

# ATLAS 2018 tape recalls analysis

German Cancio [german.cancio@cern.ch](mailto:german.cancio@cern.ch) (<mailto:german.cancio@cern.ch>)

This workbook contains a quick analysis of recalled files on the ATLAS stager between January and December 2018, across both service classes (“default” and “t0atlas”).

## Some basic statistics and histograms

### General:

- Total files recalled: 4196637
- Total data volume recalled: 7862.9365475 TB
- Median file size: 2513.705165 MB, average: 1873.6279901 MB
- Median tape transfer speed (including positioning): 360.5690171 MB, average: 336.6568009 MB
- Median service speed (including positioning): 217.5440363 MB, average: 208.6465443 MB

### “default” svcclass:

- Total files recalled: 1868871
- Total data volume recalled: 3717.7541889 TB
- Median file size: 2471.02326 MB, average: 1989.3048738 MB
- Median tape transfer speed: 360.0106055 MB, average: 330.1087342 MB
- Median service speed (including positioning): 185.5000011 MB, average: 184.0809096 MB

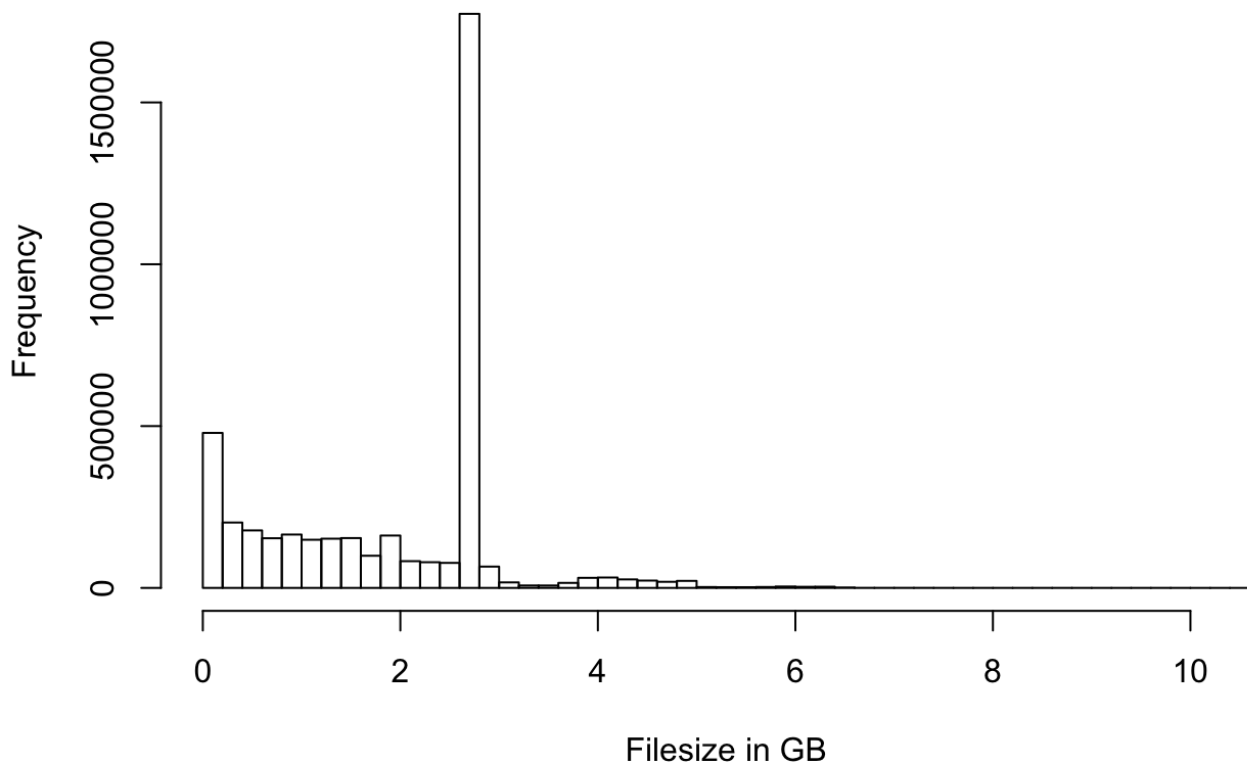
### “t0atlas” svcclass:

- Total files recalled: 2327766
- Total data volume recalled: 4145.1823586 TB
- Median file size: 2572.6981995 MB, average: 1780.755608 MB
- Median tape transfer speed: 360.7775845 MB, average: 341.9139839 MB
- Median service speed (including positioning): 243.8739069 MB, average: 228.369318 MB

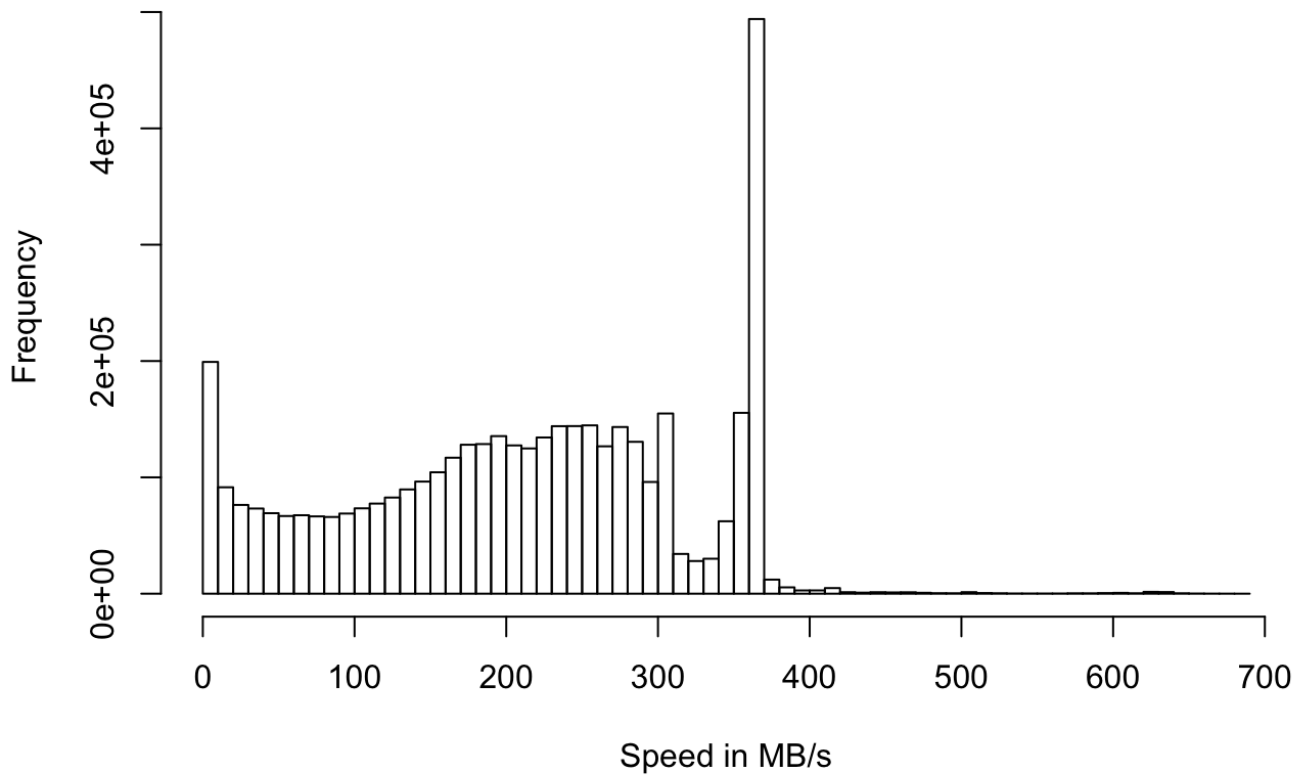
## Histograms of file size, service and transfer speed

