



TECHNICAL GUIDELINES FOR APPLICANTS TO PRACE 18th CALL (Tier-0)

The contributing sites and the corresponding computer systems for this call are:

| System | Architecture | Site (Country) | Core Hours (node hours) | Minimum request |
|---------------------------|---------------------|-----------------------|------------------------------------|--------------------------------------|
| Joliot Curie - SKL | BULL Sequana X1000 | GENCI@CEA (FR) | 142 million (3 million) | 15 million core hours |
| Joliot Curie - KNL | BULL Sequana X1000 | GENCI@CEA (FR) | 101 million (1.5 million) | 15 million core hours |
| JUWELS | BULL Sequana X1000 | GCS@JSC (DE) | 70 million (1.5 million) | 35 million core hours |
| Marconi-Broadwell | Lenovo System | CINECA (IT) | 36 million (1 million) | 15 million core hours |
| Marconi-KNL | Lenovo System | CINECA (IT) | 610 million (9 million) | 30 million core hours |
| MareNostrum | Lenovo System | BSC (ES) | 240 million (5 million) | 30 million core hours |
| Piz Daint | Cray XC50 System | CSCS (CH) | 510 million (7.5 million) | 68 million core hours Use of GPUs |
| SuperMUC-NG | Lenovo ThinkSystem | GCS@LRZ (DE) | 125 million (2.2 million) | 35 million core hours |

The site selection is done together with the specification of the requested computing time by the two sections at the beginning of the online form. The applicant can choose one or several machines as execution system, as long as proper benchmarks and resource request justification are provided on each of the requested systems. The parameters are listed in tables. The first column describes the field in the web online form to be filled in by the applicant. The remaining columns specify the range limits for each system.

A - General Information on the Tier-0 systems available for PRACE 18th Call

| | | <i>JOLIOT CURIE - SKL</i> | <i>JOLIOT CURIE - KNL</i> | <i>JUWELS</i> | <i>Marconi Broadwell</i> | <i>Marconi KNL</i> | <i>MareNostrum 4</i> | <i>Piz Daint</i> | <i>SuperMUC-NG</i> |
|-------------|-------------------------|----------------------------------|---------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|--|--|
| System Type | | Bull Sequana | Bull Sequana | Bull Sequana | Lenovo System NeXtScale | Lenovo System Adam Pass | Lenovo | Hybrid Cray xC50 | Lenovo ThinkSystem |
| Compute | Processor type | Intel Xeon Platinum 8168 2.7 GHz | Intel Knights Landing | Intel Xeon Skylake Platinum 8168 | Intel Broadwell | Intel Knights Landing | Intel Xeon Platinum 8160 2.1 GHz | Intel® Xeon® E5-2690 v3 @ 2.60GHz (12 cores) | Intel Skylake EP |
| | Total nb of nodes | 1656 | 828 | 2511 | 720 | 3 600 | 3 456 | 5 320 | 6 480 |
| | Total nb of cores | 79 488 | 52 992 | 120 528 | 25 920 | 244 800 | 165 888 | 63 840 | 31 1040 |
| | Nb of accelerators/node | n.a. | n.a. | n.a. | n.a. | n.a. | n.a. | 1 GPU per node | n.a. |
| | Type of accelerator | n.a. | n.a. | n.a. | n.a. | n.a. | n.a. | NVIDIA® Tesla® P100 16GB | n.a. |
| Memory | Memory / Node | 192 GB DDR4 | 96 GB DDR4 + 16 GB MCDRAM | 96 GB | 128 GB - DDR4 | 96 GB – DDR4 + 16 GB - MCDRAM | 96 GB (200 nodes with 384GB) | 64 GB | 96 GB |
| Network | Network Type | Infiniband EDR | BULL BXI | InfiniBand EDR | Intel Omni-Path Architecture 2:1 | Intel Omni-Path Architecture 2:1 | Intel Omni-Path Architecture | Cray Aries | Intel Omni-Path Architecture |
| | Connectivity | Fat Tree | Fat Tree | Fat Tree | Fat Tree | Fat Tree | Fat Tree | Dragonfly | Fat tree within island (512 nodes) pruned tree between islands |

