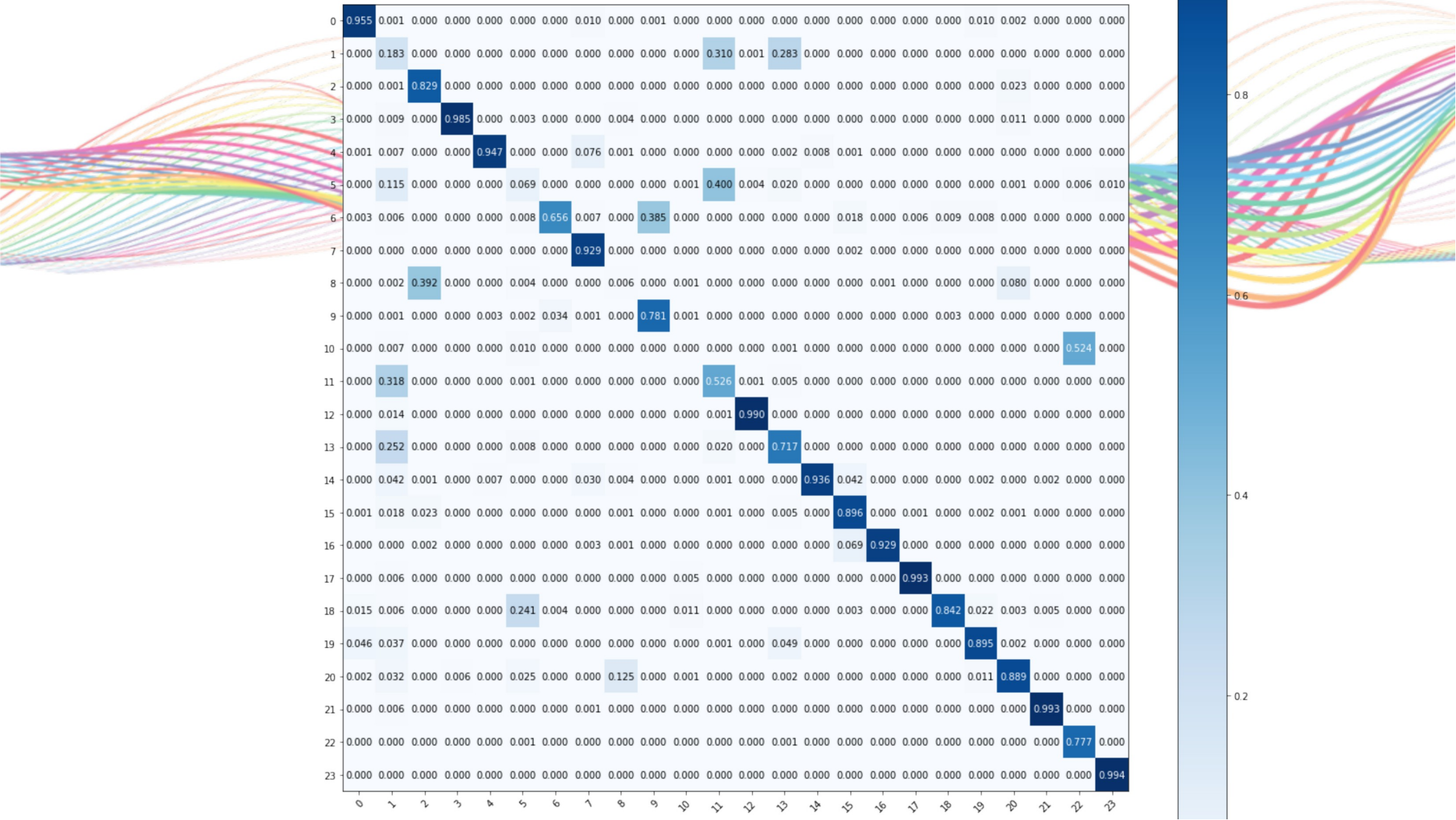
# HCA: Samusik

**original**



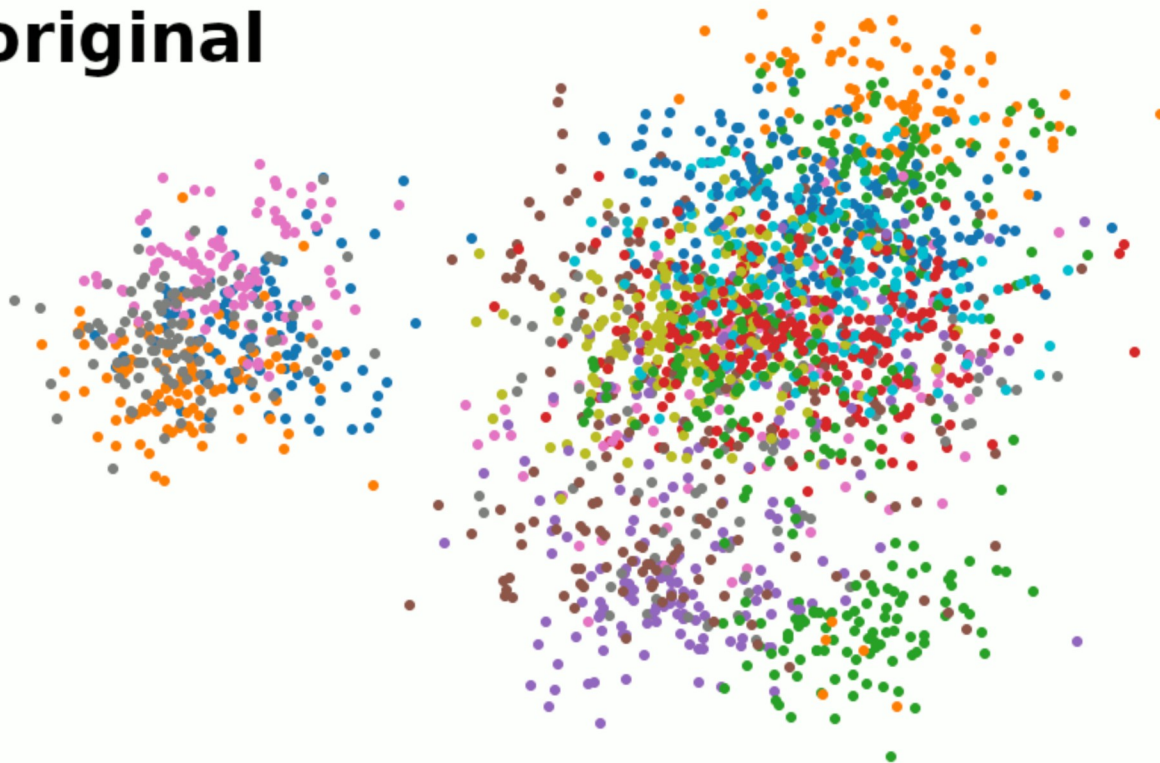
$F_1$  skóre: 0.738262

Čas: 434s (86 468 samples)