- general
- by service class

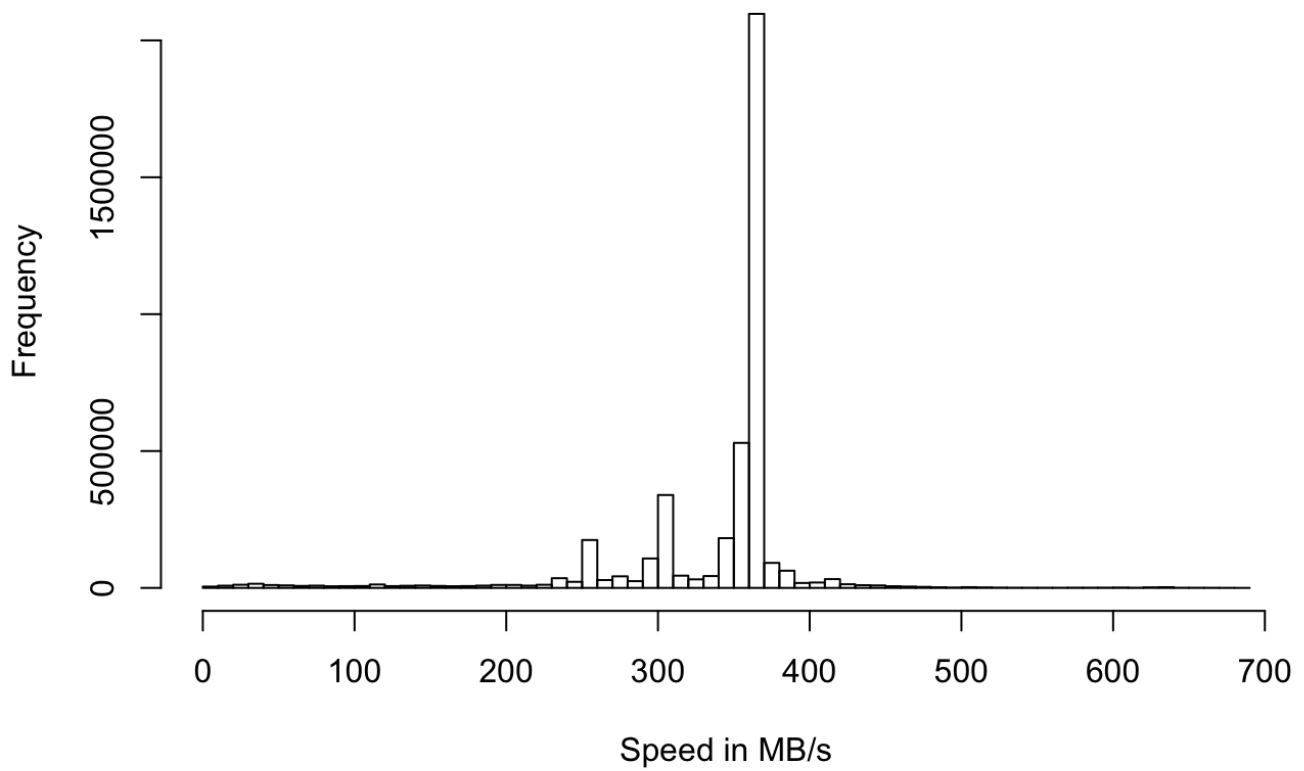
### Filesize



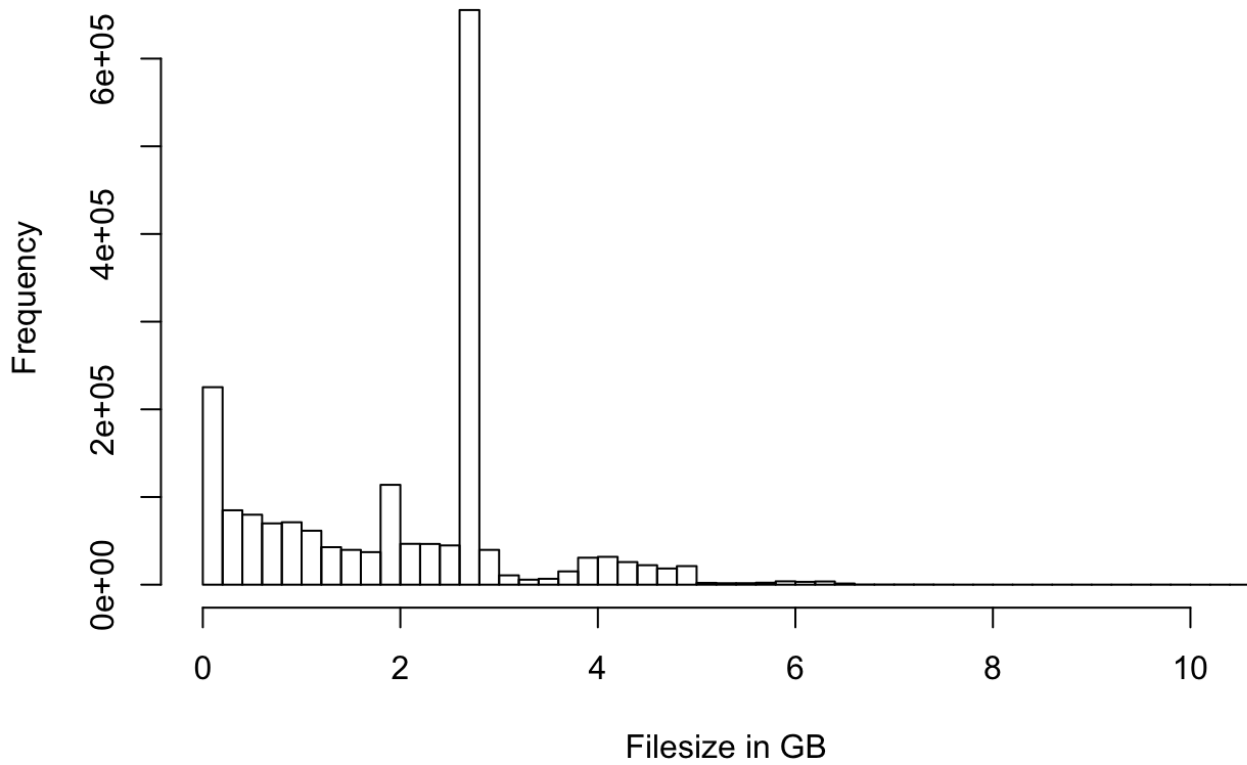
### Service speed



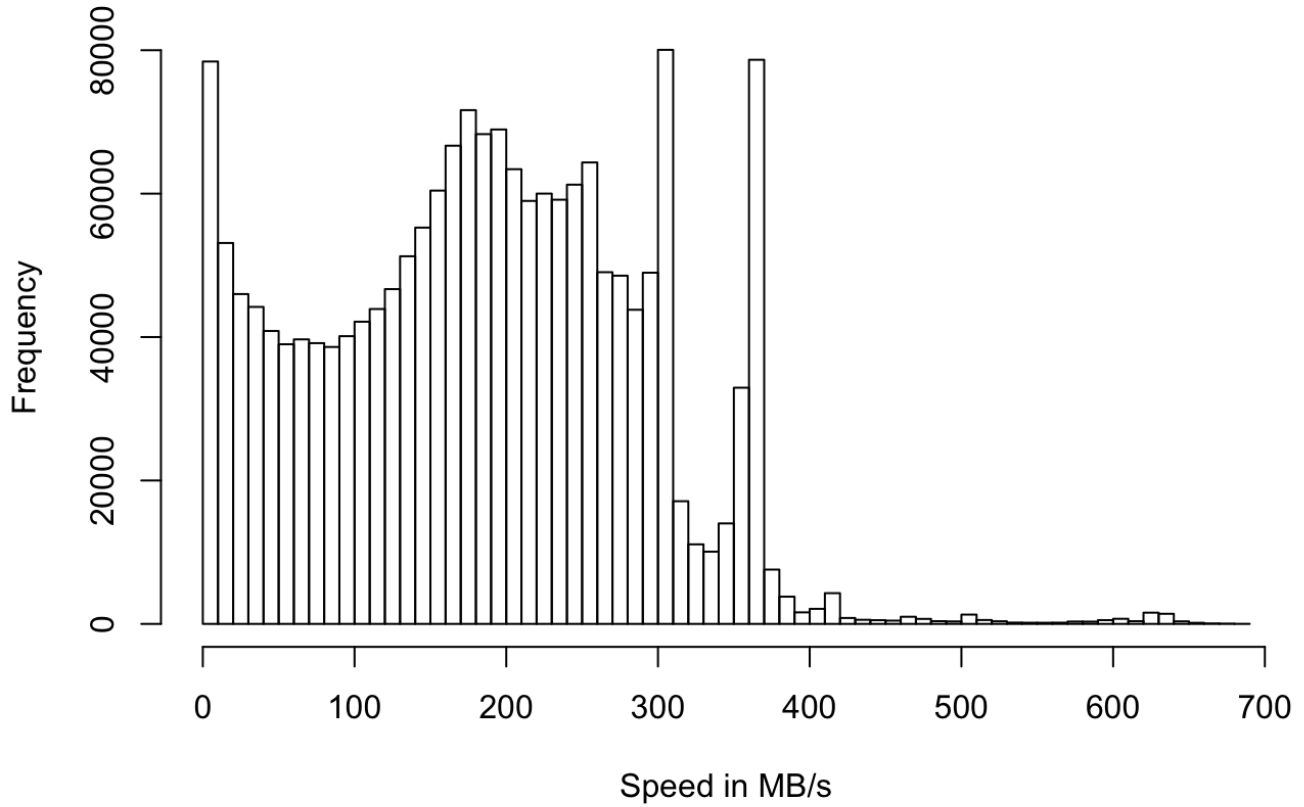
### Transfer speed



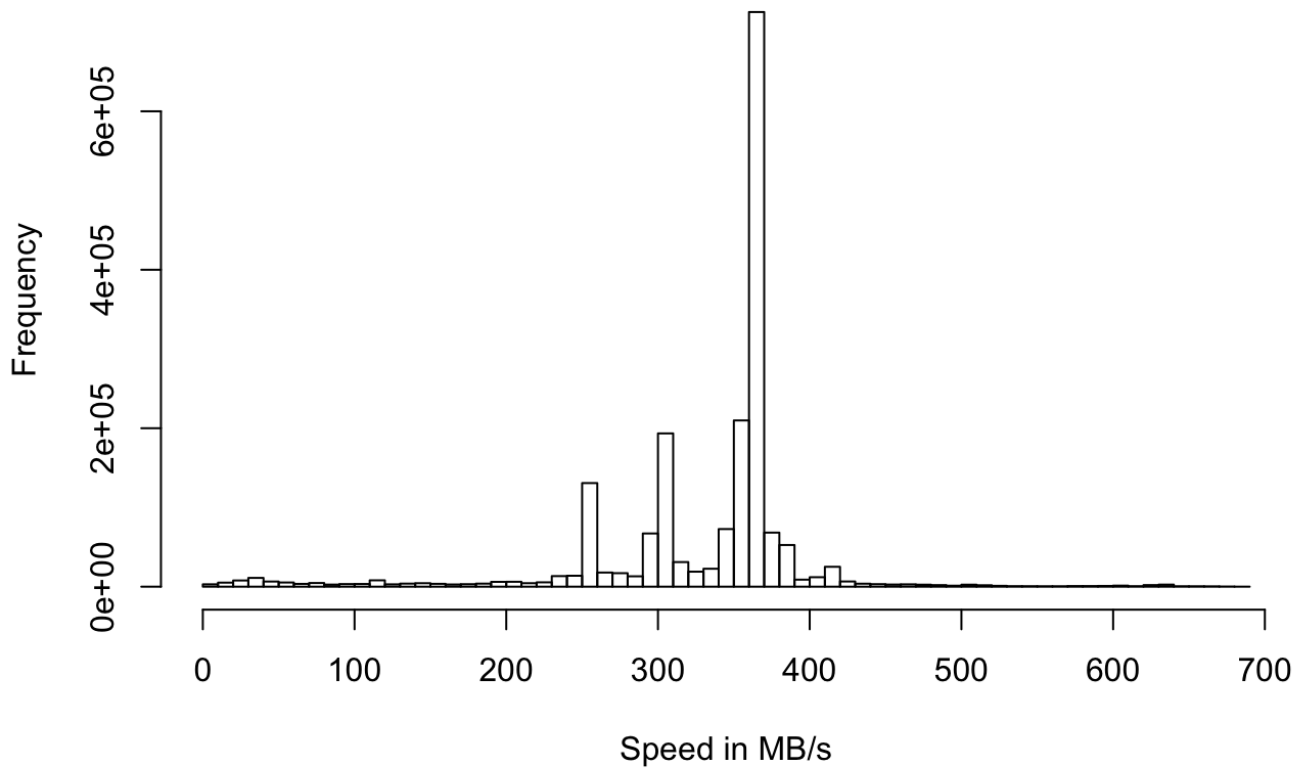
### Filesize - default svcclass



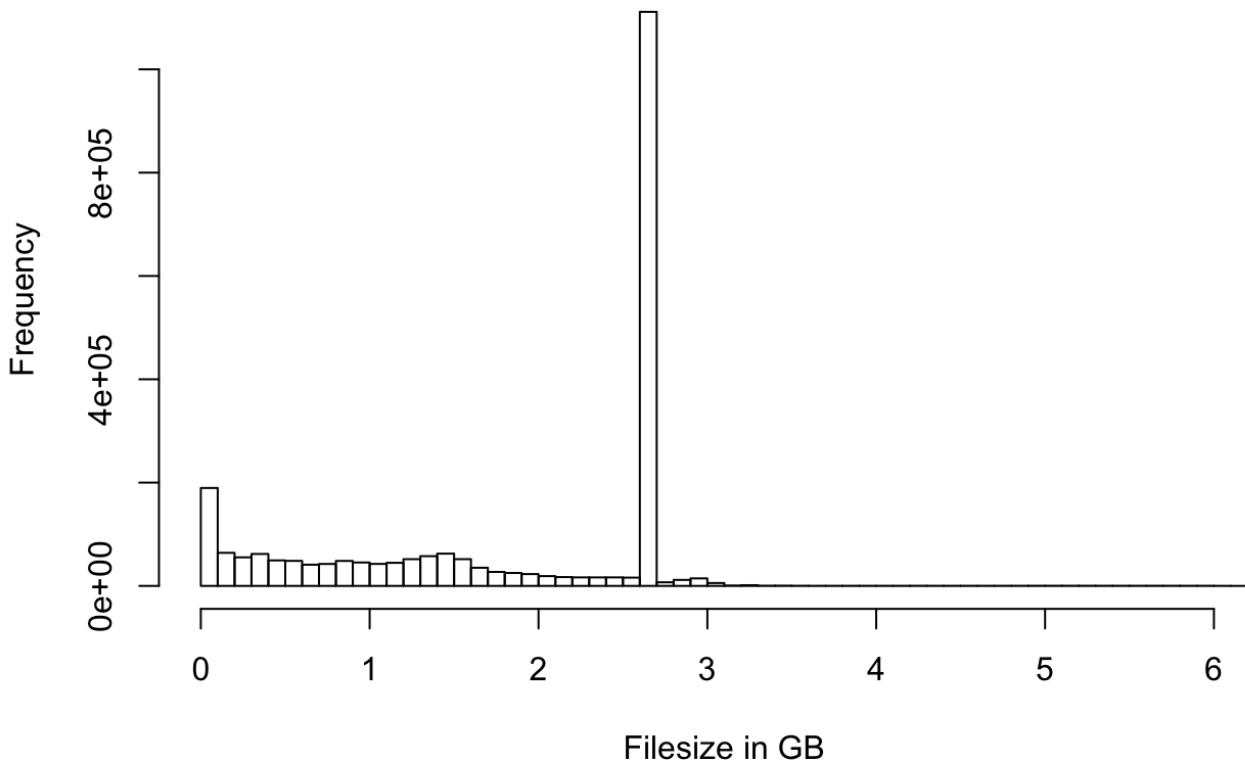
### Service speed - default svcclass



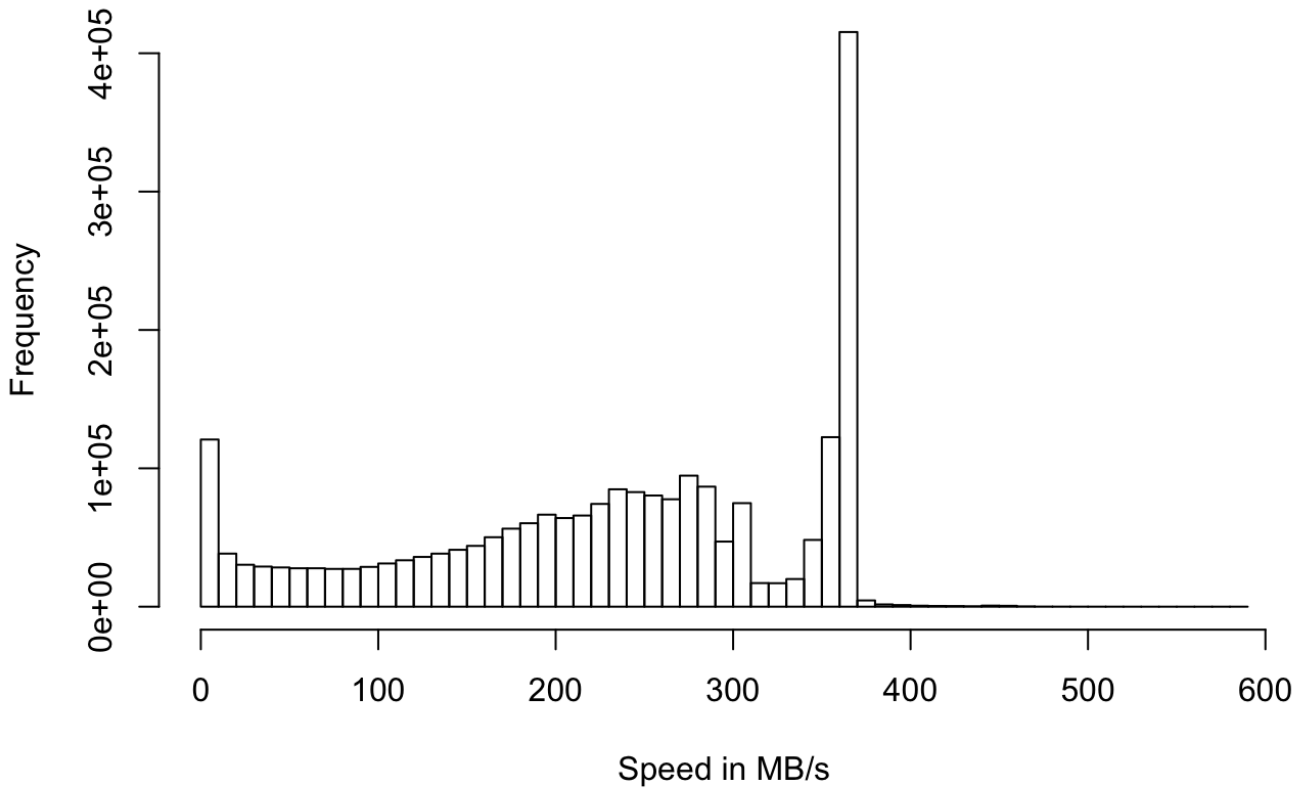
### Transfer speed - default svcclass



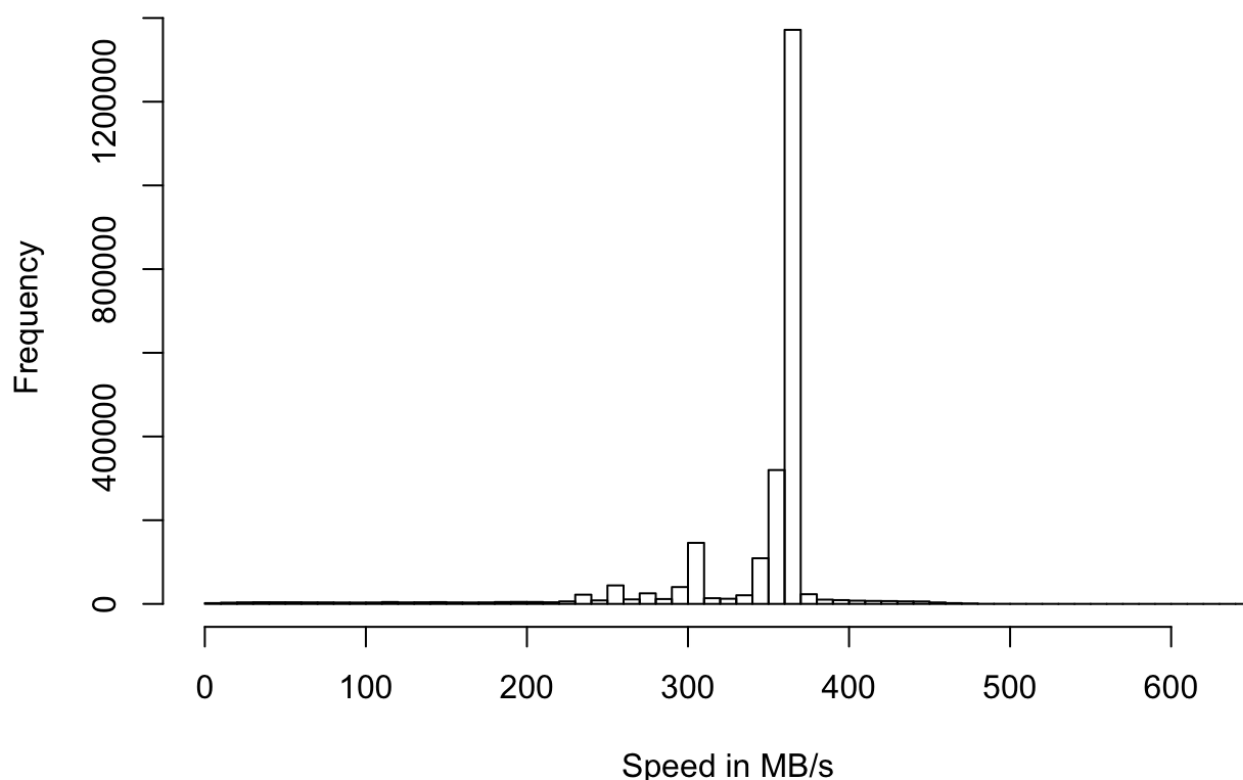
### Filesize - t0atlas svcclass



### Service speed - t0atlas svcclass



## Transfer speed - t0atlas svcclass



## Per-mount statistics

- Total tape read mounts: 81468
- Distinct tapes mounted: 7288 on 109 distinct drives.
- Average volume read per mount: 96.5156448 GB
- Average files read per mount: 51.5127044
- Number of tape mounts across multiple svcclasses: 325, percentage over total: 0.3989296 %

## default service class

- Number of tape mounts on default svcclass: 73798
- Distinct tapes mounted (default): 7126 on 109 distinct drives.
- Average volume read per mount: 50.3642247 GB
- Average files read per mount: 25.2823925

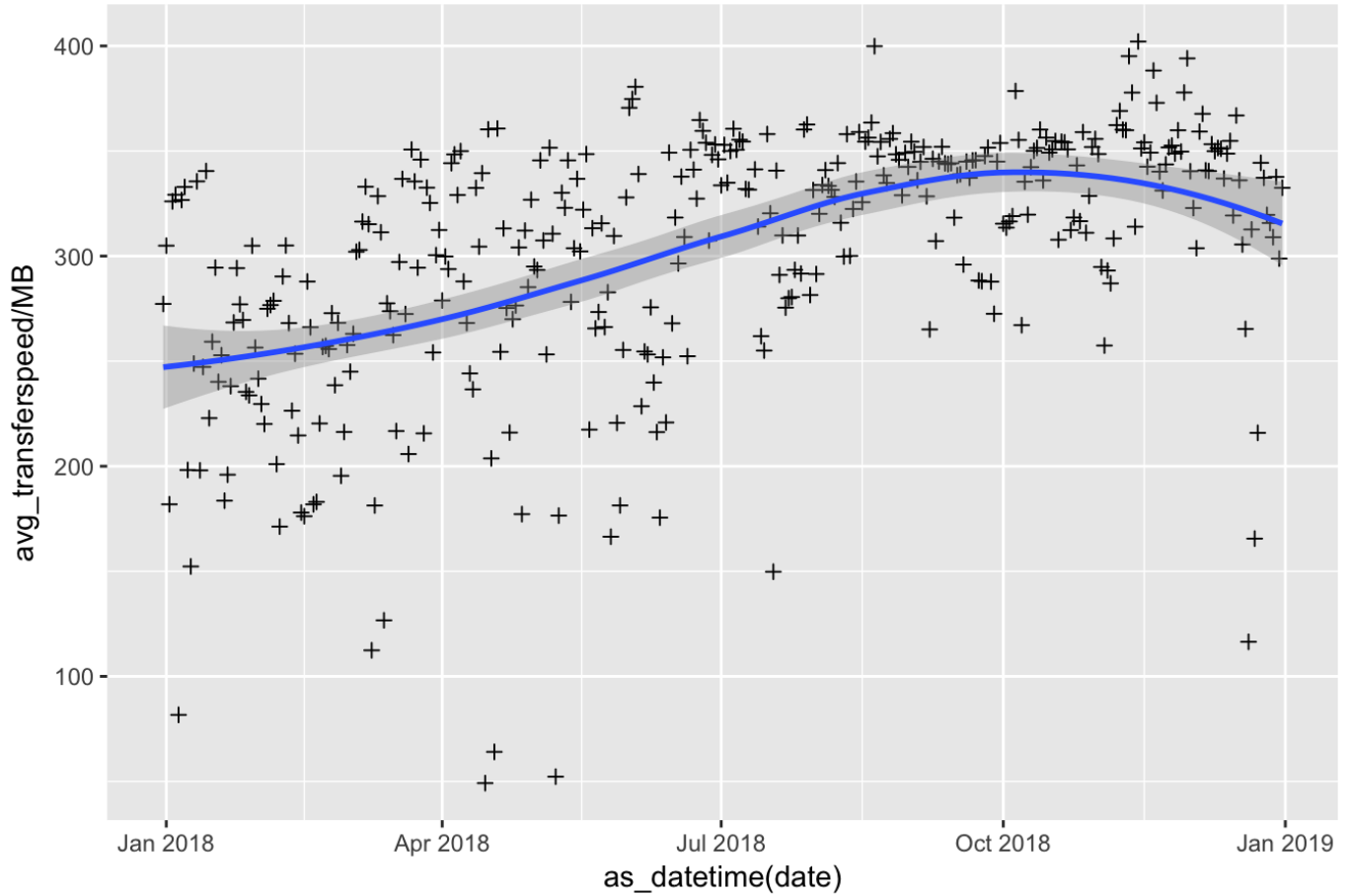
## t0atlas service class

- Number of tape mounts on t0atlas svcclass: 7670
- Distinct tapes mounted (t0atlas): 2449 on 108 distinct drives.
- Average volume read per mount: 540.5681214 GB
- Average files read per mount: 303.891395

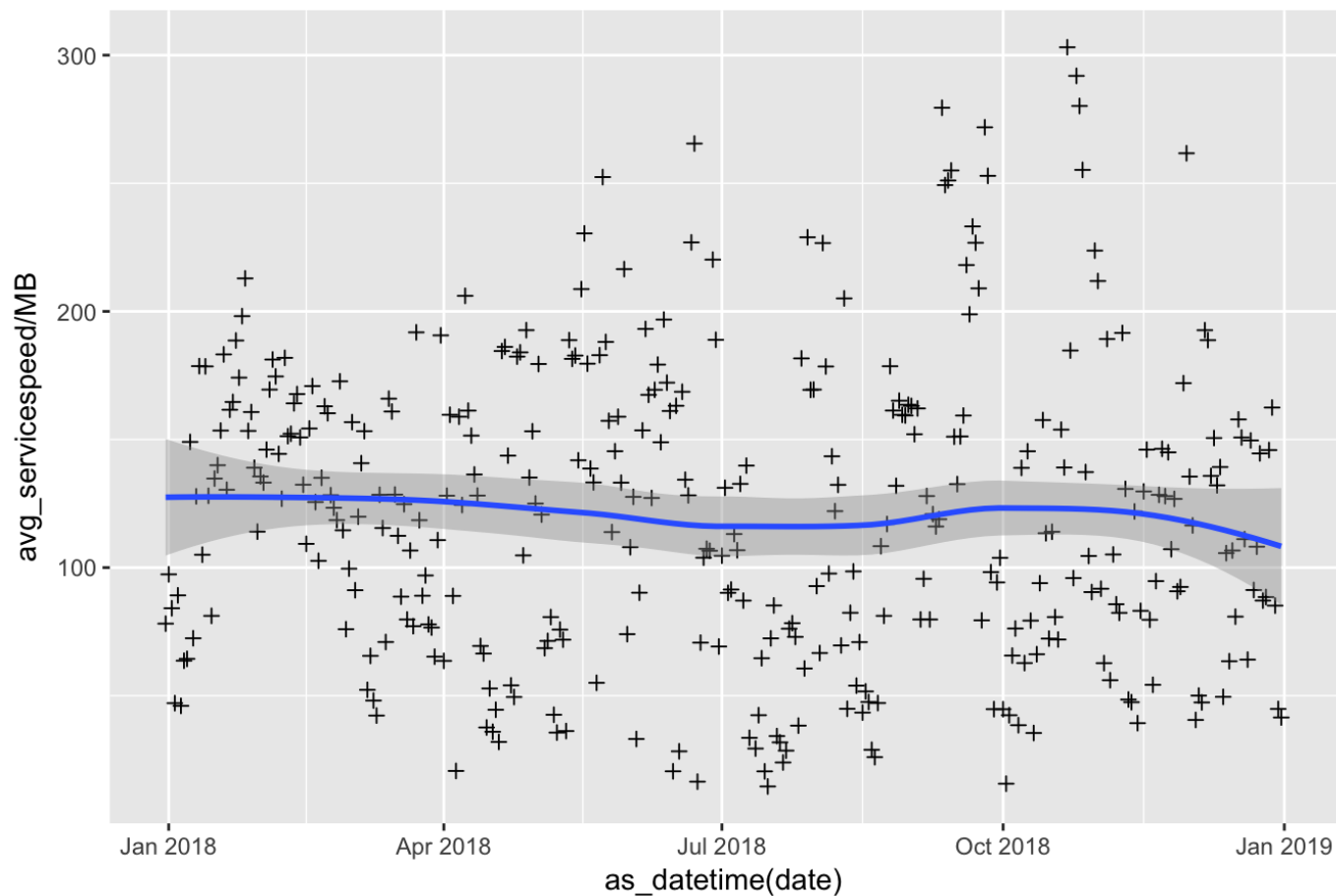
# Plots over time for daily transfer and service speeds

