

Supplementary material for “Every which way? On predicting tumor evolution using cancer progression models”

Ramon Diaz-Uriarte, Claudia Vasallo Vega
Dept. Biochemistry, Universidad Autónoma de Madrid
Instituto de Investigaciones Biomédicas “Alberto Sols” (UAM-CSIC)
Madrid, Spain*

2018-07-30 (Release: Rev: 5f59ef1)

Contents

1	Generating random fitness landscapes	3
1.1	Local maxima without reciprocal sign epistasis?	3
1.2	Rough Mount Fuji	4
2	Plots of fitness landscapes and inferred DAGs	5
3	Simulations	6
3.1	Runs until fixation	6
3.2	All genes part of lines of descent with frequency > 0.001	6
3.3	Detection regimes: sampling	6
3.4	Other parameters of the simulations	7
4	Material and methods: others	7
4.1	Terminology	7
4.2	CPM software	7
4.3	Preprocessing of data for CPMs	7
4.4	Computing probabilities of paths	8
4.5	Example where perfect recall and precision do not guarantee Jensen-Shannon divergence of 0	10
4.6	Measuring predictability: comparing paths from CPMs and LODs of different lengths	10
4.6.1	Commented example for paths of unequal length	11
4.7	Coefficients of linear models	13
5	Cancer data sets	14
5.1	Cancer data sets: sources and characteristics	14
5.2	Bootstrapping on the cancer data sets	15
6	CAPRI, CAPRESE, and paths of tumor progression	16
7	Fitness landscapes: characteristics, evolutionary predictability, clonal interference, and sampled genotypes	19

*ramon.diaz@iib.uam.es, rdiaz02@gmail.com, <http://ligarto.org/rdiaz>

8 Overall patterns for the six methods	37
9 OT and CBN, JS, weighted vs. unweighted	39
10 CAPRI and CBN, 1-precision, unweighted	40
11 CAPRESE and OT, 1-precision, unweighted	41
12 Probability of recovering the most common LOD	42
13 Number of paths inferred	44
14 Slopes of regressions of 1-recall and 1-precision on LOD diversity, S_p	45
15 Coefficient of variation of S_c	46
16 Estimated S_c by CBN	47
17 Analysis of deviance tables for fitted models	48
17.1 Models fitted to the complete data set	49
17.1.1 Two-way interactions	49
17.1.2 Three-way interactions	50
17.1.3 Four-way interactions	51
17.2 Models fitted to each combination of fitness landscape by method	53
17.2.1 Main effects	53
17.2.2 Two-way interactions	55
17.2.3 Four-way interactions	58
18 Number of mutations of local maxima and performance	63
19 LOD and CPM diversity: ratios and slopes	65
20 Regression of individual CBN unpredictability estimates on LOD diversity	65

1 Generating random fitness landscapes

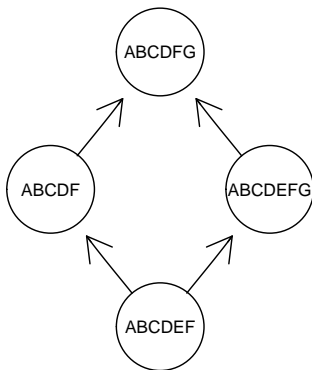
For the representable and local-maxima fitness landscapes, we started by generating random DAGs. Since no agreed upon model exists for the distribution of DAGs in CPMs, we have used two different procedures, choosing each one randomly with the same probability. One of the procedures uses the function `sim0Graph`, in the `OncoSimulR` package. To generate random DAGs with `sim0Graph` for N genes, the genes were first randomly split in a number of levels, where the number of levels used was a randomly chosen integer between 3 and $N - 1$, both included. Then, each gene from each level i was randomly connected (as descendant) to randomly chosen genes (the ancestors) from levels j , where $j < i$; the number of incoming connections of each gene is a randomly chosen integer between 1 and $maxp$ (both included), where $maxp$ is a randomly chosen integer between 2 and $N - 2$ ($maxp$ is common for all genes in a DAG, but can vary between DAGs). The final DAG is the transitive reduction of the above generated DAG. (Note that this procedure can occasionally result in star DAGs, i.e., DAGs without any dependencies; in such a case, the DAG was discarded and a new one obtained). The other procedure uses the function `random_poset` in package `MC-CBN` (<https://github.com/cbg-ethz/MC-CBN>); this function is undocumented, but it returns the transitive reduction of a randomly-filled adjacency matrix for a DAG where the initial number of non-zero connections is equal to the number of possible connections * constant; we used the default value of that constant (0.15).

1.1 Local maxima without reciprocal sign epistasis?

As explained in the paper, creating fitness landscapes with local maxima generally results in creating reciprocal sign epistasis and the number of local fitness maxima is associated with reciprocal sign epistasis —see Figures in section 7. There were, however, seven cases (out of 420) where introduction of local fitness maxima did not lead to introduction of reciprocal sign epistasis. These cases (which can be seen in the files referred to in 2) are landscapes with IDs “7E10pIyu7UguIUl8I”, “8QIFQCUIUVfC10PZr”, “bCsk2Qo5VMVm55fM”, “GedZaWDeb1029Mf88”, “hw8kQ4g44p4XAkDa”, “WpF105HbEDoECa8vs”, “t1yUXsv5fVuo10GRi”. To look at one example in detail, we will use “GedZaWDeb1029Mf88” (it is the smallest one). The fitness of four of the relevant genotypes are

ABCDFG : 2.007
ABCDEF: 1.8
ABCDF : 1.749
ABCDEF: 1.712

so there is no reciprocal sign epistasis (use, for example, the graphical criterion in [7] or [13]) and “ABCDFG” is a global maximum.



Of course, under an evolutionary model that assumes no back mutations (as is the case for CPMs), two of those transitions, those that involve losing “E” ($ABCDEF \rightarrow ABCDF$ and $ABCDEFG \rightarrow ABCDFG$) are not allowed, leading to two local maxima.

Note also that here, for that set of four genotypes, mutating gene E decreases fitness. But mutating E increases fitness in genotypes “ABC” or “ABCD”. Thus, this fitness landscape does not fulfill either the assumption that a mutation never decreases the probability of acquiring other mutations (even if the fraction of pairs of genotypes with reciprocal sign epistasis is 0). Regardless, one can also simply focus on the fact that this fitness landscape contains local maxima (and is missing paths relative to the corresponding fitness graph from the DAG of restrictions).

1.2 Rough Mount Fuji

In the Rough Mount Fuji random fitness landscapes the reference genotype (i.e., the genotype with maximum fitness) was randomly chosen (setting `reference = 'random'` in the `rfitness` function in OncoSimulR). The standard deviation, sd , of the random normal variate was set to 0.2 and the decrease in fitness (strictly, birth rate) of a genotype per each unit increase in Hamming distance from the reference genotype, c , was chosen from a uniform $U(0, 0.2)$ distribution. This gives a wide variety of fitness landscapes that encompass from close to additive (large values of c) to House of Cards (c close to 0), with maximum fitness (birth rate) comparable to those of the representable and local-peaks fitness landscapes.

The generated RMF fitness landscape was checked to ensure that all seven or ten genes were present in at least one accessible genotype; if they were not, a new fitness landscape was generated (with, possibly, different values of c and reference genotype). Function `rfitness` from the OncoSimulR package [9] was used.

2 Plots of fitness landscapes and inferred DAGs

Files `f1-fg-7.pdf` and `f1-fg-10.pdf` show the 1260 fitness landscapes used. Each `f1-fg-x.pdf` shows the 630 fitness landscapes used for x genes. In each PDF, the first 210 fitness landscapes are fully representable fitness landscapes, the next 210 (pages 211 to 420) are the “Local maxima” fitness landscapes. The next 210 (pages 421 to 630) correspond to Rough Mount Fuji (RMF) models.

In each page, the following figures/tables are shown:

DAG of restrictions (top left) The true DAG of restrictions. This only applies properly to the representable fitness landscapes. In the local maxima fitness landscapes, not all paths between genotypes are available. (See below). For the RMF landscapes this, of course, is not available, since there is no underlying DAG of restrictions.

Fitness landscape (bottom left) As it says, the fitness landscape. Boxes surround the fitness maxima in the seven-genes landscapes. To minimize clutter, genotype labels are not shown in the 10-genes landscapes.

ID and landscape characteristics (center) The ID of the fitness landscape (a random string that matches the value of ID in the data tables), the number of accessible genotypes, the number of fitness maxima or peaks. Values for “removed edges” and “prop rem edges” denote, respectively, the number and proportion of edges in the fitness graph that were removed; these only applies to the “Local maxima” landscapes; thus, these values are 0 for “Representable” fitness landscapes, and NA for RMF fitness landscapes.

Fitness graph from DAG of restrictions (top right) The fitness graph implied by the DAG of restrictions. Not available for RMF (since there is no DAG of restrictions).

True fitness graph (bottom right) The actual fitness graph that corresponds to the fitness landscape. For “Representable” fitness landscapes this would be the same as the Fitness graph from DAG of restrictions, so we do not show it (to make the size of the files smaller).

3 Simulations

3.1 Runs until fixation

Simulations were run until fixation of a genotype, where the genotype was one of the genotypes among the local maxima (or the single global maximum). We used OncoSimulR, with the `fixation` option (introduced in version 2.9.8 of the program). A genotype was considered to have been fixated if it maintained a proportion ≥ 0.98 during 15000 consecutive sampling periods (this means that if after reaching a minimum frequency ≥ 0.98 , at any time the proportion became smaller than 0.98 the counter of successive periods was reset to 0). We cannot ask for a fixation with a proportion of 1.0 because for local maxima, if mutation rate is larger than 0 and neighboring genotypes have non-zero birth rate, the fixated genotype can occasionally generate descending genotypes that exist, with small frequencies, for short periods of time. Using much shorter number of consecutive sampling periods such as 1000 or 5000 did not produce different results over using 15000 in trial runs; however, to err on the safe side and make sure fixation had been established, we used that overly long period.

These 15000 periods were excluded from the computation of clonal interference statistics.

3.2 All genes part of lines of descent with frequency > 0.001

When the 20000 simulations were completed, we verified that the frequency of all genes in the last genotypes (i.e., the fixated genotypes or the final genotypes of the LODs) were at least 0.001. If they were not, a new fitness landscape was generated and the processes started again. In other words, we avoided fitness landscapes that have a nominal number of, say, 10 genes, but where a smaller number of genes were effectively ever part of the paths of tumor progression (this issue can affect the local maxima and RMF landscapes). With the threshold of 0.001, in a sample of 4000 individuals, the probability that the gene with smallest frequency is never part of a LOD is about 0.018 ($= (1 - 0.001)^{4000}$), so less than 2%.

3.3 Detection regimes: sampling

For each detection regime, we generated 20000 random deviates (called r , below) from the specified beta distribution ($B(1, 1)$, $B(5, 3)$, and $B(3, 5)$) (for uniform, large, and small, respectively).

Using the those random deviates, we defined the target size of each sample as $t = \exp(r (\ln(M) - \ln(m)) + \ln(m))$, where M and m are the largest and smallest values, respectively, of population sizes ever attained in any of the 20000 simulations. Thus, we obtain target sizes that are uniform or biased towards large sizes or biased towards small sizes in the log scale. In the model of [18], tumor population size increases logarithmically with number of driver mutations. Therefore, uniform, small, and large biases would correspond to approximately uniform, small, or large in terms of number of driver mutations.

For each of the 20000 simulations, the actual sample was the one corresponding to the first sampling period at which the total tumor size achieved a value equal to, or larger than, t . If all values of tumor population size were $> t$, we returned the sample with the largest population size, and if all values were $< t$ the sample with smallest size.

This procedure determines at which of the sampling times we take the sample. The actual genotype returned is the single genotype with the largest frequency. Thus, we are not emulating whole-tumor sampling (bulk sequencing) but, rather, single-cell sampling, and sampling the single most common genotype.

We carried the above steps using OncoSimulR's function `samplePop`, with the values of t (thresholded as explained for $> t$ and $< t$) as arguments to `popSizeSample` and using `typeSample = 'single'`.

3.4 Other parameters of the simulations

Simulations used the implementation of the McFarland model in the OncoSimulR package [9]. In addition to the parameters specified in the main text, other parameters for the simulations on the 500 fitness landscapes were (see specific meaning in documentation of OncoSimulR [9]): $finalTime = 10000$, $keepEvery = 1$, $sampleEvery = 0.03$, $max.wall.time = 20$, $max.num.tries = 500$.

4 Material and methods: others

4.1 Terminology

The **number of local (fitness) maxima** is a static feature of the landscape (number of genotypes such that all genotypes withing a distance of one mutation have lower fitness). The number of **observed local (fitness) maxima** can be smaller, since some peaks (local maxima) might never be visited. For representable fitness landscapes, both numbers are 1. For the other two landscapes, those numbers were ≥ 2 .

4.2 CPM software

Other methods for cancer progression models have been described but either are too slow for routine use such as [22], or have dependencies on external libraries that are not open source such as DiP [12], or have no software available (e.g., [1, 6]). See further details in [8].

For CBN version 0.1.04b from March 2016, and still current as of April 2018, was downloaded from <https://www.bsse.ethz.ch/cbg/software/ct-cbn.html>. Defaults for CBN were used. I wrote a wrapper to call their code from R, and I used the default settings for temp ($-T = 1$) and steps ($-N = \text{number of nodes}^2$); I started the simulated annealing search for the best poset from an initial poset built using Oncogenetic Tress [23], as preliminary runs suggested this initial poset is as good as, or better than, the default linear poset in [15]. The parameters for the transition rates between genotypes (λ s) were obtained doing an additional run on the fitted model from the previous step, as in [15].

MCCBN was run using version 1.1.9 of the MC-CBN package, downloaded from github (<https://github.com/cbg-ethz/MC-CBN>).

OT was run using version 0.3.3 of the Onctoree package [23].

For CAPRI and CAPRESE we used version 2.11.0 of TRONCO, current as of April 2018, downloaded from the official BioConductor site. All options were left at the recommended defaults (e.g., 100 bootstrap samples for the estimation of the selective advantage scores with p-value of 0.05, and heuristic search using Hill Climbing).

Wrapper code was written for all methods to obtain the fitness graphs and, for OT, CBN, and MCCBN, the weighted predicted paths. For OT, we use `ot.fit$parent$est.weight` to obtain the probabilities of transition to each descendant genotype; if the OT fit, however, cannot return an error estimate, that operation fails and thus we use the `ot.fit$parent$obs.weight` component.

4.3 Preprocessing of data for CPMs

Before analyzing data with CPMs, data were preprocessed as follows:

- All columns that had all 0s (i.e., genes that were absent in all samples) were removed. Since these are never present in the data given to the methods, no inference can be made about the removed genes, and this necessarily decreases the dimension of the fitness landscape implied by the CPM and the length of the paths to the maximum.

- If two or more columns (genes) were identical over all individuals (i.e., were indistinguishable), all the identical replicate columns, except one, were removed from the analyses¹. This, of course, will preclude the matching of some (or all) of the true paths since we are constructing the CPM from a data set of smaller dimensionality and the CPM's paths are shorter than they ought to be (see also details in section 4.6).

Indistinguishable events will unavoidably create problems. Alternative ways of handling them are not better. If the indistinguishable columns are left in the data, some methods (e.g., CAPRI and CAPRESE) complain about it, whereas others (CBN, MCCBN) make the indistinguishable events depend on one another, with a very large λ , with the order in the DAG depending on the order on the column of data (leftmost events placed as ancestors). OT also places them as independent events (as CAPRI and CAPRESE do).

The consequence of leaving the events as in CBN would be similar to expanding the paths of progression, post-analysis, and placing the indistinguishable events one right after the other in the path. The order would, of course, have to be arbitrary and, in most cases, this would actually make matching any true LOD harder. If only one of the replicates is left in the path, the LOD needs to match, by chance, one particular order. If two or more are placed, the probability of matching decreases.

The proportion of data sets with one or more columns is about .16, .16, and .11 for the representable, local maxima, and RMF landscapes. They decrease with sample sizes, generally being under 0.05 for sample size 4000.

- Whenever one or more genes were present in all samples, to prevent the removal of these genes present in all cases, we added one case (one "pseudosample") with no mutations to the data set (this is not unlike [8], but we add only one sample, not a fixed percentage, to minimize altering any estimates of probabilities of paths). This allows us to use exactly the same data for all methods (CAPRI cannot deal with data where one or more columns are present in all subjects, OT removes them, whereas CBN can use this data).

Even more importantly, this procedure does not decrease the dimensionality of the data set and, thus, does not decrease the length of the CPM's paths to the maximum

The event that has a frequency of 1 is placed at the top of the DAG of restrictions (it is the first mutation after WT in all the paths to the maximum). The (very minor) inconvenience is that it has a minor effect on CBM's λ s estimates, but that should be inconsequential for practical purposes.

The proportion of data sets with pseudosamples added was .30, .21, and .01, for representable, local maxima, and RMF fitness landscapes.

4.4 Computing probabilities of paths

The procedure we used is as follows (it might be simpler to understand it by referring to p. i729 of Montazeri et al., 2016 [20]):

1. Obtain the set of genotypes that can exist under the poset (as we will use it in step 3 below).
2. Obtain the set of paths that can exist under the poset (to be used in step 5, below). This itself is obtained from 1.
3. Obtain the transition rate matrix between genotypes from the lambdas (e.g., what is shown in matrix S in Montazeri et al.). As explained in Montazeri et al., "the non-zero off-diagonal elements of the transition matrix are the transition rates from each genotype to

¹In more detail, the identical columns were flagged as such by "fusing" the names of the genes, so as to be able to identify them. Then, the paths were post-processed before the analysis to remove the combined name, leaving one of the two, or more, identical names.

its successive genotypes in the genotype lattice, also shown in Figure 1(b).” See also legend of Figure 1: “(b). Directed transition rates among neighboring genotypes are shown on the edges of the lattice”.

4. Set the diagonal of the previous matrix to 0² and for each row of the transition rate matrix, divide by $\sum \lambda$. Now the entries are probabilities of transition to each descendant genotype given a transition.
5. Go over the list of paths (from step 2) and for each path, obtain its probability by multiplying the probabilities of the transitions (from step 4) between the genotypes in a path.
6. (Check: verify sum of probabilities of all paths equals 1, within numerical margin of error of machine.)

In the code, steps 3 and 4 above were carried out by creating, from the set of paths to the maximum and the output from CBN, what we called a weighted fitness graph: the fitness graph of paths to the maximum with the lambdas on the edges (or the weighted adjacency matrix corresponding to paths to the maximum where weights are lambdas). This would be Figure 1b in [20]. Dividing by the row sum gives us the transition matrix between genotypes in step 4.

As an example, suppose we obtain from CBN the following DAG of restrictions and estimated lambdas:

From	To	λ
Root	A	2
Root	B	3
A	C	4
C	D	5

The paths to the maximum, with their probabilities, are:

path	probability
WT \rightarrow A \rightarrow A, B \rightarrow A, B, C \rightarrow A, B, C, D	$2/5 * 3/7 * 1 * 1$
WT \rightarrow A \rightarrow A, C \rightarrow A, B, C \rightarrow A, B, C, D	$2/5 * 4/7 * 3/8 * 1$
WT \rightarrow A \rightarrow A, C \rightarrow A, C, D \rightarrow A, B, C, D	$2/5 * 4/7 * 5/8 * 1$
WT \rightarrow B \rightarrow A, B \rightarrow A, B, C \rightarrow A, B, C, D	$3/5 * 1 * 1 * 1$

For example, from WT with a probability of 2/5 we take the path to A and with 3/5 to B. Once in genotype A, we can either add a B mutation (probability = 3/7) or a C mutation (probability = 4/7). If we add a B mutation, from genotype AB we can only move to ABC (probability 1). Etc. This procedure is equivalent to the one used by Hosseini [16].

An analogous procedure was used with OT.

²In fact, the diagonal entries are never computed explicitly and are always 0.

4.5 Example where perfect recall and precision do not guarantee Jensen-Shannon divergence of 0

Suppose the following set of two paths to the maximum, with predicted and observed probabilities as shown:

Path	CPM predicted probability	True LOD probability
WT \rightarrow A \rightarrow AB	0.99	0.01
WT \rightarrow B \rightarrow AB	0.01	0.99

The JS divergence (on a scale 0 to 1) is 0.9192 (remember 1 is the maximum value of divergence), even when 1-recall and 1-precision are both 0 (none of the CPM's paths are missing from the LODs, and none of the LODs are missing from the CPM's paths).

4.6 Measuring predictability: comparing paths from CPMs and LODs of different lengths

Let i and j denote two paths, one from the LOD and the other from the CPM, with corresponding probabilities p_i and q_j . Here, the i index refers to paths from the LOD, and the j to paths from the CPM (this is in contrast to the paper, where we did not make this specification, to keep the description general).

Let K_i, K_j denote the length of paths i and j , respectively. Note that all K_j are equal, and equal to the single K^C (as all go up to the genotype with all m genes mutated). When there is ambiguity about which K we are referring to, we will use K_j^C for the K_j from the CPM (again, $K_1^C = K_2^C = \dots = K_C$) and K_i^L for the K_i from the LOD.

The vectors P, Q , for the computation of JS will have the following types of matching pairs:

1. p_i, q_j when $K_i = K_j$
2. $p_i \frac{K_j}{K_i}, q_j$, when i is partially included in j ($K_i > K_j$),
3. $p_i, q_j \frac{K_i}{K_j}$ when j is partially included in i ($K_j > K_i$),
4. $\sum p_i \frac{K_i - K_j}{K_i}, 0$ when i with $K_i > K_j$,
5. $0, \sum q_i \frac{K_j - K_i}{K_j}, 0$ when i with $K_j > K_i$,
6. $\sum p_u, 0$, for all paths u among the paths from the LOD that do not match any j (any path from the CPM),
7. $0, \sum q_v, 0$, for all paths v in among the paths from the CPM that do not match any i (any path from the LOD).

(Where some notation above, again, can be simplified by noting that all $K_j = K_C$).

We can sum, as appropriate, the relevant entries to simplify computations as the JS is the same for the pair of vectors $P = [p_1, p_2, 0, 0, p_5, p_6], Q = [q_1, q_2, q_3, q_4, 0, 0]$ and the pair of vectors

$P' = [p_1, p_2, 0, p_5 + p_6], Q' = [q_1, q_2, q_3 + q_4, 0]$ (this follows from the definition of JS).

In other words, we have the 0 entry in Q correspond to $\sum p_i \frac{K_i - K_j}{K_i} + \sum p_u$, where the i are the paths in the LOD with some partial match among the paths in the CPM, and the u denote those paths in the LOD that do not match any path in the CPM.

The relevant entries above can be used to compute 1-recall and 1-precision. For example, for 1-recall, $P(\neg DAG|LOD)$, the sum of the probabilities of the paths in the LODs that are not among the paths allowed by the CPMs, we will use $\sum p_i \frac{K_i - K_j}{K_i} + \sum p_u$.

