

New report

O. Denas

June 5, 2017

Contents

1	Current performance – not working	1
2	Input properties	1
2.1	Node properties	1
2.2	Consecutive <code>wl()</code> calls	3
2.3	Consecutive <code>parent()</code> calls	5
3	<code>wl()</code> optimizations	7
3.1	sandbox tests	7
3.2	full tests	8
4	parent optimizations	19
4.1	sandbox tests	19
4.2	full tests	22
5	parallelization	22

1 Current performance – not working

2 Input properties

All experiments are performed on the following inputs (the mutation period refers to the index sequence):

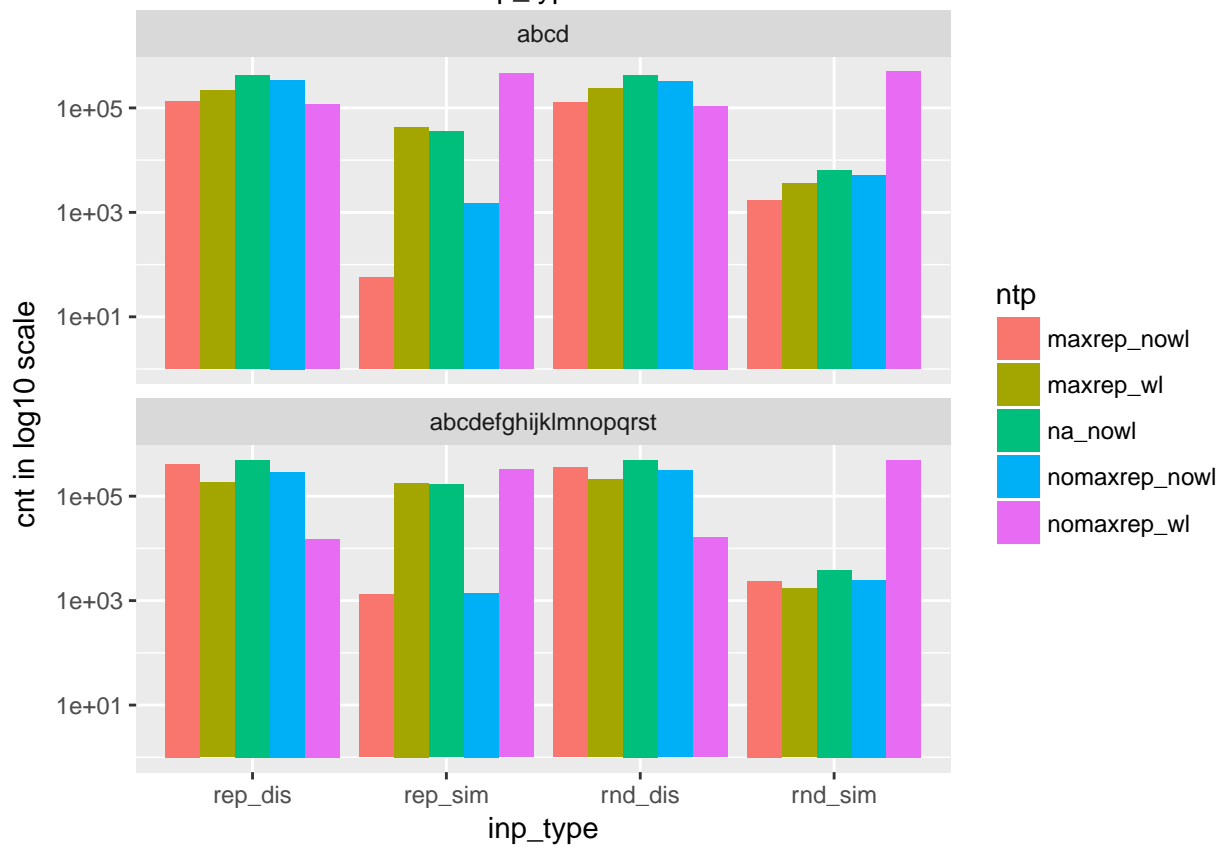
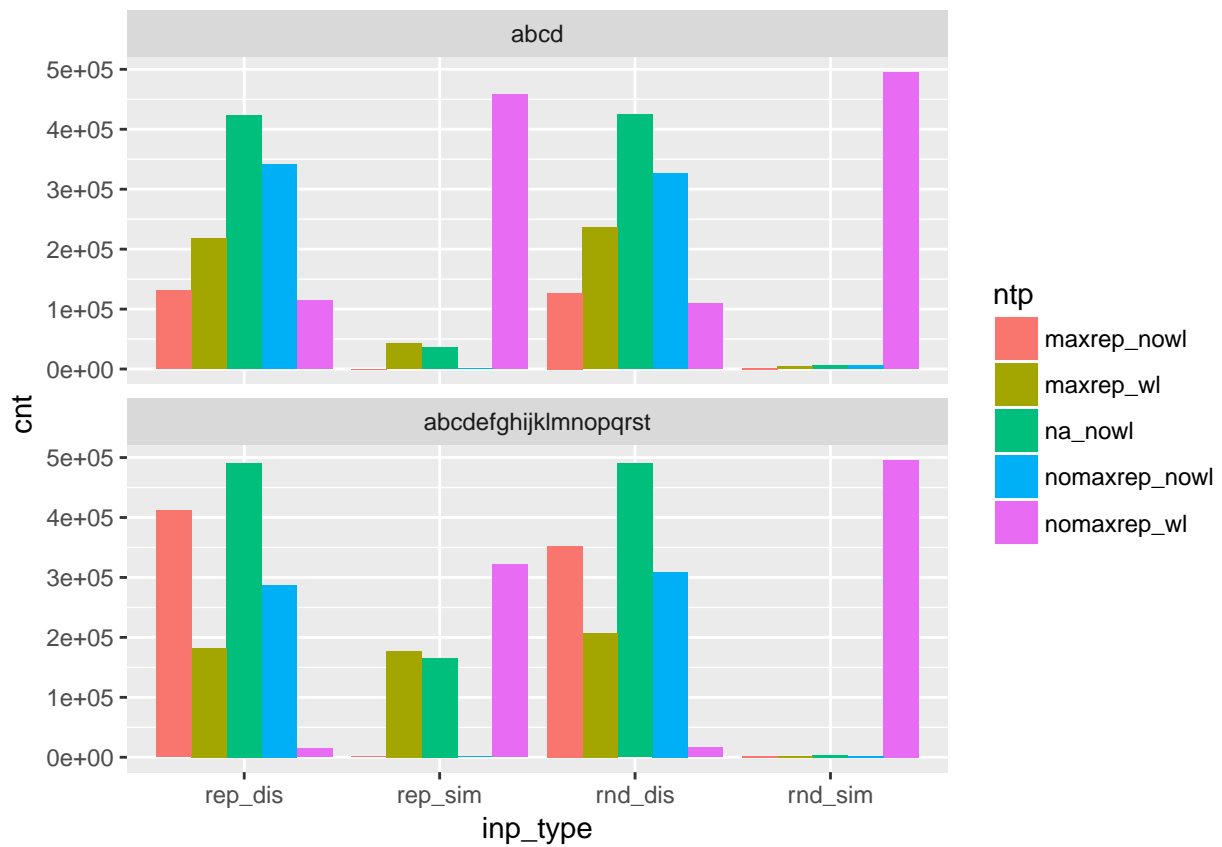
Table 1: Input datasets

s_type	t_type	slen	tlen	alp
rep	dis	100'000'000	500'000	abcdefghijklmnpqrst
rep	dis	100'000'000	500'000	abcd
rep	sim	100'000'000	500'000	abcdefghijklmnpqrst
rep	sim	100'000'000	500'000	abcd
rnd	dis	100'000'000	500'000	abcdefghijklmnpqrst
rnd	dis	100'000'000	500'000	abcd
rnd	sim	100'000'000	500'000	abcdefghijklmnpqrst
rnd	sim	100'000'000	500'000	abcd

In this section, we run the matching statistics algorithm and gather information on the way.

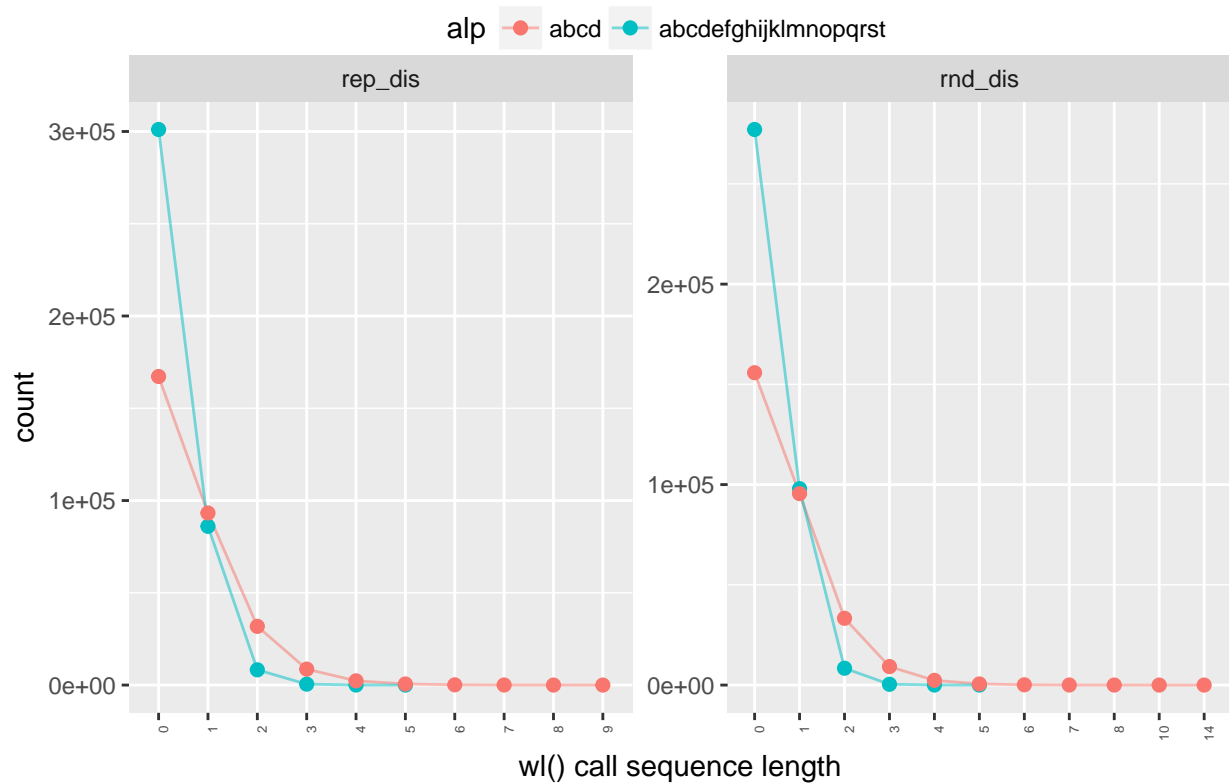
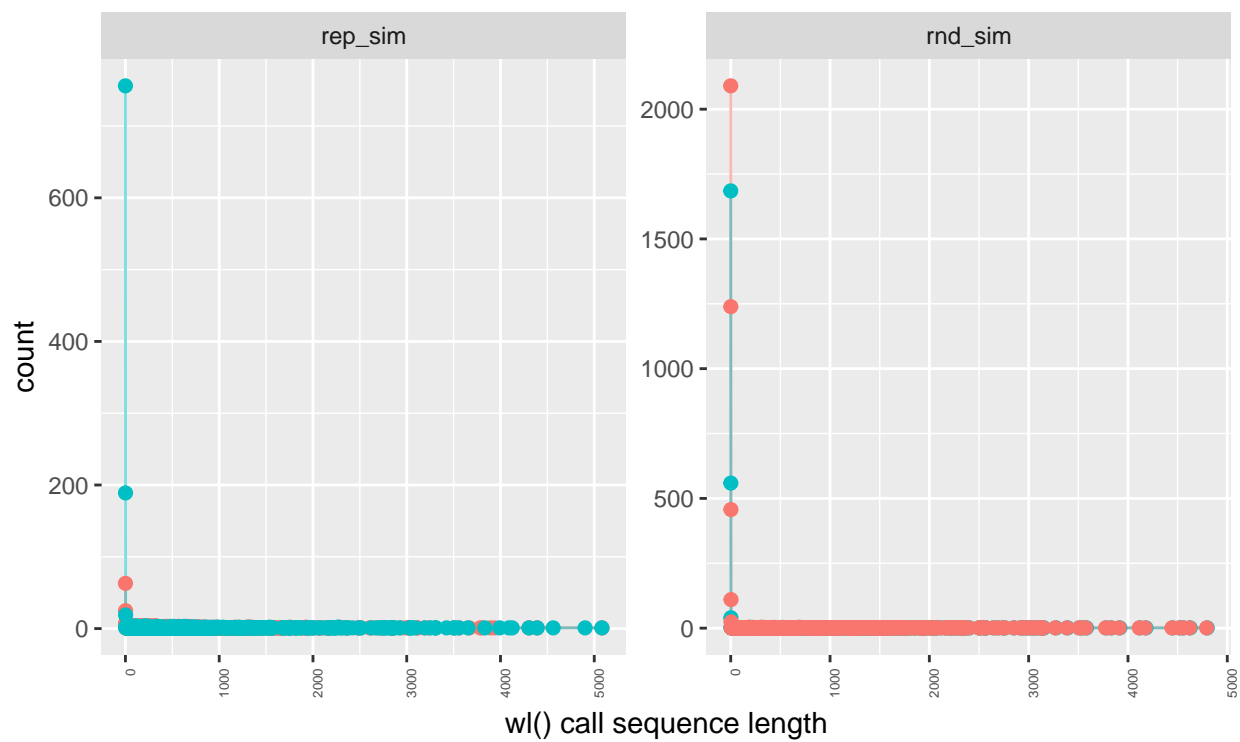
2.1 Node properties