- Transfer speed: actual drive speed
- Service speed: includes positioning time

Average transfer speed, MB/s



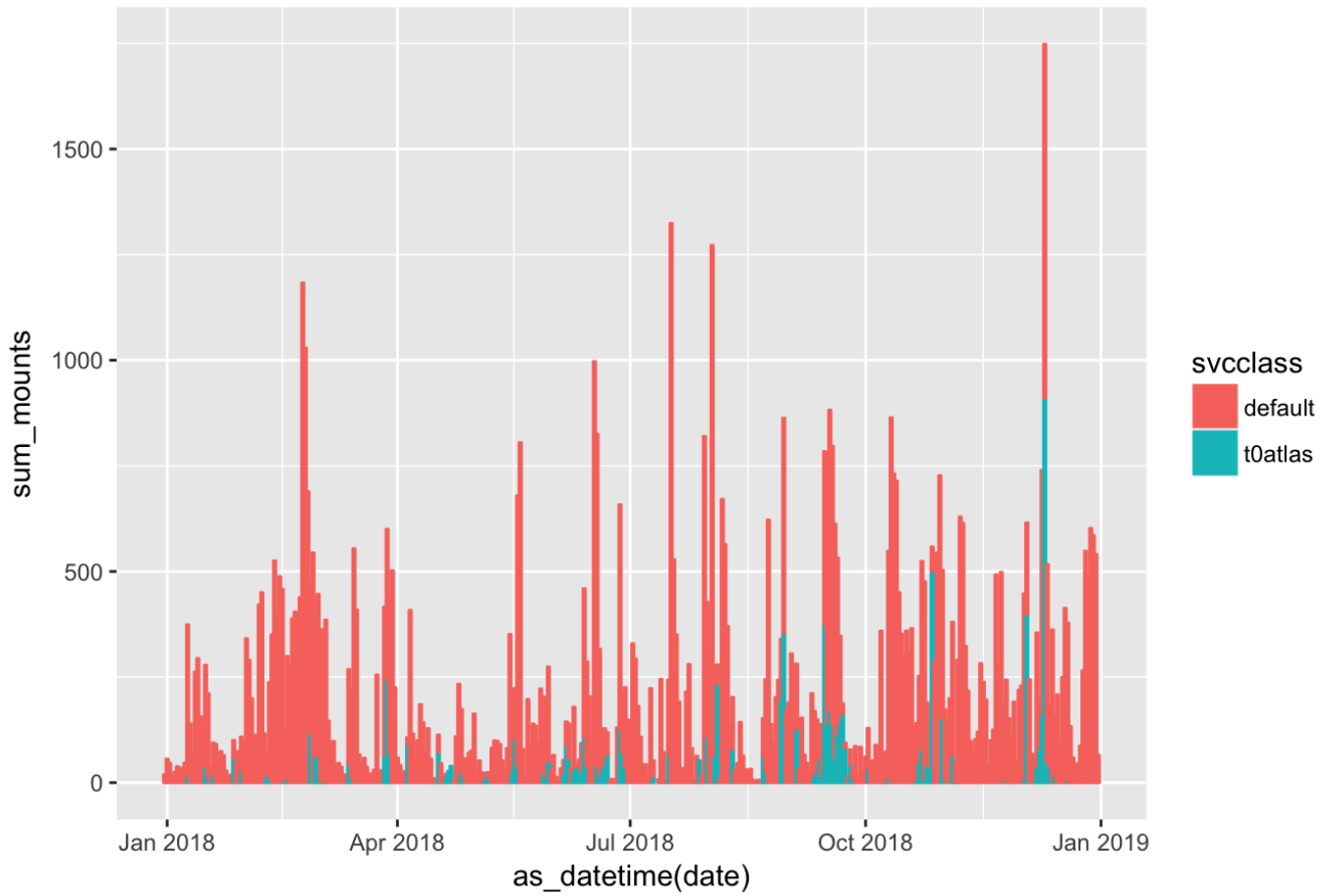
## Average service speed, MB/s



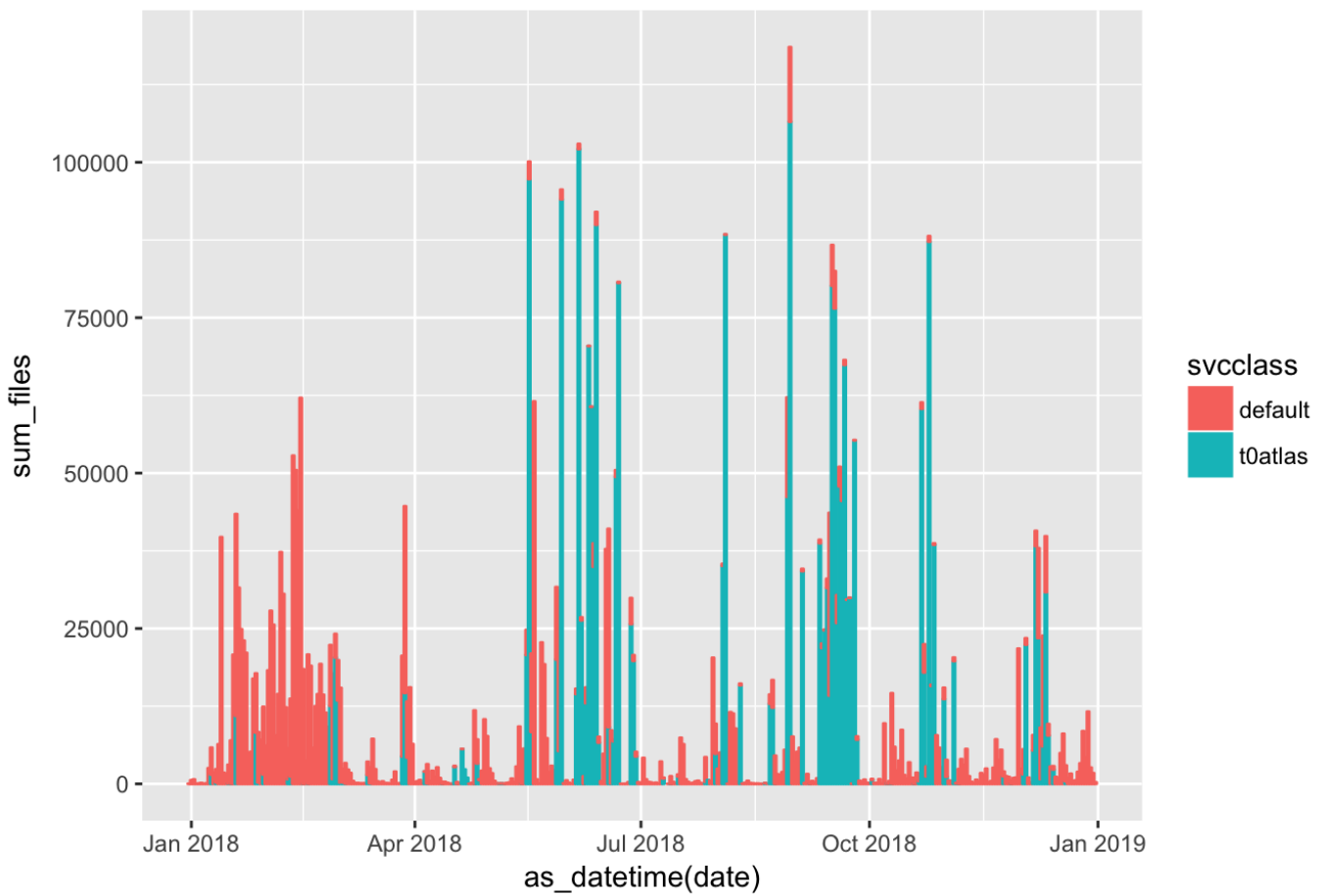
## Plots for per-svcclass daily mounts, files, data volume



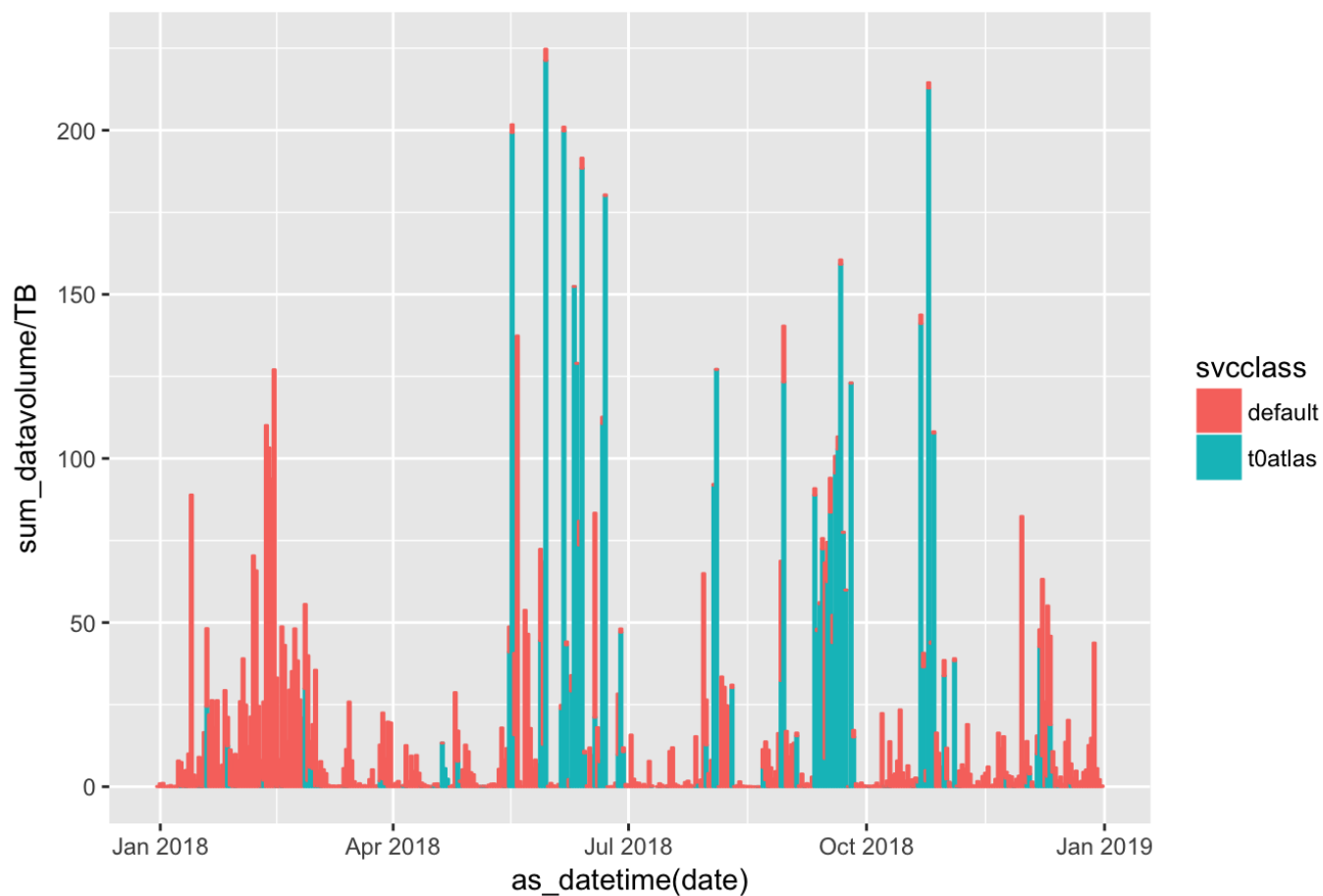
### daily mounts, stacked



### daily files read, stacked

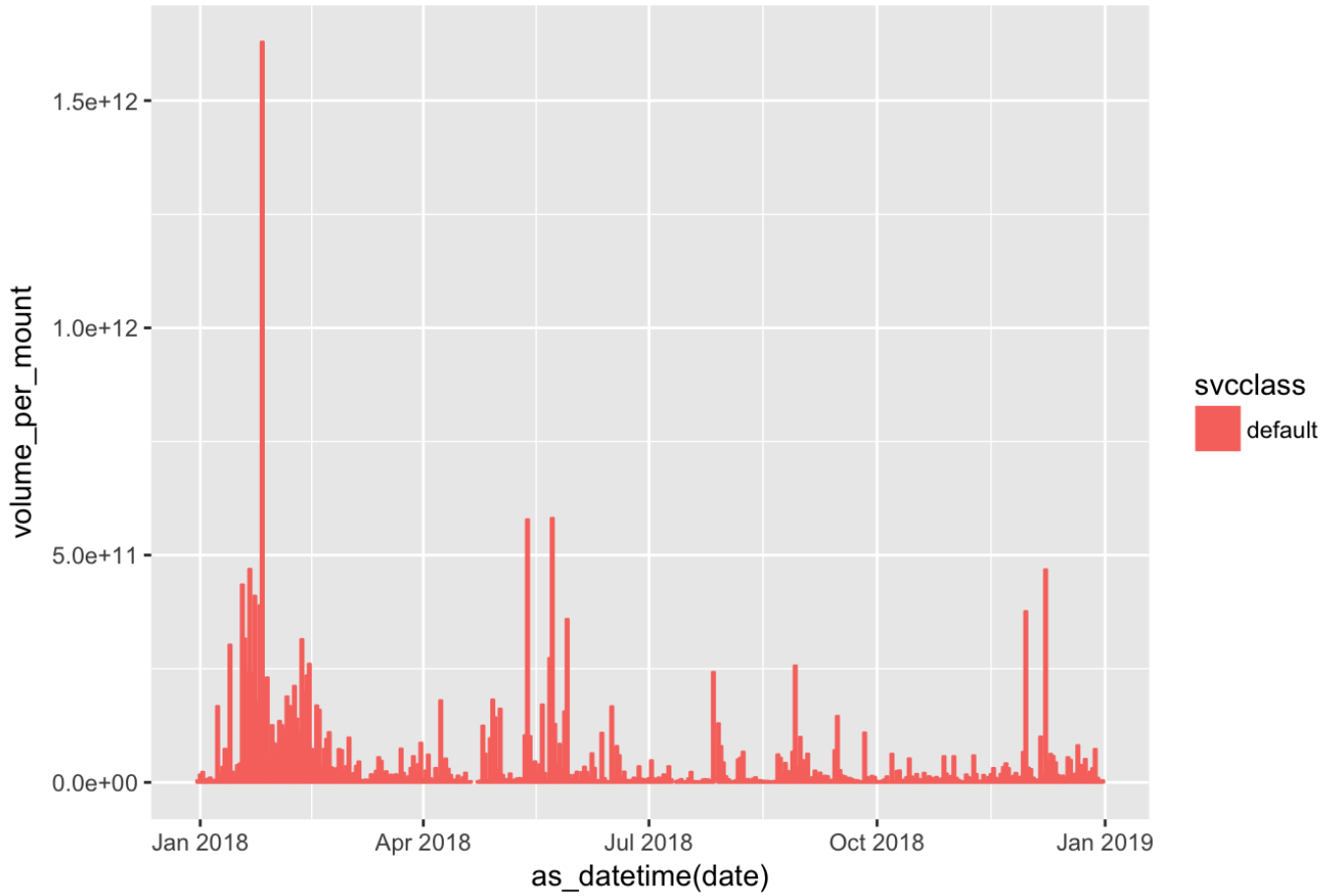


## daily volume read, in TB, stacked

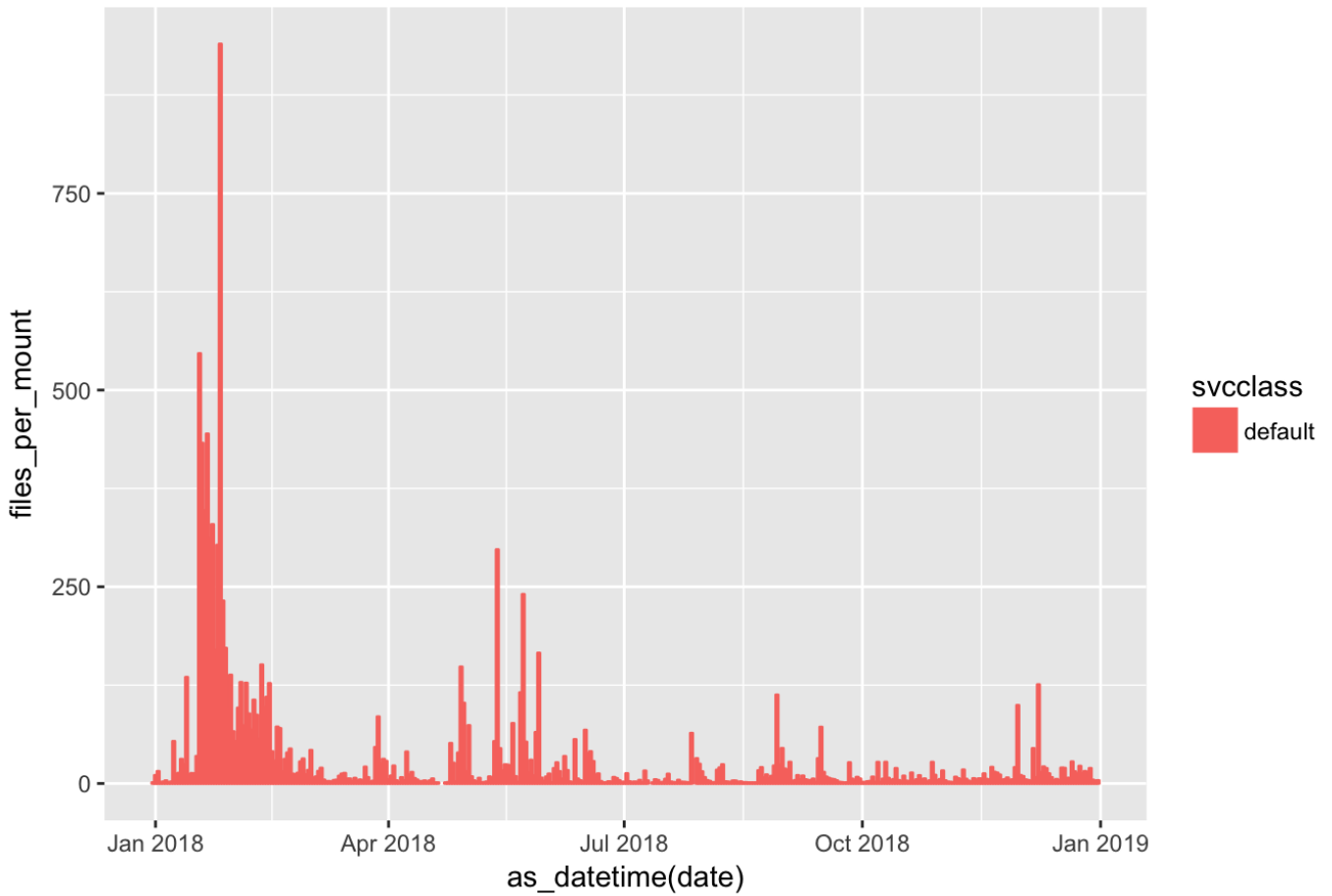


## per-svcclass volume per mount and files per mount evolution

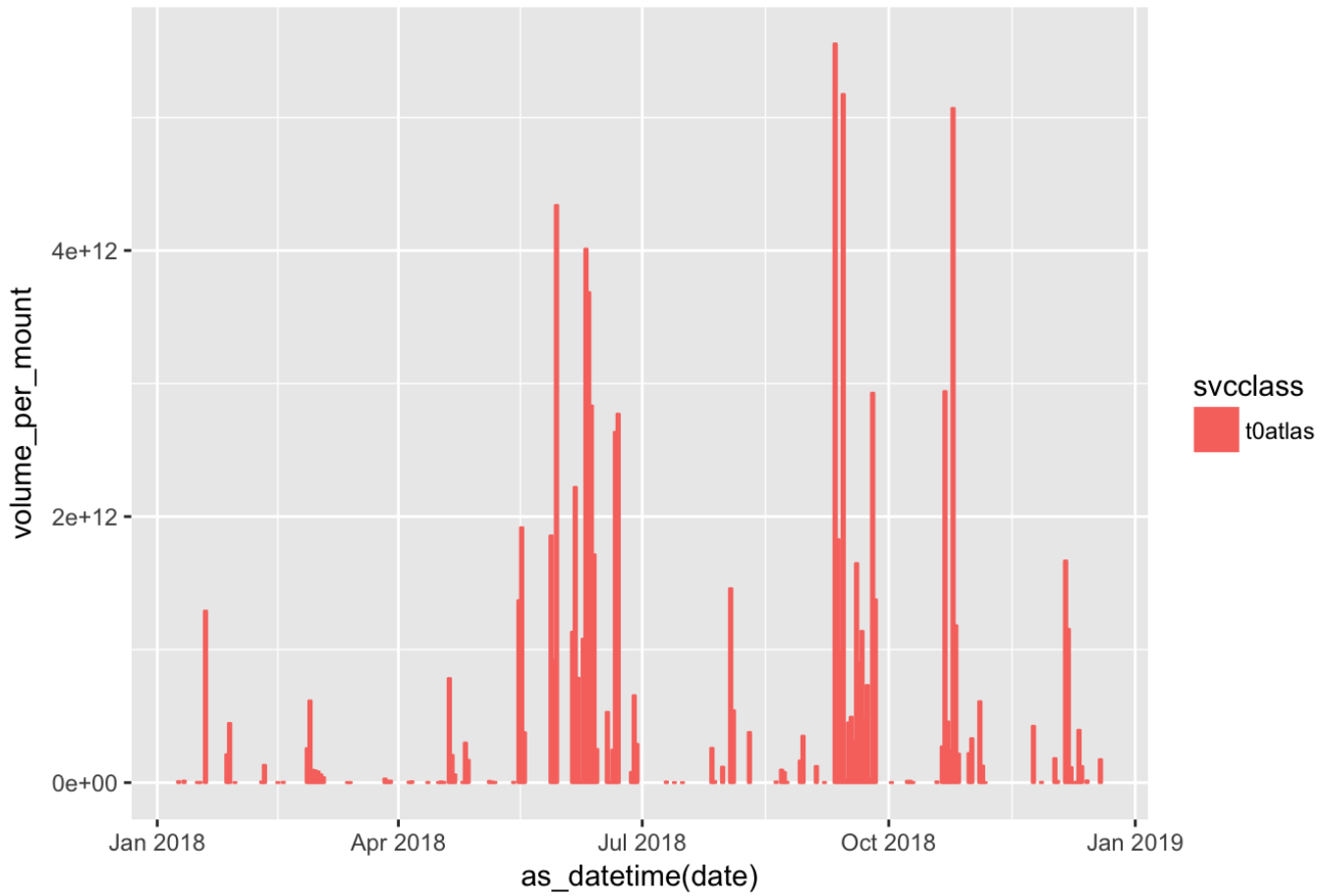
volume per mount, default



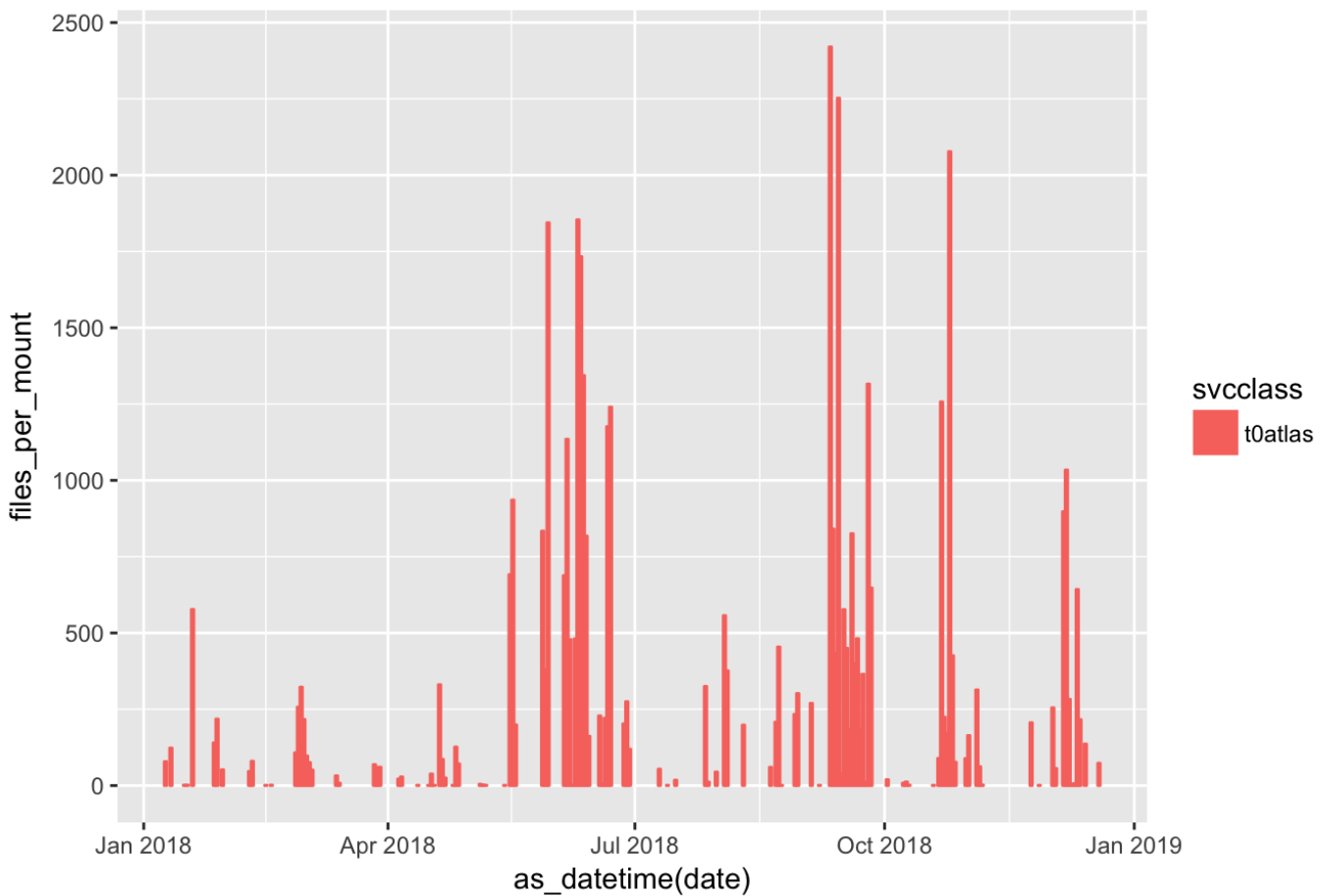
files per mount, default



### volume per mount, t0atlas



### files per mount, t0atlas



# Per-svcclass repeated mount rates

Average repeat mount rate (number of times a same tape is mounted during a time interval):

