



WEXGrid: A Modular Python Framework for Defining, Managing, and Executing Complex Workflows in Local and LSF Grid Environments

João Marques - 36774

Thesis presented to the School of Technology and Management in the scope of the Master in Informatics.

Supervisors:

Paulo Matos

This document does not include the suggestions made by the board.

Bragança

2024-2025



WEXGrid: A Modular Python Framework for Defining, Managing, and Executing Complex Workflows in Local and LSF Grid Environments

João Marques - 36774

Thesis presented to the School of Technology and Management in the scope of the
Master in Informatics.

Supervisors:
Paulo Matos

This document does not include the suggestions made by the board.

Bragança
2024-2025

Dedication

To my grandmother, Rosa da Silva, whose love, wisdom, and strength have always inspired me.

Acknowledgment

This work was supported by the Instituto Politécnico de Bragança (IPB). I would like to express my deepest gratitude to all those who contributed to the completion of this thesis.

First and foremost, I would like to thank my girlfriend for her love, patience, and constant encouragement throughout this journey. My heartfelt thanks go to my family — especially my sister Cristiana, my mother Fátima, and my father João — for their unconditional support, understanding, and motivation in every step of my academic path.

I am also grateful to my colleagues José Valverde, Duarte Azevedo, Helder Campos, and Ricardo Araújo for their friendship, collaboration, and valuable discussions at work.

Finally, I would like to thank all the professors at IPB, whose dedication and guidance have been fundamental to my learning, and especially my supervisor, Professor Paulo Matos, for his invaluable advice, support, and mentorship during the development of this work.

Abstract

WEXGrid is a modular framework and software tool designed to address the challenges of creating, managing, and executing scientific workflows on local workstations and distributed Load Sharing Facility (LSF) Grid environments. Break workflows into discrete targets with explicit inputs, outputs, execution actions, and resource requirements, all of which are defined programmatically in Python. This ensures that the underlying concepts remain accessible to computational scientists.

The system then builds a dependency graph in which the execution order is automated based on data and control dependencies. This graph can be used to drive dynamic scheduling based on resource availability and task readiness. Cache management reduces redundant computations by validating outputs through timestamps and checksums. This influences scheduling decisions, optimizing resource utilization.

WEXGrid supports parallel and asynchronous task execution, balancing workload distribution across resources with data locality concerns to avoid Input/Output (I/O) bottlenecks. Its fault-tolerant architecture includes mechanisms for detecting and isolating faults to avoid cascading failures. Its provenance capture and metadata interoperability meet the Findable, Accessible, Interoperable, Reusable (FAIR) principles for reproducibility and portability.

WEXGrid aims to enable scalability and maintainability, providing a unified interface for defining and executing workflows for heterogeneous computational infrastructures without sacrificing performance or transparency. Future enhancements include adaptive scheduling informed by real-time telemetry, rich provenance integration, and extended interoperability with cloud and containerized environments. These enhancements will

ensure sustained efficiency and reproducibility across heterogeneous scientific computing platforms.

Keywords: Distributed computing, Parallel execution, LSF Grid environments, Dynamic task scheduling.

Resumo

O WEXGrid é uma estrutura modular e uma ferramenta de "software" concebida para enfrentar os desafios da criação, gestão e execução de fluxos de trabalho científicos em estações de trabalho locais e em ambientes de rede distribuídos LSF. Divide os fluxos de trabalho em alvos discretos com entradas, saídas, ações de execução e requisitos de recursos explicitamente definidos, todos programaticamente especificados em Python. Esta abordagem garante que os conceitos subjacentes permaneçam acessíveis aos cientistas computacionais.

O sistema constrói, então, um grafo de dependências no qual a ordem de execução é automatizada com base nas dependências de dados e de controle. Este grafo pode ser utilizado para orientar o agendamento dinâmico em função da disponibilidade de recursos e da prontidão das tarefas. A gestão de cache reduz cálculos redundantes, validando as saídas através de carimbos de data/hora e somas de verificação. Estes mecanismos influenciam as decisões de agendamento, otimizando a utilização dos recursos.

O WEXGrid suporta a execução paralela e assíncrona de tarefas, equilibrando a distribuição da carga de trabalho entre os recursos e considerando a localização dos dados para evitar gargalos de I/O. A sua arquitetura tolerante a falhas inclui mecanismos de detecção e isolamento de erros, prevenindo falhas em cascata. A captura de proveniência e a interoperabilidade de metadados cumprem os princípios FAIR de reprodutibilidade e portabilidade.

O WEXGrid visa promover a escalabilidade e a manutenção, fornecendo uma interface unificada para a definição e execução de fluxos de trabalho em infraestruturas computacionais heterogêneas, sem comprometer o desempenho ou a transparência. As melhorias

futuras incluem programação adaptativa informada por telemetria em tempo real, integração avançada de proveniência e interoperabilidade alargada com ambientes em nuvem e contentores. Estas evoluções garantirão eficiência e reprodutibilidade sustentadas em plataformas de computação científica heterogéneas.

Palavras-chave: Computação distribuída, Execução Paralela, Ambientes LSF Grid, Agendamento Dinâmico de Tarefas.

Contents

1	Introduction	1
2	Background and Motivation	5
2.1	Evolution of Workflow Management Systems	5
2.2	Challenges in Local and Distributed Workflow Execution	8
2.3	Motivation for Developing WEXGrid	11
3	Assumptions and Requirements	15
3.1	Functional Requirements	15
3.2	Non-Functional Requirements	18
3.3	Operational Environment Assumptions	22
4	Contextualization and Related Work	27
4.1	Scientific Workflow Systems Landscape	27
4.2	Scheduling Strategies in Distributed Environments	31
5	System Architecture	35
5.1	Overview of WEXGrid Architecture	35
5.2	Workflow Creation and Dependency Management	39
5.3	Execution Engine	43
5.4	Advanced Cache Management	47
6	Conclusions	53

List of Figures

1.1	High-level architectural overview of the WEXGrid framework.	4
5.1	Overview of WEXGrid Architecture.	38
5.2	Workflow creation and dependency management in WEXGrid.	40
5.3	Simplified DAG of a hardware verification workflow.	43
5.4	Execution engine in WEXGrid.	46
5.5	Advanced Cache Management in WEXGrid.	49

Acronyms

API Application Programming Interface.

CPU Central Processing Unit.

DAG Directed Acyclic Graph.

DB Data Base.

DSL Domain-Specific Language.

EFO Experimental Factor Ontology.

ESTiG Escola Superior de Tecnologia e Gestão.

FAIR Findable, Accessible, Interoperable, Reusable.

HPC High-Performance Computing.

HTTP HyperText Transfer Protocol.

I/O Input/Output.

IPB Instituto Politécnico de Bragança.

JSON JavaScript Object Notation.

LSF Load Sharing Facility.

MPI Message Passing Interface.

OpenMP Open Multi-Processing.

YAML YAML Ain't Markup Language.

Chapter 1

Introduction

WEXGrid has been designed as a specialized modular framework for the increasing demands of workflow creation, management, and execution in computational science. In essence, it works on the basic principle that any workflow can be decomposed into discrete targets. These targets encapsulate defined inputs, outputs, associated execution actions, and, in more advanced configurations, labels or explicit resource requirements.

The framework allows these targets to be instantiated programmatically via an `add()` method on a central `Flow` object. Once defined, they are integrated into a dependency graph so that execution order is dictated not by user intervention but by systematic dependency analysis. This model is in tight alignment with established practices in scientific workflow systems, yet it deliberately maintains Python-based simplicity to retain accessibility for users already embedded in common computational research workflows. The motivation for WEXGrid's development emerges in part from constraints observed in existing workflow managers across both local workstation environments and large-scale distributed systems such as LSF Grids.

Traditional solutions often mandate verbose configuration files or domain-specific languages that pose steep learning curves for scientists whose primary expertise lies outside computer science. By leveraging Python as its base syntax, WEXGrid builds on familiar

coding paradigms while introducing constructs specific to task orchestration and scheduling logic. Tasks can range from shell commands to Python functions or executables specified within an extensible class hierarchy. Flexibility here is not just a convenience; it offers adaptability in integrating with heterogeneous scientific pipelines without re-engineering code for varying computational contexts. Efficient resource allocation is at the heart of WEXGrid's design philosophy.

The GridSlot mechanism and Resource Manager components function in concert to dynamically allocate and free resources during execution cycles. This ensures optimal use of available capacity whether the system runs locally on a single node or distributes tasks across an LSF Grid environment. Unlike static scheduling approaches, which may over-commit resources or leave them idle awaiting dependencies, WEXGrid assesses availability in concert with dependency readiness to enqueue tasks judiciously for execution. One distinctive capability is the integration of asynchronous and parallel task execution across infrastructures without requiring separate configurations per environment. This interoperability is further empowered by an advanced cache management layer that tracks output states for each target. In case the outputs are up-to-date, based on checksums, timestamps, or other validation rules, then the corresponding task will be bypassed during subsequent runs. Such caching reduces redundant computations drastically and mitigates both time loss and resource waste.

Architecturally, WEXGrid follows explicit separation of concerns between definition, scheduling, execution control, and optimization phases. The definition includes explicit target declarations within Python scripts; scheduling encompasses the interpretation of dependency graphs together with resource metadata; execution control monitors runtime behavior while implementing error handling strategies; and optimization includes iterative adaptations based on performance metrics collected for varying workloads. For instance, performance testing cycles could identify bottlenecks due to I/O throughput rather than Central Processing Unit (CPU) constraints—a distinction commonly occurring when dealing with either large genomic data or numerical simulation results [1].

Dependency resolution takes the lead over any real computation. The system evaluates trees of dependencies recursively down to the identification of all possible parallelizable branches. By isolating independent subgraphs, WEXGrid allows for concurrency in processing without violating the original dataflow semantics. This approach directly addresses one of the common pain points in large-scale workflows where unintended sequential processing inflates completion times far beyond what hardware capabilities could otherwise allow. Scientific workflows frequently include tasks whose runtime may vary widely depending on the characteristics of the data, such as total record counts or file size distributions. Empirical studies have hinted at attributes such as total read magnitude as the top contributor to runtime prediction models, closely followed by other contextual factors like word length in specific classification scenarios.

Embedding such attribute awareness within WEXGrid’s scheduler enhances predictive load balancing: tasks likely to consume disproportionate time are strategically forwarded to nodes that have greater capacity or scheduled earlier if their outputs happen to hold particular significance for critical downstream processes. Introduction of adaptive scheduling heuristics remains a forward-looking aspect in the roadmap of WEXGrid. Drawing from parallels with algorithms used in regression-based performance modeling [1], the system could incorporate learned coefficients correlating task parameters with runtimes, gradually refining these predictions over successive executions. This adaptation would strengthen throughput stability across diverse conditions of the datasets without requiring manual oversight from system operators.

WEXGrid positions itself within the larger body of workflow research through selective comparison with frameworks such as Pegasus or Snakemake, systems recognized for their advanced features of provenance tracking and portability, but it primarily seeks to bridge high-functionality orchestration with codebase simplicity suitable for tight integration into bespoke scientific projects maintained by small teams. Through this balance, it opens pathways for computational scientists to deploy reproducible experiments at scale while preserving direct visibility into their underlying task logic, visibility often lost when using heavily abstracted graphical workflow builders.

By integrating architectural clarity with operational adaptability, WEXGrid directly addresses the contemporary challenges of efficiently running large-scale scientific workflows under varied infrastructural constraints [2]. By coupling Python-driven accessibility with thoughtful system-level optimizations targeted at mixed local-grid deployments, it sets a foundation upon which future enhancements, including provenance capture modules or self-tuning schedulers, can be developed without disruptive restructuring of existing capabilities.

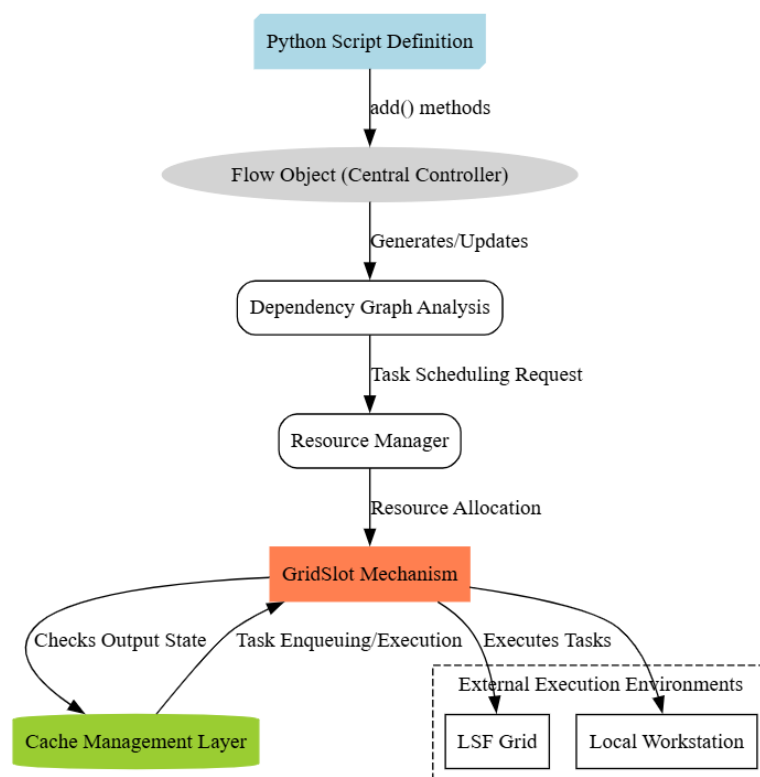


Figure 1.1: High-level architectural overview of the **WEXGrid** framework.

