# Fractal: Fault-Tolerant Shell-Script Distribution

Zhicheng Huang\* Brown University Ramiz Dundar\*

Brown University

Yizheng Xie Brown University

Konstantinos Kallas University of California, Los Angeles Nikos Vasilakis Brown University

### **Abstract**

This paper presents FRACTAL, a new system that offers fault tolerant distributed shell script execution with unmodified scripts. FRACTAL first identifies recoverable regions from side-effectful ones, and augments them with additional runtime support aimed at fault recovery. It employs precise dependency and progress tracking at the subgraph level to offer sound and efficient fault recovery. It minimizes the number of upstream regions that are re-executed during recovery and ensures exactly-once semantics upon recovery for downstream regions. Evaluation on 4- and 30-node clusters indicates average failure-free speedups of (1)  $>9.6\times$ over Bash, a single-node shell-interpreter baseline, (2)  $> 5.5 \times$ over Hadoop Streaming, a MapReduce system that supports language-agnostic third-party components, and (3) 17% over DISH, a state-of-the-art failure-intolerant shell-script distribution system—all while recovering  $7.8-16.4\times$  faster than Hadoop Streaming in cases of failures.

#### 1 Introduction

The Unix shell remains the 8<sup>th</sup> most popular language on GitHub in 2024 [20], widely used for a variety of workloads [17, 27, 29, 50, 57]. Its popularity can be attributed to several characteristics, including (1) language-agnosticism, flexibly composing an arsenal of task-specific components available in a variety of languages, and (2) dynamism, providing features such as command substitution, variable expansion, and file system reflection.

Unfortunately, these characteristics complicate *fault-tolerant* shell-script scale out. The black-box nature of third-party components complicates recovery after node failures by hindering internal state tracking and limiting scale-out opportunities. Dynamic behaviors and arbitrary side effects make re-executing script failed fragments challenging, affecting the correctness of re-executed scripts. Tolerating faults is

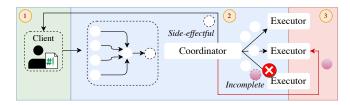


Fig. 1: FRACTAL's high-level workflow. FRACTAL (1) isolates side-effectful regions from recoverable regions; (2) executes recoverable subgraphs on nodes, tracking locality, dependencies, progress, and health; (3) detects failures, re-scheduling the minimal set of unfinished subgraphs for re-execution.

often at odds with retaining the shell's expressiveness without requiring users to modify existing (often legacy) scripts.

Typical approaches such as checkpointing [4,5,7,35,60], barriers [10,11,25], and lineage [28,66] are ill-suited for the setting of the shell (*viz.* §2). As a result, while resaerch on and around the shell is exploding [23,24,32,38,47,49,62], currently no system tolerates failures during the distributed shell-script execution.

**Fault-tolerant distribution with Fractal**: FRACTAL is a system supporting fault-tolerant shell-script distribution: it operates on existing shell scripts without modification, supports all of the shell's dynamic features, allows for language-agnostic composition of black-box components, and is able to recover from node failures.

FRACTAL (Fig. 3) begins by building a dataflow graph of the user's POSIX shell script via PaSh-JIT [32]. Its coordinator then uses command annotations to identify side-effectful commands unsafe for fault-tolerant distribution. Next, it wraps inter-subgraph edges lightweight remote-pipe primitive that records byte-level progress and enforces exactly-once semantics. FRACTAL's per-subgraph heuristic decides whether to persist outputs locally to balance runtime overhead against recovery speed. At runtime, after subgraphs are scheduled to executor nodes, FRACTAL's health and progress monitors continuously track inter-subgraph dependencies and

<sup>&</sup>lt;sup>1</sup>Authors contributed equally. Zhicheng is now with the University of California, San Diego. Ramiz is now with Google.

TEL 1 C '	C C 1, , 1	1 ' 1	1 '1 ' 1 11 ' '
Iah. I: Comparis	on of fault folerance	mechanisms across ke	ev desiderata in shell scripts.
I about 1 Companie	on or runt torerunce	meenamonio aeross ne	y desiderata in shen sempts.