4.6.1 Commented example for paths of unequal length

To give a specific example, suppose the following paths from a LOD and a CPM:

Path (i)	LOD frequency	p_i
1	WT → A	0.1
2	WT → B → AB	0.3
3	WT → B → AB → ABC	0.1
4	WT → B → BC → ABC	0.2
5	WT → B → BC → ABC → ABCD	0.3

Path (j)	CPM predicted probability	q_j
1	WT → A → AB	0.6
2	WT → B → AB	0.4

Where $i = 1, 2, 3, 4, 5$ are the four LODs and $j = 1, 2$ the two paths to the maximum from the CPM. To compute JS, 1-recall and 1-precision, it is much simpler to use an algorithm that splits the cases to be considered into three:

1. $K_i < K_j$ (i.e., $K_i < K_C$), those paths from the LOD that are shorter than the paths from the CPM;
2. $K_i = K_j$ (i.e., $K_i = K_C$), those paths from the LOD that have the same length as the paths from the CPM;
3. $K_i > K_j$ (i.e., $K_i > K_C$), those paths from the LOD that are larger than the paths from the CPM;

We can iterate over all distinct k for $K_i < K_j$, and weight the output by w_k , the sum of all paths from the LOD that end at k mutations; computations for $K_i = K_j$ can be subsumed into those for $K_i < K_j$. Thus, one part of the algorithm iterates over all $k \leq K_j = K_C$ (remember all K_j are identical and equal to K_C , the single, unique K at which the paths from the CPM stop). Computations for $K_i > K_j$ can be done in one iteration.

It might also be helpful to think about a cut operation on a path. For instance, for each k where $K_i < K_j$, we can cut the paths from the CPM at k mutations, leaving only the paths from WT up to k mutations (and collapsing, as appropriate, any collection of now indistinguishable subsets of paths, summing their probabilities).

In the exposition below, some computations could be further simplified; they are left as they are for clarity (e.g., we multiply by total frequencies of mutations for the weights when we have previously scaled the total probability by it, so it is 1 in each k , etc).

JS

1. LOD $i = 1$ finishes at one mutated gene. Cutting the CPM path at $k = 1$, the JS for $k = 1$ is obtained from the vectors of probabilities $P = [1, 0, 0]$ (from the LOD) and

$Q = [0.6, 0.5, 0.4, 0.5, 0.5]$ from the CPM. The last entry in Q is the sum of all the flow through the paths of the CPM that cannot be matched because the length of the CPM paths is $K_C = 2$. And the 1 in P comes from $p_1 / \sum p_i$ for all i that end in $k = 1$ which is only p_1 .

The first and second entries are q_1^1 and q_2^1 multiplied by k/K_C , i.e., the probabilities of the fractions of paths from the CPM cut at $k = 1$ mutations.

The weight, w_k , for this value is 0.1 (the frequency of LODs that finish at $k = 1$).

2. LOD $i = 2$ finishes at $k = 2$. Here the comparison is the immediate one for equal length paths and the vectors used for JS are: $P = [1, 0]$ and $Q = [0.4, 0.6]$. $w_2 = 0.3$.
3. Paths $i = 3, 4, 5$ are longer than the CPM paths. The flow $AB \rightarrow ABC$ for $i = 3$, the flow $BC \rightarrow ABC$ for $i = 4$, etc, cannot be matched by the CPM.

Here the total amount of evolutionary flow through the LOD that cannot be captured by the CPM, because the CPM ends prematurely, is $\sum_i p_i (K_i - K) / K_i = (0.1 * (1/3) + 0.2 * (1/3) + 0.3 * (1/2)) * (1/0.6) = 0.25/0.6$, for $i = 3, 4, 5$, where the $1/0.6$ scales relative to the total probability in paths $i = 3, 4, 5$. Then, to compute JS, the two vectors of probabilities would be $P = [0, (2/3) (0.1/0.6), (2/3) (0.2/0.6), (1/2) (0.3/0.6), 0.25/0.6]$, from the LOD, and $Q = [0.6, 0.4, 0, 0, 0]$, for the CPM.

But that is equivalent to using the two vectors

$P = [0, (2/3) (0.1/0.6), (0.5/0.6) + (1/3) (0.1/0.6)]$, $Q = [0.6, 0.4, 0]$. The last entry in P might be easier to see from adding LODs $i = 4, 5$ and the unmatched portion of $i = 3$. Here $w_i = 0.6$

4. We can now add all the JS with their corresponding weights.

1-recall For 1-recall, the sum of the probabilities of the paths in the LODs that are not among the paths allowed by the CPMs, we have: $p_5 + p_4 + p_3 \frac{(K_3 - K_C)}{K_3} = 0.3 + 0.2 + (1/3) 0.1$, where $K_C = 2$ (all CPM paths end at $k = 2$). Note that these same values can be obtained by iterating over k as above, and doing a weighted sum, but in this example it is much simpler to use the computation directly. Had we used weighted sums, we would have got: $0 w_1 + 0 w_2 + (p_5 + p_4 + p_3 \frac{K_3 - K_C}{K_3}) \frac{1}{0.6} w_{\geq 3} = 0.8889 \cdot 0.6 = 0.533 = 0.3 + 0.2 + (1/3) 0.1$ (where the $\frac{1}{0.6}$ scales so that the probabilities considered when $k \geq 3$ add to 1 —and, sure, we are dividing by 0.6 only to multiply by it because the scaling factor is the weight of the k_{ge3} stratum).

1-precision For 1-precision, the sum of the probabilities of the paths in the CPM's paths that are not among the LOD paths (the paths followed by evolution), it is simpler to use a weighted sum:

$$((q_1 \frac{K_C - K_1^L}{K_C}) + q_2) w_1 + q_1 w_2 + q_1 w_{\geq 3}.$$

When $k = 1$, $i = 1$ is $WT \rightarrow A$, and thus all of q_2 is not captured, and half of q_1 is captured; at $k = 2$, all of q_2 but none of q_1 is captured; for $k \geq 3$ again path $j = 1$, with $q_1 = 0.6$ is not captured. Thus, we have $(0.6 * (1/2) + 0.4) * 0.1 + 0.6 * 0.3 + 0.6 * 0.6$.

Note that for 1-precision we do not need to re-scale so that probabilities always add up to 1 because they already do (we consider all the j). This was not the case for 1-recall (where only some of the i might be considered in turn when we iterate over k).

To recap, as stated in the paper, we want to compute JS, 1-recall, and 1-precision taking into account that:

1. Any LOD that finishes at k mutations and matches a CPM path up to k mutations is a perfect match, up to k ; this is the reason we match each LOD with the fitness graph from CPMs cut at the number of mutations of the final genotype of the LOD.
2. Any set of LODs that finishes at k mutations when the CPM goes to K with $K > k$ necessarily misses $(K - k)/K$ of the total evolutionary flow, all that which goes from $k + 1$ to K . This is why we use a category unmatchable by construction.

Without this, it would be possible to obtain perfect JS from LODs that missed most of the evolutionary flow, for instance very short LODs that finished at one mutation.

This part of the procedure, thus, accounts for that part of the evolutionary process that the CPM predicts and is not matched by stopping evolution at a local fitness maximum; remember, again, that by construction the CPMs predict that the evolutionary process should go all the way to a global maximum with all genes mutated.

3. A similar reasoning applies when the paths from the CPM are shorter than the paths from the LOD.

4.7 Coefficients of linear models

Coefficients from the generalized linear mixed-effects models shown in the text are from over-parameterized models, and those are obviously not the models fitted. What I have done is fit the models several times, always with sum-to-zero contrasts, but changing the level of the factor set to $-\Sigma$ rest of levels so as to explicitly obtain the coefficients and standard errors for all levels of all terms (e.g., the coefficients that correspond to “Detection, Uniform”, “Detection, Small” and “Detection, Large”).

5 Cancer data sets

5.1 Cancer data sets: sources and characteristics

Name	Source	Original source	Number of genes or pathways	Number of subjects
All Pathways	[15]	(From sources for colon, glioblastoma, and pancreas genes data sets)	12	268
Colon Genes	[15]	[24]	8	95
Colon Pathways	[15]	[24]	10	95
Glioblastoma Genes (Gliob Genes)	[15]	[21]	8	78
Glioblastoma Pathways (Gliob Pathways)	[15]	[21]	10	78
Lung	[19]	[11]	51	161
MSI	[5]	[4]	30	27
MSS	[5]	[4]	34	152
Ovarian	[19]	[2]	192	326
Pancreas Genes	[15]	[17]	7	90
Pancreas Pathways	[15]	[17]	7	90

These are further details about how the data were obtained:

All Pathways and Colon, Glioblastoma, Pancreas pathways The mapping from genes to pathways was done by [15], from the original papers with data sets. Our scripts to reproduce the analysis are provided with the code. Note that for Pancreas pathways we eliminate the four pathways that were present in all subjects (see also [15] and notes in the code for details). For Glioblastoma pathways, two pathways had identical patterns (Apoptosis and Small GTPase-dependent signaling (other than KRAS)) and only one was used.

What we call “All Pathways” here, for brevity, is called “All cancer types” in [15].

Lung Available as a text file from the supplementary material of [19] (file “BMLv1.tar.gz”) as file Lung_SM4.

Ovarian Available as a text file from the supplementary material of [19] (file “BMLv1.tar.gz”) as file OV_SM5.

MSI Colorectal cancer, microsatellite unstable tumors. From [5]. The original data (as well as MSS) come from COADREAD [4] and were splitted by tumour subtype MSI and MSS.

We used GIMP to open the pdf file page (Figure 3 on page 6 of [5]) where the figure was and cropped the grid of the figure and exported it as JPEG with high resolution. Then we imported it in ImageJ(Fiji) (<https://fiji.sc/>), converted it to 8-bit, applied threshold option and set it to B/W, then exported it as text image (matrix as txt). Then we imported the text image in R and used the code in `fig_to_matrix_capri_pnas.R` to convert the text image into a matrix of genotypes. The data were checked against the original figures.

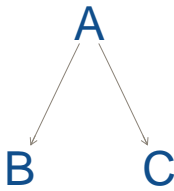
MSS Colorectal cancer, microsatellite stable tumors. From [5]. Same process as for MSI; the figure is Figure S5 from page 16 of the supplementary material to [5]. The authors explain that “Events selected for reconstruction are those involving genes altered in at least 5% of the cases, or part of group of alterations showing an exclusivity trend (see Figure S4).” There are 18 genes with a frequency $> 5\%$, and the largest number we analyzed was 13. The original data for MSI and MSS contain, for some genes, events for both deletion and mutation that affect the same gene (e.g., NRAS or ACVR1). Obviously, they both cannot happen in the same cell, and would constitute a violation of CPM assumptions. However, this does not affect the data we used because none of the deletion (or amplification) events were selected at the 13 events thresholds (and obviously neither at the 10 or 7). (The alternative would have been to remove some of the offending events or collapse into a single one the relevant columns, but since this becomes a non-issue with up to 13 events, it makes no difference.)

5.2 Bootstrapping on the cancer data sets

If the bootstrapping process resulted in a feature becoming absent from the data, or two or more features having identical patterns (i.e., one feature being identical to another) we discarded the bootstrap sample and obtained a new one; this is done to ensure that all bootstrapped data sets have paths of identical length (see also section 4.3). This, therefore, leads to JS values that are more optimistic (smaller).

6 CAPRI, CAPRESE, and paths of tumor progression

With both OT and CBN if we see a DAG such as



this is saying that, except for errors (errors in the model and observational errors) the genotypes that can exist under the model are only (with our usual notation of using a capital letter to denote that the gene is mutated, and no letter to denote absence of mutation)

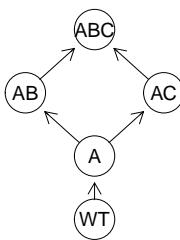
WT [genotype without driver mutations]

A
AB
AC
ABC

and, consequently, there are only two paths to the maximum:

WT → A → AB → ABC
WT → A → AC → ABC

or



Which means that, for example, under OT and CBN the following paths to the maximum are not allowed under the model:

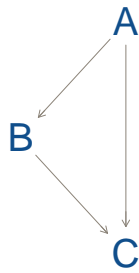
WT → B → BC → ABC
WT → B → AB → ABC

This interpretation, however, does not necessarily follow with CAPRI. CAPRI returns, as output, a DAG and a set of conditional probability tables (CPTs) associated to each node. What CAPRI seems to say, as can be checked from non-zero entries in the CPTs for B without A or C without A in a DAG as above, is that the DAG says that the most likely mutational paths are

WT → A → AB → ABC
WT → A → AC → ABC

but other mutational paths could also take place under the model. (This follows directly from seeing CAPRI return a CPT where, say, $P(B|\neg A) = 0.4$, i.e., mutating B without A is certainly not a rare event, even when an arrow exists from A to B). In fact, any path might happen (with some conditional probability tables), and one would seem to need to look at the CPTs to understand which ones can or cannot happen (but see below). And we are given no clear definition of what most likely paths really mean (e.g., “trends of selective advantage among genomic alterations” or “most common evolutionary trajectories” in the wording of 5) in terms of what probabilities are considered or not.

We see a similar issue with the following two DAGs:



because, of course, the transitive reduction of the first DAG is the second DAG. So from the point of view of paths from the non-mutated to the fully mutated genotype, both DAGs have the same meaning, as both imply that the same order of events needs to happen (or the same restrictions in the accumulation of mutations hold). Yet CAPRI seems to make a distinction between them. A distinction that could only be disentangled, presumably, by looking at the CPTs.

However, no information is available on how to go from CPTs to the conditional probabilities of genotypes implied by the model. In fact, directly using the CPT information seems discouraged and not something that users are supposed to do: access to the CPT has to be done via `TRONCO:::as.bnlearn.network(some_tronco_model)`, i.e., using the `:::` which denotes that we are accessing a non-exported function from the software. (This was the case with v. 2.11.0 and is the case as of July-2018 with version v. 2.13.0).

This contrasts with, say, CBN (and MCCBN) and OT which provide, respectively, estimated λ s and edge weights (`object$parent$est.weight`). These conditional probabilities are what we use to obtain the probability-weighted paths implied by the model.

This is a key difference between OT and CBN on the one hand and CAPRI (and CAPRESE) on the other: OT and CBN incorporate errors in their models, but they return the estimates of the parameters of their models, i.e., the λ s or the `est.weights`, that map directly into what can happen under the model, what genotypes can arise and from which other genotypes. (This is analogous to a simple linear regression: there is an error component in the model $y = \alpha + \beta x + \epsilon$, the ϵ , often assumed normally distributed with mean 0 and variance σ^2 , etc, but the method, and the software, return the α and β which allow us to predict the expected value of y given x , under the model).

In contrast, with CAPRI we can (again, using a non-exported function) obtain the CPTs that correspond to the DAG returned. Remember that essentially what CAPRI is doing is fitting a Bayesian Network to the observational data, with the DAG built so that arrows respect the temporal priority and probability raising restrictions. But the CPTs themselves are not estimated parameters of a model that could be mapped into probabilities of paths. The CPTs of CAPRI seem to be the conditional probabilities of observing what we observe under the DAG and they incorporate errors (model errors and noise in the data), and can contain non-zero entries for child nodes when their parents are absent. Obtaining probabilities of paths from these CPTs is, thus, impossible.

An additional issue that has been noted with CAPRI in the paper is its behavior with in-

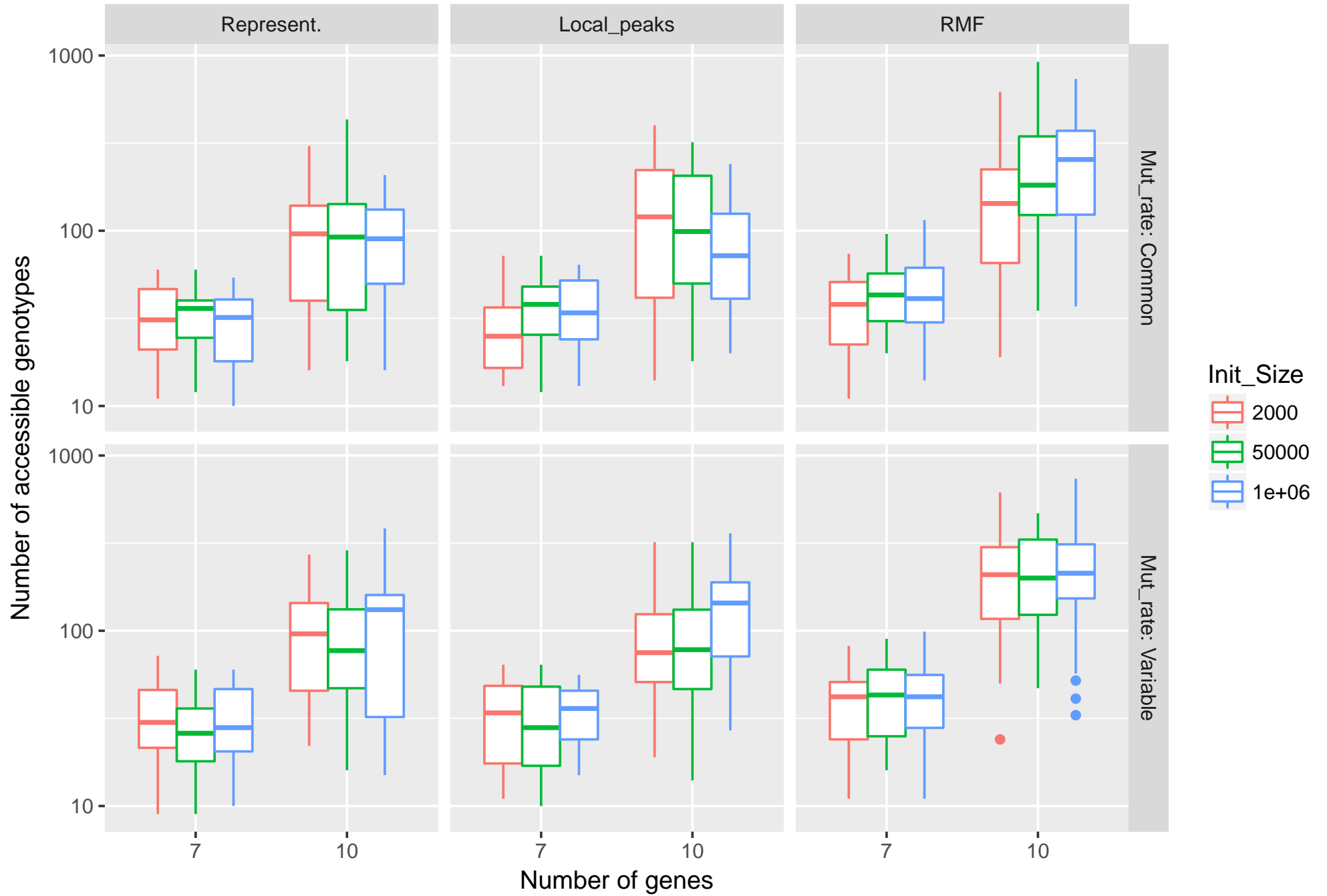
creasing sample sizes. That was also noted in [10] and seems related to the penalization in the fits: CAPRI does not seem suited to deal with large data sets as it tends to allow only one, or a few, paths to the maximum when N is 4000. The (transitive reduction of) the DAGs returned from CAPRI is often a linear sequence. This behavior did not change considerably whether we used AIC or BIC.

CAPRESE does not return DAGs, but trees, and in that sense it is simpler to deal with. But still, as with CAPRI, there is no information available on how to go from CPTs to the conditional probabilities of genotypes implied by the model (and accessing the CPT does not seem encouraged) and, as for CAPRI, the CPTs seem to mix the underlying model with the error model, and obtaining probabilities of paths from these CPTs is, thus, impossible.

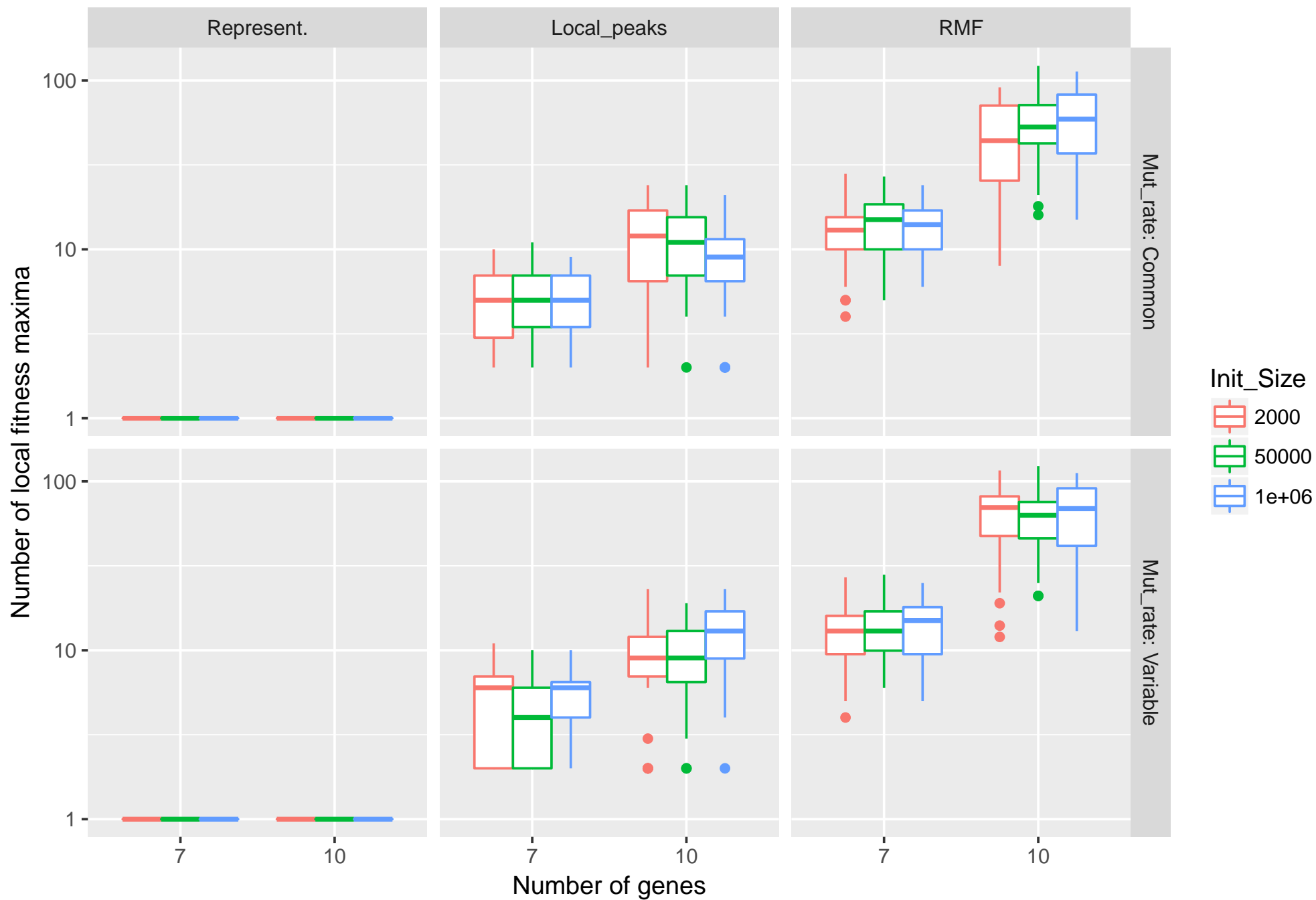
7 Fitness landscapes: characteristics, evolutionary predictability, clonal interference, and sampled genotypes

The following figures show the main fitness landscape characteristics and the resulting variation in evolutionary predictability, clonal interference, and sampling characteristics, between types of fitness landscapes and simulation conditions (initial population size and mutation rate). Note that the “static fitness landscape” characteristics do not depend on the simulations.

