

Data
Collection



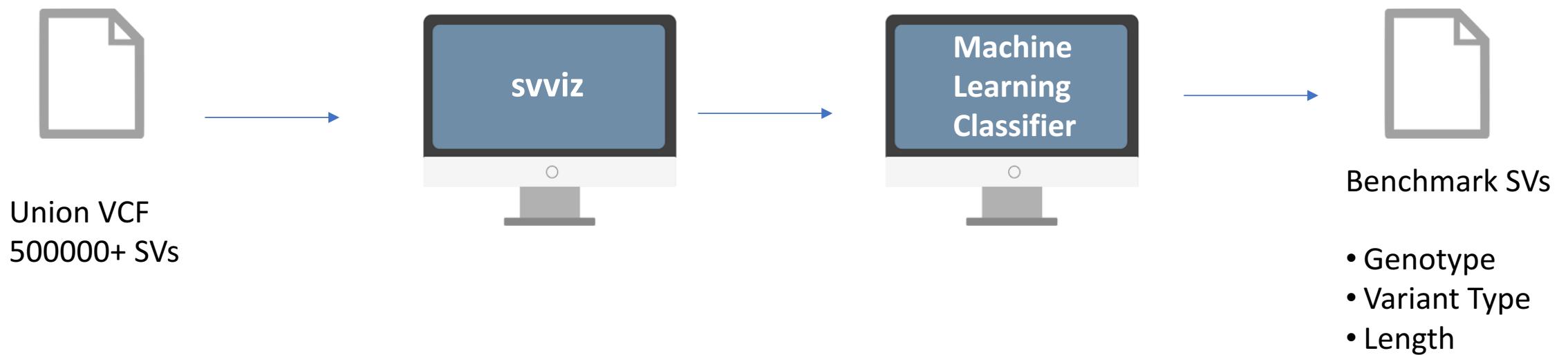
Feature
Generation

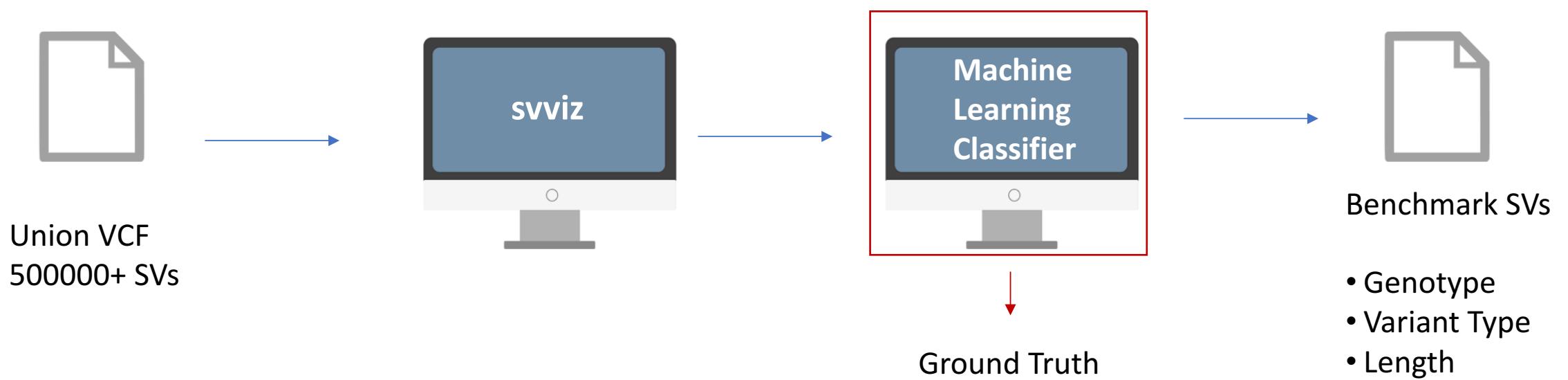


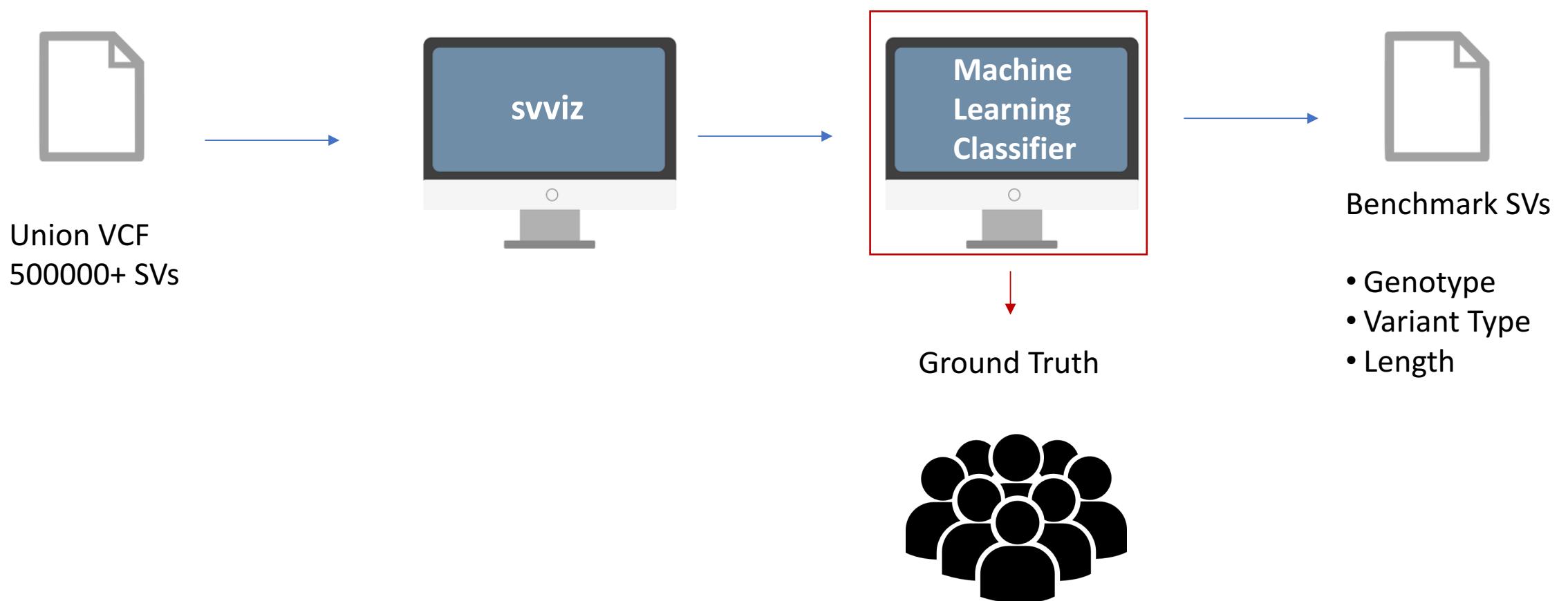