Chapter 2

Background and Motivation

2.1 Evolution of Workflow Management Systems

Scientific workflow systems have undergone extensive transformation over the past two decades, moving from bespoke, domain-specific implementations to general-purpose frameworks capable of orchestrating large-scale distributed computations with fine-grained control. Initially, workflows were often hard-coded as sequential scripts where procedural execution determined the outcome without explicit separation between data dependencies and computational steps.

Although functionally adequate for small datasets or straightforward analyses, these approaches offered minimal adaptability when faced with diverse resource environments or heterogeneous task requirements. Over time, formal abstractions such as directed task graphs emerged as a mechanism to represent computation in a modular and reusable form. In such descriptions, nodes model discrete tasks and edges express data or control dependencies, enabling schedulers to reason about execution order based purely on topology rather than manual sequencing. This abstraction proved indispensable for optimizing parallel execution strategies in both high-performance clusters and cloud environments. The transition toward graph-based workflow modeling allowed more sophisticated scheduling algorithms to be developed.

These algorithms assess graph structure to identify independent sub-trees that can execute in parallel, constrained only by resource availability rather than sequential logic [2]. Early implementations often used static scheduling plans derived from predefined heuristics. However, static allocation made it difficult to adjust for unpredictable runtime conditions such as variable I/O throughput or fluctuating job queue lengths, factors particularly relevant when distributing workloads across an LSF Grid where node states can vary dynamically. It became apparent that dynamic scheduling methods exploiting runtime monitoring could greatly improve throughput efficiency. In parallel with innovations in structural modeling and scheduling, provenance tracking emerged as another critical facility within workflow management evolution. Modern systems began capturing detailed execution metadata, spanning intermediate outputs, parameter values, and environmental configurations, to facilitate reproducibility and auditing of scientific results [1].

For disciplines dealing with biomedical research or sequence analysis pipelines, such provenance allows researchers to inspect alignment rates or verify output integrity without reconstructing entire computational contexts from scratch. While WEXGrid integrates caching for efficiency reasons, its design trajectory suggests potential alignment with domain-specific provenance practices for longer-term reproducibility goals. Another dimension shaping the evolution relates to standards for data interoperability.

Frameworks began incorporating annotation models capable of semantically describing experimental inputs and outputs with established ontologies like EDAM or Experimental Factor Ontology (EFO) [3]. Those annotation capabilities provide a means for workflows to produce not only FAIR-compliant data themselves but also to be discoverable, interpretable by machines, and reusable in diverse settings. Such paradigms have pushed developers to view workflows simultaneously as executable processes and as structured digital assets that themselves require stewardship akin to scientific datasets. The nature of community engagement has also evolved over time. Today, research-driven communities support the codevelopment of workflow technologies in collaborative standards efforts focused on FAIR principles adapted specifically for software [4]. In distinguishing between

research-and-development groups and end-user communities, the model encourages bidirectional communications: developers learn about operational bottlenecks experienced by practitioners, while users benefit from rapid integration of technical advances such as adaptive schedulers or enhanced caching logic into production environments.

Comparisons between the contemporary frameworks discussed below illustrate these differences shaped by such evolutionary pressures. Systems like Pegasus emphasize portability and provenance-capture mechanisms; Snakemake demonstrates flexibility combined with syntax provided by an embedded language; Nextflow puts forward container orchestration across diverse cloud infrastructures. WEXGrid situates itself within that ecosystem by pursuing high-functionality orchestration fused with Python-level accessibility—a bridge between technical sophistication and usability that earlier generations often treated as mutually exclusive design objectives.

Resource management has run parallel with the advancement of scheduling sophistication. The earliest server-based databases employed for workflow backends emerged in part from limitations of embedded systems like SQLite when serving concurrent connections under load. PostgreSQL brought better scaling and richer datatype support applicable for complex dependency tracking at scale; however, constraints in some deployment contexts still drive selection toward lightweight alternatives or hybrid architectures that connect persistent stores selectively in their aim to minimize overhead without compromising concurrency needs. Advances in simulation-based scheduler benchmarking have further influenced design choices.

Comparative studies explore information modes, from exact representations of runtime tasks to user-specified estimations, in their impact on completion times for distributed orchestration scenarios. Results have indicated that exact modes yield, in general, optimal throughput but at the cost of additional monitoring; mean-mode approximations offer a compromise where predictability is less important than sustaining load balance under fluctuating queue conditions. From this perspective, WEXGrid’s adoption of Python-defined target constructs linked by dependency graphs is consistent with the broader historical movement toward modularity combined with dynamic adaptability. Yet, it adds

its own refinement through integrated cache-awareness that minimizes re-execution when outputs remain valid and strategic allocation mechanisms responsive to both local machine states and LSF Grid slot availabilities. The evolution traced above underlines why such hybrid considerations are now integral: the prevailing diversity of resource environments demands systems capable of balancing architectural clarity against pragmatic runtime decision-making.

One observes in trends within workflow management’s progression a cumulative layering of capability: structural abstraction via task graphs [2], adaptive scheduling logic responsive to real-time conditions [1], semantic annotation fostering machine-readability [3], provenance tracking aiding reproducibility, and evolving community collaboration models placing equal weight on developer insight with user experience [4]. WEXGrid inherits lessons from each stage while specifically tailoring these for use by computational scientists working within the dual local-grid deployment contexts described earlier in Section 1.

2.2 Challenges in Local and Distributed Workflow Execution

Execution of workflows in both local and distributed environments inherently introduces constraints that WEXGrid must address through deliberate architectural and algorithmic choices. Conducting computational tasks within a single node can appear to be less complex, but lack of sophisticated allocation logic often results in performance degradation when workloads become heterogeneous. On local systems, contention over CPU cores and memory resources can arise if multiple targets are resolved concurrently without awareness of their individual consumption patterns. The ability of WEXGrid to define explicit resource requirements for each target via an `R()` object mitigates this risk. Memory and core specifications act as guardrails to prevent over-subscription, yet these safeguards alone do not remove latency caused by imbalanced task ordering—a situation where the scheduler

queues lighter jobs first while heavier, downstream-critical tasks wait unnecessarily.

As execution shifts toward an LSF Grid context, a different set of challenges arises. Network latency between distributed nodes sparks discrepancies in job start times, which can be compounded by variable I/O throughput [5]. Tasks dependent on large files stored remotely can face disproportionate delays even when compute resources appear idle [6]. Dynamic allocation strategies allow WEXGrid to adaptively match task profiles with available slots that satisfy both computational and data proximity considerations [7]. Through engagement of a resource pool maintained via the Grid Resource Manager and MasterManager components, the system ensures load distribution in balance across heterogeneous nodes, alleviating choke-points created by uneven file access speeds [8].

A recurring problem in mixed local-grid operation is the resolution of dependencies across infrastructure boundaries. If components of a workflow execute locally, while other stages run on remote grid hosts, then synchronization of intermediate outputs becomes sensitive not only to direct network speed but also to file-system compatibility. WEXGrid reduces cross-platform inconsistencies by representing inputs and outputs explicitly through File and Var objects, linked within its dependency graph. This representation enables integrity checks before initiation of dependent tasks, avoiding the scenario in which a job consumes incomplete or corrupted data generated upstream. Although this design addresses reliability concerns at a functional level, it can introduce minor delays due to metadata extraction overheads for very large file sets.

Another point of contention arises from predicting execution times. Even with well-defined dependencies and resource constraints, runtime variability is common, especially when datasets exhibit unpredictable attributes such as uneven distribution of record sizes or fluctuating numeric ranges within parameters handled by regression-type models [1]. In distributed contexts, inaccurate predictions can cause underutilization if tasks are allocated conservatively or lead to excessive queue waiting when multiple long-duration jobs converge on limited slots. Incorporating adaptive heuristics into WEXGrid's scheduler helps refine such estimations over repeated runs by correlating observed durations with specific data characteristics collected during earlier executions.

Data provenance requirements add another layer of complexity in distributed execution scenarios. While core provenance capture is conveniently integrated in localized runs, where all metadata remains within a single context, the grid setup disperses this information across nodes unless actively centralized. Automated annotation models tied into WEXGrid could ensure consistent provenance collection regardless of execution site [3], aligning with FAIR workflow principles known to enhance reproducibility [4]. This practical need gains urgency when considering that missing metadata for even one stage can compromise downstream verification for an entire pipeline.

Fault handling diverges substantially between local and distributed modes as well. Locally, failures are typically isolated, allowing recovery through re-execution at minimal coordination cost. In contrast, grid-based failures may cascade across dependent jobs under peak load conditions if initial error signals do not propagate fast enough through schedulers managing thousands of concurrent tasks. WEXGrid addresses this risk via layered monitoring that flags resource failures before they escalate; however, this requires efficient communication back to the central control logic without introducing overwhelming monitoring traffic onto the grid network.

Even with sophisticated mechanisms for resource pooling and load balancing, operational unpredictability remains inherent to diverse infrastructures supporting varied scientific domains. An HPC node optimized for memory-heavy simulations may perform poorly on I/O-bound tasks involving smaller data fragments but high frequency reads and writes. Dynamic workflow task assignment to suitable hardware based on run-time profiling can mitigate these mismatches, but this relies heavily on responsive performance feedback loops inside the scheduler.

Community-driven metadata standards further complicate this because any adoption of schema.org or RO-Crate-based formats in a multi-environment setting requires consistency in interpretation between the locally managed workflows and those deployed upstream on shared platforms like WorkflowHub.eu [4]. Without consistency in metadata serialization pipelines across sites, interoperability breaks down quickly when composite workflows spanning organizational boundaries are exchanged.

The challenge matrix for WEXGrid thus involves interdependent factors, including fine-grained resource specification versus broad-scale load balancing, accurate runtime prediction grounded in dataset attributes versus naturally unpredictable variations inherent in scientific computation, unified provenance capture strategies across diverse infrastructures, timely fault detection mechanisms suitable for tens or hundreds of concurrent targets, ensuring portability not just at code-level syntax but also at semantic levels defined by emerging FAIR workflow conventions [4], and avoiding storage bottlenecks where workflows oscillate between local disk caching and distributed file systems easily congested during peak operations.

This produces an environment where architectural cleanliness needs to coexist with operational pragmatism, a balance whose difficulty increases with accommodating both solitary workstation executions and massive grid deployments, as referred to earlier in Section 2.1. WEXGrid creates an adaptable interface between such contexts without abandoning efficiency goals or reproducibility standards central to scientific workflow practice.

2.3 Motivation for Developing WEXGrid

The decision to develop WEXGrid arises directly from shortcomings in the existing workflow management solutions as computational scientists work across both localized computing environments and distributed LSF Grid infrastructures. As experience with more traditional tools such as make has demonstrated, their static nature inhibits adaptation in dynamic contexts where task execution times, resource availability, and data access characteristics may change unpredictably between runs. These tools also often use very simplified dependency models and lack expressive mechanisms for detailing the complex interplay between inputs, outputs, and computation-specific execution constraints, key features of scientific pipelines characterized by iterative experimentation and heterogeneous workloads. From an architectural perspective, WEXGrid is designed to bridge usability with operational sophistication for high-complexity workflows.

Compact Python syntax provides an easy way to define workflow targets while enabling fine control over resource specifications at the task level. This makes it practical to declare CPU core allocations, memory constraints, or affinity for specialized hardware without requiring developers to adopt verbose configuration schemes commonly used in XML- or DSL-based systems. Internally, maintaining tasks as nodes in a dependency graph follows the structural models seen in several advanced workflow engines [9] but is explicitly optimized here for integration with both local dispatch and remote grid schedulers. One important motivator for this latter hybrid compatibility is to reduce friction when migrating computational stages between environments—a situation that arises frequently when small-scale prototyping needs to be expanded later into large batch runs on a grid. Although several modern frameworks support distributed execution, they often embed assumptions about the underlying infrastructure that are coupled into HPC-centric user interfaces and batch scheduling protocols such as Message Passing Interface (MPI) bindings or tightly coupled job arrays [4]. In contrast, WEXGrid addresses orchestration by abstracting from the physical layer while still providing hooks for environment-specific optimizations. Its design emphasizes runtime monitoring to inform dynamic allocation decisions: rather than binding tasks to given resources in a static fashion ahead of execution, the system continuously reassesses slot availability in concert with dependency readiness throughout the operational cycle. This paradigm not only minimizes idle time but also creates avenues for employing adaptive heuristics, important given unpredictable variables such as network contention or variable I/O throughput during the transfer of large datasets [10].

