



Memorial Sloan Kettering
Cancer Center

Variant classifier

Developing a knowledge-based
approach using IMPACT data

November 8, 2018

Papaemmanuil Lab

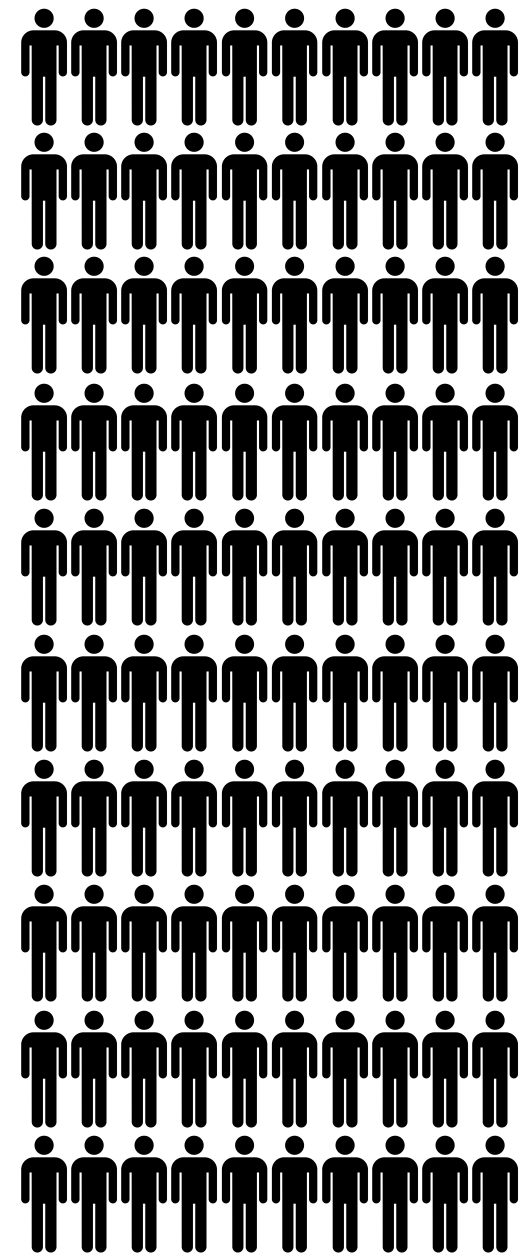
Pierre Guilmin | Elsa Bernard

In collaboration with A. Zehir, R. Ptashkin and C. Debyani



A variant classifier, why is it important?

100 patients



A variant classifier, why is it important?

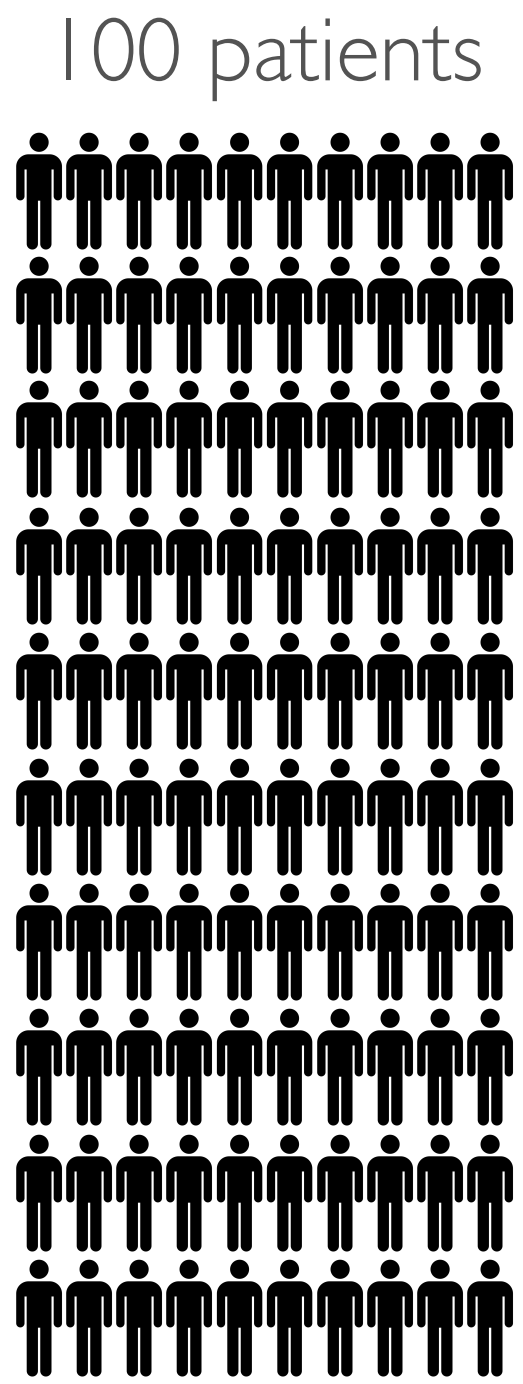
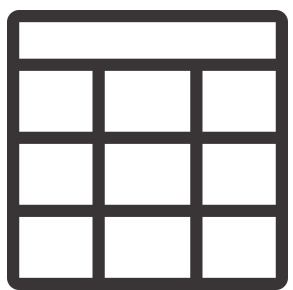
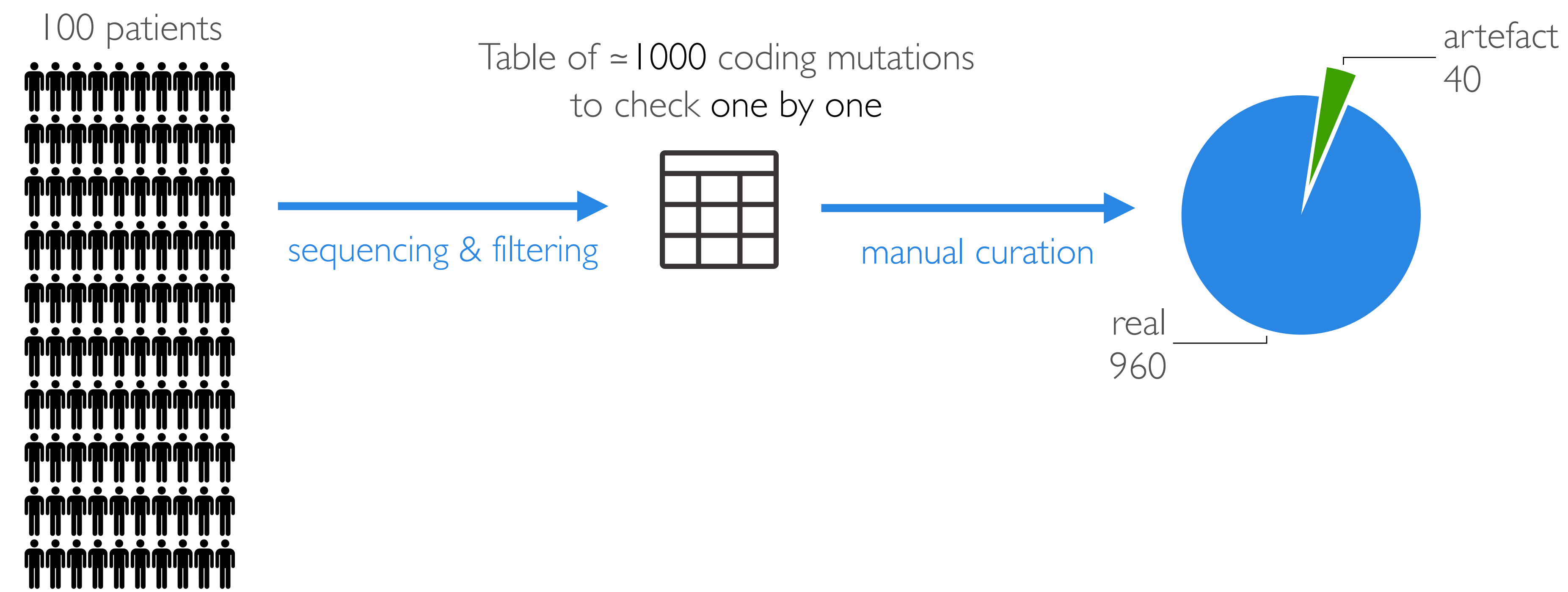


Table of ≈ 1000 coding mutations
to check one by one



A variant classifier, why is it important?



A variant classifier, why is it important?