| | | <i>JOLIOT CURIE</i> | <i>JUWELS</i> | <i>Marconi Broadwell</i> | <i>Marconi KNL</i> | <i>MareNostrum 4</i> | <i>Piz Daint</i> | <i>SuperMUC-NG</i> |
|---------------------------|-------------|---------------------|---------------|--------------------------|--------------------|----------------------|------------------|--------------------|
| Home file system | type | NFS | GPFS | GPFS | GPFS | GPFS | GPFS | NAS * |
| | capacity | TBA | 2.3 PB | 200 TB | 200 TB | 32 TB | 86 TB | 256 PB |
| Work file system | type | Lustre | n.a. | GPFS | GPFS | GPFS | GPFS | GPFS * |
| | capacity | TBA | n.a. | 7.1 PB | 7.1 PB | 4.3 PB | 5.7 PB | 20 PB |
| Scratch file system | type | Lustre | GPFS | GPFS | GPFS | GPFS | Lustre | GPFS * |
| | capacity | 4.6 PB | 9.1 PB | 2.5 PB | 2.5 PB | 8.7 PB | 6.2 PB | 30 |
| Archive | capacity | Unlimited | On demand | On demand | On demand | NA | NA | On demand |
| Minimum required job size | Nb of cores | 1 024 | | 540 | 2 040 | 1 024 | 6 nodes | 512 |

IMPORTANT REMARKS:

- Applicants are strongly advised to apply for PRACE Preparatory Access to collect relevant benchmarks and technical data for the system they wish to use through Project Access (note: if requesting resources on Piz Daint, it is mandatory to show benchmarking on this system in your submission). Further information and support from high performance computing (HPC) Technical teams can be requested during the preparation of the application through PRACE Peer Review at peer-review@prace-ri.eu or directly at the centres.
- Please contact the peer review office of PRACE at peer-review@prace-ri.eu in order to request assistance from the high-level support team, at least 1 month before the submission deadline.

More details on the website of the centers:

JOLIOT CURIE:

<http://www-hpc cea.fr/en/complexe/tgcc-irene.htm>

JUWELS:

http://www.fz-juelich.de/ias/jsc/EN/Expertise/Supercomputers/JUWELS/JUWELS_node.html

Marconi:

<http://www.hpc.cineca.it/hardware/marconi>

<https://wiki.u-gov.it/confluence/display/SCAIUS/UG3.1%3A+MARCONI+UserGuide>

MareNostrum:

<https://www.bsc.es/marenostrum/marenostrum/technical-information>

<https://www.bsc.es/user-support/mn4.php>

Piz Daint:

<http://user.cscs.ch/index.html>

SuperMUC-NG:

<https://www.lrz.de/services/compute/supermuc/supermuc-ng/>

Subsection for each system

JOLIOT CURIE, GENCI@CEA

JOLIOT CURIE is a BULL Sequana system X1000 based on 9 compute cells integrated into 2 partitions:

- **The SKL partition** is composed of 6 cells, each containing 272 compute nodes with two 24-core Intel Skylake Platinum 8168 processors 2.7 GHz, 4 GB/core (192 GB/node). These 6 cells are interconnected by an Infiniband EDR 100 Gb/s high speed network;
- **The KNL partition** is composed of 3 cells, each containing 276 nodes with one Intel Knights Landing 68-core 7250 1.4 GHz manycore processor with 16 GB of high-speed memory (MCDRAM) and 96 GB of main memory. These 3 cells are interconnected by a BULL BXI 100 Gb/s high speed network. A KNL node provides 64 cores for user jobs and keeps 4 cores for the system. A node is configured in quadrant for the cluster node and in cache mode for the memory. The threadmultiple mode for the hybrid program is not yet supported.