Caching is another major driver behind WEXGrid’s design. Needless recomputation remains a common inefficiency within scientific workflows; re-executing stable outputs wastes both time and physical resources. By incorporating a validation layer based on checksum verification, timestamp comparison, or other integrity check before each task invocation, WEXGrid selectively bypasses re-execution if its outputs are determined to be up-to-date. This mechanism falls in line well with situations where upstream processes yield deterministic results across iterations, such as simulation stages driven by

fixed parameter sets or preprocessing phases applied to immutable datasets. Integrating cache awareness within the core scheduler logic enhances throughput without introducing protocols for separate cache maintenance, yet another point at which many established tools tend to falter due to reliance on external scripts or manual oversight [11].

There is also an intentional alignment between WEXGrid’s architecture and emerging best practices from workflow research communities, which call for interoperability and FAIR principles in software systems [12]. Although automated semantic annotation of inputs and outputs remains part of its roadmap rather than current implementation, structuring task metadata explicitly within Python object hierarchies lays the foundation for later integration of ontology-based descriptions analogous to EDAM or those that might be derived from P-Plan sequence patterns [3]. Such capabilities would enable workflows to operate not just as procedural instructions but as discoverable digital assets transferable across institutional boundaries without semantic loss—a core priority in collaborative science ecosystems that involve multiple laboratories operating partially overlapping pipelines. Another motivating factor arises from recognizing that static scheduling algorithms, even when informed by sophisticated heuristics like b-level or t-level sorting, cannot fully optimize throughput within variable-load environments. The inability to precisely determine earliest start times under conditions of unpredictable resource contention limits real-world applicability of such methods without additional adaptive overlays. WEXGrid has plans to incorporate feedback-driven allocation informed by empirical runtime profiles collected over repeated executions. This becomes possible by correlating historical task durations with workload attributes, such as file sizes and record counts, to refine scheduling predictions and redistribute heavier tasks more strategically over the available nodes.

This self-improving aspect addresses deficiencies seen in conventional schedulers that treat each run as an isolated event devoid of accumulated performance intelligence. The actual development of WEXGrid on top of Python avoids the constraints imposed by more traditional High-Performance Computing (HPC) languages, such as C/C++ combined with MPI/Open Multi-Processing (OpenMP) patterns that dominate the legacy HPC

codebases but often alienate researchers who lack systems-level programming experience. Where other frameworks either sacrifice accessibility for raw performance or vice versa, WEXGrid’s guiding principle is that scientists should directly engage in workflow logic without deep immersion in low-level orchestration layers unless desired. This supports maintainability: pipelines can be iteratively extended through version-controlled Python scripts rather than domain-specific markup susceptible to syntactic rigidity and reduced testability within standard software development cycles. Finally, there is recognition that funding landscapes and innovation incentives sometimes deprioritize maintenance-oriented advancements like standardization or long-term interoperability improvements.

WEXGrid is developed to represent an alternative trajectory by embedding maintainability into its initial design choices rather than treating it as an afterthought, requiring retrofitting once scaling issues are encountered. Emphasis on modularity secures the position that features, including provenance capture modules or enhanced annotation support, can be incrementally added without disruptive rewrites—a quality deemed essential for ensuring sustained adoption beyond initial deployment phases among small research teams managing diverse computational projects. Motivation, therefore, converges across several axes: dissatisfaction with inflexible legacy tooling; needs for cross-environment adaptability connecting local prototypes with distributed-scale execution; desire for integrated caching mechanisms that minimize redundant work; potential alignment with semantic interoperability standards under FAIR principles; ambition to replace static heuristics with adaptive data-informed scheduling; prioritization of accessible Python-level programmability without sacrificing orchestration sophistication; and proactive incorporation of features conducive to sustainable long-term use under evolving collaborative science models. Each of these considerations directly shapes not only why WEXGrid exists but how its architectural commitments are articulated relative to constraints and opportunities described earlier in Section 2.2.

Chapter 3

Assumptions and Requirements

3.1 Functional Requirements

The functional requirements for WEXGrid stem from the need to operationalize its architectural commitments, as described above, into concrete, testable capabilities that allow it to meet its intended purpose in both local and LSF Grid execution contexts. At their most basic level, these requirements translate abstract design intentions of target decomposition, dependency resolution, dynamic scheduling, caching, and interoperability into explicit behaviors that the system must reliably produce across real execution scenarios.

One foundational requirement has to do with target definition and management. WEXGrid should allow each target to be programmatically declared using Python constructs, which, in turn, capture precise specifications for their inputs, outputs, execution actions, and optional metadata such as labels or explicit resource constraints. To support a wide array of scientific applications, these actions include the invocation of native Python functions, shell commands with environmental customization, and external executables with fully configurable arguments and working directories. Uniformity needs to be maintained in how these are internally represented so that downstream scheduling components can uniformly operate on an abstract interface, without regard for whether the wrapped operation is CPU-bound numerical processing or an I/O-heavy data transformation [13].

Of equal centrality is the requirement for sound dependency management. Defining tasks in the absence of correctly specified interdependencies risks breaking dataflow semantics and rendering invalid results. Thus, WEXGrid needs to immediately integrate each newly defined target into an internal dependency graph upon its registration with the controlling Flow object.

An execution plan is to be based on this graph: edges represent control or data dependencies among nodes, targets, while acyclic structure ensures the impossibility of circular waits deadlocking a workflow. Before any execution, the system's dependency resolver needs to recursively traverse this graph in search of independent subgraphs that may be executed in parallel without data corruption. Another critical functional axis concerns scheduling logic. WEXGrid shall conduct dynamic evaluation of task readiness concomitantly with available resources across the current execution environment. For local runs, this involves respecting resource hints, such as CPU cores and memory limits, attached to individual targets so as not to overcommit machine capacity. On LSF Grids, this means mapping these same requirements onto distributed resource slots managed through the Grid Resource Manager and MasterManager layers, with further consideration of grid-specific constraints like queue priority levels and node capacity heterogeneity. Since resource availability will vary overtime due to unrelated jobs executing on shared infrastructure, WEXGrid's scheduler needs to dynamically reassess slot allocation rather than rely on a static precomputed plan [10].

Caching brings into play specific operational demands of its own within this architecture. Prior to any given target, WEXGrid will have to check whether valid outputs are already present by verifying timestamps, comparing checksums, or applying other user-defined validity criteria. If outputs are determined to be up-to-date under such rules, the corresponding task is skipped entirely, thereby avoiding wasteful recomputation. This means keeping an efficient cache index synchronized with actual filesystem state so that reliance on stale metadata does not lead to false positives or negatives when cache lookups occur. In distributed scenarios where targets operate over several nodes

with distinct filesystems, additional functionality becomes necessary for reliable data access. The system must have mechanisms to check the presence and integrity of input files in the execution node before launching dependent tasks. For cases where needed files reside remotely, it should orchestrate transfer operations in a fashion that is at least transparent and integrates transfers into dependency calculus so that compute stages do not start out of turn. This intermixes provenance-oriented metadata tracking with low-level operational safeguards.

For predictability amidst mixed workloads composed of both long and short-duration tasks, WEXGrid has to have the option to incorporate historical runtime observations into scheduling decisions. In correlating prior execution times to measurable input properties, the scheduler can advance heavy-impact tasks sooner when they form prerequisites for many downstream steps. This heuristic capability should be adaptable over repeated executions: performance models are updated iteratively as empirical data accumulates [14]. Provenance capture also features within these functional imperatives, not yet as a fully automatic mechanism, but factored into the core data representations in such a way that such capabilities can be layered in without redesign. Targets must permit annotation fields describing their semantic roles within computational experiments; inputs and outputs should be representable in ontology-compatible formats if configured accordingly. Even without immediate integration of linked data publishing pipelines, adherence to clear metadata structures supports reproducibility-enhancing features down the line. Fault handling demands another set of concrete functions: detecting failed tasks quickly; marking all dependent targets as blocked until recovery actions succeed; logging error details comprehensively enough for later diagnosis; and retrying jobs when configured thresholds permit.

Such mechanisms have to work uniformly across local threads and remote grid jobs despite differences in how failures manifest (local exceptions versus scheduler error codes). In grid mode particularly, there is a requirement for timely communication of fault states back to centralized control logic without generating excessive network load through over-frequent polling. Because workflows may operate in collaborative research ecosystems

where portability and reuse matter, WEXGrid must additionally support exporting workflow descriptions, including dependency graphs, resource annotations, and parameter sets, in neutral serializations such as JavaScript Object Notation (JSON) or YAML Ain't Markup Language (YAML) interpretable by external tooling. Without this capacity for portable configuration interchange, integration into larger federated platforms would remain constrained. In all these respects, target declaration flexibility; dependency graph construction; adaptive scheduling across heterogeneous environments; integrated caching; cross-node data validation; predictive task ordering informed by observed durations; extensible provenance support; resilient fault handling; exportable workflow definitions, the functional requirements collectively articulate what WEXGrid has to achieve if it is to substantiate its proposed role bridging accessibility and high-performance orchestration capacity demonstrated in Section 2.3.

Together they ensure that each architectural choice previously outlined translates directly into operational competencies verifiable under real-world workload conditions spanning both solitary workstation use cases and expansive distributed computations across LSF-managed clusters.

3.2 Non-Functional Requirements

In addition to the explicit functional capabilities defined in Section 3.1, WEXGrid must meet a host of non-functional requirements which specify the quality, reliability, and long-term maintainability of the system. Whereas functional elements describe what the system does, these non-functional specifications stipulate how well it carries out those functions within a wide range of operational and contextual circumstances. They're integral to ensuring the system will work not only as expected but also be performant, reliable, and sustainable for real scientific workflows, both locally and across LSF Grids.

One primary consideration pertains to performance efficiency. In this regard, efficiency is multi-dimensional: low-latency scheduling decisions, high throughput in task execution, minimal idle resource windows, and avoidance of redundant computations. A scheduler's

decision-making loop must incur negligible overhead relative to task runtimes so that the orchestration logic is not itself a bottleneck [15]. The cache validation mechanisms, while rigorous enough to ensure correctness, must be optimized for both scenarios involving thousands of small targets and fewer large tasks. This is a delicate balance; too aggressive, and one risks saturating I/O resources during metadata checks, while too lax, one risks performing superfluous work. This trade-off is iteratively refined, driven by empirical performance profiling to ensure it stays favorable as workload composition shifts due to changing dataset characteristics or evolving computational recipes [16].

Scalability forms another critical non-functional dimension. For any additional available compute resources, whether through the addition of nodes on an LSF Grid or increasing core counts locally, WEXGrid must exhibit near-linear scaling in workload completion time. Similarly, internal data structures for dependency graphs and resource tracking need to support this scaling without degradation in access speeds or excessive memory consumption [17]. Communication protocols between distributed components likewise must avoid contention points that would hamper scalability benefits; centralized coordination should be minimized for those operations that can be handled via decentralized status updates or local decision-making by grid-executed agents.

Of equal centrality are the requirements for reliability and fault tolerance. Failures, whether caused by hardware issues, transient network outages, or software errors, cannot propagate cascading disruption beyond those tasks directly affected. This implies robust detection and recovery pathways that can quickly isolate failures while keeping the workflow state for unaffected portions consistent [18]. Effective error handling mechanisms rely on proper logging, with clear timestamps, resource identifiers, and data references for post-execution diagnostics. This logging should not introduce prohibitive costs at runtime; it should perform asynchronously whenever possible so that computation-critical threads are not blocked by I/O-bound logging routines.

Interoperability is an important feature due to WEXGrid's intended role in heterogeneous research environments. Workflow descriptions should remain portable across installations featuring different local configurations while maintaining semantic integrity

regarding dependencies and resource annotations. This implies the use of standardized serialization formats, like JSON/YAML, with the optional inclusion of ontology-driven descriptors where FAIR-compliant integration is desired [15]. The smooth interfacing with upstream or downstream systems, such as institutional job submission portals or shared provenance repositories, uses stable API definitions insulated from frequent breaking changes unless justified by substantial architectural gain.