Before making a `wl(v, c)` call register whether `v` is a maximal repeat, and whether the Weiner link exists or not.



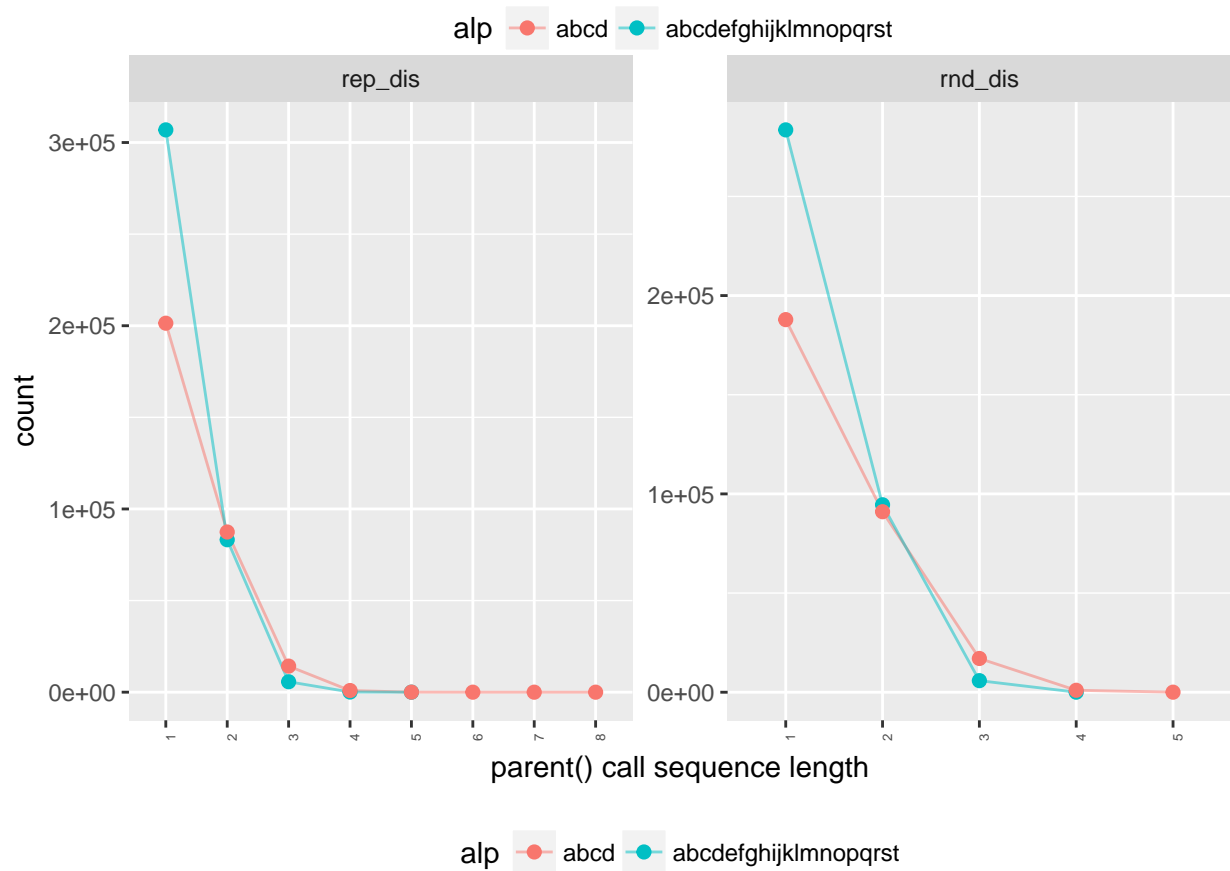
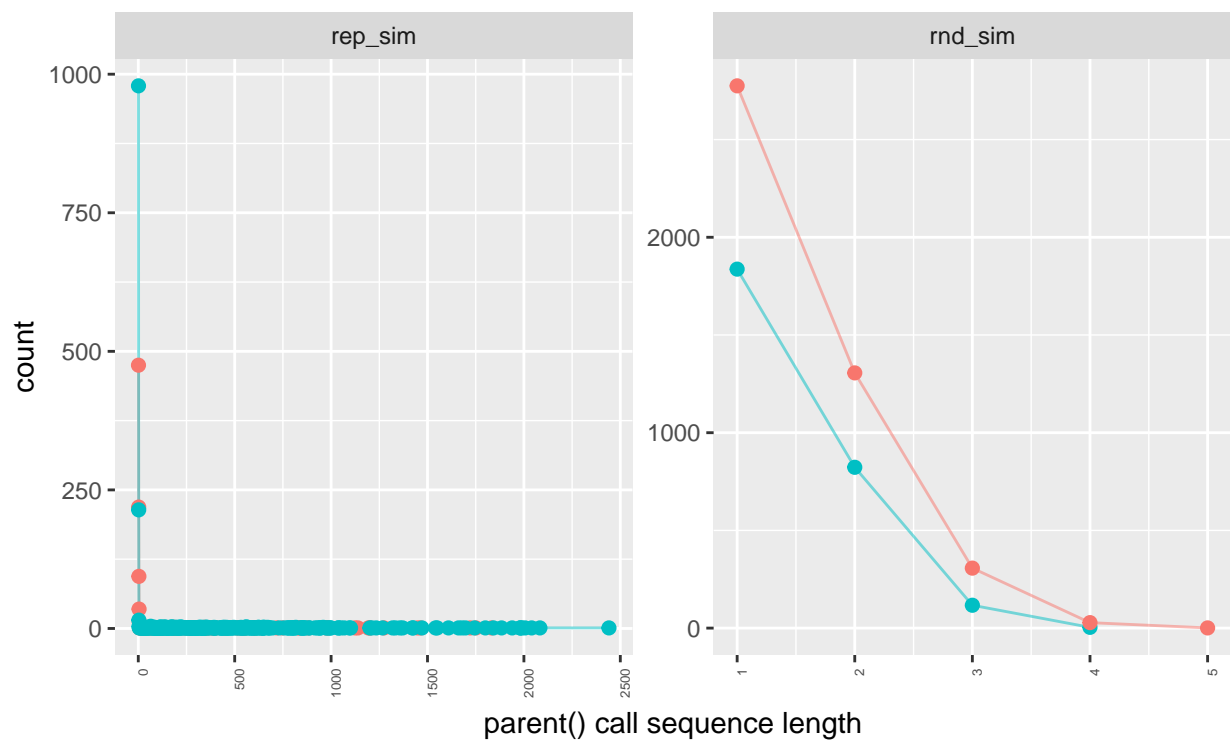
2.2 Consecutive `wl()` calls

During the MS vector construction, count the number of consecutive `wl()` calls due to matches between reversed indexed string and the query. In other words count the k -length iterations of the while cycle.



2.3 Consecutive `parent()` calls

During the MS vector construction, count the number of consecutive `parent()` calls after a failed `w1()` call.