# KMeans: Levine13

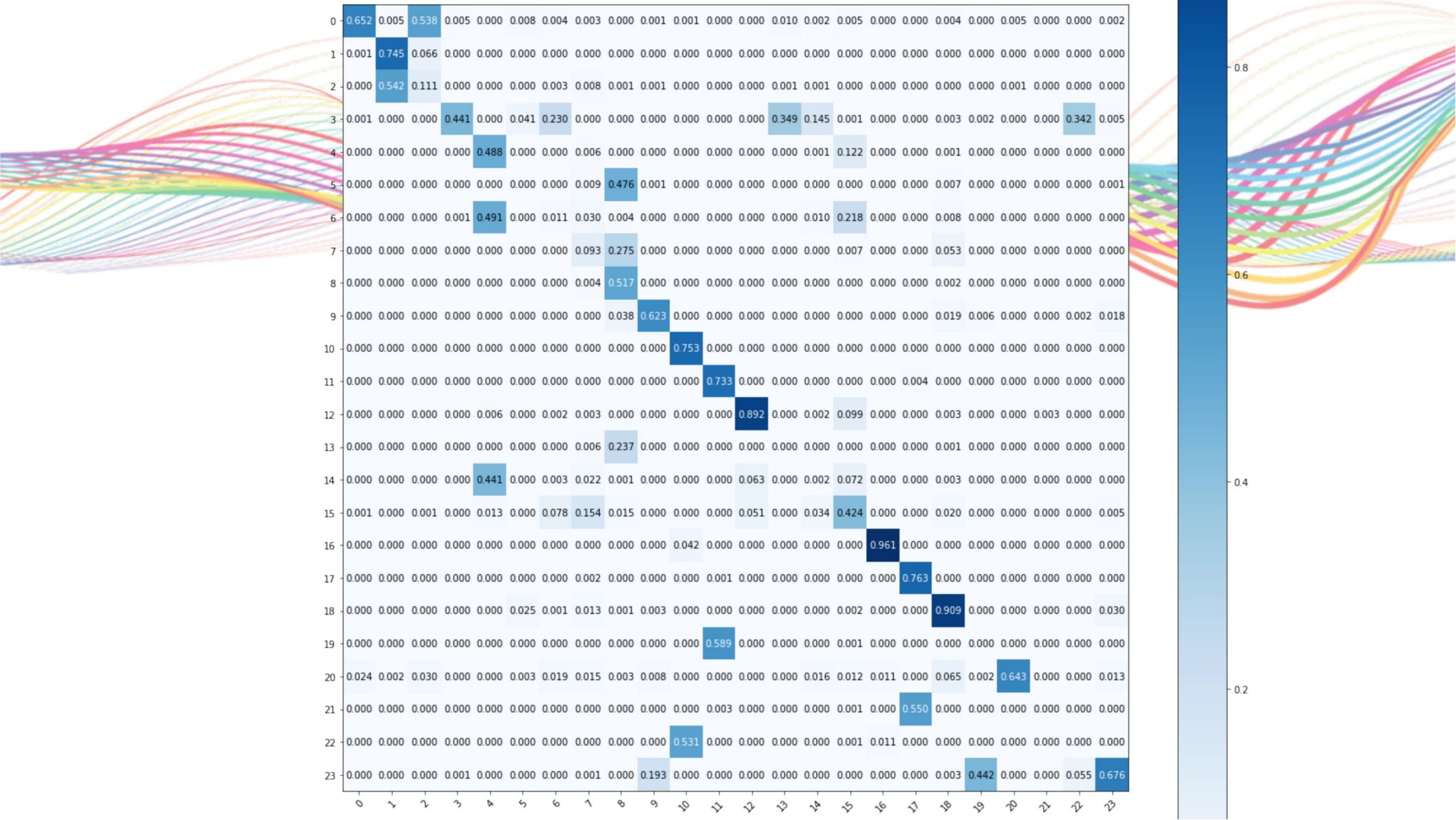
**original**



$F_1$  skóre: 0.434779

Čas: 11s

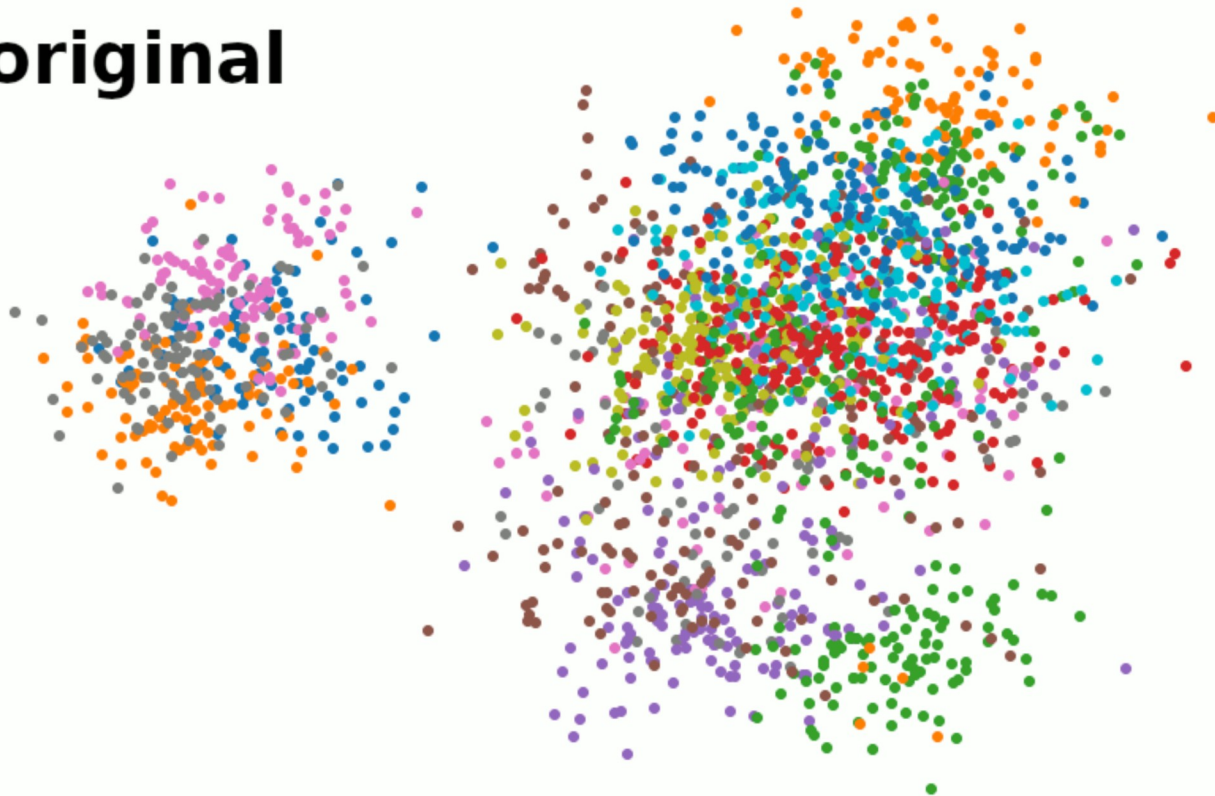






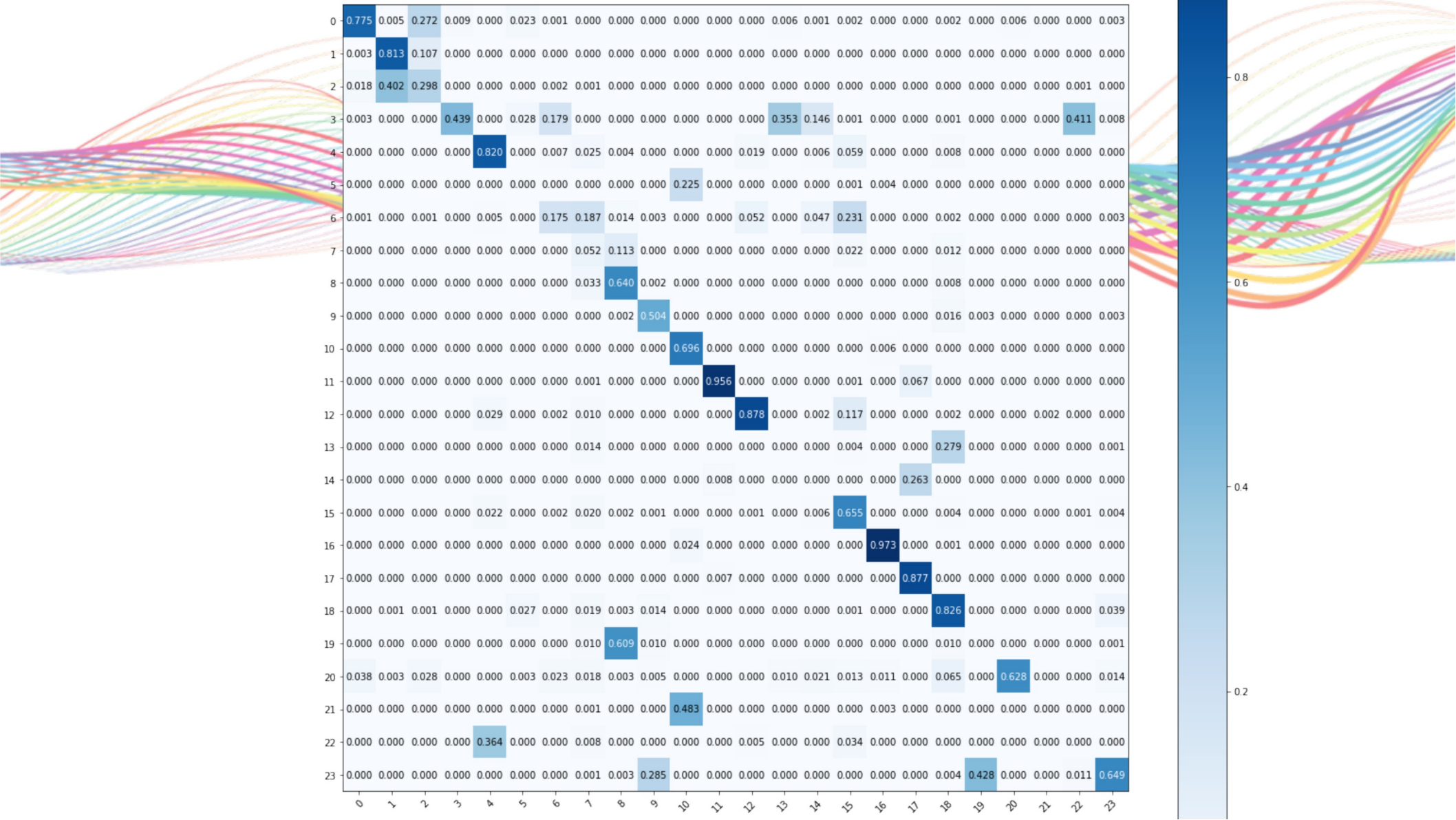
# HCA: Levine13

**original**



$F_1$  skóre: 0.485611

Čas: 207s



# Záver

## Levine\_13dim

- nepodarilo sa prekročiť  $F_1$  skóre z originálneho článku (0.518)
- HCA s 0.485 by bol tretí

## Samusik\_all

- najlepšie  $F_1$  je v originálnom článku 0.702
- HCA: 0.738
- neviem ale, či by bol rýchlejší :-)



# Future: Mahalanobis Clustering?



- chcel by som vyskúšať aglomeratívny clustering s Mahalanobisovou metrikou
- predstavujem si, že zhluky sa môžu vo flow cytometry dátach tvarovať do elipsoidov
- zatiaľ to nie je naimplementované v SciPy / scikit-learn (skúšal som to dohackovať, ale nepodarilo sa)





# Zdroje a odkazy

Videa: <https://blackhole.sk/~kabel/doc/bioinf>

Comparison of clustering methods for high-dimensional single-cell flow and mass cytometry data,  
<https://doi.org/10.1002/cyto.a.23030>