Generate
Labels



Machine
Learning

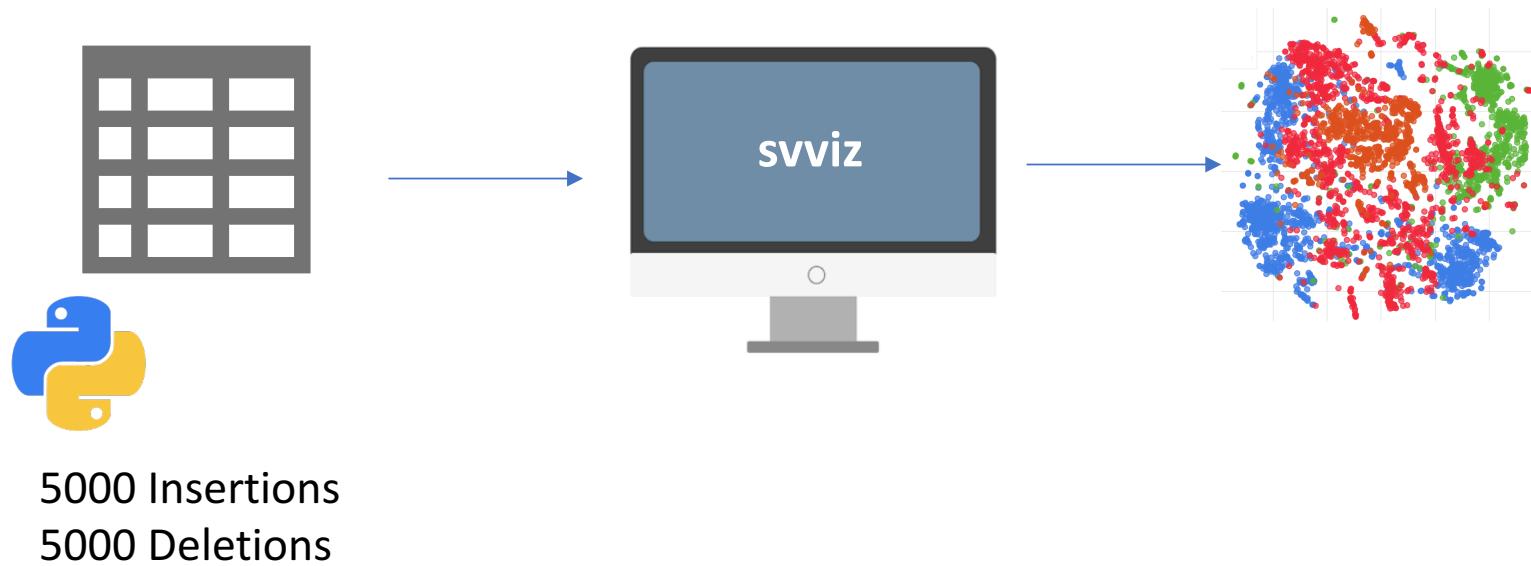




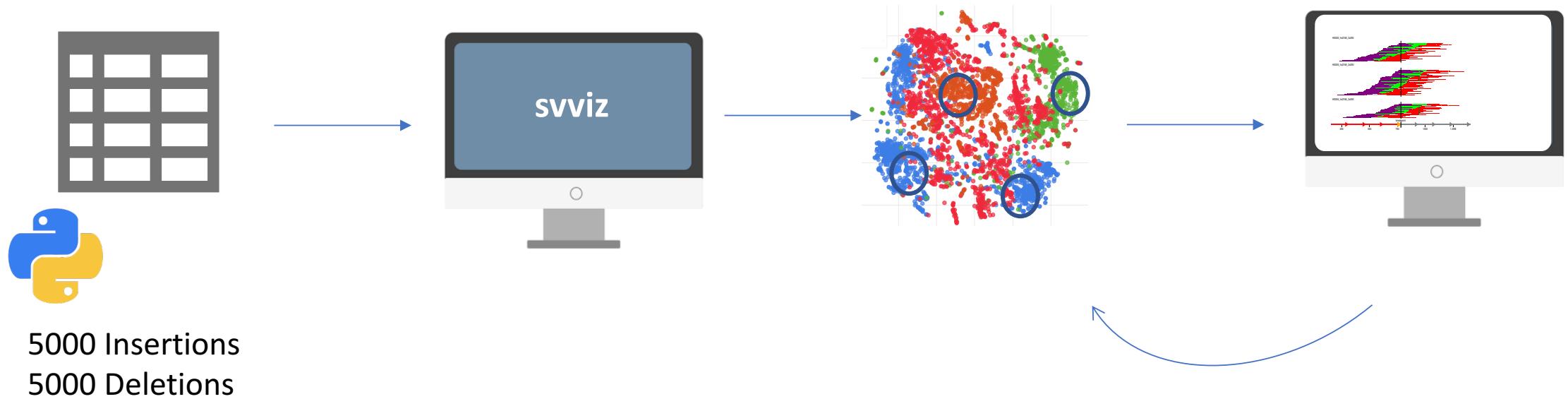


Generate labeled data

