Memorial Sloan Kettering
Cancer Center

# IMPACT annotator

September 26, 2018

Pierre Guilmin

Elsa Bernard

# Create a tool that classifies variant automatically

- somatic vs non-somatic   OR   driver vs passenger

- using Supervised Machine Learning Classification

- on the IMPACT dataset

matched normal
variant calling

588,547 mutations
23,162 patients

the IMPACT dataset

Manually curated dataset
for somatic vs non-somatic:
- OK
- UNLIKELY
- UNKNOWN

IMPACT dataset

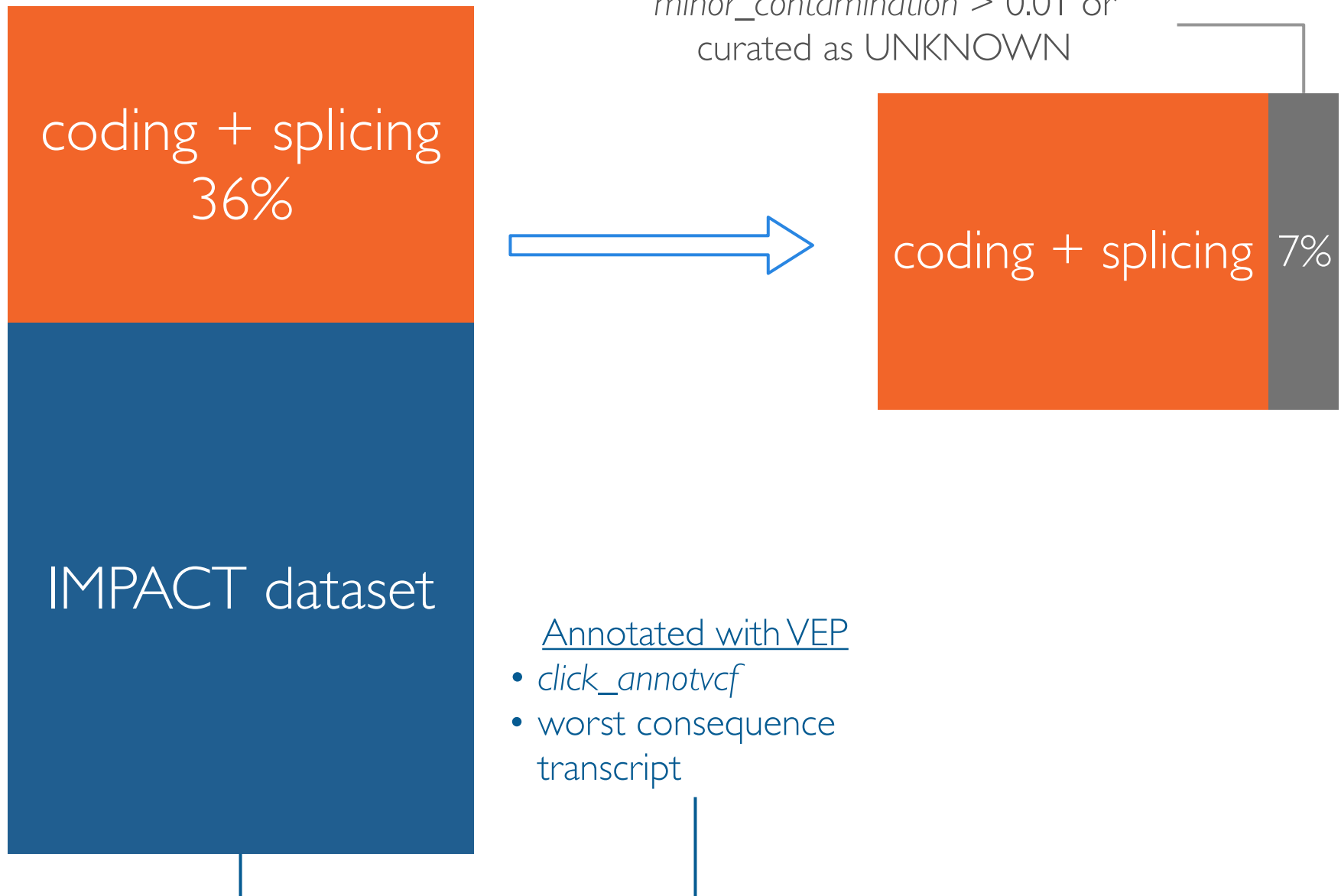Annotated with VEP
- *click_annotvcf*
- worst consequence transcript