This configuration is completed with 5 fat nodes for pre/post processing (3 TB of memory each and a fast local storage based on NVMe) and 20 hybrid nodes used for remote visualisation.

These resources are federated across a multi-layer shared Lustre parallel filesystem with a first level (/scratch) of more than 5 PB at 300 GB/s.

The peak performance of this system is 9 petaflops.

A second tranche is planned to be installed during Q2 2019.

JUWELS, GCS@JSC

JUWELS (Jülich Wizard for European Leadership Science) is designed as a modular system. The JUWELS Cluster module, supplied by Atos, based on its Sequana architecture, consists of about 2500 compute nodes, each with two Intel Xeon 24-core Skylake CPUs and 96 GiB of main memory. The compute nodes are interconnected with a Mellanox EDR InfiniBand interconnect. The peak performance of this CPU based partition is 10.4 petaflops. A booster module, optimized for massively parallel workloads, is currently scheduled for the beginning of 2020.

Marconi, CINECA

Marconi system consists of three partitions, from which two will be available for Call 18:

- **Marconi - Broadwell** consists of 10 Lenovo NeXtScale racks with 72 nodes per rack. Each node contains 2 Broadwell processors each with 18 cores and 128 GB of DDR4 RAM.¹
- **Marconi - KNL** consists of 3600 Intel server nodes integrated by Lenovo. Each node contains 1 Intel

¹ Marconi system will be upgraded in 2019. This upgrade will affect particularly the Broadwell partition. After this upgrade the network infrastructure of this partition could be different.

Technical Guidelines for Applicants 18th Call PRACE Project Access

Knights Landing processor with 68 cores, 16 GB of MCDRAM and 96 GB of DDR4 RAM.

All the nodes (for both systems) are connected via Intel Omni-Path network.

The aggregate peak performance of this system is more than 20 petaflops.

MareNostrum 4, BSC

MareNostrum 4 consists of 48 Compute Racks with 72 compute nodes per rack. Each node has two Intel Xeon Platinum 8160 processors with 2.1 GHz, 24 cores per socket (48 cores/node) and 96 GB of main memory (2 GB/core), connected via Intel Omni-Path fabric at 100 Gbits/s.

There are a subset of 200 fat nodes available that have 384 GB of main memory (8 GB/core). Their use is restricted to a maximum of 50% of their hours for all projects combined during each PRACE call.

Piz Daint, CSCS

Named after Piz Daint, a prominent peak in Grisons that overlooks the Fuorn pass, this supercomputer is a hybrid Cray XC50 system and is the flagship system for national HPC Service. The compute nodes are equipped with Intel® Xeon® E5-2690 v3 @ 2.60GHz (12 cores) and NVIDIA® Tesla® P100 16GB, and 64 GB of host memory.

The nodes are connected by the "Aries" proprietary interconnect from Cray, with a dragonfly network topology.

SuperMUC-NG, GCS@LRZ

SuperMUC-NG will start operation in Q1/2019. It will provide 6480 Lenovo ThinkSystem dual socket nodes equipped with 24 core Intel Skylake EP processors and 96 GB of main memory. A subset of 144 fat nodes will be equipped with 768 GB of main memory. The nodes are connected via a fat-tree Omni-Path network. The Peak Performance will be at 26.7PF.

B – Guidelines for filling-in the online form

Resource Usage

Computing time

To apply for PRACE Tier-0 resources there is a minimum amount of core hours for all systems. Proposals which do not comply with this requirement should apply in the Tier-1 national calls.

The amount of computing time has to be specified in core hours (or alternatively node hours) (wall clock time [hours]*physical cores (nodes) of the machine applied for). It is the total number of core (node) hours to be consumed within the twelve months period of the project.

Please justify the number of core (node) hours you request by providing a detailed work plan and the appropriate technical data on the systems of interest. Applicants are strongly invited to apply for PRACE Preparatory Access.