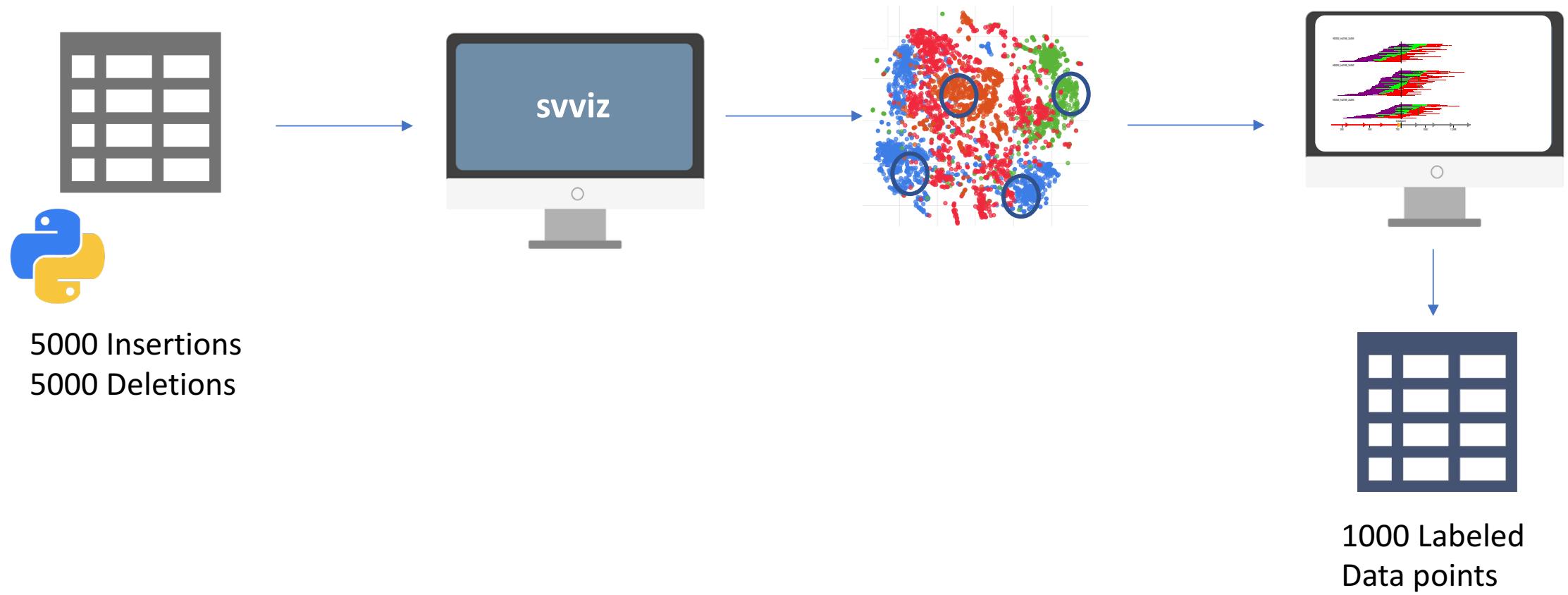
Overview



Overview

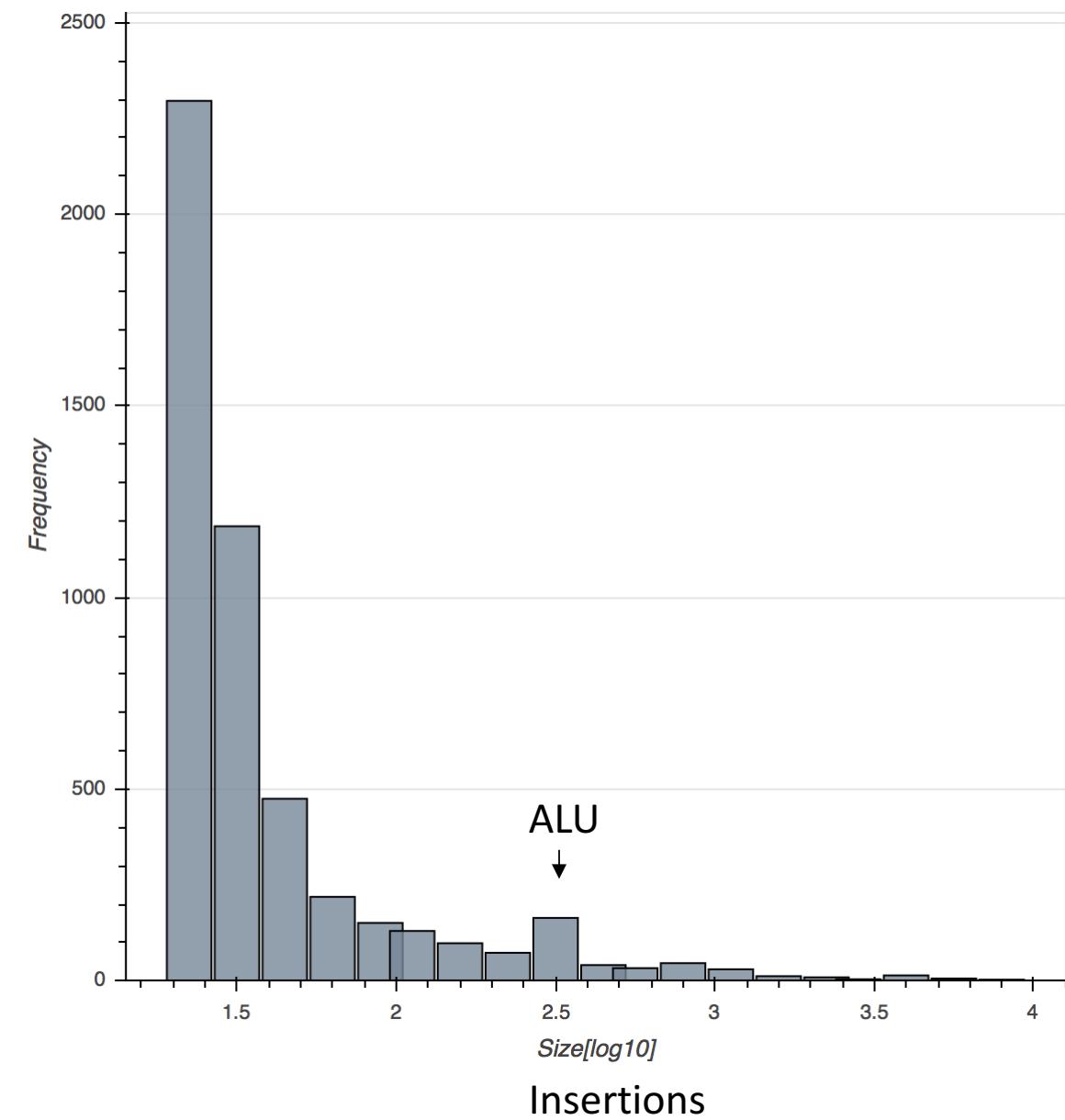
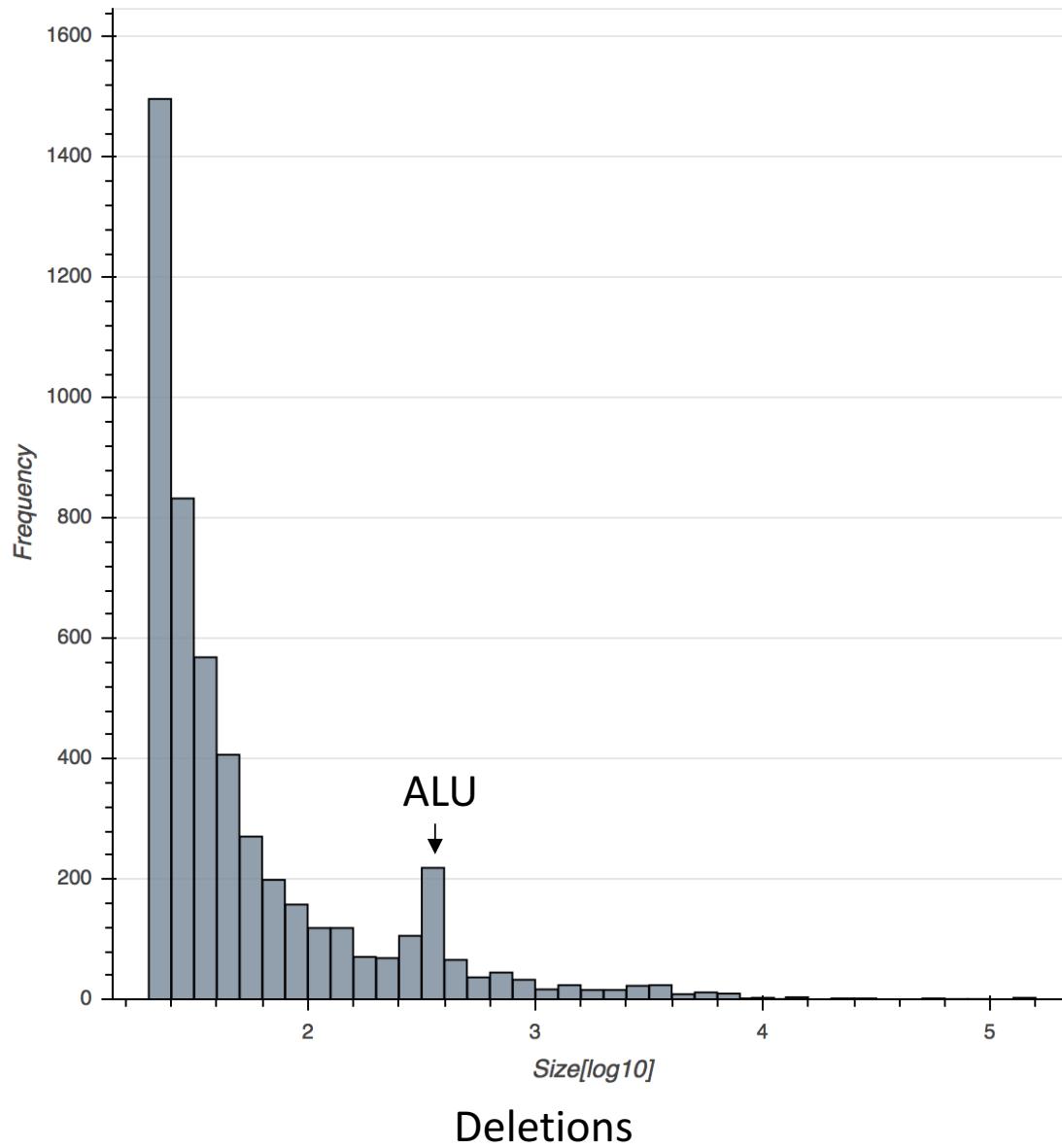


