

Glossary of Workflow Terms

General Definitions

Participant: Participant is an object that transforms inputs into outputs. Example of participants are: dCache (dccb), PBS (qsub), user applications and shell scripts. All Participants are considered to be atomic operations from the executing workflow's point of view.

Participant Product: Participant product is the output generated by a participant including provenance data (participant inputs and parameters).

Workflow: Workflow are the procedures whereby data and control are passed among participants according to a defined set of rules (e.g. data and control dependencies) to achieve a specific goal.

Workflow management system: A workflow management system manages and executes workflows on computing and storage resources. It is responsible for resolving dependencies, keeping track of data products, scheduling and fault tolerance.

Reliability and Fault Tolerance: A system is said to *fail* when it cannot meet its promises. An *error* is a part of a system's state that may lead to a failure. The cause of error is called a *fault*. *Reliability* is a measure of the continuous service accomplishment (or, equivalently, of the time to failure) from an initial reference point. *Fault tolerance* means that a system can provide its services even in the presence of faults.

Quality of Service: Quality of service is defined as the level of performance in a computing system. For example, to finish a job in a given time limit. A workflow could have different performance constraints, such as time, cost, fidelity, reliability, and security. The primary focus of the LQCD workflow will be *time to complete a participant*.

Checkpoint: Checkpointing is obtaining a snapshot of the system's current state that is stored for later use. Checkpointing could be at the system level or at the application level. An application level checkpoint records the present state of a user job (process).

Restart: Restart is defined as starting from the originally defined starting point.

Resume: Resume is defined as starting from the last checkpoint (or last milestone in terms of LQCD computation).

Workflow Types Definitions

Template: The workflow template is a pattern or parameterized description of how a particular problem is solved.

Ideal: An ideal workflow assumes unlimited resources and it contains only information relevant to solving the problem in the best possible circumstances. In other words, no resource constraints are included.

Instantiated: An instantiated workflow is the result of the application of relevant input files and parameters to a template by a user. It could include multiple versions of the participants.

Workflow Instance = workflow template + applications + input parameters and data + output data

Executable: An executable workflow is produced by applying cluster resource usage policies and constraints to an instantiated workflow by the workflow system. In other words, the workflow graph is transformed from ideal to something that can be run on an existing cluster.

Domain Specific Definitions

Lattice: A lattice is a cubic four dimensional space-time grid. Boundary conditions are normally periodic in the spatial directions and (anti)periodic in the time direction. *Sites*, the vertices of the grid, are referred to by coordinates (x,y,z,t). There are eight *links* connecting a site to its nearest neighbors. Lattices are typically cubic: $L_s = L_x = L_y = L_z$. Some lattice dimensions now in use include $L_s^3 \times L_t = 16^3 \times 48$, $20^3 \times 64$, $28^3 \times 96$, $40^3 \times 96$ and $48^3 \times 128$.

(Gauge) Configuration: A gauge configuration is a four dimensional (space-time) snapshot of the *gluon* field. On each *link* of the lattice, the gluon variable is represented by an SU(3) (complex 3 x 3) matrix. Configurations are stored in configuration files containing $3 \times 3 \times \text{sizeof}(\text{complex}) \times 4 \times L_s^3 \times L_t$ bytes.

Ensemble: An ensemble is an ordered collection of gluon configurations sharing the same physics parameters e.g. lattice spacing (or QCD coupling strength), number of sea quarks and sea quark masses. Configurations are generated in a Markov sequence. At each step, the last configuration serves as the starting gluon configuration which is evolved forward in simulation time by Monte Carlo techniques. At regular intervals in Monte Carlo simulation time, gluon configuration snapshots are saved. Configurations within an ensemble are labeled by a unique sequence number or *configuration number* which is typically the number of steps in simulation time. An ensemble may contain a fork where more than one Monte Carlo evolution sequence was started from the same input configuration. Forks are identified by a unique *series label*.

Configuration generation: Configuration generation is a campaign describing the creation of an *ensemble*. Each step of the workflow(s) of this ensemble generates a gauge configuration which depends upon the output gauge configuration of the previous step.

Quark Propagator: The quark propagator is the field that describes how a quark propagates or hops from site-to-site of the lattice. A quark is represented by complex matrix at every site of the lattice. *Clover quarks* are represented by a 4 x 4 (spin) x 3 x 3 (color) matrix while *staggered quarks* are represented by a 3 x 3 (color) matrix at every lattice site. Quark propagators are the solution of a sparse matrix problem. The sparse matrix has dimension equal to the number of sites on the lattice times the size of the quark matrix. Conjugate gradient techniques are used to solve for quark propagators. Quark propagators are produced (and saved) as intermediate results during a campaign. During a campaign, quark propagators are computed for every configuration in an ensemble for every given set of physics input parameters e.g. quark mass, quark source type.

Meson 2- and 3-point functions: Meson 2- and 3-point functions are generated by connecting together quark propagators to form operators that create and annihilate mesons. These n-point functions describe how mesons, particles composed of quarks, propagate on the lattice. Meson n-point functions typically have one or more time coordinates fixed and they are normally summed over the spatial dimensions before they are stored. A meson 2-point function is typically stored as L_t complex numbers per configuration. A campaign might generate hundreds of distinct meson n-point functions.

Ensemble average: Ensemble averages are the physical results that are calculated in a statistical analysis of meson *n-point functions* that have been averaged over every configuration in an ensemble.

Campaign: A *campaign* is a coordinated set of calculations aimed at determining a set of specific physics quantities - for example, predicting the mass and decay constant of a specific particle determined by computing *ensemble averaged 2-pt functions*. A typical campaign consists of taking an *ensemble* of vacuum gauge *configurations* and using them to create intermediate data products (e.g. *quark propagators*) and computing *meson n-point functions* for every configuration in the ensemble. An important feature of such a campaign is that the intermediate calculations done for each configuration are independent of those done for every other configuration.

An example campaign could consist of a single workflow, where the intermediate products are used immediately, or it could be broken up into multiple workflows, with the intermediate products stored for later use.

Milestone: A milestone is the persistent state of the last intermediate result reached by a workflow instance. A milestone can be used for recovering a workflow instance or to extend a campaign. A milestone is typically the output of a participant. The milestone has information about the workflow instance that generated it, including parameters and input files that were used (provenance).

Extend Campaign: A previously executed campaign can have new input files added, an act that requires further workflow execution. It could also use milestones from a previous campaign and generate additional output files.

Campaign = workflow instance(s) = template(s) + ensemble + applications + parameters
+ desired outputs