Daily:

- default: 1.2532334
- t0atlas: 1.1107892

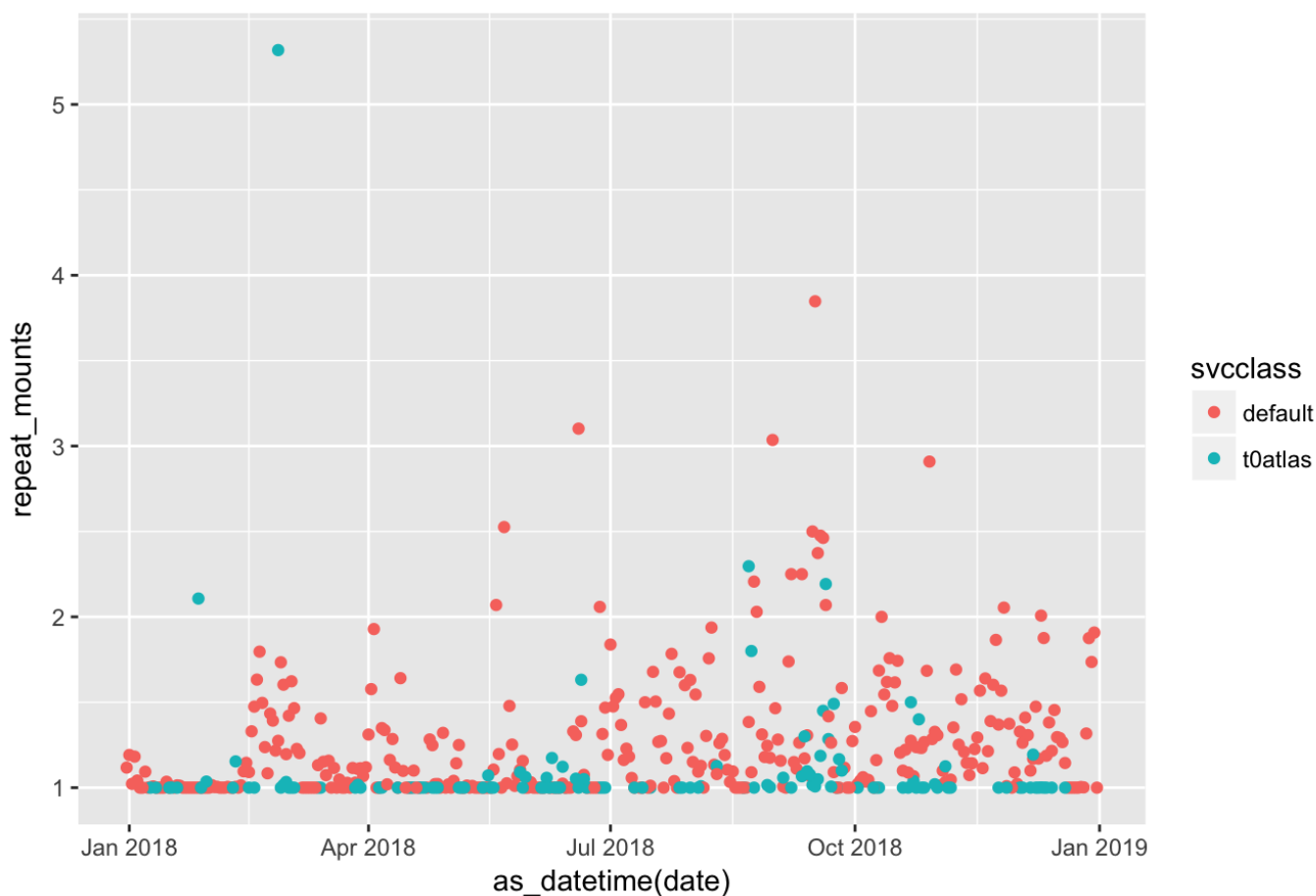
3 days:

- default: 1.6843782
- t0atlas: 1.1766316

1 week:

- default: 2.0192878
- t0atlas: 1.4035372

average daily repeated mounts



## Positioning times

- Median positioning times: 2.823471s, average: 4.6154742s
- Percentage of positioning times > 20s: 3.3433199 %
- Percentage of contiguously read files (positioning times < 0.1s): 25.7029855 %

**default service class:**

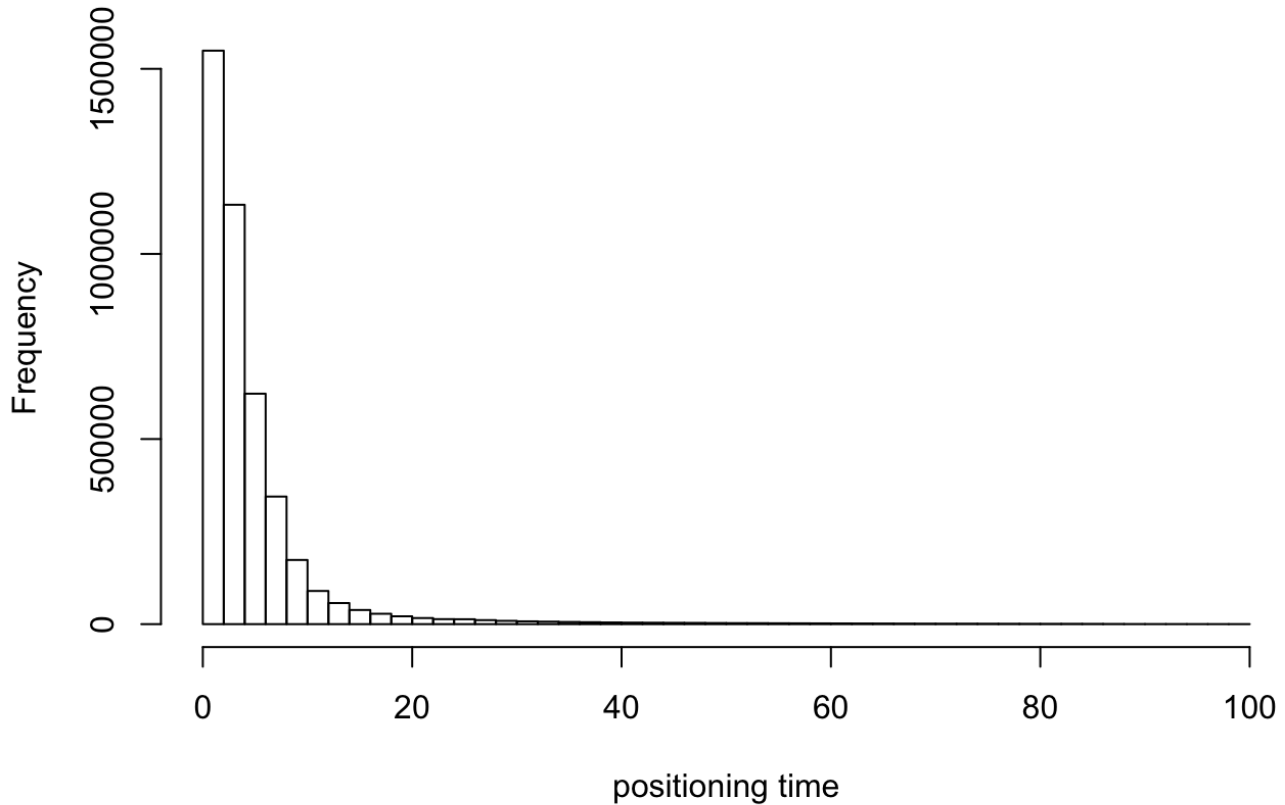
- Median positioning times: 4.110638s, average: 6.6432804s

- Percentage of positioning times > 20s: 6.0390471 %
- Percentage of contiguously read files (positioning times < 0.1s): 18.0604226 %

## t0atlas service class:

- Median positioning times: 2.1973055s, average: 2.987429s
- Percentage of positioning times > 20s: 1.1790274 %
- Percentage of contiguously read files (positioning times < 0.1s): 31.8388962 %

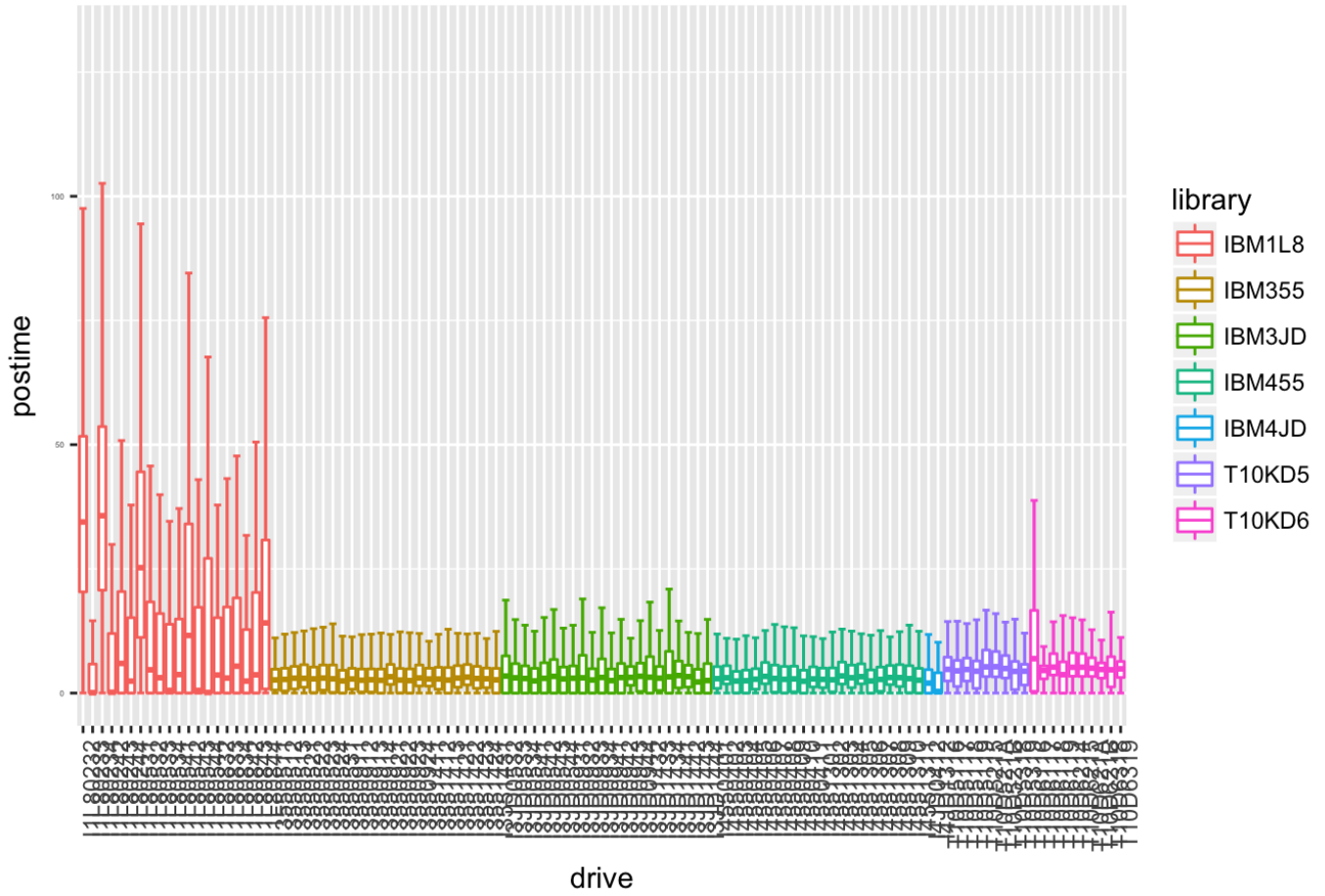
### positioning time per file



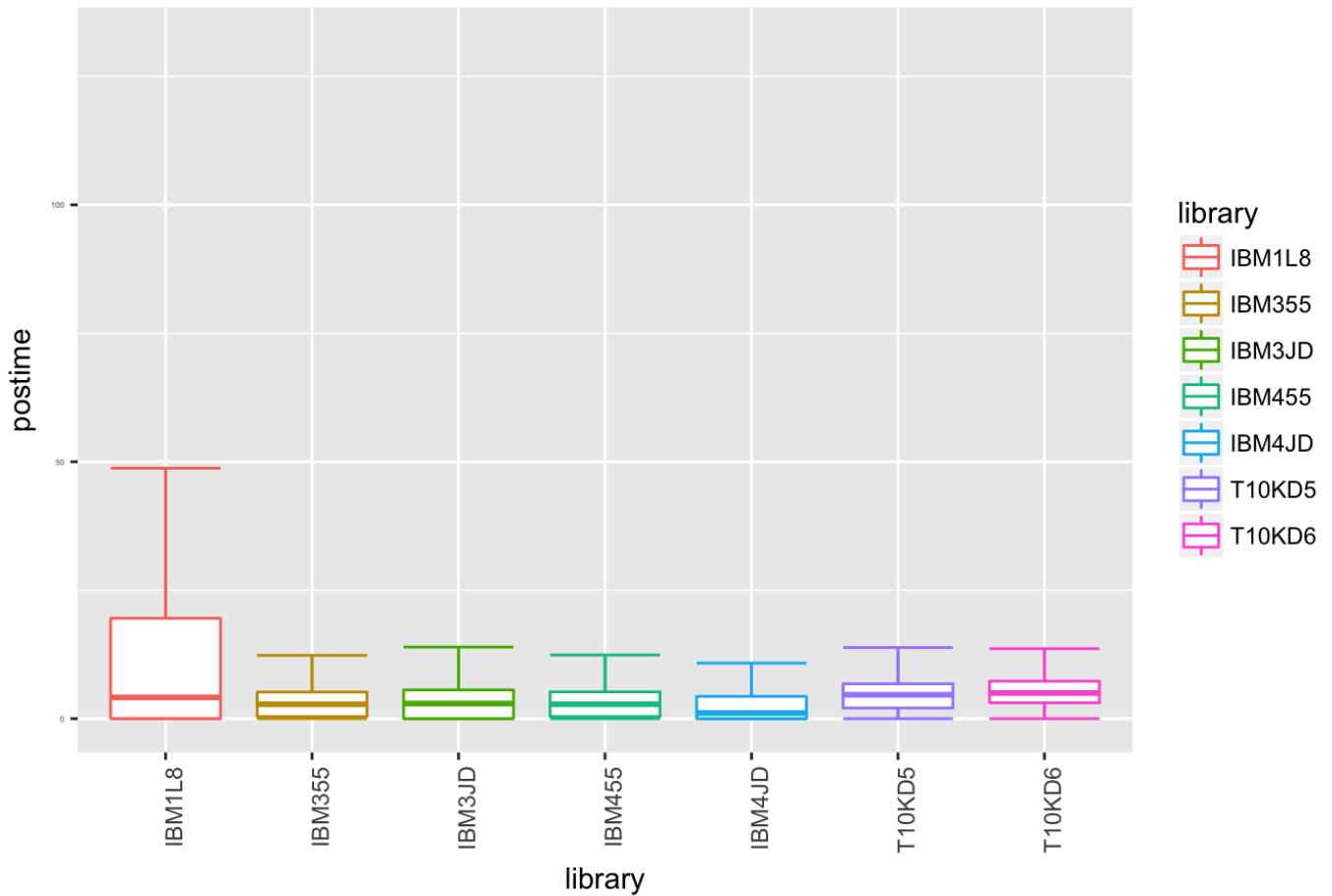
## Per-drive, per-library performance

Check for tape drive and overall library performance figures (positioning time, transfer speed)