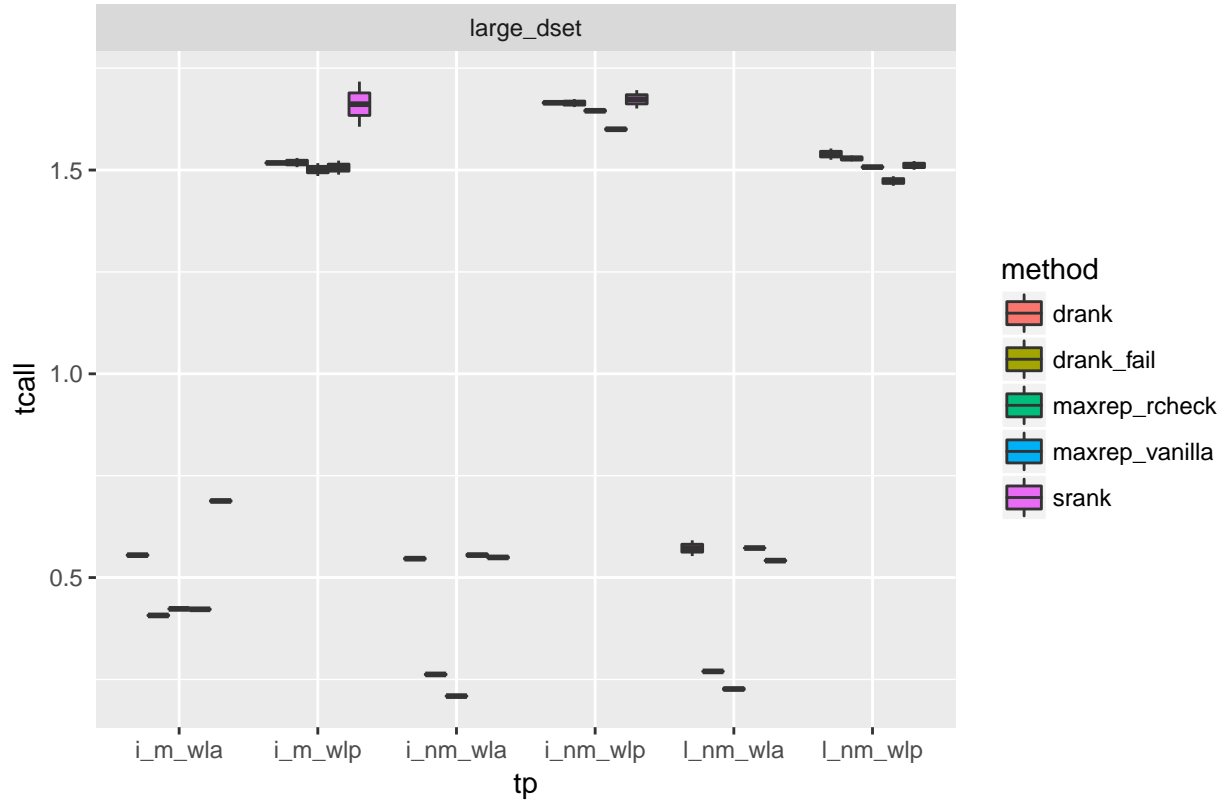
3 wl() optimizations

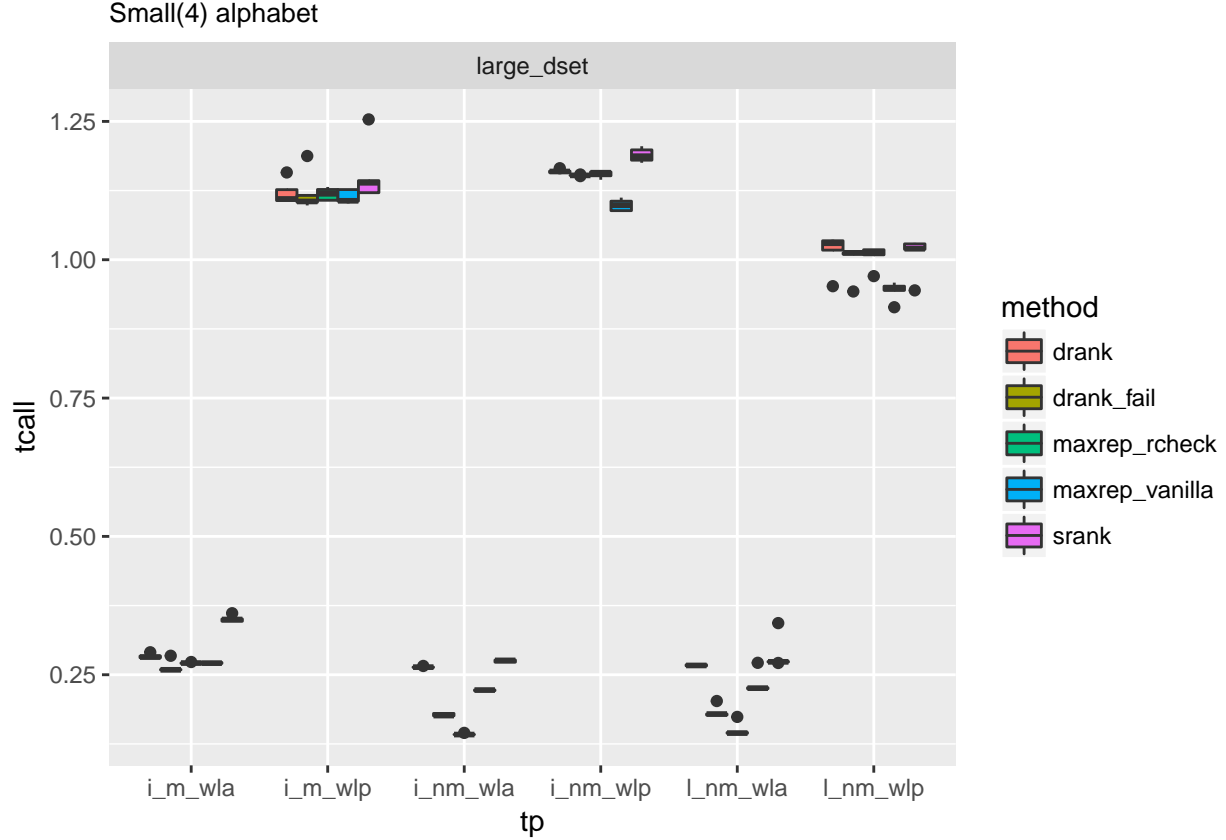
3.1 sandbox tests

The 4 plots correspond to combinations of index size (one 100Mb and one 1Mb) and alphabet size (one 4 and one 20). Each plot shows 6 measurements based on

- node type: leaf (l) or internal (i)
- maximality: maximal (m) or non-maximal (nm)
- WL presence: present (wlp) or absent (wla)

Large(20) alphabet





3.2 full tests

I run the program several times with and without the optimization. The pointrange plots report (median, with quartile ranges) the relative difference of each optimized time from the average non-optimized time in the construction time of the MS vector which is

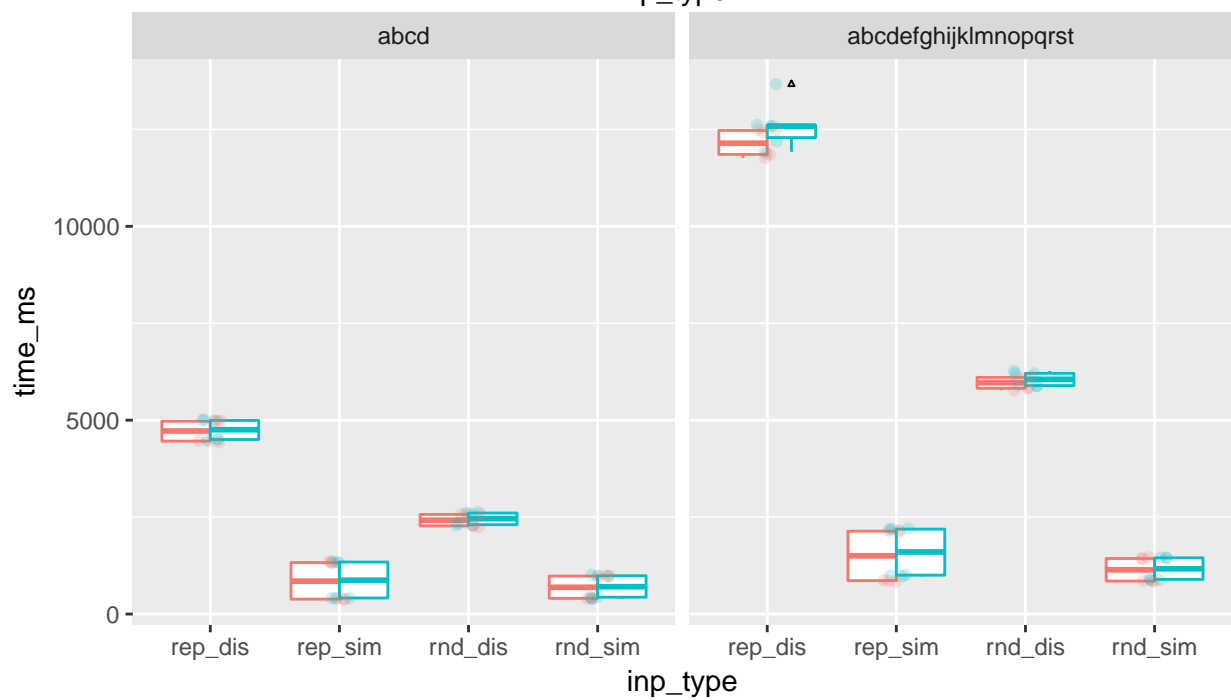
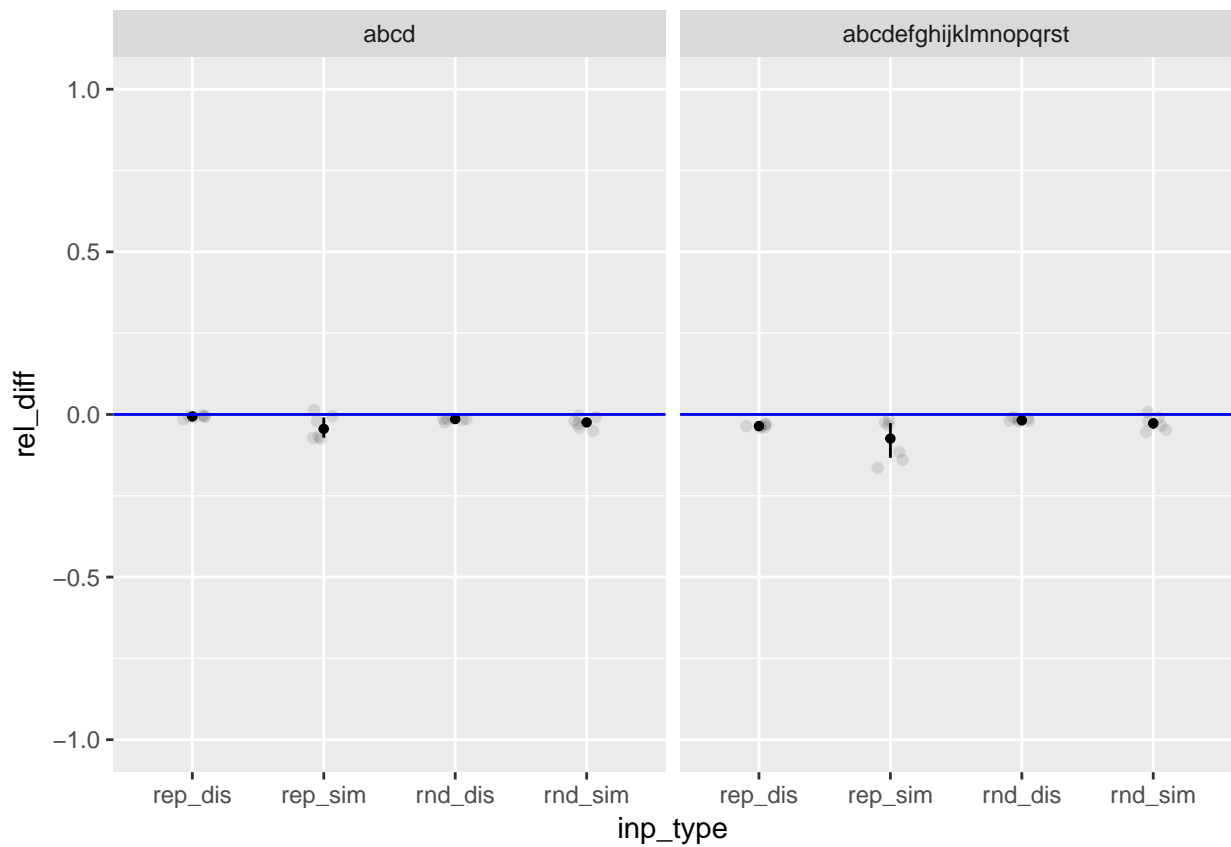
$$d^{(i)} = \frac{t_{\text{opt}}^{(i)} - \bar{t}_{\text{non_opt}}}{\max\{t_{\text{opt}}^{(i)}, \bar{t}_{\text{non_opt}}\}}$$

with $\bar{t}_{\text{non_opt}} = 1/n \sum t_{\text{non_opt}}^{(i)}$. Hence negative values indicate a speedup by the optimization over the average non-optimized time: -0.5 is a 2x speedup etc, in general the speed up is $-1/d^{(i)}$.