coding + splicing 36%

IMPACT dataset

Annotated with VEP
- *click_annotvcf*
- worst consequence transcript

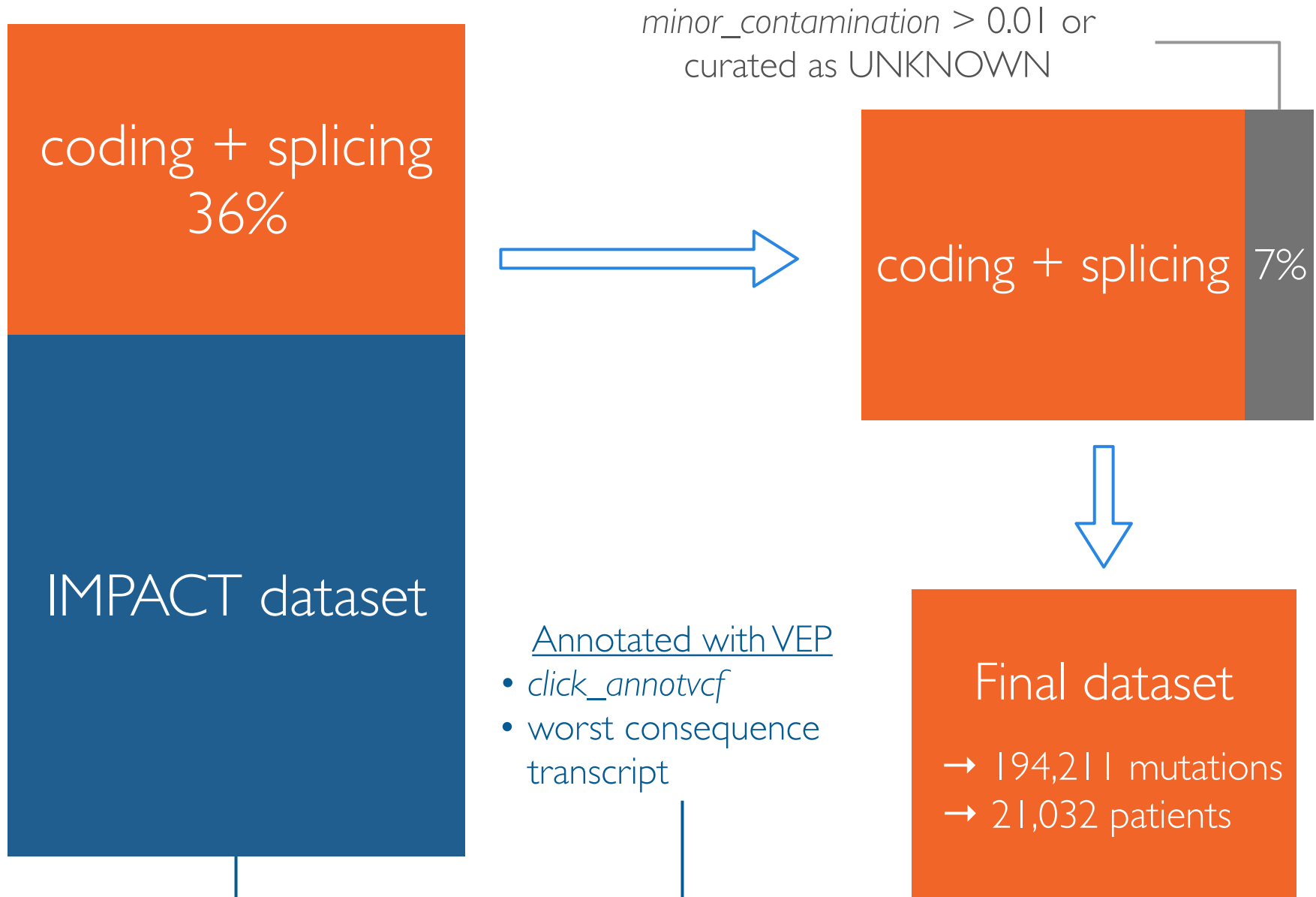# The IMPACT sub-dataset used

coding + splicing
36%

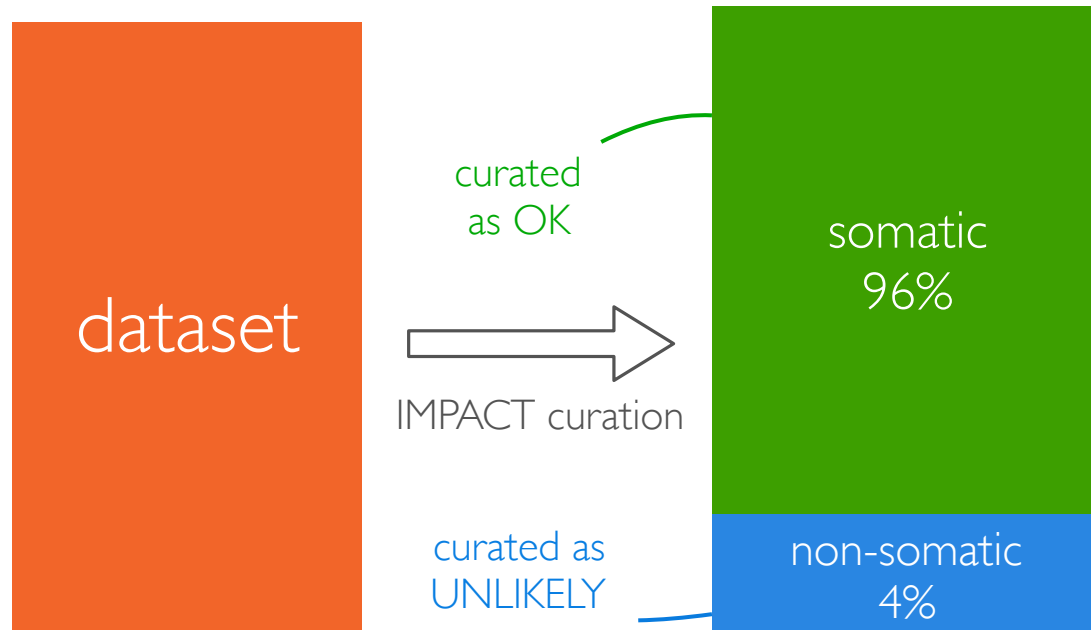IMPACT dataset

*minor_contamination* > 0.01 or
curated as UNKNOWN

coding + splicing  7%

<u>Annotated with VEP</u>
- *click_annotvcf*
- worst consequence transcript

# The IMPACT sub-dataset used

coding + splicing