Overview

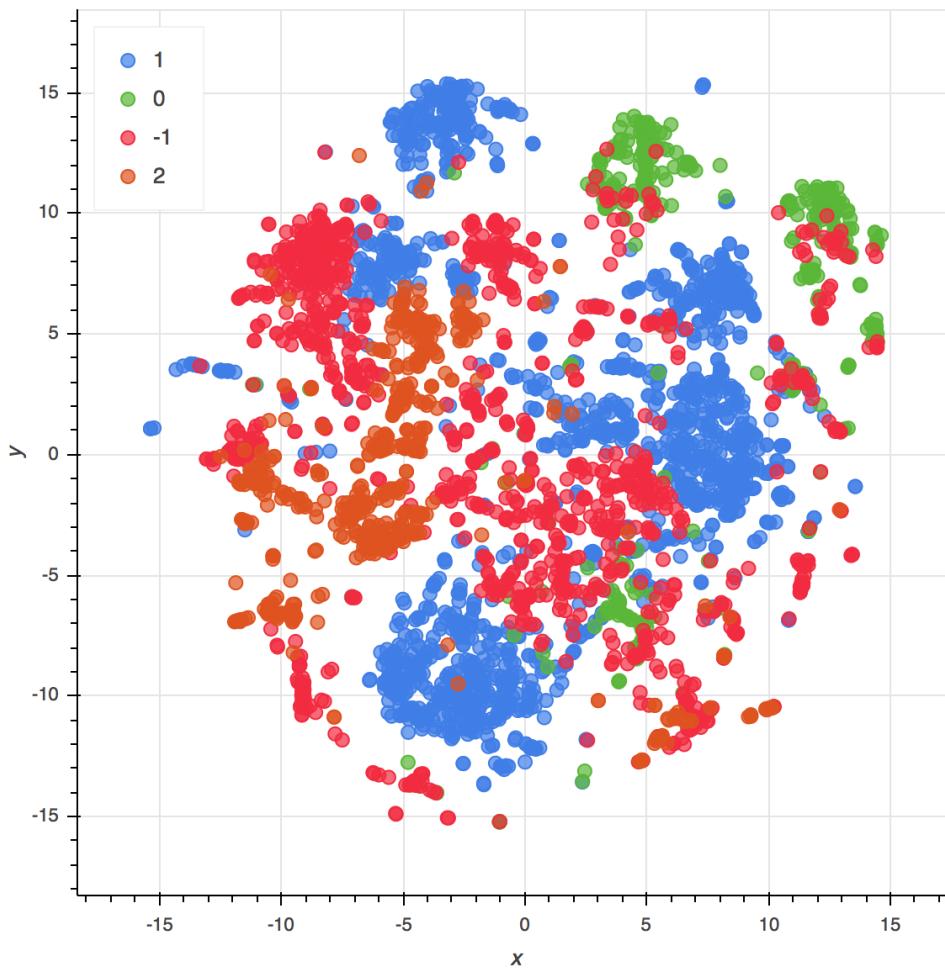


Results

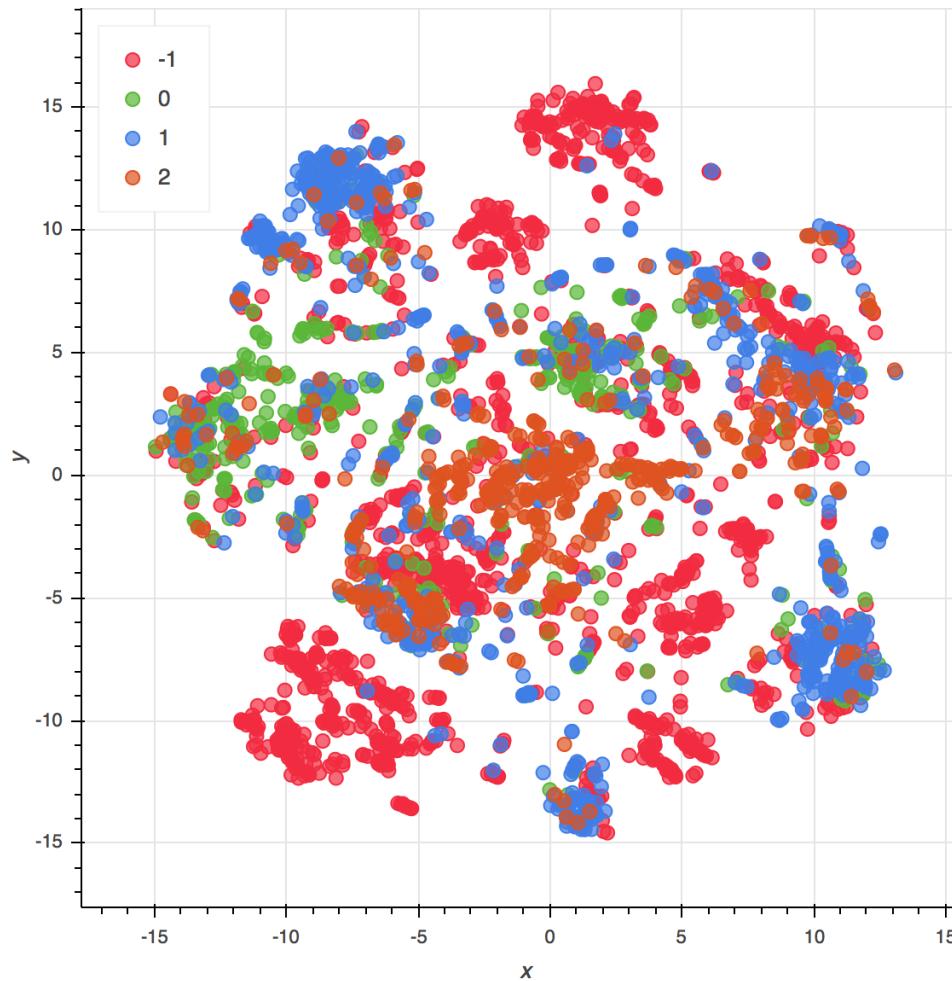
Size



Estimated Consensus GT



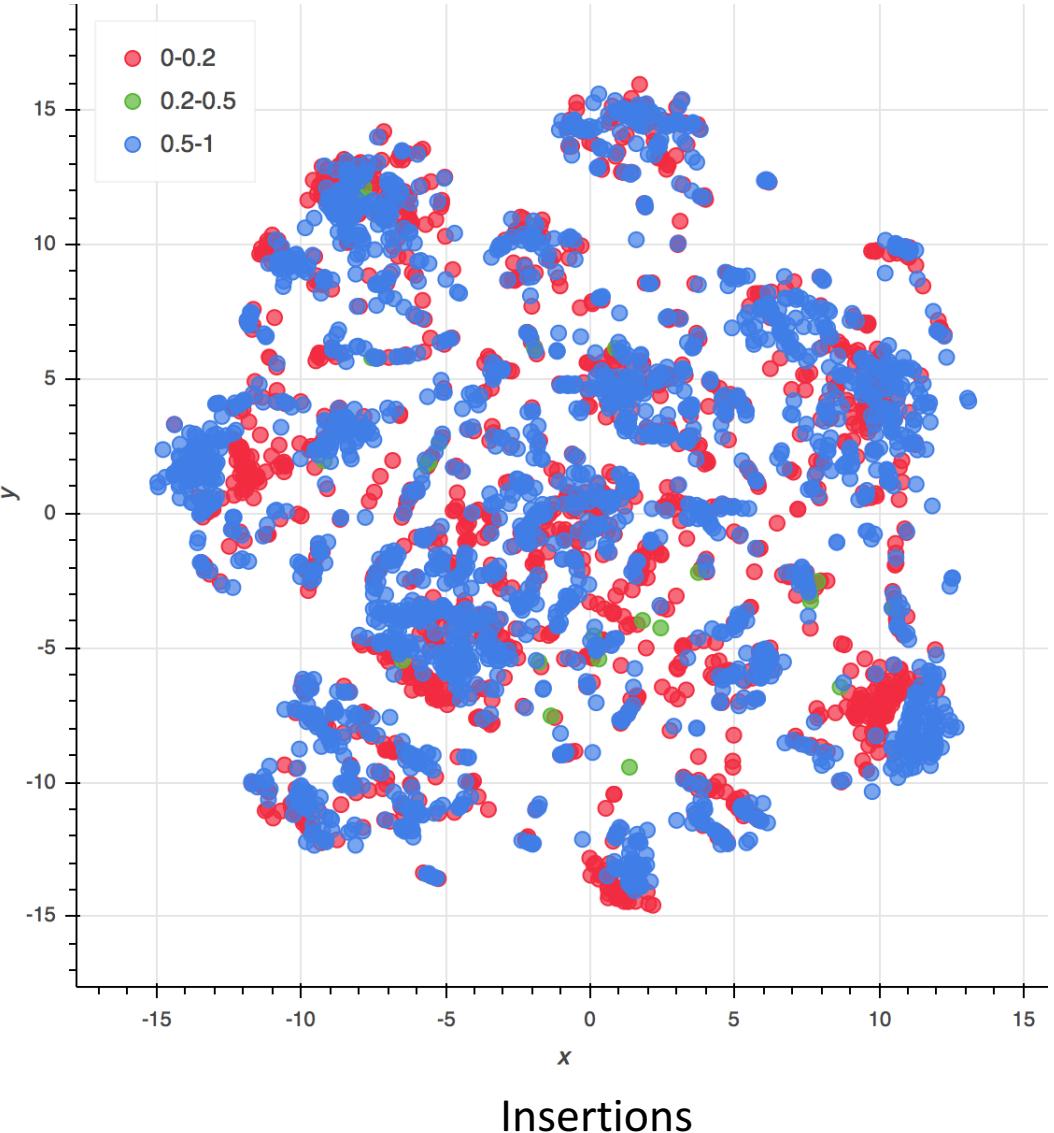
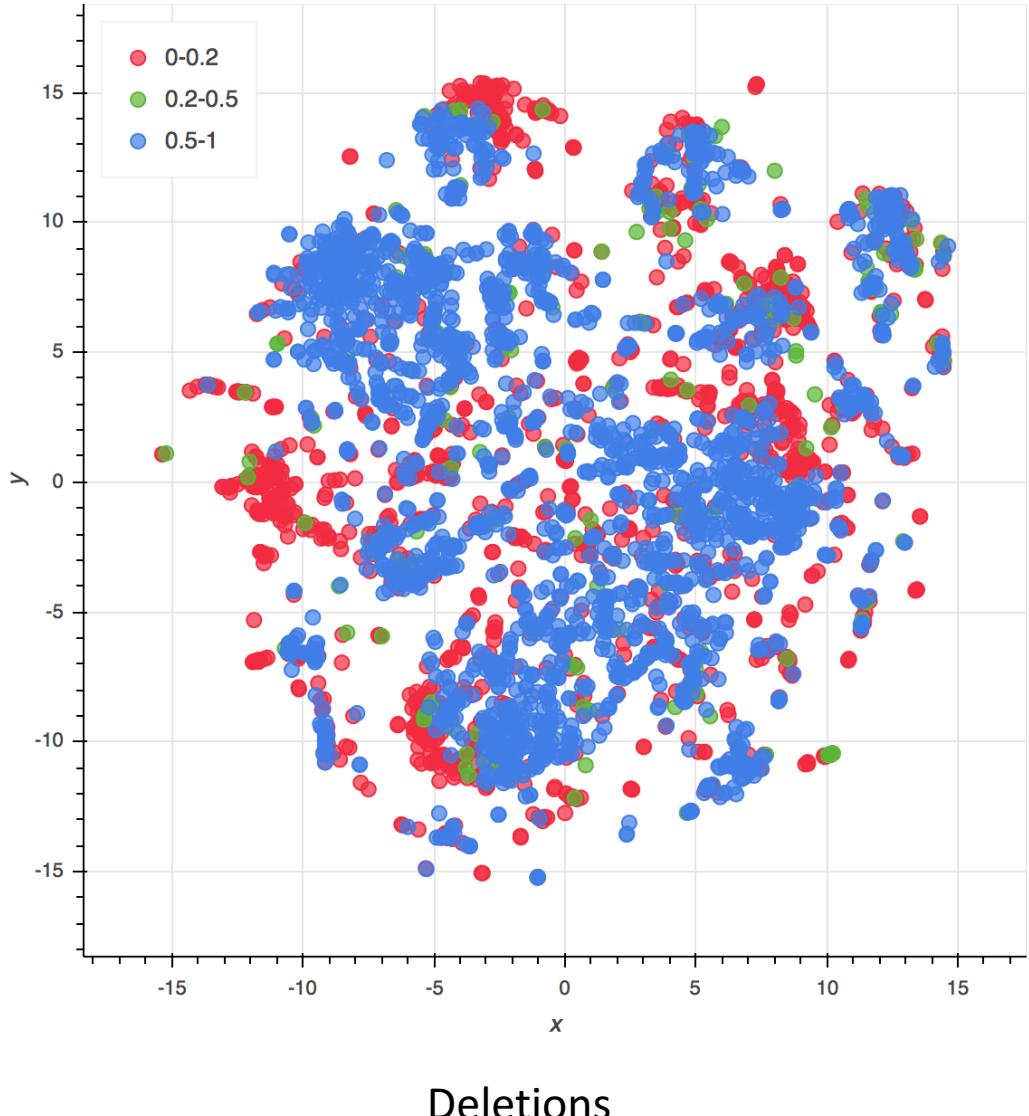
Deletions