The boxplots report the raw times for the MS construction.

3.2.1 drank

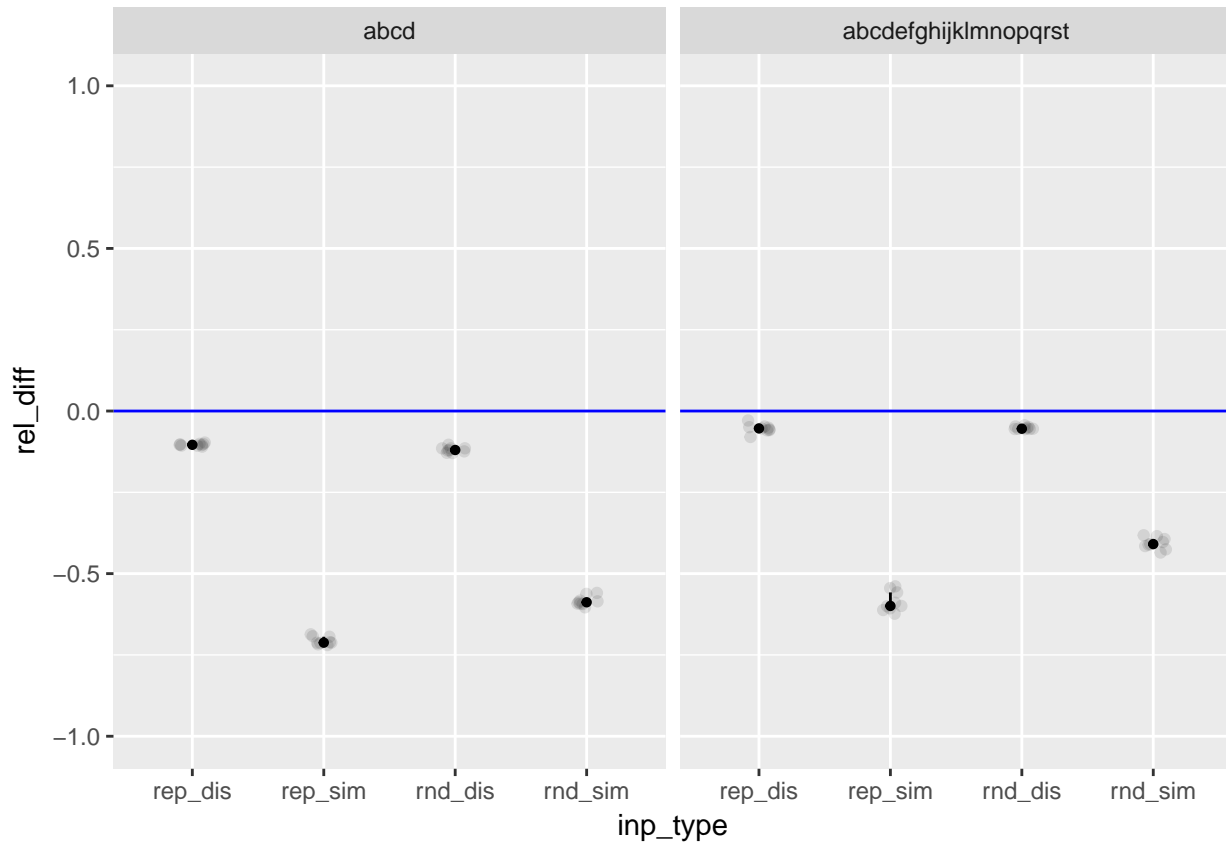
```
## # A tibble: 2 x 2
##   drank      a
##   <chr> <int>
## 1 drank    48
## 2 srnk     48
```

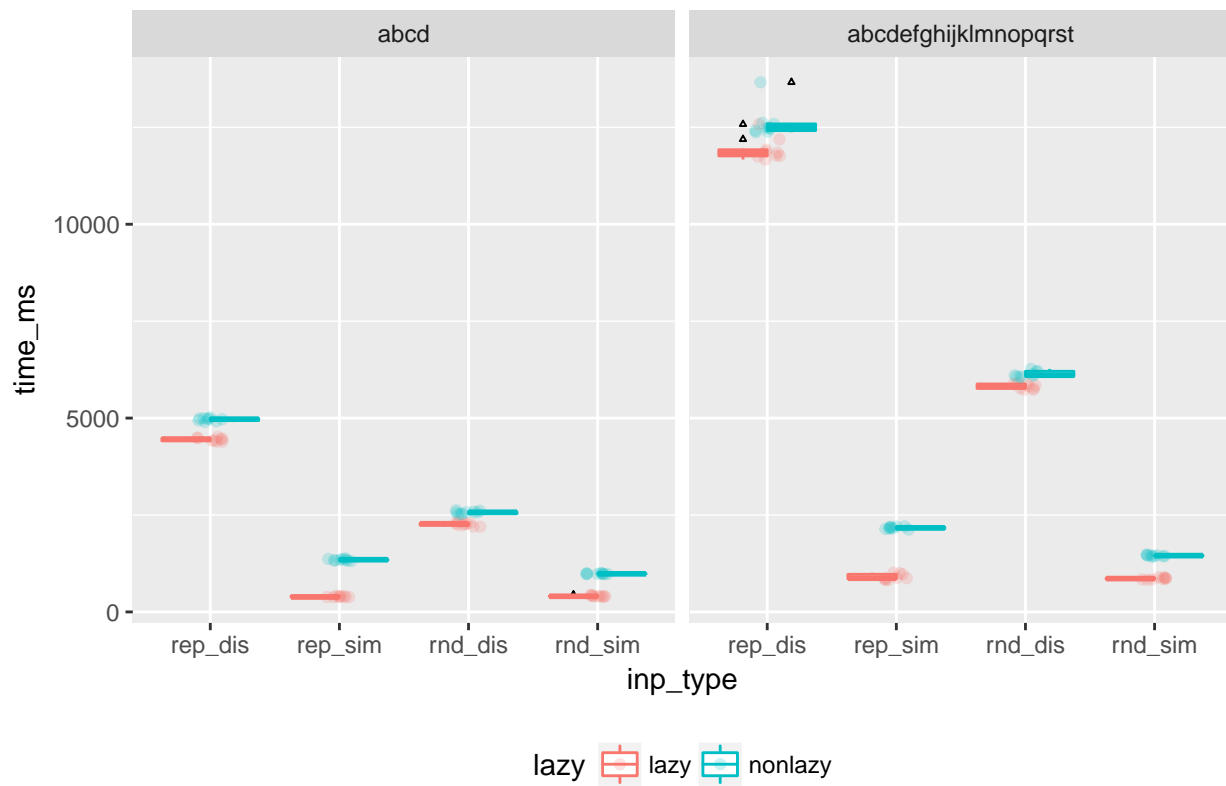



drank drank srank

3.2.2 Lazy

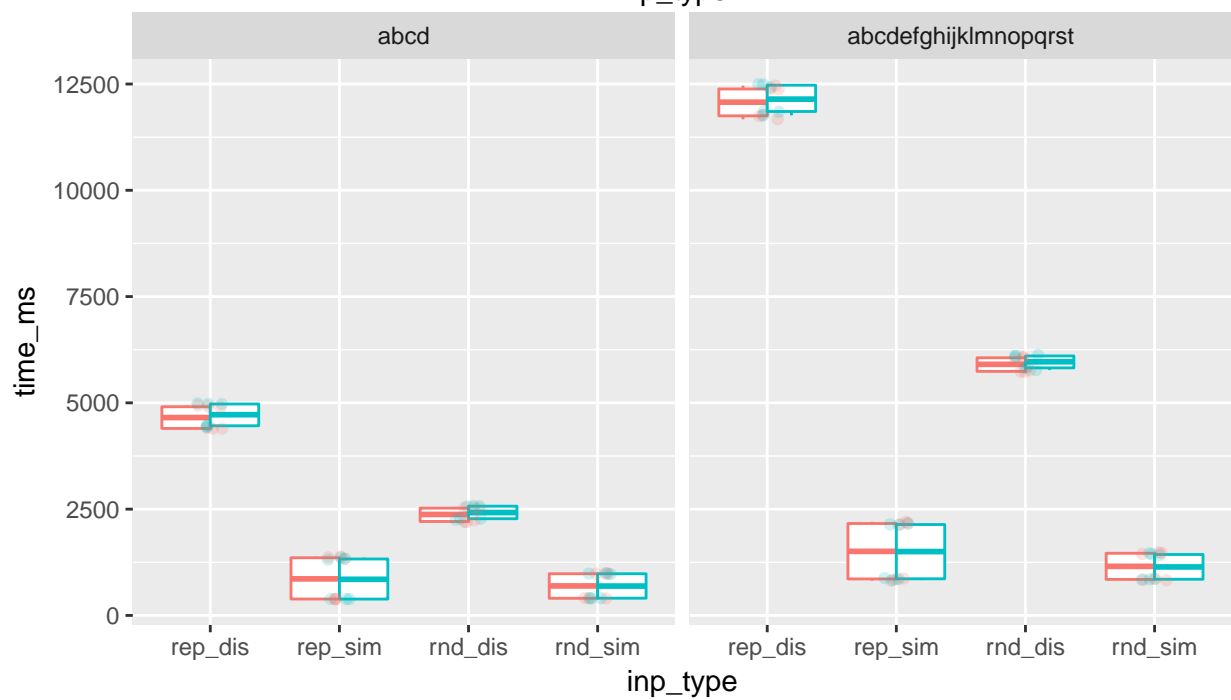
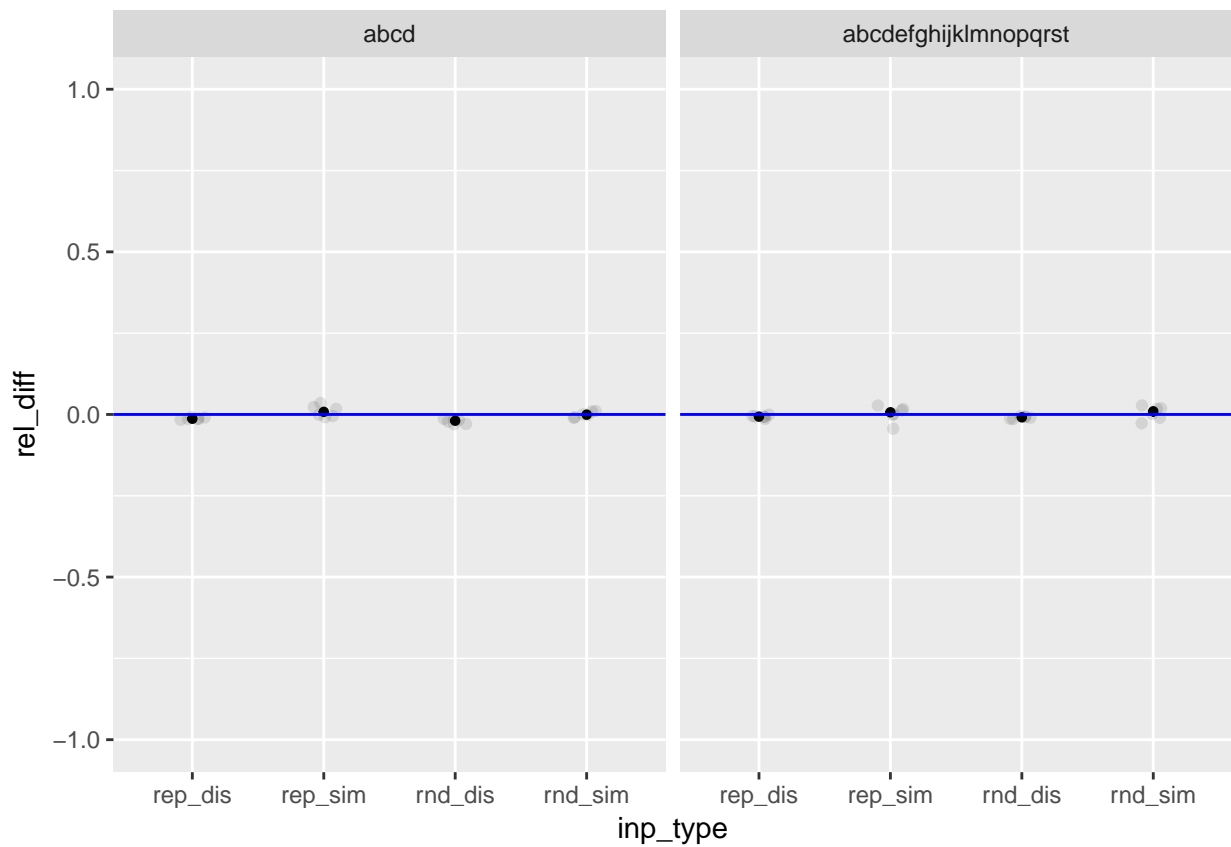
```
## # A tibble: 2 x 2
##   lazy      a
##   <chr> <int>
## 1 lazy    72
## 2 nonlazy 72
```