Deside	ratum	<b>Checkpointing</b> [4,5,7,35,60]	Barrier-based [10, 11, 25]	Lineage-based [28,66]	FRACTAL
D1 Ha	andles Black-Box State	No	Yes	No	Yes
D2 Ad	d-Hoc Pipe Streaming Integrity	No	No	No	Yes
D3 Sid	de-Effect Management	Partial	No	No	Yes
D4 Dy	ynamism Compatibility	Partial	No	No	Yes
D5 Re	ecovery Granularity	Coarse	Coarse	Fine	Fine
D6 No	o Script Modification	Partial	Partial	No	Yes

byte-level delivery to detect failures. When a failure occurs, its coordinator computes exactly which subgraphs—and any upstream fragments whose outputs were not persisted—must be replayed. By re-executing only this minimal set of fragments, FRACTAL eliminates redundant work and guarantees exactly-once semantics across all downstream regions.

FRACTAL also introduces a new subsystem, **frac** (for *fracture*), for injecting runtime faults for automatically tuning key parameters and aiding recovery characterization—potentially useful to other distributed systems that combine black-box components and thus released as a separate tool.

**Key results**: In failure-free scenarios, FRACTAL achieves substantial performance improvements, delivering an average speedup of  $9.6\times$  over Bash, a standard shell interpreter,  $5.5\times$  over Hadoop Streaming (AHS), cluster-computing system incorporating black-box Unix commands, and 17% over DISH, a recent *failure-intolerant* scaleout system.

FRACTAL recovers from faults within  $1.26\times$  of the script's fault-free runtime, achieving a  $9.3\times$  speedup over AHS—while supporting improved expressiveness and no manual modifications to the source programs.

**Paper outline and contributions**: The paper starts with a discussion about the design lanscape for fault-tolerant shell-script distribution (§2). With a motivating example script example reifying challenges in tolerating faults (§3), it introduces the FRACTAL design overview. It then proceeds with FRACTAL's key subsystems (§4–§6):

- Execution engine (§4): FRACTAL's remote-pipe instrumentation, progress and health monitors, and the executor runtime work in synergy to provide efficient and precise recovery.
- Performance optimizations (§5): FRACTAL's targeted optimizations to minimize overhead on the critical path, including event-driven execution, buffered I/O, and batched scheduling.
- Fault injection (§6): FRACTAL's **frac** tool enables precise, large-scale fault injections to characterize recovery behavior under real-world conditions.

The paper then presents FRACTAL's evaluation (§7), related work (§8), and conclusion (§9).

**Availability**: Upon acceptance, FRACTAL's implementation and evaluation—currently hosted in a private repository—will

be made publicly accessible as an MIT-licensed open-source project at github/blind/fractal.

### **2** Fault Tolerance for Shell Script

This section begins by outlining the desiderata for fault-tolerant shell-script distribution (Table 1, col. 1), derived from the shell's unique characteristics. It then examines the limitations of existing fault-tolerance mechanisms in meeting these desiderata (Table 1, cols. 2-4). Finally, it presents FRACTAL's design for meeting these desiderata (Table 1, col. 5).

We assume that *worker nodes* may crash in either fail-stop or fail-restart fashion, mirroring typical large-scale deployments on commodity hardware. The *control plane* or the *client node* resilience lies outside the scope for this discussion and can be mitigated with established techniques including durable state logging, consensus protocols, and leader election via ZooKeeper [26].

#### 2.1 Desiderata

Shell scripts uniquely blend diverse commands, streaming pipelines, flexible control flow, and dynamic expansion at runtime. While this provides unmatched expressiveness and simplicity, it complicates fault-tolerant execution significantly. To guide our design, we identify six key fault-tolerance desiderata that any robust shell-script distribution mechanism must satisfy.

- **D1 Black-box state handling**: Shell pipelines invoke external binaries (*e.g.*, sort, grep, unzip) whose internal state cannot be inspected or checkpointed. Such commands may hold gigabytes of data internally without any API for partial snapshots, making it impossible to selectively roll back and resume. A fault-tolerance scheme must recover progress without requiring analysis of or hooks in these opaque commands.
- **D2** Ad-Hoc pipe streaming integrity: Shell commands communicate via unstructured byte streams over ad-hoc UNIX pipes, with arbitrary buffering, chunking, and transformation semantics that vary by command. Failures leave no record of how many bytes or which logical "records" were consumed, and replaying an opaque stream risks duplicating or dropping data. Under failure, the system must guarantee

exactly-once delivery so that no data is lost or duplicated despite retries.