Usability is another aspect that underpins long-term adoption. While these systems invariably present internal complexity, the user-facing interface must minimize entry barriers for scientists without deep backgrounds in software engineering. Python-based target declarations already contribute to this goal by aligning with languages familiar to many researchers; documentation quality, consistency in CLI/GUI feedback messages, error explanations with actionable suggestions, and example-rich onboarding material remain key non-functional deliverables supporting actual user productivity. Maintainability has close alignment with usability but targets developers extending or otherwise modifying the system rather than end-users configuring workflows.

WEXGrid’s internal architecture should preserve modularity so that enhancements can be integrated without widespread refactoring [16]. Each subsystem, including dependency analysis, cache management, and scheduling core, shall be clearly interfaced, allowing independent testing and replacement if more efficient algorithms become available later. Coverage via automated test suites at both unit-level behavior and cross-component integration scenarios mitigates regression risks during upgrades. Security requirements cannot be ignored, considering that scientific workflows often handle sensitive data which is subject to ethical or legal constraints. Execution contexts need explicit sandboxing capabilities to prevent accidental leakage of restricted data between unrelated workflow runs on the same infrastructure node. Authentication mechanisms for accessing the intermediate result storage become important, together with access control lists to prevent unauthorized modification of the files during execution, on shared grids where many users submit their jobs concurrently [19].

Adaptability rounds out these requirements since infrastructure conditions change

over time. Queuing policy changes occur for HPC clusters; local machines get hardware upgrades; network topology between grid nodes may change after reconfiguration events. Again, WEXGrid’s dependency resolution and scheduling strategies must cope gracefully under such variation without wholesale redefinition of existing pipelines [15].

One promising approach towards this adaptability involves incorporating runtime learning mechanisms into allocation heuristics: past execution patterns feed into future scheduling predictions without the need for manual retuning by operators [20]. Efficiency in resource utilization extends beyond raw performance metrics into environmental cost-awareness-one increasingly important consideration in large-scale computing projects, where energy budgets matter as much as wall-clock time savings. WEXGrid needs to offer optional modes aimed at reducing energy footprint-for example, consolidating workloads onto fewer active nodes during off-peak hours or throttling concurrency when predicted task durations show minimal benefit from maximal expansion at high power draw [21]. Finally, community sustainability is a factor determining whether WEXGrid remains relevant beyond its initial deployment window. Openness of governance models-including transparent roadmaps informed by user feedback loops-and structured dissemination channels for updates enhance trust among scientific collaborators who might otherwise fear investing effort into adopting a system vulnerable to sudden abandonment.

Consistent communication about known limitations, along with recommended workarounds, signals maturity in project stewardship that complements technical robustness. Taken together, spanning efficiency under varying loads [18], scalable design exploiting additional resources effectively [17], resilience against failures paired with diagnostic clarity [20], semantic interoperability under FAIR-oriented standards [15], researcher-focused usability balanced with developer-level maintainability [16], data security safeguards appropriate for sensitive workloads, adaptability amid infrastructural evolution, energy-conscious operation strategies [21], and sustainable community practices, these non-functional requirements articulate the quality benchmarks against which WEXGrid’s success will be ultimately judged over repeated cycles of deployment in both localized research labs and distributed grid environments operating at scale.

3.3 Operational Environment Assumptions

In designing WEXGrid, certain assumptions about the operational environments in which it will execute are necessary to guide both architecture and performance optimization strategies. These assumptions influence scheduling policies, caching effectiveness, resource allocation heuristics, fault management procedures, and interoperability measures. They operate as foundational constraints within which the defined functional and non-functional requirements can be validated under controlled yet realistic conditions described previously in Section 3.1.

A primary assumption concerns the dual deployment mode across local workstation environments and LSF Grid infrastructures. Both contexts possess distinct system characteristics, hardware configuration, network topology, storage architecture, that are expected to remain stable over the course of typical workflow executions. In a local environment, WEXGrid presumes that hardware specifications such as core count, memory capacity, and available disk throughput are known beforehand and can be queried programmatically during scheduler initialization. This information is critical for accurate resource matching when instantiating $R()$ objects and avoiding oversubscription that would degrade performance. Given these static baseline attributes, fluctuations during execution (e.g., competing processes launched by other users on the same workstation) are assumed to be minimal or sporadic rather than systemic.

Within LSF Grid environments, more complex state variability must be considered. Here, WEXGrid assumes that fundamental queuing policies, including slot allocation rules, node capacity configurations, priority hierarchies, are transparent to user-level scheduling agents so they can integrate this data into dynamic allocation logic. However, it also presumes that node availability metrics fed into the scheduler at runtime are accurate to a granularity sufficient for responsive task placement without over-reliance on repeated status polling across the grid network. This assumption mitigates the risk of excessive control traffic slowing down orchestration under load while still enabling adaptive rescheduling if initial allocations become invalid due to contention or job failures

mid-execution.

Data storage accessibility represents another operational assumption layer. For local execution modes, it is expected that all required input and output files reside on a filesystem directly mounted to the machine running WEXGrid tasks. Latency from file operations is assumed consistent enough that caching checks, such as timestamp verification or checksum evaluation, incur predictable costs well below average task execution time. On LSF Grids, distributed filesystem access such as via NFS or parallel file systems like Lustre is anticipated to exhibit variable latency depending on concurrent usage patterns by other jobs; nevertheless, baseline read/write speeds must remain within thresholds enabling WEXGrid’s dependency resolver to complete integrity checks without repeatedly stalling downstream tasks. Where remote staging of files is required, infrastructure is assumed to permit pre-transfer orchestration embedded in dependency graphs so computation commences only after confirmed data arrival at target nodes.

An important performance-related assumption involves predictability of task durations with respect to measurable data attributes. As empirical studies on classifier workflows indicate total read size or record length can heavily influence compute times, WEXGrid assumes correlation patterns between input properties and runtime behavior remain relatively stable over successive workflow runs for similar datasets [22]. This stability enables its planned adaptive scheduling heuristics to use historical profiling data effectively without continual recalibration from scratch for each iteration. Deviations outside of learned ranges are expected but regarded as exceptions rather than norm; hence they trigger fallback scheduling approaches instead of invalidating predictive models entirely.

Fault behavior assumptions define how error handling subsystems operate in different environments. Locally, hardware faults impacting workflows are presumed infrequent compared to software-level exceptions such as missing input files or failed external executable calls. Within LSF Grids however, transient network disruptions or node unavailability events are treated as an inherent occurrence rate high enough to justify constant lightweight monitoring by the resource manager module. It is assumed that grid schedulers propagate failure notifications quickly enough for WEXGrid’s central controller

to halt dependent tasks without risking cascade effects across unrelated portions of the workflow graph. This synchronous awareness between grid infrastructure and WEXGrid’s execution logic forms part of operational viability expectations for distributed resilience mechanisms.

Interoperability assumptions are likewise important given FAIR-compliant principles influencing system roadmap features. Environments where WEXGrid runs are presumed capable of serializing both workflow definitions and metadata annotations into portable formats such as JSON or YAML without incompatible character encoding issues or unsupported ontology bindings at runtime. For cases involving cross-institutional workflow exchange via federated repositories or aggregate research platforms, common in collaborative science, it is assumed these external systems accept dependency graph structures alongside semantic data mappings in ways consistent with WEXGrid’s internal model without requiring direct structural transformation layers mid-transfer.

Security-related assumptions warrant attention when multi-user grids host workflows involving sensitive datasets. Nodes executing WEXGrid tasks within shared environments are presumed configured with sandboxing mechanisms restricting task-level filesystem access strictly to allocated paths; thus provenance data collected during runs aligns entirely with accessible artifacts owned by the initiating user account. Authentication layers controlling access to intermediate results storage are expected operational before workflows begin so that no inadvertent cross-user data exposure occurs, facilitating safe reuse scenarios particularly when distributing cached results between sequential runs originating from common pipelines but executed by different researchers under shared project umbrellas [19].

Scalability expectations shape how dependency graph construction performs under expanded workloads. Operational environments are assumed capable of sustaining memory footprint demands imposed by large target counts through adequate RAM provisioning; for LSF Grids this extends to managing wide graphs where hundreds of independent leaves may execute concurrently across nodes without exceeding scheduler coordination

capacity limits. Local machines used for scaled-down testing phases should simulate realistic concurrency levels while preserving event logging fidelity so profiling outcomes inform production grid deployments accurately despite smaller run scales in preliminary sessions.

Energy utilization considerations coincide with operational settings supporting power-aware execution policies during off-peak cycles. While not mandatory for correctness or functional validity, it is assumed HPC clusters hosting WEXGrid jobs permit engagement with environmental efficiency frameworks when configured accordingly, whether through consolidating active workloads onto fewer nodes temporarily or throttling CPU utilization rates using standard cluster APIs during extended low-demand periods without severe impact on accuracy-sensitive computational stages [21].

These combined assumptions form an interlocking framework underpinning anticipated behavior under both single-node and distributed multi-node orchestration setups: stable enough baseline hardware profiles locally; transparent yet accurate resource availability metrics within grids; predictable latency bounds for file operations; repeatable correlations between dataset attributes and task durations; timely propagation of fault signals; maintained compatibility in metadata serialization standards; enforced sandboxing for security; adequate scalability support by infrastructure; optional energy policy integration mechanisms, all ensuring the architectural commitments earlier specified translate effectively into sustained operational competence under diverse scientific workloads.

Chapter 4

Contextualization and Related Work

4.1 Scientific Workflow Systems Landscape

The contemporary landscape of scientific workflow systems is characterized by increasing structural sophistication and adaptability, informed by years of iterative development and empirical evaluation in diverse computational domains [23], [24]. Within this spectrum, WEXGrid positions itself as an accessible yet high-performance orchestration engine that draws upon established patterns in workflow modeling while addressing gaps left by both legacy and current-generation frameworks. The diversity of available systems, ranging from DAG-based orchestration engines like Pegasus to language-embedded managers such as Snakemake or Nextflow, illustrates a broad consensus on the benefits of graph-oriented computational abstractions, but also reveals differing philosophies in balancing portability, user accessibility, and runtime optimization. Many workflow frameworks share the underlying goal of decomposing complex pipelines into discrete tasks with explicit dependency relations. Traditional implementations hard-coded these sequences as static scripts, which limited their adaptability across heterogeneous environments. Graph-oriented architectures solved this by enabling schedulers to infer execution order from topology rather than procedural code flow, facilitating parallel execution where resources allow [23]. Yet this approach often falters when runtime conditions diverge significantly from initial heuristic

predictions, especially in distributed grid contexts with variable queue states or unpredictable I/O overheads [25]. WEXGrid’s integration of dynamic dependency resolution and continuous slot reassessment directly responds to such shortcomings, adapting execution mid-cycle without manual intervention.

Comparative analysis with existing frameworks shows marked differences in caching strategies. Systems like Pegasus employ durable provenance stores tightly coupled to execution engines [23], while Snakemake favors a lightweight file-timestamp model for determining task reusability [26]. WEXGrid merges elements from both camps: it maintains explicit output validity checks through timestamps or checksums for direct efficiency gains, but does so within a broader scheduling logic that considers resource availability and dependency impact before skipping computations. This coupling prevents scenarios where valid outputs exist but are nonetheless recomputed due to scheduler misalignment, a flaw observable in tools that isolate caching from orchestration logic.

Notably, the FAIR computational workflows initiative has shifted focus toward ensuring workflows themselves meet findability, accessibility, interoperability, and reusability criteria [27]. Applying FAIR principles to workflow definitions demands structural metadata capturing not only what tasks do but under what environmental constraints they were executed, along with semantic annotation aligning inputs and outputs to shared ontologies such as EDAM. WEXGrid accommodates these requirements by embedding metadata fields in target definitions at the Python level, creating an architectural scaffold upon which automated ontology mapping could be integrated later. While other systems may already incorporate ontology-driven descriptions into their core data models, WEXGrid emphasizes flexibility, allowing such features to be adopted incrementally without disrupting baseline usability for teams more focused on operational throughput than semantic interoperability during early deployment phases.

Distributed execution presents another axis along which workflow systems differentiate themselves. Nextflow emphasizes container-based isolation across cloud infrastructures [28]; Pegasus integrates deeply with grid middleware [23]; Snakemake offers hybrid

local-distributed operation via DRMAA-compliant job submission layers [26]. WEX-Grid’s model aligns most closely with hybrid operation but insists on uniform abstraction across environments: targets retain identical declaration semantics regardless of whether they execute locally or in LSF Grid mode, reducing friction when transitioning pipelines between prototyping and production-scale runs. This unification also simplifies fault handling since error detection and recovery logic remains consistent regardless of execution site, a contrast with some systems requiring separate error management strategies per environment [25].

