

Figure 4. Feature Set Analysis

Jordan A. Lee

Oregon Health & Science University

Portland, Oregon

12/18/20



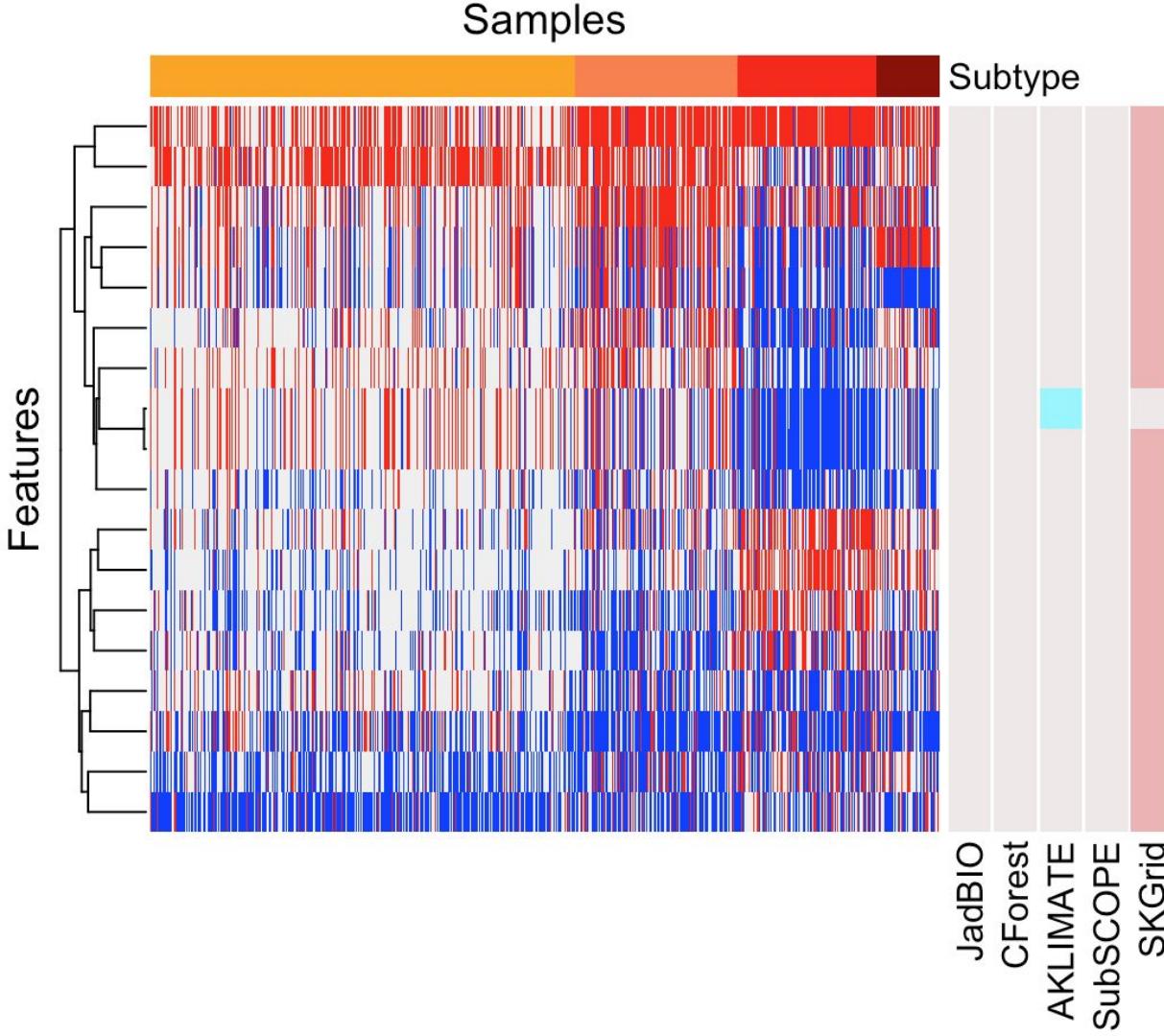
Developed in collaboration with Chris, Christina, Brian

Data Analysis

Jordan, Chris, Christina, Brian

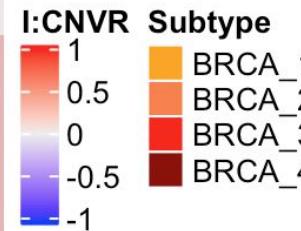
- Based on V8 tarball, these results are before the miRNA correction to certain cancers
- If interest, will rerun once all v9 results from classifier groups in
- Process:
 - Select a cancer cohort
 - Select best model per team (overall weighted F1)
 - Pull corresponding feature set for each model
 - Map back to molecular values (in v8 tarball)
 - Scale data if appropriate (z-scores for METH, GEXP, MIR)
 - Heatmaps
 - Plot scaled/raw data
 - Cluster feature rows
 - No clustering on sample/subtype columns

All heatmaps clustered using complete linkage method (max dist between clusters before merging) and euclidean distance



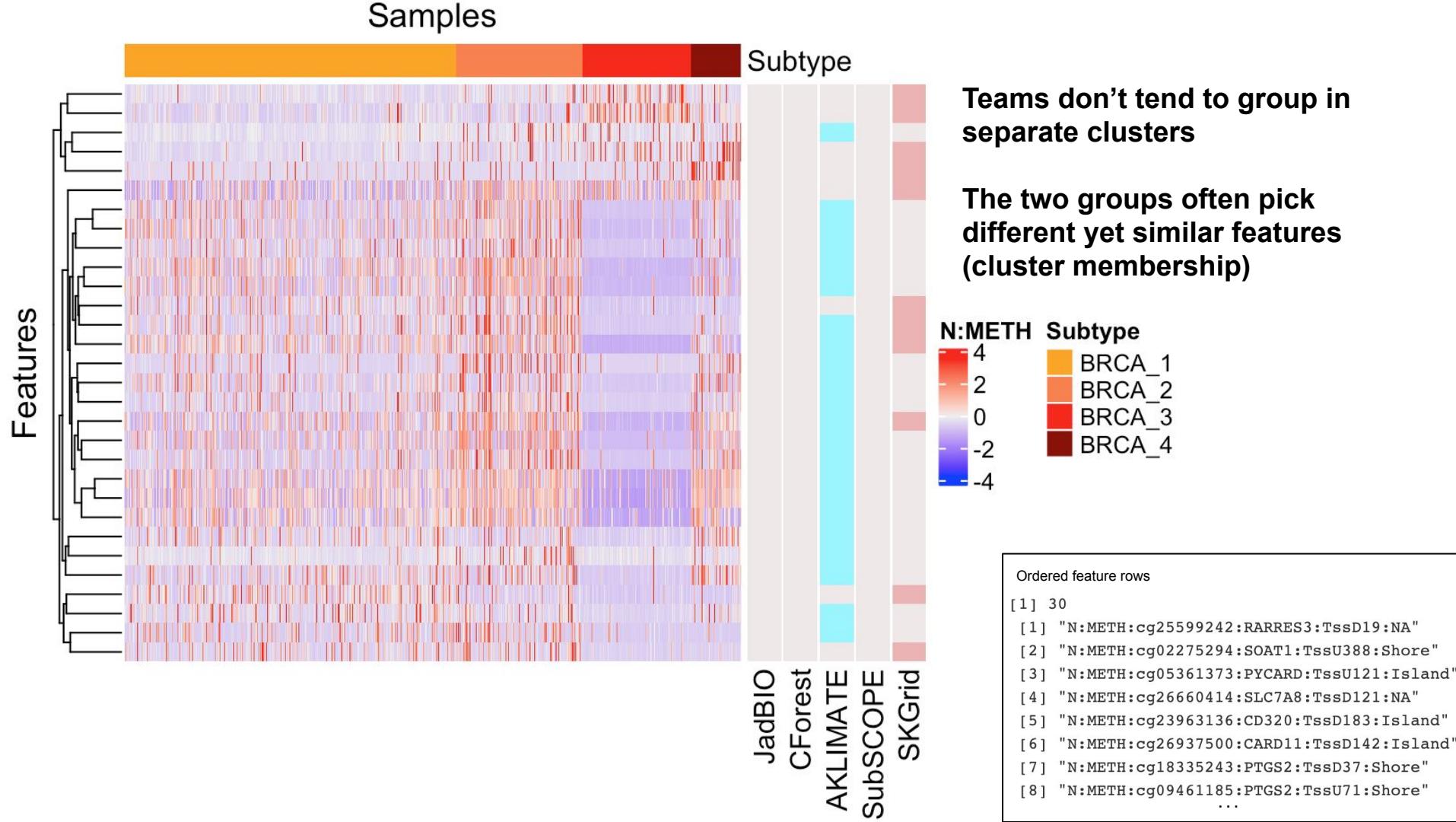
The one AKLIMATE feature is different than any picked from SKGrid