36%

IMPACT dataset

*minor_contamination* > 0.01 or
curated as UNKNOWN

coding + splicing  7%

## Annotated with VEP
- *click_annotvcf*
- worst consequence
  transcript

## Final dataset

→ 194,211 mutations
→ 21,032 patients

dataset

# Different classes



dataset

IMPACT curation

curated as OK

curated as UNLIKELY

somatic 96%

non-somatic 4%
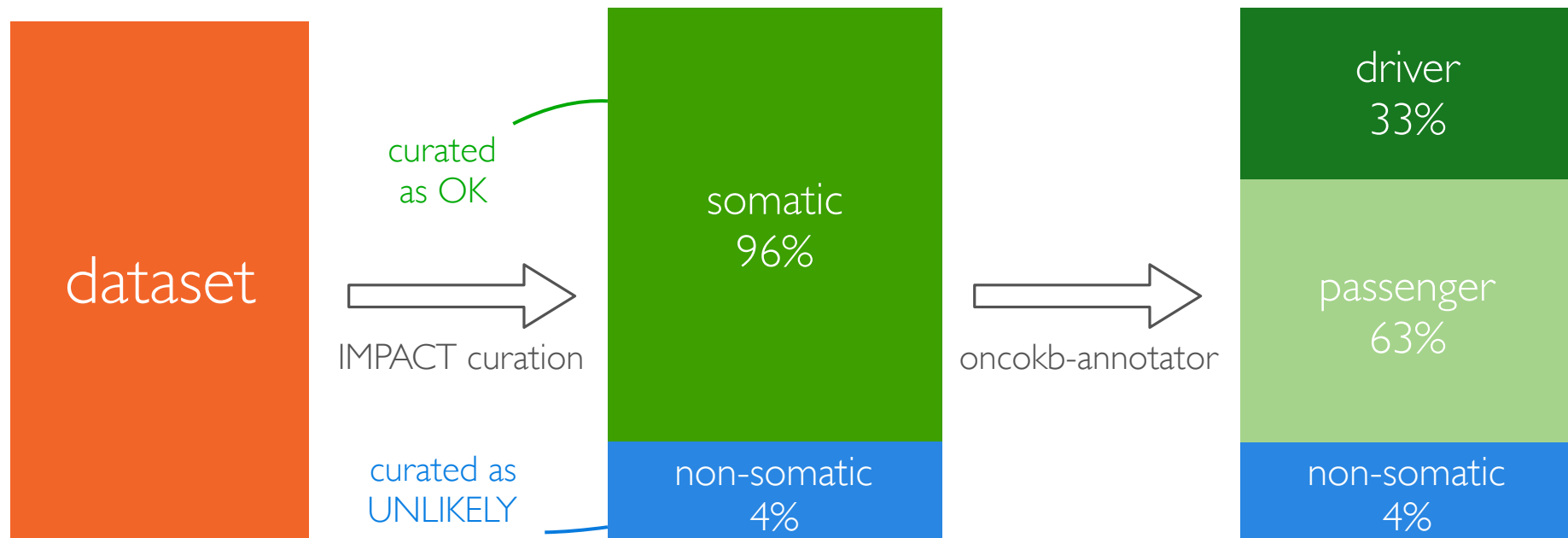
Tumor VAF for somatic vs non-somatic