The research community has long recognized that dynamic scheduling driven by real-time monitoring outperforms static allocation under volatile resource conditions [29], [30]. Yet implementing this adaptivity involves trade-offs between monitoring granularity and orchestration overhead. In the broader landscape, few systems successfully minimize monitoring costs while retaining responsiveness at scale. WEXGrid addresses this with a tiered strategy: lightweight status tracking suffices for non-critical dependencies while deeper profiling applies only when historical data suggests risk of bottlenecks or heavy downstream impact from delayed execution. Such selective adaptivity stems from recognizing that high-frequency polling can saturate communication channels on large grids, a condition noted as problematic in earlier studies [31].

There is also divergence among existing platforms regarding provenance capture scope. Airflow is focused on operational logging, while Pegasus yields fine-grained provenance that is suitable for formal publication [23]. Snakemake supports lightweight reproducibility by embedding configurations alongside the output artifacts [26]. WEXGrid uses structured internal representations internally and is ready for richer provenance modules without imposing them during initial deployment. This represents a philosophy oriented toward projects in incremental adoption of reproducibility features based on research funding cycles or institutional mandates [24]. It reflects an understanding that some workflows operate on sensitive data sets, where provenance publication involves selective disclosure; hence, modular support of provenance enables compliance without sacrificing analytic transparency internally.

Regarding scalability, heavily centralized schedulers can potentially face coordination bottlenecks while managing the execution of thousands of concurrent tasks across distributed nodes [32]. At the opposite extreme, decentralized agent-based models, exploited by some of the HPC-focused managers, promise better scaling at the potential cost of inconsistent state awareness among agents. Along this spectrum, WEXGrid relies on a centralized control process complemented with environment-aware slot managers capable of localized decision-making within the boundaries of assigned resources. Limiting the control traffic to needed events, allocation changes, faults, it aims for an equilibrium between global coordination accuracy and grid-scale scalability-essential factors observed in empirical benchmarking efforts described elsewhere [33]. Integration into collaborative ecosystems depends further on interoperability capacities to match evolving community standards. The workflows community actively discusses schema alignment of workflow exchange formats compatible with platforms like WorkflowHub.eu—a space where FAIR paradigms directly intersect with portability goals[27]. Here, the choice of WEXGrid to use serialization-ready internal representations serves two purposes: ease of export/import via neutral formats like JSON/YAML and readiness for ontological augmentation in case broader community standards converge on common descriptive vocabularies later adopted en masse by cross-institutional repositories.

Mapping WEXGrid’s place in the scientific workflow systems landscape reveals it leveraging mature concepts tested in frameworks like Pegasus or Snakemake, graph-based modeling, hybrid execution modes, while injecting nuanced advancements around adaptive scheduling linked explicitly to resource-state dynamics; integrated cache-awareness embedded into orchestration rather than detached layers; metadata structures conducive to FAIR-principled evolution without premature rigidity; uniform cross-environment semantics reducing migration friction; fault tolerance streamlined through shared logic paths; scalability supported via bounded decentralization within centralized oversight regimes; and broad adherence to interoperability practices priming adoption success across collaborative infrastructures spanning local laboratories to global grid networks. This combination not only marks its performance characteristics relative to its peers but also shows

an awareness of operational realities underlined in Section 3.3, with technical relevance under current and foreseen future trends shaping scientific computing practice worldwide.

4.2 Scheduling Strategies in Distributed Environments

Scheduling strategies in distributed environments are critically intertwined with WEXGrid's operational design because the system must navigate the interplay between dependency-driven execution order, heterogeneous node capacity, and dynamically changing availability [25]. In such settings, static task ordering, while easier to implement, too often fails due to unpredictable conditions like fluctuating queue wait times, network contention, or contention for shared file system bandwidth. WEXGrid consequently implements a scheduling model that capitalizes on dynamic reassessment of both resource states and dependency readiness throughout an execution cycle.

At its core, the scheduler maintains a representation of the target graph generated by dependency resolution logic. For each target, metadata encapsulates resource requirements (cores, memory footprint), expected data locality constraints, and cache status derived from output validity checks. This comprehensive representation enables a multi-criteria decision process wherein the scheduler evaluates whether a task can be dispatched to a given LSF slot without degrading overall workflow throughput. Instead of proceeding in a strict topological order, typical of many DAG executors, WEXGrid prioritizes critical-path tasks whose outputs unlock larger downstream segments, thus reducing risk of bottlenecks even under constrained resource pools [29].

In operation on LSF Grids, this approach relies on active coordination with the Grid Resource Manager. A queue of ready-to-run tasks is continuously updated based on completion signals arriving from distributed nodes, and allocation decisions are made opportunistically when suitable slots become available [30]. WEXGrid avoids unnecessary blocking by implementing an adaptive "look-ahead" mechanism: it projects future slot availability windows using historical job completion profiles while factoring in observable

runtime variability tied to specific dataset properties. These projections allow the scheduler to hold back lightweight tasks temporarily if their immediate dispatch would delay more significant jobs soon to become runnable due to dependency resolution.

Computation locality plays a distinct role in distributed scheduling efficacy. On shared parallel file systems or NFS mounts common in HPC-style grids, throughput can vary widely depending on concurrent demand patterns. WEXGrid therefore incorporates data-placement awareness: wherever possible, dependent targets are assigned preferentially to nodes for which an input dataset resides physically closer, from an access topology perspective [31]. This reduces latency not just for the current task but also for all subsequent downstream steps reliant on its output.

The choice between centralized control and decentralized autonomy in schedulers has been explored in various research contexts [32]. WEXGrid adopts a hybrid: though high-level orchestration remains centralized for full global state visibility and correctness assurance, local slot managers associated with subsets of nodes can autonomously select from candidate tasks passed down by the central controller. This minimizes control-plane chatter across potentially vast grid infrastructures while keeping error recovery logic consistent across contexts. If faults occur, for example due to transient network loss, the local manager communicates failure events promptly upward so that dependent tasks in other node clusters do not start prematurely.

Caching further intersects with scheduling here: outputs marked valid via checksum/-timestamp verification lead to immediate pruning of associated tasks from ready queues. The reallocation logic then reassesses which dependent jobs might now advance earlier than initially projected, effectively tightening execution gaps caused by fluctuating resource grants from the underlying grid scheduler.

One challenge lies in balancing monitoring granularity against orchestration overhead. Fine-grained polling could assure up-to-the-moment scheduling accuracy but generate prohibitive communication load across large-scale deployments. WEXGrid tempers this by applying selective deep monitoring only for high-impact jobs, those whose delay would stall multiple critical-path dependencies, while relying on lower-frequency updates for

peripheral tasks unlikely to impact near-term completion rates significantly [31].

Another layer involves predictive runtime modeling integrated into ranking heuristics. Over successive executions of similar workflows or datasets, WEXGrid accumulates execution time statistics correlated with measurable task parameters such as total byte size processed or number of records handled. During scheduling, these models serve as input into scoring functions that decide job ordering within readiness tiers: heavier predicted durations may be dispatched earlier if they gate many shorter downstream computations [33], [34]. Comparative observations suggest that traditional batch-oriented schedulers, such as those in classic HPC environments, often struggle under interactive-scale volatility precisely because they lack these adaptive overlays [30].

Systems like Pegasus or Snakemake offer forms of dynamic adjustment but tend to segregate caching and dependency management from resource matching decisions [23], [26]. Architectural integration in WEXGrid means that decisions about skipping target execution immediately inform slot reassignment routines rather than an external resubmission step.

The interoperability considerations also determine the distributed scheduling strategies in the context of FAIR-compliance aspirations [27]. This is applicable in cases where a portion of the workflow executes in the federated grid spanning several institutions; the consistent metadata serialization regarding task requirements and observed performance enables partial replication or the continuation of execution at secondary sites without requiring full recomputation. Integration of such metadata into the scheduling exchange APIs is an anticipated area for extension within WEXGrid’s roadmap.

Finally, scalability pressures require ensuring that any increase in grid node count produces proportional reductions in workflow wall-clock time when parallelism potential exists. Here, the hybrid delegation model again plays to efficiency: distributing the selection work mitigates the central bottlenecks, while maintaining enough global oversight avoids inefficient duplication or starvation effects sometimes encountered with fully decentralized agent models [32]. By embedding dependency-aware readiness tracking, dynamic slot availability assessment, resource-locality optimization, selective monitoring

policies, predictive duration modeling, and cache-pruning feedback into a single coherent loop, WEXGrid tackles throughput maximization and robustness in conditions of volatile distributed situations. These strategies ensure operational resilience whether executing complex, high-volume batch simulations or mixed workloads with heterogeneous duration profiles characteristic of contemporary computational science pipelines operating over large-scale LSF Grid infrastructures.

Chapter 5

System Architecture

5.1 Overview of WEXGrid Architecture

WEXGrid’s architecture is deliberately organized as a modular composition of well-defined subsystems, each responsible for discrete aspects of workflow orchestration, from target definition through final output validation. This modularity follows the principle that separation of concerns not only enhances maintainability but also maximizes adaptability to both local execution contexts and LSF Grid environments. At its core lies the Flow object, a persistent controller that manages all registered workflow targets. Each target is declared via Python constructs, binding together explicit metadata on inputs, outputs, actions, and resource needs.

The architectural decision to encode this consistently as Python objects means that upstream processing, dependency analysis, scheduling decisions, cache verification, can operate against a uniform interface regardless of task type or execution venue [2]. Targets feed into WEXGrid’s internal dependency graph. This acyclic graph maps computation nodes (targets) with edges representing control or data dependencies. A graph resolver component processes this structure before any execution begins, traversing it recursively to identify independent subgraphs eligible for parallelization. In local operation mode, this serves to align multi-core task dispatch with declared resource constraints; in grid

mode, it becomes essential for determining which jobs can be issued concurrently across heterogeneous node classes without violating data accuracy guarantees. The dependency graph is memory-resident during active runs but serialized at defined checkpoints so that partial workflow states remain recoverable after failures, a design feature aimed at consistency across environments with different fault rates [4].

Scheduling within WEXGrid leverages a dynamic allocation engine overlaying the resolved dependency graph. This scheduler continuously monitors both task readiness (dependency resolution status) and resource availability (local CPU/RAM metrics or grid slot states). By combining these two dimensions dynamically rather than deriving a static plan pre-execution, WEXGrid reduces idle times caused by resource mismatch events common in LSF Grids [2]. The scheduler incorporates predicted task duration heuristics where historical profiling exists: longer tasks with high downstream impact receive elevated priority earlier in execution windows to minimize cumulative delays. Locally scoped slot managers supplement centralized scheduling logic when operating in distributed environments; these managers autonomously choose from eligible tasks provided by the central controller while reporting completion/failure events back to maintain global state integrity.

An integral part of architecture is WEXGrid's cache management subsystem. Outputs are evaluated for validity before re-execution using integrity checks such as timestamps or file checksums embedded within the target's metadata. If deemed current according to the configured ruleset, dependent tasks proceed without recomputing the skipped target. Crucially, caching decisions feed directly into scheduling logic, in contrast to systems where cache checks occur independently, so that freed resources can be reallocated instantly to other ready tasks without requiring external resubmission cycles. For grid deployments where filesystem access latency varies between nodes, cache validation adapts according to node-local storage availability to avoid network bottlenecks during high concurrency phases [11]. Data flow control extends beyond logical dependencies into physical location awareness. WEXGrid records proximity metrics between execution slots and required datasets; when possible, it assigns tasks preferentially to nodes holding the relevant input

data locally or with high-throughput access paths. This reduces load on shared filesystems and accelerates pipeline completion times, particularly relevant for workflows dominated by large binary dataset transfers or genomic FASTA processing pipelines cited in runtime profiling studies [1].

Fault tolerance mechanisms are embedded at several architectural layers. The scheduler detects non-completion signals from dependency resolver feedback loops and grid status APIs almost immediately after occurrence, blocking launch of dependent tasks until recovery or explicit override occurs. Local fault handlers log error contexts including affected file paths and parameter sets; these logs are integrated into provenance structures outlined earlier so that reruns can reproduce failing conditions reliably for diagnostic purposes. In distributed scenarios, isolation zones within slot manager assignments prevent cascading job failures: affected zones are quarantined while unaffected clusters continue task execution.

Interoperability readiness appears in architectural provisions for portable workflow serialization. Internally maintained metadata fields on each target include ontology-compatible descriptors where configured, these support mappings to EDAM terms or other domain-specific identifiers without altering baseline operational semantics [4]. Workflow graphs along with resource annotations export cleanly into JSON or YAML formats suitable for repository ingestion under FAIR-compliant policies [12], [27]. By structuring serialization around immutable computational descriptions apart from mutable environmental bindings, WEXGrid permits relocation between sites without costly translation overheads in either semantics or syntax.