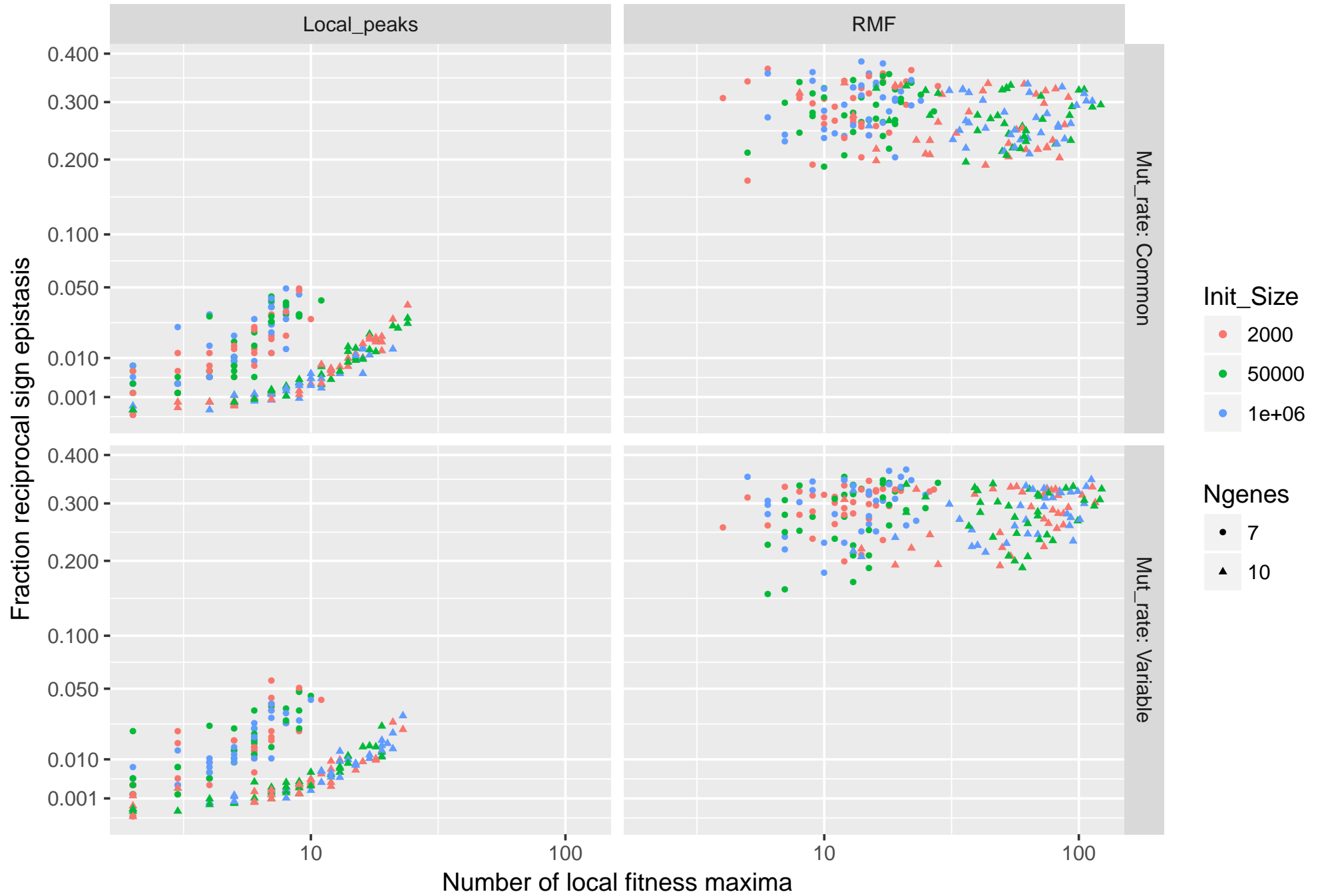
Static fitness landscape characteristics



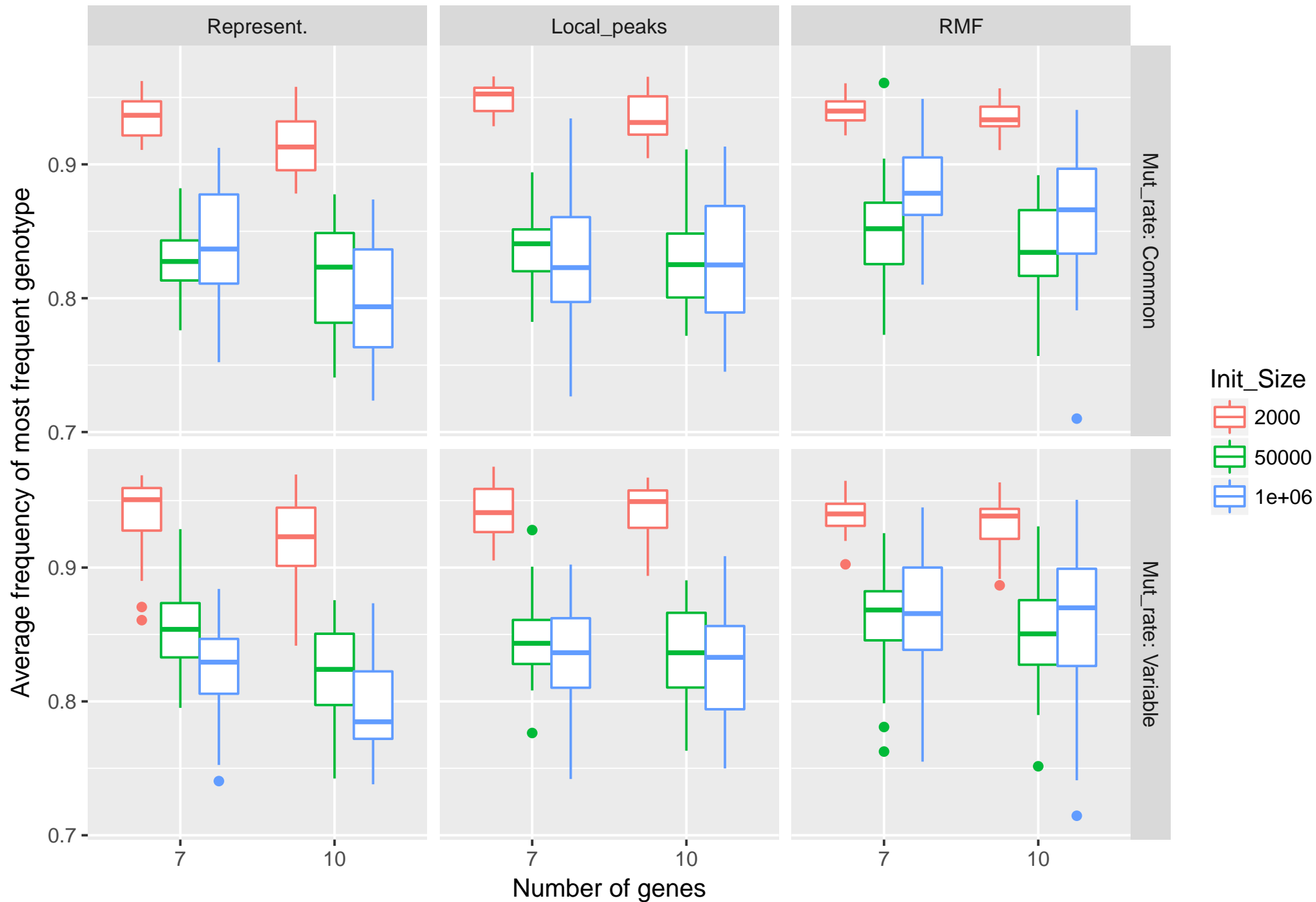
Static fitness landscape characteristics



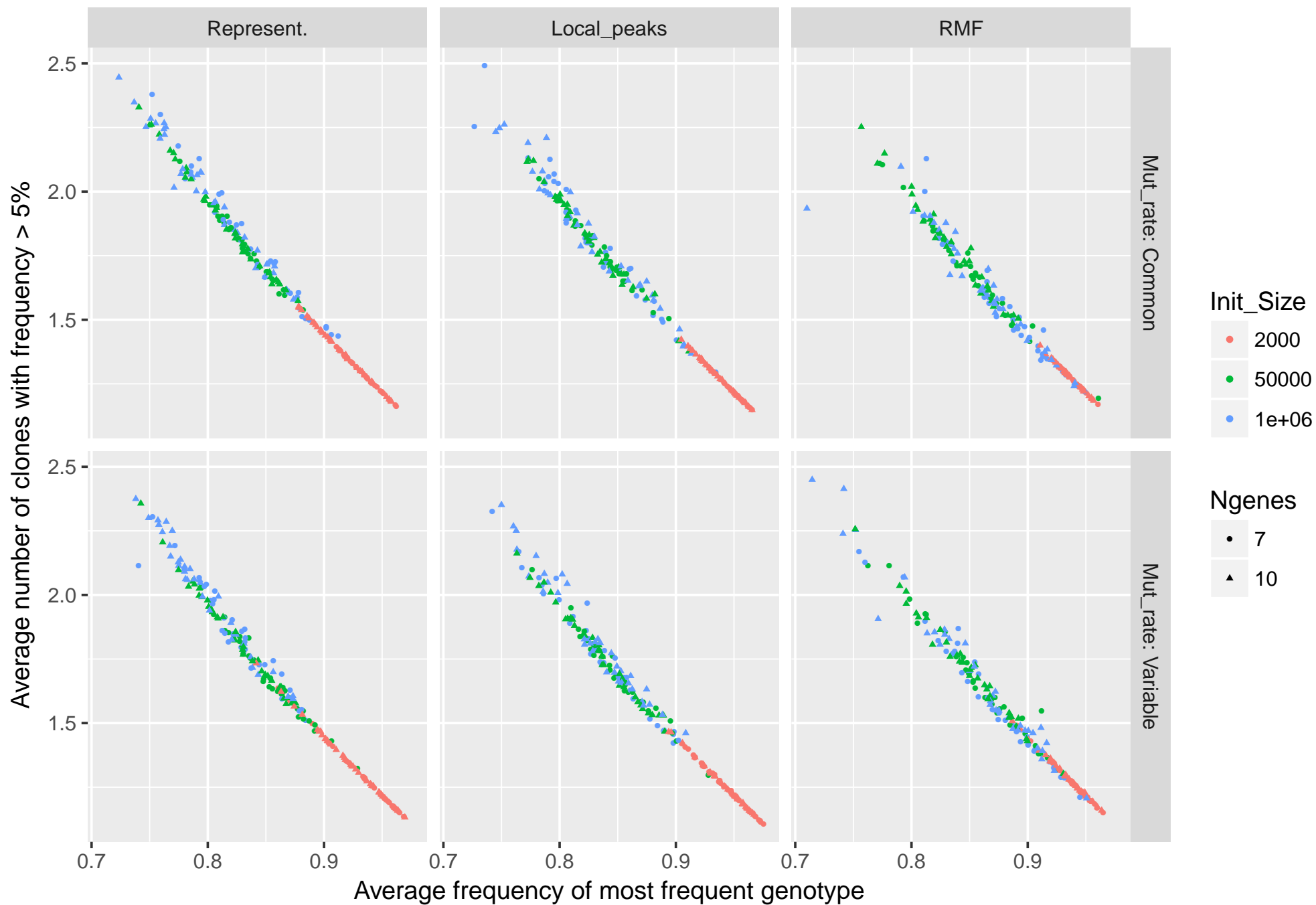
Static fitness landscape characteristics