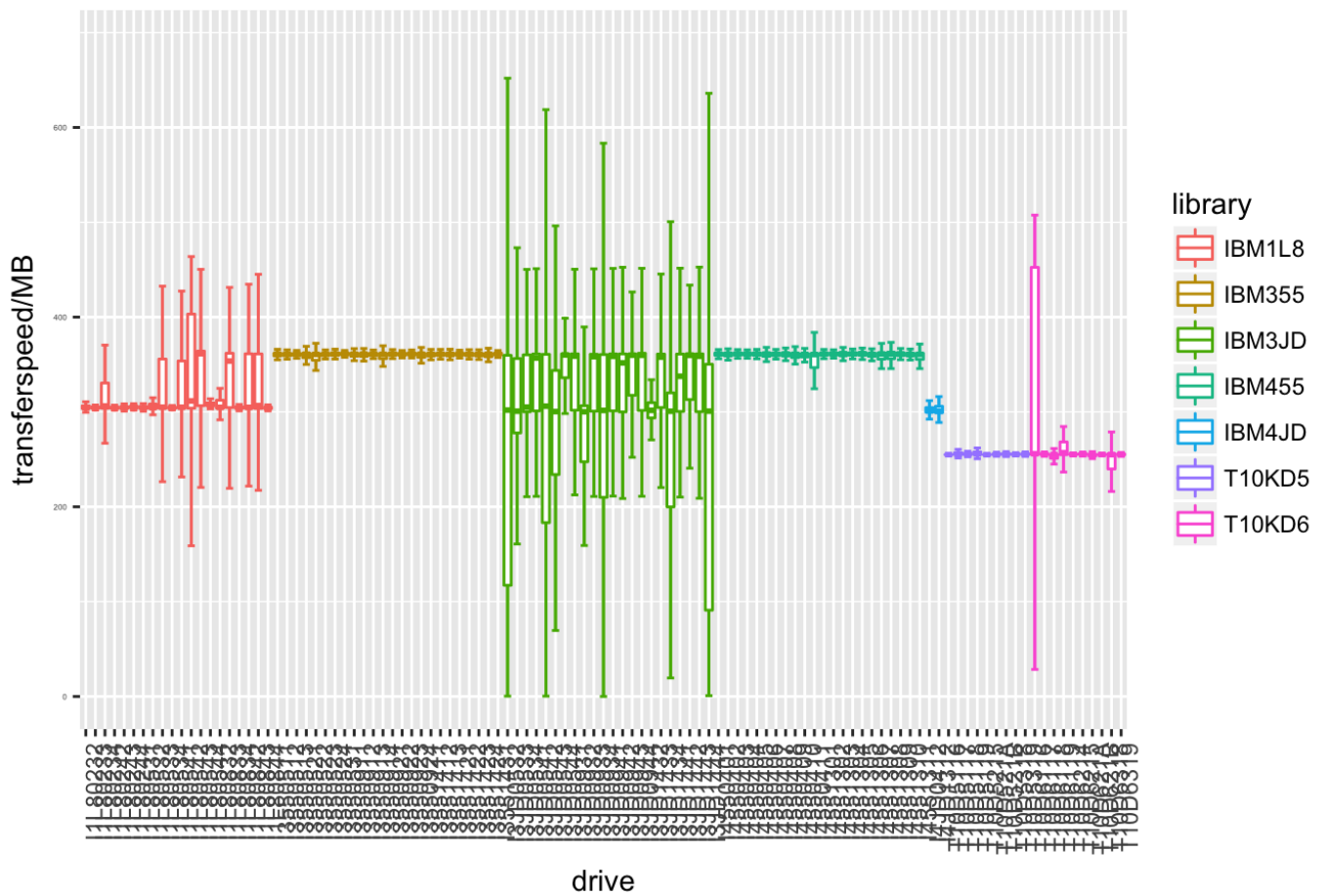
### positioning time, by drive



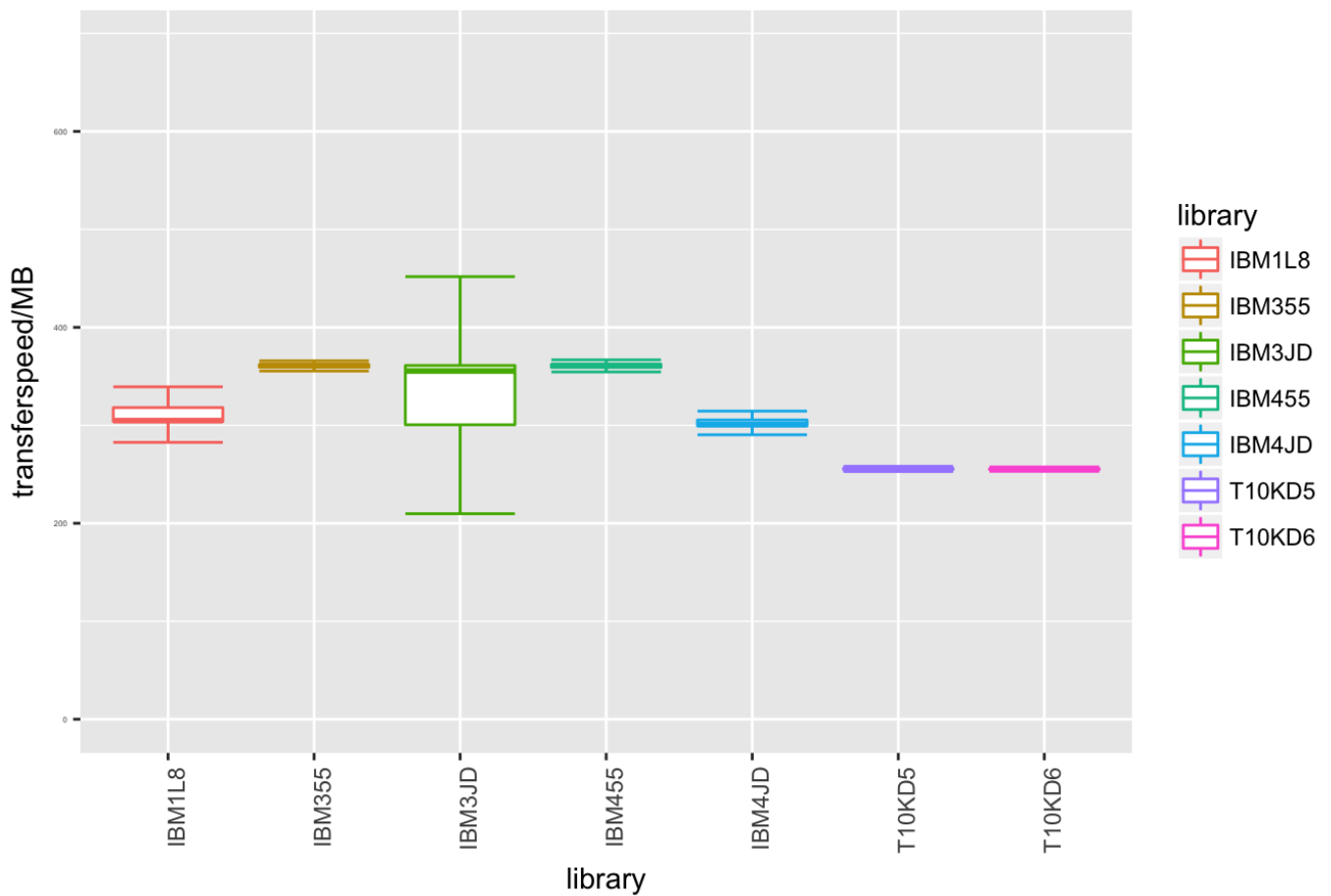
### positioning time, by library



### transfer speed, by drive

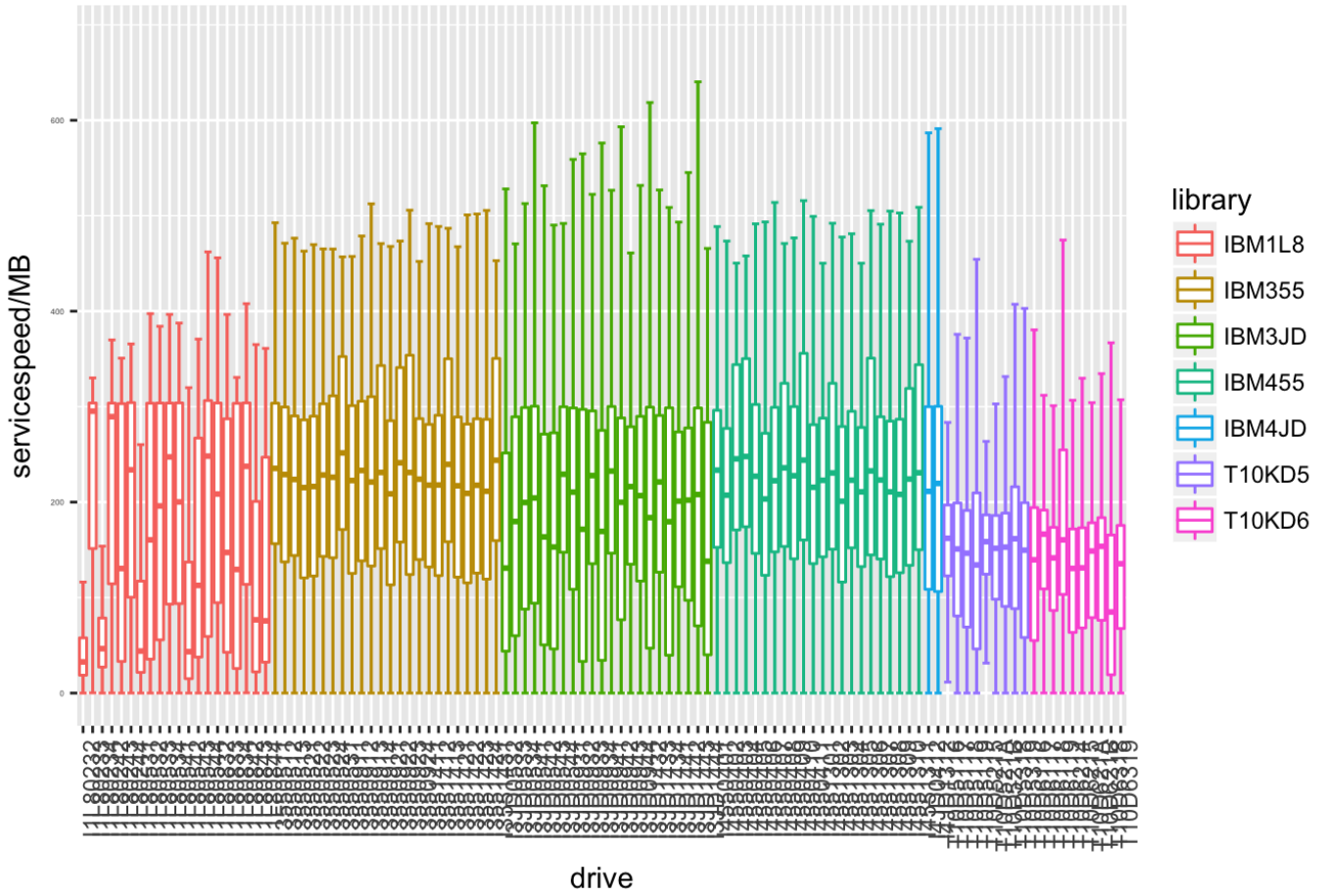


### transfer speed, by library

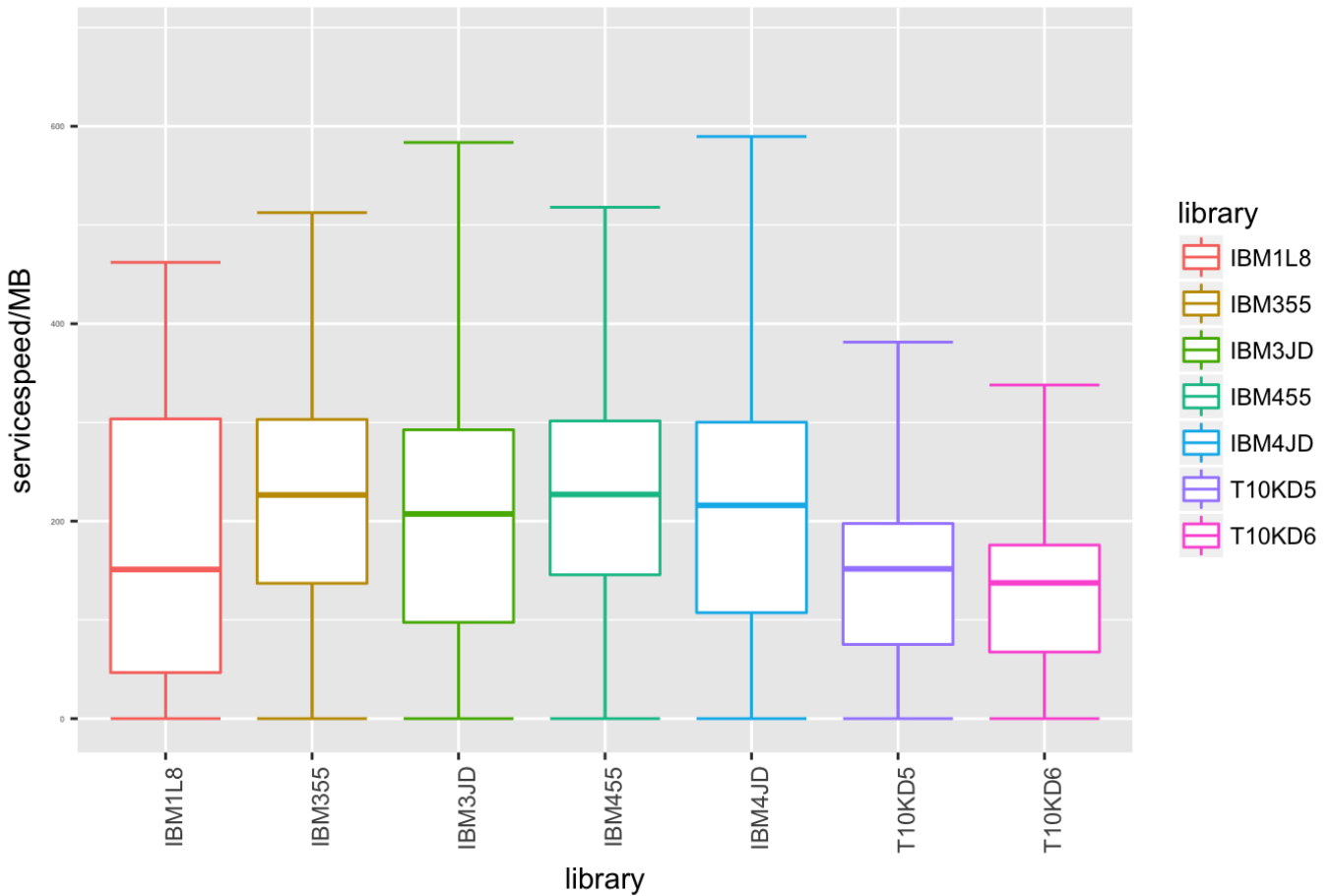




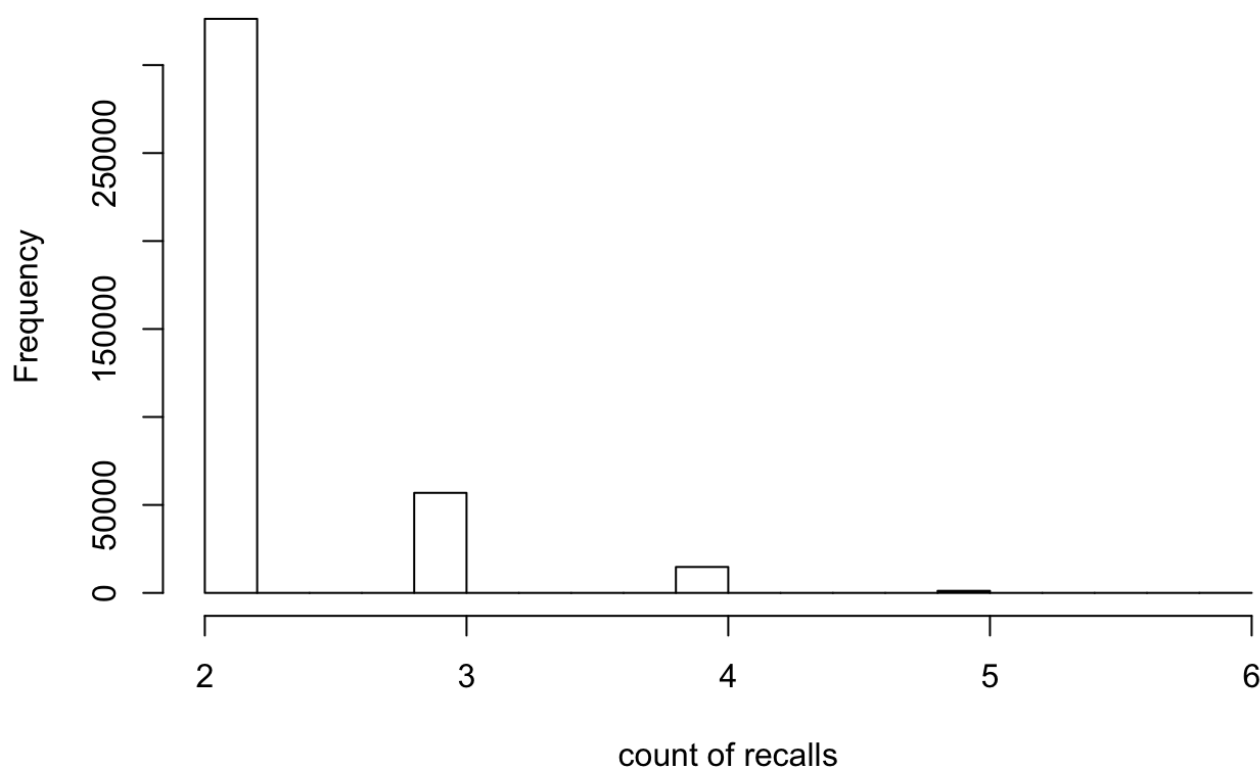
### service speed, by drive



### service speed, by library



## Files recalled more than once



## Repeated file access

- how many files have been recalled multiple times? What was the time interval distribution between recalls for these?
- how does this differ between “default” and “t0atlas”?

### General:

From the 4196637 files read, 489086 files (11.6542365 %) have been recalled more than once. Out of these, 83160 files (1.9815867 %) have been recalled from both service classes.

### default service class:

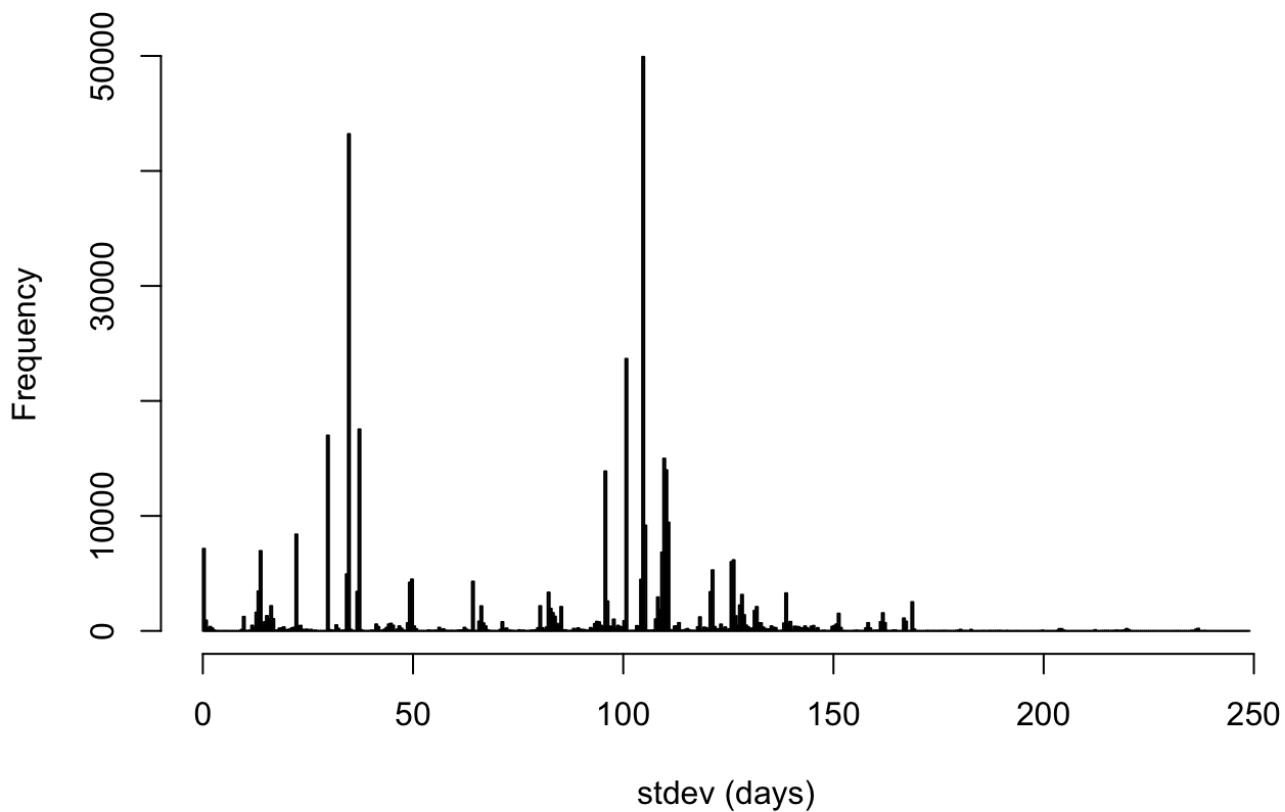
From the 1868871 files read, 86665 files (4.6372917 %) have been recalled more than once.

### t0atlas service class:

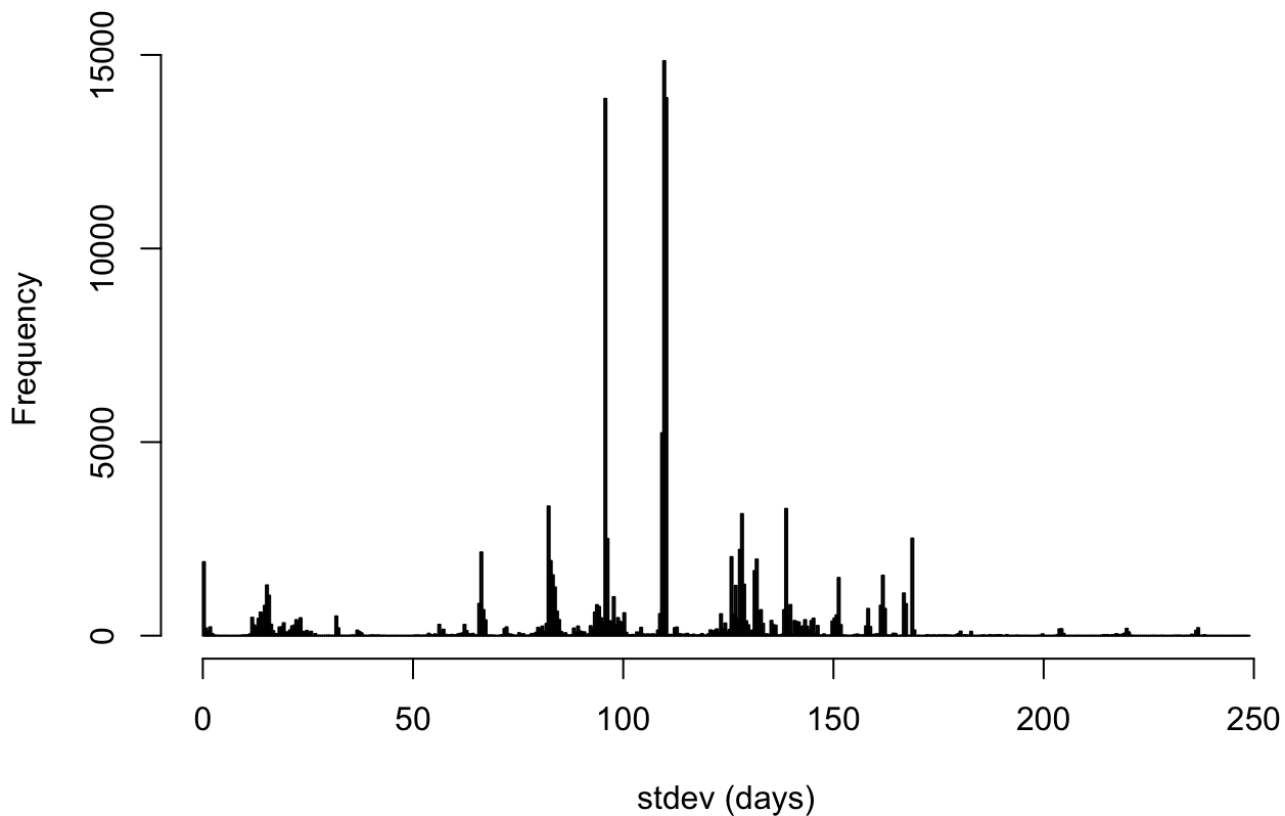
From the 2327766 files read, 319261 files (13.7153391 %) have been recalled more than once.

- For the files recalled more than once within 2018, what is the histogram of standard deviation in repeated access time?
- For the files recalled more than once within 24h, how does the histogram look like?
- Are there any differences between t0atlas and default?