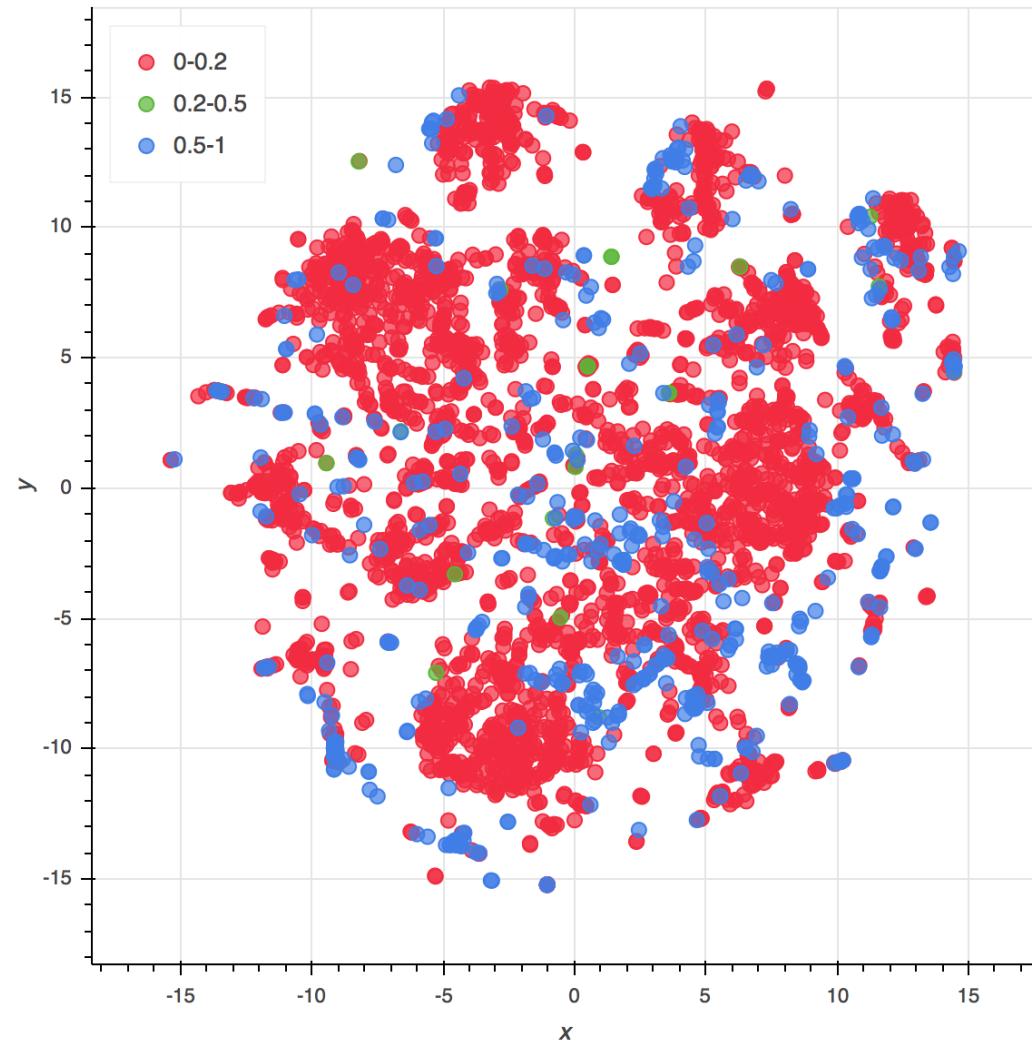
Insertions

Key	
0	Homozygous Reference
1	Heterozygous Variant
2	Homozygous Variant
-1	Undetermined