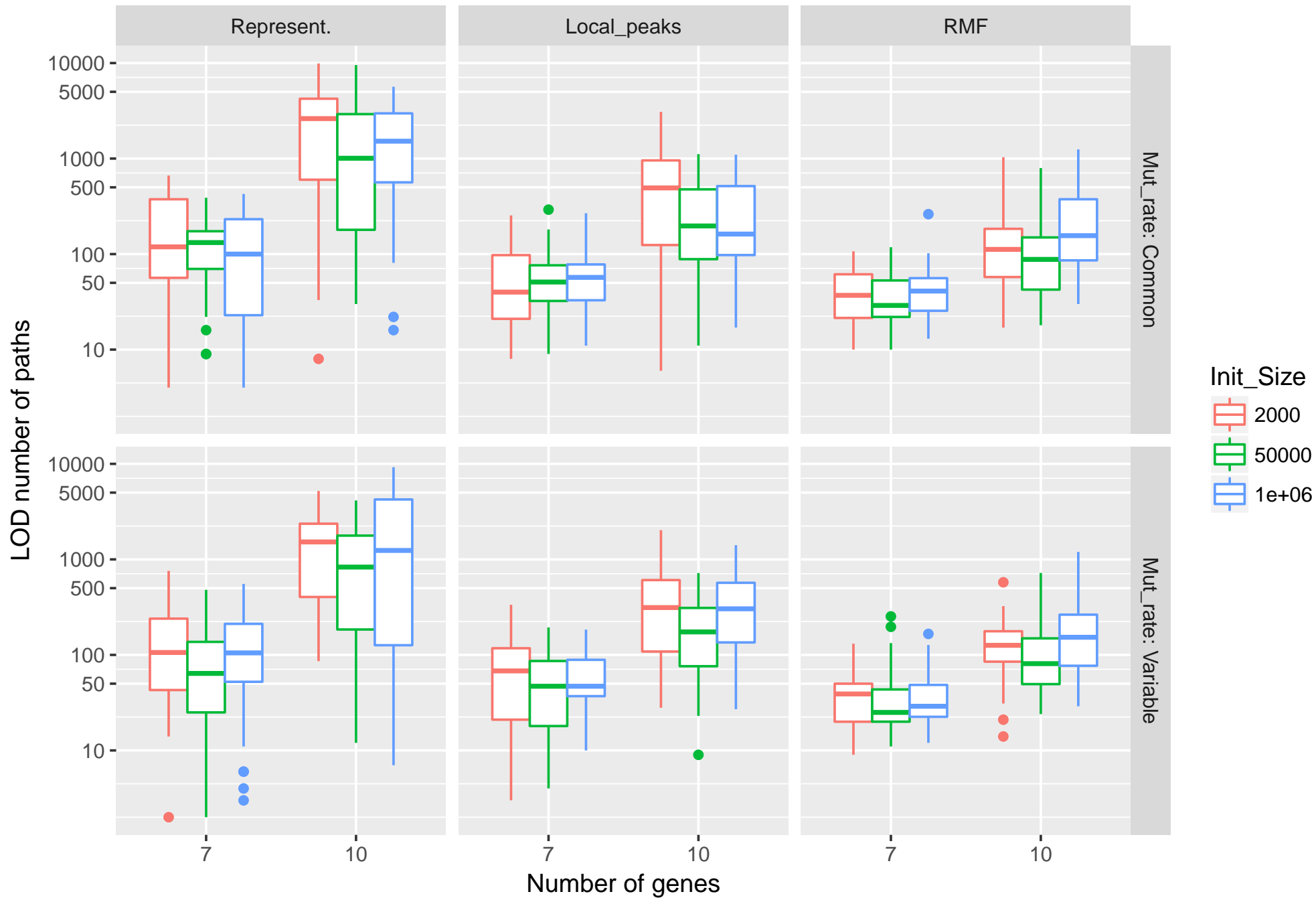
Clonal interference



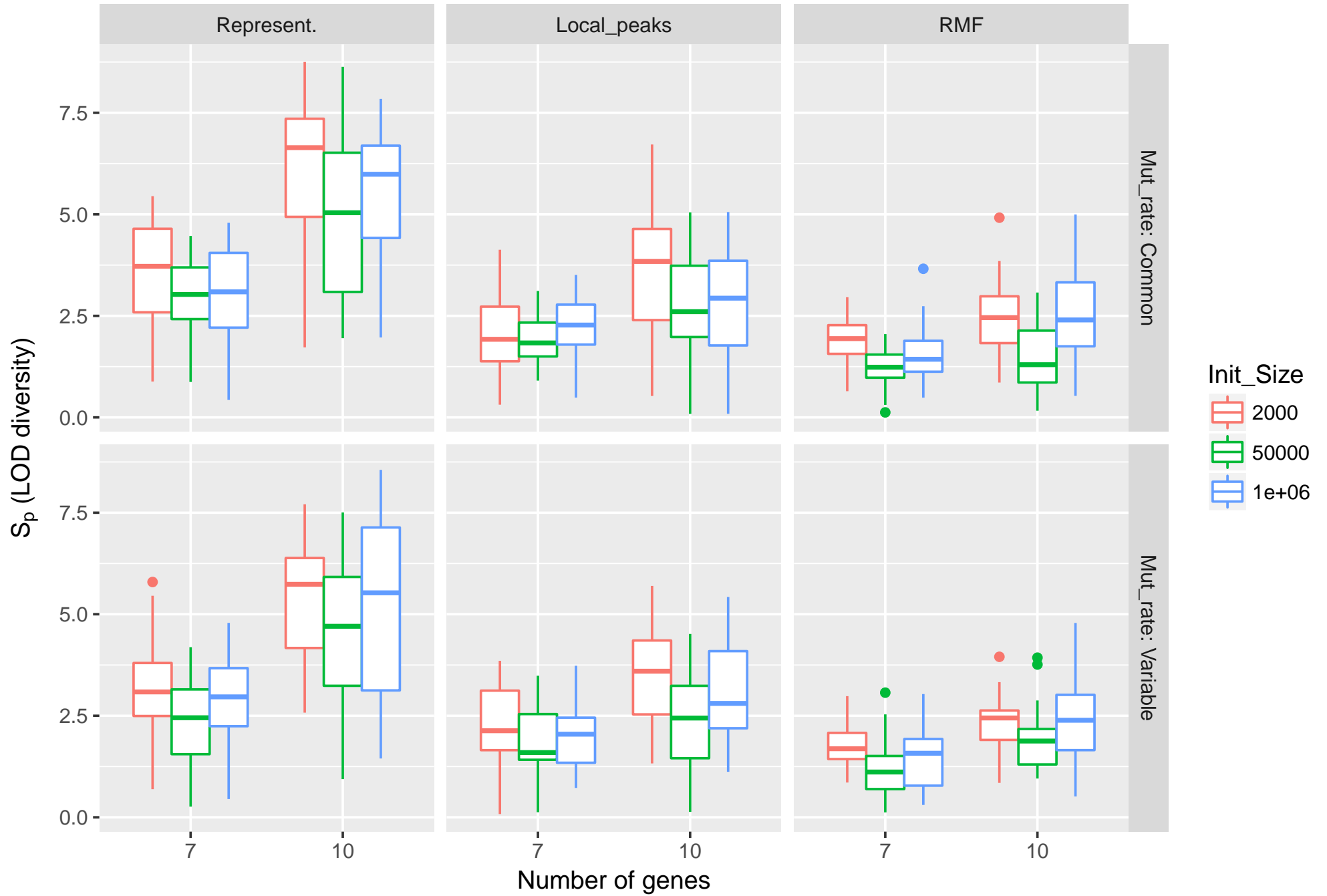
Clonal interference



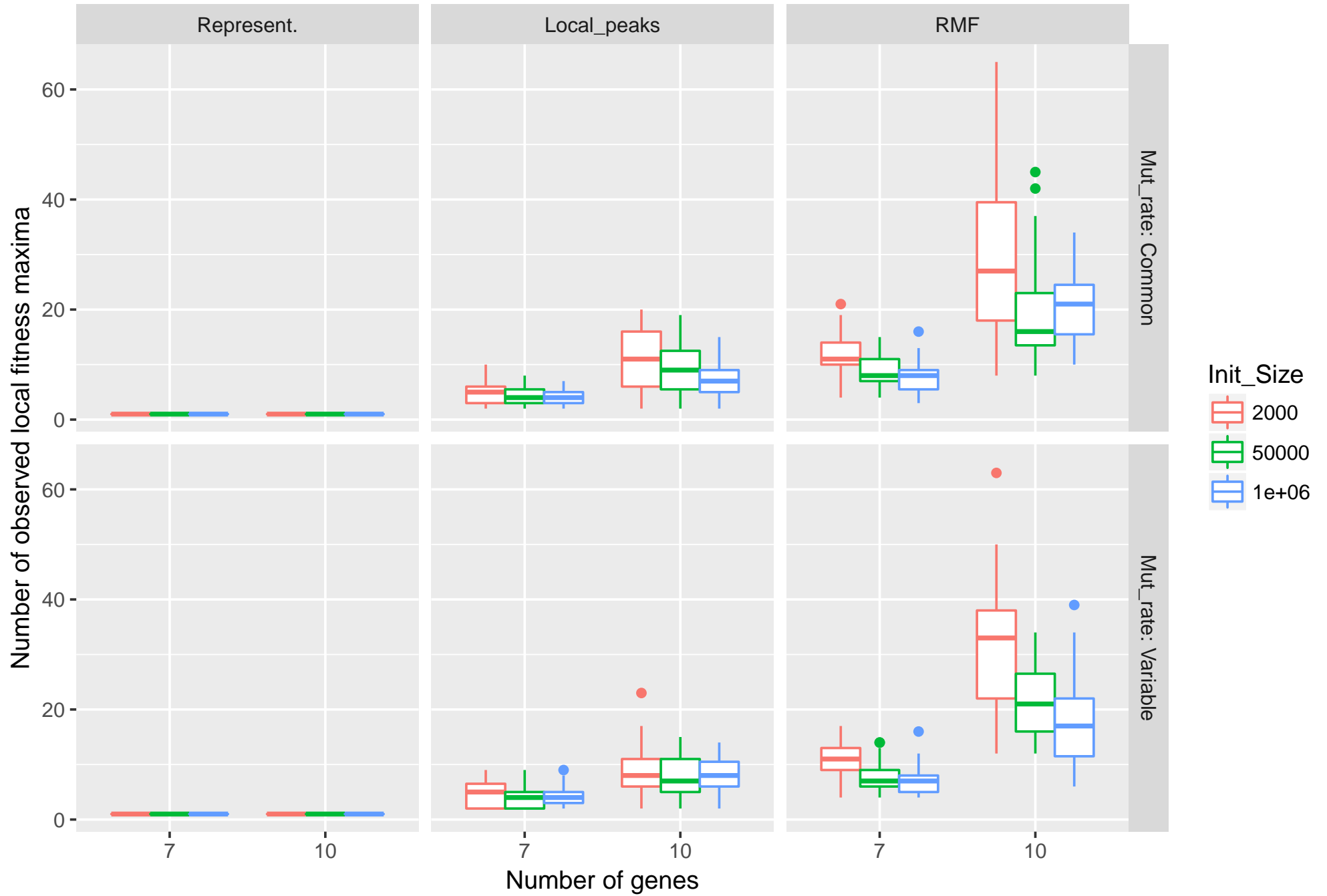
Evolutionary predictability of paths



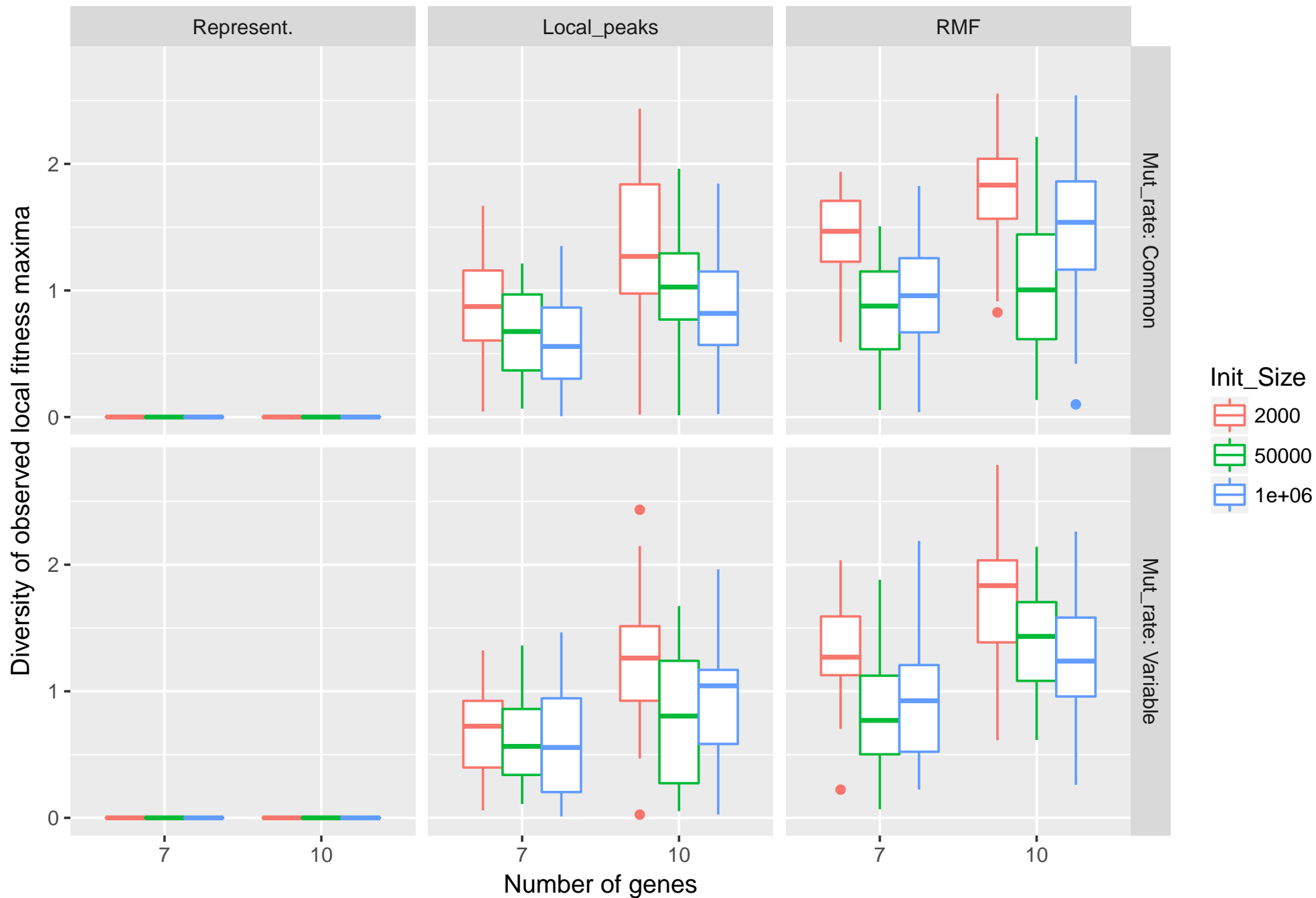
Evolutionary predictability of paths



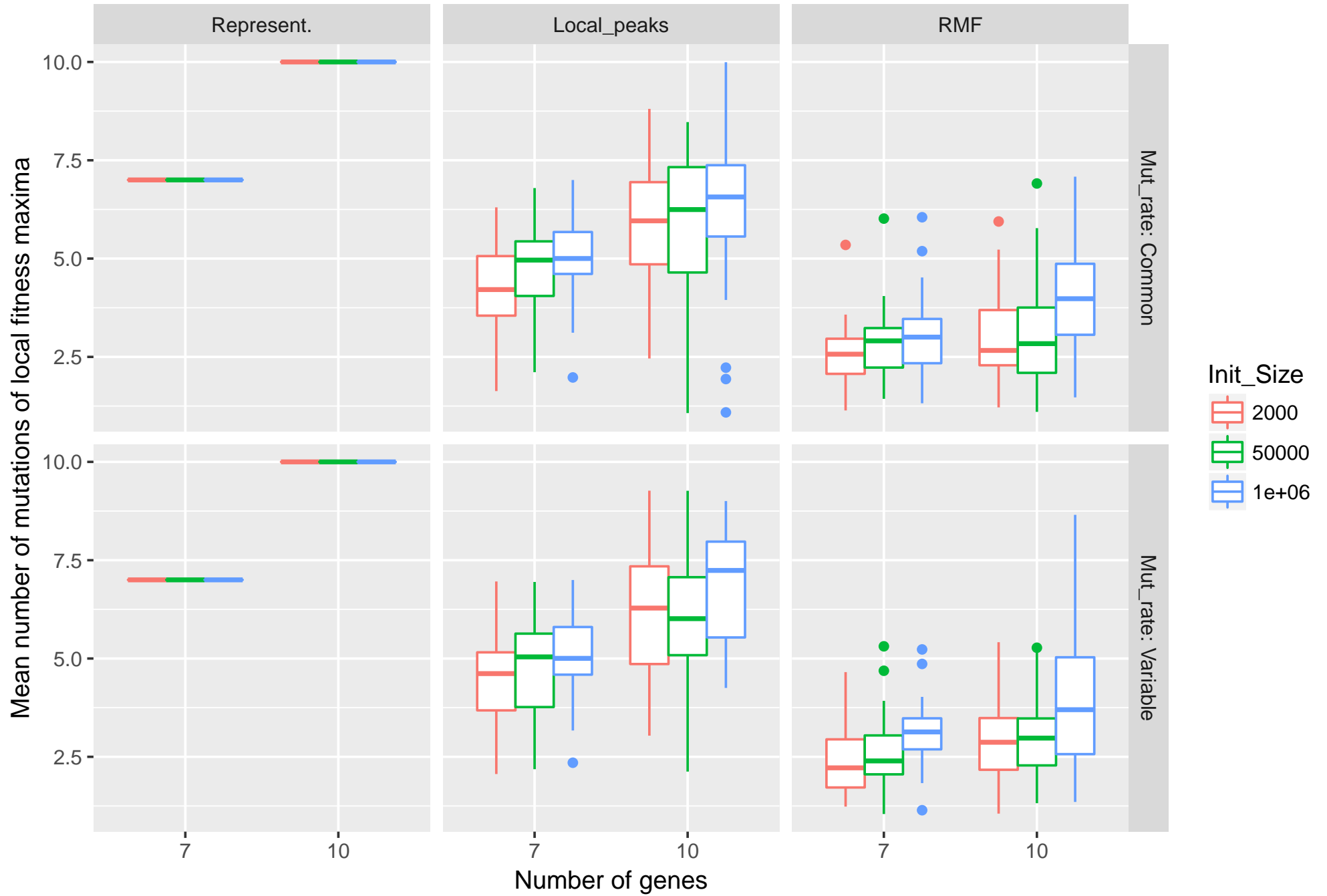
Evolutionary predictability of paths



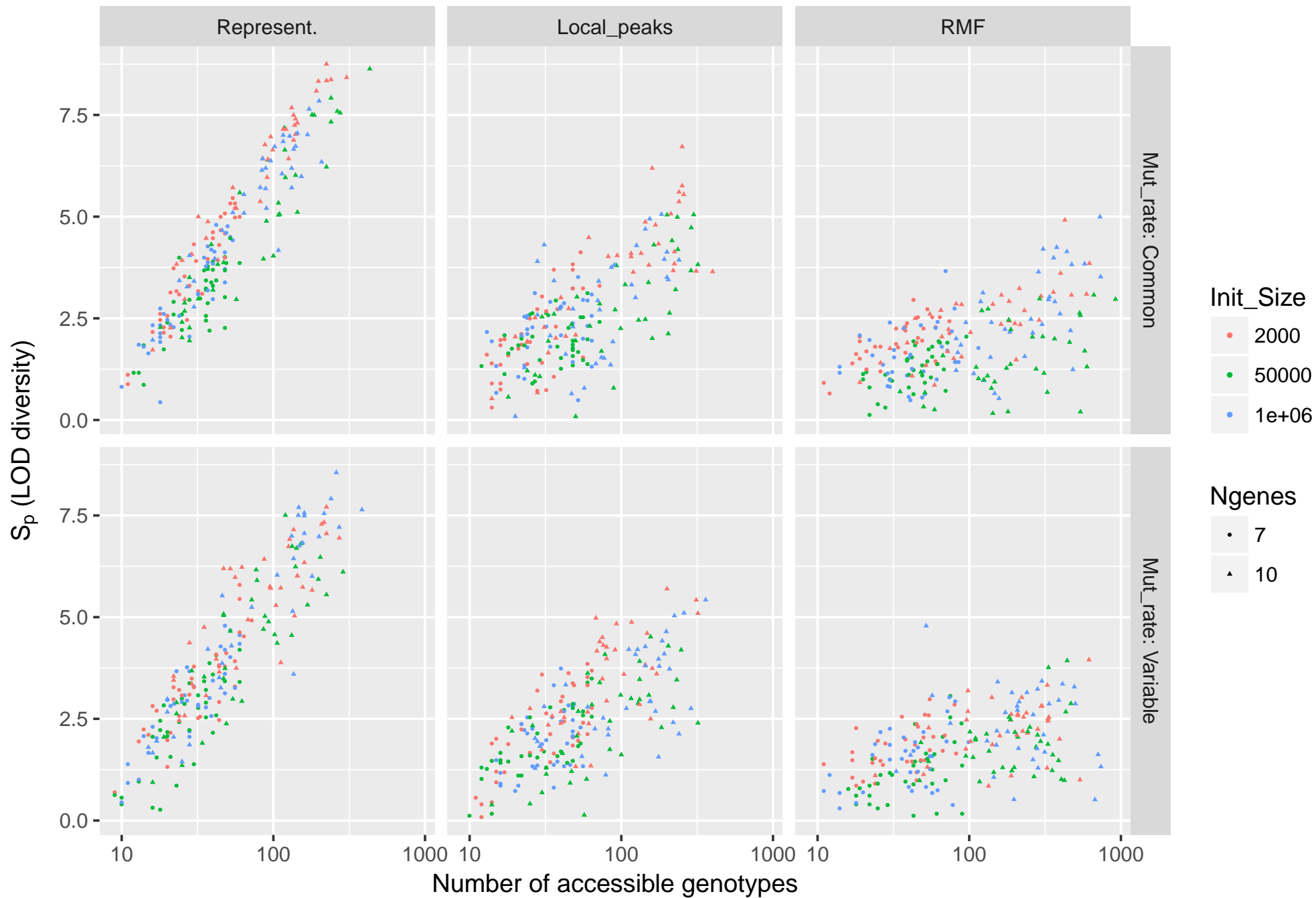
Evolutionary predictability of paths



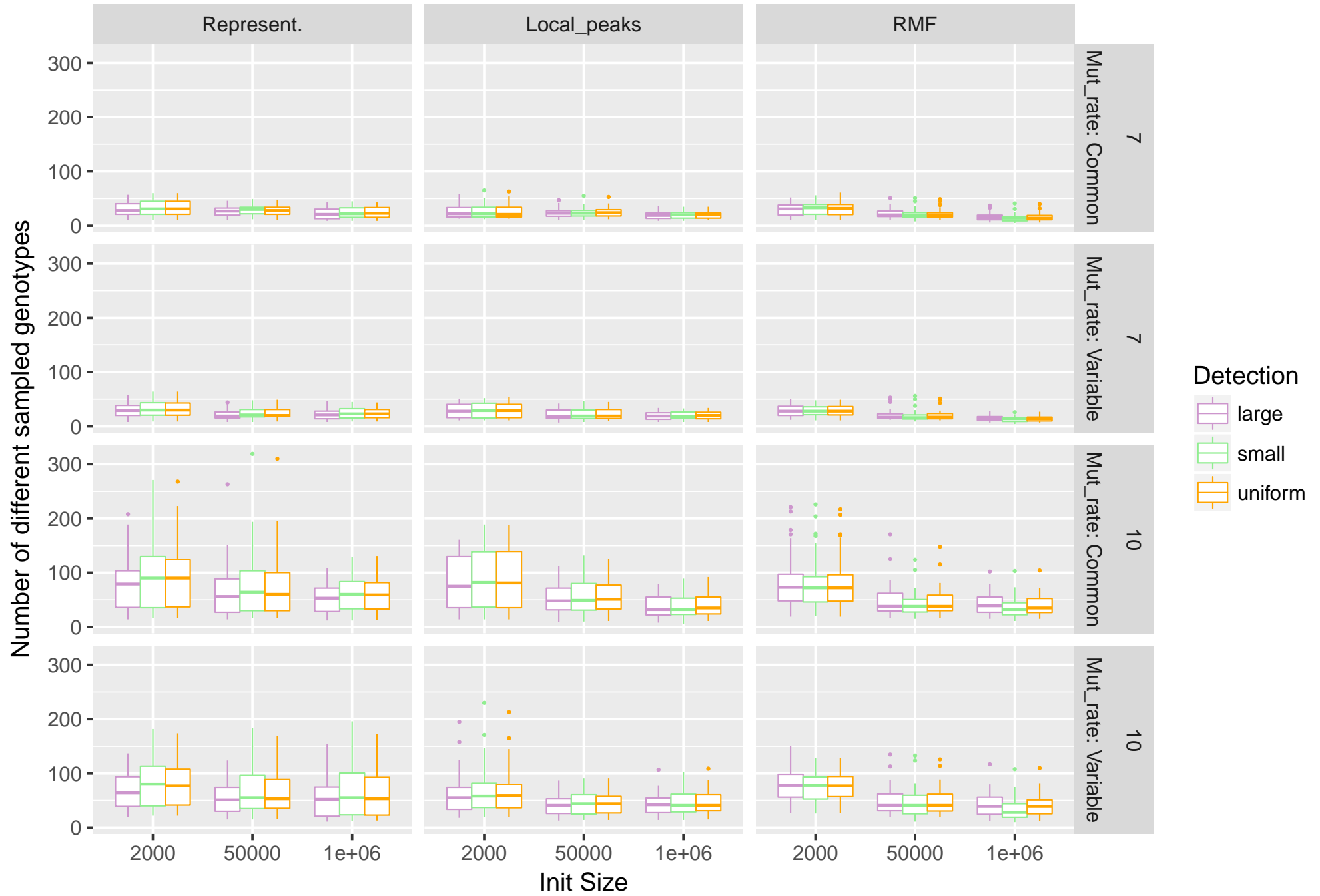
Evolutionary predictability of paths



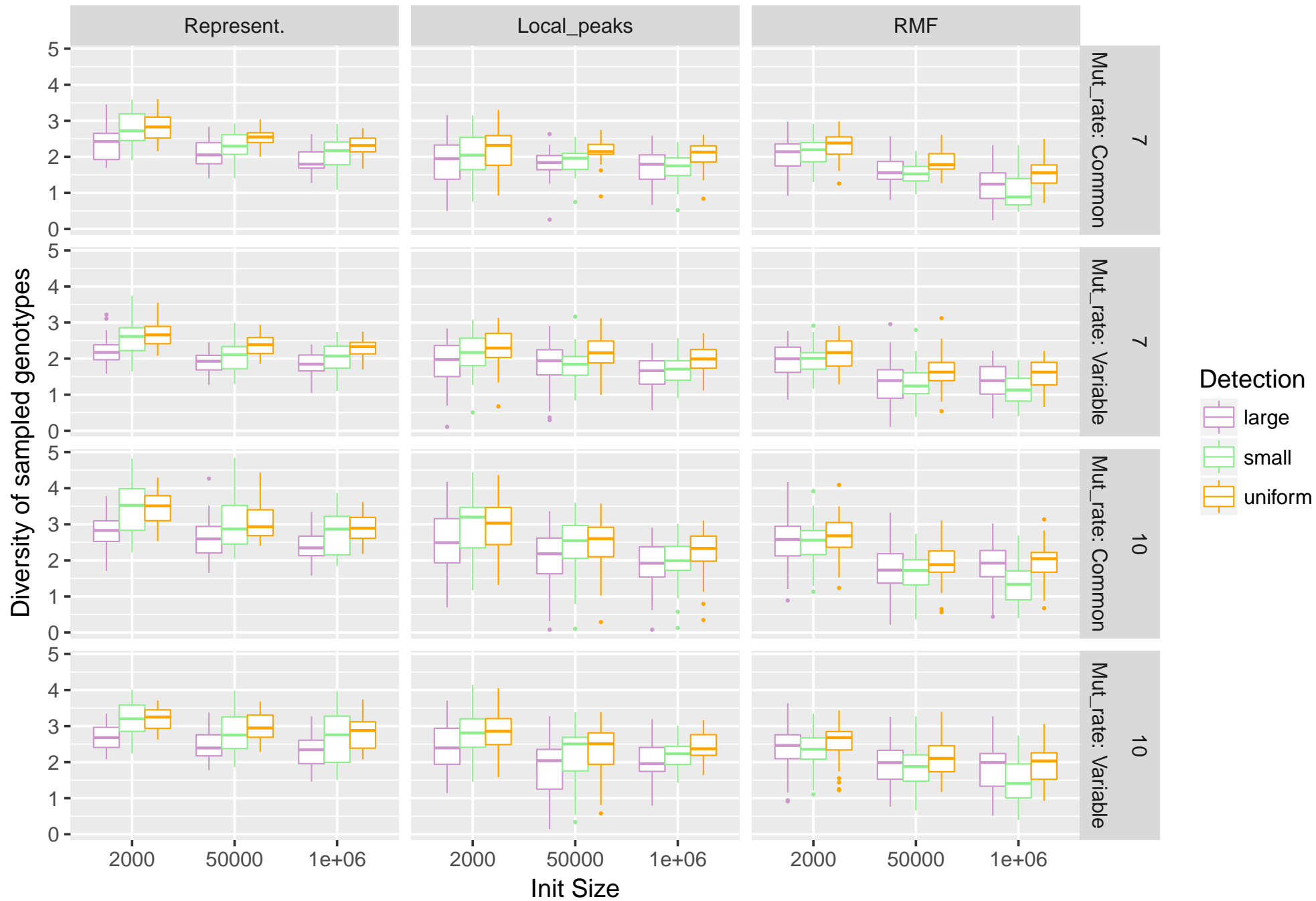
Evolutionary predictability of paths



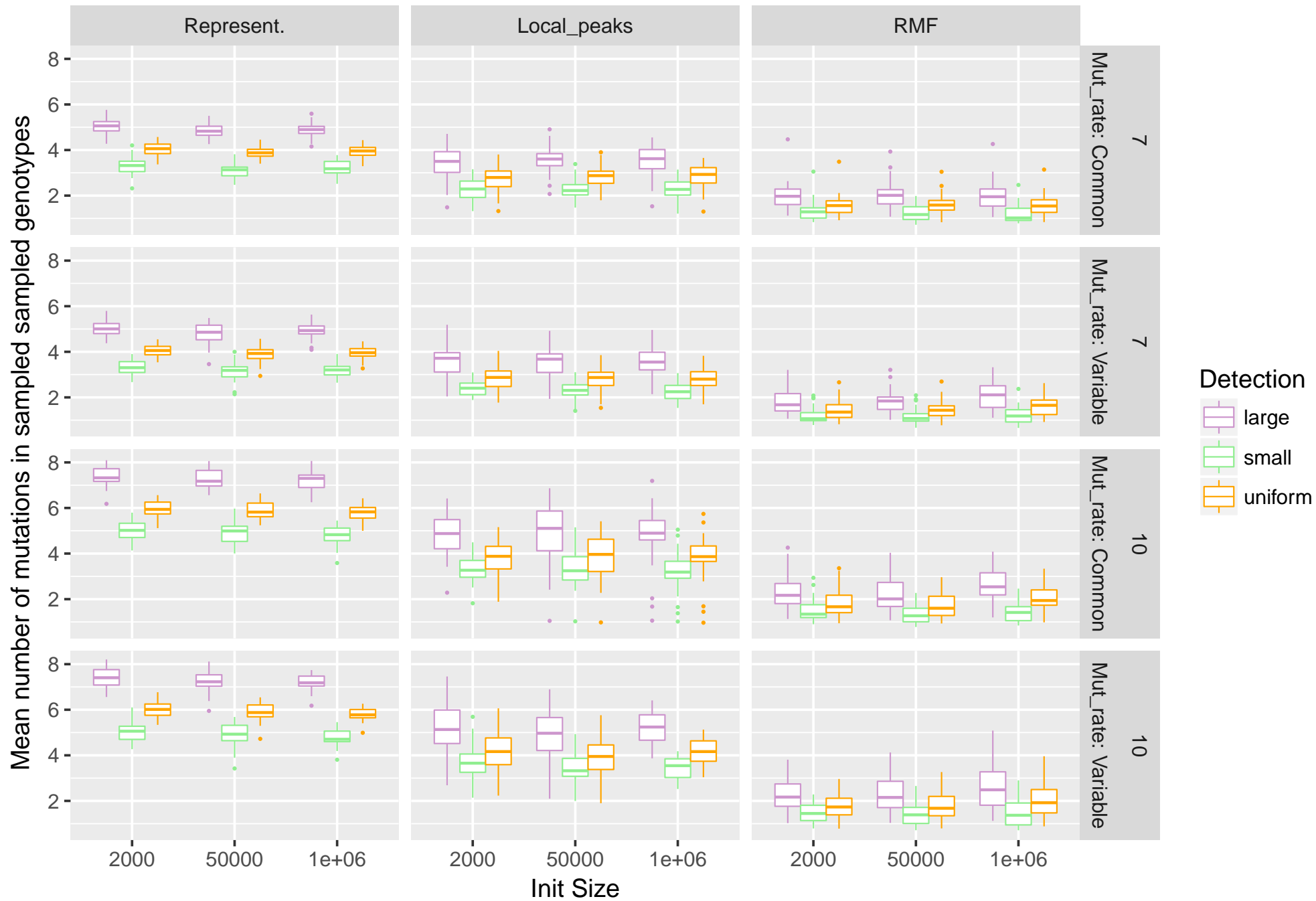
Sample's characteristics.