3.2.3 Fail

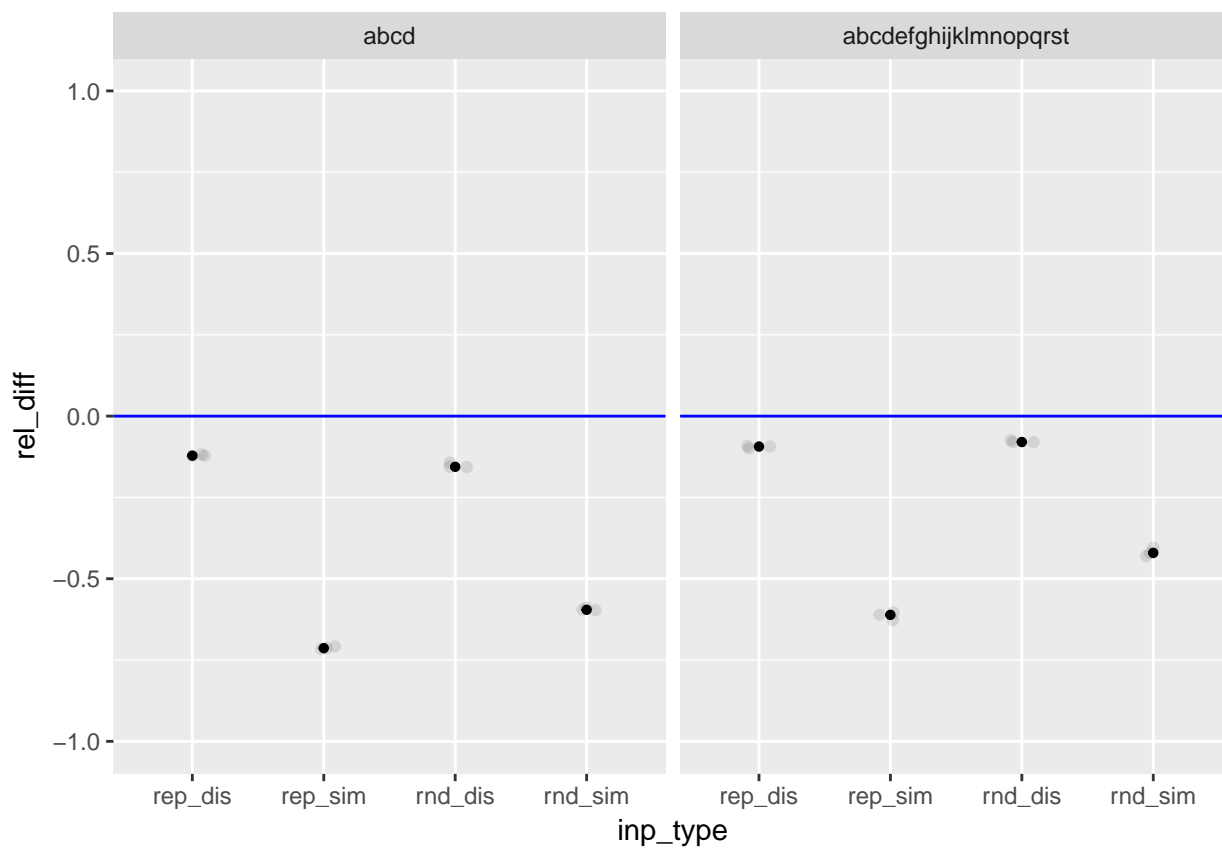
```
## # A tibble: 2 x 2
##   fail      a
##   <chr> <int>
## 1   fail    48
## 2 nonfail  48
```

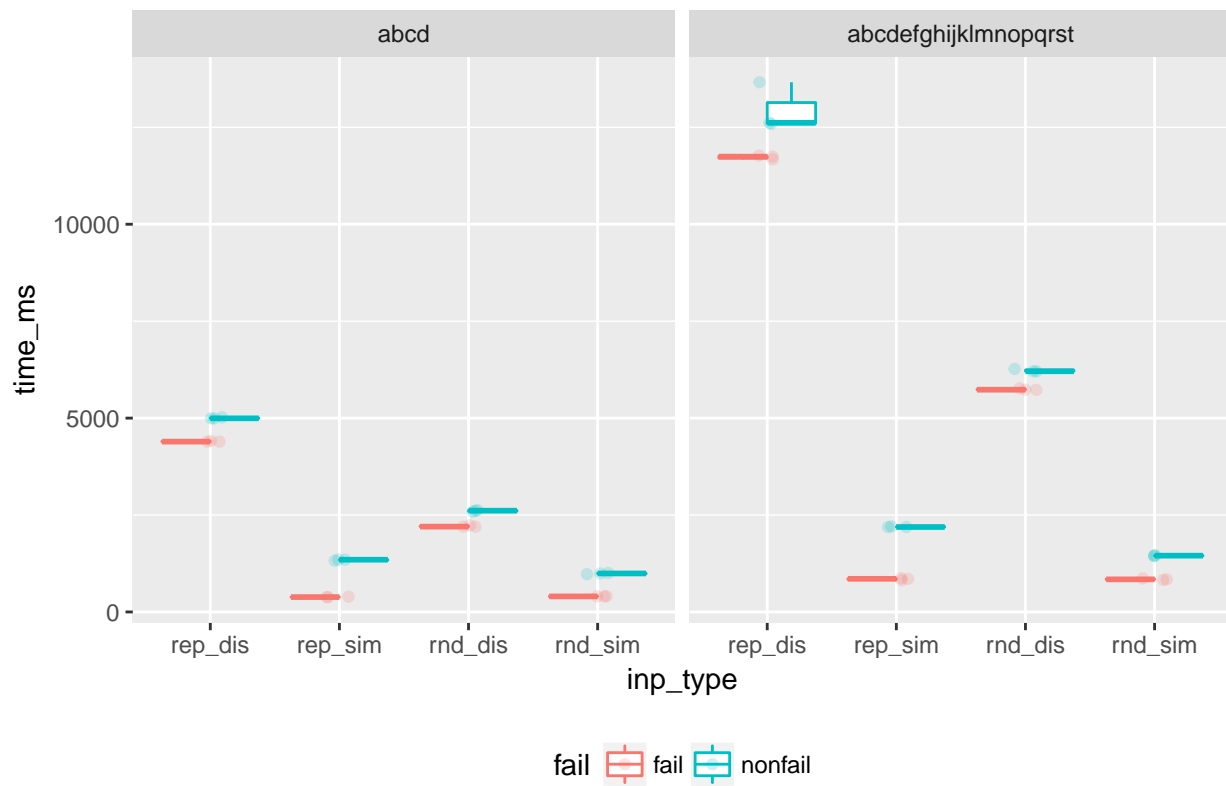


fail fail nonfail

3.2.4 all flags (lazy, fail, double_rank) vs. single rank

```
## # A tibble: 2 x 2
##   drank      a
##   <chr> <int>
## 1 drank    24
## 2 srnk     24
```