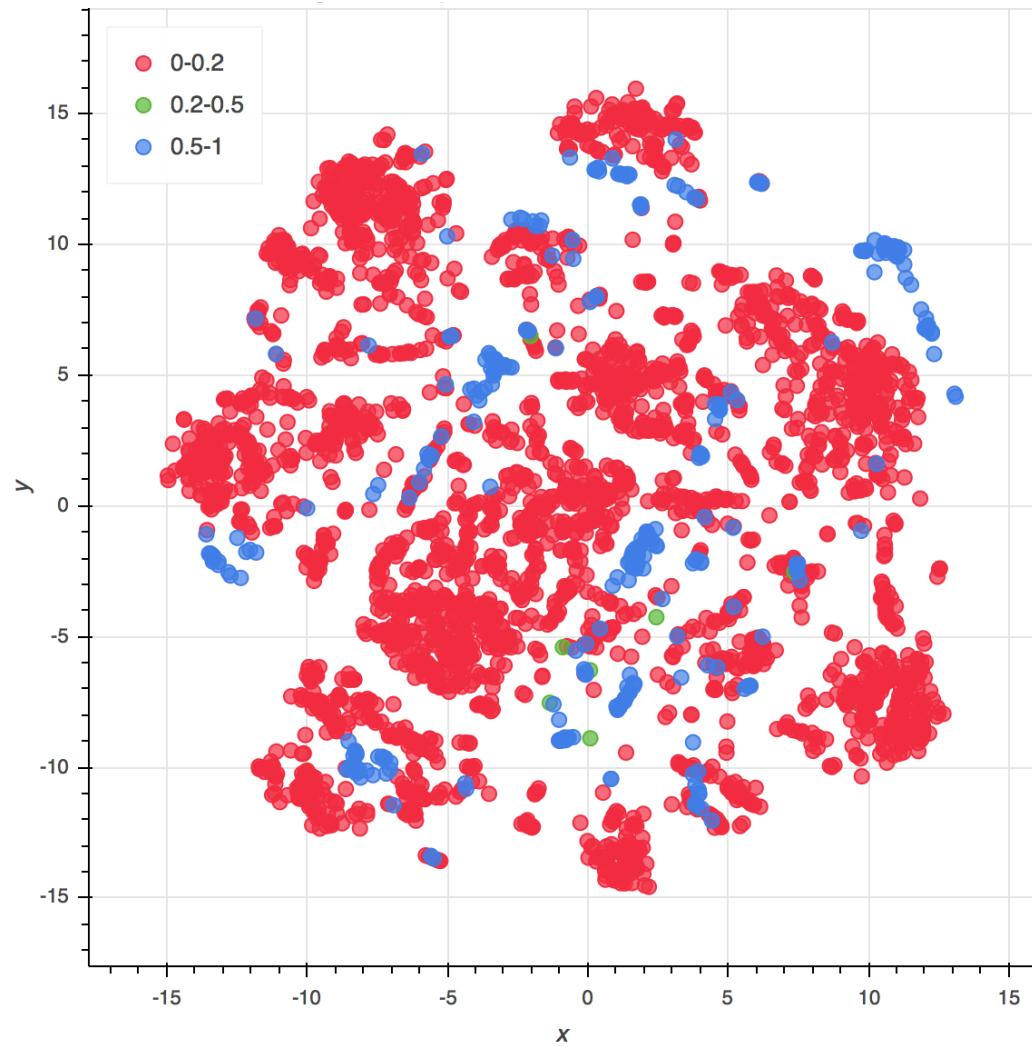
Tandem Repeats



Segmental Duplications

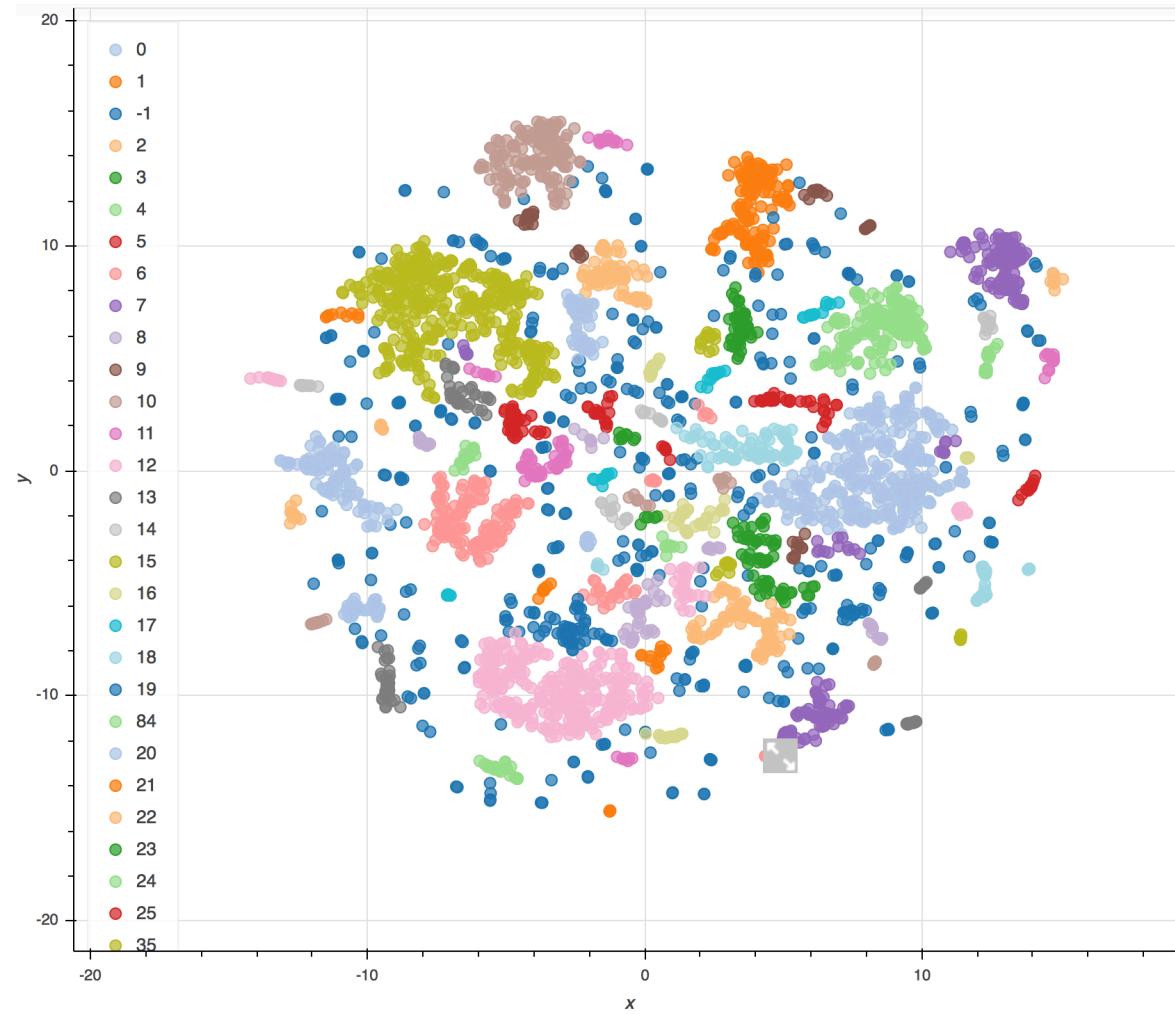


Deletions

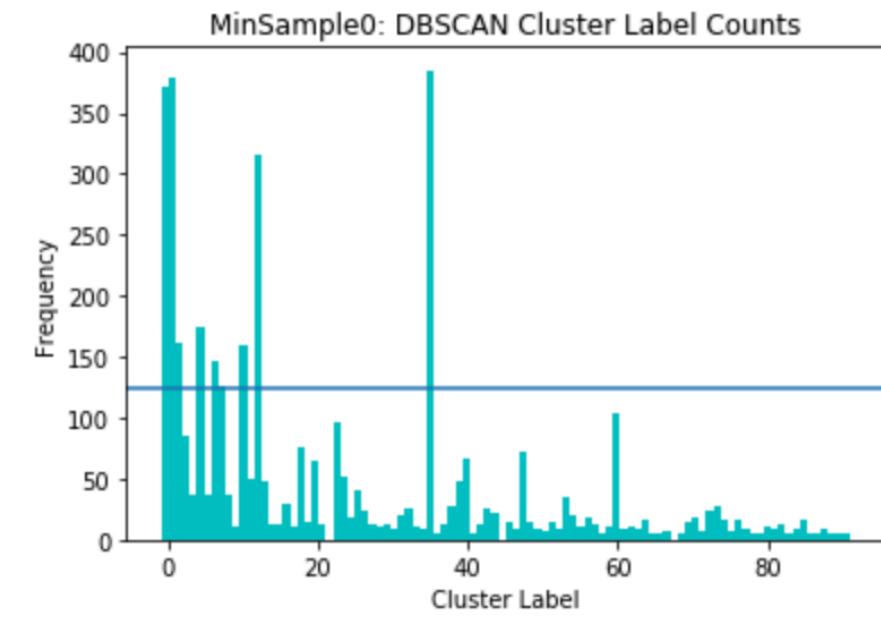


Insertions

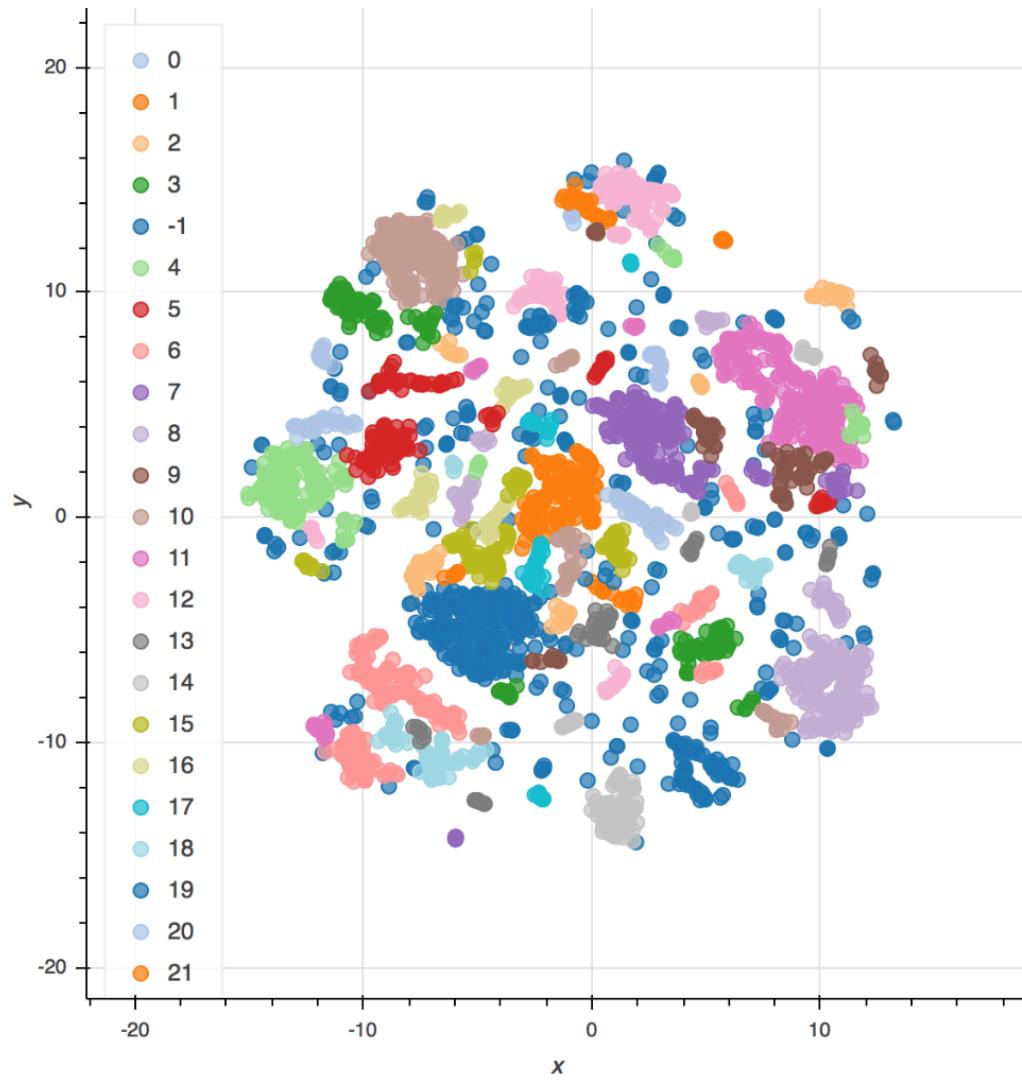
DBSCAN



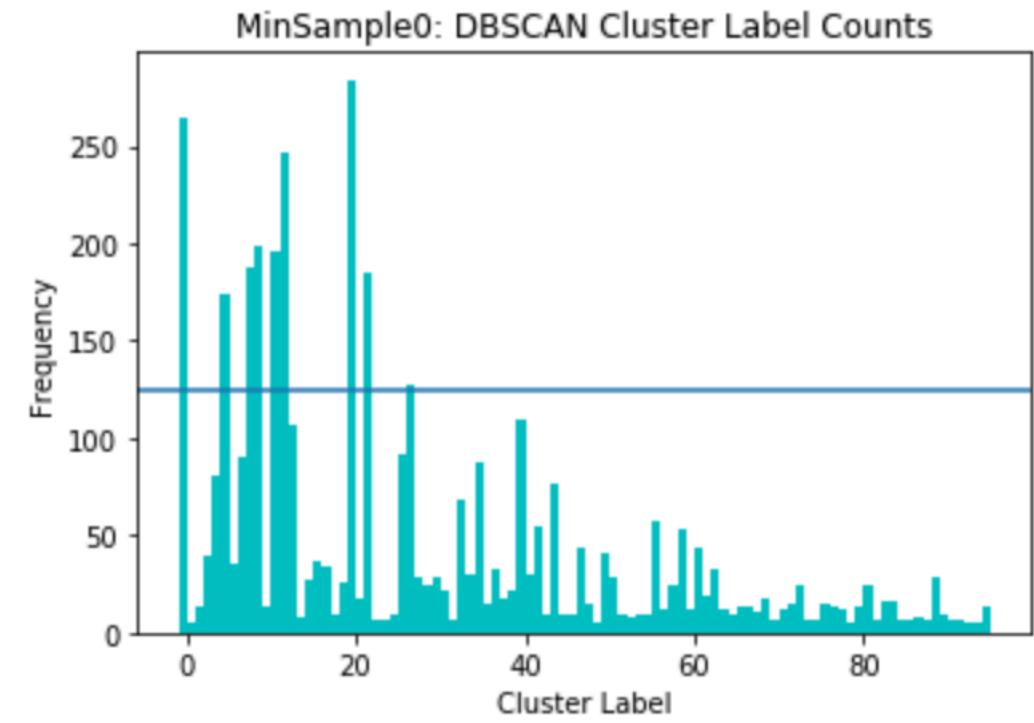
Deletions

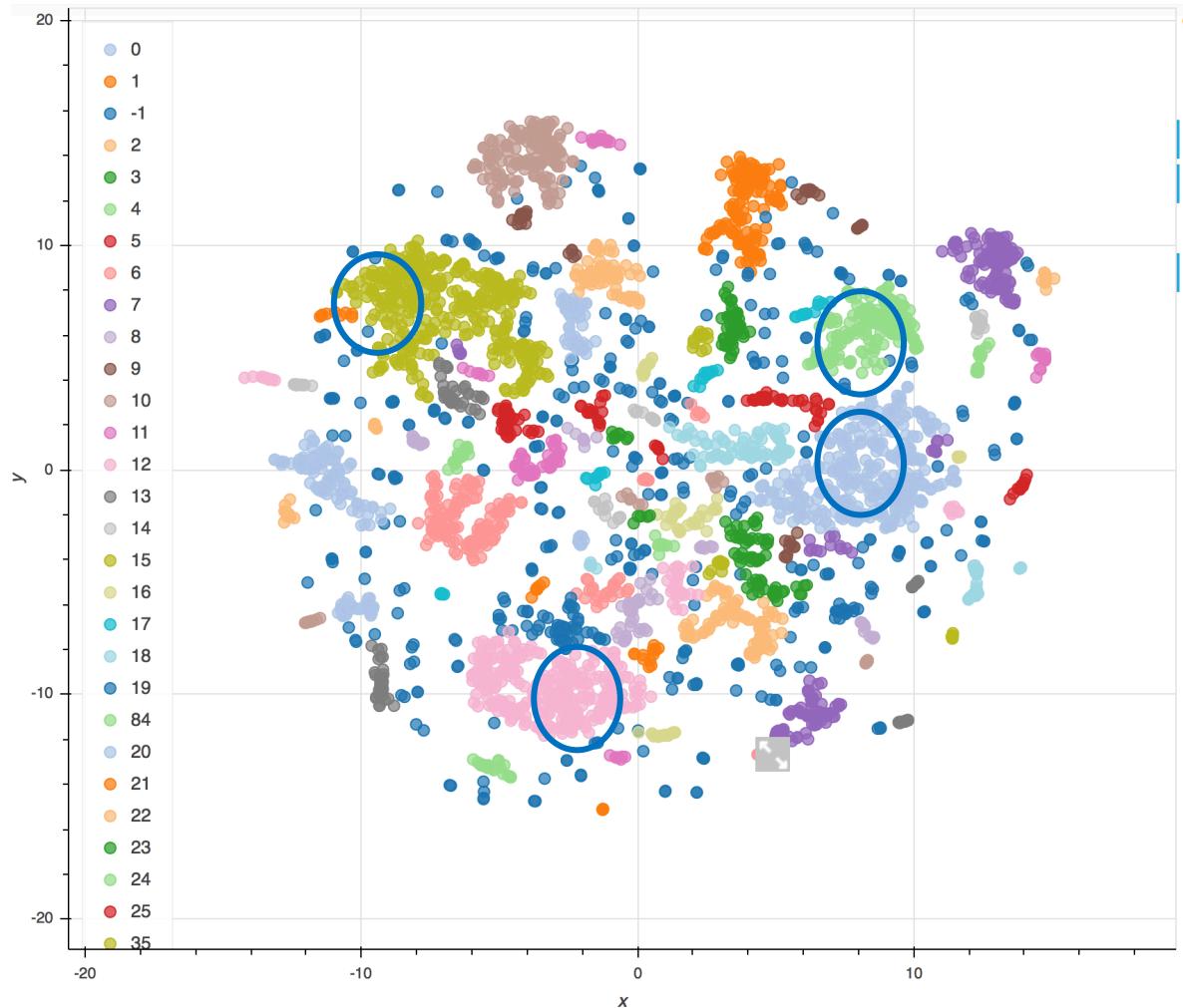


DBSCAN

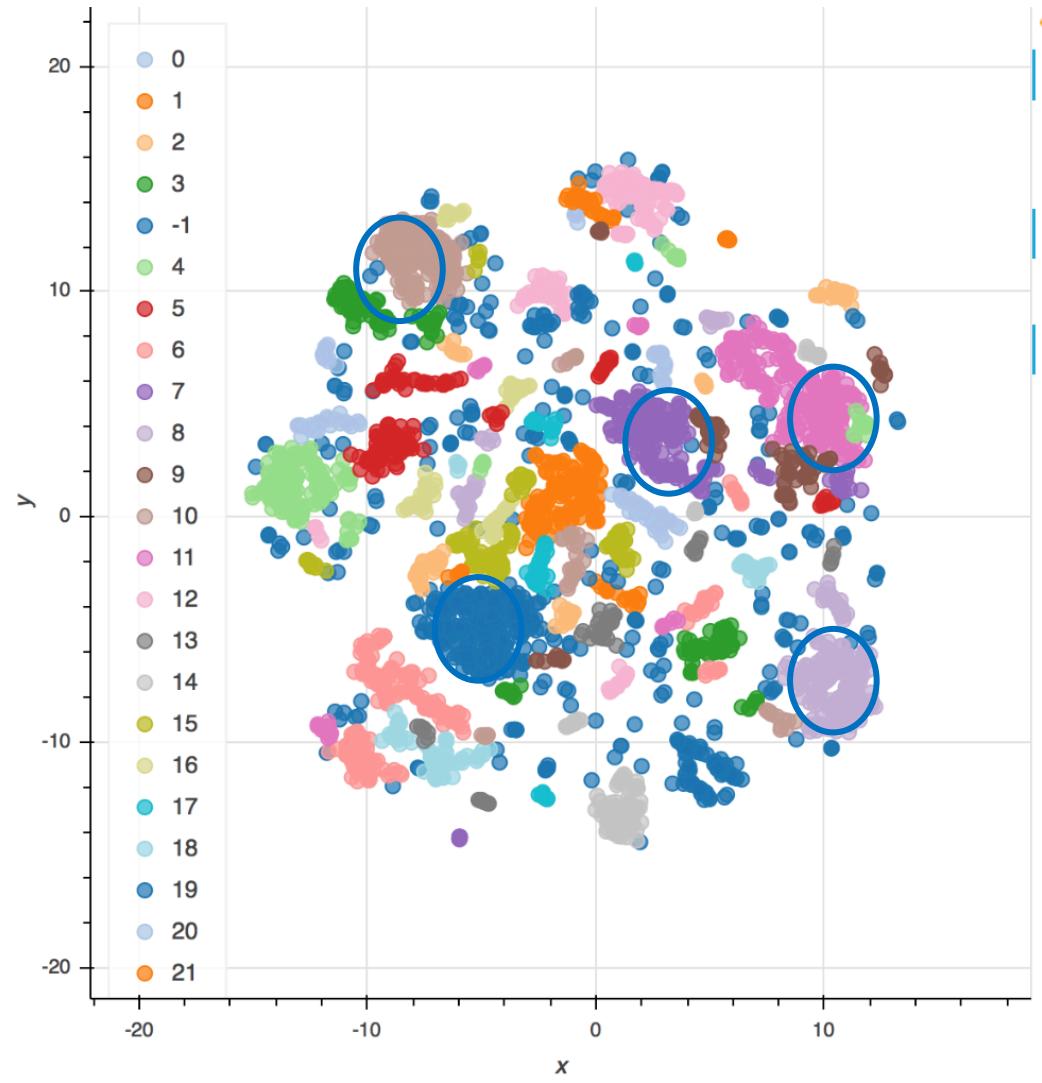


Insertions





Deletions



Insertions