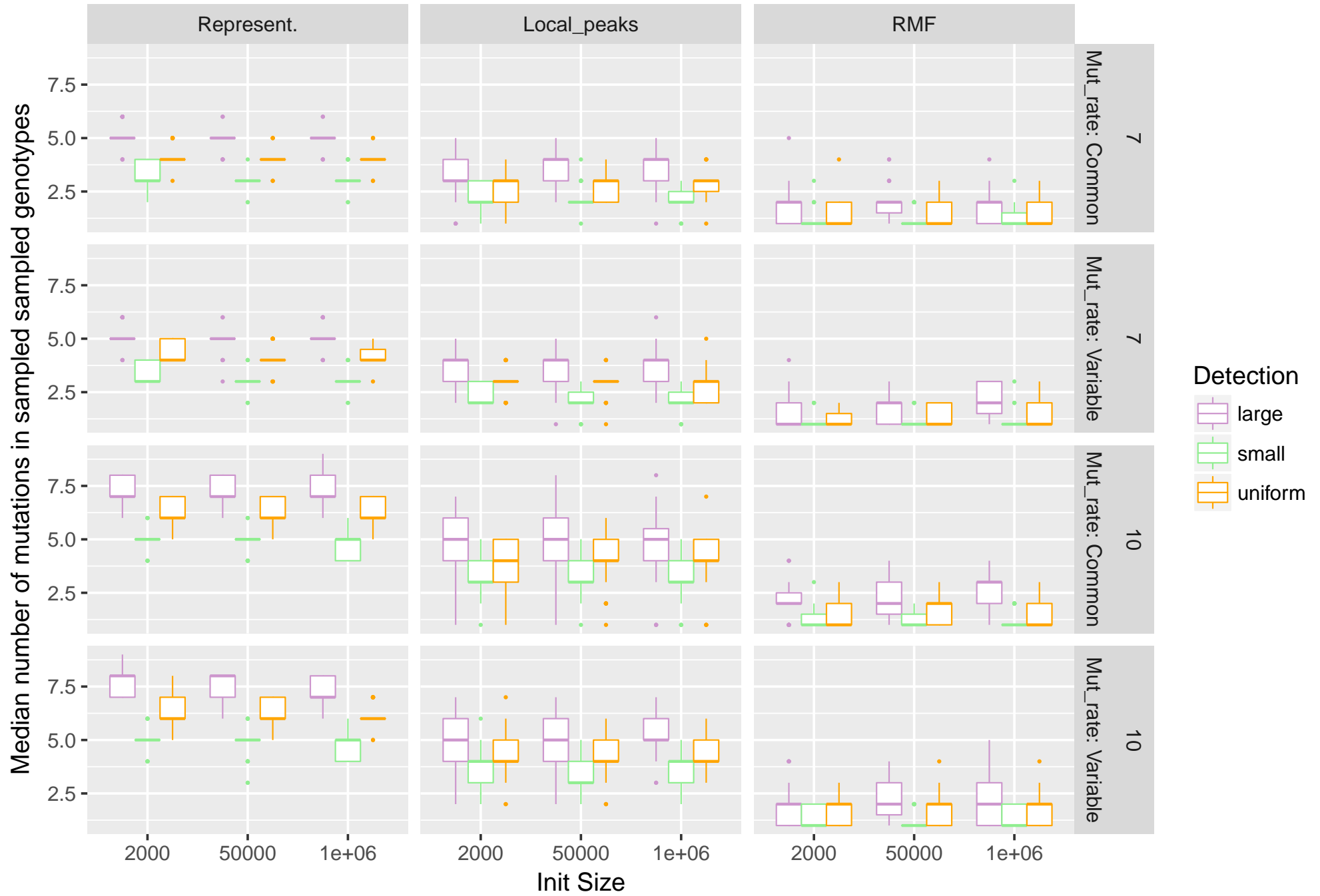
Sample's characteristics.



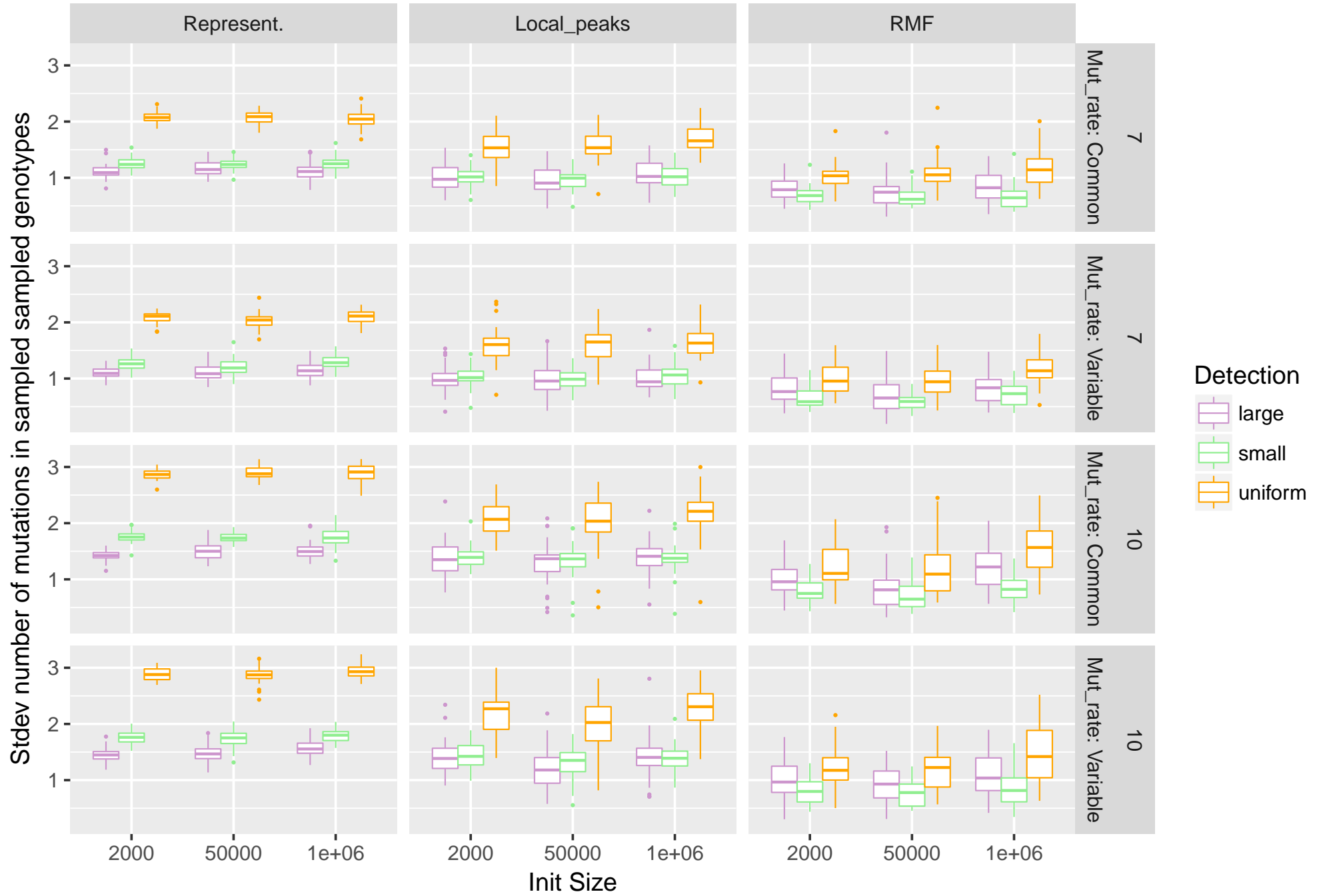
Sample's characteristics.



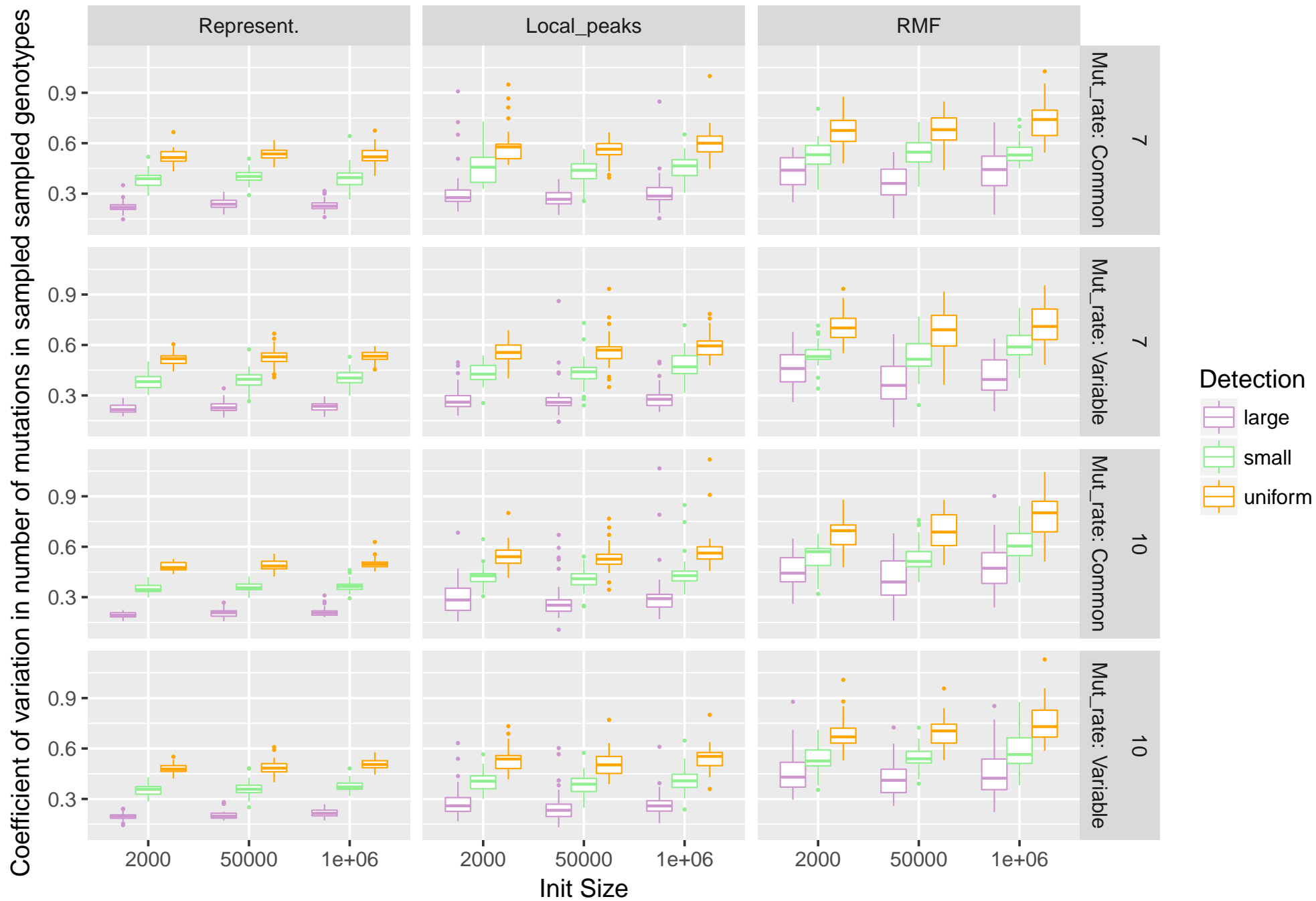
Sample's characteristics.



Sample's characteristics.



Sample's characteristics.



8 Overall patterns for the six methods

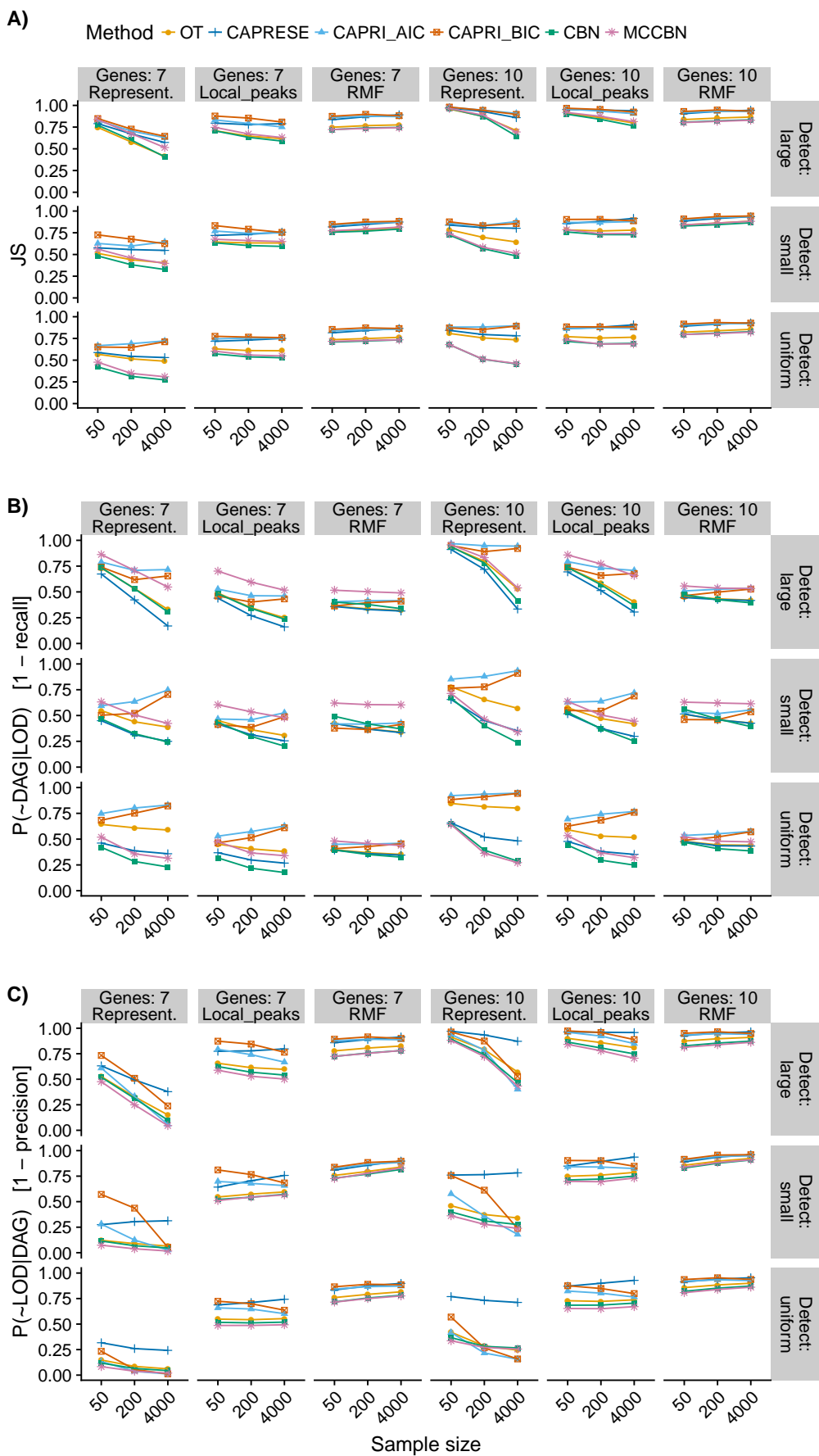


Figure 1: Summary performance measures for all six methods for all combinations of sample size by type of landscape by detection regime by number of genes. For all measures, smaller is better. For OT, CBN, and MCCBN, Jensen-Shannon entropy and 1-precision use probability-weighted predicted paths (see text). Each point represented is the average of 210 points (35 replicates of each one of the six combinations of 3 initial size by 2 mutation rate regimes; we are thus marginalizing over initial size by mutation rate; each one of the 210 points is, itself, the average of five runs on different partitions of the simulated data).

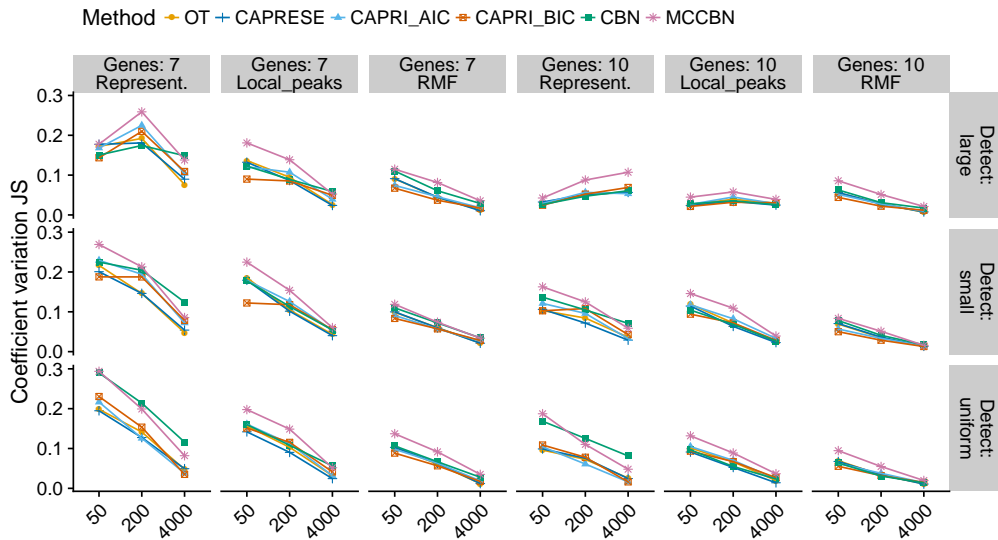


Figure 2: Coefficient of variation (standard deviation/mean) of JS for each combination of method and type of fitness landscape. The coefficient of variation has been computed from the five runs for each landscapes on each combination of sample size and detection regime.

9 OT and CBN, JS, weighted vs. unweighted

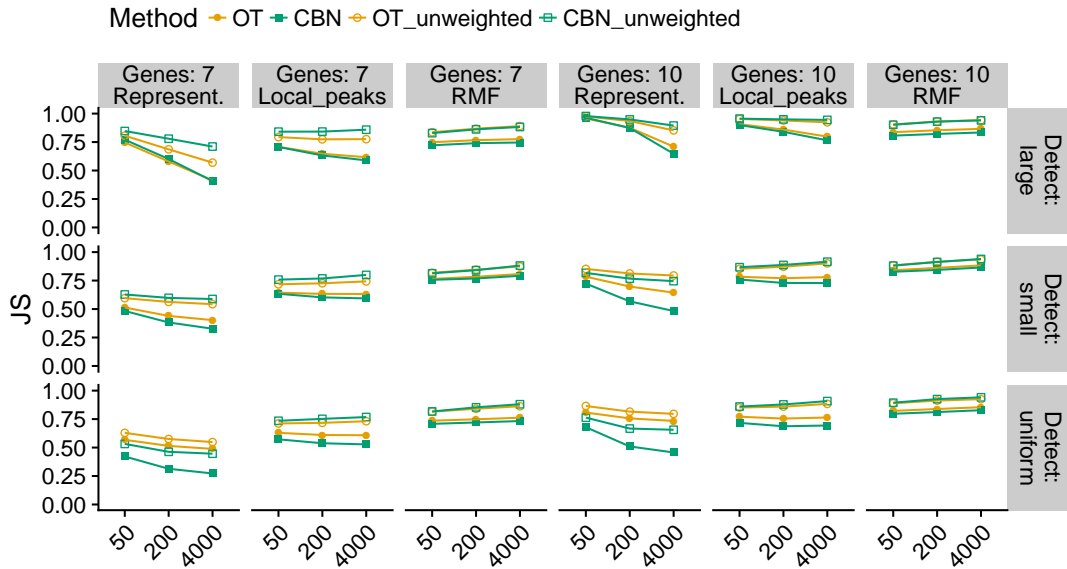


Figure 3: Comparison of the performance of OT and CBN using weighted and unweighted probabilities of paths to the maximum.

10 CAPRI and CBN, 1-precision, unweighted

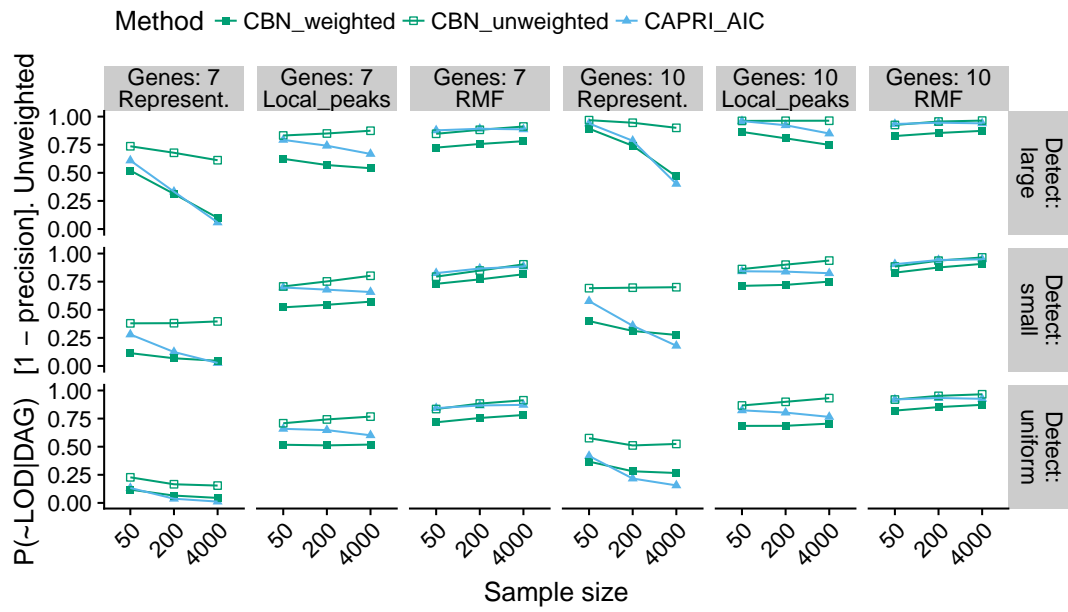


Figure 4: Comparison of the performance of CAPRI with CBN using weighted and unweighted probabilities of paths to the maximum.

11 CAPRESE and OT, 1-precision, unweighted

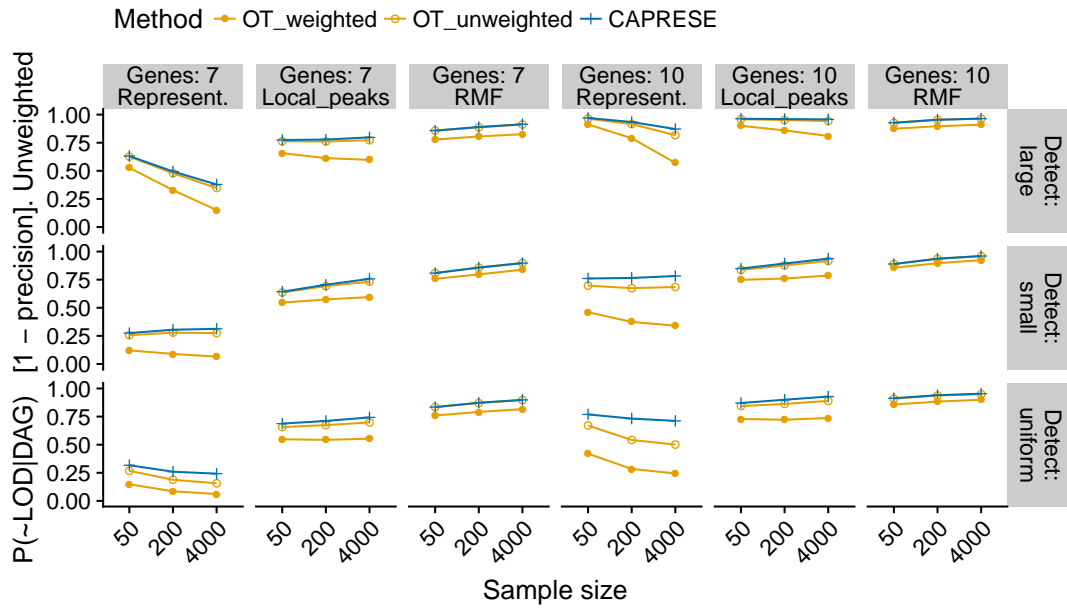


Figure 5: Comparison of the performance of CAPRESE with OT using weighted and unweighted probabilities of paths to the maximum.