100 patients

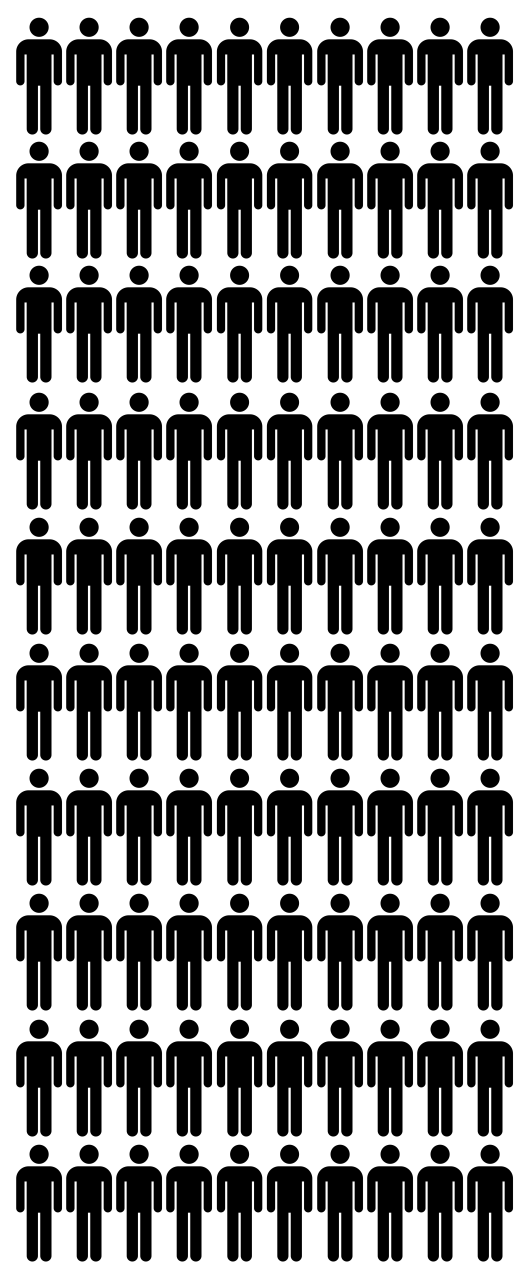
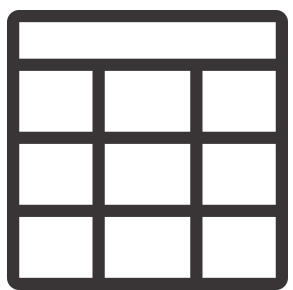
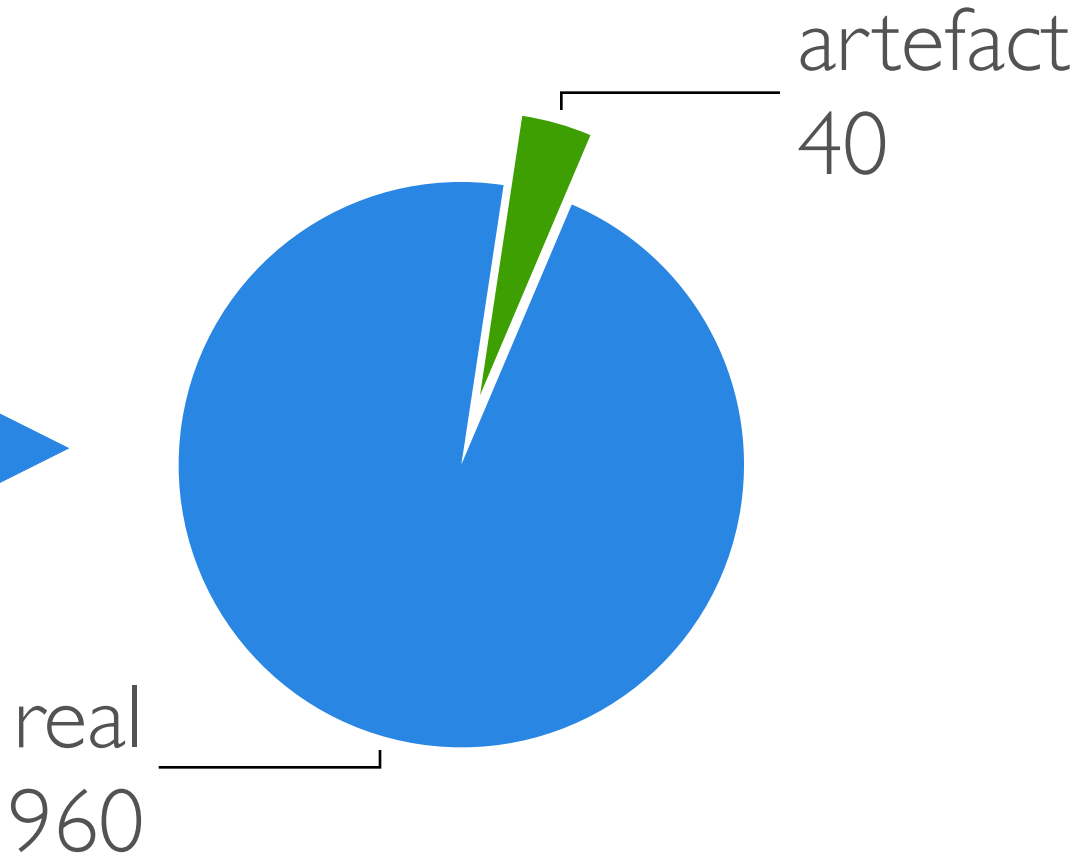