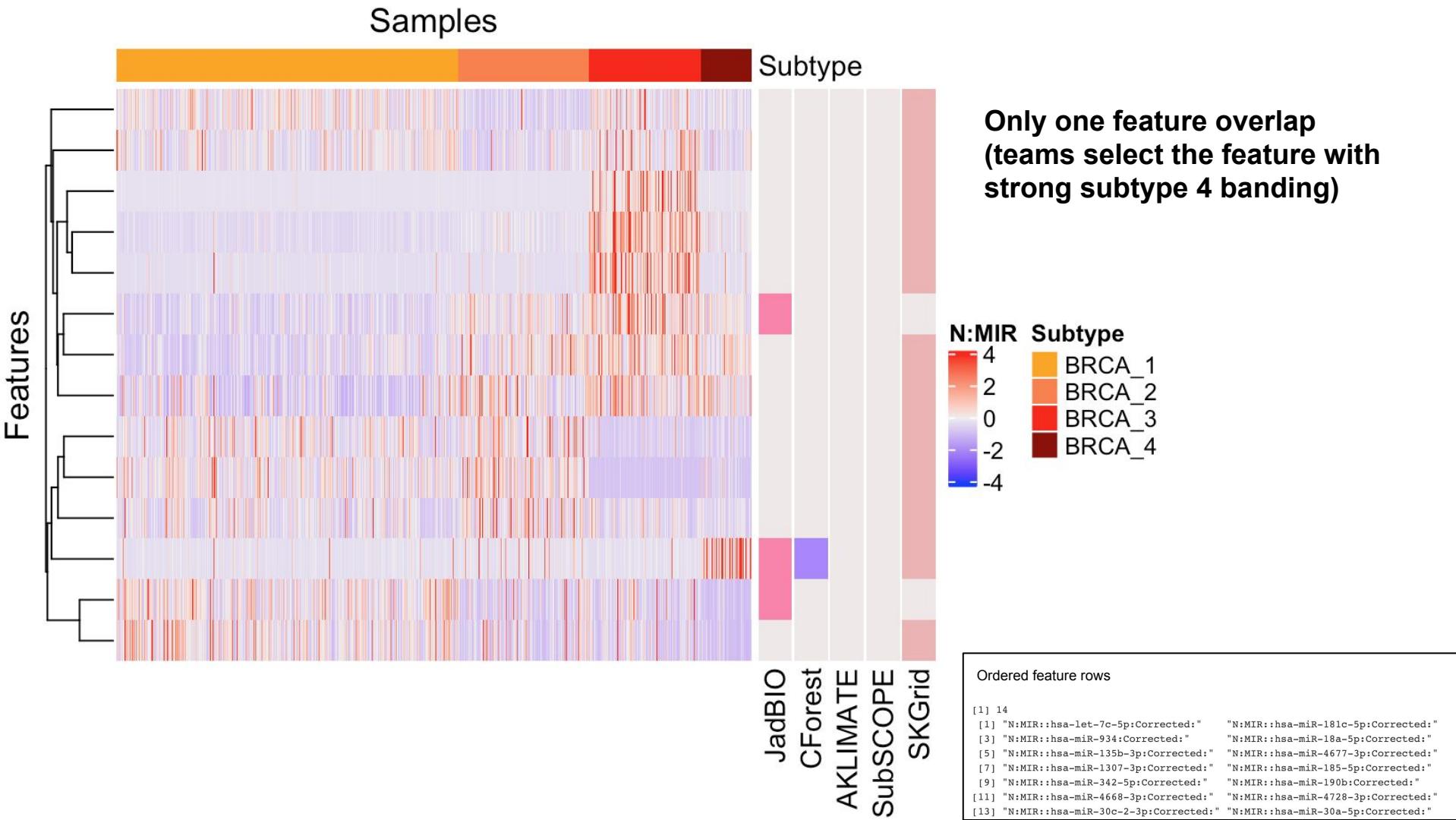
But appears to be very similar to a SKGrid feature