The results of OT here are remarkable because OT can only build trees, and therefore cannot reflect the dependency of a mutation on two or more upstream mutations so it is prone to allow more paths to the maximum. The results of OT contrasts with those of CAPRESE, the other method that only builds trees. CAPRESE is building DAGs of restrictions that have too few restrictions and, therefore, allow for too many paths to the maximum. One notable difference between the two methods is that with OT it is relatively simple to use a measure of 1-precision that weights by the probability of each path. The performance of OT, even if we use unweighted probabilities of paths, is much better than that of CAPRESE but improves even further when using weighted paths, again highlighting the usefulness of weighting paths to obtain more accurate predictions.

12 Probability of recovering the most common LOD

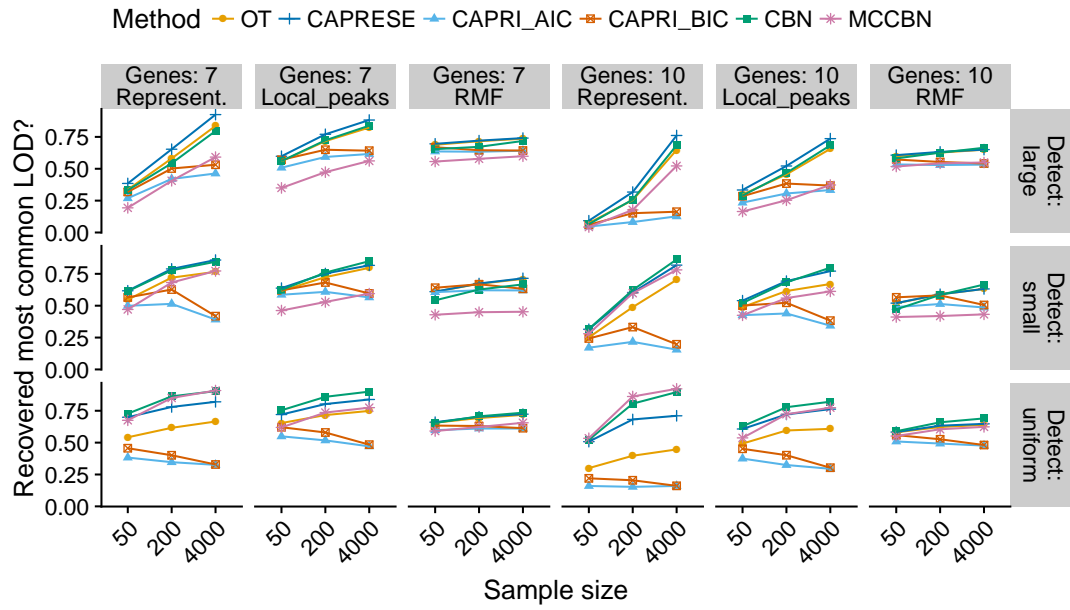


Figure 6: Probability of recovering the most common LOD: probability that the most common observed path to the maximum is among the paths allowed by the CPMs.

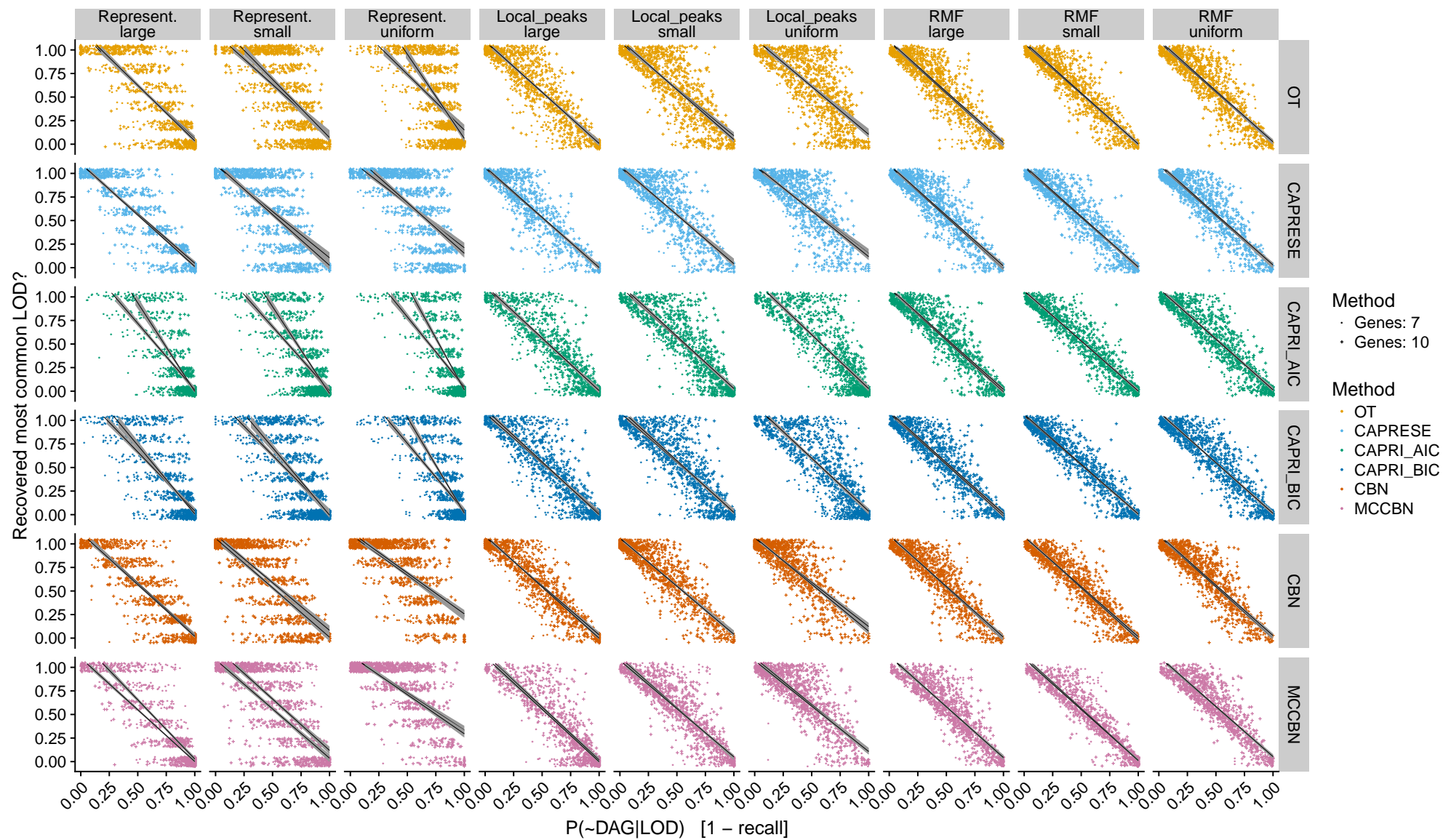


Figure 7: Probability of recovering the most common LOD and 1-recall: relationship.

13 Number of paths inferred

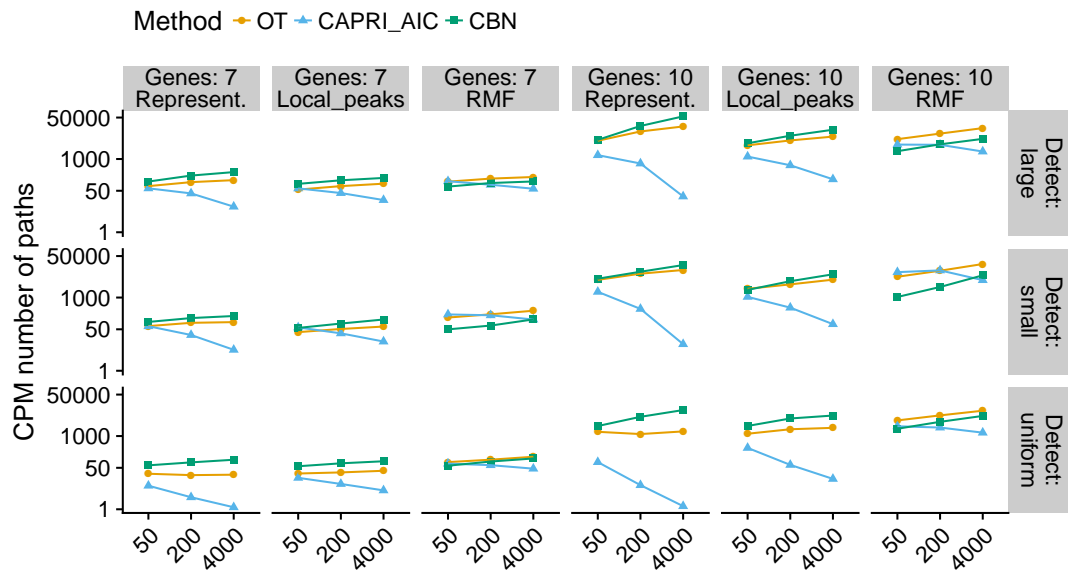


Figure 8: Number of paths to the maximum according to the CPMs.

14 Slopes of regressions of 1-recall and 1-precision on LOD diversity, S_p

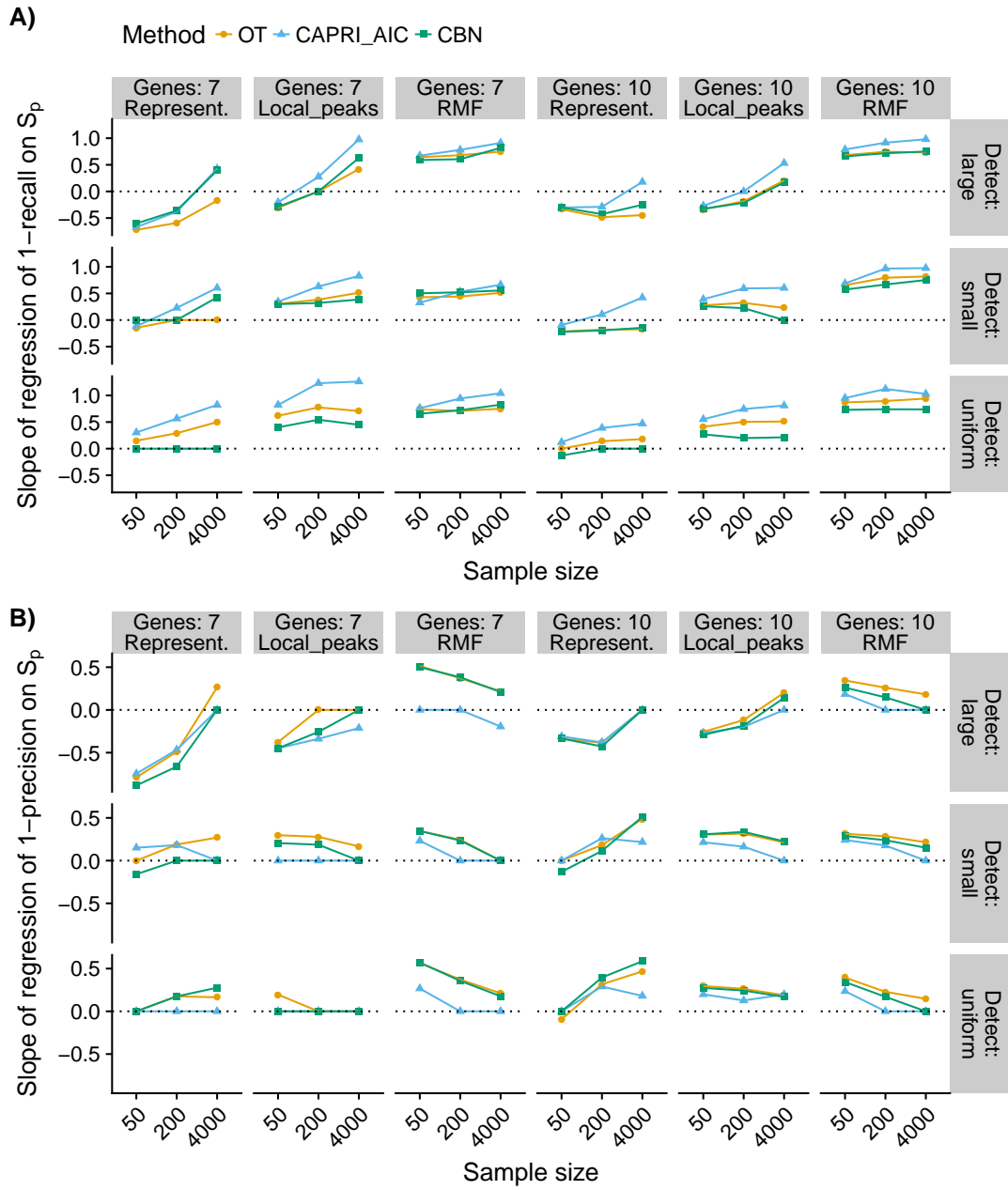


Figure 9: Slopes of regressions of 1-recall and 1-precision on LOD diversity, S_p

15 Coefficient of variation of S_c

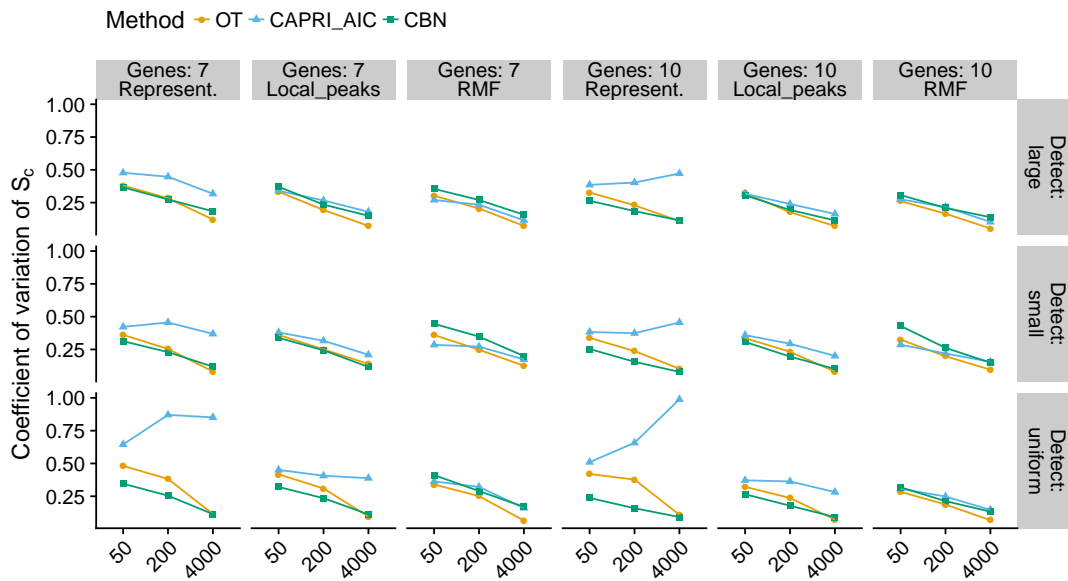


Figure 10: Coefficient of variation (standard deviation/mean) of S_c for each combination of method and type of fitness landscape. The coefficient of variation has been computed from the five runs for each landscapes on each combination of sample size and detection regime. For OT and CBN, it is computed using the probability-weighted predicted paths (see text). Each point plotted is the average of 210 points.

16 Estimated S_c by CBN

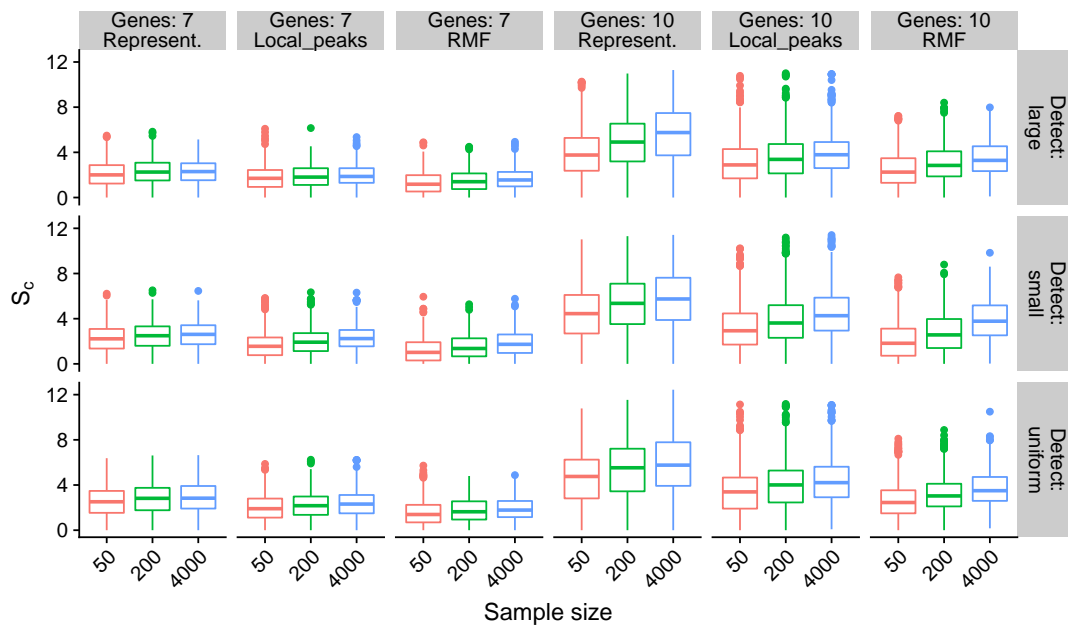


Figure 11: Estimated S_c by CBN for all combinations of sample size by type of landscape by detection regime by number of genes. Each box plot shows 1050 points.

17 Analysis of deviance tables for fitted models

Analysis of deviance tables for the generalized (beta regressions) linear mixed effects models. Models where fitted using the R package glmmTMB [3]. Analysis of deviance tables from package car [14]. All analysis of deviance tables use Type II Wald chi-square tests.

17.1 Models fitted to the complete data set

(Notice the strong evidence we need the three and four and even five way interactions.)

17.1.1 Two-way interactions

	Chisq	Df	Pr(>Chisq)
Num. Genes	223.34	1.00	< .0001
Method	9459.99	2.00	< .0001
Landscape	322.37	2.00	< .0001
Detection	4768.56	2.00	< .0001
Sample Size	1599.17	2.00	< .0001
LOD diversity	139.98	1.00	< .0001
Num. Genes:Method	22.98	2.00	< .0001
Num. Genes:Landscape	91.35	2.00	< .0001
Num. Genes:Detection	2318.60	2.00	< .0001
Num. Genes:Sample Size	381.80	2.00	< .0001
Num. Genes:LOD diversity	8.72	1.00	0.0032
Method:Landscape	192.88	4.00	< .0001
Method:Detection	678.41	4.00	< .0001
Method:Sample Size	818.96	4.00	< .0001
Method:LOD diversity	170.20	2.00	< .0001
Landscape:Detection	6534.14	4.00	< .0001
Landscape:Sample Size	3110.94	4.00	< .0001
Landscape:LOD diversity	78.21	2.00	< .0001
Detection:Sample Size	2420.56	4.00	< .0001
Detection:LOD diversity	3303.42	2.00	< .0001
Sample Size:LOD diversity	939.31	2.00	< .0001

Table 1: Full model, 2-way interactions

17.1.2 Three-way interactions

	Chisq	Df	Pr(>Chisq)
Num. Genes	197.19	1.00	< .0001
Method	11331.58	2.00	< .0001
Landscape	411.00	2.00	< .0001
Detection	4217.33	2.00	< .0001
Sample Size	1536.56	2.00	< .0001
LOD diversity	146.46	1.00	< .0001
Num. Genes:Method	52.74	2.00	< .0001
Num. Genes:Landscape	80.98	2.00	< .0001
Num. Genes:Detection	2394.91	2.00	< .0001
Num. Genes:Sample Size	354.72	2.00	< .0001
Num. Genes:LOD diversity	8.15	1.00	0.0043
Method:Landscape	285.47	4.00	< .0001
Method:Detection	707.50	4.00	< .0001
Method:Sample Size	842.82	4.00	< .0001
Method:LOD diversity	117.61	2.00	< .0001
Landscape:Detection	6924.74	4.00	< .0001
Landscape:Sample Size	3407.98	4.00	< .0001
Landscape:LOD diversity	77.13	2.00	< .0001
Detection:Sample Size	2507.80	4.00	< .0001
Detection:LOD diversity	4229.09	2.00	< .0001
Sample Size:LOD diversity	1120.46	2.00	< .0001
Num. Genes:Method:Landscape	90.47	4.00	< .0001
Num. Genes:Method:Detection	36.64	4.00	< .0001
Num. Genes:Method:Sample Size	10.43	4.00	0.0338
Num. Genes:Method:LOD diversity	41.35	2.00	< .0001
Num. Genes:Landscape:Detection	1302.06	4.00	< .0001
Num. Genes:Landscape:Sample Size	347.91	4.00	< .0001
Num. Genes:Landscape:LOD diversity	3.70	2.00	0.1572
Num. Genes:Detection:Sample Size	553.21	4.00	< .0001
Num. Genes:Detection:LOD diversity	3.91	2.00	0.1417
Num. Genes:Sample Size:LOD diversity	56.88	2.00	< .0001
Method:Landscape:Detection	250.27	8.00	< .0001
Method:Landscape:Sample Size	192.75	8.00	< .0001
Method:Landscape:LOD diversity	605.91	4.00	< .0001
Method:Detection:Sample Size	19.44	8.00	0.0127
Method:Detection:LOD diversity	126.85	4.00	< .0001
Method:Sample Size:LOD diversity	117.71	4.00	< .0001
Landscape:Detection:Sample Size	2084.94	8.00	< .0001
Landscape:Detection:LOD diversity	867.54	4.00	< .0001
Landscape:Sample Size:LOD diversity	1163.94	4.00	< .0001
Detection:Sample Size:LOD diversity	736.44	4.00	< .0001

Table 2: Full model, 3-way interactions

17.1.3 Four-way interactions

	Chisq	Df	Pr(>Chisq)
Num. Genes	193.68	1.00	< .0001
Method	11619.48	2.00	< .0001
Landscape	426.49	2.00	< .0001
Detection	4059.78	2.00	< .0001
Sample Size	1500.59	2.00	< .0001
LOD diversity	150.27	1.00	< .0001
Num. Genes:Method	54.25	2.00	< .0001
Num. Genes:Landscape	77.19	2.00	< .0001
Num. Genes:Detection	2370.70	2.00	< .0001
Num. Genes:Sample Size	357.14	2.00	< .0001
Num. Genes:LOD diversity	8.47	1.00	0.0036
Method:Landscape	298.90	4.00	< .0001
Method:Detection	716.32	4.00	< .0001
Method:Sample Size	841.73	4.00	< .0001
Method:LOD diversity	119.30	2.00	< .0001
Landscape:Detection	6697.76	4.00	< .0001
Landscape:Sample Size	3459.44	4.00	< .0001
Landscape:LOD diversity	77.85	2.00	< .0001
Detection:Sample Size	2466.93	4.00	< .0001
Detection:LOD diversity	4132.54	2.00	< .0001
Sample Size:LOD diversity	1103.02	2.00	< .0001
Num. Genes:Method:Landscape	101.44	4.00	< .0001
Num. Genes:Method:Detection	35.22	4.00	< .0001
Num. Genes:Method:Sample Size	11.99	4.00	0.0174
Num. Genes:Method:LOD diversity	43.71	2.00	< .0001
Num. Genes:Landscape:Detection	1262.35	4.00	< .0001
Num. Genes:Landscape:Sample Size	350.64	4.00	< .0001
Num. Genes:Landscape:LOD diversity	3.55	2.00	0.1699
Num. Genes:Detection:Sample Size	532.95	4.00	< .0001
Num. Genes:Detection:LOD diversity	2.93	2.00	0.231
Num. Genes:Sample Size:LOD diversity	50.25	2.00	< .0001
Method:Landscape:Detection	229.42	8.00	< .0001
Method:Landscape:Sample Size	215.46	8.00	< .0001
Method:Landscape:LOD diversity	627.29	4.00	< .0001
Method:Detection:Sample Size	21.27	8.00	0.0065
Method:Detection:LOD diversity	131.95	4.00	< .0001
Method:Sample Size:LOD diversity	88.97	4.00	< .0001
Landscape:Detection:Sample Size	2216.86	8.00	< .0001
Landscape:Detection:LOD diversity	867.13	4.00	< .0001
Landscape:Sample Size:LOD diversity	1175.97	4.00	< .0001
Detection:Sample Size:LOD diversity	819.51	4.00	< .0001
Num. Genes:Method:Landscape:Detection	12.73	8.00	0.1214
Num. Genes:Method:Landscape:Sample Size	7.44	8.00	0.4898
Num. Genes:Method:Landscape:LOD diversity	14.84	4.00	0.0051
Num. Genes:Method:Detection:Sample Size	25.87	8.00	0.0011
Num. Genes:Method:Detection:LOD diversity	53.99	4.00	< .0001
Num. Genes:Method:Sample Size:LOD diversity	2.51	4.00	0.6425
Num. Genes:Landscape:Detection:Sample Size	321.87	8.00	< .0001
Num. Genes:Landscape:Detection:LOD diversity	74.72	4.00	< .0001
Num. Genes:Landscape:Sample Size:LOD diversity	101.05	4.00	< .0001

Num. Genes:Detection:Sample Size:LOD diversity	115.78	4.00	< .0001
Method:Landscape:Detection:Sample Size	24.77	16.00	0.0739
Method:Landscape:Detection:LOD diversity	19.36	8.00	0.013
Method:Landscape:Sample Size:LOD diversity	23.93	8.00	0.0024
Method:Detection:Sample Size:LOD diversity	10.02	8.00	0.2634
Landscape:Detection:Sample Size:LOD diversity	237.96	8.00	< .0001

Table 3: Full model, 4-way interactions

17.2 Models fitted to each combination of fitness landscape by method

Remember each model uses 3780 observation: 35 replicates, 3 mutation rates, 2 variance settings, 2 number of genes, 3 detection regimes, 3 sample sizes. These correspond to 420 different fitness landscapes: 35 by 3 by 2 by 2. Each observation is itself the average of five different splits of the set of 20000 simulations.