non-somatic = artefact + germline

# Different classes



dataset → IMPACT curation (curated as OK / curated as UNLIKELY) → somatic 96% / non-somatic 4% → oncokb-annotator → driver 33% / passenger 63% / non-somatic 4%
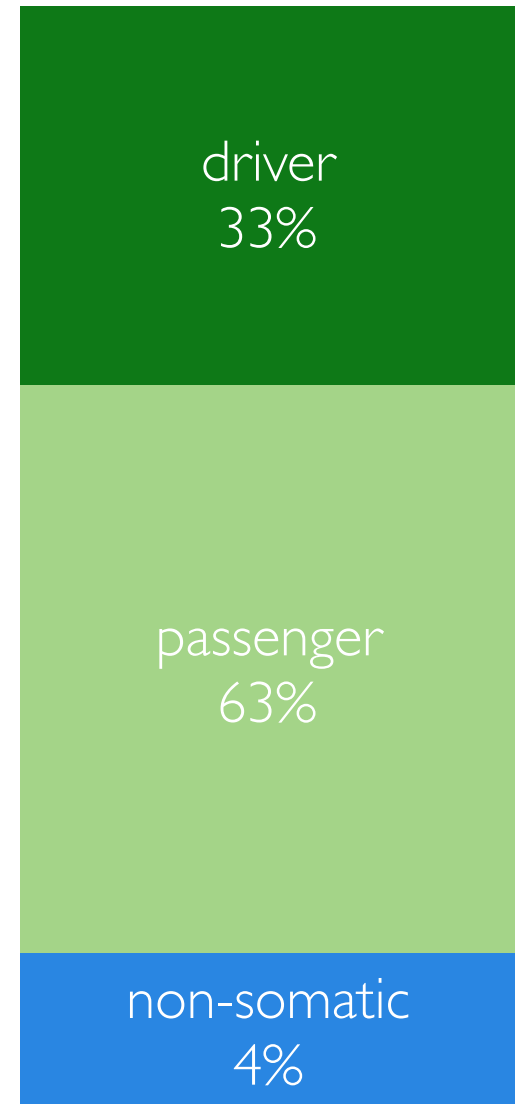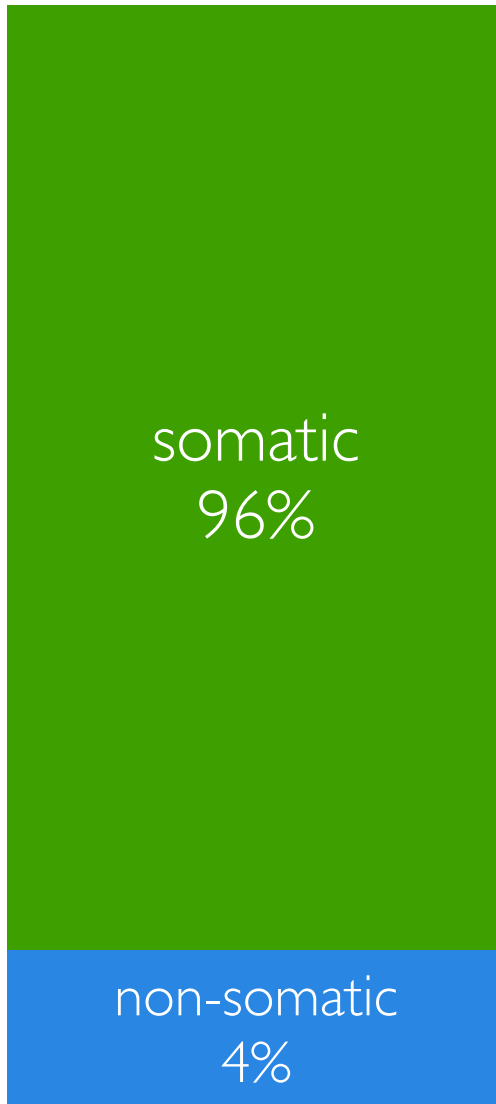
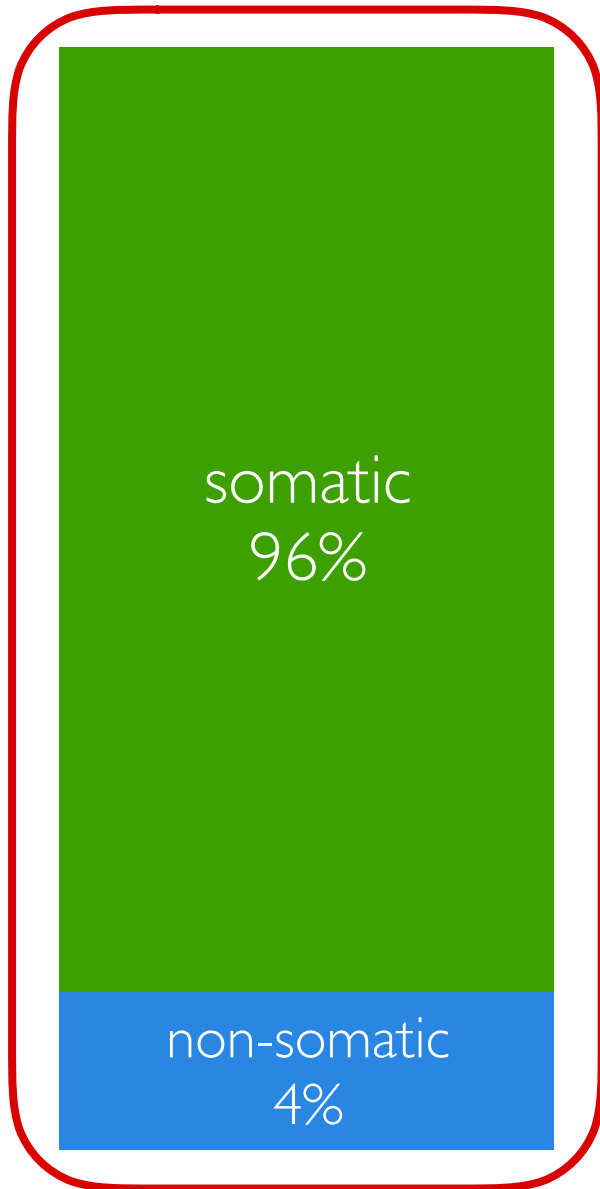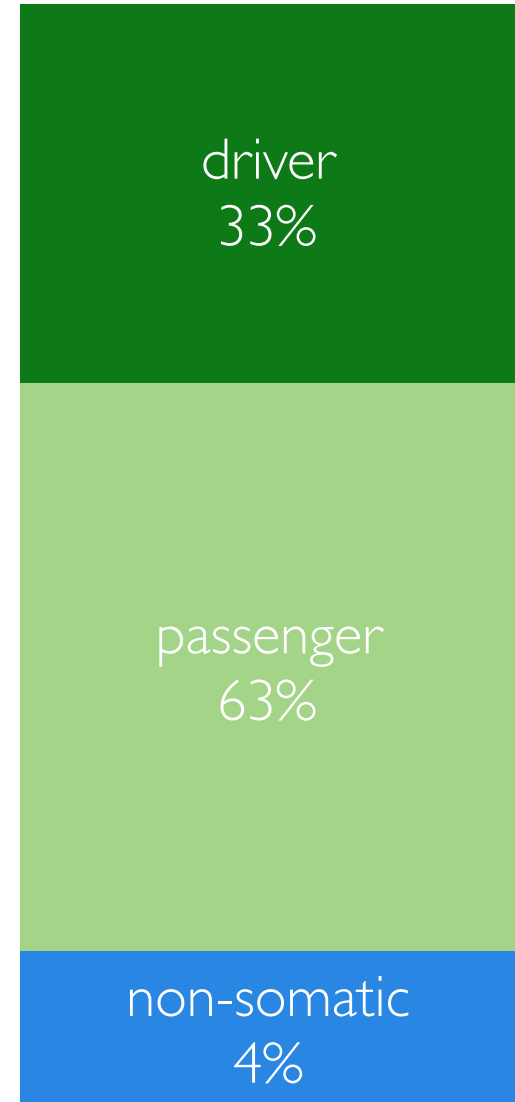Tumor VAF for somatic vs non-somatic