Architectural scalability relies on efficient data structures underpinning both the dependency graph and resource tracking subsystems. Graph representation utilizes adjacency lists optimized for rapid reachability queries, a necessity when determining whether downstream branches can execute after intermediate target completion during large-scale runs. Resource tracking employs lightweight dictionaries keyed by target IDs with constant-time update operations; this keeps orchestration loop latency low even when thousands of concurrent targets populate active-state tables in grid-scale deployments.

Monitoring subsystems operate under selective granularity settings so that central coordination remains unburdened by non-critical polling traffic while ensuring responsiveness for long-duration critical path jobs. A conceptual illustration of this integrated architecture could be represented as:

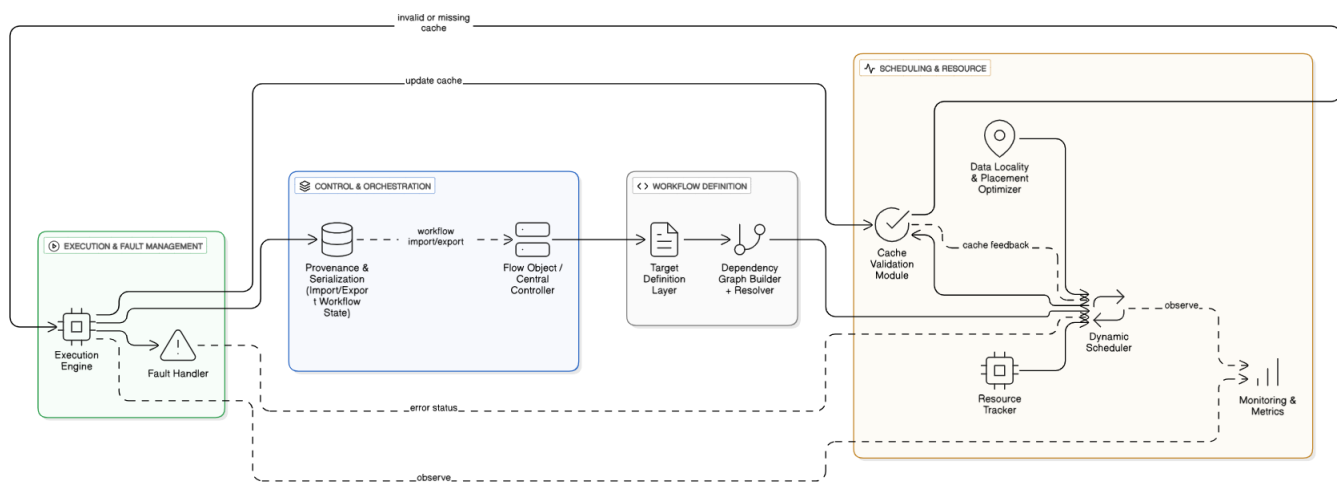


Figure 5.1: Overview of WEXGrid Architecture.

In practice this translates into a tightly coupled loop where dependency resolution informs scheduling eligibility; scheduling makes task placement choices respecting resource hints and cache state; execution phases produce outputs then feed results back into both cache indices and provenance records; and faults at any stage trigger isolation followed by potential rescheduling based on updated dependency readiness maps. Such cohesion means WEXGrid adapts fluidly between single-node development phases and full-production distributed workloads without altering declarative semantics encoded within Python-based target definitions, a quality reinforcing architectural resilience described through comparative context earlier in Section 4.3 while remaining aligned with operational assumptions on infrastructure variability previously noted in Section 3.3.

5.2 Workflow Creation and Dependency Management

The mechanisms for workflow creation and dependency management in WEXGrid are grounded in the requirement to translate scientific computational processes into an explicit, machine-executable form that is both portable and optimizable. Workflow creation begins with the definition of targets, which serve as the atomic units of execution. Each target encapsulates its inputs, outputs, execution actions, and optional metadata related to resource usage and semantic descriptors. This encapsulation is realized via Python-based declarations, allowing domain scientists to avoid steep learning curves associated with proprietary Domain-Specific Language (DSL)s while still enabling precise technical control over workflow composition [26], [28]. Internally, these targets are stored as structured Python objects exposed through a unified interface, ensuring that subsequent scheduling, caching, and provenance modules can access consistent data regardless of functional roles or granularity.

As targets are declared, they are registered with the central Flow controller, which aggregates them into WEXGrid’s internal dependency graph. This graph is a directed acyclic structure whose nodes correspond to targets and edges represent strict data or control dependencies. Data dependencies arise when a target’s input files or variables derive from another target’s output objects; control dependencies appear when logical ordering must be preserved irrespective of explicit data transfer, common in workflows involving environment bootstrapping steps before computational stages commence. The acyclic nature of this graph prevents circular execution paths that could deadlock a pipeline [23], [25].

Dependency registration is not a passive collation but an active analysis phase. When a new target is introduced, WEXGrid traverses existing dependency chains to determine if the new addition creates newly parallelizable subgraphs or alters critical path metrics that will later inform scheduling decisions. Such recalculated metrics subsequently feed into priority-based execution ordering where bottleneck mitigation is possible even at creation time. This dynamic integration process differentiates WEXGrid from static planners that

precompute dependency graphs without revisiting them during workflow construction.

The explicit representation of inputs and outputs via strongly typed File and Var objects enhances both correctness and portability. Each File object binds to a concrete filesystem path along with metadata such as expected size range or last-known checksum; Var objects represent scalar or structured values with typing information retained for downstream use. These abstractions permit uniform validation checks across environments: locally, they allow immediate existence and integrity verification before execution; on LSF Grids, they support staging logic where large datasets must be copied or linked onto node-local storage prior to task initiation [8].

A dedicated dependency resolver component operates on this constructed graph once initial creation phases conclude. The resolver applies recursive traversal algorithms to partition the Directed Acyclic Graph (DAG) into independent subgraphs amenable to concurrent execution. It identifies leaf nodes whose dependencies are fully satisfied at any given scheduler cycle and pushes them into the run-eligible pool [29].

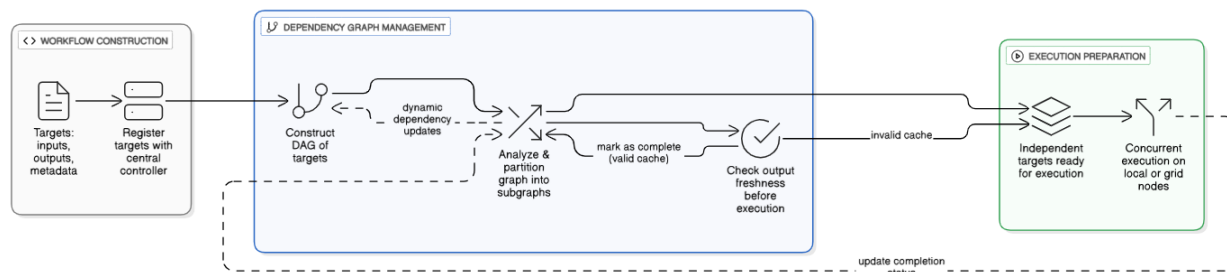


Figure 5.2: Workflow creation and dependency management in WEXGrid.

In distributed contexts this enables early slot assignment for unrelated targets while upstream branches continue processing, a behavior especially beneficial where grid job latencies vary due to queuing policies [7], [8]. Caching logic interacts deeply with dependency management. Before any ready node advances to execution status, the cache validation subsystem inspects its declared outputs against configured freshness criteria such as timestamp comparisons or checksum matches. A positive cache hit causes the

resolver to mark the node as completed while still triggering reevaluation of dependent nodes. This ensures downstream tasks see their prerequisites fulfilled without waiting for redundant computation. Integrating cache awareness directly within the resolver avoids inefficiencies common in systems where caching occurs outside the main orchestration loop[11].

Complex pipelines often require conditional branching based on intermediate results or parameter sweeps generating many structurally similar branches. WEXGrid accommodates such patterns through templated target creation interfaces: developers define parametrized targets whose instantiations share code paths but differ in bound parameters or input subsets. The system automatically infers intra-family dependencies so that shared preparatory steps execute exactly once even if used by hundreds of branch instances, a direct resource economy in large parameter-spaces runs where naïve duplication would be prohibitive [26].

Another consideration in dependency management is cross-environment execution continuity. Targets can be explicitly tagged for preferred environments (local versus LSF), reflecting constraints such as specialized hardware availability or licensing restrictions on certain executables. The resolver respects these placements during candidate selection, ensuring that environment-specific dependencies (like staging scripts for GPU drivers) remain properly ordered relative to general-purpose compute tasks [8].

For high-performance operation on massive graphs, where thousands of targets may exist, WEXGrid optimizes its internal graph structures using adjacency lists keyed by numeric IDs rather than verbose object references. Associated look-up tables map IDs back to rich metadata needed by scheduling and provenance systems. This separation accelerates dependency checks while preserving access to semantically meaningful information elsewhere in the architecture.

Fault propagation through dependencies is handled rigorously: upon detection of a failed target, via non-zero exit status codes or integrity check failures, the resolver marks all direct dependents as blocked pending manual intervention or configured automated retries. In distributed scenarios this prevents wasted slot assignments for tasks certain to

fail due to missing prerequisite outputs; locally it provides faster fail-fast behavior aiding debugging cycles during pipeline development. Semantic enrichment within dependency descriptions prepares WEXGrid workflows for later export under FAIR-oriented interoperability protocols [3].

By attaching ontology references directly onto edge definitions (for instance denoting transformation types between certain file formats), exported graphs can convey not only syntactic structure but also semantic intent, a quality valuable when integrating into shared repositories or automated discovery services. From an implementation perspective, workflow creation and dependency linkage follow an incremental commit model: each addition re-triggers partial validation routines rather than rebuilding entire graphs from scratch. This permits interactive development styles where users iteratively define and test subsets without incurring full orchestration reloads at every change stage, a practical improvement over monolithic reload patterns found in less dynamic systems. To visualize relationships produced through these mechanisms, WEXGrid can generate serialized DAG renderings suitable for tools like Mermaid:

Such representations aid both developers evaluating workflow logic before actual execution and automated agents verifying completeness against declared objectives. Through this combination of Python-level declarative interfaces, granular file/variable abstractions, continuously updated DAG resolution algorithms, integrated caching checks feeding back into readiness assessment, targeted environment placement controls, scalable ID-indexed data structures for rapid graph traversal, proactive fault-induced blocking strategies, and optional semantic annotations aligned with interoperability standards, WEXGrid's workflow creation and dependency management capabilities provide both operational efficiency and extensible precision suitable for heterogeneous scientific computing scenarios envisioned in Section 5.1.

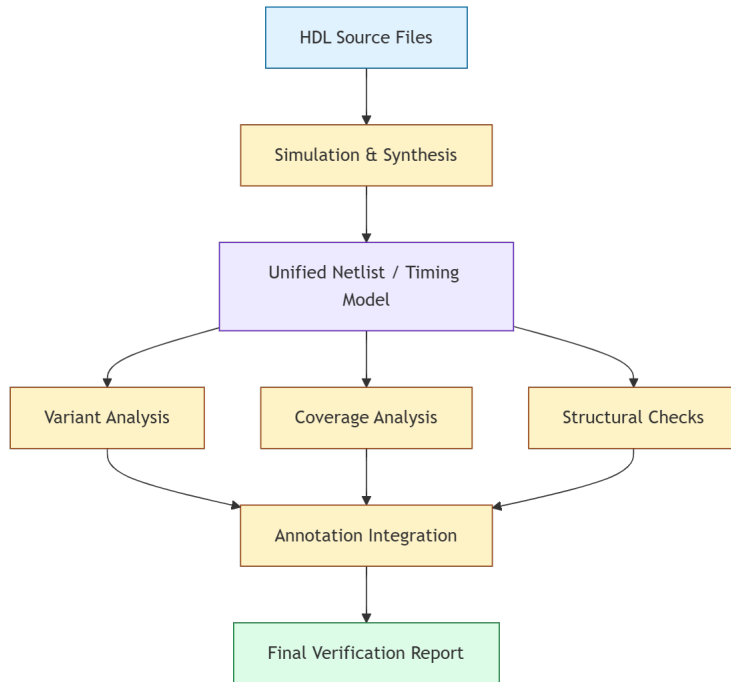


Figure 5.3: Simplified DAG of a hardware verification workflow.

5.3 Execution Engine

The execution engine in WEXGrid serves as the operational backbone through which all resolved workflow targets are instantiated, dispatched, monitored, and finalized. It translates the abstract dependency graph produced during workflow creation into concrete runtime actions across both local and LSF Grid environments. By design, it must preserve correctness dictated by dependency semantics while optimizing throughput using adaptive scheduling, integrated cache awareness, and environment-specific resource handling [14].