17.2.1 Main effects

	Chisq	Df	Pr(>Chisq)
Num. Genes	325.32	1.00	< .0001
Sample Size	797.56	2.00	< .0001
Detection	1101.14	2.00	< .0001
LOD diversity	2.41	1.00	0.1206

Table 4: Represent..OT

	Chisq	Df	Pr(>Chisq)
Num. Genes	84.38	1.00	< .0001
Sample Size	179.66	2.00	< .0001
Detection	470.52	2.00	< .0001
LOD diversity	63.71	1.00	< .0001

Table 5: Local Peaks.OT

	Chisq	Df	Pr(>Chisq)
Num. Genes	9.65	1.00	0.0019
Sample Size	199.97	2.00	< .0001
Detection	206.52	2.00	< .0001
LOD diversity	162.60	1.00	< .0001

Table 6: RMF.OT

	Chisq	Df	Pr(>Chisq)
Num. Genes	206.70	1.00	< .0001
Sample Size	62.82	2.00	< .0001
Detection	692.92	2.00	< .0001
LOD diversity	79.61	1.00	< .0001

Table 7: Represent..CAPRI_AIC

	Chisq	Df	Pr(>Chisq)
Num. Genes	110.80	1.00	< .0001
Sample Size	48.42	2.00	< .0001
Detection	383.72	2.00	< .0001
LOD diversity	81.86	1.00	< .0001

Table 8: Local Peaks.CAPRI_AIC

	Chisq	Df	Pr(>Chisq)
Num. Genes	35.67	1.00	< .0001
Sample Size	144.12	2.00	< .0001
Detection	48.52	2.00	< .0001
LOD diversity	70.86	1.00	< .0001

Table 9: RMF.CAPRI.AIC

	Chisq	Df	Pr(>Chisq)
Num. Genes	157.75	1.00	< .0001
Sample Size	1331.50	2.00	< .0001
Detection	2493.82	2.00	< .0001
LOD diversity	0.35	1.00	0.5538

Table 10: Represent..CBN

	Chisq	Df	Pr(>Chisq)
Num. Genes	71.83	1.00	< .0001
Sample Size	315.26	2.00	< .0001
Detection	823.90	2.00	< .0001
LOD diversity	64.31	1.00	< .0001

Table 11: Local Peaks.CBN

	Chisq	Df	Pr(>Chisq)
Num. Genes	3.82	1.00	0.0507
Sample Size	100.38	2.00	< .0001
Detection	320.43	2.00	< .0001
LOD diversity	151.14	1.00	< .0001

Table 12: RMF.CBN

17.2.2 Two-way interactions

	Chisq	Df	Pr(>Chisq)
Num. Genes	229.87	1.00	< .0001
Sample Size	1084.79	2.00	< .0001
Detection	1145.03	2.00	< .0001
LOD diversity	0.02	1.00	0.8958
Num. Genes:Sample Size	168.67	2.00	< .0001
Num. Genes:Detection	705.00	2.00	< .0001
Num. Genes:LOD diversity	4.46	1.00	0.0347
Sample Size:Detection	726.35	4.00	< .0001
Sample Size:LOD diversity	308.31	2.00	< .0001
Detection:LOD diversity	1019.12	2.00	< .0001

Table 13: Represent..OT

	Chisq	Df	Pr(>Chisq)
Num. Genes	66.47	1.00	< .0001
Sample Size	163.41	2.00	< .0001
Detection	449.46	2.00	< .0001
LOD diversity	74.54	1.00	< .0001
Num. Genes:Sample Size	34.28	2.00	< .0001
Num. Genes:Detection	443.26	2.00	< .0001
Num. Genes:LOD diversity	0.02	1.00	0.895
Sample Size:Detection	307.15	4.00	< .0001
Sample Size:LOD diversity	21.38	2.00	< .0001
Detection:LOD diversity	403.27	2.00	< .0001

Table 14: Local Peaks.OT

	Chisq	Df	Pr(>Chisq)
Num. Genes	9.65	1.00	0.0019
Sample Size	213.84	2.00	< .0001
Detection	215.42	2.00	< .0001
LOD diversity	155.94	1.00	< .0001
Num. Genes:Sample Size	31.75	2.00	< .0001
Num. Genes:Detection	5.39	2.00	0.0676
Num. Genes:LOD diversity	0.05	1.00	0.8191
Sample Size:Detection	7.68	4.00	0.1041
Sample Size:LOD diversity	190.09	2.00	< .0001
Detection:LOD diversity	9.29	2.00	0.0096

Table 15: RMF.OT

	Chisq	Df	Pr(>Chisq)
Num. Genes	144.41	1.00	< .0001
Sample Size	70.06	2.00	< .0001
Detection	720.04	2.00	< .0001
LOD diversity	89.03	1.00	< .0001
Num. Genes:Sample Size	154.58	2.00	< .0001
Num. Genes:Detection	574.17	2.00	< .0001
Num. Genes:LOD diversity	5.16	1.00	0.0231
Sample Size:Detection	691.60	4.00	< .0001
Sample Size:LOD diversity	734.28	2.00	< .0001
Detection:LOD diversity	958.49	2.00	< .0001

Table 16: Represent..CAPRI_AIC

	Chisq	Df	Pr(>Chisq)
Num. Genes	90.86	1.00	< .0001
Sample Size	36.29	2.00	< .0001
Detection	351.07	2.00	< .0001
LOD diversity	89.21	1.00	< .0001
Num. Genes:Sample Size	8.98	2.00	0.0112
Num. Genes:Detection	309.60	2.00	< .0001
Num. Genes:LOD diversity	0.01	1.00	0.91
Sample Size:Detection	262.51	4.00	< .0001
Sample Size:LOD diversity	80.39	2.00	< .0001
Detection:LOD diversity	429.26	2.00	< .0001

Table 17: Local Peaks.CAPRI_AIC

	Chisq	Df	Pr(>Chisq)
Num. Genes	33.69	1.00	< .0001
Sample Size	164.54	2.00	< .0001
Detection	50.28	2.00	< .0001
LOD diversity	65.36	1.00	< .0001
Num. Genes:Sample Size	34.37	2.00	< .0001
Num. Genes:Detection	2.23	2.00	0.3273
Num. Genes:LOD diversity	1.24	1.00	0.2663
Sample Size:Detection	15.51	4.00	0.0037
Sample Size:LOD diversity	154.86	2.00	< .0001
Detection:LOD diversity	40.59	2.00	< .0001

Table 18: RMF.CAPRI_AIC

	Chisq	Df	Pr(>Chisq)
Num. Genes	111.35	1.00	< .0001
Sample Size	1727.19	2.00	< .0001
Detection	2829.04	2.00	< .0001
LOD diversity	2.47	1.00	0.116
Num. Genes:Sample Size	252.55	2.00	< .0001
Num. Genes:Detection	583.85	2.00	< .0001
Num. Genes:LOD diversity	11.81	1.00	6e-04
Sample Size:Detection	583.96	4.00	< .0001
Sample Size:LOD diversity	367.88	2.00	< .0001
Detection:LOD diversity	634.67	2.00	< .0001

Table 19: Represent..CBN

	Chisq	Df	Pr(>Chisq)
Num. Genes	56.06	1.00	< .0001
Sample Size	313.03	2.00	< .0001
Detection	836.24	2.00	< .0001
LOD diversity	76.98	1.00	< .0001
Num. Genes:Sample Size	35.78	2.00	< .0001
Num. Genes:Detection	413.32	2.00	< .0001
Num. Genes:LOD diversity	0.30	1.00	0.5868
Sample Size:Detection	296.90	4.00	< .0001
Sample Size:LOD diversity	19.84	2.00	< .0001
Detection:LOD diversity	292.03	2.00	< .0001

Table 20: Local Peaks.CBN

	Chisq	Df	Pr(>Chisq)
Num. Genes	3.82	1.00	0.0506
Sample Size	107.26	2.00	< .0001
Detection	331.99	2.00	< .0001
LOD diversity	144.95	1.00	< .0001
Num. Genes:Sample Size	21.97	2.00	< .0001
Num. Genes:Detection	4.24	2.00	0.12
Num. Genes:LOD diversity	1.16	1.00	0.2823
Sample Size:Detection	4.15	4.00	0.3856
Sample Size:LOD diversity	181.70	2.00	< .0001
Detection:LOD diversity	1.36	2.00	0.5071

Table 21: RMF.CBN

17.2.3 Four-way interactions

	Chisq	Df	Pr(>Chisq)
Num. Genes	236.38	1.00	< .0001
Sample Size	1105.42	2.00	< .0001
Detection	1046.08	2.00	< .0001
LOD diversity	0.84	1.00	0.3593
Num. Genes:Sample Size	193.98	2.00	< .0001
Num. Genes:Detection	778.62	2.00	< .0001
Num. Genes:LOD diversity	6.38	1.00	0.0116
Sample Size:Detection	722.33	4.00	< .0001
Sample Size:LOD diversity	322.04	2.00	< .0001
Detection:LOD diversity	1027.32	2.00	< .0001
Num. Genes:Sample Size:Detection	96.81	4.00	< .0001
Num. Genes:Sample Size:LOD diversity	28.18	2.00	< .0001
Num. Genes:Detection:LOD diversity	26.39	2.00	< .0001
Sample Size:Detection:LOD diversity	116.00	4.00	< .0001
Num. Genes:Sample Size:Detection:LOD diversity	76.54	4.00	< .0001

Table 22: Represent..OT

	Chisq	Df	Pr(>Chisq)
Num. Genes	57.63	1.00	< .0001
Sample Size	166.26	2.00	< .0001
Detection	380.88	2.00	< .0001
LOD diversity	72.98	1.00	< .0001
Num. Genes:Sample Size	35.33	2.00	< .0001
Num. Genes:Detection	409.75	2.00	< .0001
Num. Genes:LOD diversity	0.04	1.00	0.8394
Sample Size:Detection	326.02	4.00	< .0001
Sample Size:LOD diversity	24.07	2.00	< .0001
Detection:LOD diversity	412.53	2.00	< .0001
Num. Genes:Sample Size:Detection	198.32	4.00	< .0001
Num. Genes:Sample Size:LOD diversity	1.85	2.00	0.3967
Num. Genes:Detection:LOD diversity	5.28	2.00	0.0715
Sample Size:Detection:LOD diversity	219.97	4.00	< .0001
Num. Genes:Sample Size:Detection:LOD diversity	8.77	4.00	0.0672

Table 23: Local Peaks.OT

	Chisq	Df	Pr(>Chisq)
Num. Genes	9.69	1.00	0.0019
Sample Size	215.06	2.00	< .0001
Detection	215.20	2.00	< .0001
LOD diversity	155.03	1.00	< .0001
Num. Genes:Sample Size	32.40	2.00	< .0001
Num. Genes:Detection	5.40	2.00	0.0673
Num. Genes:LOD diversity	0.05	1.00	0.8191
Sample Size:Detection	7.55	4.00	0.1096
Sample Size:LOD diversity	191.25	2.00	< .0001
Detection:LOD diversity	9.34	2.00	0.0094
Num. Genes:Sample Size:Detection	0.26	4.00	0.9925
Num. Genes:Sample Size:LOD diversity	1.73	2.00	0.4219
Num. Genes:Detection:LOD diversity	32.05	2.00	< .0001
Sample Size:Detection:LOD diversity	1.65	4.00	0.799
Num. Genes:Sample Size:Detection:LOD diversity	1.63	4.00	0.8035

Table 24: RMF.OT

	Chisq	Df	Pr(>Chisq)
Num. Genes	163.81	1.00	< .0001
Sample Size	53.20	2.00	< .0001
Detection	661.33	2.00	< .0001
LOD diversity	100.78	1.00	< .0001
Num. Genes:Sample Size	170.47	2.00	< .0001
Num. Genes:Detection	654.22	2.00	< .0001
Num. Genes:LOD diversity	5.37	1.00	0.0205
Sample Size:Detection	691.21	4.00	< .0001
Sample Size:LOD diversity	750.44	2.00	< .0001
Detection:LOD diversity	972.35	2.00	< .0001
Num. Genes:Sample Size:Detection	48.16	4.00	< .0001
Num. Genes:Sample Size:LOD diversity	24.57	2.00	< .0001
Num. Genes:Detection:LOD diversity	42.41	2.00	< .0001
Sample Size:Detection:LOD diversity	116.19	4.00	< .0001
Num. Genes:Sample Size:Detection:LOD diversity	53.68	4.00	< .0001

Table 25: Represent..CAPRI_AIC

	Chisq	Df	Pr(>Chisq)
Num. Genes	87.57	1.00	< .0001
Sample Size	35.06	2.00	< .0001
Detection	329.14	2.00	< .0001
LOD diversity	88.97	1.00	< .0001
Num. Genes:Sample Size	9.80	2.00	0.0074
Num. Genes:Detection	309.31	2.00	< .0001
Num. Genes:LOD diversity	0.03	1.00	0.8725
Sample Size:Detection	272.21	4.00	< .0001
Sample Size:LOD diversity	87.34	2.00	< .0001
Detection:LOD diversity	437.56	2.00	< .0001
Num. Genes:Sample Size:Detection	63.36	4.00	< .0001
Num. Genes:Sample Size:LOD diversity	1.44	2.00	0.4879
Num. Genes:Detection:LOD diversity	4.16	2.00	0.1248
Sample Size:Detection:LOD diversity	73.43	4.00	< .0001
Num. Genes:Sample Size:Detection:LOD diversity	2.63	4.00	0.6213

Table 26: Local Peaks.CAPRI.AIC

	Chisq	Df	Pr(>Chisq)
Num. Genes	33.91	1.00	< .0001
Sample Size	165.45	2.00	< .0001
Detection	50.85	2.00	< .0001
LOD diversity	65.50	1.00	< .0001
Num. Genes:Sample Size	34.42	2.00	< .0001
Num. Genes:Detection	2.30	2.00	0.3163
Num. Genes:LOD diversity	1.25	1.00	0.2641
Sample Size:Detection	15.66	4.00	0.0035
Sample Size:LOD diversity	155.70	2.00	< .0001
Detection:LOD diversity	41.21	2.00	< .0001
Num. Genes:Sample Size:Detection	1.76	4.00	0.7805
Num. Genes:Sample Size:LOD diversity	9.35	2.00	0.0093
Num. Genes:Detection:LOD diversity	3.06	2.00	0.2169
Sample Size:Detection:LOD diversity	2.52	4.00	0.6414
Num. Genes:Sample Size:Detection:LOD diversity	0.83	4.00	0.9342

Table 27: RMF.CAPRI.AIC

	Chisq	Df	Pr(>Chisq)
Num. Genes	109.99	1.00	< .0001
Sample Size	1798.63	2.00	< .0001
Detection	2781.42	2.00	< .0001
LOD diversity	4.29	1.00	0.0384
Num. Genes:Sample Size	292.35	2.00	< .0001
Num. Genes:Detection	633.77	2.00	< .0001
Num. Genes:LOD diversity	13.30	1.00	3e-04
Sample Size:Detection	577.47	4.00	< .0001
Sample Size:LOD diversity	383.39	2.00	< .0001
Detection:LOD diversity	647.75	2.00	< .0001
Num. Genes:Sample Size:Detection	91.32	4.00	< .0001
Num. Genes:Sample Size:LOD diversity	48.42	2.00	< .0001
Num. Genes:Detection:LOD diversity	16.31	2.00	3e-04
Sample Size:Detection:LOD diversity	97.29	4.00	< .0001
Num. Genes:Sample Size:Detection:LOD diversity	98.90	4.00	< .0001

Table 28: Represent..CBN

	Chisq	Df	Pr(>Chisq)
Num. Genes	46.47	1.00	< .0001
Sample Size	326.11	2.00	< .0001
Detection	762.69	2.00	< .0001
LOD diversity	74.96	1.00	< .0001
Num. Genes:Sample Size	35.59	2.00	< .0001
Num. Genes:Detection	377.15	2.00	< .0001
Num. Genes:LOD diversity	0.43	1.00	0.5135
Sample Size:Detection	315.03	4.00	< .0001
Sample Size:LOD diversity	21.87	2.00	< .0001
Detection:LOD diversity	296.12	2.00	< .0001
Num. Genes:Sample Size:Detection	213.53	4.00	< .0001
Num. Genes:Sample Size:LOD diversity	2.25	2.00	0.3248
Num. Genes:Detection:LOD diversity	20.79	2.00	< .0001
Sample Size:Detection:LOD diversity	233.91	4.00	< .0001
Num. Genes:Sample Size:Detection:LOD diversity	4.47	4.00	0.3461

Table 29: Local Peaks.CBN

	Chisq	Df	Pr(>Chisq)
Num. Genes	3.85	1.00	0.0499
Sample Size	107.17	2.00	< .0001
Detection	332.57	2.00	< .0001
LOD diversity	143.99	1.00	< .0001
Num. Genes:Sample Size	22.27	2.00	< .0001
Num. Genes:Detection	4.25	2.00	0.1195
Num. Genes:LOD diversity	1.16	1.00	0.2811
Sample Size:Detection	4.19	4.00	0.3803
Sample Size:LOD diversity	182.19	2.00	< .0001
Detection:LOD diversity	1.29	2.00	0.5253
Num. Genes:Sample Size:Detection	0.42	4.00	0.9807
Num. Genes:Sample Size:LOD diversity	0.44	2.00	0.8014
Num. Genes:Detection:LOD diversity	24.00	2.00	< .0001
Sample Size:Detection:LOD diversity	8.52	4.00	0.0744
Num. Genes:Sample Size:Detection:LOD diversity	1.49	4.00	0.8284

Table 30: RMF.CBN

18 Number of mutations of local maxima and performance

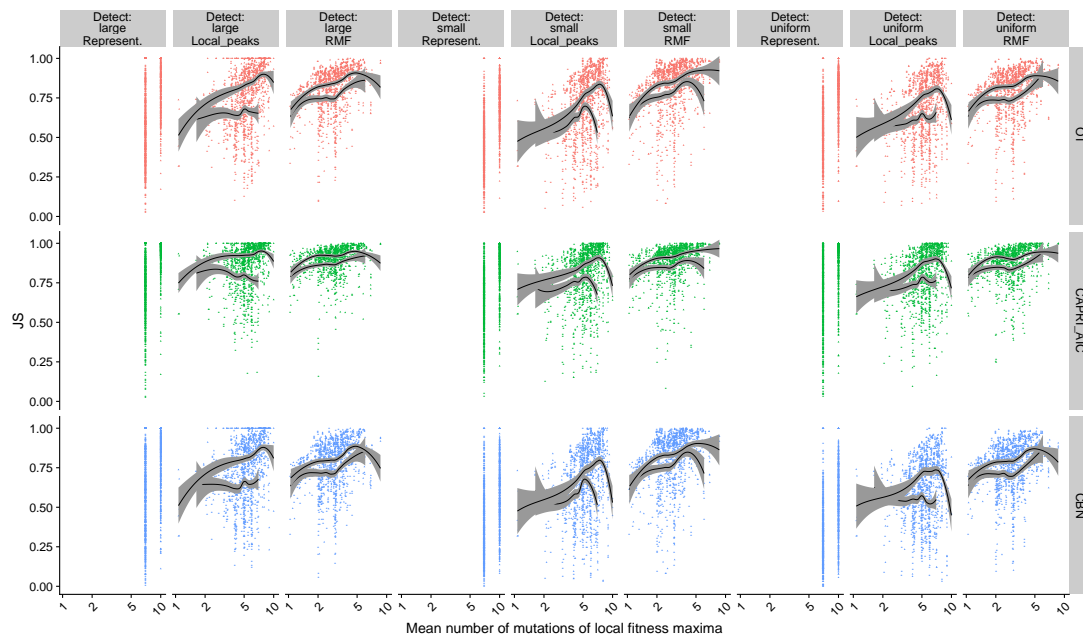


Figure 12: Mean number of mutations of local maxima and JS

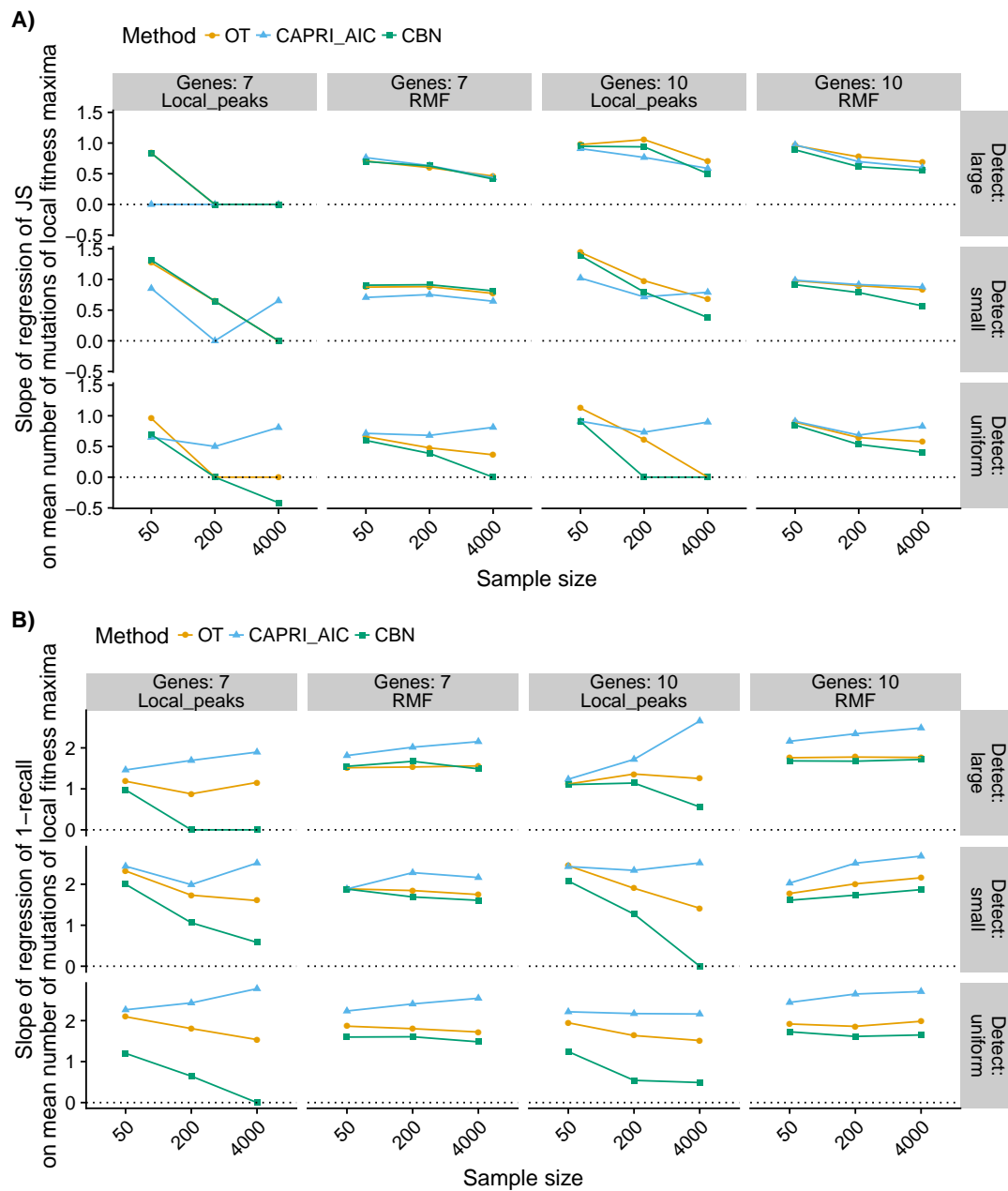


Figure 13: Slopes of regression of JS and 1-recall on mean number of mutations of local maxima.