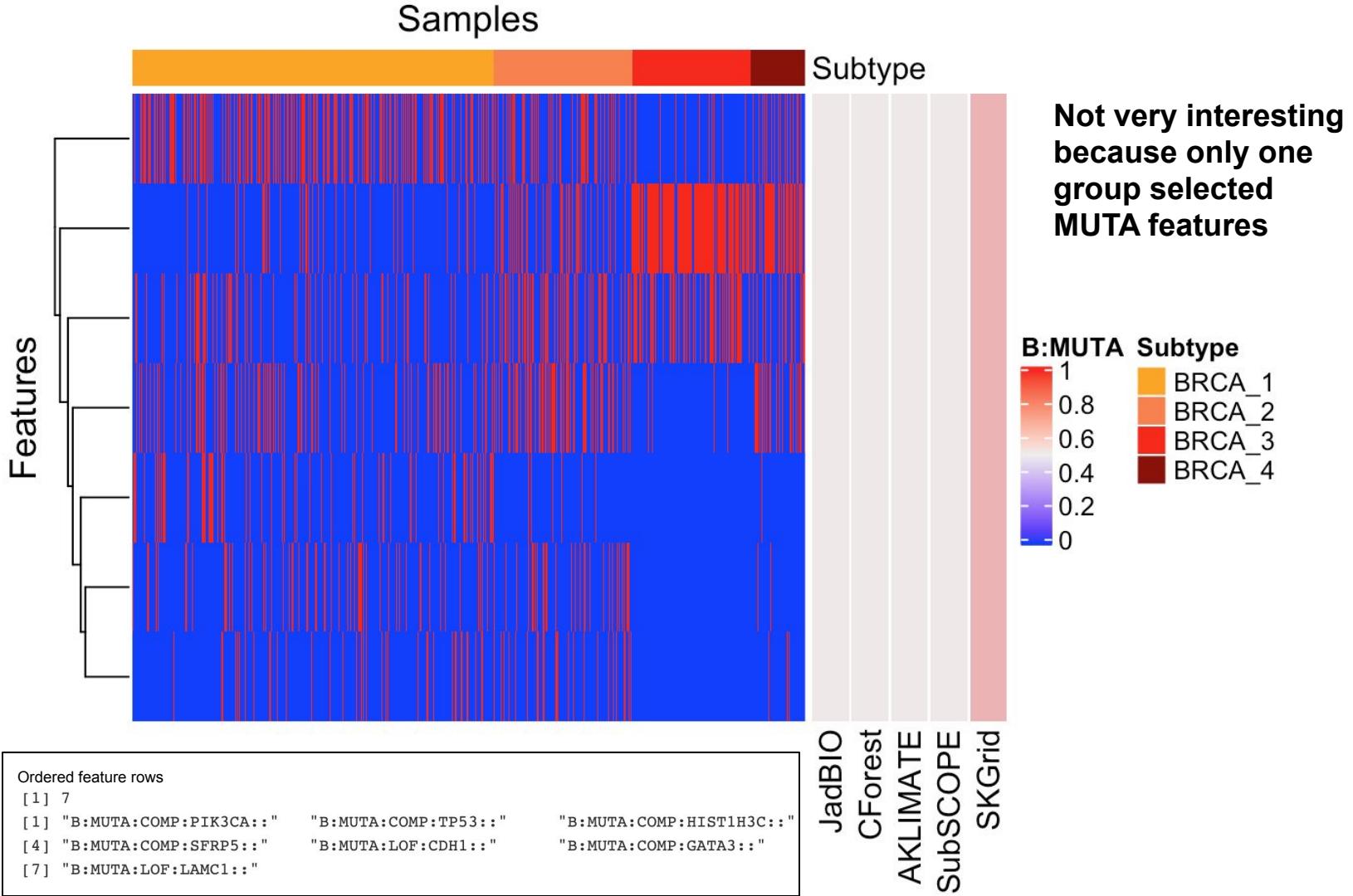


Ordered feature row

```
[1] 18
[1] "I:CNVR::hsa-mir-548d-1:-1533:" "I:CNVR::BRICD5:283870:"
[3] "I:CNVR::PSMD12:5718:" "I:CNVR::MIR4728:100616132:"
[5] "I:CNVR::SLC4A1:6521:" "I:CNVR::FOXA1:3169:"
[7] "I:CNVR::OR6C1:390321:" "I:CNVR::LOC644936:644936:"
[9] "I:CNVR::CRSP8P:441089:" "I:CNVR::RFC1:5981:"
[11] "I:CNVR::TREM1:54210:" "I:CNVR::ACTR2:10097:"
[13] "I:CNVR::CYP4A22:284541:" "I:CNVR::TTC39B:158219:"
[15] "I:CNVR::DOHH:83475:" "I:CNVR::LZTS1:11178:"
[17] "I:CNVR::ZNRF3:84133:" "I:CNVR::GOT2:2806:"
```

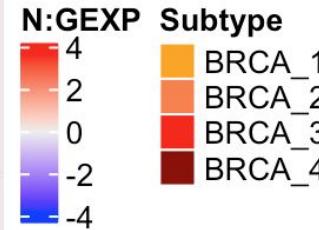
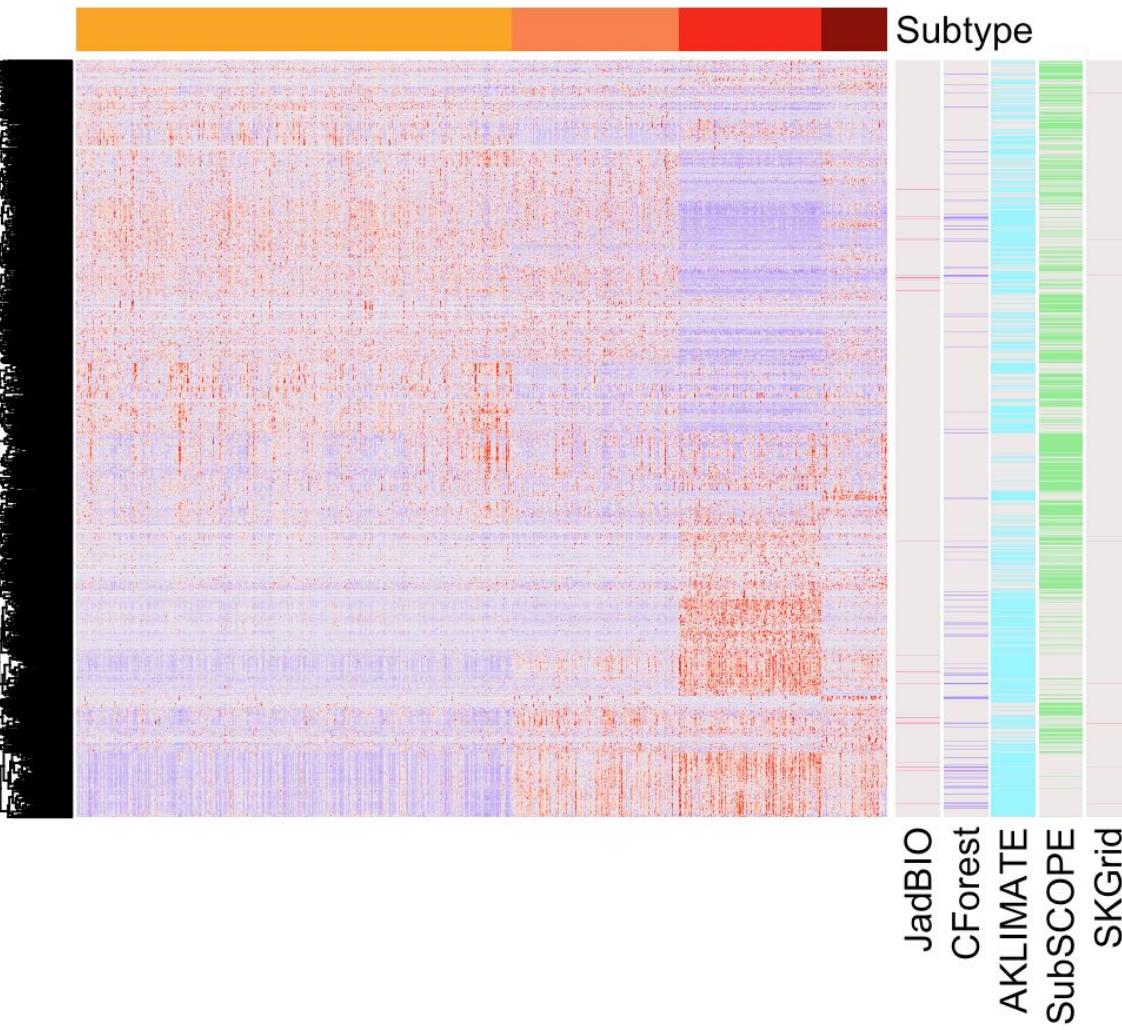






Features

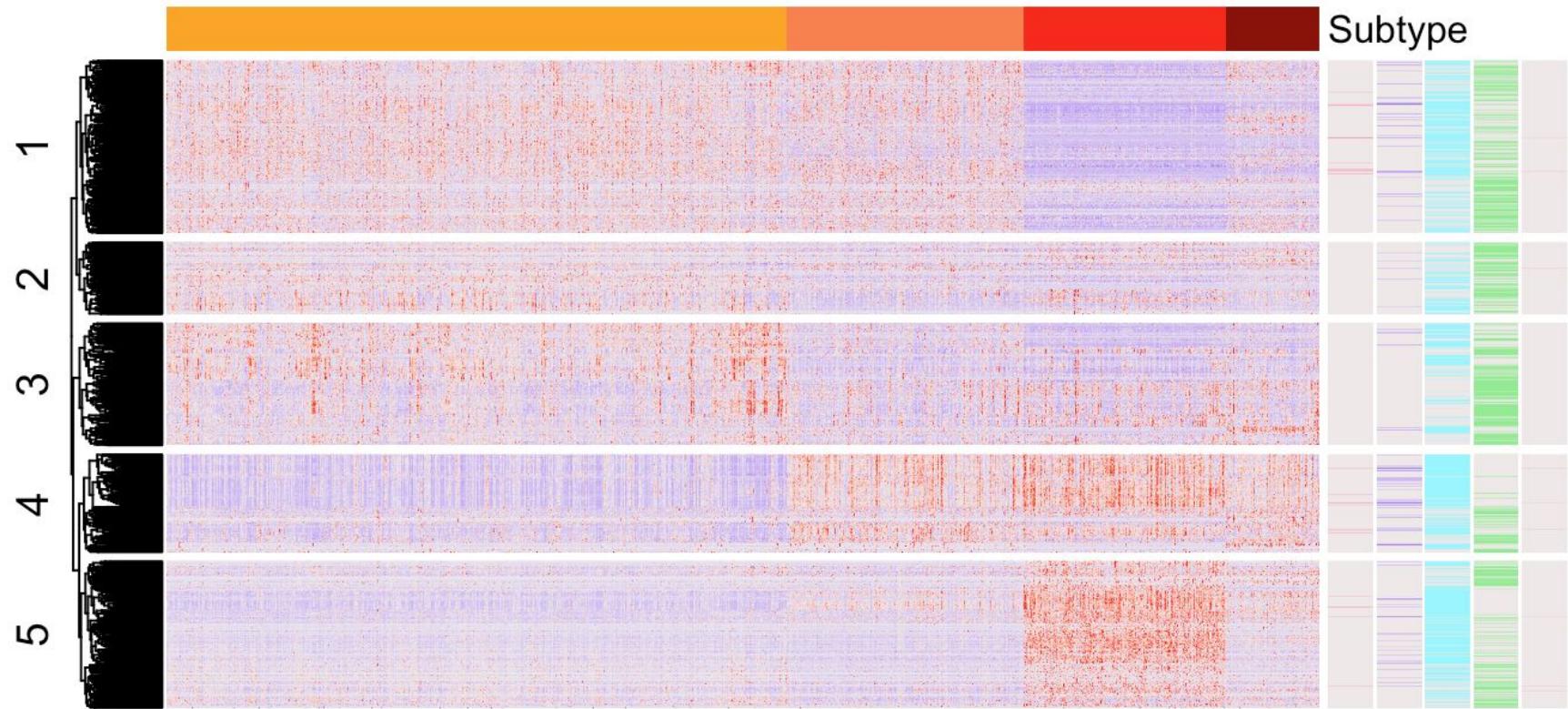
Samples



**CForest, AKLIMATE, and
SubSCOPE tend to pick features
across each cluster**

**JadBIO and SKGrid much more
sparse across clusters. They do
appear to pick features within the
same clusters as each other**

