# Profiling Contur

# Tools Used (1)

- We use cPython to profile the code which produces a .prof file. The profiling can be done with a small adjustment to the contur script in bin and then just running contur as normal from terminal.

- We use two separate packages to view the results in the .prof file, Snakeviz and gprof2dot.

- Snakeviz produces icicle and sunburst plots, it is browser based so has a dashboard style format where the user can dynamically play with the plots. See https://jiffyclub.github.io/snakeviz/ for further info.

# Tools Used (2)

- gprof2dot produces static plots which show the flow of the code and the proportion of time spent in each part of the code, see https://github.com/jrfonseca/gprof2dot.

- All profiling results can be found in the following repository https://github.com/SeanBrayUCL/contur_profiles .

- Each profiling result in the repository consists of the .prof file which Snakeviz can be run on to see the profiling results and the png image which is the output of gprof2dot.

# Changes Made So Far (1)

1. ErrMap Function ([https://gitlab.com/hepcedar/contur/-/commit/5df35e7e2e3419fb650e3a2611338a9119865425](https://gitlab.com/hepcedar/contur/-/commit/5df35e7e2e3419fb650e3a2611338a9119865425) ) : We were calling the ErrMap function inside a for-loop, the update calls ErrMap just once outside the loop, saves the results in a dictionary and this dictionary is then called in the loop.

2. Test to statistic to confidence level ([https://gitlab.com/hepcedar/contur/-/commit/0b80789580cd0529bfdce28d2bae252755c45a63](https://gitlab.com/hepcedar/contur/-/commit/0b80789580cd0529bfdce28d2bae252755c45a63) ):  To convert the test statistic first to a p value we use scipy's survival function. We were passing test statistics to this function individually, the change first collects test statistics into a numpy array and passes this numpy array into the scipy survival function.

# Changes Made So Far (2)

3.  Sort blocks, replacing list with dataframe (https://gitlab.com/hepcedar/contur/-/commit/407b0618b8538fe6491a5358b078900e7f01d1e9 ): For each pool we wanted the maximum confidence level of the analysis objects in the pool. We were doing this with a list and an if statement to restrict ourselves to analysis objects in a specified pool and then taking the max of the list. We were running this list concatenation inside another loop though which made things slow. Instead now we start with a dataframe of analysis objects pools and confidence level, group the dataframe by pool (using pandas groupby method) taking the max confidence level by pool.

# Ongoing work

1. The change we made in update 2, passing numpy arrays instead of individual test statistics to the survival function can be taken further. Right now for each analysis object in a yoda file we will call this survival function three times, so for a whole yoda file we call it c.a. 3xnum_analysis_objects, two possible changes we can make here:

   - Reduce the number of calls in each analysis object from 3 to either 2 or 1, gains from doing this are likely very small, but may makes things cleaner;
   - Add an extra dimension to the numpy array storing test statistics for each analysis object in the numpy array and then passing this 3 dimensional numpy array to calculate p values for all yoda files in one call to the survival. This change likely offers a material performance improvement but requires more work to implement.
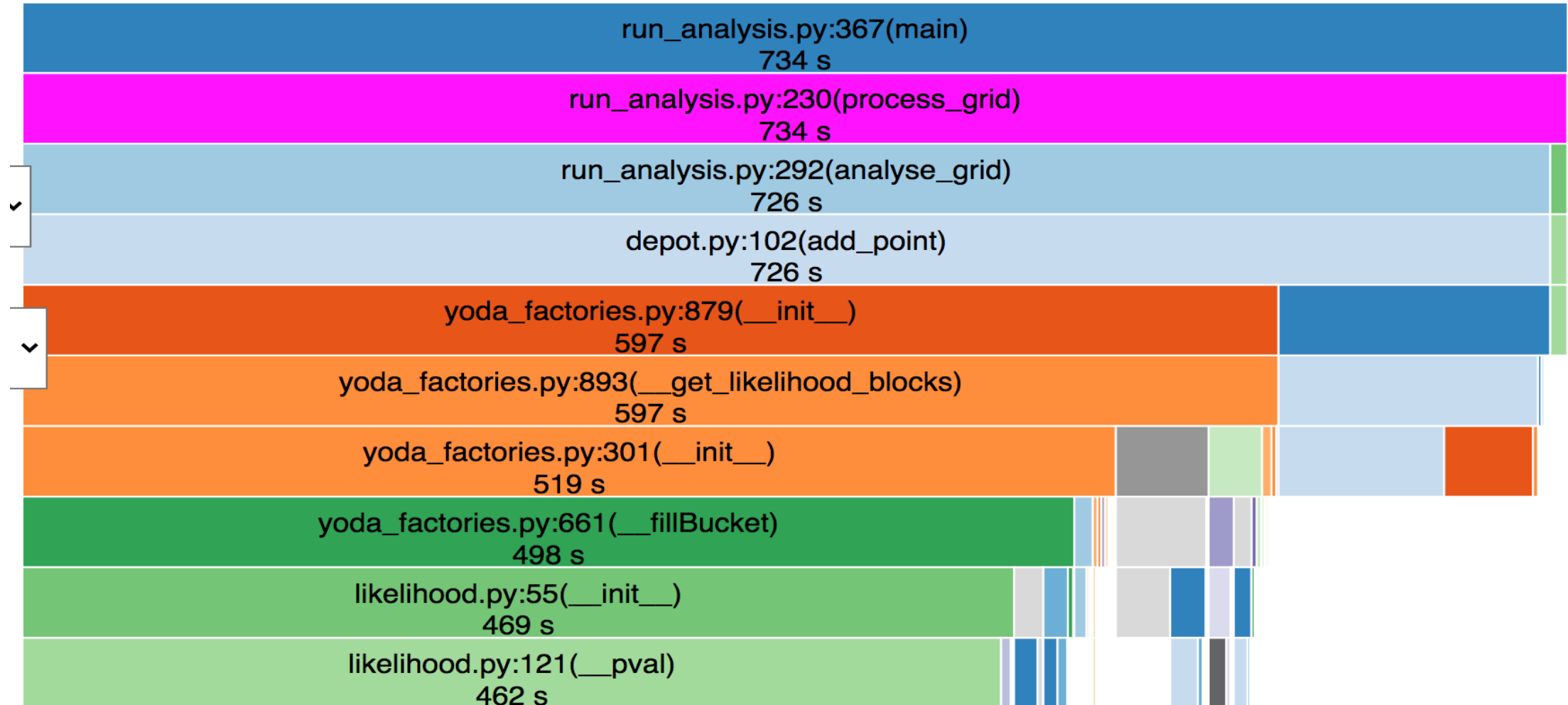
# Ongoing work