- **D3 Side-effect management**: Shell commands often perform non-idempotent external actions (*e.g.*, appending to files, making network calls, updating system states) that modify the environment. Simply re-running a partially completed side-effectful command can append data again or re-trigger external actions. Recovery must prevent repeated side effects or orphaned partial writes.
- **D4 Dynamism compatibility**: Shell scripts resolve control flow and command invocations only at runtime via loops, conditionals, and variable expansions. A fault-tolerance design must support these on-the-fly pipelines and not assume a static operator graph.
- D5 Fine recovery granularity: Shell pipelines often chain many long-running commands so re-executing whole stages or other coarse-grained recovery units wastes substantial work. Achieving fine-grained replay is further complicated by blackbox commands with opaque internal state and by ephemeral outputs already consumed downstream. Therefore, a robust recovery mechanism must isolate and replay only the minimal affected fragment of the workflow.
- D6 No script modification: Shell scripts are notoriously difficult to program and maintain [18, 21], and modifying or re-implementing legacy scripts can be costly and errorprone [14, 15]. Thus, forcing rewrites is prohibitive. An ideal solution should preserve existing scripts unchanged, transparently adding recovery support.

### 2.2 Existing Approaches

Here, we analyze existing fault-tolerance mechanisms and their limitations in meeting the desiderata for shell-script distribution, using three representative paradigms: checkpointing, barrier-based, and lineage-based systems (some systems incorporates facets from multiple paradigms). It is worth noting that, in practice, failing to meet even one of the desiderata may be enough to disqualify a system as a viable solution for shell scripts.

Checkpointing: Checkpointing-based systems capture periodic snapshots of process or operator state [4,5,7,35,60]. This model succeeds when the runtime exposes hooks into each operator, as in streaming engines or controlled process trees. D1 checkpointing systems require components to implement APIs such as get-processing-state [5] and getState [44] to retrieve internal state, but this approach is not viable due to the opaque nature and language agnosticism of shell commands. Incremental or per-operator snapshots reduce overhead but still require instrumentation inside each command, which is impossible for opaque shell binaries. D2 Frameworks like Flink [4], Storm [60], and Kafka [35] embed barriers or use offset-tracked logs, but UNIX pipes

lack any barrier semantics or offset markers, making it infeasible to checkpoint and resume a byte stream without rewriting the pipeline around an external broker. D3 Transactional sink APIs (e.g., Flink's TwoPhaseCommitSinkFunction) manage side-effect writes within the framework, but shell scripts perform arbitrary file and network I/O outside any transaction boundary. D4 Traditional checkpointing assumes a fixed set of operators known ahead of time, complicating recovery when new processes appear mid-execution. CRIU only snapshots processes it has been explicitly told to monitor, whereas shell loops and eval spawn new binaries at runtime without notifying CRIU—capturing those on-the-fly children would require continuous shell-level hooks to register each new process, which is impractical. D5 Full snapshots capture entire pipelines or process trees, imposing high overhead and causing unnecessary re-execution for chains of short-lived commands; per-command checkpointing would introduce prohibitive runtime overhead and complex coordination. D6 Adopting checkpointing for shell scripts demands wrapping or replacing every invocation, violating the no-modification requirement for legacy or ad-hoc scripts.

Barrier-based: Barrier-based systems such as MapReduce [10, 11] achieve fault tolerance by retrying entire map or reduce tasks upon failures, relying on a static task graph. D1 While this model supports black-box tasks, it lacks the ability to resume partially completed computation: failed components must restart from scratch, even if most work had completed. D2 Barrier-based models are not ideal for streaming data; their retry model simply replays upstream outputs, leading to potential duplication and breaking exactly-once guarantees. Streaming extensions (e.g., Kafka Streams [35]) guarantee exactly-once by buffering entire micro-batches or writing to durable topics, but to adapt raw UNIX pipes one must replace each | with a brokered topic or RDD stage. This forces serialization and network hops for every pipeline edge and introduces head-of-line blocking at batch boundaries, undermining the shell's low-latency, in-memory streaming model. D3 Barrier-based retries rerun every command in a failed task including any non-idempotent side-effectful commands such as file appends or HTTP calls. Because there is no transactional or deduplication API at the shell level, each retry re-issues external side-effects, and preventing duplicates requires invasive wrappers or bespoke idempotency logic around every shell command, violating transparent execution. D4 The static task graph must be fully specified before execution; shell scripts that spawn new commands via loops or eval cannot be dynamically incorporated, leaving on-thefly pipelines untracked. D5 Recovery granularity is fixed at task boundaries; defining each shell command as its own task could narrow scope but forces scripts to be restructured into dozens or hundreds of map/reduce jobs, incurring prohibitive scheduling overhead. D6 D6 Finally, while MapReduce does not mandate full rewrites, it still requires structuring logic into map and reduce phases—limiting flexibility and

imposing extra effort when adapting existing or evolving shell scripts. Hadoop Streaming [25] supports arbitrary binaries but still forces explicit mapper/reducer wrappers, violating the no-modification desideratum.