Samples



Zoom in
on GEXP

The ordering of feature rows may change from the heatmap on the last slide

But the hierarchical structure remains unchanged (same relationships and clustering distances)

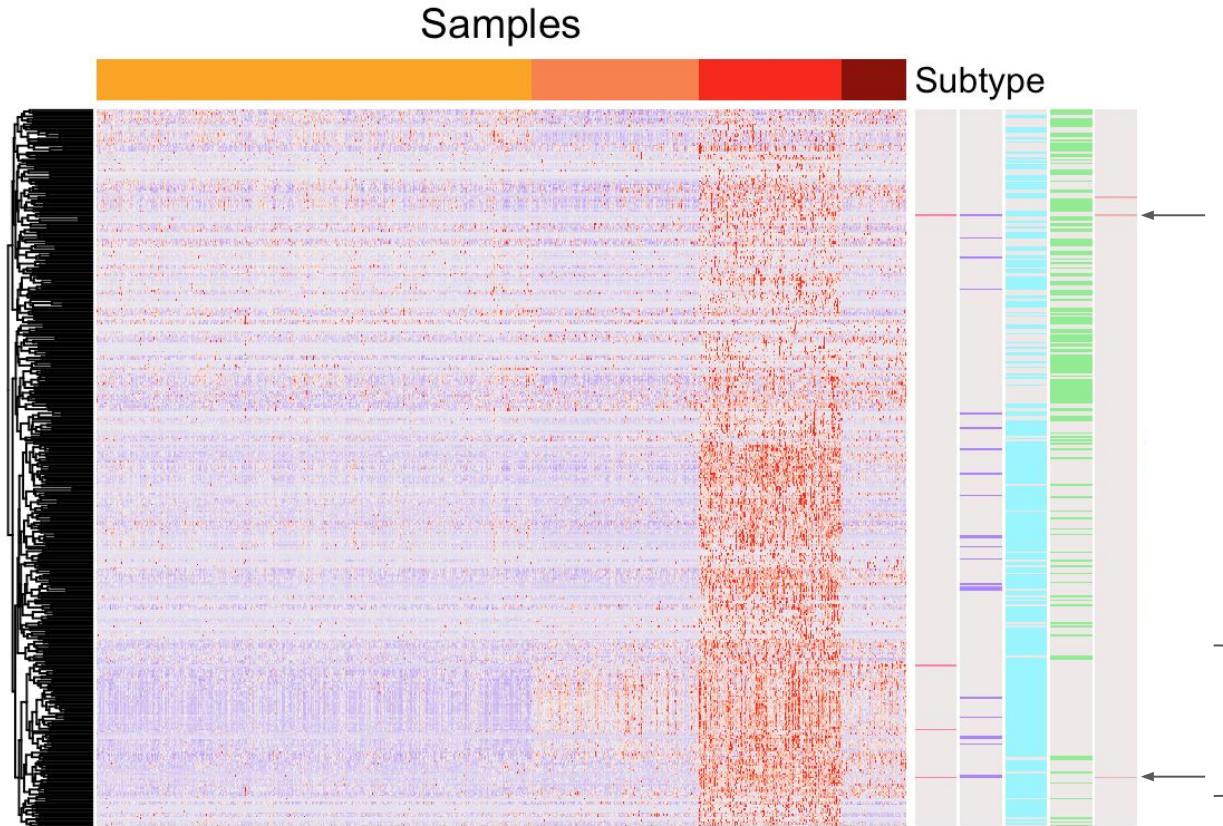
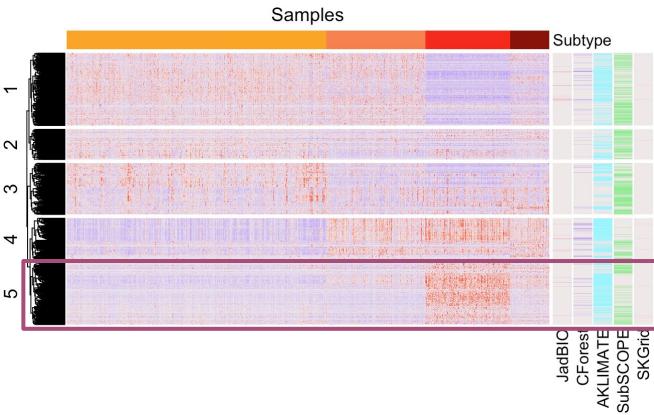
JadBIO
CForest
AKLIMATE
SubSCOPE
SKGrid

Cluster includes all 5 teams

JadBIO and SKGrid are very sparse compared to other groups

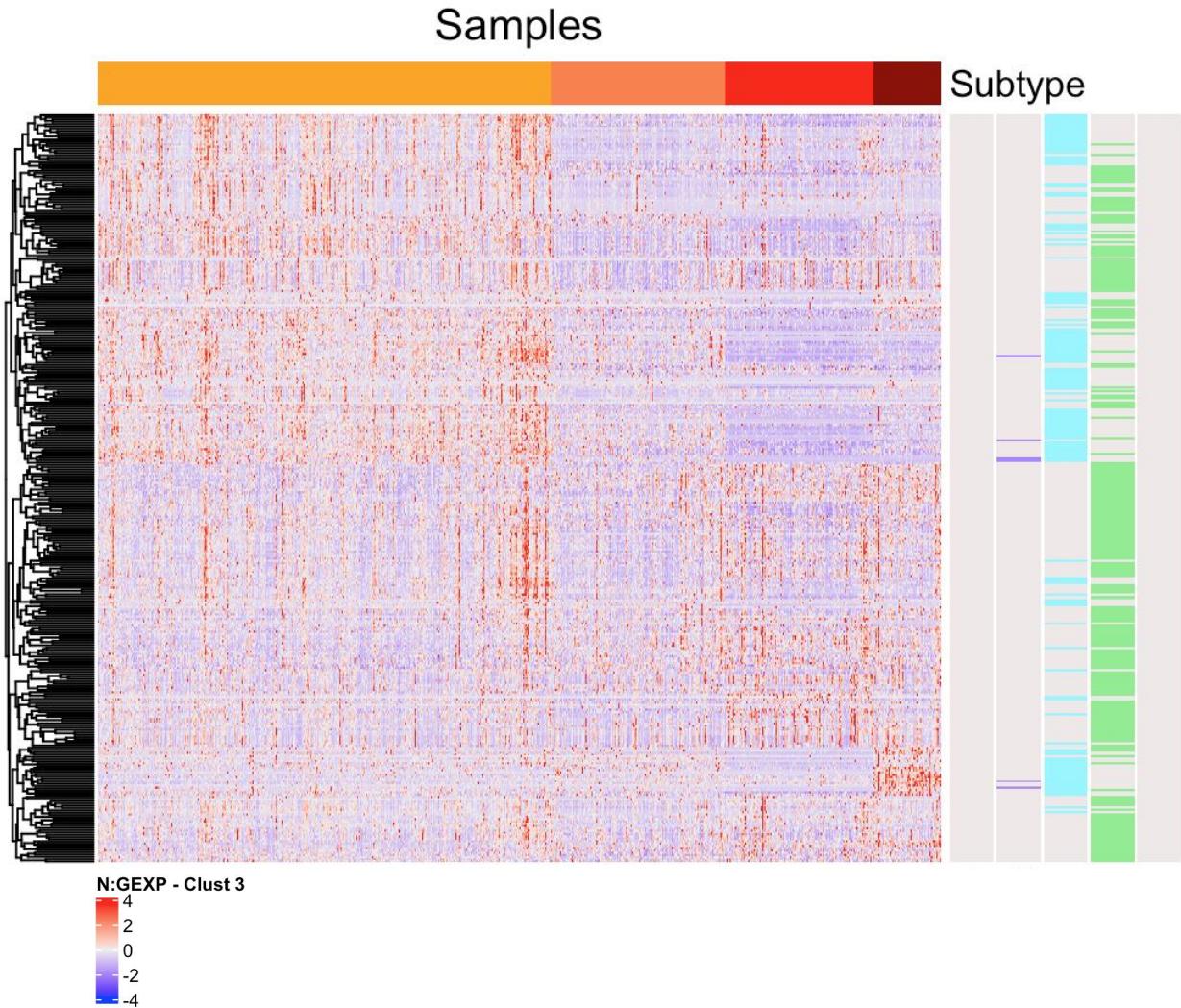
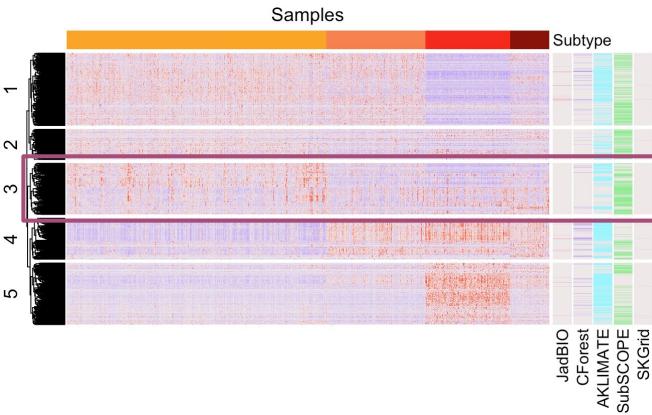
Some overlap at arrows (exact fts and/or similar fts)