non-somatic = artefact + germline

# Supervised classification

# Supervised classification



2 class:
somatic | non-somatic

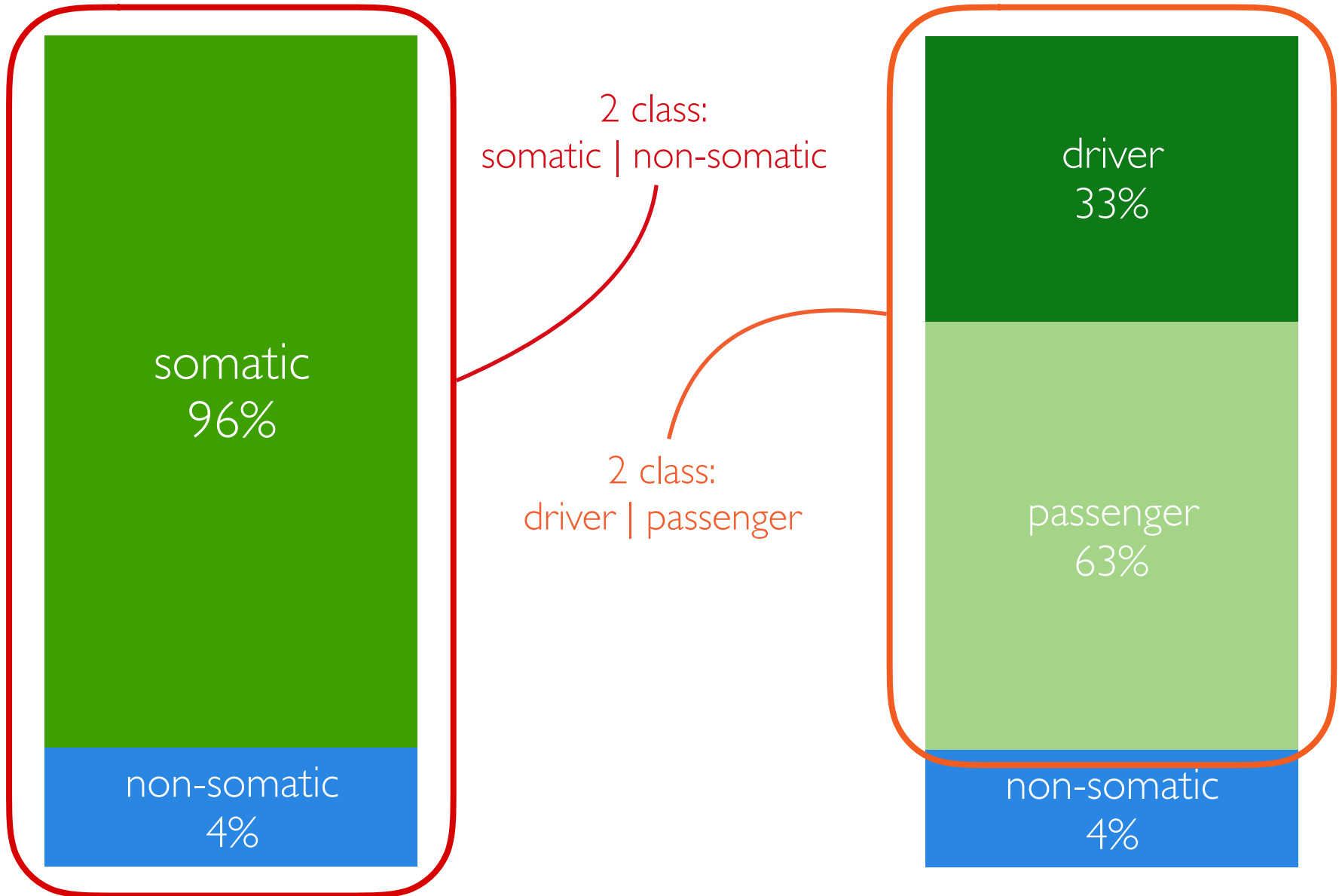somatic
96%

non-somatic
4%

driver
33%

passenger
63%

non-somatic
4%

# Supervised classification



2 class:
somatic | non-somatic

driver
33%

passenger
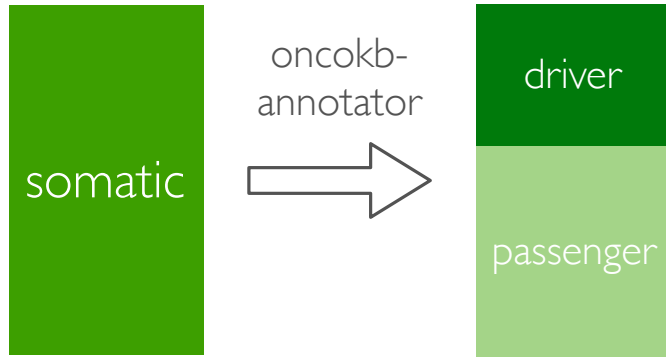63%

non-somatic
4%

2 class:
driver | passenger

somatic
96%

non-somatic
4%

# Supervised classification



somatic
96%

non-somatic
4%

2 class:
somatic | non-somatic

2 class:
driver | passenger

3 class:
driver | passenger
| non-somatic

driver
33%

passenger
63%

non-somatic
4%

## Variant annotation databases

somatic → oncokb-annotator → driver / passenger

OncoKB driver
34,4 %

CanDL, DoCM,
CIViC or CGI driver
0,3 %

passenger
65,3 %

# Literature

## Variant annotation databases

somatic

oncokb-annotator →

driver

passenger



OncoKB driver 34,4 %

CanDL, DoCM, CIViC or CGI driver 0,3 %

passenger 65,3 %

## Existing algorithms

→ SIFT, PolyPhen-2: predicts whether an amino-acid substitution affects protein function
  • not cancer specific

→ CHASM, FATHMM, CanDrA, CScape, rDriver, …: supervised learning for driver classification
  • no sequencing features (features: SIFT, PolyPhen-2, conservative features, genetic features, other prediction algorithm features…)
  • smaller and not "real" dataset (driver: COSMIC / passenger: synthetic, dbSNP, …)
  • never done on somatic vs non-somatic: applied on <u>real somatic mutations only</u>

non-synonymous SNVs only (no indels)