Lineage-based: Lineage-based fault tolerance mechanisms, as in Dryad [28] and Spark [66], record a DAG of operator dependencies and recover by replaying only failed tasks. While effective for deterministic dataflows within a single framework, they struggle with the ad-hoc, mixed environment of shell pipelines. D1 Lineage frameworks assume each operator is a pure function of its visible inputs and outputs, but black-box shell commands (e.g., sort, uniq) buffer and transform data internally without producing retrievable artifacts, so their progress cannot be reconstructed from lineage alone. D2 Spark Streaming enforces exactly-once by slicing streams into micro-batches with checkpointed offsets and write-ahead logs, but ad-hoc UNIX pipes are unbounded byte streams with no batch boundaries or offset metadata, making transparent mid-stream resume or replay impossible without rewriting each pipe as a Spark streaming stage. D3 Lineage frameworks mitigate side-effects via transactional sinks only if every write goes through their API, but shell commands perform arbitrary I/O (e.g., », mv) outside any transaction boundary, requiring invasive wrappers to prevent duplicates. D4 Dynamic DAG registration in systems like Spark Structured Streaming still requires user callbacks (e.g., writeStream), whereas shell loops and eval spawn processes silently, leaving lineage unnotified and unprepared to recover new branches. D5 Lineage-based models can recompute with a fine granularity as long as the tasks and dependencies are explicitly captured within the lineage framework, where they recompute only the failed partition(s) and their direct dependencies. D6 Lastly, lineage-based approaches impose significant restructuring burdens on script authors; shell scripts must be rewritten into deterministic, functional transformations conforming to the lineage system's programming model, a requirement particularly cumbersome for legacy or rapidly evolving scripts.

#### 2.3 Our approach

At its core, FRACTAL treats a program fragment, as the atomic unit of computation. This design avoids costly instrumentations of every individual shell command while remaining finegrained enough to avoid large-scale re-execution. At fragment boundaries, FRACTAL injects minimal runtime primitives that transparently track progress and enable precise recovery without affecting the internal logic of any black-box command.

This model addresses the key challenges of shell-script fault tolerance. D1 By tracking only inputs and outputs for each fragment, we never peek inside a command's memory or file descriptors. Each command remains unmodified; recovery works solely from its byte-stream boundaries. D2 Byte-level progress tracking guarantees no data loss or duplication, even

```
#1/bin/bash
   in=${in:-$TOP/log-analysis/nginx-logs}
   out=${out:/outputs}
   bots='Googlebot|Bingbot|Baiduspider|Yandex|'
   mkdir $out && hdfs dfs -mkdir /log-analysis
   # 1. Download and store nginx logs to HDFS
   wget "$SOURCE/data/nginx-logs.zip"
   unzip nginx-logs.zip && rm nginx-logs.zip
   hdfs dfs -put nginx-logs /log-analysis/nginx-logs
11
12
   # 2. Analyze log files
   for log in $ (hdfs dfs -ls -C $in); do
13
     name="$out/$(basename "$log".log)'
     # 3. Identify bot IPs by visit frequency
     hdfs dfs -cat "$log" | grep -E $bots | cut -d" " -f1 |
16
           sort | uniq -c | sort -rn >> "${name}.out"
17
     # Further analysis omitted for brevity...
18
   done
```

**Fig. 2:** Log analysis script (*Cf*.§3). The script downloads Nginx logs, stores them on a distributed filesystem, and analyzes them to extract traffic statistics—slightly modified from POSH [49] to highlight idiomatic shell challenges.

when a fragment mixes streaming filters with blocking operators. D3 Commands with non-idempotent side effects remain in a special fragment under user control; distributed fragments perform only pure data transformations or write to isolated files that can be atomically swapped in upon success. D4 Fragments are derived at compile time from the AST, so any new commands—spawned via loops, conditionals, or environment expansions—automatically become first-class fault-tolerance units without requiring a static representation. D5 FRACTAL recovers at the fragment level, not per-command. Command-level recovery would incur prohibitively high scheduling and bookkeeping overhead. Fragment-level recovery strikes a sweet spot: small enough to avoid re-doing large amounts of work, yet coarse enough to amortize the runtime instrumentation cost. D6 All faulttolerance logic is injected by the compiler. Users run unmodified POSIX shell scripts under FRACTAL, without needing to reexperss their script in constrained APIs.