Table of ≈ 1000 coding mutations
to check one by one

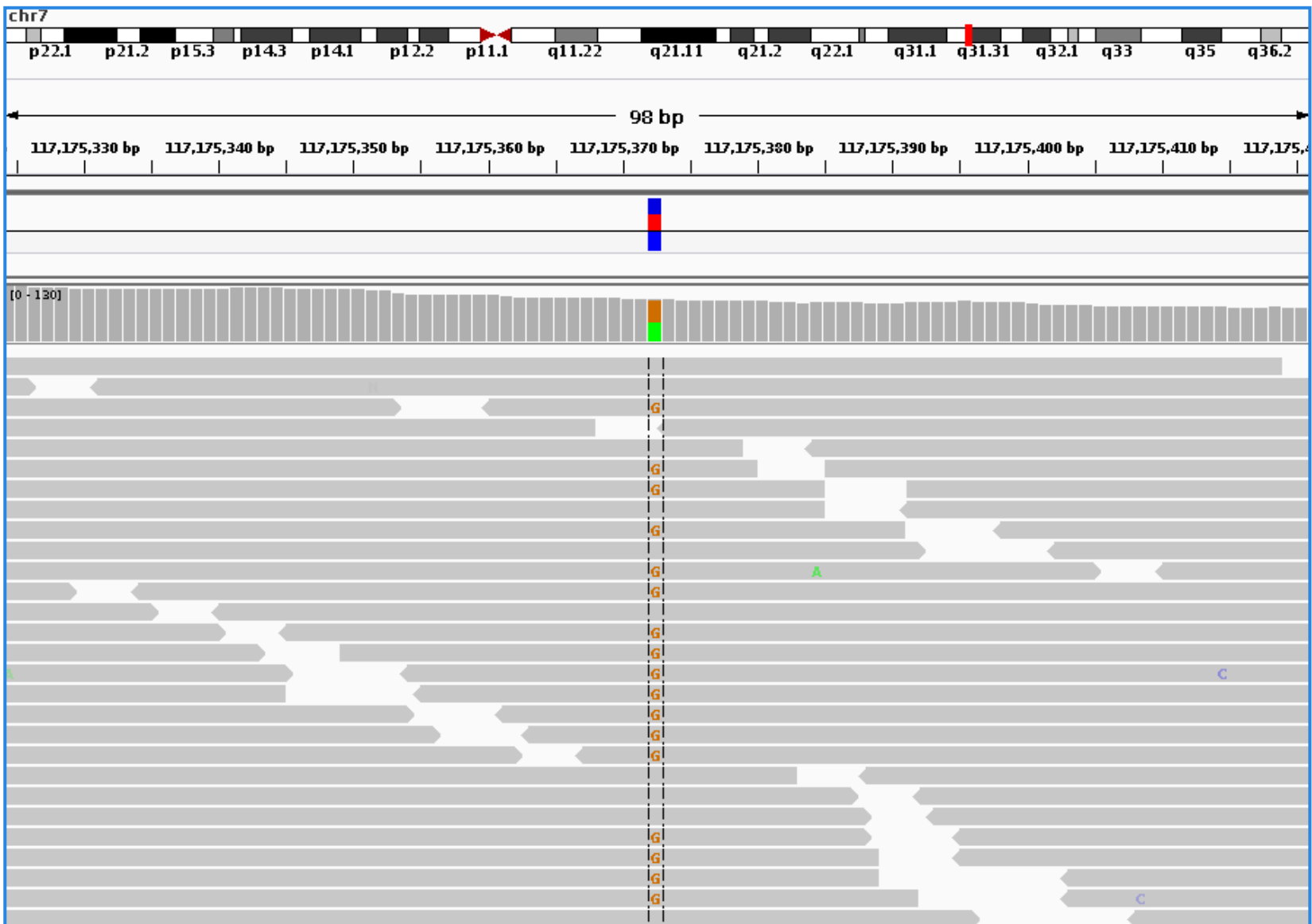
sequencing & filtering



manual curation

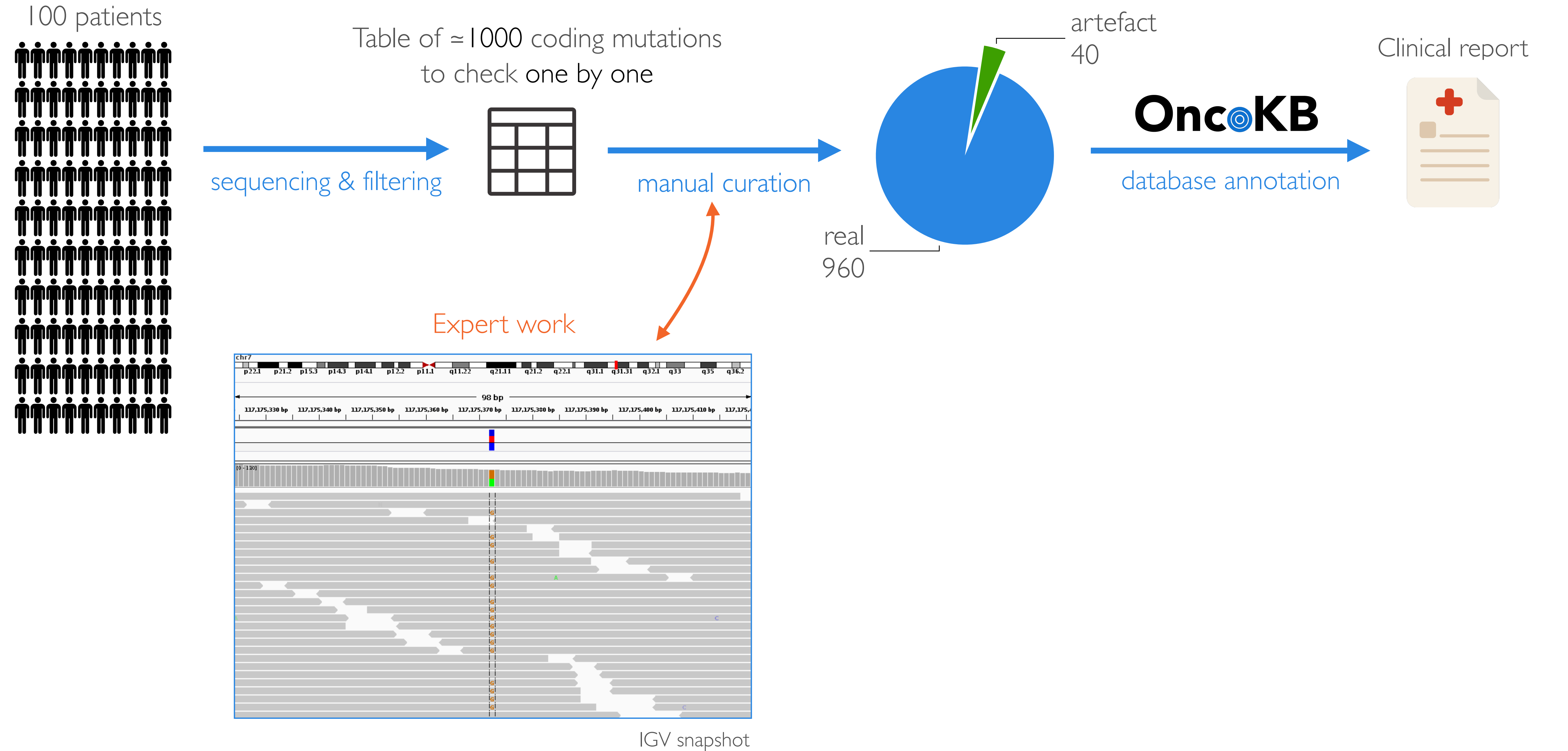


Expert work

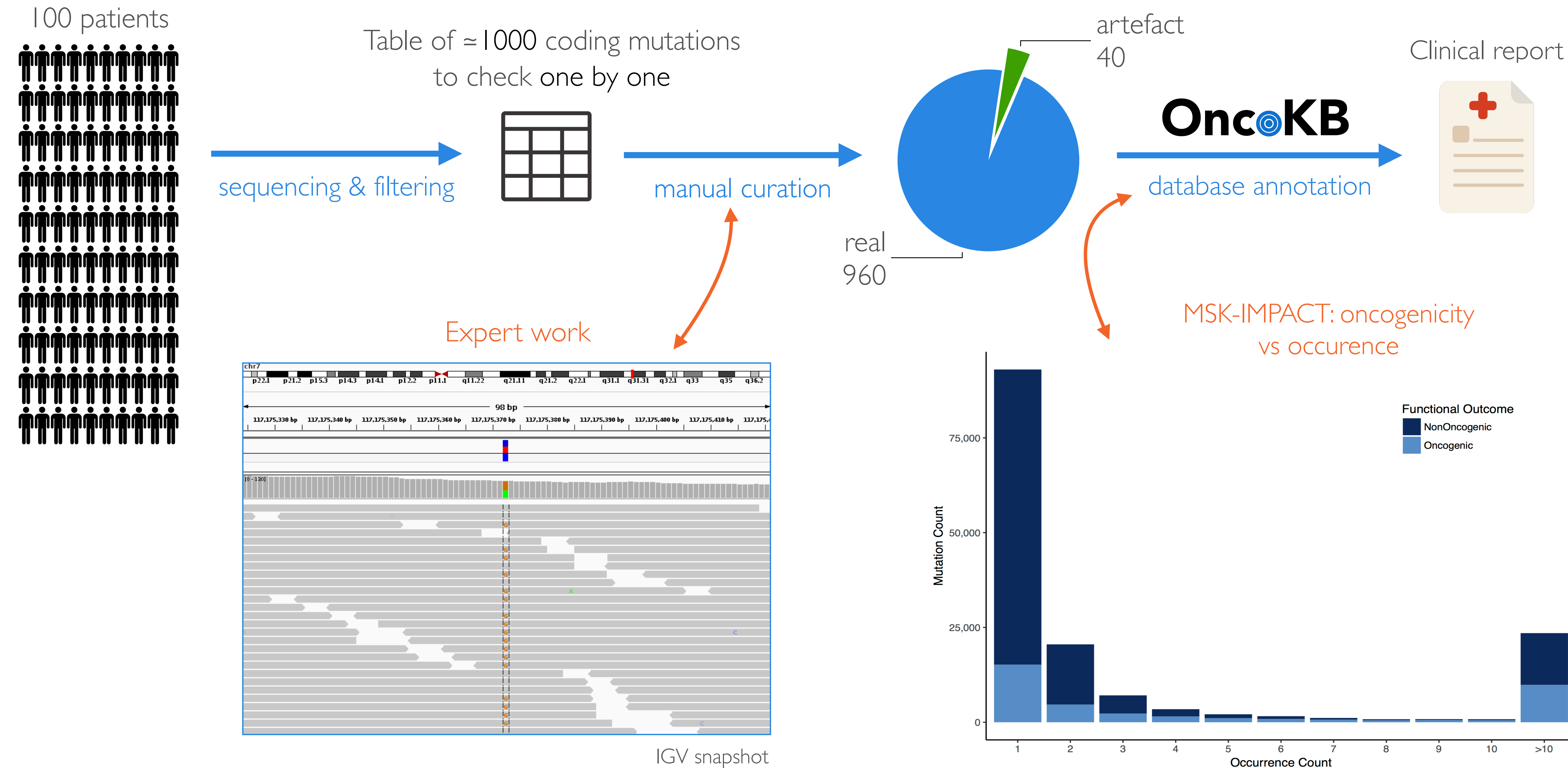


IGV snapshot

A variant classifier, why is it important?



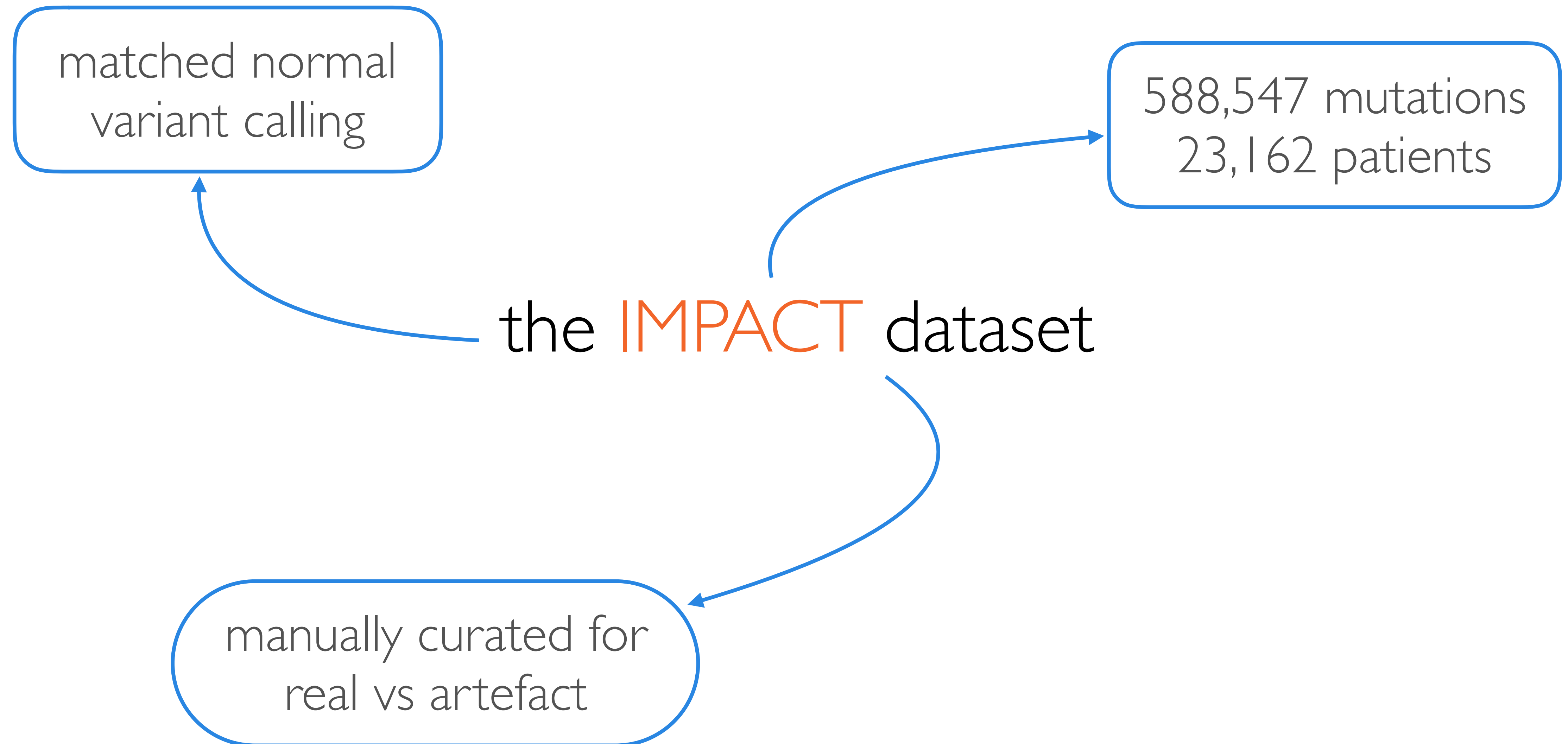
A variant classifier, why is it important?