Once allocated, the project has to be able to start immediately and is expected to use the resources continuously.

When planning for access, please take into consideration that the effective availability of the system is about 80 % of the total availability, due to queue times, possible system maintenance, upgrade and data transfer time.

Tier-0 proposals are required to respect the minimum and maximum request of resources as indicated in the Terms of Reference to be found [here](#).

Job Characteristics

This section describes technical specifications of simulation runs performed within the project.

Wall Clock Time

A simulation consists in general of several jobs. The wall clock time for a simulation is the total time needed to perform such a sequence of jobs. This time could be very large and could exceed the job wall clock time limits on the machine. **In that case the application has to be able to write checkpoints and the maximum time between two checkpoints has to be less than the wall clock time limit on the specified machine.**

| <i>Field in online form</i> | <i>Machine</i> | <i>Max</i> |
|--|--|--|
| Wall clock time of one typical simulation (<u>hours</u>) <number> | SuperMUC-NG Marconi Other systems Piz Daint | 2500 hours * 4000 hours * < 10 months - |
| Able to write checkpoints <check button> | All | Yes (apps checkpoints) |

Technical Guidelines for Applicants 18th Call PRACE Project Access

| | | |
|---|----------------------------------|----------------------|
| Maximum time between two checkpoints (= maximum wall clock time for a job) (hours) <number> | SuperMUC-NG All other systems | 48 hours 24 hours |
|---|----------------------------------|----------------------|

* This is the time a job really can use the CPUs. This limited time is mainly due to waiting times in the queue especially for jobs with a limited scalability (using less than 10,000 cores).

** This might be changed during project runtime, guaranteed minimum is the value in brackets.

Number of simultaneously running jobs

The next field specifies the number of independent runs which could run simultaneously on the system during normal production conditions. This information is needed for batch system usage planning and to verify if the proposed work plan is feasible during project run time.

| Field in online form | Machine | Max |
|--|--|---|
| Number of jobs that can run simultaneously <number> | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | 25 (1 024 cores), 4 (8 192 cores) 25 (1 024 cores), 4 (8 192 cores) for the SKL partition and 10 (1 024 cores), 2 (8 192 cores) for the KNL partition 3 (more on demand) 2-20 (depending on the job size) Dynamic* No shared nodes: 1 job per node maximum tba |

* Depending on the amount of PRACE projects assigned to the machine, this value could be changed.

Job Size

The next fields describe the job resource requirements, which are the number of cores (nodes) and the amount of main memory. These numbers have to be defined for three different job classes (with minimum, average, or maximum number of cores/nodes).

Please note that the values stated in the table below are absolute minimum requirements, allowed for small jobs, which should only be applicable to a small share of the requested computing time. **Typical production jobs should run at larger scale.**

Job sizes must be a multiple of the minimum number of cores (nodes) in order to make efficient use of the architecture.

IMPORTANT REMARK

*Please provide explicit scaling data of the codes you plan to work with in your project at least up to the minimum number of physical cores required by the specified site (see table below) using input parameters comparable to the ones you will use in your project (a link to external websites, just referencing other sources or "general knowledge" is not sufficient). **Generic scaling plots provided by vendors or developers do not necessarily reflect the actual code behavior for the simulations planned. Scaling benchmarks need to be representative of your study case and need to support your resource request on every system of interest (application to PRACE preparatory access project is strongly recommended and mandatory on Piz Daint). Missing technical (scaling, etc.) data may result in rejection of the proposal.***