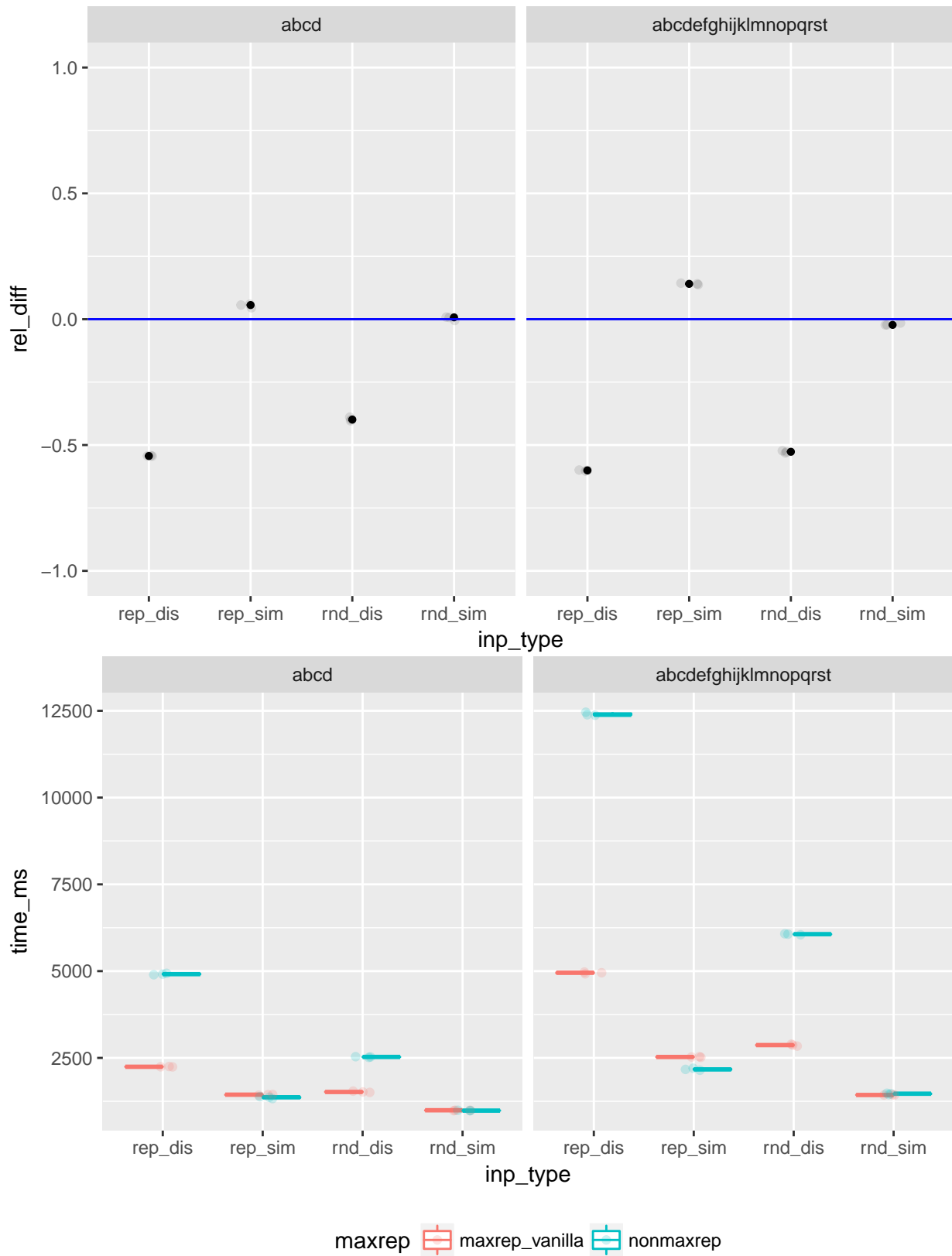




3.2.5 maxrep

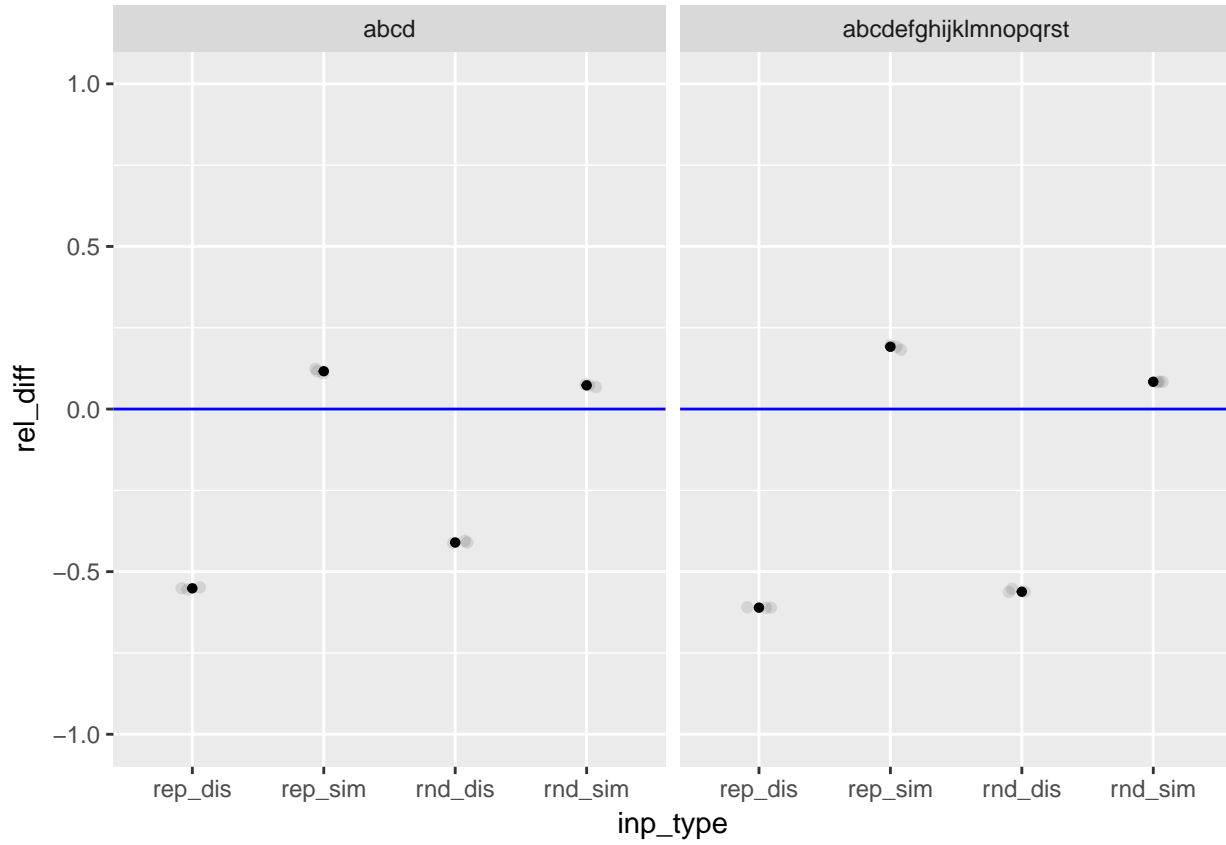
3.2.5.1 maxrep_vanilla vs. non-maxrep

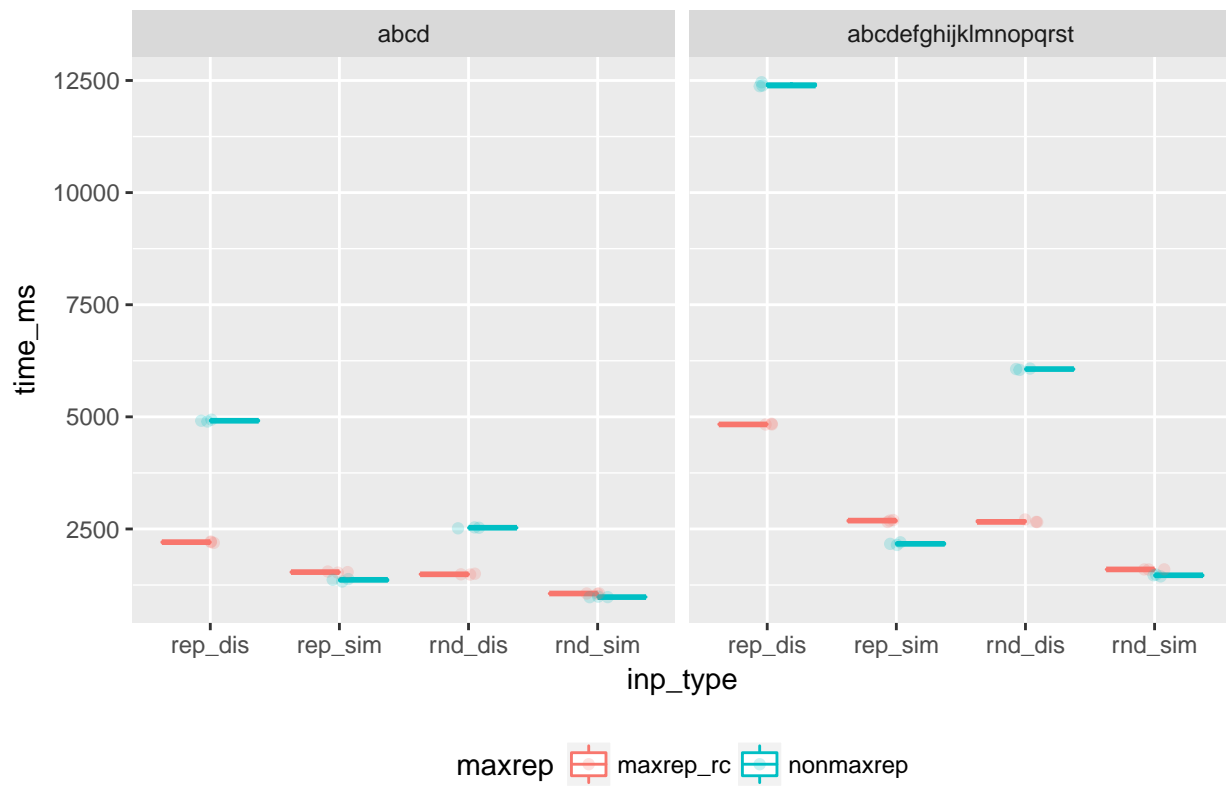
```
## # A tibble: 2 x 2
##       maxrep      a
##       <chr> <int>
## 1 maxrep_vanilla    24
## 2      nonmaxrep    24
```