# The features

**NGS features**
- VAF
- strand ratio
- …

**Frequency in normals**

**Population AF**
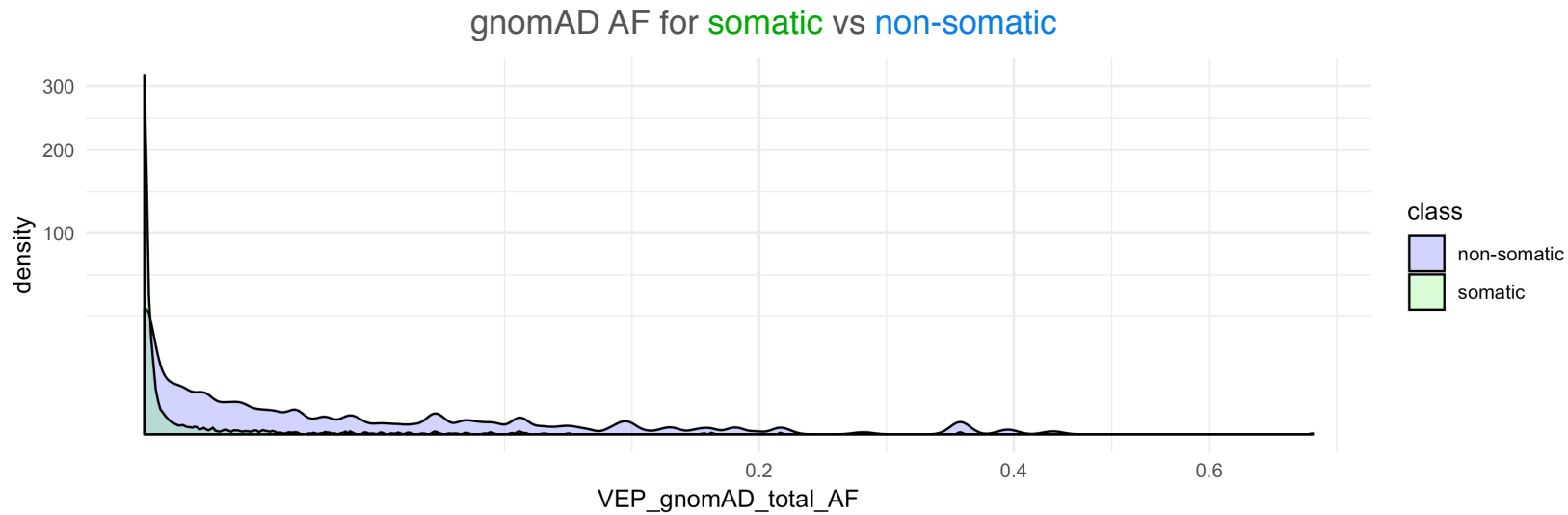- Kaviar
- gnomAD

**COSMIC score**

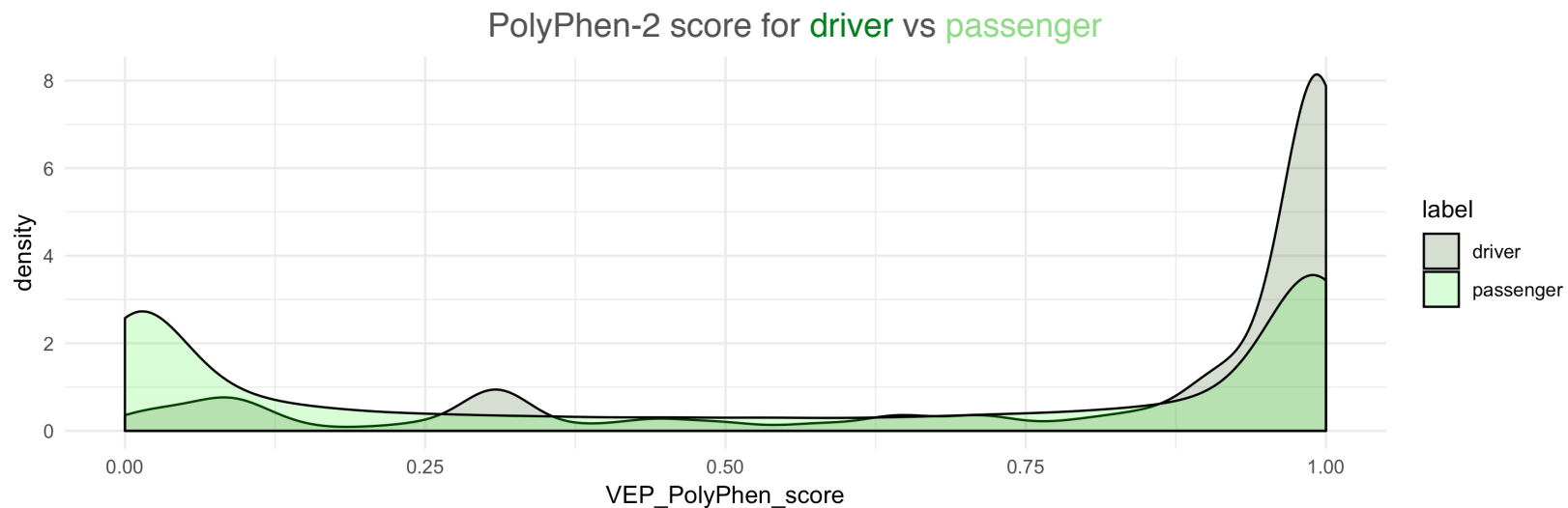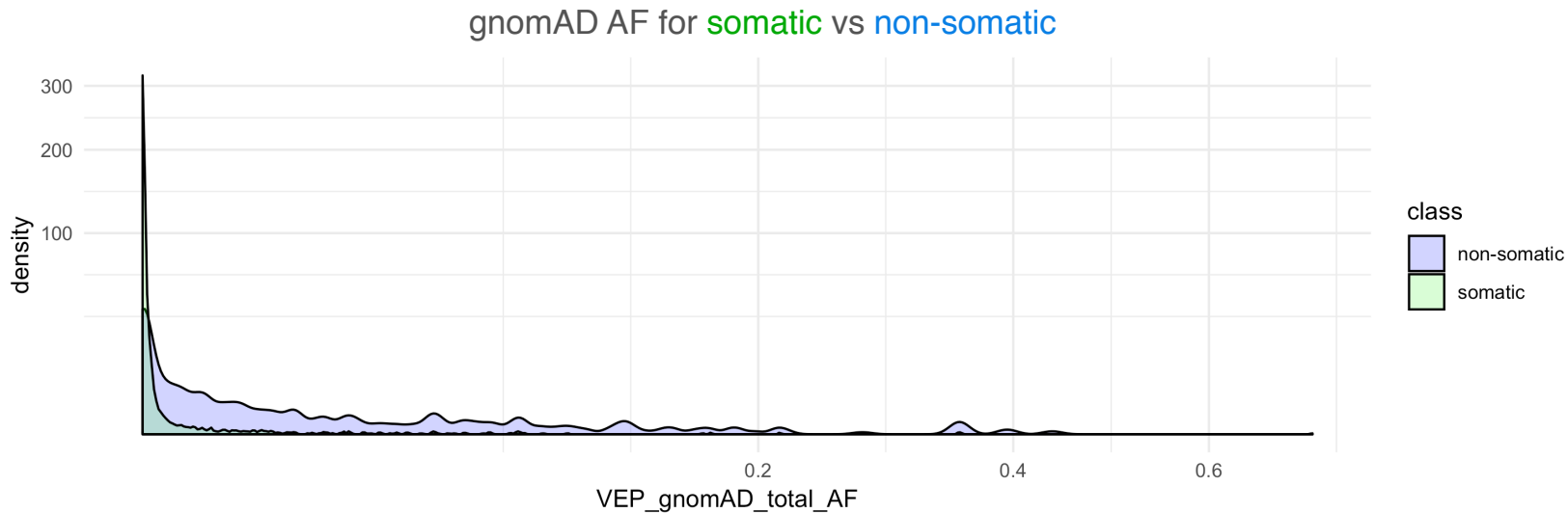**Impact scores**
- SIFT
- PolyPhen-2
- …