Create a tool that classifies variant automatically

- real vs artefact OR driver vs passenger
- using Supervised Machine Learning Classification
- all cancer, all mutation type
- on the **IMPACT** dataset

IMPACT, the dataset



Two steps classification

IMPACT
dataset

coding + splicing
(194,211 mutations
= 36%)

Two steps classification

IMPACT
dataset

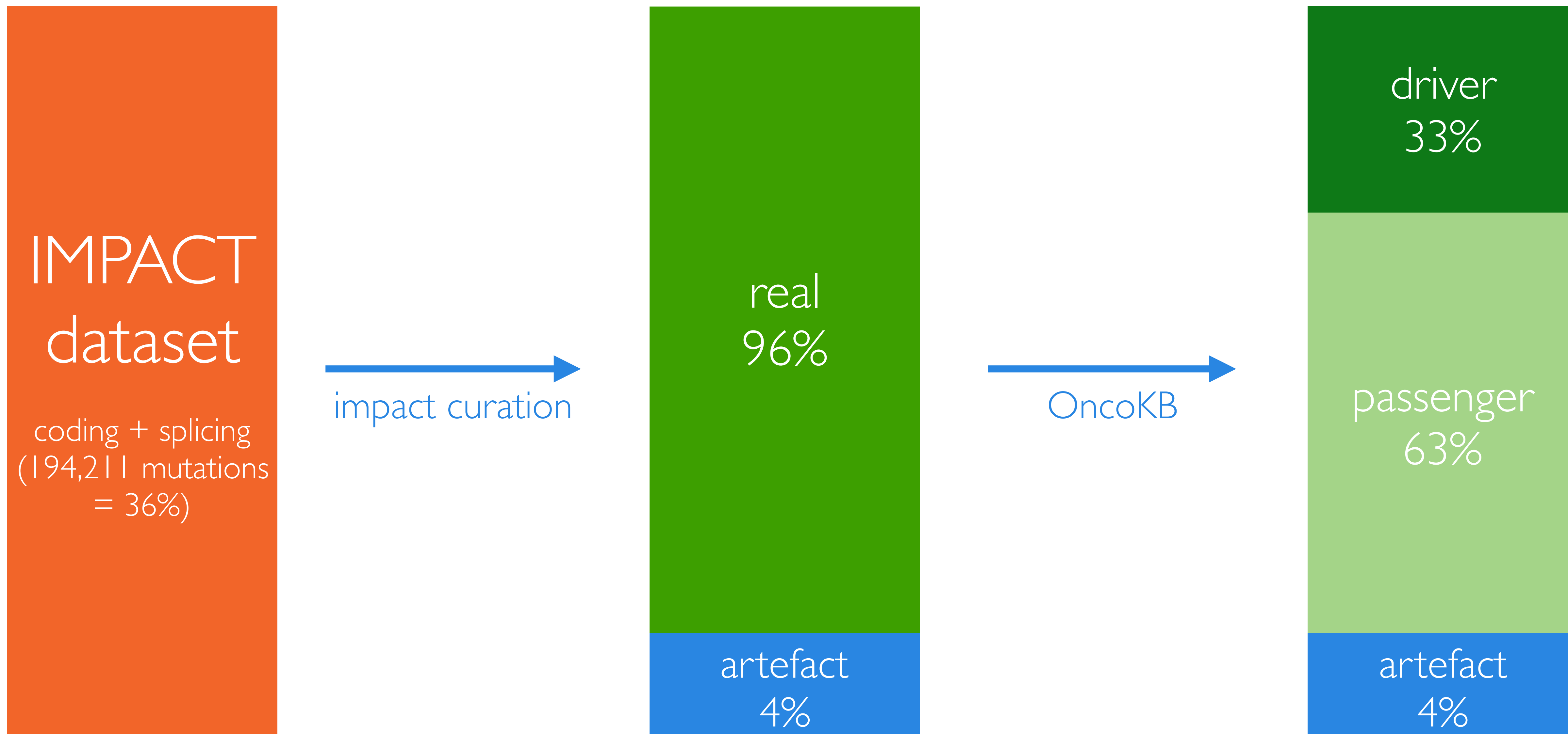
coding + splicing
(194,211 mutations
= 36%)