While the focus is correct and efficient recovery, FRAC-TAL also aims to deliver near state-of-the-art performance in failure-free executions.

### 3 Example and Overview

Scripts that process large datasets usually need to interact with distributed file systems such as HDFS [52], NFS, or Alluxio [37], as their input data does not fit on a single computer. FRACTAL scales out the computation to facilitate data locality, data parallelism, and pipeline parallelism, while ensuring recoverability when a participating node fails.

**Example script and problem**: Fig. 2 presents an example shell script analyzing log files generated by Nginx, divided

into three parts: (1) setup ( $L_{7-10}$ ), downloading 150GB of log data and storing them on HDFS; (2) driver ( $L_{13-14}$ ,  $L_{19}$ ), iterating over the HDFS directory, piping log files to the analysis pipeline and appending results to a dynamically determined local file; and (3) analysis ( $L_{15-18}$ ), identifying known bot IPs by visit frequency.

A developer opting for distributed execution, either manual or more recently automated [47, 49], is left with only one option when a node—*i.e.*, part of Fig. 2—fails: to restart the entire computation. Unfortunately, such a restart impacts both *performance*, as a full rerun will waste over 3 hours, and *correctness*, as the script appends to a file and thus upon failure may result in partial outputs—worse even, potentially mixed with correct results from earlier failure-free executions.

**Challenges for fault tolerance**: The script in Fig. 2 illustrates why fault tolerance is so challenging: it invokes blackbox commands (e.g. uniq -c at  $L_{17}$ ) whose internal counters cannot be inspected or checkpointed (D1), passes data through chains of ephemeral UNIX pipes (e.g. from grep to sort to uniq) with no built-in barriers or offsets (D2), and relies on side-effectful operations (e.g. the append operator >> at  $L_{17}$ ) that risk duplication or partial writes upon retry (D3). Control-flow constructs such as for log in \$IN/\*.log spawn commands dynamically based on variable values, preventing any static view of the computation graph (D4). Reexecuting the entire script on 150GB of logs takes over 3 hours, so coarse-grained restart is prohibitively expensive (D5). Finally, these are often legacy or incrementally maintained scripts, so any fault-tolerance scheme must operate transparently on unmodified POSIX shell programs (D6).

**Fractal overview**: Fig. 3 presents an overview of FRAC-TAL. FRACTAL builds on a DFG representation of a POSIX shell script via PaSh-JIT [32]. FRACTAL's coordinator then leverages command annotations to partition subgraphs into one of three types with different fault-recovery semantics (Fig. 3 A1): (1) main, which contains the AST region previously deemed as "unsafe" for distribution and is executed on a node containing the authoritative shell state and broader environment—typically, the client node from which the computation is initiated, (2) regular, which does not include an aggregator vertex, and (3) merger, which includes an aggregator vertex, such as sort -m, responsible for merging the outputs of multiple upstream subgraphs, It then instruments every inter-subgraph edge with lightweight communicative primitives (Fig. 3 A2) that record delivered byte offsets and enforce exactly-once semantics over ad-hoc UNIX pipes. Next, a persubgraph, heuristic—based component (Fig. 3 A3) decides whether to persist each subgraph's outputs locally, recovery speed with failure-free overhead. Once prepared, the coordinator schedules subgraphs across executor nodes and relies on progress and health monitors to track execution progress and detect failures (Fig. 3 A4-6). On node failure, it identifies and re-executes only the minimal downstream subgraphs that

did not complete, ensuring both correctness and efficiency in recovery. Executors (Fig. 3 B1-4) receive their assigned subgraphs, reconstruct them into shell scripts, and run them in a tight, non-blocking event loop that maximizes CPU utilization without oversubscription.

**Results**: On a 30-node Cloudlab cluster (§7), FRACTAL executes Fig. 2's script in 220s (speedup:  $40\times$ ). Upon failure at 50% of the execution, FRACTAL executes only the necessary fragments—outputting correct results across all local and HDFS files in 330s (27.1×).

## 4 System Design

This section presents FRACTAL's fault recovery design. It then introduces FRACTAL's core components that drive execution and recovery.