Visually easier to see heatmap blocks in last several rows (note all teams select features in this cluster region)



Cluster contains no features selected by JadBIO or SKGrid

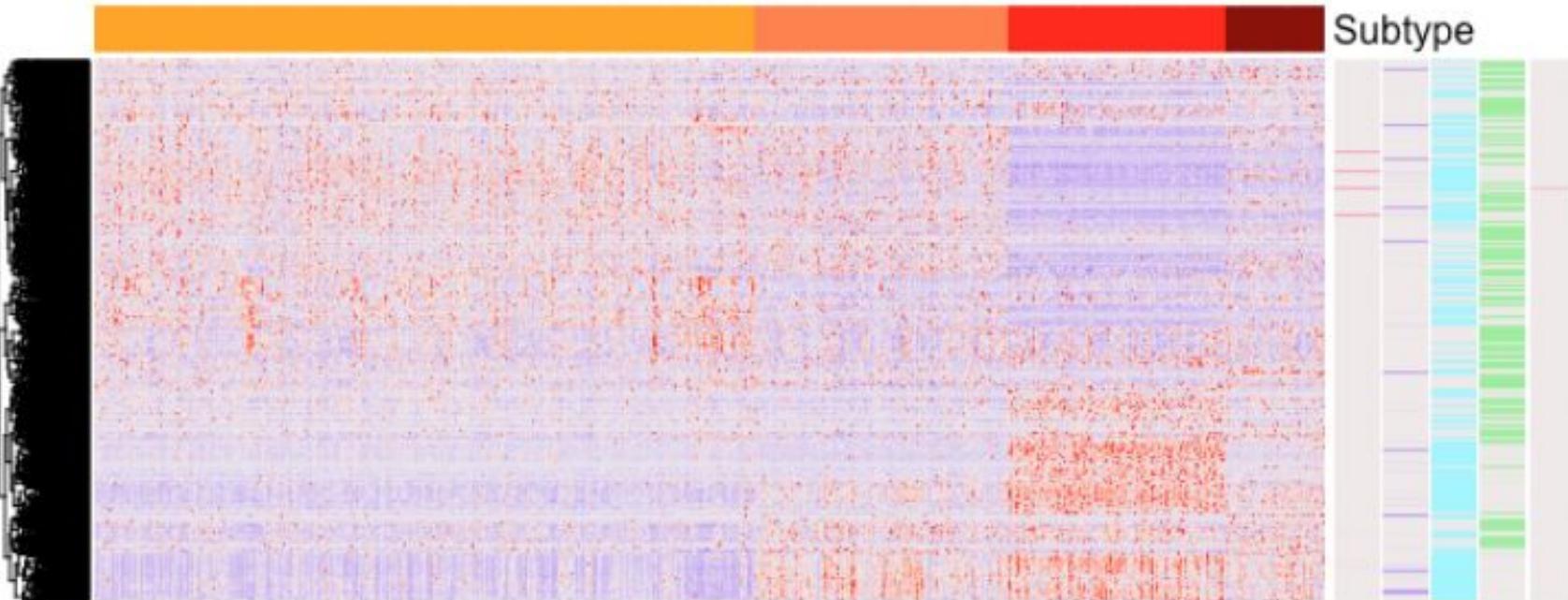
AKLIMATE and SubSCOPE tend to select related features



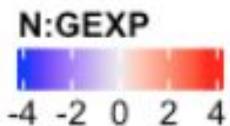
Additional Slides

Samples

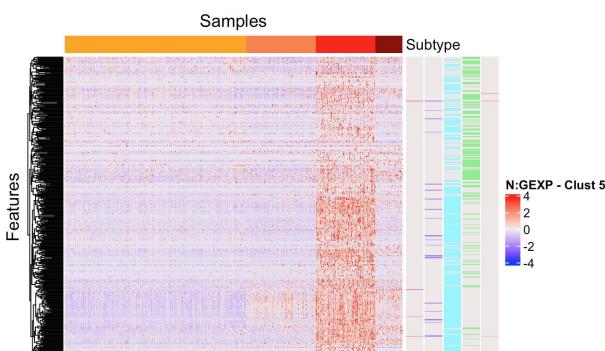
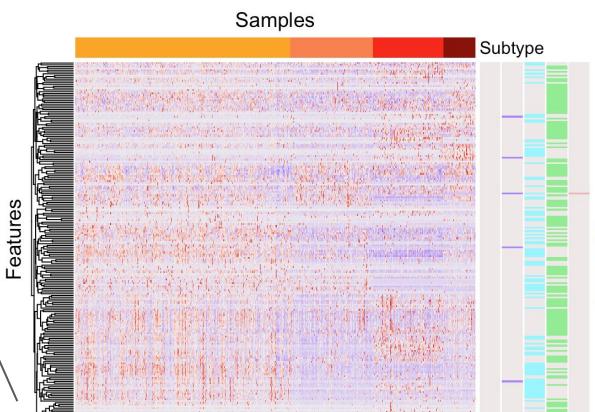
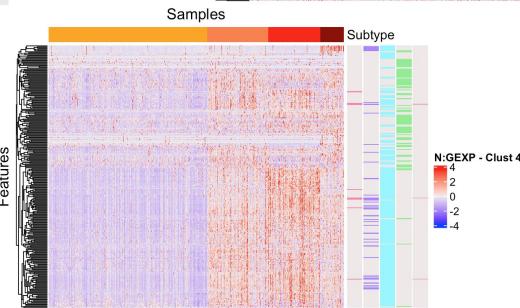
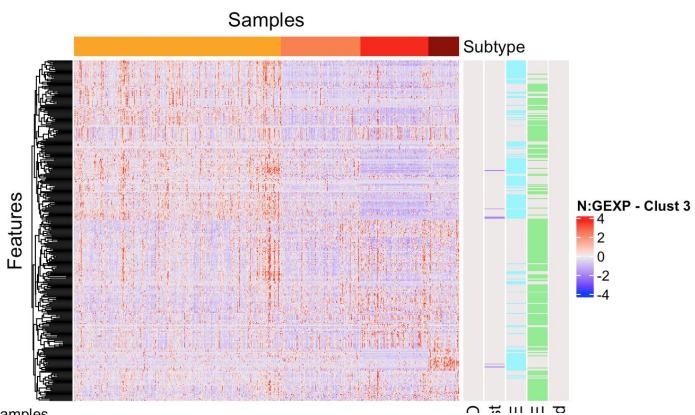
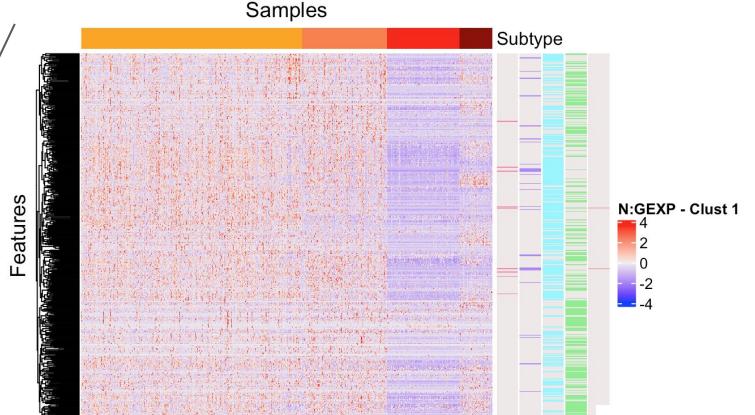
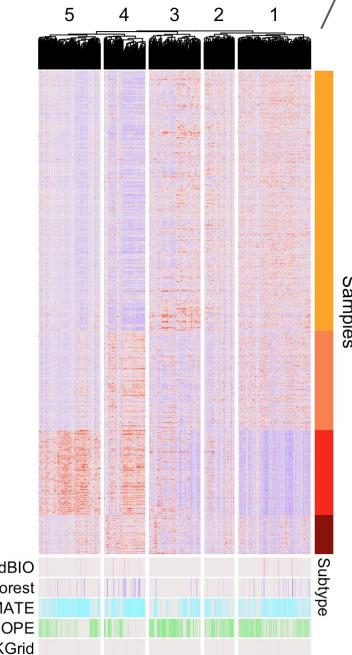
Features