19 LOD and CPM diversity: ratios and slopes

The following R code will, via a simple example, show that it is easy to have data where the average of the ratios is larger than one whereas the slope of the regression is negative:

```
a <- 10
n <- 100
sd <- 0.5
x <- runif(n, min = 1, max = 5)
y <- -1 * x + a + rnorm(n, mean = 0, sd = sd)
plot(y ~ x)
summary(lm(y ~ x))
mean(y/x)
```

20 Regression of individual CBN unpredictability estimates on LOD diversity

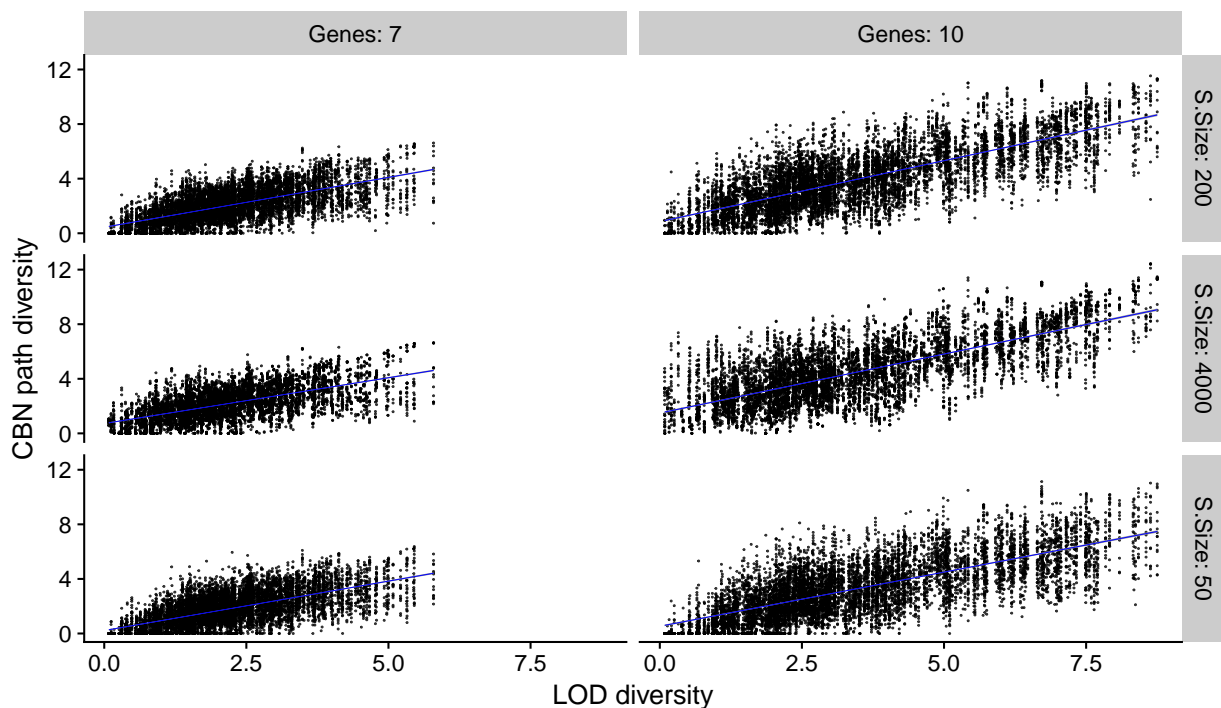


Figure 14: Scatterplot of diversity of paths to the maximum from CBN on true LOD diversity. Each panel contains 9450 points: 3 fitness landscapes * 3 initial sizes * 2 mutation regimes * 3 detection regimes * 5 replicate splits.

References

- [1] Attolini, C., Cheng, Y., Beroukhim, R., Getz, G., Abdel-Wahab, O., Levine, R. L., Mellinghoff, I. K., Michor, F., 2010. A mathematical framework to determine the temporal sequence of somatic genetic events in cancer. *Proceedings of the National Academy of Sciences*, **107**(41):17604–17609. doi:10.1073/pnas.1009117107/-/DCSupplemental.www.pnas.org/cgi/doi/10.1073/pnas.1009117107. URL <http://www.pnas.org/content/107/41/17604.short>.

- [2] Bell, D., Berchuck, A., Birrer, M., Chien, J., Cramer, D. W., Dao, F., Dhir, R., DiSaia, P., Gabra, H., Glenn, P., Godwin, a. K., Gross, J., Hartmann, L., Huang, M., Huntsman, D. G., Iacocca, M., Imielinski, M., Kalloger, S., Karlan, B. Y., a. Levine, D., Mills, G. B., Morrison, C., Mutch, D., Olvera, N., Orsulic, S., Park, K., Petrelli, N., Rabeno, B., Rader, J. S., Sikic, B. I., Smith-McCune, K., Sood, a. K., Bowtell, D., Penny, R., Testa, J. R., Chang, K., Dinh, H. H., a. Drummond, J., Fowler, G., Gunaratne, P., Hawes, a. C., Kovar, C. L., Lewis, L. R., Morgan, M. B., Newsham, I. F., Santibanez, J., Reid, J. G., Trevino, L. R., Wu, Y.-Q., Wang, M., Muzny, D. M., a. Wheeler, D., a. Gibbs, R., Getz, G., Lawrence, M. S., Cibulskis, K., Sivachenko, a. Y., Sougnez, C., Voet, D., Wilkinson, J., Bloom, T., Ardlie, K., Fennell, T., Baldwin, J., Gabriel, S., Lander, E. S., Ding, L., Fulton, R. S., Koboldt, D. C., McLellan, M. D., Wylie, T., Walker, J., O’Laughlin, M., Dooling, D. J., Fulton, L., Abbott, R., Dees, N. D., Zhang, Q., Kandoth, C., Wendl, M., Schierding, W., Shen, D., Harris, C. C., Schmidt, H., Kalicki, J., Delehaunty, K. D., Fronick, C. C., Demeter, R., Cook, L., Wallis, J. W., Lin, L., Magrini, V. J., Hodges, J. S., Eldred, J. M., Smith, S. M., Pohl, C. S., Vandin, F., Raphael, B. J., Weinstock, G. M., Mardis, E. R., Wilson, R. K., Meyerson, M., Winckler, W., Verhaak, R. G. W., Carter, S. L., Mermel, C. H., Saksena, G., Nguyen, H., Onofrio, R. C., Hubbard, D., Gupta, S., Crenshaw, A., Ramos, a. H., Chin, L., Protopopov, A., Zhang, J., Kim, T. M., Perna, I., Xiao, Y., Zhang, H., Ren, G., Sathiamoorthy, N., Park, R. W., Lee, E., Park, P. J., Kucherlapati, R., Absher, D. M., Waite, L., Sherlock, G., Brooks, J. D., Li, J. Z., Xu, J., Myers, R. M., Laird, P. W., Cope, L., Herman, J. G., Shen, H., Weisenberger, D. J., Noushmehr, H., Pan, F., Triche Jr, T., Berman, B. P., Van Den Berg, D. J., Buckley, J., Baylin, S. B., Spellman, P. T., Purdom, E., Neuvial, P., Bengtsson, H., Jakkula, L. R., Durinck, S., Han, J., Dorton, S., Marr, H., Choi, Y. G., Wang, V., Wang, N. J., Ngai, J., Conboy, J. G., Parvin, B., Feiler, H. S., Speed, T. P., Gray, J. W., Socci, N. D., Liang, Y., Taylor, B. S., Schultz, N., Borsu, L., Lash, a. E., Brennan, C., Viale, A., Sander, C., Ladanyi, M., a. Hoadley, K., Meng, S., Du, Y., Shi, Y., Li, L., Turman, Y. J., Zang, D., Helms, E. B., Balu, S., Zhou, X., Wu, J., Topal, M. D., Hayes, D. N., Perou, C. M., Zhang, J., Wu, C. J., Shukla, S., Sivachenko, A., Jing, R., Liu, Y., Noble, M., Carter, H., Kim, D., Karchin, R., Korkola, J. E., Heiser, L. M., Cho, R. J., Hu, Z., Cerami, E., Olshen, A., Reva, B., Antipin, Y., Shen, R., Mankoo, P., Sheridan, R., Ciriello, G., Chang, W. K., a. Bernanke, J., Haussler, D., Benz, C. C., Stuart, J. M., Benz, S. C., Sanborn, J. Z., Vaske, C. J., Zhu, J., Szeto, C., Scott, G. K., Yau, C., Wilkerson, M. D., Zhang, N., Akbani, R., a. Baggerly, K., Yung, W. K., Weinstein, J. N., Shelton, T., Grimm, D., Hatfield, M., Morris, S., Yena, P., Rhodes, P., Sherman, M., Paulauskis, J., Millis, S., Kahn, A., Greene, J. M., Sfeir, R., a. Jensen, M., Chen, J., Whitmore, J., Alonso, S., Jordan, J., Chu, A., Zhang, J., Barker, A., Compton, C., Eley, G., Ferguson, M., Fielding, P., Gerhard, D. S., Myles, R., Schaefer, C., Mills Shaw, K. R., Vaught, J., Vockley, J. B., Good, P. J., Guyer, M. S., Ozenberger, B., Peterson, J., Thomson, E., 2011. Integrated genomic analyses of ovarian carcinoma. *Nature*, **474(7353)**:609–615. doi:10.1038/nature10166. URL <http://www.nature.com/doifinder/10.1038/nature10166>.
- [3] Brooks, M. E., Kristensen, K., van Benthem, K. J., Magnusson, A., Berg, C. W., Nielsen, A., Skaug, H. J., Maechler, M., Bolker, B. M., 2017. glmmTMB Balances Speed and Flexibility Among Packages for Zero-inflated Generalized Linear Mixed Modeling. *The R Journal*, **9(2)**:378–400. URL <https://journal.r-project.org/archive/2017/RJ-2017-066/index.html>.
- [4] Cancer Genome Atlas Network, 2012. Comprehensive molecular characterization of human colon and rectal cancer. *Nature*, **487(7407)**:330–337. doi:10.1038/nature11252. URL <https://www.nature.com/articles/nature11252>.
- [5] Caravagna, G., Graudenzi, A., Ramazzotti, D., Sanz-Pamplona, R., Sano, L. D., Mauri, G., Moreno, V., Antoniotti, M., Mishra, B., 2016. Algorithmic methods to infer the evolutionary trajectories in cancer progression. *PNAS*, **113(28)**:E4025–E4034. doi:10.1073/pnas.1520213113. URL <http://www.pnas.org/content/113/28/E4025>.

- [6] Cheng, Y.-K., Beroukhi, R., Levine, R. L., Mellinghoff, I. K., Holland, E. C., Michor, F., 2012. A mathematical methodology for determining the temporal order of pathway alterations arising during gliomagenesis. *PLoS computational biology*, **8(1)**:e1002337. doi:10.1371/journal.pcbi.1002337. URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3252265&tool=pmcentrez&rendertype=abstract>.
- [7] Crona, K., Greene, D., Barlow, M., 2013. The peaks and geometry of fitness landscapes. *Journal of Theoretical Biology*, **317**:1–10. doi:10.1016/j.jtbi.2012.09.028. URL <http://www.sciencedirect.com/science/article/pii/S0022519312005061>.
- [8] Diaz-Uriarte, R., 2015. Identifying restrictions in the order of accumulation of mutations during tumor progression: Effects of passengers, evolutionary models, and sampling. *BMC Bioinformatics*, **16(41)**:0–36. doi:doi:10.1186/s12859-015-0466-7. URL <http://www.biomedcentral.com/1471-2105/16/41/abstract>.
- [9] Diaz-Uriarte, R., 2017. OncoSimulR: Genetic simulation with arbitrary epistasis and mutator genes in asexual populations. *Bioinformatics*, **33(12)**:1898–1899. doi:10.1093/bioinformatics/btx077. URL <https://academic.oup.com/bioinformatics/article/33/12/1898/2982052/OncoSimulR-genetic-simulation-with-arbitrary>.
- [10] Diaz-Uriarte, R., 2018. Cancer progression models and fitness landscapes: A many-to-many relationship. *Bioinformatics*, **34(5)**:836–844. doi:10.1093/bioinformatics/btx663. URL <https://academic.oup.com/bioinformatics/article/34/5/836/4557185>.
- [11] Ding, L., Getz, G., Wheeler, D. A., Mardis, E. R., McLellan, M. D., Cibulskis, K., Sougnez, C., Greulich, H., Muzny, D. M., Morgan, M. B., Fulton, L., Fulton, R. S., Zhang, Q., Wendl, M. C., Lawrence, M. S., Larson, D. E., Chen, K., Dooling, D. J., Sabo, A., Hawes, A. C., Shen, H., Jhangiani, S. N., Lewis, L. R., Hall, O., Zhu, Y., Mathew, T., Ren, Y., Yao, J., Scherer, S. E., Clerc, K., Metcalf, G. A., Ng, B., Milosavljevic, A., Gonzalez-Garay, M. L., Osborne, J. R., Meyer, R., Shi, X., Tang, Y., Koboldt, D. C., Lin, L., Abbott, R., Miner, T. L., Pohl, C., Fewell, G., Haippek, C., Schmidt, H., Dunford-Shore, B. H., Kraja, A., Crosby, S. D., Sawyer, C. S., Vickery, T., Sander, S., Robinson, J., Winckler, W., Baldwin, J., Chirieac, L. R., Dutt, A., Fennell, T., Hanna, M., Johnson, B. E., Onofrio, R. C., Thomas, R. K., Tonon, G., Weir, B. A., Zhao, X., Ziaugra, L., Zody, M. C., Giordano, T., Orringer, M. B., Roth, J. A., Spitz, M. R., Wistuba, I. I., Ozenberger, B., Good, P. J., Chang, A. C., Beer, D. G., Watson, M. A., Ladanyi, M., Broderick, S., Yoshizawa, A., Travis, W. D., Pao, W., Province, M. A., Weinstock, G. M., Varmus, H. E., Gabriel, S. B., Lander, E. S., Gibbs, R. A., Meyerson, M., Wilson, R. K., 2008. Somatic mutations affect key pathways in lung adenocarcinoma. *Nature*, **455(7216)**:1069–1075. doi:10.1038/nature07423.
- [12] Farahani, H. S., Lagergren, J., 2013. Learning oncogenetic networks by reducing to mixed integer linear programming. *PLoS ONE*, **8(6)**:e65773. doi:10.1371/journal.pone.0065773. URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3683041&tool=pmcentrez&rendertype=abstract>.
- [13] Ferretti, L., Schmiegel, B., Weinreich, D., Yamauchi, A., Kobayashi, Y., Tajima, F., Achaz, G., 2016. Measuring epistasis in fitness landscapes: The correlation of fitness effects of mutations. *Journal of Theoretical Biology*, **396**:132–143. doi:10.1016/j.jtbi.2016.01.037. URL <http://www.sciencedirect.com/science/article/pii/S0022519316000771>.
- [14] Fox, J., Weisberg, S., 2011. *An R Companion to Applied Regression, 2nd Ed.* Sage, Thousand Oaks, CA.
- [15] Gerstung, M., Eriksson, N., Lin, J., Vogelstein, B., Beerenwinkel, N., 2011. The Temporal Order of Genetic and Pathway Alterations in Tumorigenesis. *PLoS ONE*, **6(11)**:e27136. doi:10.1371/journal.pone.0027136. URL <http://dx.plos.org/10.1371/journal.pone.0027136%0020http://www.bsse.ethz.ch/cbg/software/ct-cbn>.

- [16] Hosseini, S.-R., 2018. Quantifying the predictability of cancer progression using Conjunctive Bayesian Networks. M.Sc. Thesis, Swiss Federal Institute of Technology, Zürich.
- [17] Jones, S., Zhang, X., Parsons, D. W., Lin, J. C.-H., Leary, R. J., Angenendt, P., Mankoo, P., Carter, H., Kamiyama, H., Jimeno, A., Hong, S.-M., Fu, B., Lin, M.-T., Calhoun, E. S., Kamiyama, M., Walter, K., Nikolskaya, T., Nikolsky, Y., Hartigan, J., Smith, D. R., Hidalgo, M., Leach, S. D., Klein, A. P., Jaffee, E. M., Goggins, M., Maitra, A., Iacobuzio-Donahue, C., Eshleman, J. R., Kern, S. E., Hruban, R. H., Karchin, R., Papadopoulos, N., Parmigiani, G., Vogelstein, B., Velculescu, V. E., Kinzler, K. W., 2008. Core signaling pathways in human pancreatic cancers revealed by global genomic analyses. *Science (New York, N.Y.)*, **321(5897)**:1801–6. doi:10.1126/science.1164368. URL <http://www.ncbi.nlm.nih.gov/pubmed/18772397>.
- [18] McFarland, C. D., Korolev, K. S., Kryukov, G. V., Sunyaev, S. R., Mirny, L. A., 2013. Impact of deleterious passenger mutations on cancer progression. *Proceedings of the National Academy of Sciences of the United States of America*, **110(8)**:2910–5. doi:10.1073/pnas.1213968110. URL <http://www.ncbi.nlm.nih.gov/pubmed/23388632>.
- [19] Misra, N., Szczurek, E., Vingron, M., 2014. Inferring the paths of somatic evolution in cancer. *Bioinformatics (Oxford, England)*, **30(17)**:2456–2463. doi:10.1093/bioinformatics/btu319. URL <http://www.ncbi.nlm.nih.gov/pubmed/24812340>.
- [20] Montazeri, H., Kuipers, J., Kouyos, R., Böni, J., Yerly, S., Klimkait, T., Aubert, V., Günthard, H. F., Beerenwinkel, N., Study, T. S. H. C., 2016. Large-scale inference of conjunctive Bayesian networks. *Bioinformatics*, **32(17)**:i727–i735. doi:10.1093/bioinformatics/btw459. URL <http://bioinformatics.oxfordjournals.org/content/32/17/i727>.
- [21] Parsons, D. W., Jones, S., Zhang, X., Lin, J. C.-H., Leary, R. J., Angenendt, P., Mankoo, P., Carter, H., Siu, I.-M., Gallia, G. L., Olivi, A., McLendon, R., Rasheed, B. A., Keir, S., Nikolskaya, T., Nikolsky, Y., Busam, D. A., Tekleab, H., Diaz, L. A., Hartigan, J., Smith, D. R., Strausberg, R. L., Marie, S. K. N., Shinjo, S. M. O., Yan, H., Riggins, G. J., Bigner, D. D., Karchin, R., Papadopoulos, N., Parmigiani, G., Vogelstein, B., Velculescu, V. E., Kinzler, K. W., 2008. An Integrated Genomic Analysis of Human Glioblastoma Multiforme. *Science*, **321(5897)**:1807–1812. doi:10.1126/science.1164382. URL <http://science.sciencemag.org/content/321/5897/1807>.
- [22] Sakoparnig, T., Beerenwinkel, N., 2012. Efficient sampling for Bayesian inference of conjunctive Bayesian networks. *Bioinformatics (Oxford, England)*, **28(18)**:2318–24. doi:10.1093/bioinformatics/bts433. URL <http://www.ncbi.nlm.nih.gov/pubmed/22782551> <http://www.bsse.ethz.ch/cbg/software/bayes-cbn>.
- [23] Szabo, A., Pappas, L., 2013. Oncotree: Estimating oncogenetic trees. R package version 0.3.3. URL <http://cran.r-project.org/package=Oncotree>.
- [24] Wood, L. D., Parsons, D. W., Jones, S., Lin, J., Sjoblom, T., Leary, R. J., Shen, D., Boca, S. M., Barber, T., Ptak, J., Silliman, N., Szabo, S., Dezso, Z., Ustyanksky, V., Nikolskaya, T., Nikolsky, Y., Karchin, R., Wilson, P. A., Kaminker, J. S., Zhang, Z., Croshaw, R., Willis, J., Dawson, D., Shipitsin, M., Willson, J. K. V., Sukumar, S., Polyak, K., Park, B. H., Pethiyagoda, C. L., Pant, P. V. K., Ballinger, D. G., Sparks, A. B., Hartigan, J., Smith, D. R., Suh, E., Papadopoulos, N., Buckhaults, P., Markowitz, S. D., Parmigiani, G., Kinzler, K. W., Velculescu, V. E., Vogelstein, B., 2007. The Genomic Landscapes of Human Breast and Colorectal Cancers. *Science*, **318(5853)**:1108–1113. doi:10.1126/science.1145720. URL <http://dx.doi.org/10.1126/science.1145720>.