App Development



Preliminary Machine Learning

CrowdSourced Labels

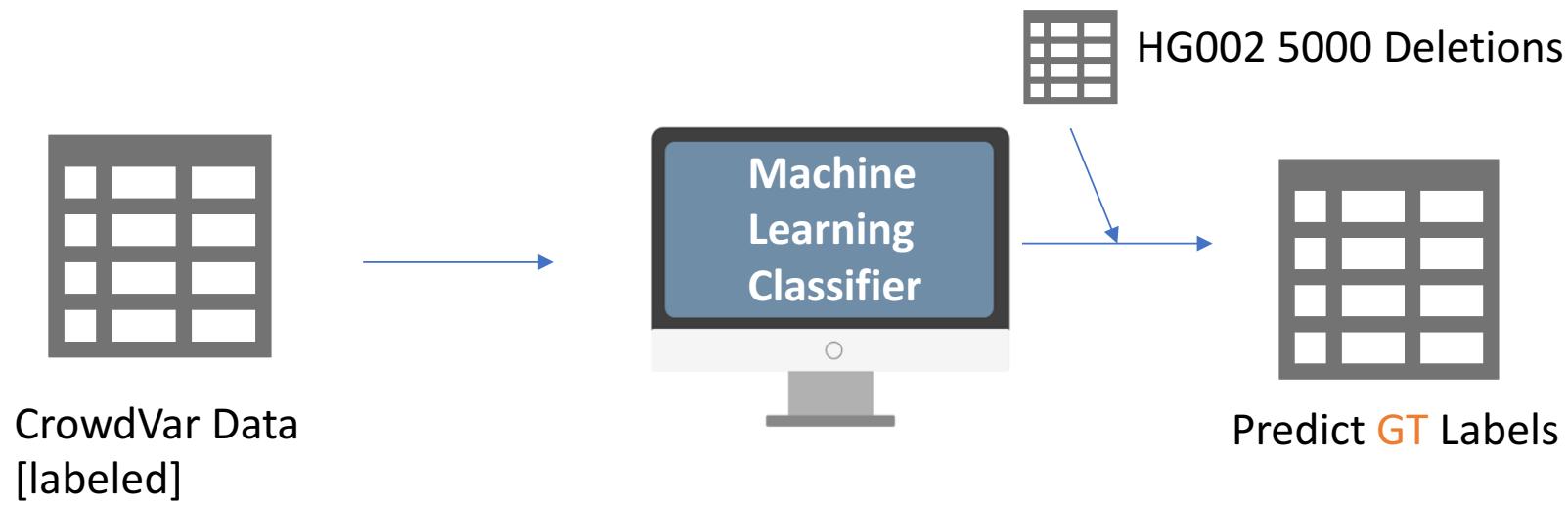


- 1700 + DEL HG002
- 8 Hom Ref
- 1097 Het Var
- 647 Hom Var

New Results

CrowdVariant: a crowdsourcing approach to classify copy number variants

Peyton Greenside, Justin Zook, Marc Salit, Ryan Poplin,
Madeleine Cule, Mark DePristo



Train Random Forest Classifier

- Train-test split [70% Train 30% Test]

Precision Score