→
impact curation

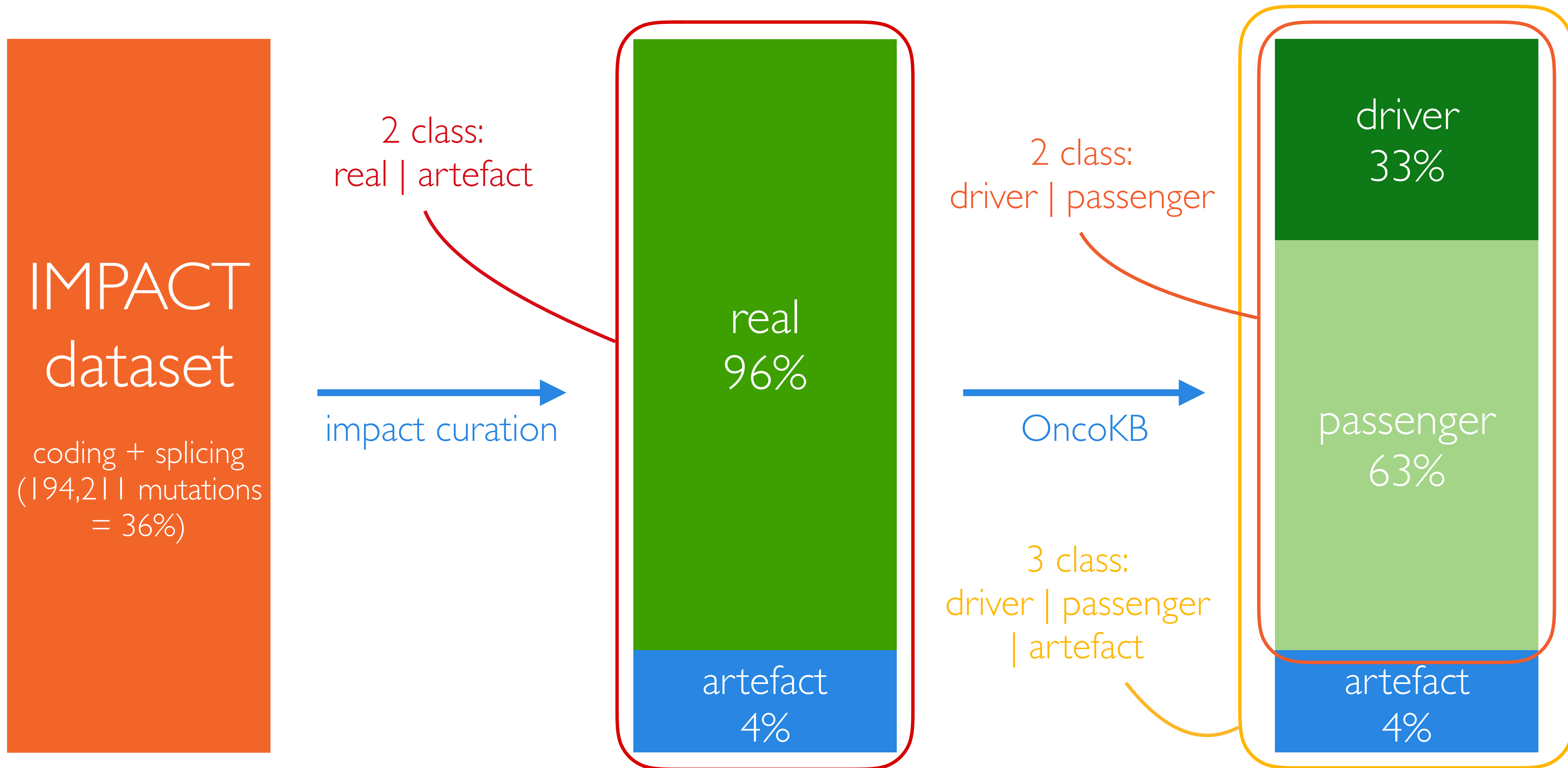
real
96%

artefact
4%

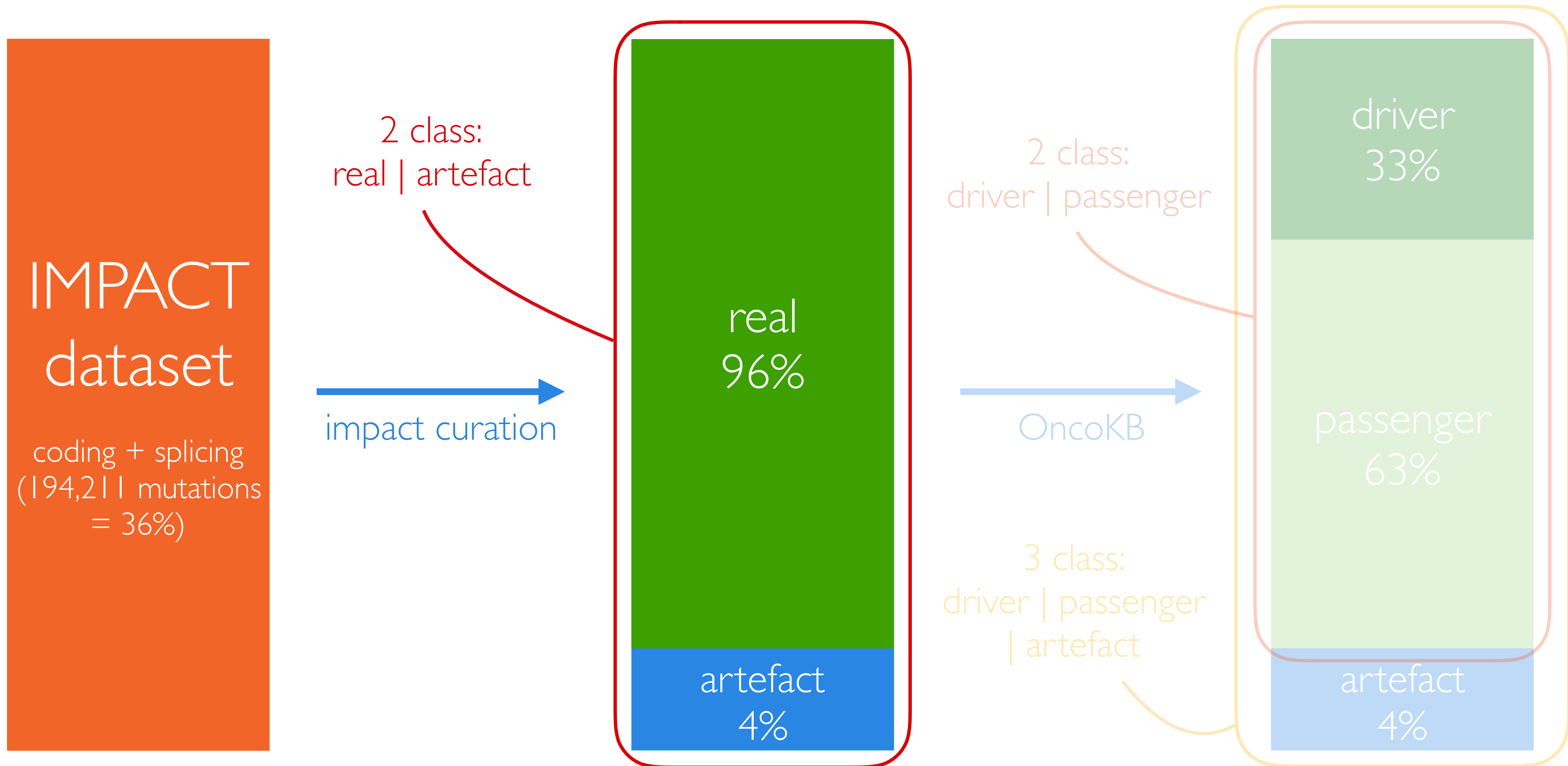
Two steps classification



Two steps classification



Two steps classification



The features used in our model



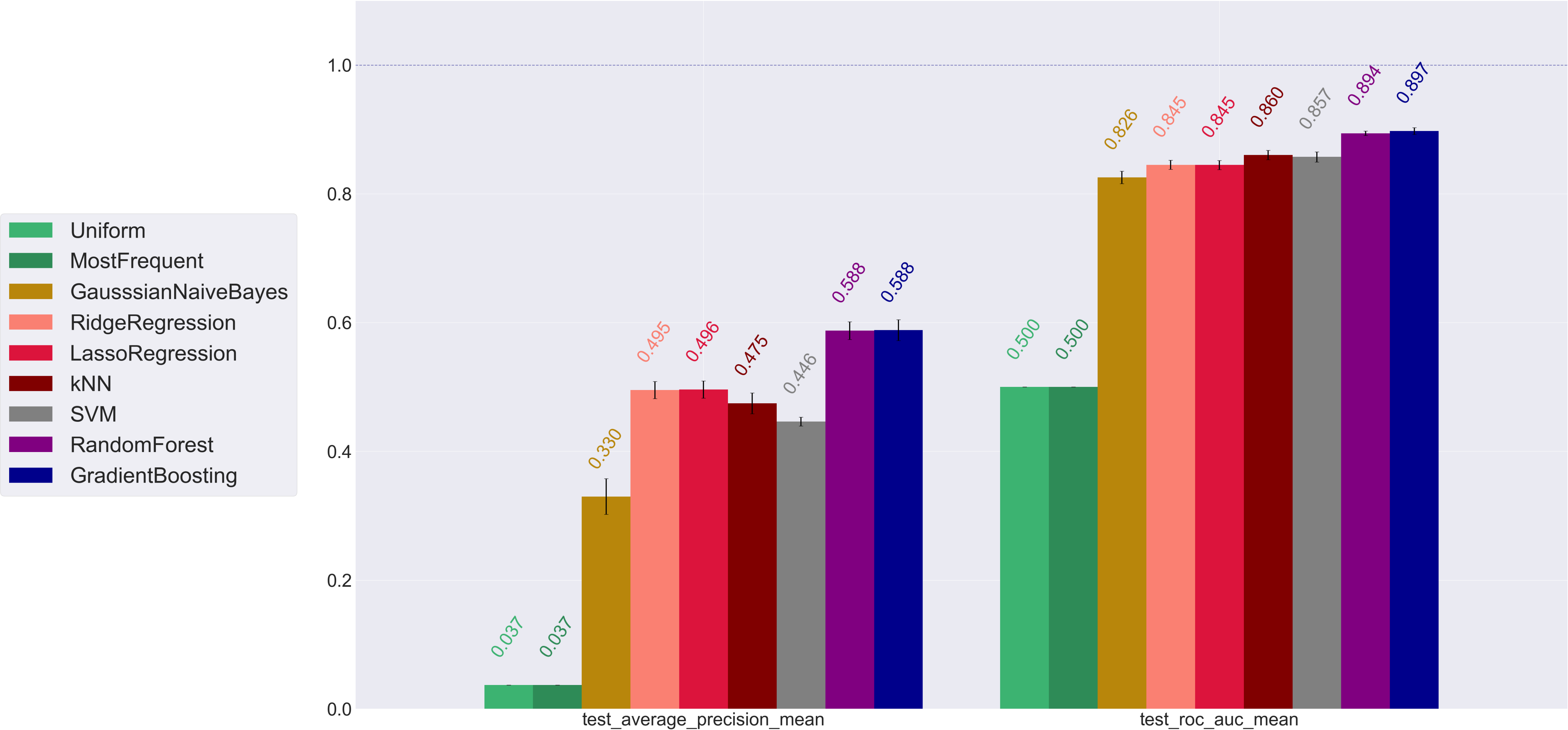
- Sequencing features ($n = 11$)
Tumor VAF, tumor depth
- Genomic coordinates ($n = 3$)
Chromosome, Hugo Symbol
- Control population ($n = 12$)
GnomAD allele frequency
- Cancer population ($n = 4$)
COSMIC count, OncoKB
- Normal control ($n = 1$)
Frequency in normals
- Mutation consequence ($n = 6$)
Protein effect, SIFT & PolyPhen class

Algorithm comparison

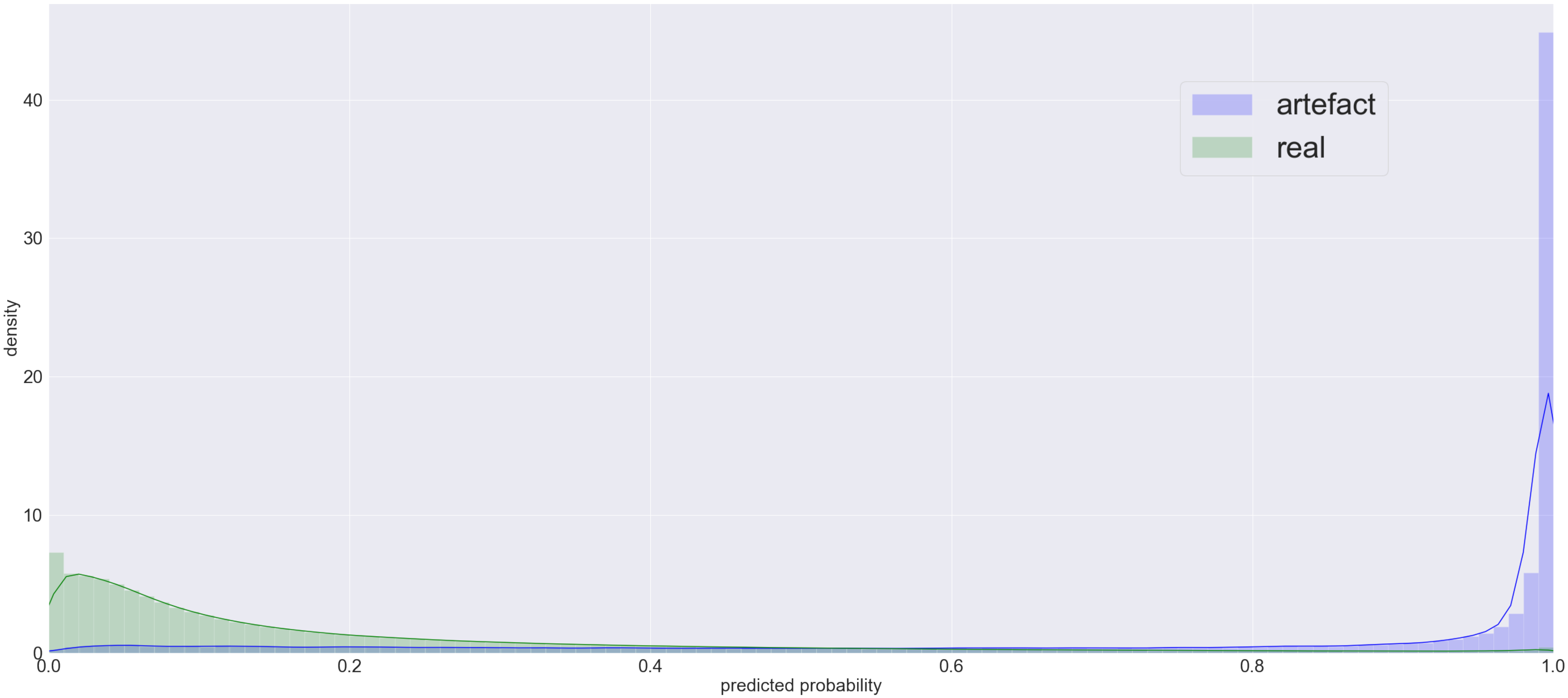


- Uniform
- MostFrequent
- GaussssianNaiveBayes
- RidgeRegression
- LassoRegression
- kNN
- SVM
- RandomForest
- GradientBoosting