Technical Guidelines for Applicants 18th Call PRACE Project Access

| <i>Field in online form</i> | <i>Machine</i> | <i>Min (cores)</i> |
|---|--|---|
| Expected job configuration (Minimum) <number> | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | 1 024 4 608 540 (Broadwell) 2 040 (KNL) 1 024 6 nodes 960 |
| Expected number of cores (Average) <number> | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | 4 096 9 216 540 (Broadwell), 4 080 (KNL) 4 096 6 to 2 400 nodes 12288 |
| Expected number of cores (Maximum) <number> | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | 40 000 (full machine on demand) 40 000 for the SKL partition and 20 000 for the KNL partition (full machine on demand) 24 576 (full machine on demand) 3 240 (Broadwell), 68 000 (KNL) 32 000 (for exceptional applications the usage of the full machine is possible) 4 400 nodes 147 456 (half machine; full machine on demand) |

Virtual cores (SMT is enabled) are not counted. Accelerator based systems (GPU, Xeo, Phi, etc.) need *special rules*.

Additional information:

Marconi

The minimum number of (physical) cores per job is 2040 on KNL partition and 540 on Broadwell.

However, this minimum requirement should only be requested for a small share of the requested computing time and it is expected that PRACE projects applying for Marconi can use at least a doubled value on average (some exceptions could be made for Broadwell partition, especially if the applicants will be able to launch several jobs at the same time: in this case a minimum number of 540 cores per job will be considered as sufficient).

The maximum number of (physical) cores per job is 68000 (on Broadwell partition is 3240). Larger jobs are possible in theory (only if requested sending an email to superc@cinca.it), but the turnaround time is not guaranteed.

Please provide explicit scaling data of the codes you plan to work within your project. On KNL partition the scaling behavior at least up to 4080 physical cores must be shown and it must demonstrate a good scalability at least up to 2040 physical cores. In any case, the applicant will have also to demonstrate that their codes can scale on a KNL system equipped with an Omnipath network at least up to the maximum

Technical Guidelines for Applicants 18th Call PRACE Project Access

number of required cores using a setup having size similar to the one proposed for the project. For hybrid (MPI+OpenMP) codes it is strongly recommended that applicants show scaling data for different numbers of threads per task in order to exploit the machine most efficiently. Providing such kind of data will be favorably considered for the technical evaluation of the project.

Proving to be able to scale on other architectures will be favorably considered for the technical evaluation of the project.

Piz Daint

Technical data needs to be provided on the Cray XC50, Piz Daint. To apply for Piz Daint use of **GPUs is a must**. Scalability, performance and technical data have to be sufficient to justify the resource request (≥ 1 million node hours). All technical data on Piz Daint must be in **node hours**.

SuperMUC-NG

The minimum number of (physical) cores per job is 960. However, it is expected that PRACE projects applying for this system can use more than 6144 physical cores per job. When running several jobs simultaneously filling complete islands should be possible, but this is not mandatory.

Job Memory

The next fields are the total memory usage over all cores of jobs.

| <i>Field in online form</i> | <i>Machine</i> | <i>Max</i> |
|---|--|---|
| Memory (Minimum job) <number> | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | <ul style="list-style-type: none"> - Jobs should use a substantial fraction of the available memory - No requirement - <118 GB per node (Broadwell), <86 GB per node (KNL) - 2 GB * #cores or 8 GB * #cores (max 200 fat nodes) - nodes are not shared - Jobs should use a substantial fraction of the available memory |
| Memory (Average job) <number> | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | <ul style="list-style-type: none"> - Jobs should use a substantial fraction of the available memory - No requirement - <118 GB per node (Broadwell); <86 GB per node (KNL) - 2 GB * #cores or 8 GB * #cores (max 200 fat nodes) - 5.3 GB per core or 64 GB per node (nodes are not shared) - Jobs should use a substantial fraction of the available memory |

Technical Guidelines for Applicants 18th Call PRACE Project Access

| Memory (Maximum job <number> | | |
|---|---------------|--|
| | Joliot-Curie | - 4 GB* #cores (or 64 GB * #nodes for CURIE and 192 GB * # nodes for Joliot Curie) |
| | JUWELS | - <92 GB per node |
| | Marconi | - <118 GB per node (Broadwell), <86 GB per node (KNL) |
| | MareNostrum 4 | - 2 GB* #cores or 8 GB * #cores (max 200 fat nodes) |
| | Piz Daint | - 5.3 GB per core or 64 GB per node (nodes are not shared) |
| | SuperMUC-NG | - 2.0 GB* #cores or 96 GB* #nodes) |

The memory values include the resources needed for the operating system, i.e. the application has less memory available than specified in the table.