3.2.5.2 maxrep+rank_and_check vs. non-maxrep

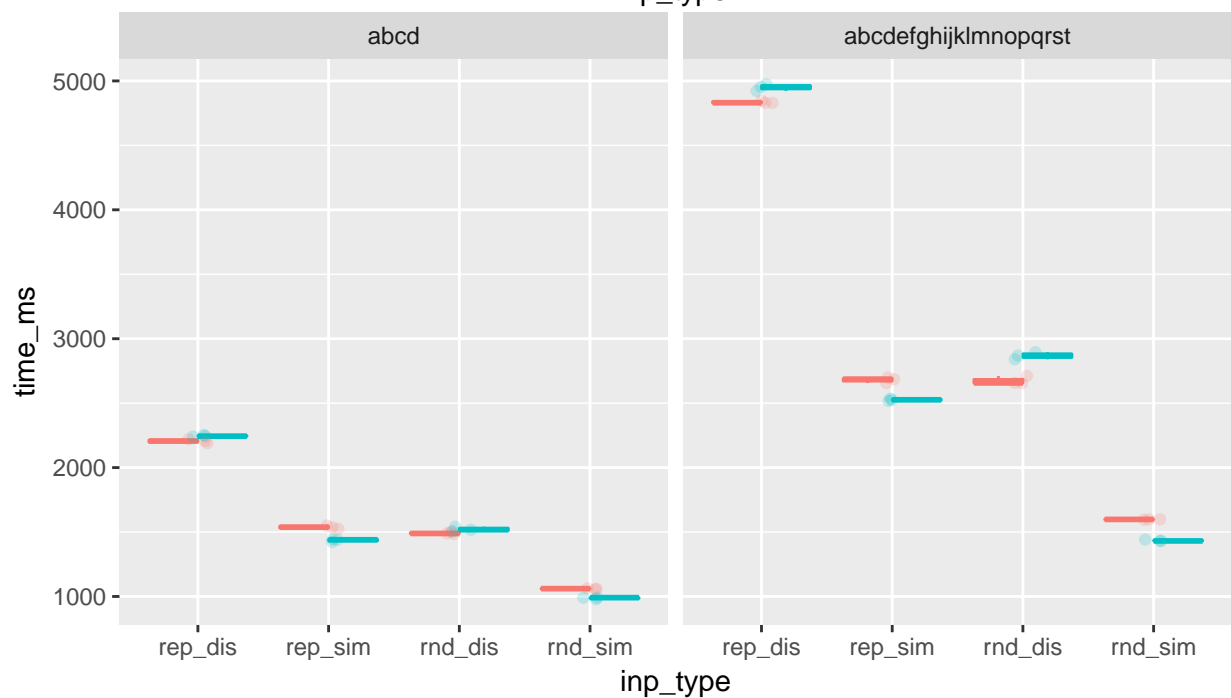
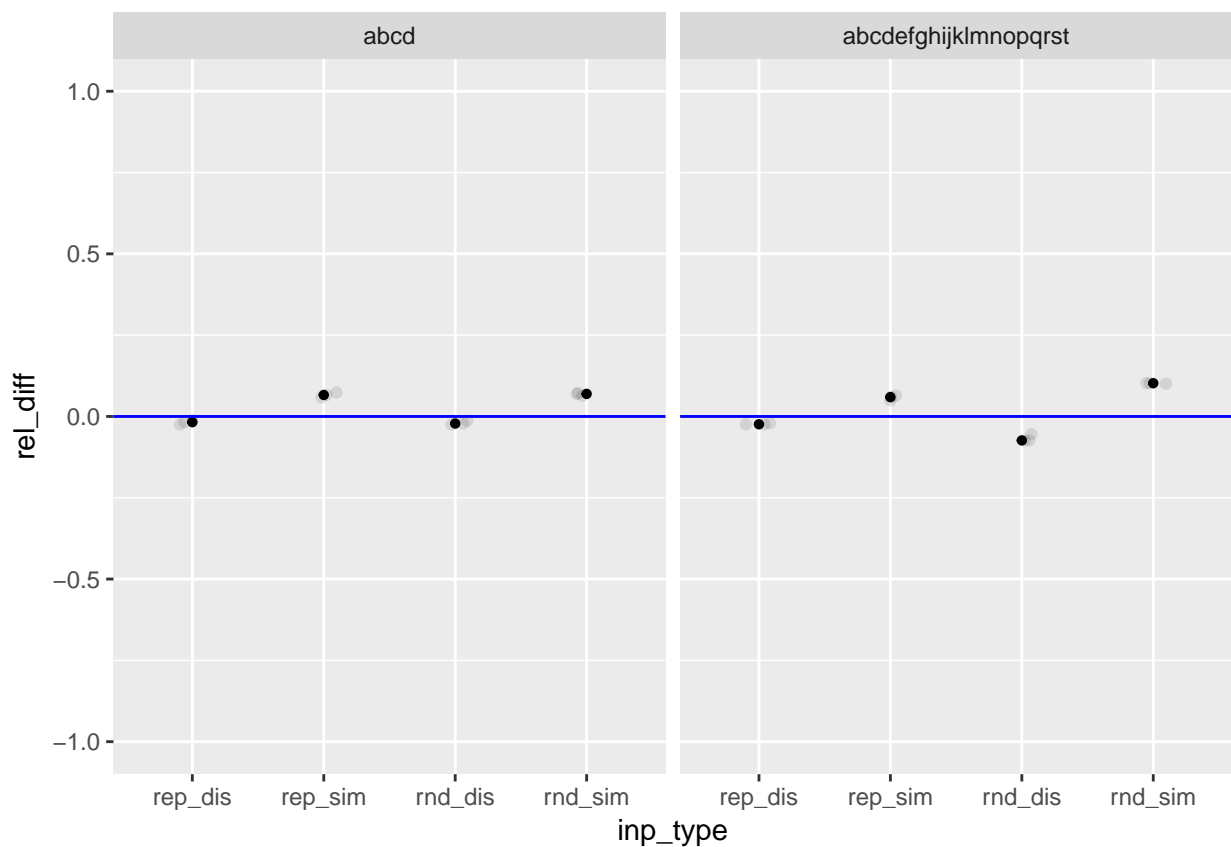
```
## # A tibble: 2 x 2
##   maxrep a
##   <chr> <int>
## 1 maxrep_rc 24
## 2 nonmaxrep 24
```







3.2.5.3 maxrep+rank_and_check vs. maxrep_vanilla

```
## # A tibble: 2 x 2
##       maxrep      a
##       <chr> <int>
## 1 maxrep_rc    24
## 2 maxrep_vanilla 24
```



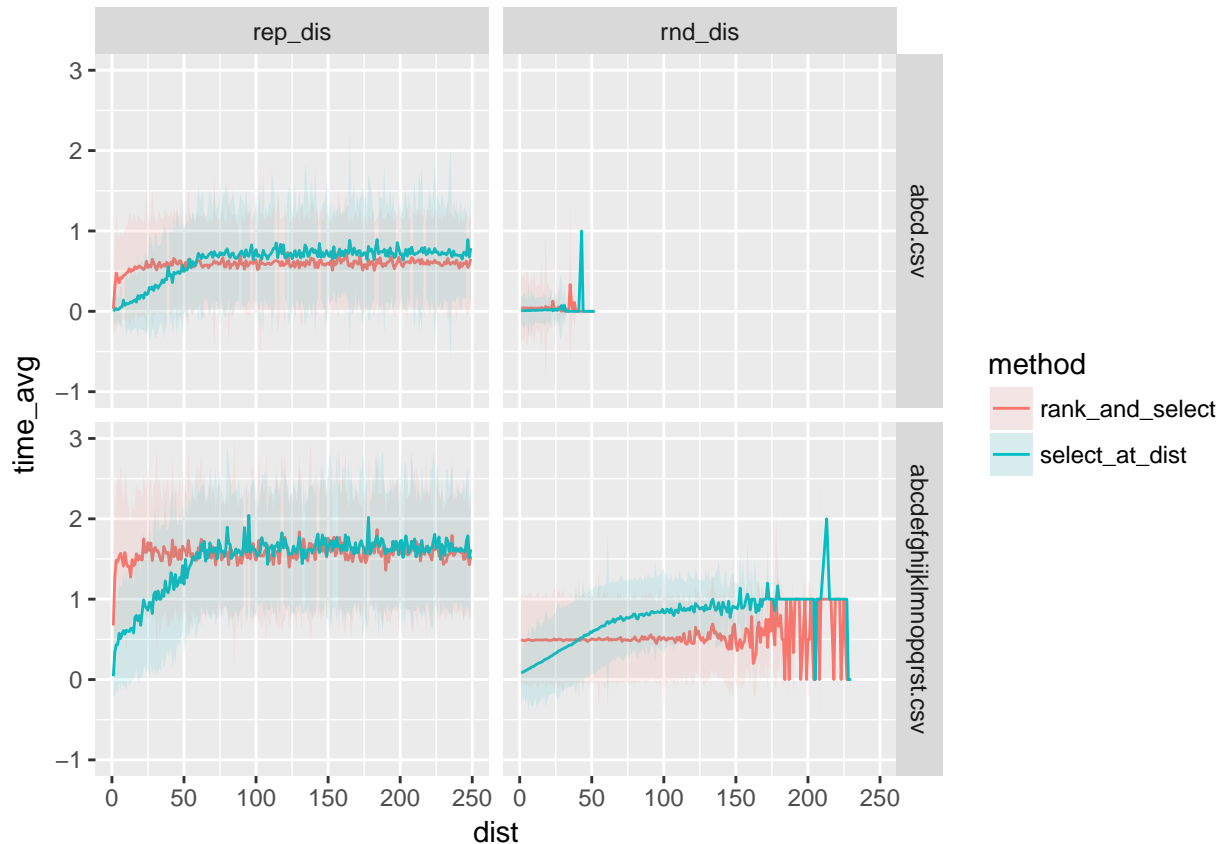
maxrep  maxrep_rc  maxrep_vanilla

4 parent optimizations

4.1 sandbox tests

4.1.1 select at dist

Traverse the tree and time the call `select_at_dist` (labeled f) and the call `select(rank())` (labeled s) on each node.



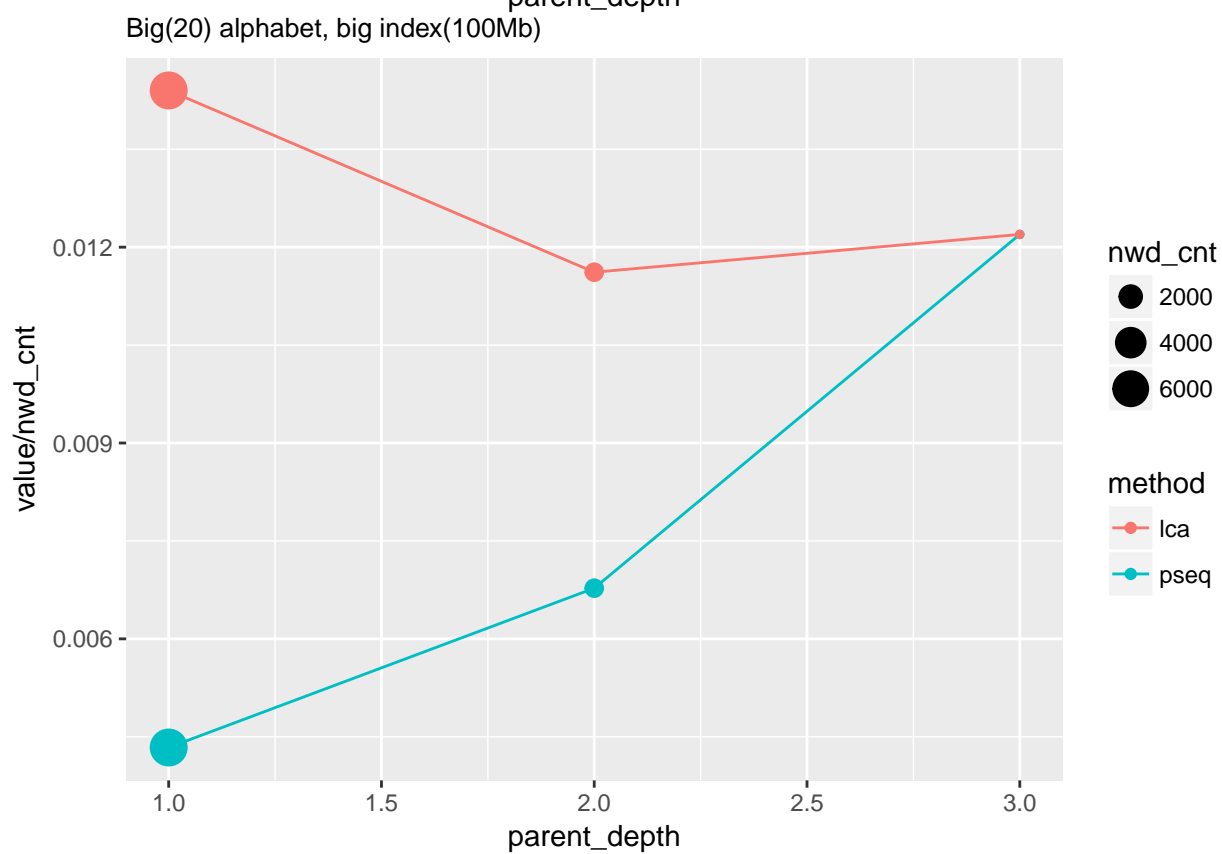
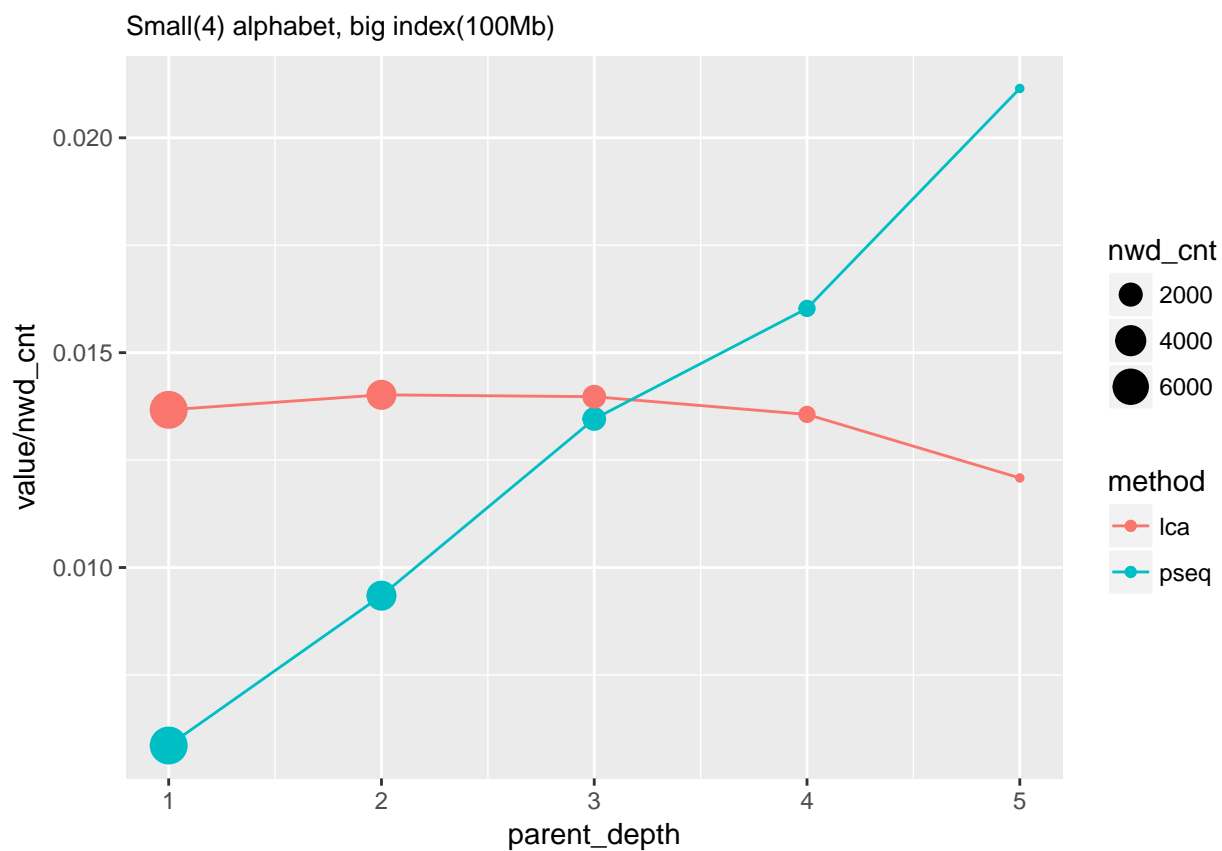
4.1.2 lowest maximal ancestor