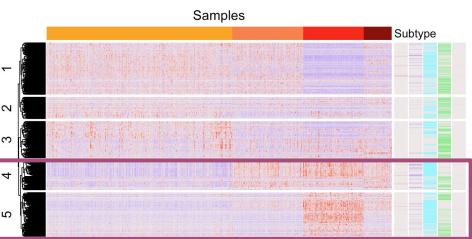
Full GEXP heatmap



All clusters when
have 5 total
clusters

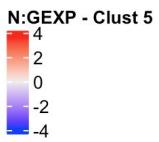
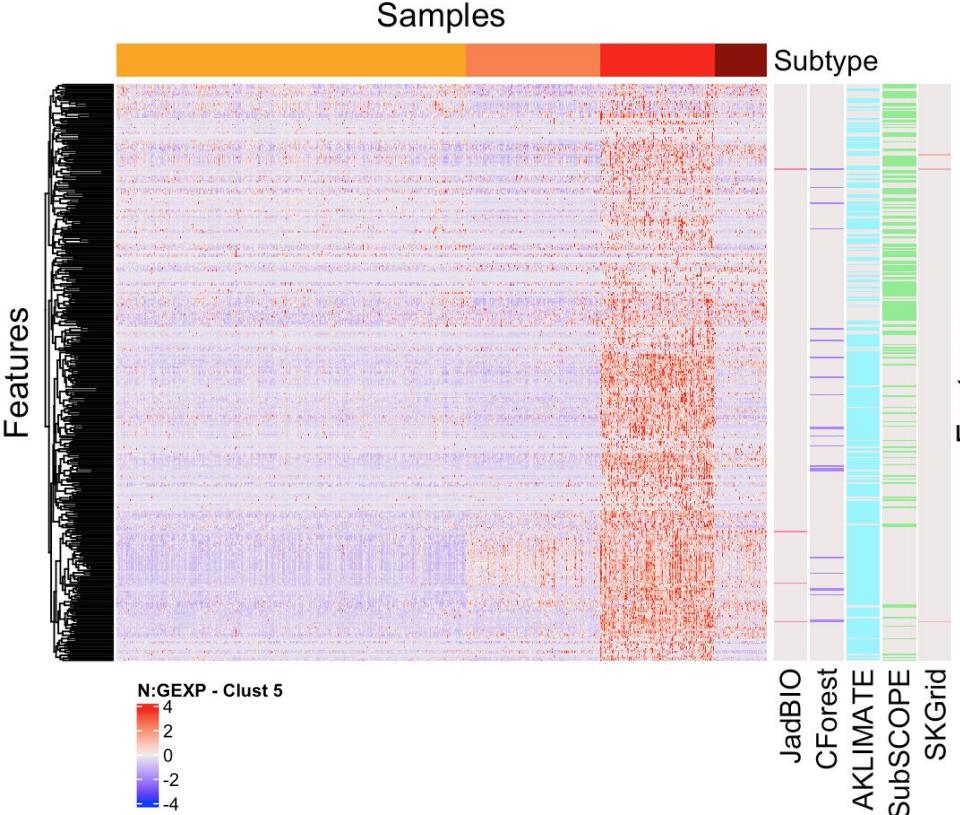


With 5 groupings



Samples

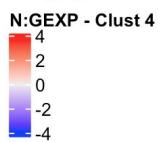
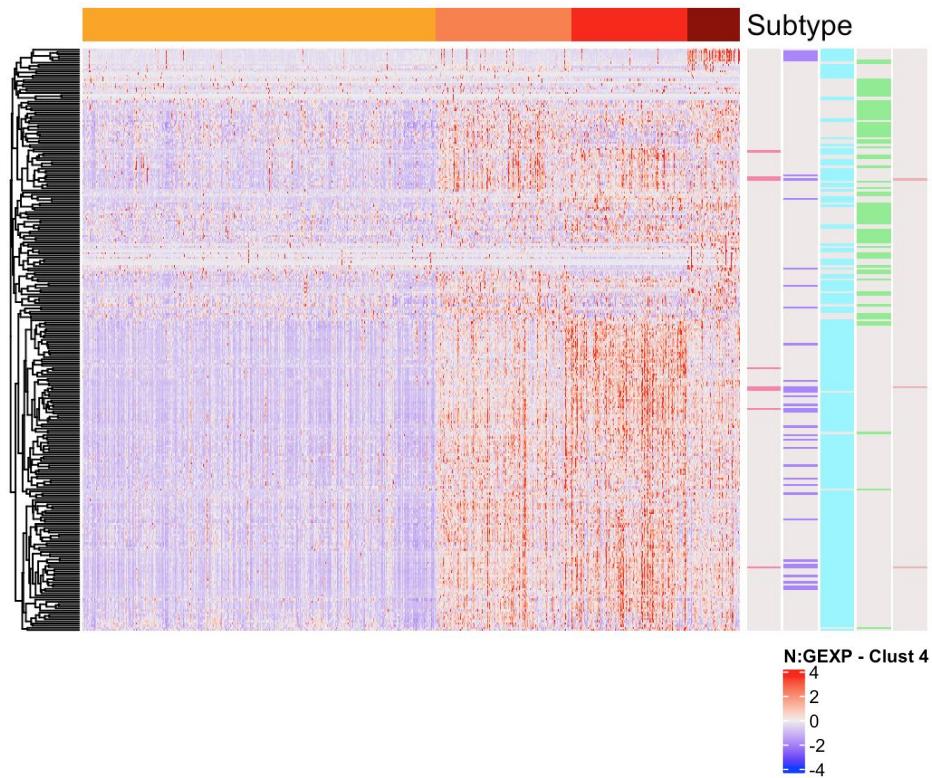
Subtype