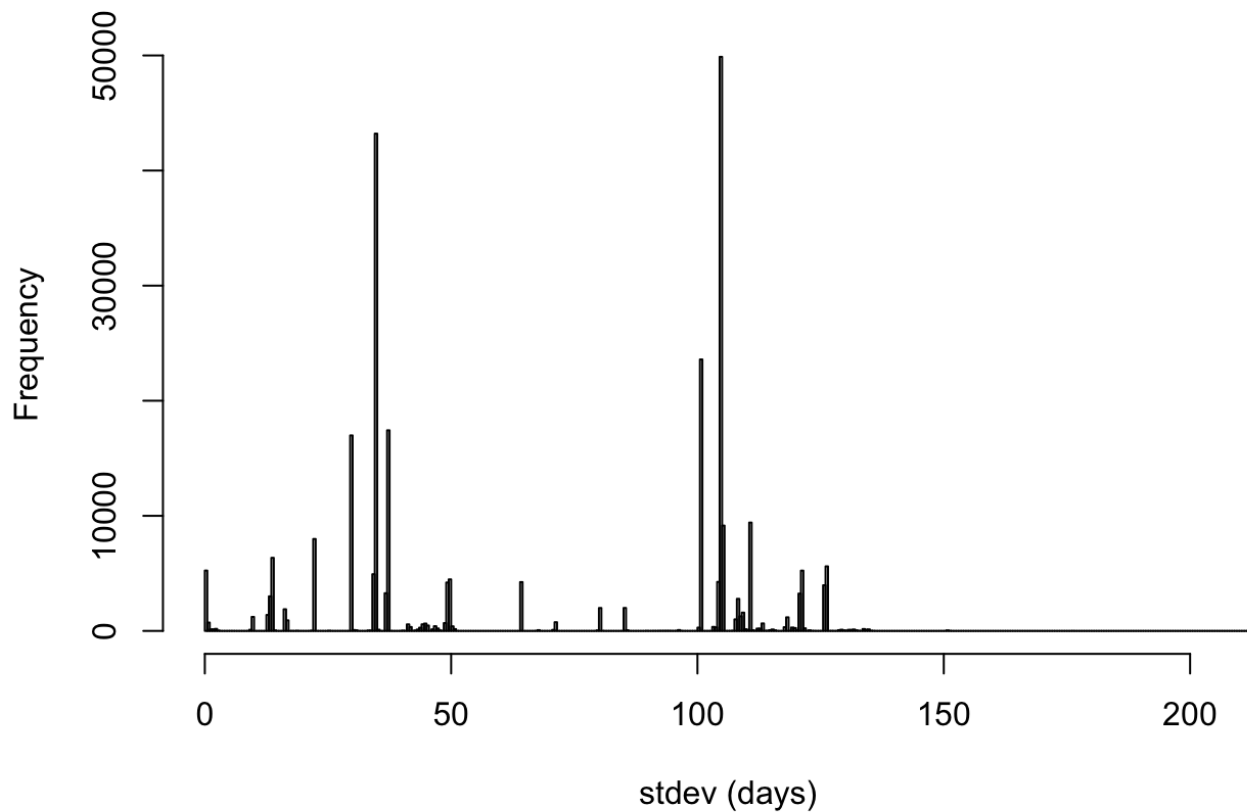
### stdev (days) between repeated file recalls



### stdev (days) between repeated file recalls - default



## stdev (days) between repeated file recalls - t0atlas



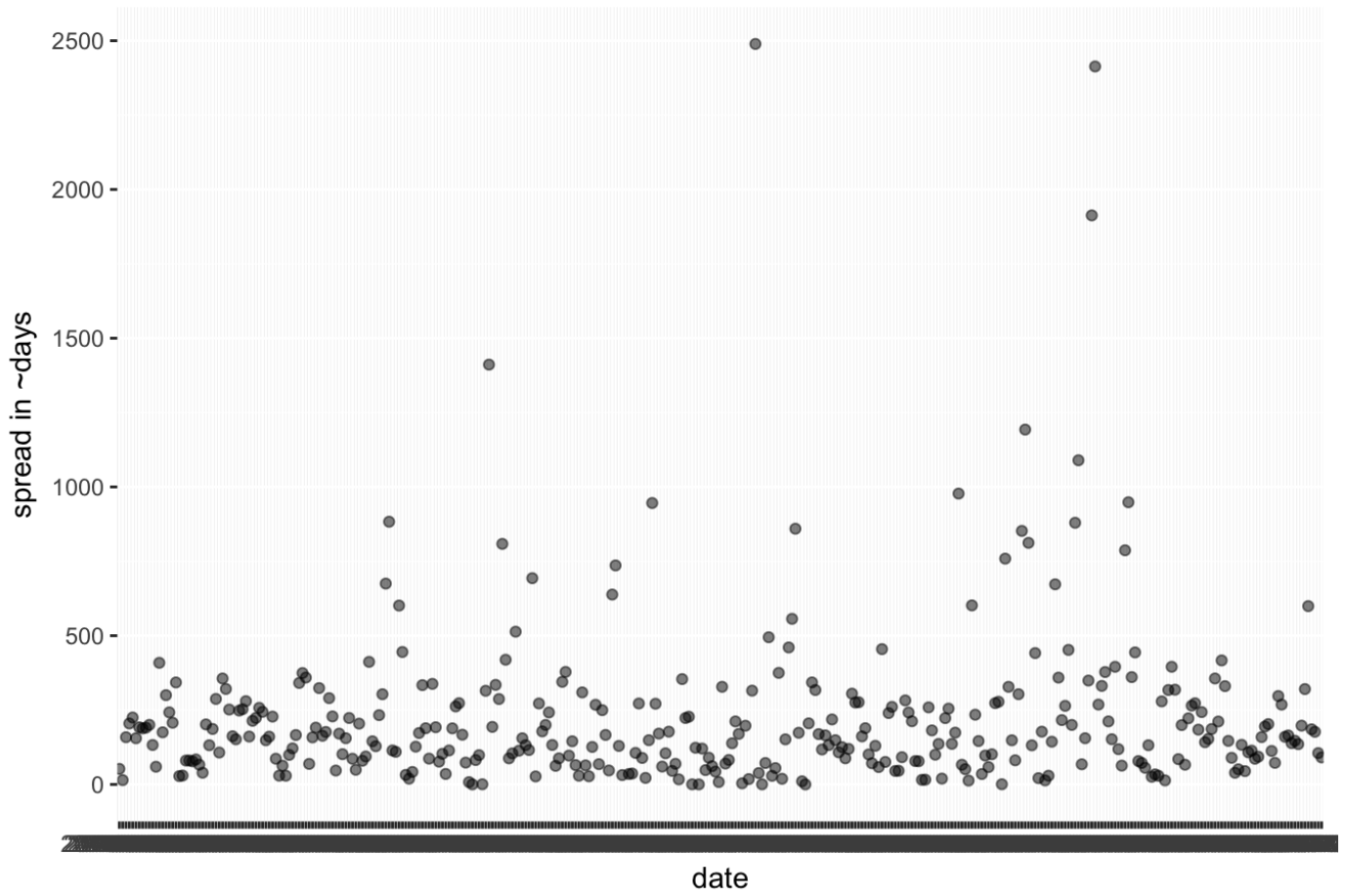
## File retrieval cohesion

- What is the spread in creation time of files retrieved within a single day? How collocated are recalls with regard to creation time?

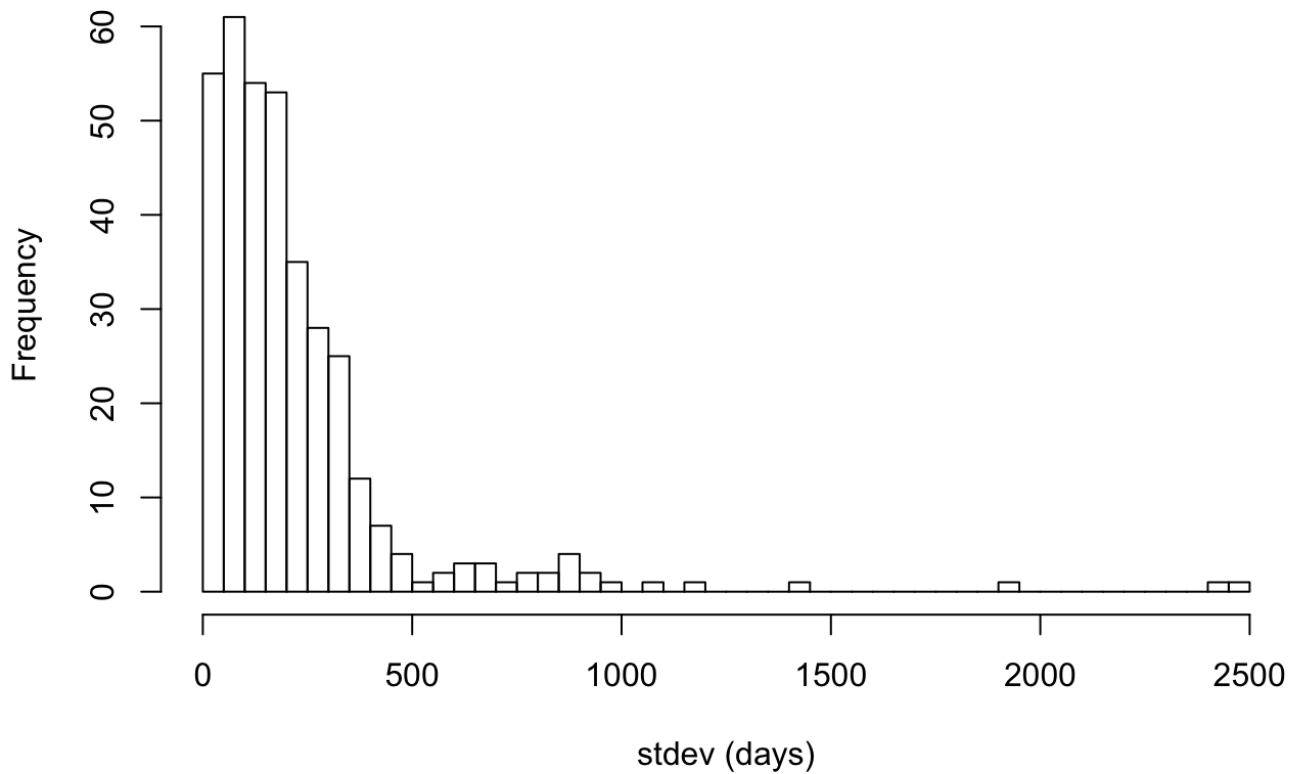
Basis is the daily standard deviation of the nsfileID of retrieved files. Average nsfileID creation/day (across all of CASTOR): ~230K files (stable over last 4 years), therefore showing distance divided by 230K files, so roughly corresponding to the number of days between the creation time of these files.

- standard deviation (in ~days) for “default” svcclass (daily plot, histogram)
- standard deviation (in ~days) for “t0atlas” svcclass (daily plot, histogram)

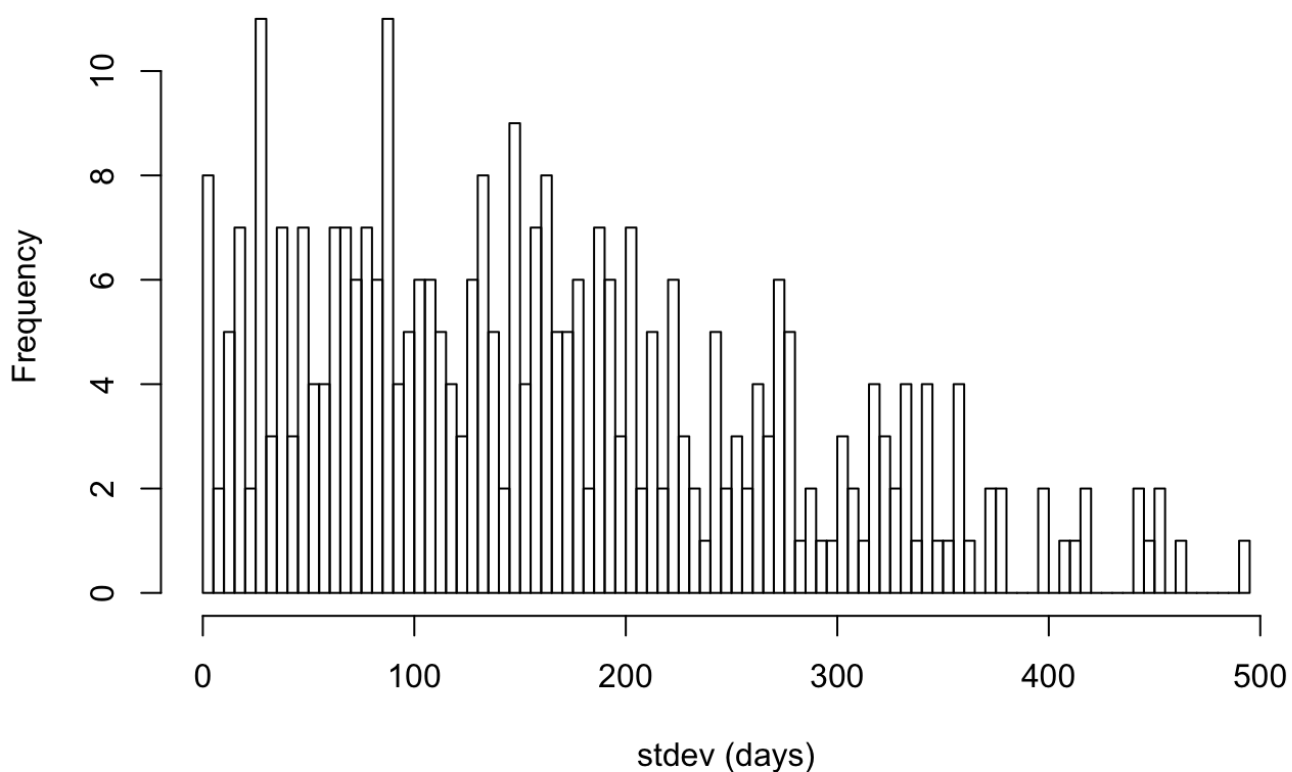
### Spread in creation time of recalled files, default svcclass



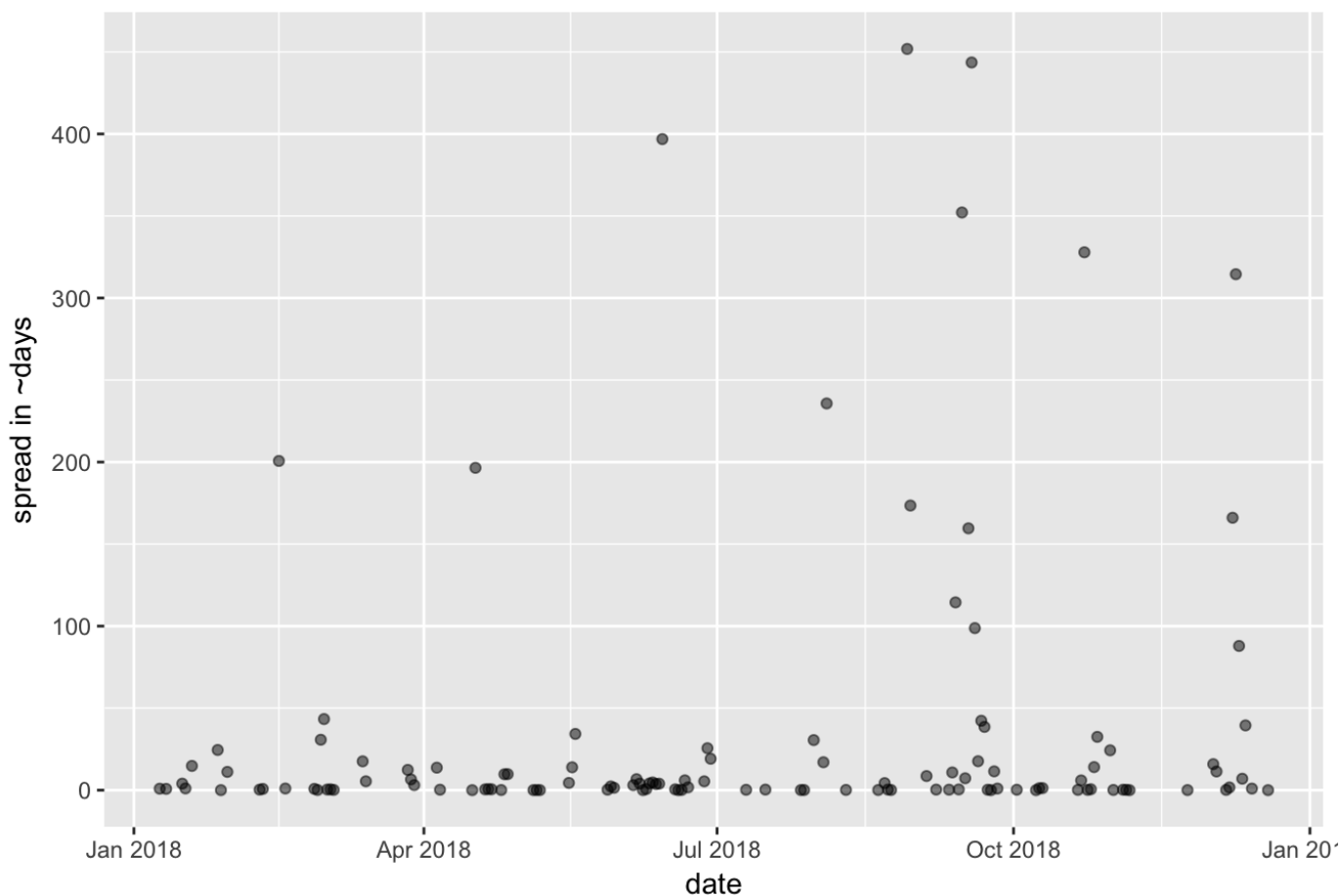
### Spread in creation time of recalled files, default svcclass



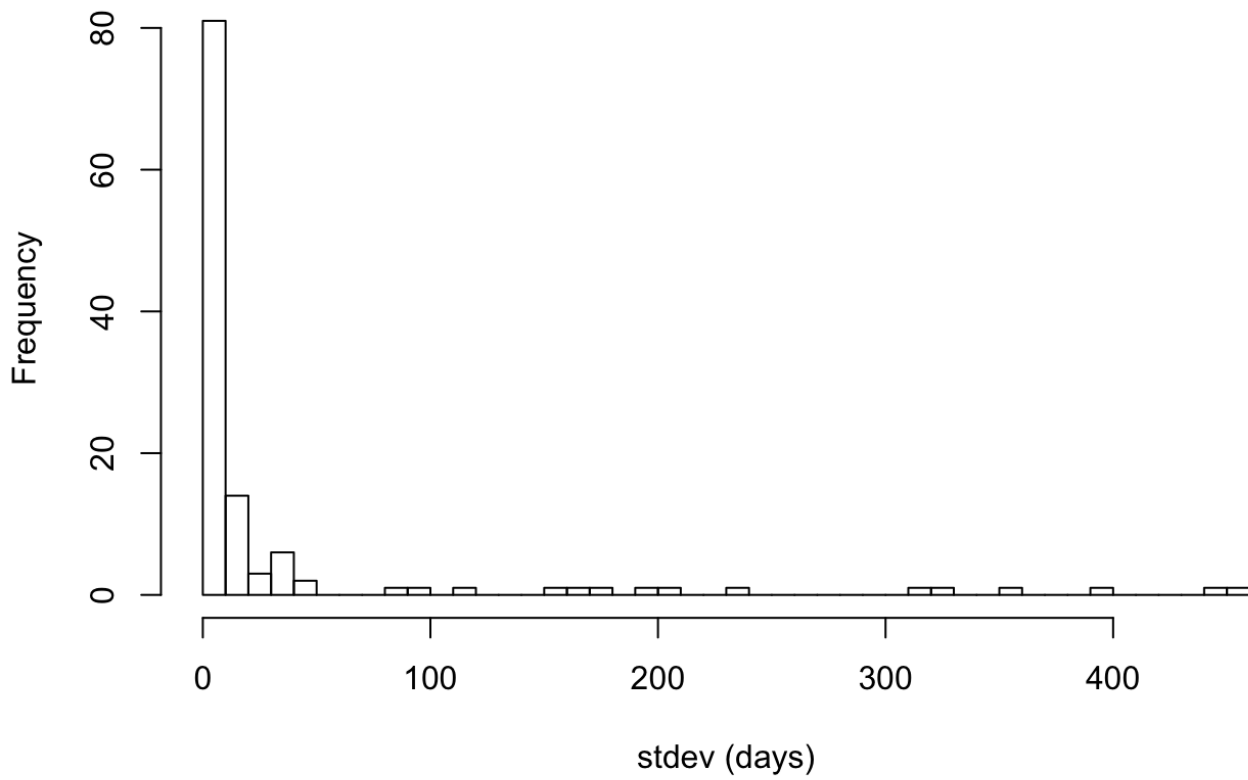
### Spread in creation time of recalled files, default svcclass, 500d