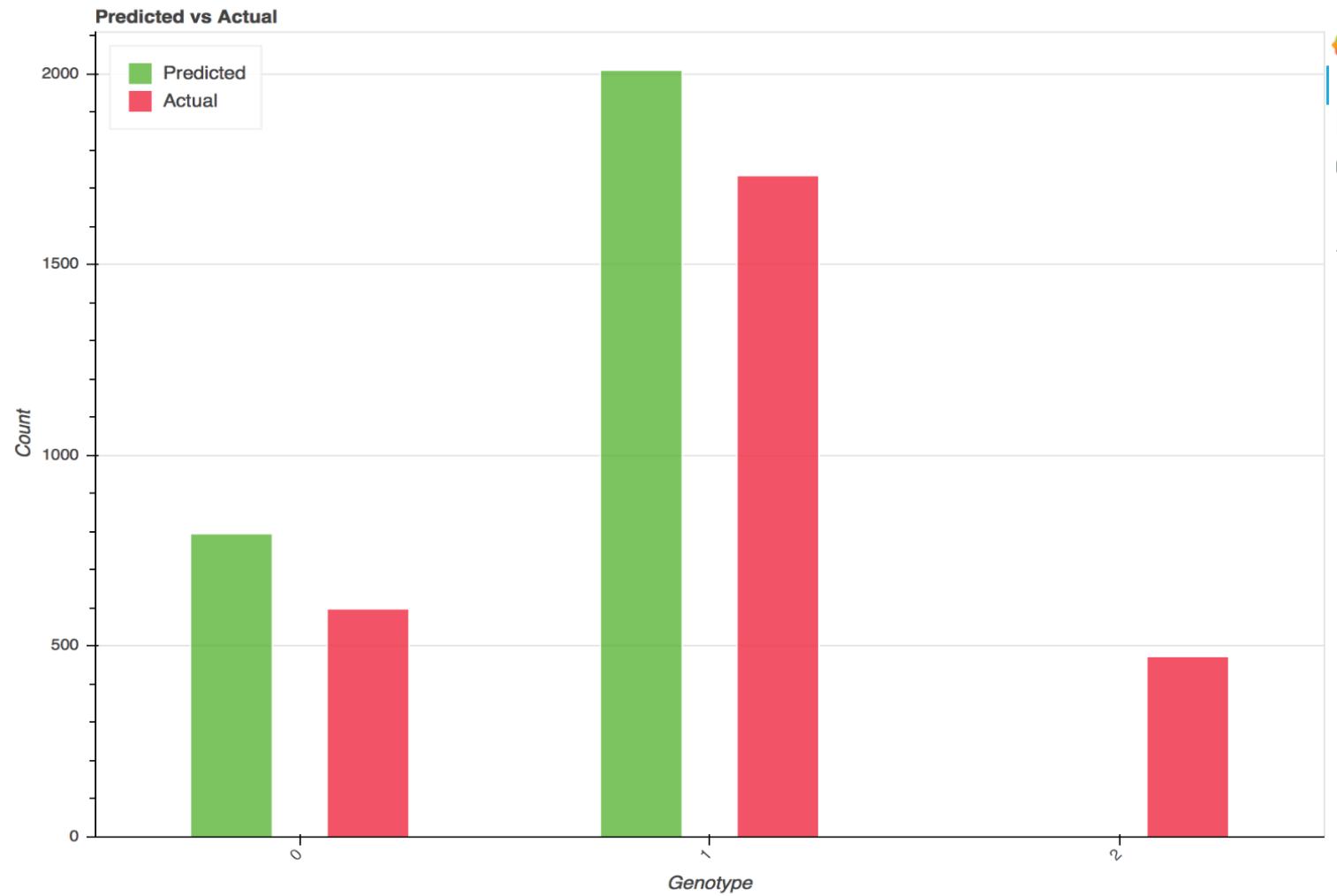
0.981619887

Predict on 5000 Randomly Selected DEL

```
Out[109]: 0.76185383244206772
```

Predict on 5000 Randomly Selected DEL

```
Out[109]: 0.76185383244206772
```



- Trained RF classifier
- Train-Test Split
- (overall prediction accuracy)
- Predict on HG002 given svviz features
- Prediction accuracy – 70%, why? Only trained on a few references
- Show examples of good and bad matches