JadBIO
CForest
AKLIMATE
SubSCOPE
SKGrid

Samples

Subtype

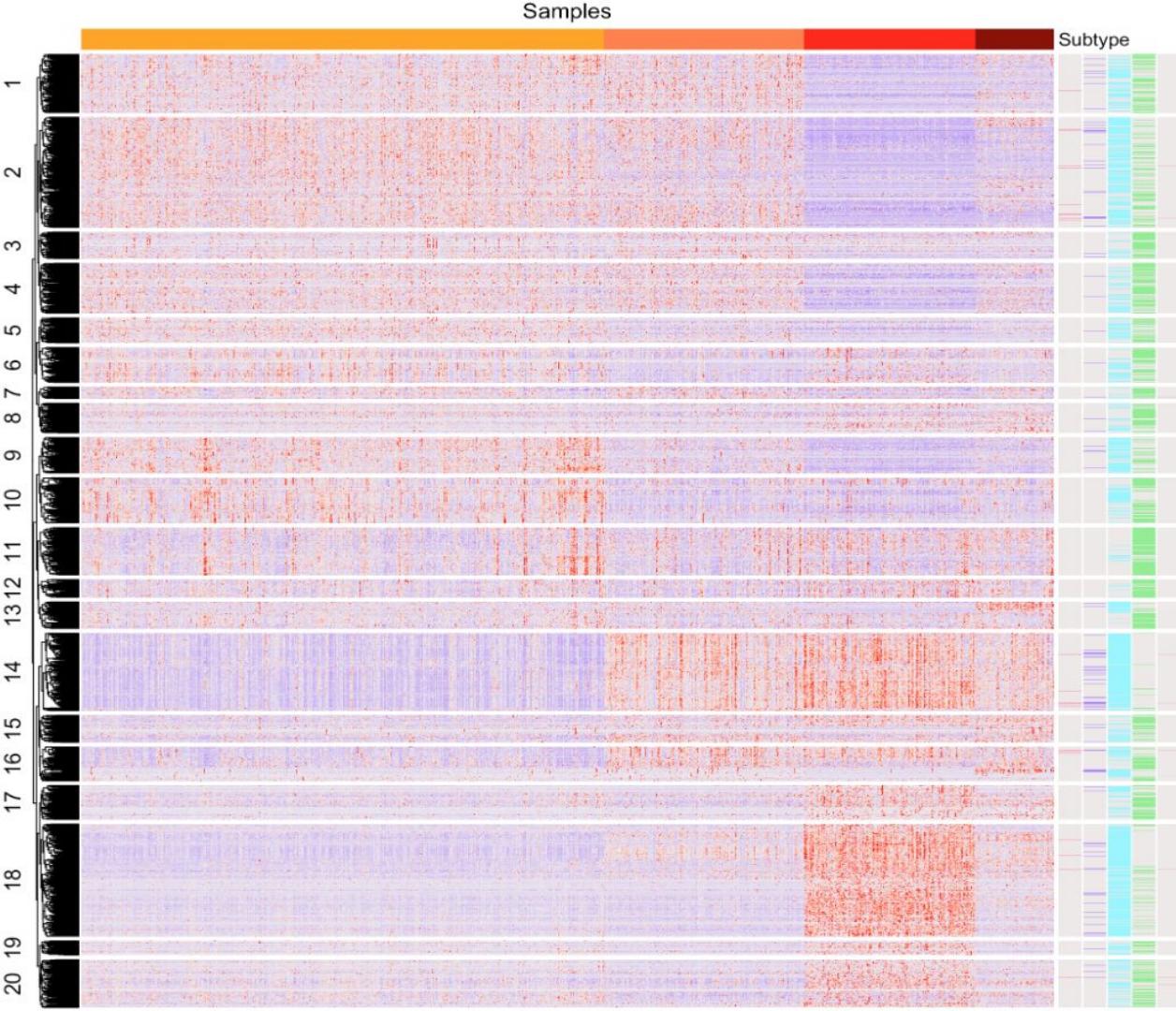


Zoom in on GEXP

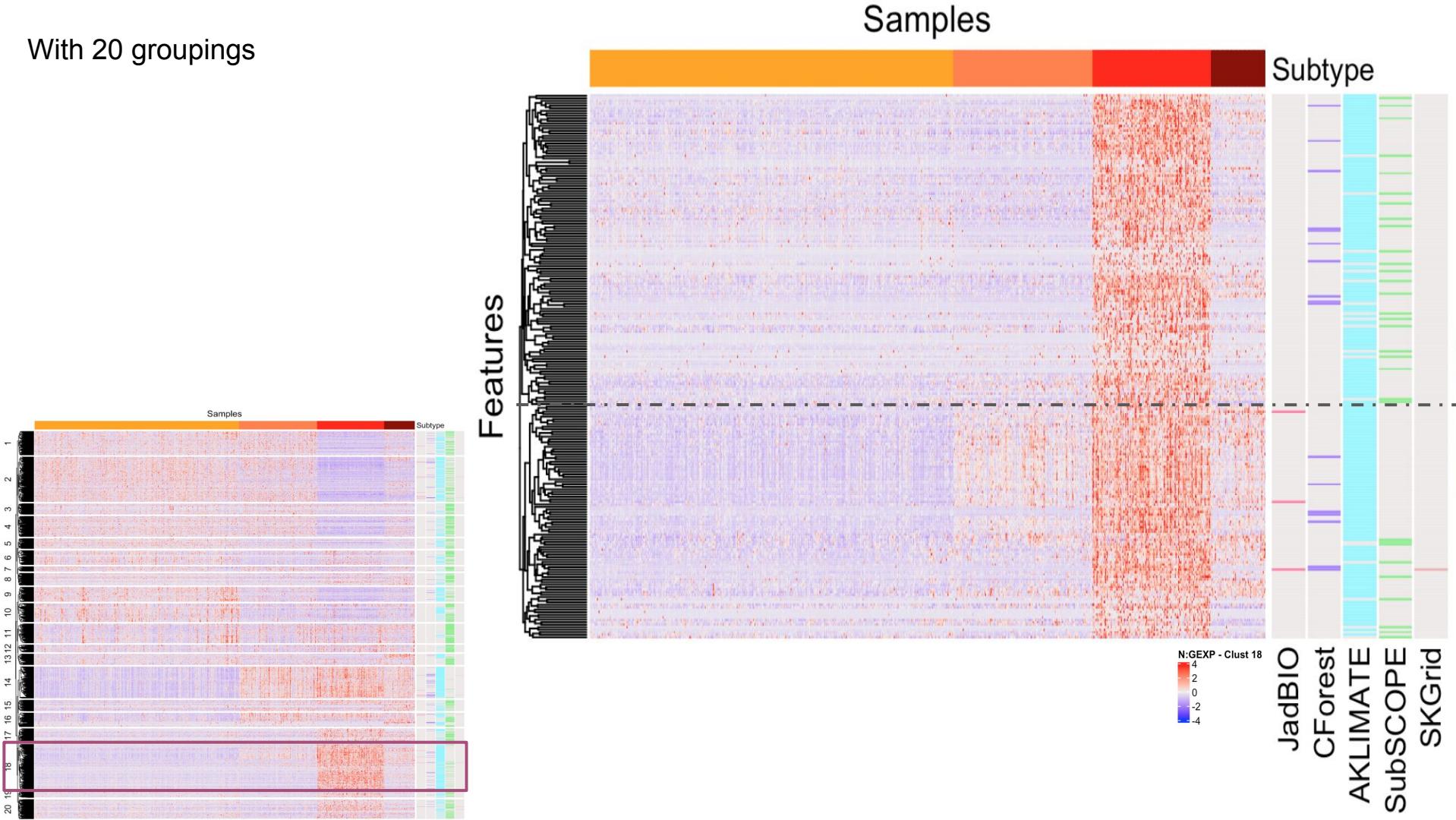
The ordering of feature rows may change from the heatmap on the last slide

But the hierarchical structure remains unchanged (same relationships and clustering distances)

