

TCGA somatic mutations (32 cancer types)

**Train**

missense mutations

**Train**

Semi-supervised labels:  
2,051 likely-driver  
623,996 likely-passenger

CHASmplus

**Predict**

CHASmplus score

$\times$

20/20+

**Predict**

gene driver score

$=$

gwCHASmplus score

**Gene hold-out cross-validation**

Train Test

fold 1

fold 2

fold 3

fold 4

fold 5

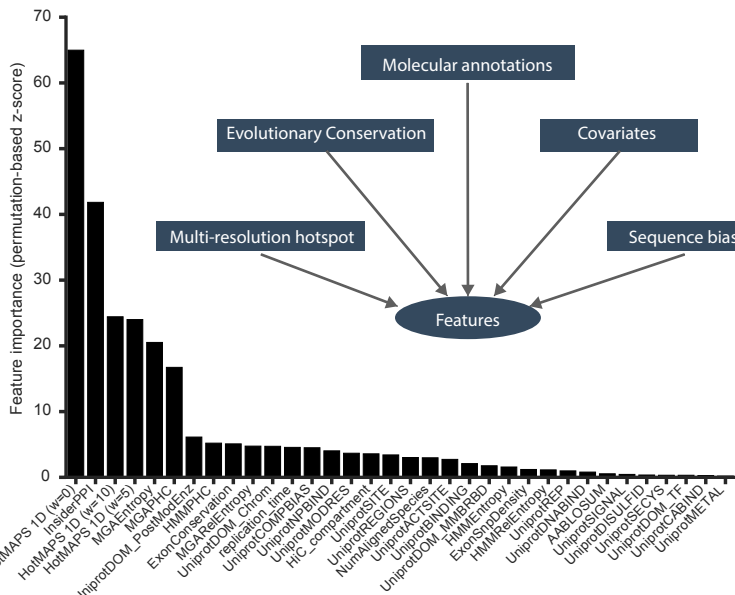
fold 6

fold 7

fold 8

fold 9

fold 10



The flowchart illustrates the gwCHASMplus simulation and analysis pipeline. At the top, a box labeled "Individual cancer types and Pan-cancer" contains "Cancer type 1", three ellipses, and "Cancer type 32".

The pipeline branches into two main paths:

- Left Path (Real Data):** An arrow points down from the top box to an oval labeled "Features". A thick vertical arrow labeled "Predict" points down from "Features" to the text "gwCHASMplus scores". Another thick vertical arrow points down from "gwCHASMplus scores" to "score p-value". A thick vertical arrow labeled "Benjamini-Hochberg" points down from "score p-value" to "Significant mutations ( $q \leq 0.01$ )".
- Right Path (Simulation):** An arrow labeled "Simulation controlling for sequence aspects" points down from the top box to a box labeled "Simulated mutations". This is followed by "x 10". An arrow points down from "Simulated mutations" to an oval labeled "Simulated Features". A thick vertical arrow labeled "Predict" points down from "Simulated Features" to the text "simulated gwCHASMplus scores".

A horizontal arrow points from "score p-value" to the "simulated gwCHASMplus scores" plot.

The "simulated gwCHASMplus scores" plot shows a 1-CDF curve. The y-axis is labeled "1-CDF" and ranges from 0.0 to 1.0. The x-axis is labeled "gwCHASMplus" and ranges from 0.0 to 1.0. The curve starts at (0, 1) and drops sharply to near zero by x=0.2.