**Genetic features**
- gene and gene type
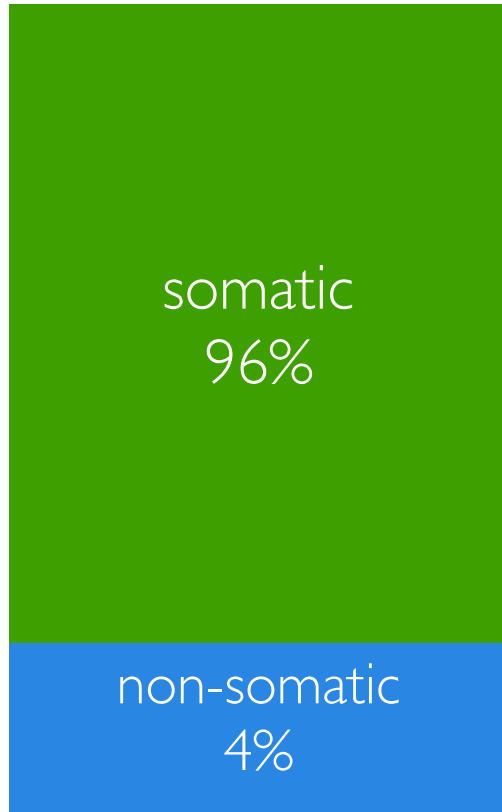- mutation effect
- …

Any ideas?

# Example of informative feature

somatic

non-somatic



gnomAD AF for somatic vs non-somatic

# Example of informative feature



somatic

non-somatic

## gnomAD AF for somatic vs non-somatic

class
- non-somatic
- somatic

driver

passenger

## PolyPhen-2 score for driver vs passenger

label
- driver
- passenger

# Dataset complexity for machine learning

Inconsistencies

somatic
96%

non-somatic
4%

# Dataset complexity for machine learning

Inconsistencies



somatic
96%

non-somatic | 25% | 10% | 8%

classified at least once as somatic

OncoKB driver

both

# Dataset complexity for machine learning

# Dataset complexity for machine learning

## Inconsistencies

somatic
96%

non-somatic

25%

10%

8%

classified at least
once as somatic

OncoKB driver

both

## Imbalanced dataset

somatic
96% (187,012)

non-somatic
4% (7,199)

## Highly recurrent mutations

78% of the driver mutations are <u>unique</u> across the dataset

BUT

100 mutations are shared by 45 patients or more
(ex: KRAS p.G12D, >1000 patients)

- Overfitting on hotspots?
- Performance on « rare » drivers?

# First results

- <u>not representative yet</u>

- logistic regression (l2-regularization | ridge regression)

- 10,000 samples

- 5-fold stratified cross-validation

- driver vs passenger