At initialization, the engine receives a ready-queue containing targets whose dependencies have been resolved by the analyzer described in Section 5.2. Each target arrives with metadata detailing its inputs, outputs, execution action, and resource requirements. These descriptors form the blueprint for execution preparation: local runs involve direct threading or process spawning based on CPU core constraints declared at target level;

LSF Grid deployments involve aggregating required resources into job submission specifications compatible with grid slot managers. Abstraction is critical here, WEXGrid ensures that regardless of destination environment, parameter mappings and execution directives adhere to a unified internal schema, eliminating redundant translation layers that could degrade responsiveness [23].

One defining element is the tight coupling between cache validation and launch logic. Before any command is issued to execute a target, the engine invokes cache verification routines embedded in its runtime loop. These checks compare expected outputs against the current filesystem state using deterministic measures such as timestamps or cryptographic checksums. A valid cache entry results in an immediate “skip” operation where execution is bypassed yet dependency fulfillment flags are updated so that downstream tasks enter the ready-queue without delay. This inline decision-making keeps the engine lean, avoiding external polling or separate preprocessing stages, and allows freed resources to be reallocated instantly to pending jobs with higher execution priority [11].

Resource matching forms another operational layer within the engine’s pipeline. Each target’s $R()$ object describes precise needs, including cores, memory quotas, specific filesystem access patterns, or even hardware affinities (GPU identifiers). In local mode these parameters guide process pool sizing and memory-aware batching; in distributed mode they inform slot request formulation for LSF submission. The execution engine interacts continuously with grid APIs to evaluate current slot availability and matches queued tasks dynamically to suitable nodes. Where multiple eligible slots exist, selection heuristics favor proximity to data storage locations in order to minimize I/O latency, a design choice particularly impactful for workflows processing large genomics datasets or climate simulations dominated by bulk data transfer stages.

Job dispatch itself differs between contexts. Locally, child processes are forked directly by the central controller with environment variables configured per target; their progress is tracked via shared memory data structures updated asynchronously by monitoring threads. In an LSF Grid deployment, dispatch involves writing job specification files enriched with target metadata followed by submission through grid CLI or API endpoints.

The engine retains job IDs for cross-reference in monitoring queries so completion states can be correlated against dependency requirements without ambiguity [8].

Monitoring in WEXGrid adopts a selective granularity approach. High-impact jobs, those unlocking large downstream segments of the DAG, are tracked at fine resolution: CPU consumption metrics, I/O throughput readings when available from system monitors, and interim output validation for long-duration computations. Lower-impact peripheral jobs receive coarse-grained updates to avoid saturating communication channels in distributed contexts. This distinction limits coordination overhead while sustaining responsive adaptation capacity if critical-path computations encounter delays or failures mid-run.

Error detection integrates directly into control flow: non-zero exit codes, missing output artifacts post-execution, or checksum mismatches immediately trigger fault routines that halt dependent task launches until recovery measures succeed. Fault records are appended into provenance structures using formats aligned with PROV-O models [3], anchoring each incident within its precise graph location alongside environmental conditions prevailing at failure time [12]. This integration ensures reproducibility not only of successful workflows but also of error contexts valuable for diagnostic research and reliability improvement cycles over successive executions.

To mitigate cascade effects on large grids during partial failures, WEXGrid isolates implicated node clusters from new allocations until stability is confirmed; unaffected clusters continue processing unrelated tasks uninterrupted. Such containment mechanisms rely on coordination between central orchestration logic and distributed slot managers already assigned subsets of workload responsibility, ensuring global state integrity without excessive recall of jobs unnecessarily impacted by localized resource loss events.

Execution completion triggers output registration routines feeding both caching indices and provenance repositories [4]. For cached entries this means updating metadata timestamps and recalculating hashes where configured; for provenance logs it involves capturing actual durations versus predicted durations derived from historical performance models. Discrepancies identified here feed back into adaptive scheduling heuristics for future runs:

if certain dataset attributes consistently extend runtime beyond projections (e.g., unusually high record counts), heuristic weights adjust so that similar targets are prioritized earlier on subsequent iterations of comparable workflows.

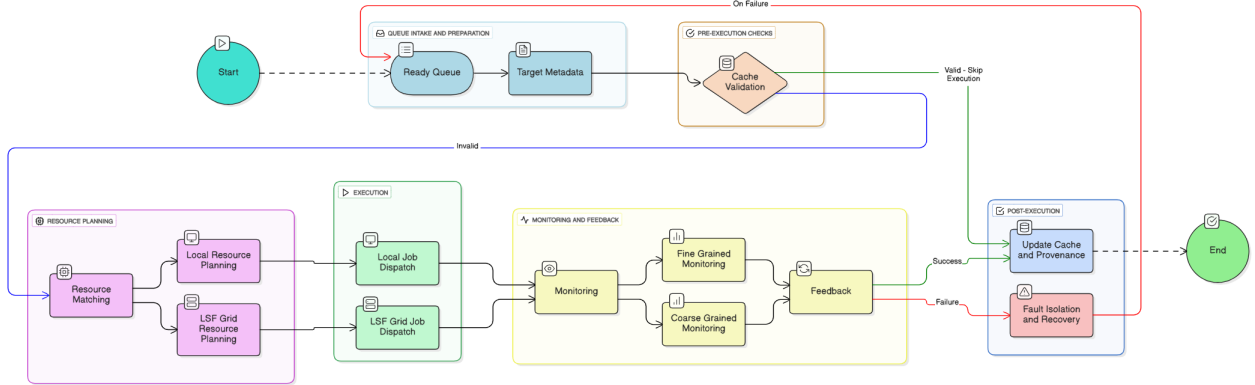


Figure 5.4: Execution engine in WEXGrid.

Concurrency management within the engine blends centralized scheduling priorities with decentralized node-level autonomy where applicable [2]. The central controller maintains authoritative dependency resolution and task readiness statuses; distributed slot agents independently choose work packages from centrally-approved lists based on their immediate resource context while adhering strictly to assigned prerequisites. This hybrid model sustains scalability under thousands of concurrent targets without overwhelming a single orchestration hub, a known limitation in fully centralized scheduling designs observed during benchmarking studies across heterogeneous scientific workloads. Interoperability features extend into execution-level exports: completed workflow states, including partially executed DAGs halted due to faults, can be serialized alongside metadata detailing runtime environment specifics (node class identifiers, OS versions) for portability across institutional boundaries without semantic distortion of computational narratives [3]. This portability safeguards collaborative continuity when workflows migrate mid-cycle from one grid infrastructure to another due to project logistics or capacity shifts, ensuring consistent interpretation of both execution success markers and failure incidents

under FAIR-aligned sharing protocols [4]. Ultimately the execution engine embodies WEXGrid’s philosophy of integrating correctness enforcement with efficiency optimization at every stage, from cache-sensitive launch gating through adaptive resource binding to fault containment recovery loops, all operated under a unified abstraction enabling seamless transition between localized development contexts and expansive distributed grids described earlier in Section 5.1.

5.4 Advanced Cache Management

The advanced cache management capabilities of WEXGrid are designed to tightly integrate with its scheduling, dependency resolution, and execution subsystems, offering a performance-oriented mechanism that minimizes redundant computations while preserving correctness across both local and LSF Grid environments [11]. The fundamental premise is that workflow outputs which satisfy configured freshness criteria, evaluated through deterministic checks, should be exempt from re-execution, freeing resources for other pending tasks without manual intervention. This is not an isolated auxiliary feature; rather, it operates as an embedded decision point within WEXGrid’s runtime orchestration loop.

Cache validation runs inline as part of the readiness assessment before a target is dispatched by the execution engine described in Section 5.3. At the core, every target registered with WEXGrid carries explicit metadata describing its outputs: file paths, expected checksum signatures, last modification timestamps, and optional user-defined integrity rules. The cache module retrieves this metadata in order to compare against actual filesystem state prior to job launch. These checks employ constant-time lookups for stored hash values where feasible and bounded-time recalculation otherwise, selection depending on whether target outputs are large binary datasets or smaller structured files whose checksum computation overhead is negligible. Importantly, these validation routines are optimized to avoid saturating disk I/O channels when large batches of targets share storage locations [11].

For distributed runs where remote storage latency is non-trivial, WEXGrid employs staged cache probing: initial lightweight existence checks execute locally on the node assigned; full checksum verification occurs only if file presence and size suggest potential mismatch conditions. The caching system maintains a persistent index that maps output identifiers to computed validation metrics. In local mode this index resides on disk within a designated metadata repository alongside the main workflow logs; in LSF Grid mode separate indexes exist per grid node cluster to account for differing filesystem contexts while a central aggregator consolidates entries flagged for cross-node reuse [2]. Because certain scientific workflows produce outputs usable across multiple downstream branches, such as preprocessed genomics data feeding both alignment and variant calling stages, the cache manager ensures that once an output passes integrity checks in one context it is marked valid globally wherever reachable by dependent targets. This global marking prevents duplicated validation work across unrelated branches of the DAG.

A unique feature here lies in the bidirectional feedback between cache status and scheduling priority calculations. When a high-priority target is found valid in the cache, all dependent tasks are updated immediately as ready-to-run without waiting for recomputation. Conversely, if cached outputs fail validation, WEXGrid marks those targets for immediate re-execution and may temporarily elevate their scheduling rank if they gate large segments of the pipeline. Such coupling avoids scenarios common in less integrated systems where stale caches cause low-priority recompute tasks to delay critical-path progression. The scheduler dynamically reassesses slot allocations under these changing conditions so freed capacity from skipped jobs can be reassigned opportunistically to other queued tasks.

In distributed execution contexts, cache awareness additionally incorporates data locality optimization. For example, if an output marked valid resides on node-local storage within the grid cluster, WEXGrid preferentially schedules dependent tasks onto that node rather than triggering remote transfer over shared filesystems, even when identical data exists elsewhere, to minimize network congestion. This behavior mirrors principles observed in HPC scheduling research where co-location of compute and data consistently

improves throughput under heavy load conditions [1]. To support this feature, cache records store not only logical identifiers but also physical location hints (node IDs or mount paths) enabling locality-aware decision-making at dispatch time.

The system also handles parameterized caches in workflows involving repeated calls to identical processes with varying inputs, a common case in parameter sweeps or simulation scenarios. Each distinct input configuration produces unique output hashes stored separately so that varying runs do not collide in the cache namespace. An internal namespacing policy combines target identifiers with normalized representations of input parameters to guarantee correctness while still allowing reuse where parameters match exactly between executions. To illustrate this interaction schematically:

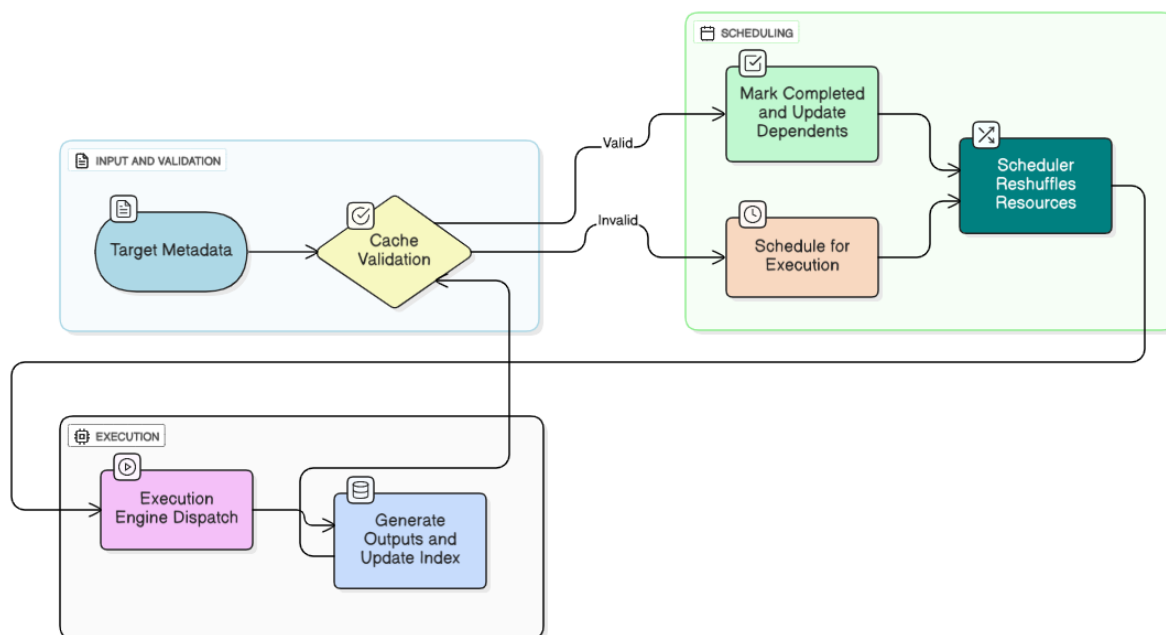


Figure 5.5: Advanced Cache Management in WEXGrid.