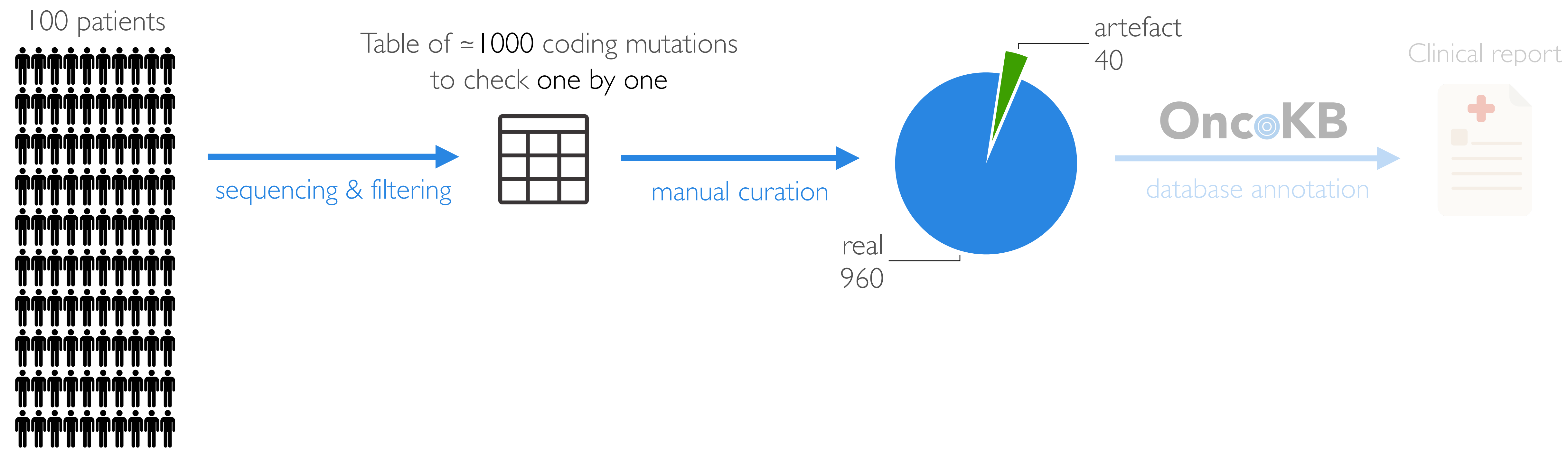
Algorithm comparison



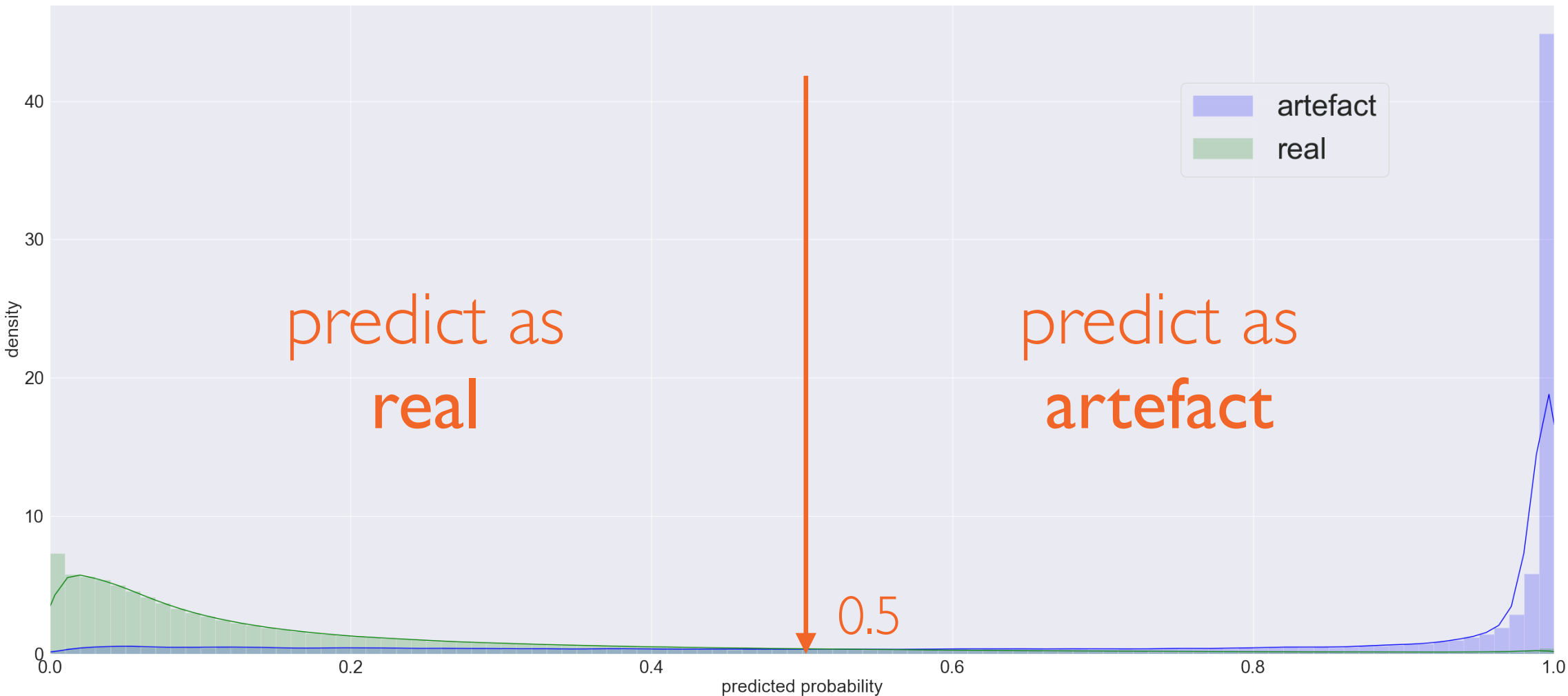
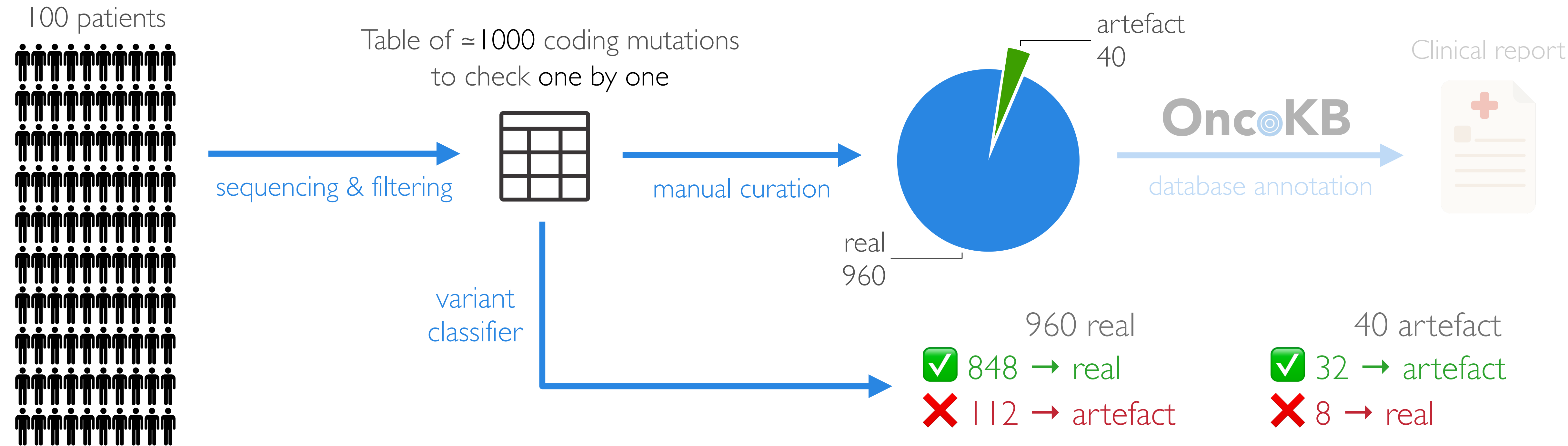
Best algorithm probability output



