

## Upcoming changes

A new reservation has been created, "safe". Jobs running within this reservation will be safe from automatic job termination. However, this reservation will be intentionally under-provisioned, so expect longer queue times. Add `--reservation=safe` to your slurm invocation to use the reservation. With the creation of this reservation, we will apply stricter rules to jobs outside of this reservation. Within this reservation we may impose limits on job width.

- Job termination will no longer be based on average wattage since the start of the job. Instead, we will use an exponential moving average [1].
- The power limit will be gradually increased to at least 100W.
- Interactive jobs have so far been fully exempted from automatic job termination. In the future, interactive jobs will only be exempt from automatic job termination for the first 8 hours of wall time.

The `safe` reservation is not intended to be a long-term solution for projects, but rather a stop-gap solution while looking to improve code. Efficient use of resources is considered by the Berzelius allocation staff.

[1] We will leave the exact parameters open for us to adjust. To start with, we aim for the "step function" for a job going from a full load down to idle to allow for the job to continue running for approximately one hour.

## Motivation

We are updating the scheduling policy in this manner based on metrics, experiences and user interactions collected since we started working on improving efficiency. So far, we see a measurable and significant improvement in how the system is used. These changes are intended to mitigate limitations imposed on users, as well as allowing for automatic job termination in some common situations where this hasn't been done so far.

The old scheduling policy places a small subset of users in a position where they can't work at all. This is partially mitigated by the `1g.10gb`-reservation, but not all users are able to use that. The `safe` reservation is intended to mitigate this, by providing a manner in which the job will always run safely. However, to provide incentive to use the system efficiently, the size of this reservation will be limited. That means that queue-times will be artificially and intentionally longer.

Basing job termination on average use works fine in most cases. "Delayed starts", where nothing happens, are terminated after one hour, as are jobs that simply don't manage to saturate the GPU. However, for "forgotten" jobs, where the GPU was used for a few hours and then went idle, we need a different averaging function to allow for a faster decay. In practice, we don't think this change will affect users noticeably.

A common pattern for interactive jobs is forgotten sessions when users allocate resources, use them for a while and then forget to terminate the job. So far, we have intentionally avoided interactive jobs at all. In the future we will terminate them according to the same policy as other jobs, but interactive jobs will have a grace time of 8 hours (= a full workday), instead of the normal one hour.

As always, please contact [berzelius-support@nsc.liu.se](mailto:berzelius-support@nsc.liu.se) with any questions or comments.