Cluster 14 and 18 on next slide



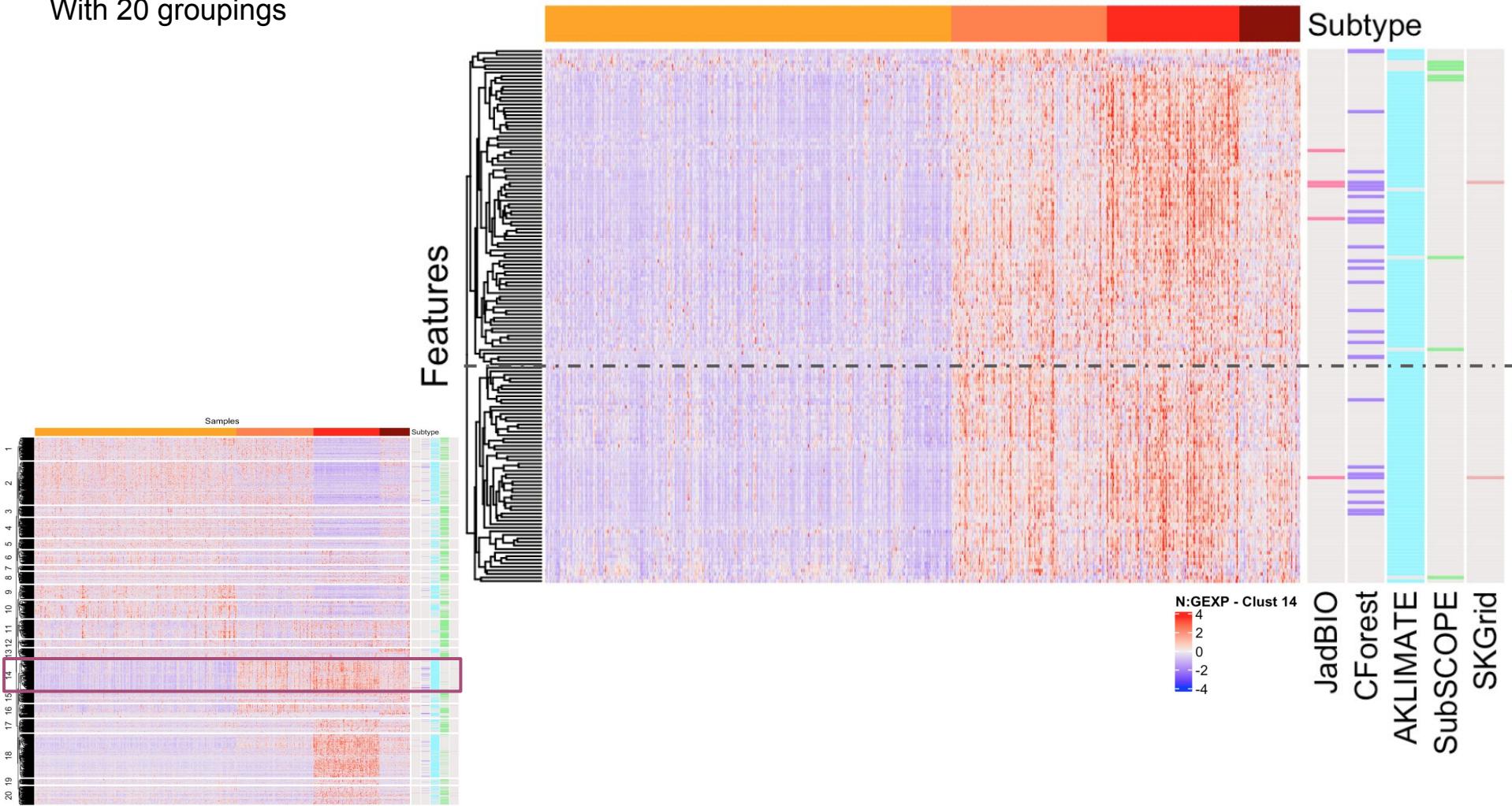
With 20 groupings

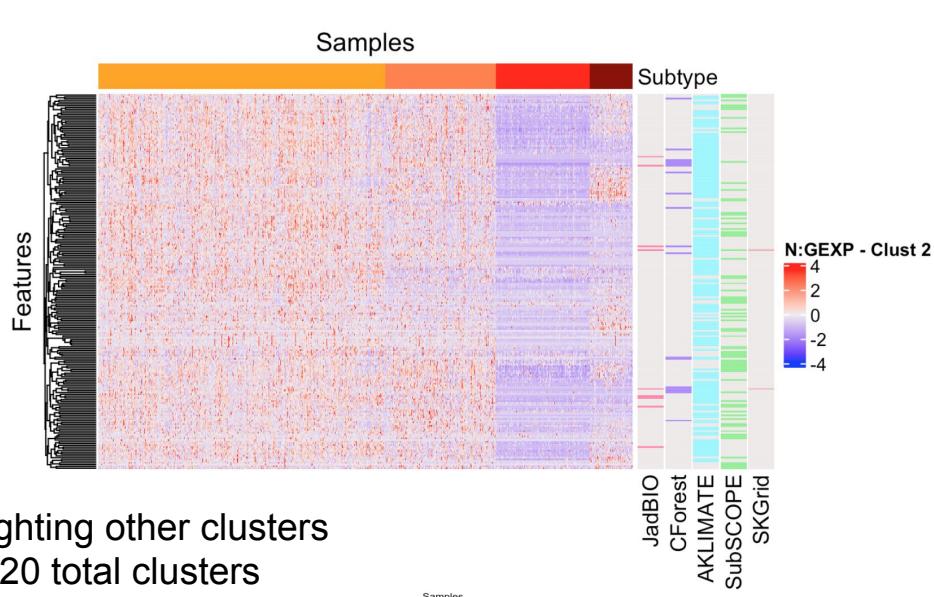
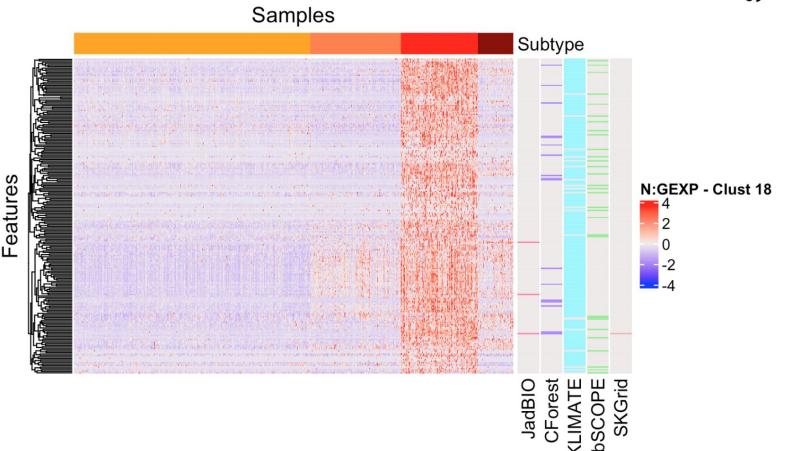
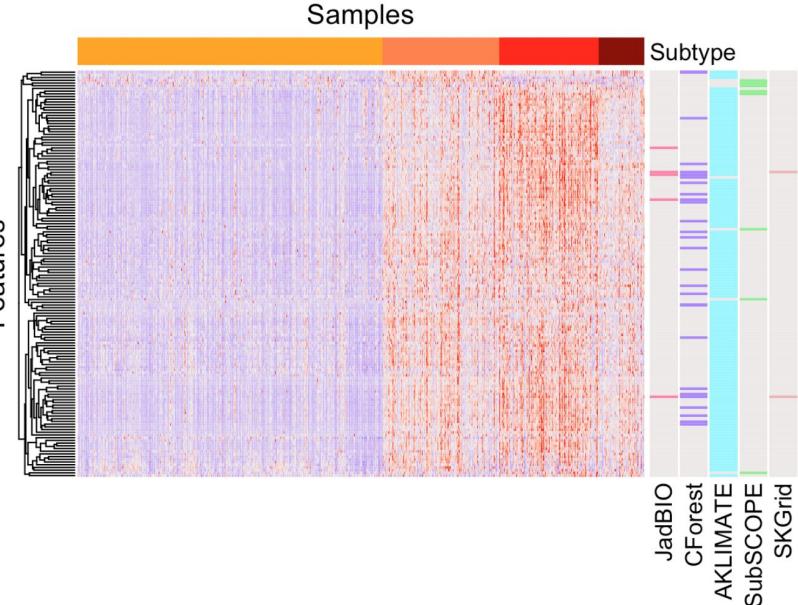


With 20 groupings

Samples

Subtype





Highlighting other clusters
when 20 total clusters