Generate all tuples (v, c, d) from a tree where, v is an internal node, c a character, and $d = \text{depth}(v) - \text{depth}(u)$ with u the lowest maximal ancestor.

Then shuffle and time the two ways of finding u : using the maxrep or with a sequence of parent calls.

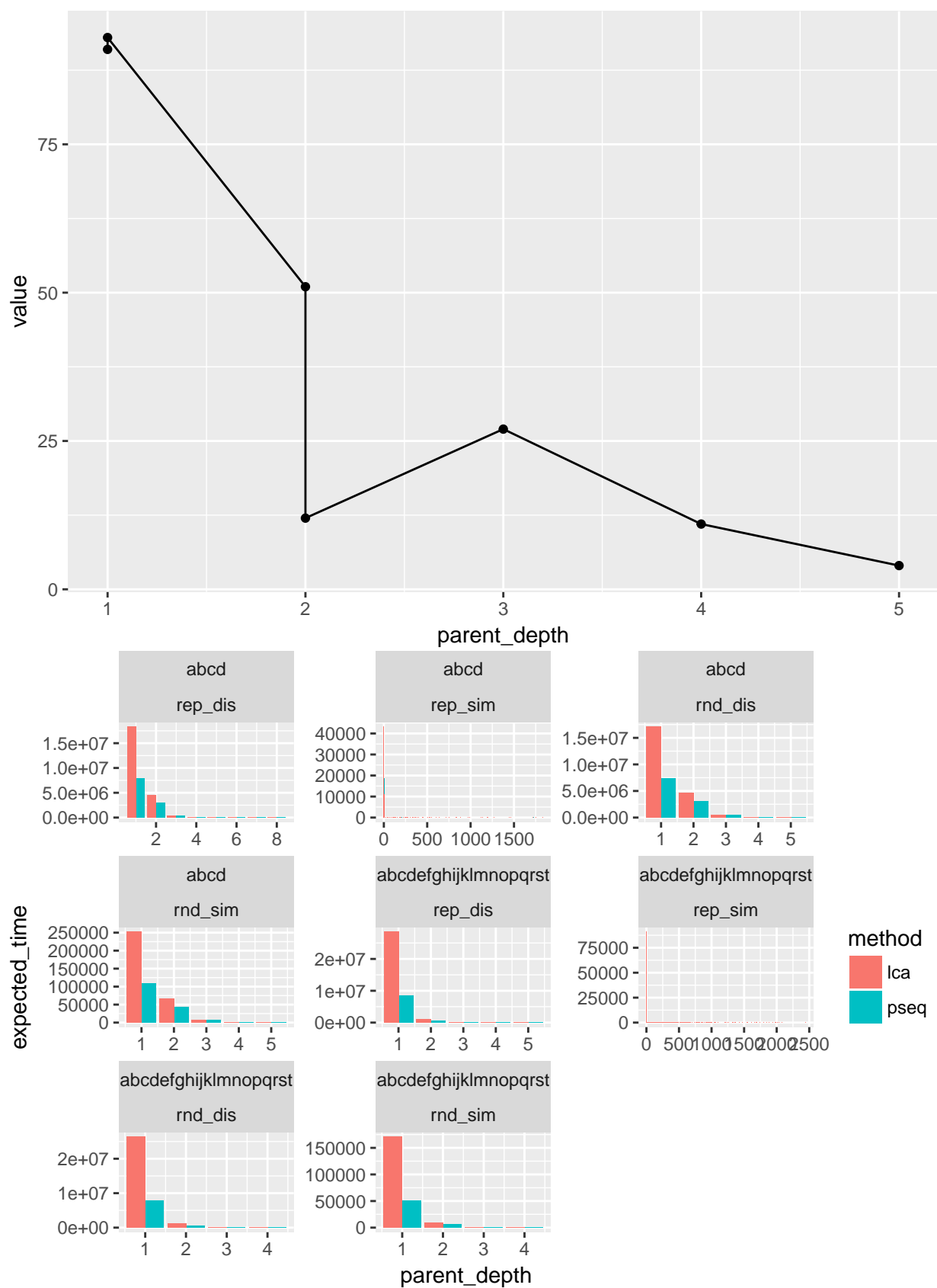
The plots show the `total time / nr of calls` for a given depth with the size of the point representing the `nr of calls` for a given depth value.

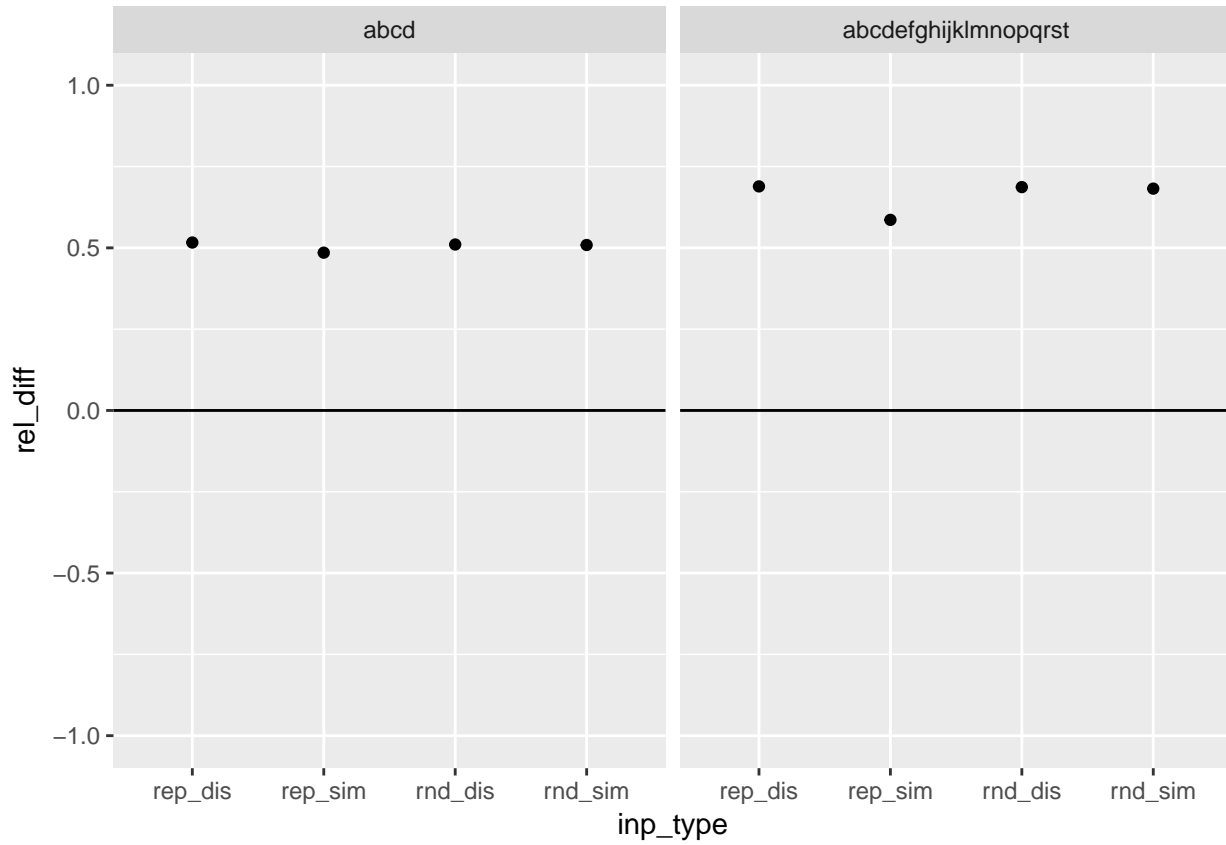
When the nr of calls is small (this happens for large depth lengths) the time resolution is not big enough and the numbers are not to be trusted – hence I have removed the data from the plots.



By counting the number of consecutive parent calls (for each number of consecutive calls) we project the

runtime of each method for an input type as below.





4.2 full tests

5 parallelization

Run the program on 1, 2, 4, 8 and 16 threads and measure the time it takes to build the RUNS and MS vectors.

TODO: there is currently a bug on the parallel version of the program for particular inputs. At the moment fixing it is not a priority.