2. Upgrade the suite of tests within Contur, this particularly important to ensure nothing is unintentionally broken when updating code to improve run time.

3. Look into parallelising the Contur grid run across yoda files. Currently the grid run effectively loops through each yoda file and combines results at the end into a single .map file.
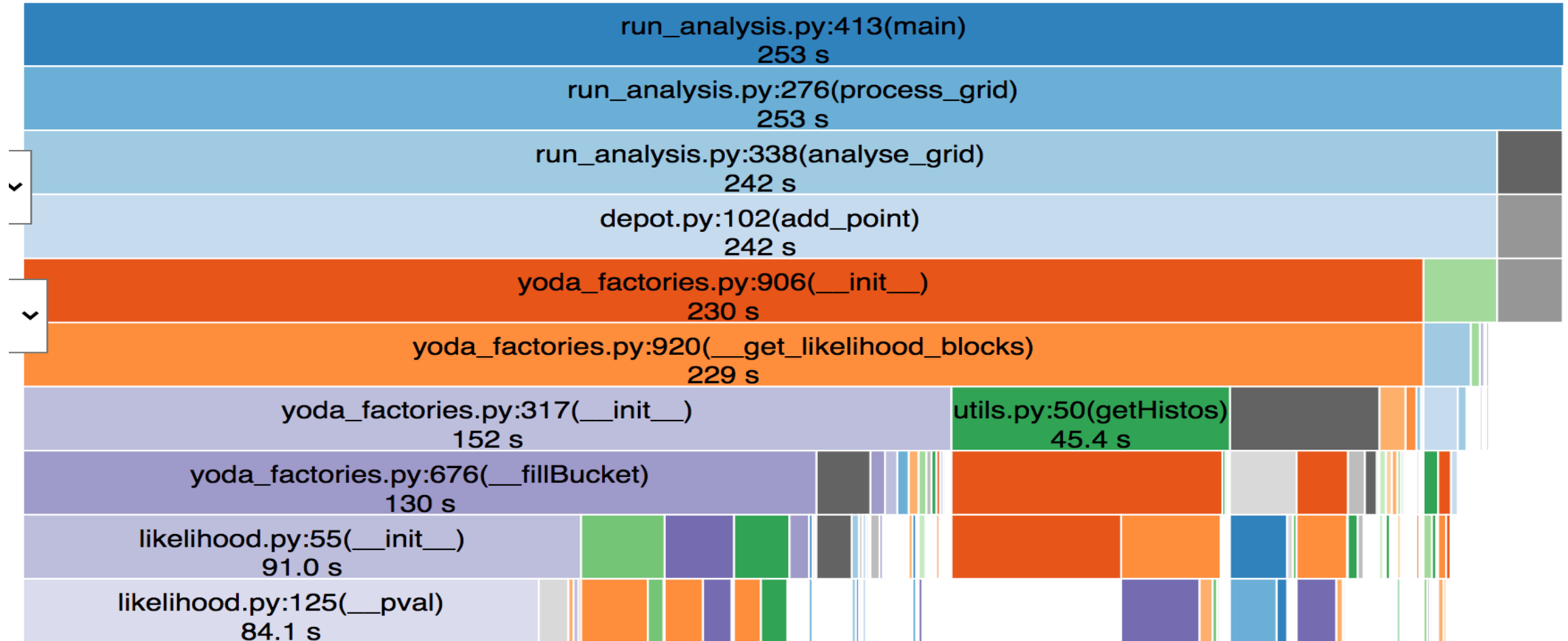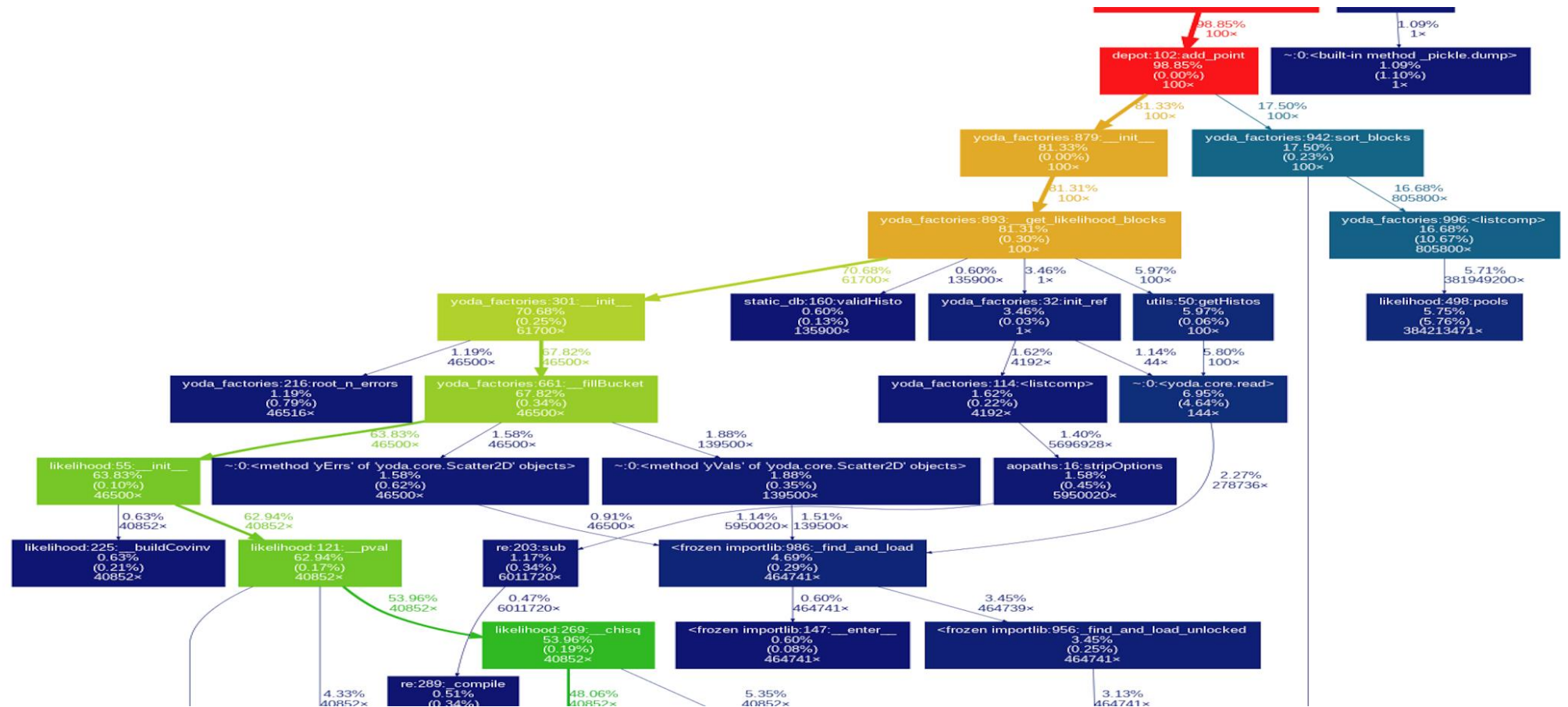
# Annex

# Grid Run (100 yoda files) Icicle Chart **before** change 2 and 3

# Grid Run (100 yoda files) Icicle Chart after change 2 and 3

# Grid Run (100 yoda files) gprof2dot Chart **before** change 2 and 3

# Grid Run (100 yoda files) gprof2dot Chart <u>after</u> change 2 and 3