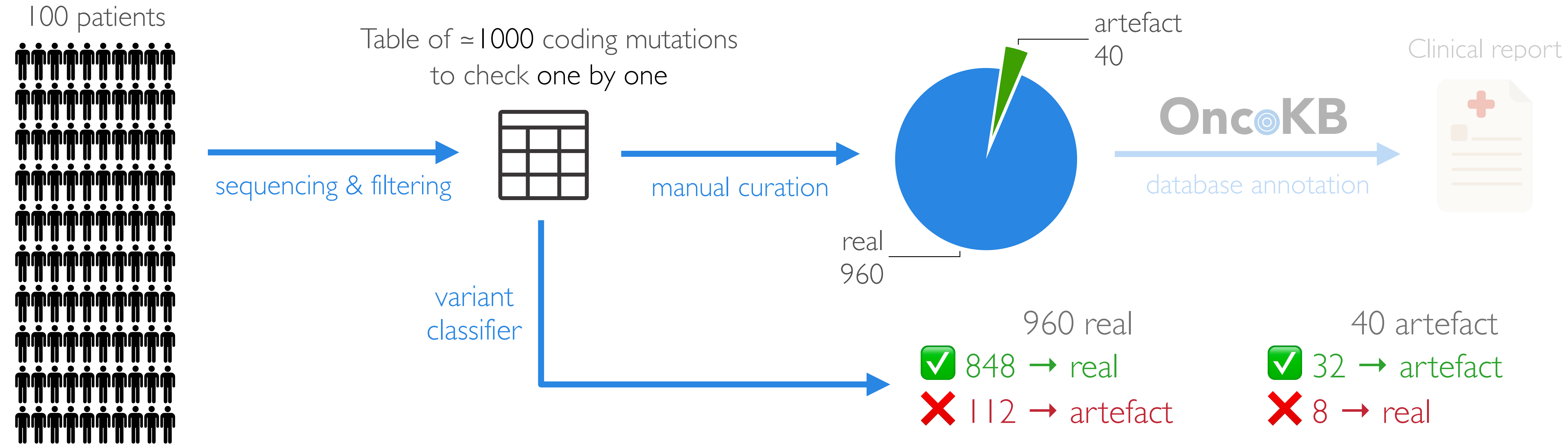
The variant classifier performances



The variant classifier performances



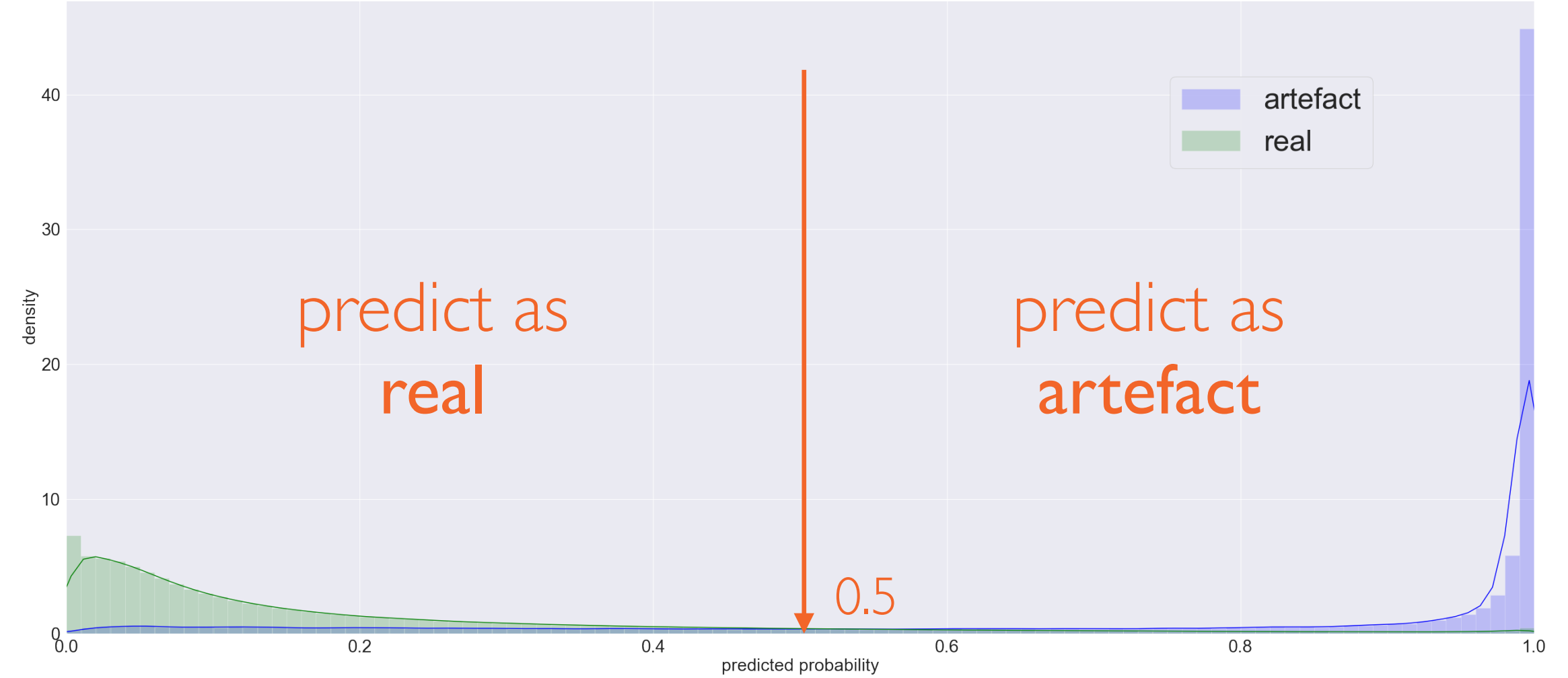
The variant classifier performances



8x less work

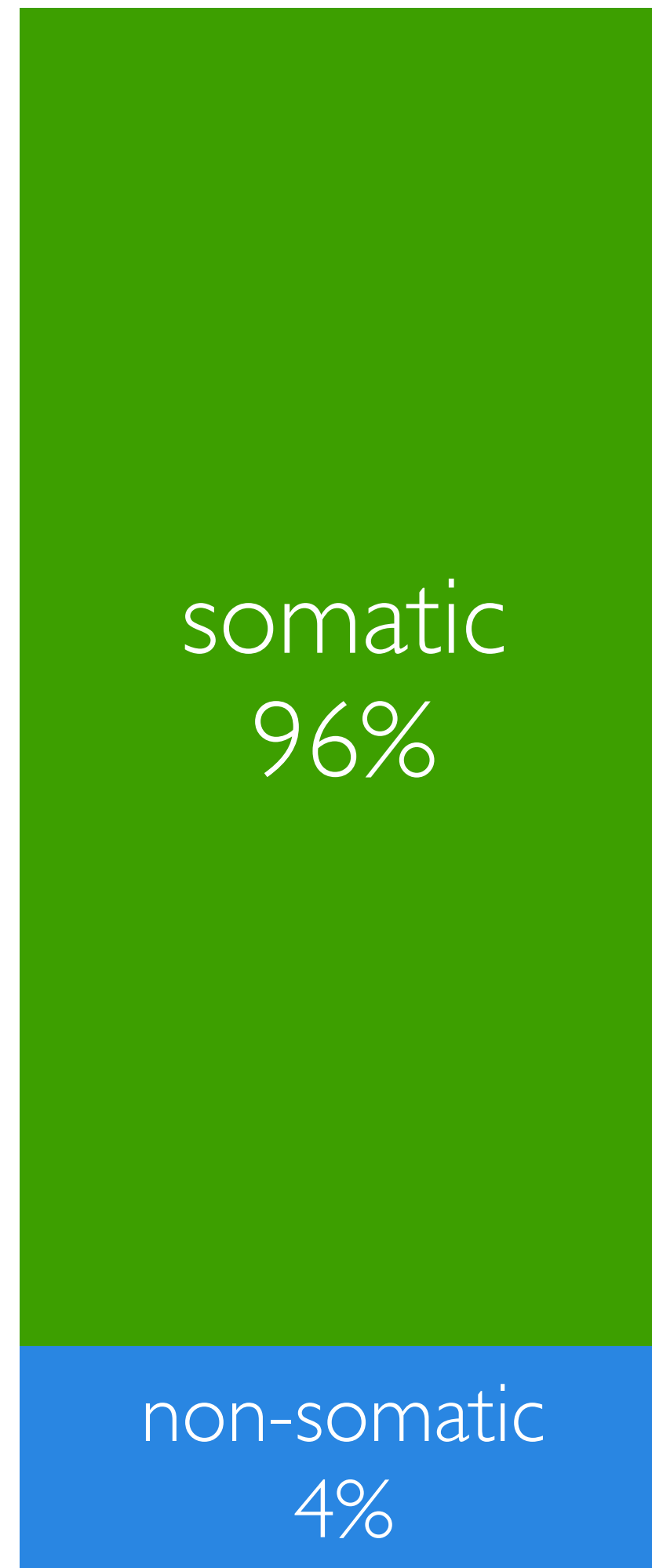
144/1000 mutations to check one by one instead of 1000/1000