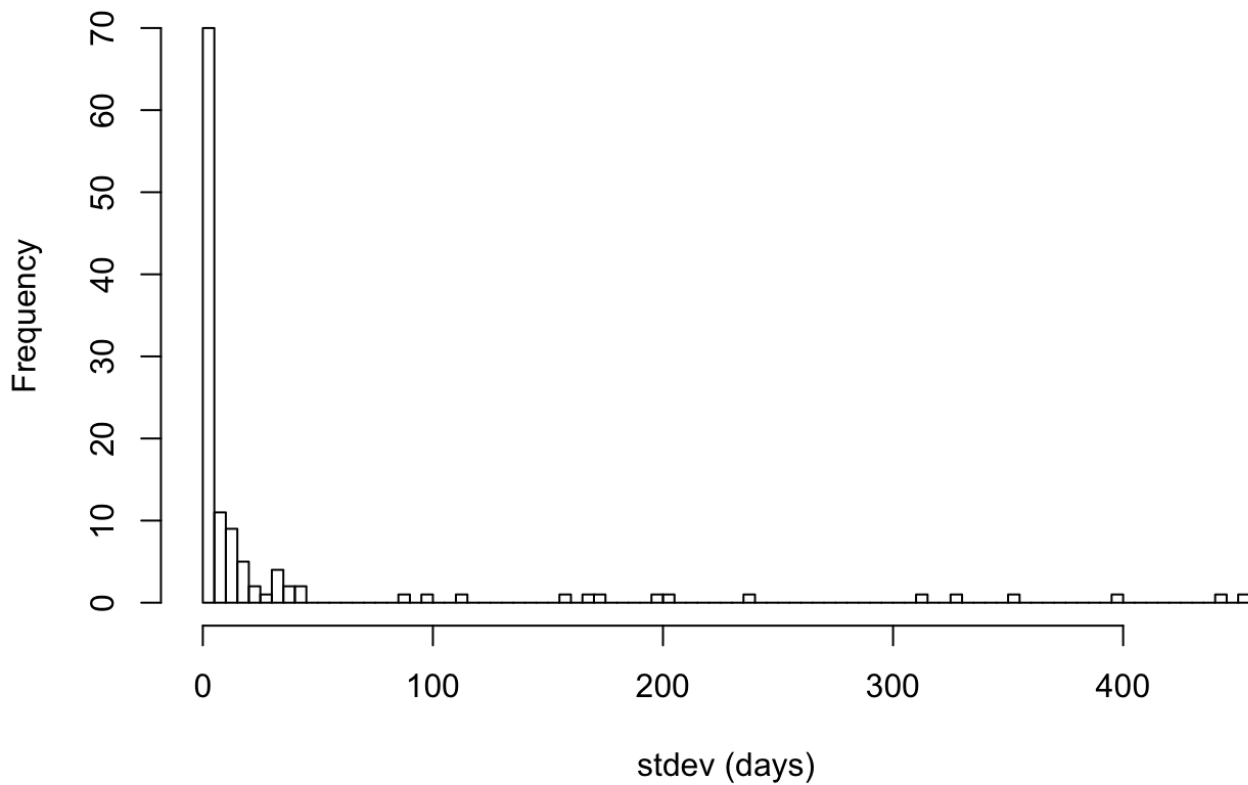
### Spread in creation time of recalled files, t0atlas svcclass



### Spread in creation time of recalled files, t0atlas svcclass



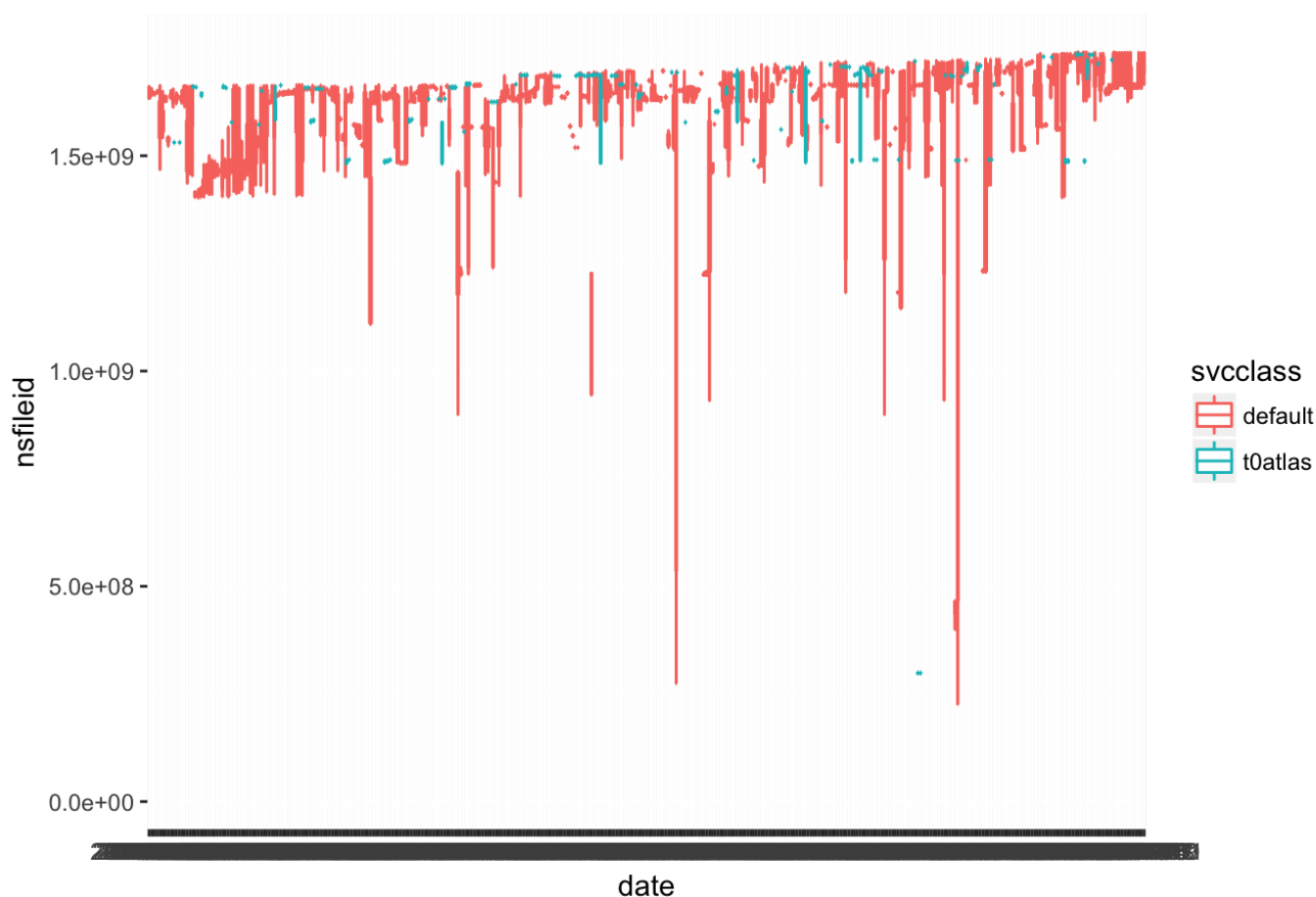
### Spread in creation time of recalled files, t0atlas svcclass, 500d



# Per-day recalled nsfileID boxplots

How far away in time are the files being recalled

nsfileID daily retrieve spread by svcclass



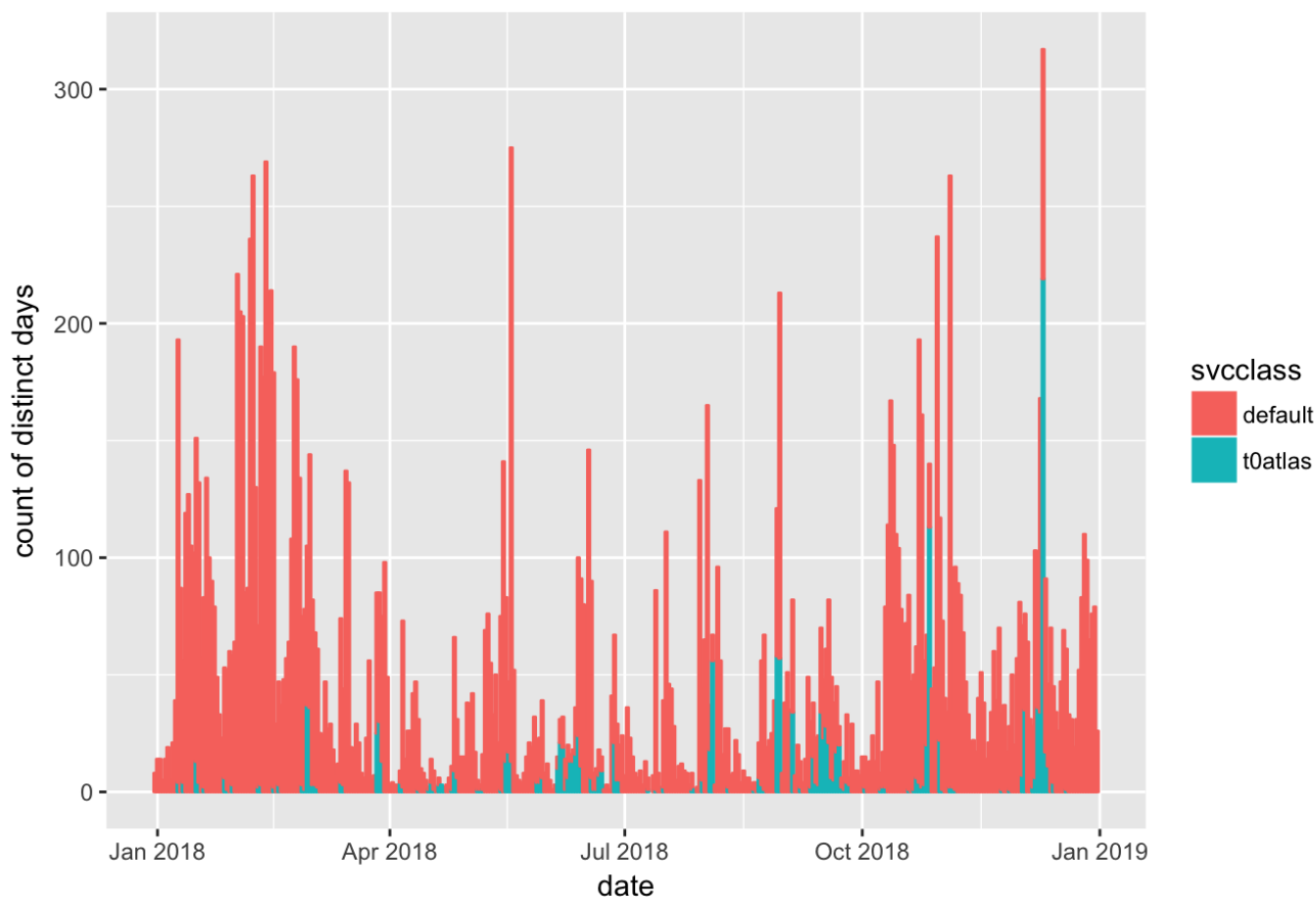
## Average number of distinct creation time days processed within a day

For a given day, in addition to the standard creation time deviation (which represents the overall spread in creation time between all files being processed) it is interesting to know from how many distinct creation days there are files being processed (which gives an idea of how grouped these files are, in addition to the overall creation time spread)

- default: average of 49.5331492 different creation days
- t0atlas: average of 11.7322835 different creation days



## distinct creation days accessed per day

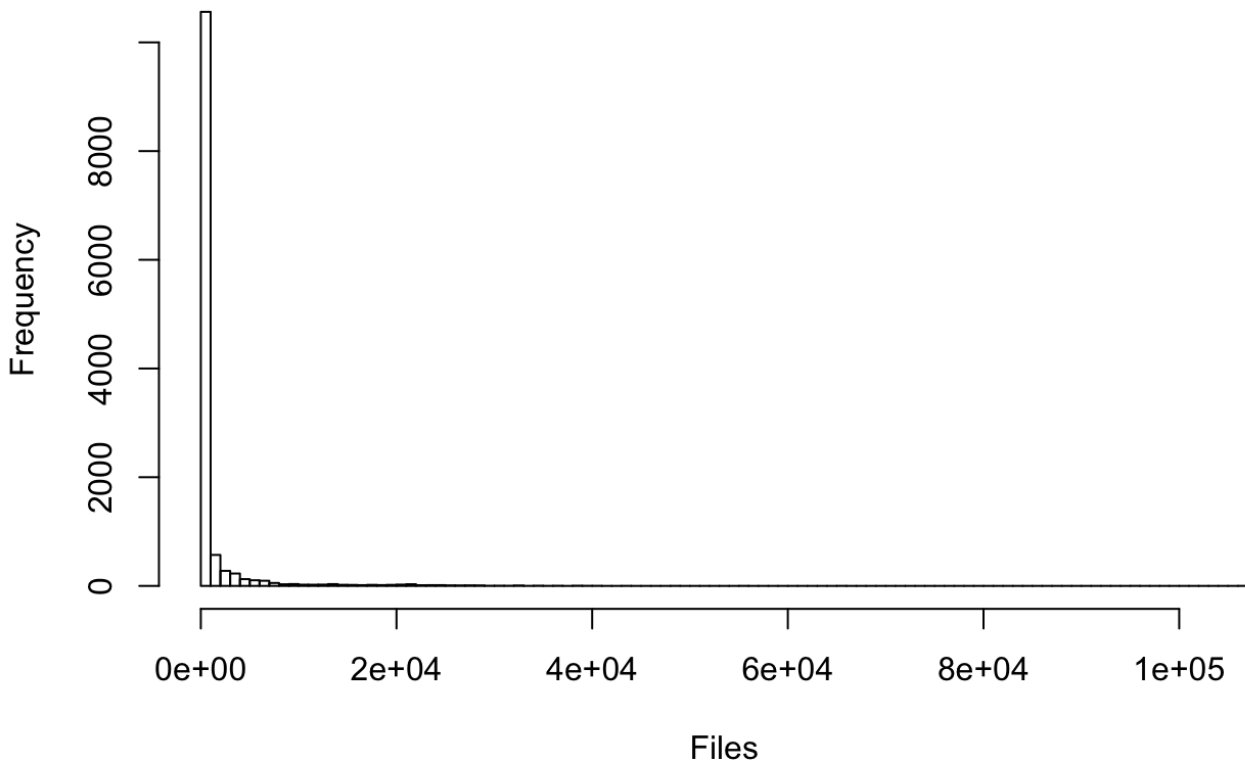


## Accessed ATLAS datasets

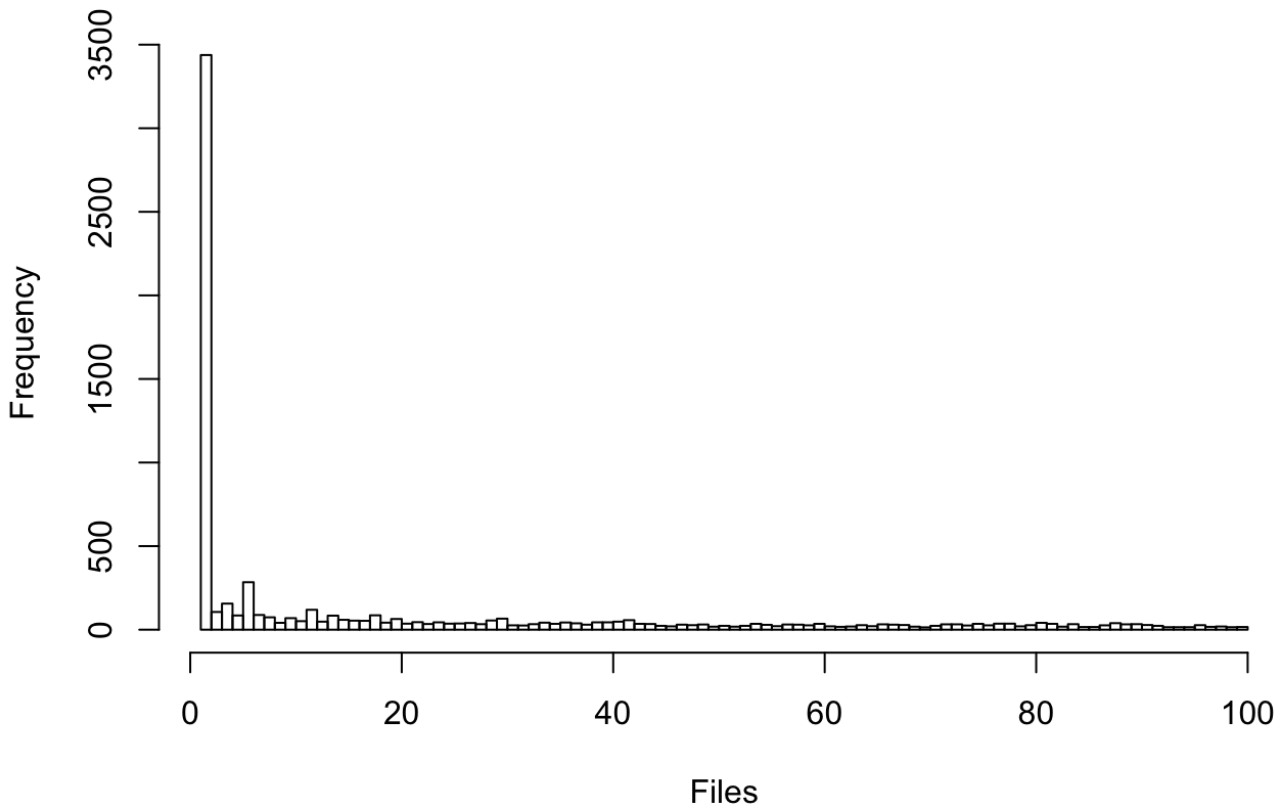
ATLAS dataset names have been identified by extracting the last directory path entry before each file name. Rucio datasets are identified by a naming convention (see <https://gitlab.cern.ch/cta/CTA/issues/461> (<https://gitlab.cern.ch/cta/CTA/issues/461>) for details). For each of the datasets that have been accessed (at least one file recalled), the following characteristics have been extracted:

- Total number of datasets: 12461, number of Rucio datasets: 12413 (99.6147982 %)
- Total number of files across all datasets: 14.265152 million
- Total volume across all datasets:  $3.052964510^4$  TB
- Average number of files per dataset: 1144.7838857; median: 55
- The mean and median differ substantially: There are many datasets with just 1-2 files (27.5900811 %) and there is a tail of a few, long datasets (see histograms below)
- Average volume per dataset: 2.4500156 TB; median: 0.0300045 TB.
- Mean and median differ again greatly as a consequence of many datasets being very small.
- Average spread in time (first to last creation time): 2.5432881 days, median: 0.0725652 days.
- Same observation as before. In addition, in the cases where a file has been re-created (e.g. re-import from T1 following a file loss), the spread will become artificially large as the creation time of the recovered file may be much younger than the rest of the files. See also the histograms below.

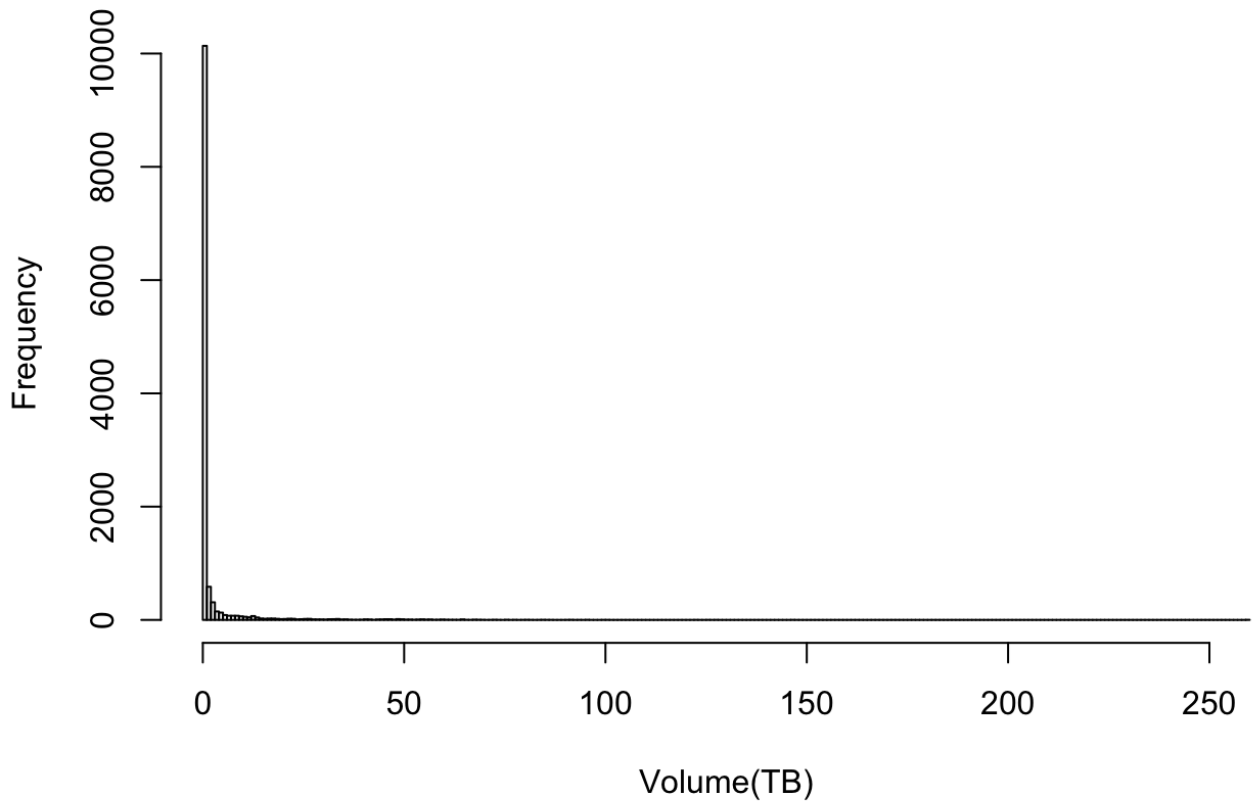
### Files per dataset



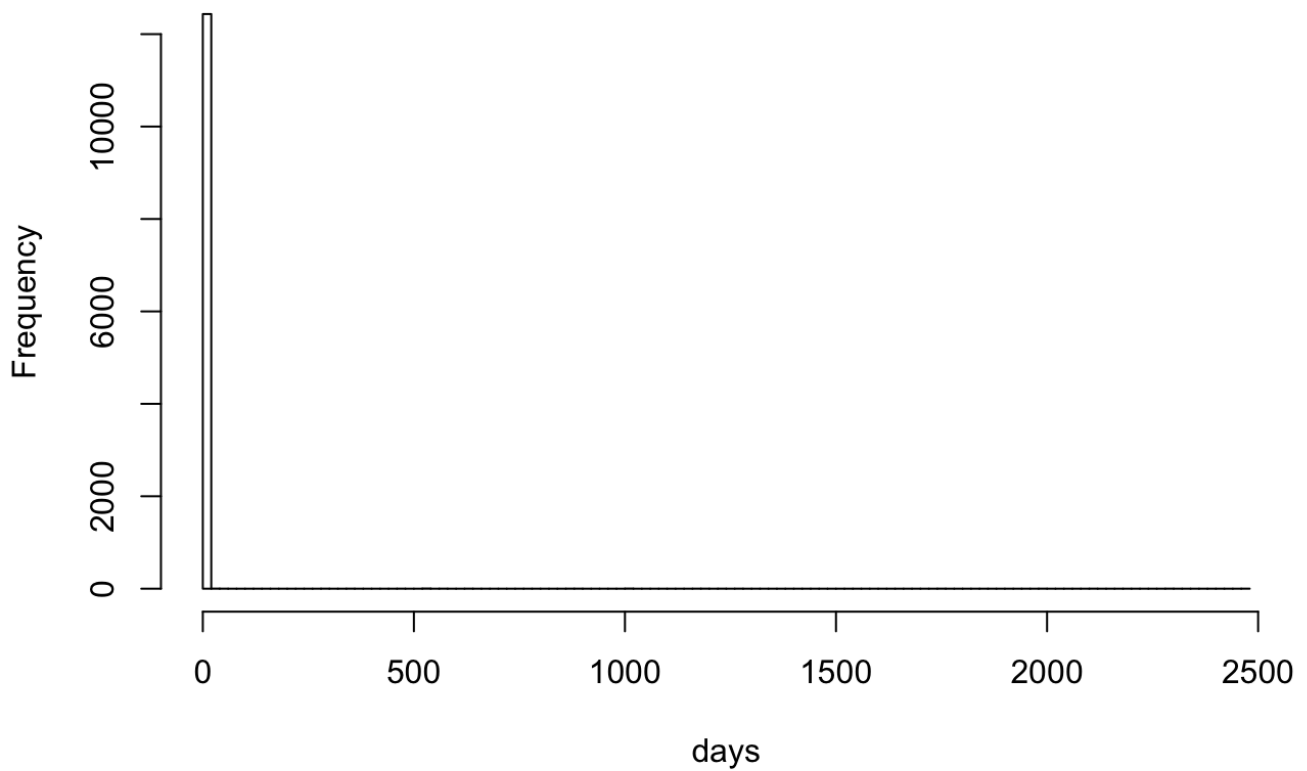
### Files per dataset (cut at 100)