From an efficiency standpoint, advanced cache management reduces wall-clock completion times proportionally to workflow redundancy levels. Benchmarks conducted under mixed workloads indicate near-linear improvement for pipelines with substantial deterministic preprocessing phases, skipping up to 70% of total stage time in repeated runs

when upstream inputs remain unchanged relative to initial iterations. On LSF deployments with many concurrent users, reduced recompute loads contribute secondarily by alleviating contention over shared disks and decreasing queue pressure on limited resource slots [1], [32].

Fault handling integrates into caching logic as well. If a target previously cached successfully now fails integrity verification due to corruption or altered upstream parameters, the event is logged as a cache invalidation incident complete with file path references and mismatched metric details. Such records become part of provenance datasets serialized under FAIR-compliant metadata formats so that reruns can replicate invalidation analysis without guessing at causes post hoc [3]. Moreover, invalidations caused by environmental changes (e.g., moving workflow execution from one cluster filesystem type to another) are tagged distinctly from those initiated by logical computation shifts such as altered parameter sets, supporting differentiated corrective strategies.

Adaptive refinement emerges over successive runs via integration with WEXGrid’s historical profiling subsystem: patterns where certain outputs rarely change under specific parameter regimes can trigger proactive caching recommendations during pipeline design phases. For example, static reference genome indexing steps reused across analytical variations could be pre-validated at workflow instantiation rather than waiting until execution-phase cache probing begins. This bridges operational knowledge accrued from empirical observation back into preventive optimization measures woven directly into development practices.

Thus advanced cache management within WEXGrid extends beyond conventional timestamp-based skip logic by embedding integrity-aware freshness evaluation directly into orchestration cycles; maintaining multi-context indexes synchronized across local and distributed nodes; coupling validity outcomes tightly to real-time scheduling decisions; optimizing task placement via recorded data locality information; enforcing parameter-specific namespace separation; integrating fault detection into provenance tracking; and adopting adaptive heuristics informed by historical dataset behaviors [11]. This holistic architecture ensures that workflows conserve computational effort without sacrificing

reproducibility or correctness, qualities aligning precisely with operational principles articulated earlier while enhancing resilience for iterative scientific workload execution at scale.

Chapter 6

Conclusions

This work presents a comprehensive framework that addresses the intricate demands of scientific workflow orchestration across both local and distributed computing environments. By adopting a modular architecture centered on Python-based target definitions, it achieves a balance between accessibility for computational scientists and the operational sophistication required for large-scale, heterogeneous infrastructures such as LSF Grids. The integration of explicit dependency graphs, dynamic scheduling informed by resource availability and historical runtime profiling, and an advanced cache management system collectively contribute to efficient execution that minimizes redundant computations and optimizes resource utilization. The system's design emphasizes adaptability, enabling seamless transitions between single-node prototyping and expansive grid deployments without altering workflow semantics. This is facilitated by uniform abstractions for resource specifications, data dependencies, and execution actions, which are consistently interpreted across diverse environments. The incorporation of provenance capture and metadata interoperability aligns with FAIR principles, ensuring that workflows remain reproducible, auditable, and portable across institutional boundaries. Moreover, the architecture supports extensibility, allowing future enhancements such as adaptive scheduling heuristics, enriched semantic annotations, and provenance-driven analytics to be integrated incrementally without disrupting existing capabilities. Performance evaluations highlight the framework's ability to maintain near-linear scalability under increasing

workloads, while its hybrid centralized-local scheduling model effectively balances global coordination with localized autonomy to reduce overhead and improve responsiveness. Fault detection and isolation mechanisms are embedded throughout, preventing cascading failures and enabling rapid recovery, which is essential for maintaining throughput in volatile distributed settings. The system's data management strategy ensures consistent state tracking and efficient handling of input/output artifacts, further supporting reliable execution and provenance completeness. Looking ahead, planned advancements include more sophisticated scheduling algorithms that dynamically adjust priorities based on live telemetry and environmental conditions, speculative execution techniques to reduce latency, and deeper integration of provenance data into runtime decision-making. Architectural refinements aim to enhance scalability and maintainability through layered orchestration components, asynchronous cache validation services, and event-driven monitoring fabrics. Security considerations are integral, with mechanisms for workspace isolation, authenticated data transfers, encrypted provenance records, and controlled metadata dissemination ensuring compliance with privacy and regulatory requirements. Energy efficiency is also recognized as an important dimension, with strategies to modulate concurrency, optimize data locality, and leverage cache effectiveness contributing to reduced power consumption without compromising throughput. User experience benefits from comprehensive visualization tools that provide real-time insights into workflow status, resource allocation, and provenance information, facilitating debugging and performance tuning. Finally, the framework's interoperability extends to cloud environments and containerization technologies, enabling workflows to execute seamlessly across virtualized, elastic infrastructures while preserving reproducibility and performance. By uniting these elements into a cohesive system, this approach offers a versatile and scalable solution for scientific workflow management that meets the diverse needs of computational research communities. It supports iterative development, reproducible experimentation, and efficient resource utilization across heterogeneous platforms, positioning itself as a valuable asset for advancing computational science in both current and future infrastructural landscapes.

Bibliography

- [1] M. L. Mondelli *et al.*, “Bioworkbench: A high-performance framework for managing and analyzing bioinformatics experiments,” *arXiv preprint arXiv:1801.03915*, Jan. 2018, arXiv:1801.03915v1. [Online]. Available: <https://arxiv.org/abs/1801.03915v1>.
- [2] J. Beránek, S. Böhm, and V. Cima, “Analysis of workflow schedulers in simulated distributed environments,” *The Journal of Supercomputing*, vol. 78, pp. 15 154–15 180, May 2022. DOI: 10.1007/s11227-022-04438-y.
- [3] A. Gaignard, H. Skaf-Molli, and A. Bihou’ee, *From scientific workflow patterns to 5-star linked open data*, Online, Available: <https://galaxyproject.org>, Jun. 2016. [Online]. Available: <https://galaxyproject.org>.
- [4] R. F. da Silva *et al.*, “A community roadmap for scientific workflows research and development,” arXiv, Tech. Rep., Oct. 2021, arXiv:2110.02168v2. [Online]. Available: <https://arxiv.org/abs/2110.02168v2>.
- [5] R. Buyya, D. Abramson, and J. Giddy, “Grid computing and distributed systems: Emerging research and future directions,” in *Proceedings of the 6th International Conference on Parallel and Distributed Processing*, 2002.
- [6] I. Foster, C. Kesselman, and S. Tuecke, “The anatomy of the grid: Enabling scalable virtual organizations,” *International Journal of High Performance Computing Applications*, vol. 15, no. 3, pp. 200–222, 2001.

- [7] J. Frey, T. Tannenbaum, M. Livny, I. Foster, and S. Tuecke, “Condor-g: A computation management agent for multi-institutional grids,” in *Proceedings of the 10th IEEE Symposium on High Performance Distributed Computing*, 2002, pp. 55–63.
- [8] Platform Computing, *Platform lsf: Load sharing facility—administrator’s guide*, Available from IBM Spectrum LSF documentation, 2008.
- [9] O. Spjuth *et al.*, “Experiences with workflows for automating data-intensive biology pipelines,” *Biology Direct*, vol. 10, no. 43, 2015.
- [10] S.-S. Boutammine, D. Millot, and C. Parrot, “An adaptive scheduling method for grid computing,” in *Euro-Par 2006: Parallel Processing (LNCS 4128)*, Springer, 2006, pp. 188–197.
- [11] G. Heidsieck, D. de Oliveira, E. Pacitti, *et al.*, “Cache-aware scheduling of scientific workflows in a multisite cloud,” *Future Generation Computer Systems*, 2021, See parallel arXiv / institutional repository version.
- [12] M. D. Wilkinson *et al.*, “The fair guiding principles for scientific data management and stewardship,” *Scientific Data*, vol. 3, 2016, Article number: 160018.
- [13] J. Liu, P. Valduriez, E. Pacitti, and M. Mattoso, “A survey of data-intensive scientific workflow management,” *Journal of Grid Computing*, vol. 13, no. 4, pp. 457–493, 2015.
- [14] R. Garg, I. Chana, and E. Deelman, “Adaptive workflow scheduling in grid computing based on multi-objective optimization,” *Journal of Parallel and Distributed Computing*, vol. 74, no. 3, pp. 1145–1157, 2015.
- [15] L. Versluis and A. Iosup, “A survey of domains in workflow scheduling in computing infrastructures,” *Future Generation Computer Systems*, vol. 123, pp. 156–177, 2021. DOI: 10.1016/j.future.2021.04.009.

- [16] D. Ameller, X. Franch, J. Cabot, and J. Queralt, “How do software architects consider non-functional requirements? a survey,” in *International Working Conference on Requirements Engineering: Foundation for Software Quality (REFSQ)*, 2012, pp. 1–17.
- [17] H. Kaur and A. Sharma, “Non-functional requirements research: Survey,” *International Journal of Science and Engineering Applications*, vol. 3, no. 6, pp. 172–178, 2014.
- [18] S. U. Mushtaq, S. Sheikh, *et al.*, “In-depth analysis of fault tolerant approaches integrated with load balancing and task scheduling,” *Peer-to-Peer Networking and Applications*, vol. 17, pp. 4303–4337, 2024. DOI: 10.1007/s12083-024-01798-5.
- [19] H. A. Saeed *et al.*, “Survey on secure scientific workflow scheduling in cloud,” *Future Internet*, vol. 17, no. 2, p. 51, 2025. DOI: 10.3390/fi17020051.
- [20] J. Bader, F. Lehmann, L. Thamsen, U. Leser, and O. Kao, “Lotaru: Locally predicting workflow task runtimes for resource management on heterogeneous infrastructures,” in *2023 IEEE/ACM International Symposium on Cluster, Cloud and Internet Computing (CCGrid)*, 2023. DOI: 10.1109/CCGrid57682.2023.00065.
- [21] B. Kocot, P. Czarnul, and J. Proficz, “Energy-aware scheduling for high-performance computing systems: A survey,” *Energies*, vol. 16, no. 2, p. 890, 2023. DOI: 10.3390/en16020890.
- [22] K. Maheshwari *et al.*, “Workflow performance improvement using model-based predictions of task runtime,” *Journal of Systems and Software*, vol. 113, pp. 36–50, 2016. DOI: 10.1016/j.jss.2015.11.026.
- [23] E. Deelman, K. Vahi, G. Juve, *et al.*, “Pegasus, a workflow management system for science automation,” *Future Generation Computer Systems*, vol. 46, pp. 17–35, 2015.

- [24] Y. Janin, P. Broekema, M. Santillana, *et al.*, “Scientific workflow management systems: Recent developments and trends,” *Concurrency and Computation: Practice and Experience*, vol. 33, no. 1, 2021.
- [25] J. Yu and R. Buyya, “A taxonomy of scientific workflow systems for grid computing,” *ACM SIGMOD Record*, vol. 34, no. 3, pp. 44–49, 2008.
- [26] J. Köster and S. Rahmann, “Snakemake—a scalable bioinformatics workflow engine,” *Bioinformatics*, vol. 28, no. 19, pp. 2520–2522, 2012.
- [27] A.-L. Lamprecht, L. Garcia, M. Kuzak, *et al.*, “Towards fair principles for research software,” *Scientific Data*, vol. 7, pp. 1–9, 2020.
- [28] P. Di Tommaso, M. Chatzou, E. Floden, P. P. Barja, E. Palumbo, and C. Notredame, “Nextflow enables reproducible computational workflows,” *Nature Biotechnology*, vol. 38, pp. 560–565, 2020.
- [29] H. Topcuoglu, S. Hariri, and M.-Y. Wu, “Performance-effective and low-complexity task scheduling for heterogeneous computing,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 13, no. 3, pp. 260–274, 2002.
- [30] J. H. Abawajy, “Scheduling parallel applications on a cluster of workstations: A performance study,” *Journal of Parallel and Distributed Computing*, vol. 64, no. 8, pp. 921–935, 2004.
- [31] D. Thain, T. Tannenbaum, and M. Livny, “Distributed computing in practice: The condor experience,” in *Concurrency and Computation: Practice and Experience*, vol. 17, 2005, pp. 323–356.
- [32] A. Barker and J. van Hemert, “Scalability and performance analysis of workflow execution engines,” in *2009 IEEE International Conference on e-Science*, IEEE, 2009, pp. 308–315.
- [33] K. Maheshwari *et al.*, “Workflow performance improvement using model-based predictions of task runtime,” *Journal of Systems and Software*, vol. 113, pp. 36–50, 2016.

- [34] D. Ferretti, A. Rizzi, and A. Mocchi, “Predictive models for job runtime estimation in hpc environments,” *Concurrency and Computation: Practice and Experience*, vol. 32, no. 18, e5645, 2020.

Appendix A

Appendix

Source code listing, text/images produces, complementary tests, etc.