✗ 8/40 artefacts considered as real

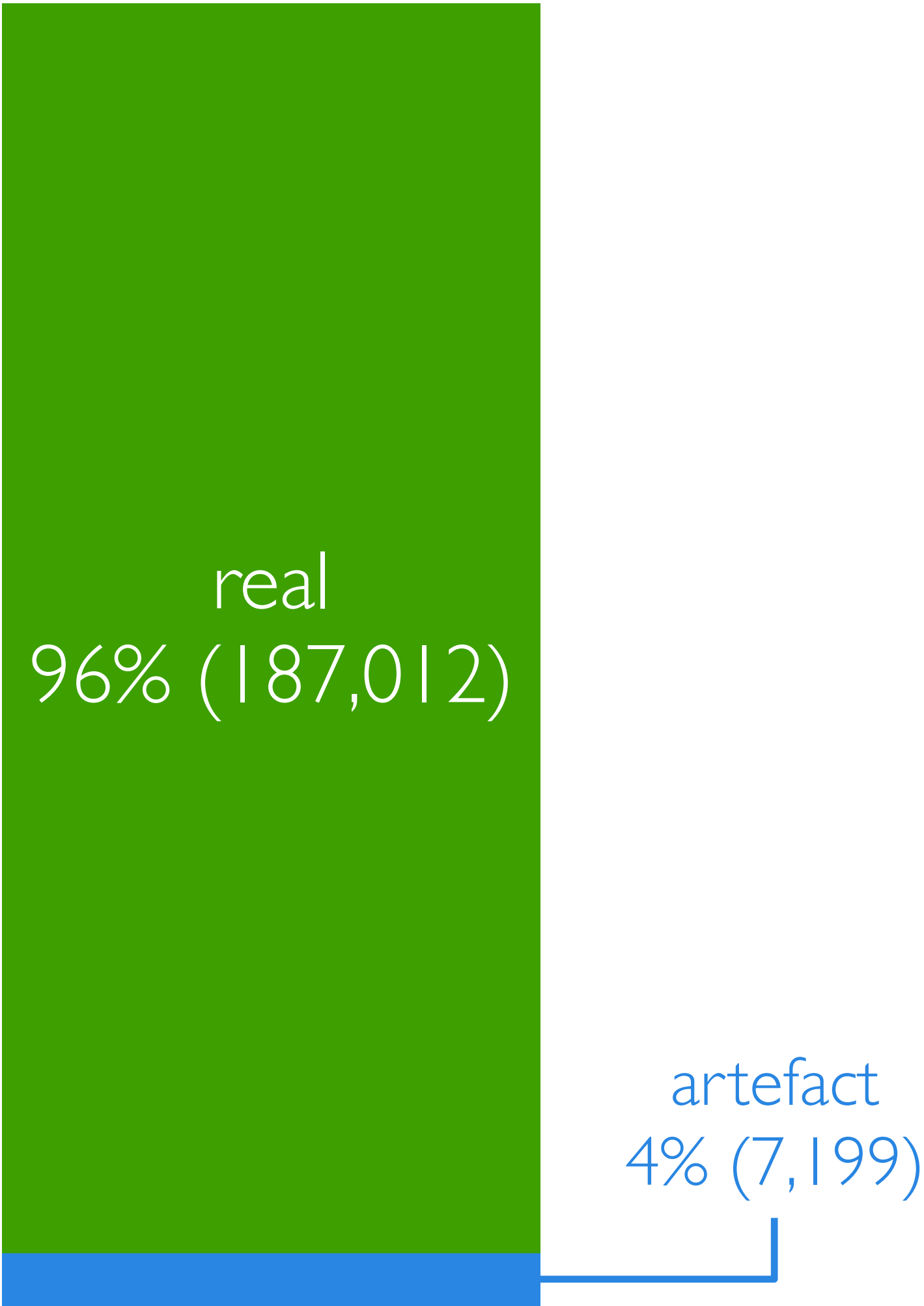


Main challenges

Imbalanced dataset



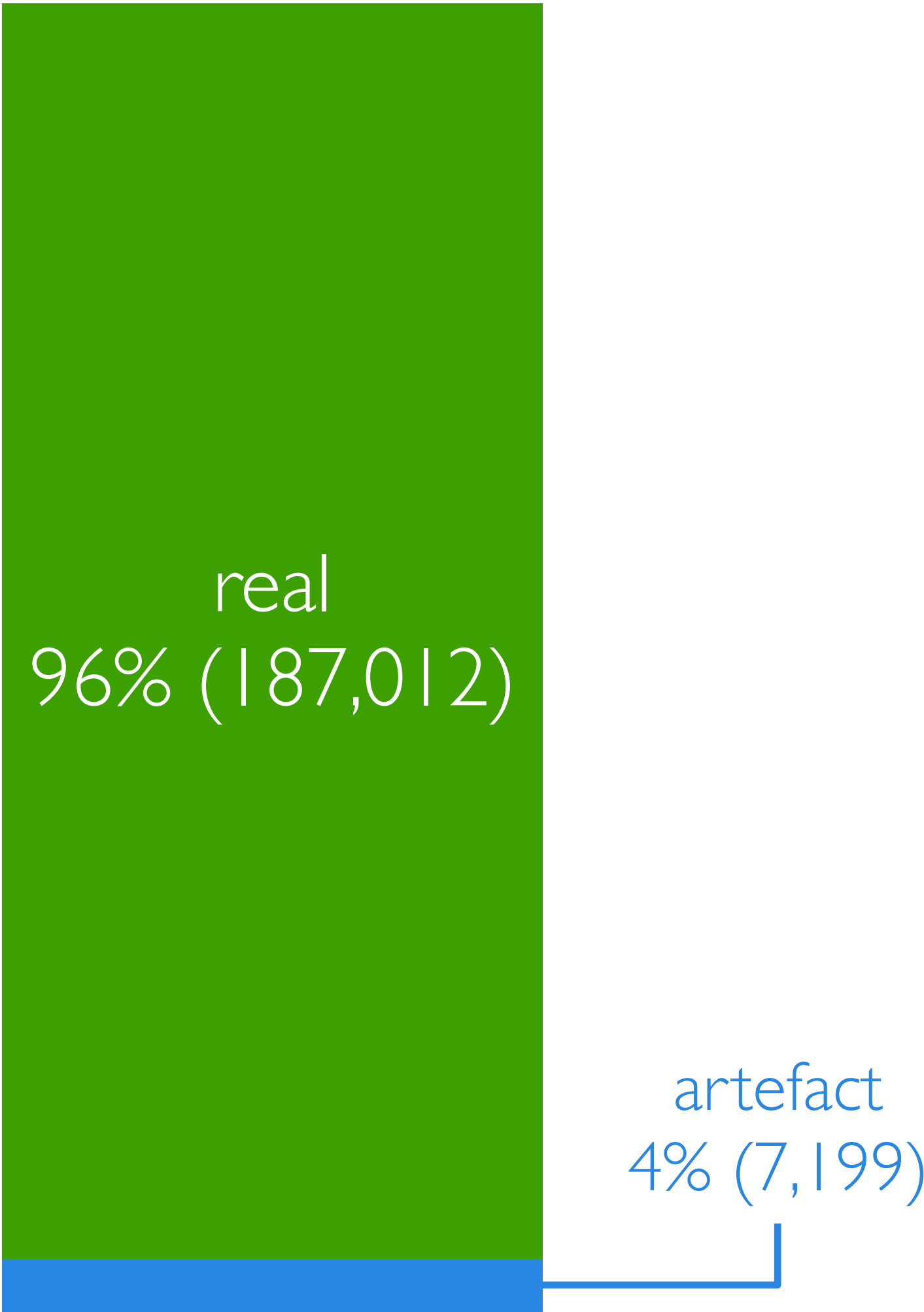
Imbalanced dataset



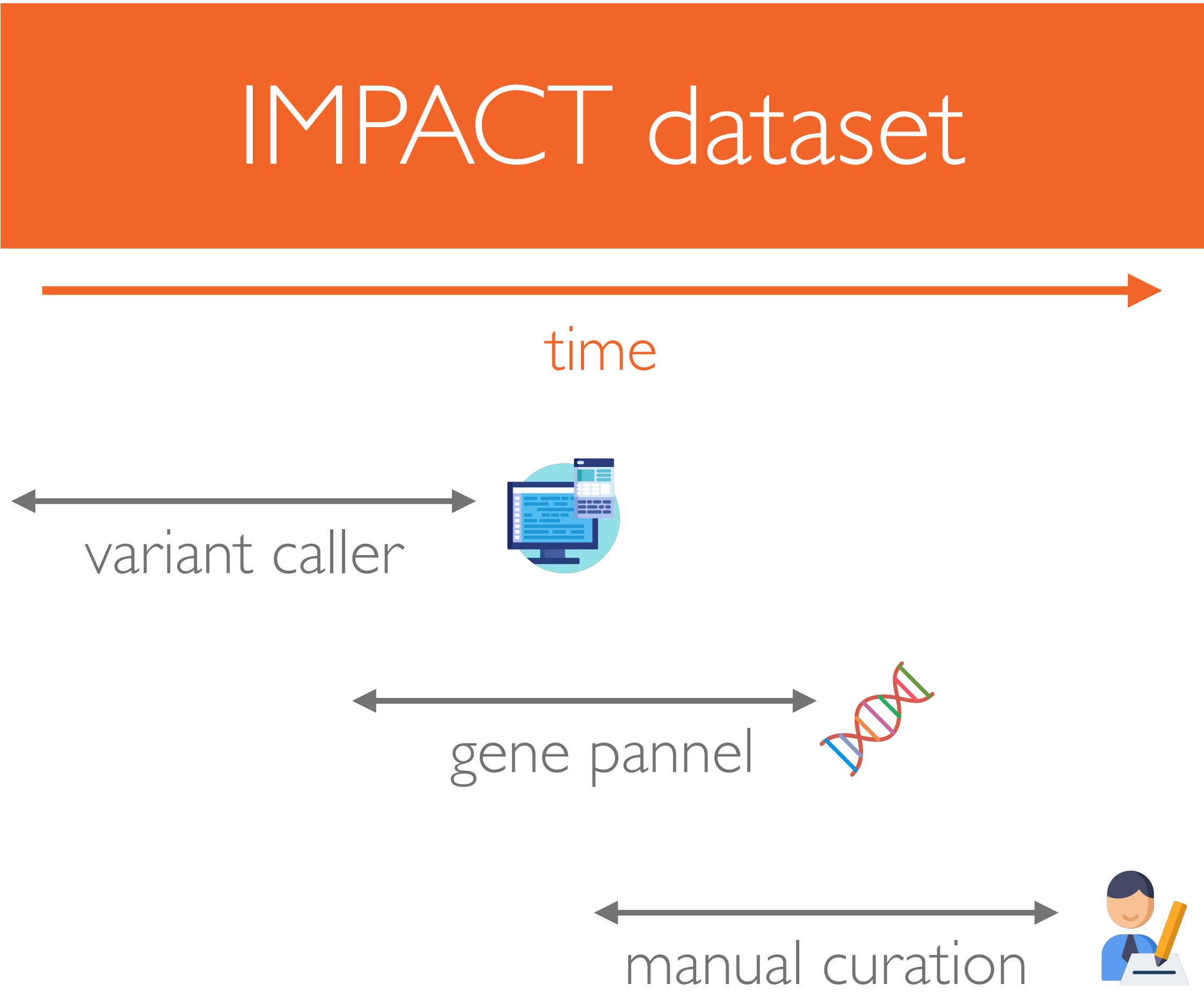
Main challenges




Imbalanced dataset



Evolution over time



Compare with new paper



TECHNICAL REPORT
<https://doi.org/10.1038/s41588-018-0257-y>

A deep learning approach to automate refinement of somatic variant calling from cancer sequencing data

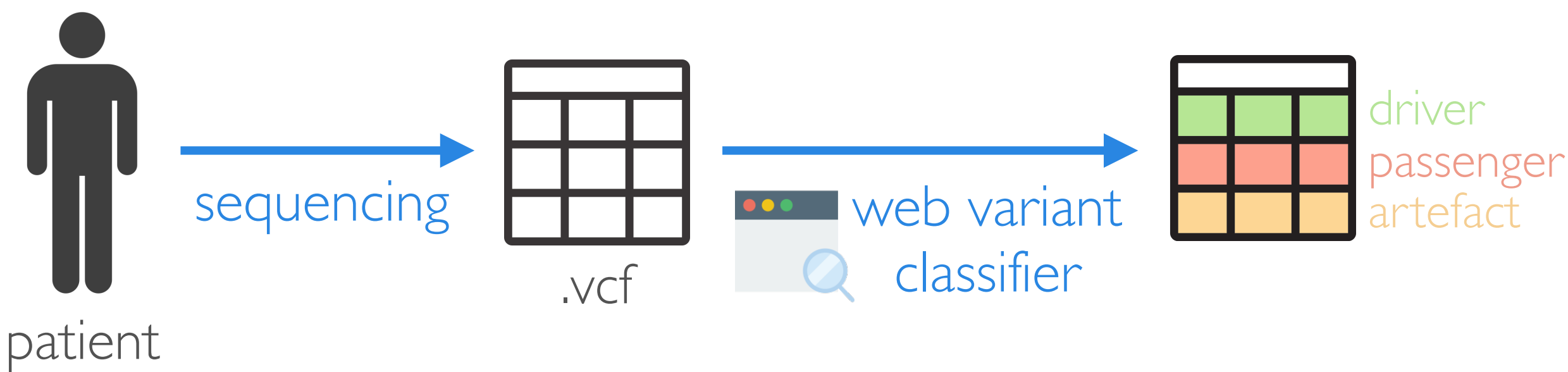
Benjamin J. Ainscough^{1,2,12}, Erica K. Barnell^{1,12}, Peter Ronning¹, Katie M. Campbell¹, Alex H. Wagner¹, Todd A. Fehniger^{2,3}, Gavin P. Dunn⁴, Ravindra Uppaluri⁵, Ramaswamy Govindan^{2,3}, Thomas E. Rohan⁶, Malachi Griffith^{1,2,3,7}, Elaine R. Mardis^{8,9}, S. Joshua Swamidass^{10,11*} and Obi L. Griffith^{1,2,3,7*}

Get detailed calling features
Variant caller flags, read mapping quality, ...

Explore new methods

- Deep learning
- Improved under-sampling strategy
- ...

Create a 2-steps web-based classifier



Optimisation of clinical heme panel