Storage

General remarks

The storage requirements have to be defined for four different storage classes (Scratch, Work, Home and Archive).

- **Scratch** acts as a temporary storage location (job input/output, scratch files during computation, checkpoint/restart files; no backup; automatic remove of old files).
- **Work** acts as project storage (large results files, no backup).
- **Home** acts as repository for source code, binaries, libraries and applications with small size and I/O demands (source code, scientific results, important restart files; has a backup).
- **Archive** acts as a long-term storage location, typically data reside on tapes. For PRACE projects also archive data have to be removed after project end. The storage can only be used to backup data (simulation results) during project's lifetime.

Data in the archive is stored on tapes. **Do not store thousands of small files in the archive, use container formats** (e.g. tar) to merge files (**ideal size of files: 500 – 1 000 GB**). Otherwise, **you will not be able to retrieve back the files from the archive within an acceptable period of time** (for retrieving one file about 2 minutes time (independent of the file size!) + transfer time (dependent of file size) are needed)!

IMPORTANT REMARK

All data must be removed from the execution system within 2 (6 on Marconi) months after the end of the project.

Total Storage

The value asked for is the maximum amount of data needed at a time. Typically, this value varies over the project duration of 12 month (or yearly basis for multi-year projects). **The number in brackets in the "Max per project" column is an extended limit, which is only valid if the project applicant contacted the center beforehand for approval.**

Technical Guidelines for Applicants 18th Call PRACE Project Access

| Field in online form | Machine | Max per project | Remarks |
|---|--|--|---|
| Total storage (Scratch) <number> Typical use: Scratch files during simulation, log files, checkpoints Lifetime: Duration of jobs and between jobs | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | 20 TB (100 TB) 20 TB (100 TB) 20 TB (100 TB)* ¹ 100 TB (more on demand) 6.2 PB 100 TB (200 TB) | - without backup, automatic clean-up procedure; - without backup, files older than 90 days will be removed automatically - without backup, clean-up procedure for files older than 50 days; - without backup, clean-up procedure. - without backup, clean-up procedure. Quota on inodes - without backup, automatic clean-up procedure. |
| Total storage (Work) <number> Typical use: Result and large input files Lifetime: Duration of project | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | 1 TB n.a. 20 TB (100 TB)* ¹ 10 TB (100TB) 250 TB (500 TB)* ² 100 TB (200 TB) | - Without backup; - With backup; - Input not readable from compute nodes; data kept only for duration of project - Without backup. |
| Total storage (Home) <number> Typical use: Source code and scripts Lifetime: Duration of project | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | 3 GB 6 TB 50 GB 20 GB 10 GB x User 100 GB | - with backup and snapshots; - with backup; - with backup; - with backup; - with backup and snapshots - with backup and snapshots. |
| Total storage (Archive) <number> | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | 100 TB * ³ 20 TB (100 TB)* ⁴ n.a. n.a. 100 TB* ⁵ | - File size > 1 GB - Ideal file size: 500 GB – 1000 GB - Typical file size should be > 5 GB |

*¹ The default value is 1 TB. Please ask to CINECA User Support (superc@cineca.it) to increase your quota after the project will start.

*² From 250 to maximum of 500 TB will be granted if the request is fully justified and a plan for moving the data is provided.

*³ Access to JUWELS' archive needs a special agreement with JSC and PRACE.

*⁴ Not active by default. Please ask to CINECA User Support after the project will start.

*⁵ Long-term archiving or larger capacity must be negotiated separately with LRZ.