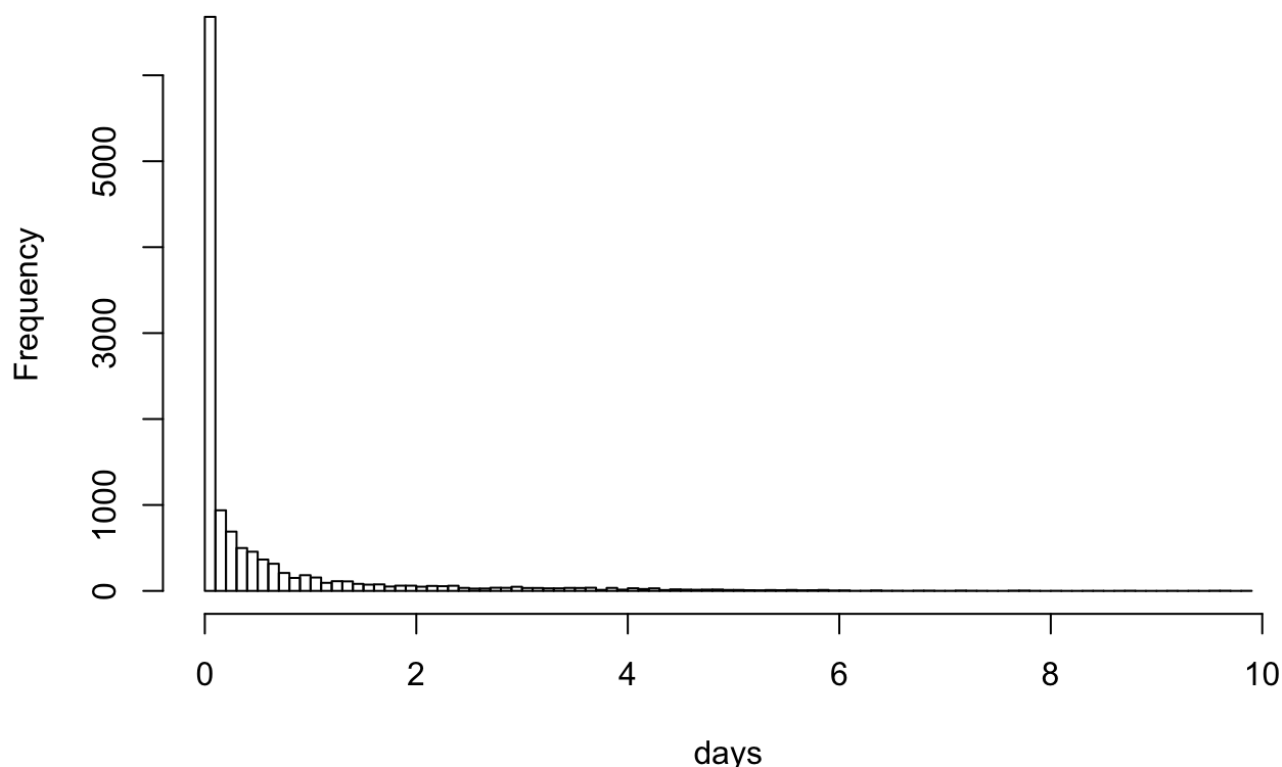
### Volume per dataset



### Dataset spread in time(days)



## Dataset spread in time(days)(cut at 10 days)



## Dataset access patterns

In the section above we have identified the ATLAS datasets that were accessed during 2018. How was the actual access to them, as a function of the service class (default vs t0atlas)?

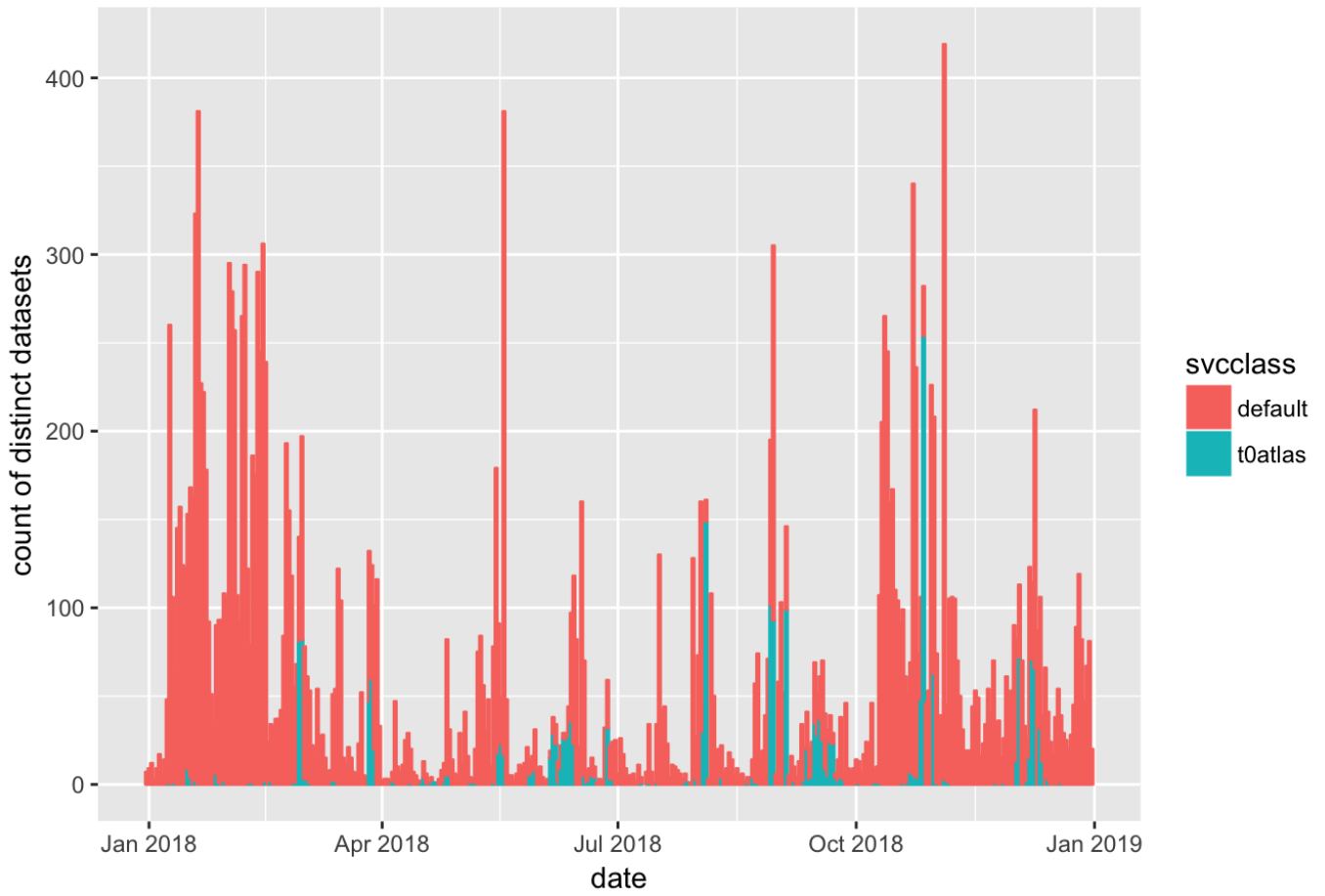
default:

- total datasets accessed by default: 9482
- total files in these datasets: 12752024; fraction accessed via default: 14.6554853 %
- total volume in these datasets: 2.797084610<sup>4</sup> TB; fraction accessed via default: 13.291533 %
- Average number of existing files per accessed dataset: 1344.8664839; median: 90

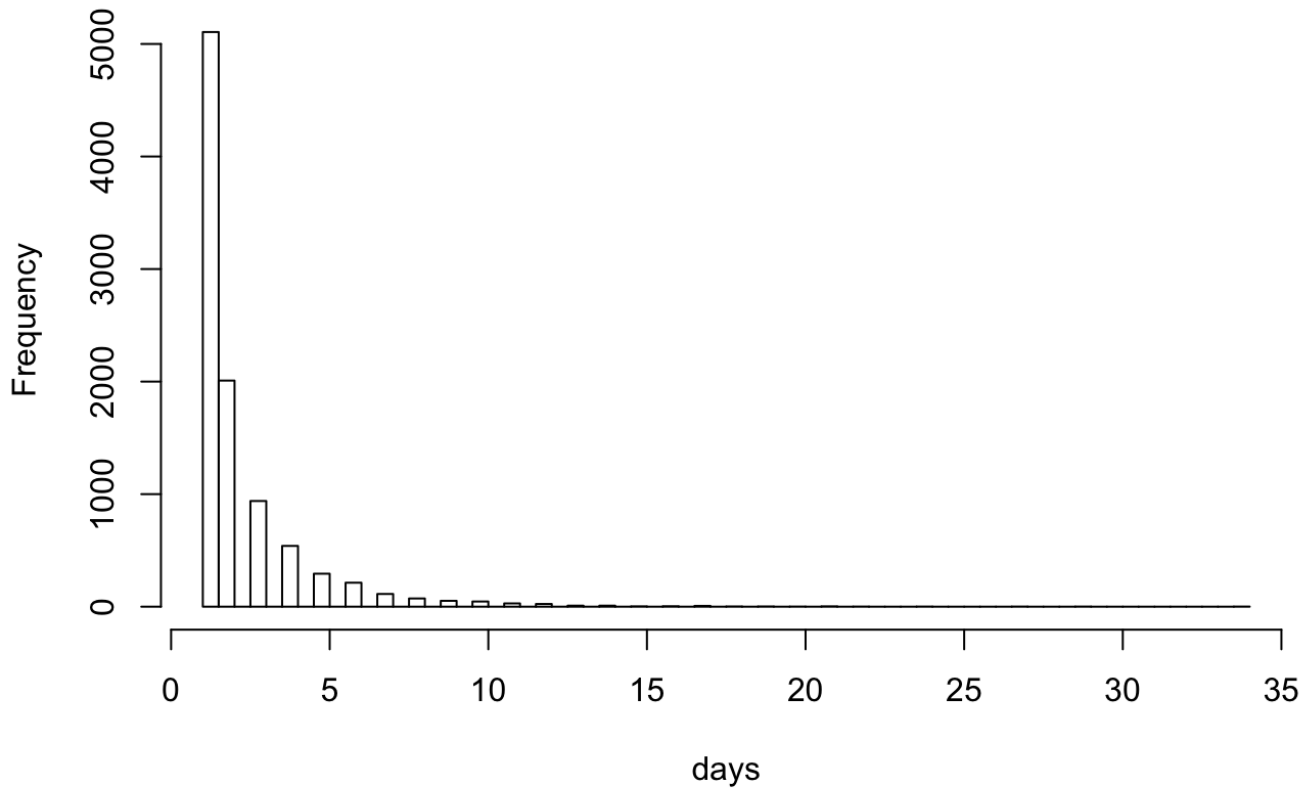
t0atlas:

- total datasets accessed by t0atlas: 3214
- total files in these datasets: 2218264; fraction accessed via t0atlas: 104.9363827 %
- total volume in these datasets: 4049.016156 TB; fraction accessed via t0atlas: 102.3750511 %
- Average number of existing files per accessed dataset: 690.1879278; median: 2. Note the low median - 50% of the t0atlas accessed datasets have less or equal than 2 files.

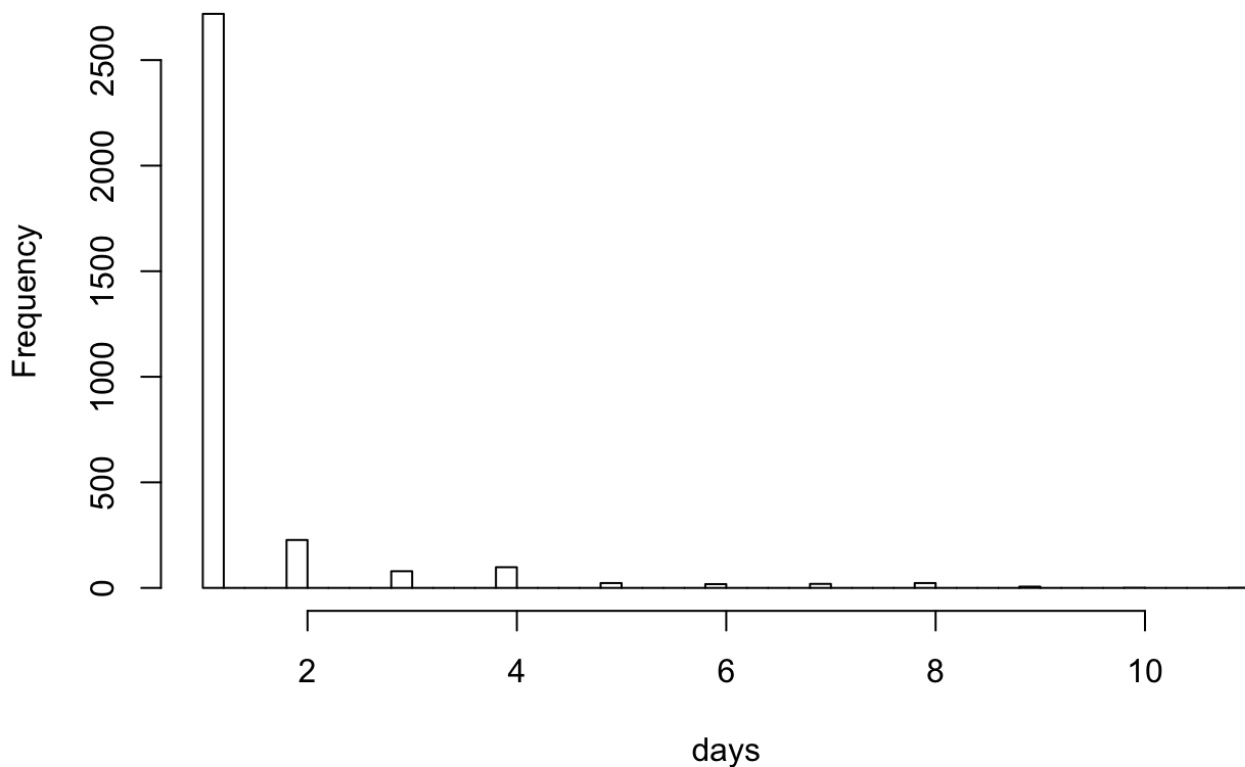
### distinct datasets accessed per day



### Distinct days on which a given dataset is accessed (default)



## Distinct days on which a given dataset is accessed (t0atlas)



## Mount excess per dataset

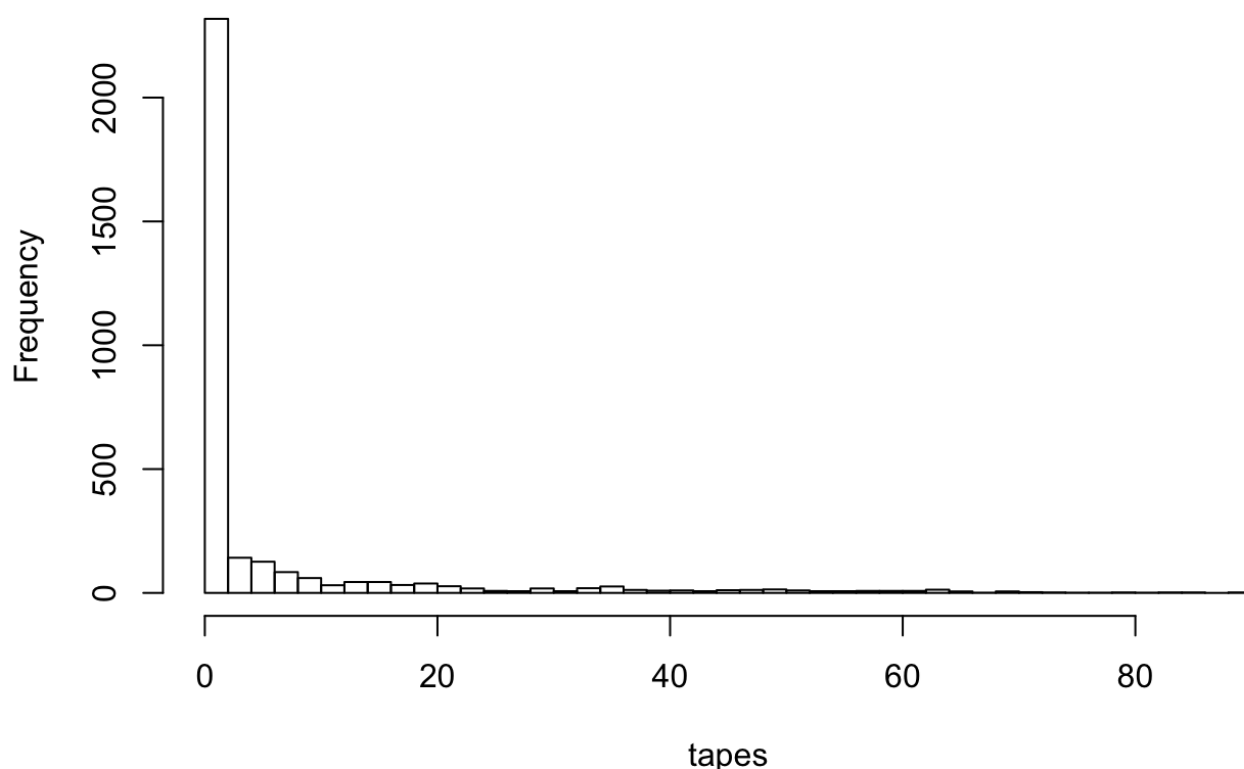
How many more mounts did we issue to what it would be necessary to read out a dataset?

Note: Both tape pools (atlas\_raw and atlas\_prod) are configured with up to 10 concurrent write drives each. So if there is sufficient data to be written, up to 10 drives will be mounted in parallel for writing. Thus datasets can be spread across 10 tapes, which means reading them back implies 10 mounts (in parallel or sequentially) - for those datasets that have sufficient files of course.

In terms of mount counts, it is difficult to go down below the count of accessed datasets, unless the datasets themselves are grouped within a single tape which is not usually expected. So what is the baseline “excess” factor?

- Number of mounts; number of datasets accessed; “excess” factor (mounts per dataset accessed):
- default: 73798 ;9482 ; 7.7829572
- t0atlas: 7670 ;3214 ; 2.3864343

## Tapes per dataset, t0atlas



## Queue Analysis (To Be Completed)

In addition to the number of tape mounts, the average queueing waiting time is a significant metric for understanding what the overall latency for data access is.

Today, in CASTOR, a two-level queueing is used: Tape read requests queue up first within the stager and then on VDQM. The stager holds back submitting tape mounts to VDQM to not exceed a given number of parallel mounts. However, the stager doesn't distinguish whether submitted jobs are running or just queued (e.g. due to busy library), which may cause a tarpit effect (jobs accumulate on busy/slow/unavailable libraries). As a workaround, jobs are occasionally (manually) released to VDQM, sometimes creating an "avalanche" effect. VDQM tape queues are processed (per library) in FIFO order across all stagers following a FIFO ordering with all VO's having the same priority.

- How much does T0 activity (t0default) have to wait as a consequence of grid (default) or other VO mounts being processed? How would a separate queueing policy help?