When requesting more than the specified scratch disk space and/or larger than 1 TB a day and/or storage of more than 4 million files, please justify this amount and describe your strategy concerning the handling of data (pre/post processing, transfer of data to/from the production system, retrieving relevant data for long-term). If no justification is given the project will be proposed for rejection.

Technical Guidelines for Applicants 18th Call PRACE Project Access

If you request more than 100 TB of disk space, please contact peer-review@prace-ri.eu before submitting your proposal in order to check whether this can be realized.

Number of Files

In addition to the specification of the amount of data, the number of files also has to be specified. If you need to store more files, **the project applicant must contact the center beforehand for approval.**

| <i>Field in online form</i> | <i>Machine</i> | <i>Max</i> | <i>Remarks</i> |
|---|--|--|--|
| Number of files (<u>Scratch</u>) <number> | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | 2 Million 4 Million 2 Million 2 Million 1 Million 1 Million | - 10 000 files max per directory, without backup, files older than 90 days will be removed automatically - Without backup, files older than 90 days will be removed automatically - Without backup, files older than 50 days will be removed automatically - No limit while running, but limits in number of files left on scratch. - Without backup, old files are removed automatically, Ideal file size: >100 GB. |
| Number of files (<u>Work</u>) <number> | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | 500 000 n.a. 2 Million 2 Million 50 000 per TB 1 Million | - Extensible on demand, 10 000 files max per directory - with backup and snapshots - Ideal file size: >100 GB |
| Number of files (<u>Home</u>) <number> | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | n.a. 12 Million 100 000 100 000 100 000 100 000 | - With backup - With backup - With backup and snapshots - With backup (snapshots) |
| Number of files (<u>Archive</u>) <number> | Joliot-Curie JUWELS Marconi MareNostrum 4 Piz Daint SuperMUC-NG | 100 000 100 000 10 000* n.a. n.a. 100 000 | - Extensible on demand, typical file size should be > 1 GB - Ideal file size: 500 GB – 1000 GB - Typical file size should be > 5 GB |

*HSM has a better performance with a small amount of very big files

Data Transfer

For planning network capacities, applicants have to specify the amount of data which will be transferred from the machine to another location. Field values can be given in Tbyte or Gbyte.

Reference values are given in the following table. *A detailed specification would be desirable: e.g. distinguish between home location and other PRACE Tier-0 sites.*

Please state clearly in your proposal the amount of data which needs to be transferred after the end of your project to your local system. Missing information may lead to rejection of the proposal.

Be aware that transfer of large amounts of data (e.g. tens of TB or more) may be challenging or even unfeasible due to limitations in bandwidth and time. Larger amounts of data have to be transferred continuously during project's lifetime.

Alternative strategies for transferring larger amounts of data at the end of projects have to be proposed by users (e.g. providing tapes or other solutions) and arranged with the technical staff.

| <i>Field in online form</i> | <i>Machine</i> | <i>Max</i> |
|--|----------------|--------------------|
| Amount of data transferred to/from production system <number> | Joliot-Curie | 100 TB |
| | JUWELS | 100 TB |
| | Marconi | 20 TB* |
| | MareNostrum 4 | 50 TB |
| | Piz Daint | Currently no limit |
| | SuperMUC-NG | 100 TB |

* More is possible, but this needs to be discussed with the site prior to proposal submission.

If one or more specifications above is larger than a reasonable size (e.g. more than tens of TB data or more than 1TB a day) the applicants must describe their strategy concerning the handling of data in a separate field (pre/post-processing, transfer of data to/from the production system, retrieving relevant data for long-term). In such a case, the application is *de facto* considered as I/O intensive.

I/O

Parallel I/O is mandatory for applications running on Tier-0 systems. Therefore, the applicant must describe how parallel I/O is implemented (checkpoint handling, usage of I/O libraries, MPI I/O, netcdf, HDF5 or other approaches). Also the typical I/O load of a production job should be quantified (I/O data traffic/hour, number of files generated per hour).