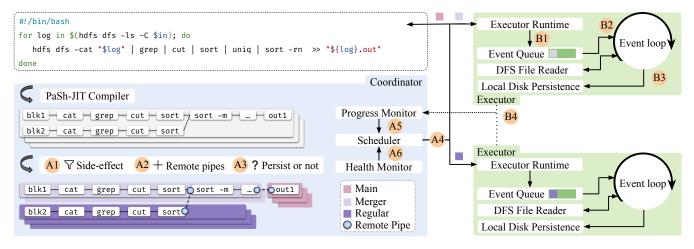
## 4.1 Fault Recovery in FRACTAL

When the health monitor alerts the coordinator about a node failure, rescheduling of the necessary subgraphs occurs in five steps where FRACTAL (1) identifies all incomplete subgraphs assigned to the crashed node and, by querying progress monitoring, their dependencies; (2) sends kill requests to subgraphs that cannot be used in the new execution plan; (3) updates the progress monitor according to the new execution plan; (4) identifies subgraphs that no longer need to be reexecuted because their results are persisted; (5) distributes the optimized list of subgraphs based on the new execution plan.

Some of these steps are different depending on whether the failed node is a merger or regular node. A faulty merger subgraph is rescheduled to a healthy executor, with incomplete upstream dependencies being re-routed to the same executor and complete dependencies having their persistent outputs transferred directly. A faulty regular subgraph is re-scheduled on a healthy executor, but its downstream merger is notified to continue reading the incomplete stream where the failed regular left off instead of re-executing the merger subgraph.

While the new execution plan is being prepared, the scheduler may receive new dataflow graphs to distribute. To avoid concurrent modifications to the progress monitor and further complications, crash handling and scheduling are performed under locks and are mutually exclusive.

When a loop is unrolled into parallel subgraphs, FRACTAL tracks read-write and write-write dependencies between iterations to establish execution order and isolation. If an iteration's corresponding subgraph fails, the scheduler applies the standard five-step crash-handling procedure only to that iteration and any upstream dependencies, while independent iterations are neither reissued nor re-executed. Completed iterations either reuse persisted outputs or are simply skipped, ensuring failures in one iteration do not force recomputation of its peers unless necessary.



**Fig. 3: FRACTAL's architecture.** From a client shell script, FRACTAL uses PaSh-JIT to build a DFG, applies annotations to isolate the unsafe main subgraph (A1), and splits the rest into regular and merger subgraphs at HDFS block boundaries. It then instruments each edge with remote-pipe primitives (A2) and uses a lightweight heuristic to persist outputs per subgraph (A3). The coordinator schedules subgraphs (A4) and leverages the progress (A5) and health (A6) monitors to re-execute only failed fragments. Executors reconstruct subgraphs into shell scripts and run them in a tight, non-blocking event loop (B1-4), streaming data via remote pipes, the distributed file reader, and local cache.

# **4.2** FRACTAL Components

**DFG augmentation**: Before scheduling, FRACTAL augments each DFG fragment with remote pipes to track execution progress and ensure exactly-once semantics during fault recovery. For instance, in Fig. 3 (A4), remote pipes are added at the boundary between the merger and regular fragments as well as between the merger and main fragment. The scheduler then assigns each fragment to executor nodes, replaces the original DFG edges with these remote pipes, and updates the progress-monitor metadata.

Remote Pipe: Efficient communication is crucial for precise recovery, since lost or duplicated streams can break correctness or incur extra work. remote pipes provide unidirectional channels between a writer (source) and a reader (destination), both identified by the same edge ID, in either transient (sockets) or persistent (files) mode under a dynamic persistence switch. If persistence is disabled, the writer opens a socket and registers its endpoint for the reader to resolve; if enabled, the writer writes to a file and exposes its path for retrieving the data during potential re-executions.

Detecting and handling failures is crucial for the remote pipe. If a connection is lost, the reader periodically queries the discovery service. When a new address is found—often due to rescheduling—a new connection is made. Since the reader knows how many bytes it has already forwarded downstream, it can discard duplicates and maintain a correct, non-repetitive data stream. This behavior is important for certain side-effectful operations (like the append operator shown in Fig. §2) to ensure that re-executed subgraphs do not produce duplicated outputs. The reader consumes the stream in buffered chunks while maintaining an 8-byte lookahead

for the EOF token. Once the unique marker is detected, it is stripped off and the complete data is delivered downstream.

**Dynamic Output Persistence**: Upon node failure, any incomplete subgraph and its upstream dependencies must be re-executed, which can be costly if those upstream tasks are expensive. To reduce this overhead, FRACTAL can persist the outputs of upstream subgraphs so that reassigned tasks can read cached results instead of recomputing them. However, writing to local storage incurs overhead during fault-free execution, and its impact varies with node hardware (e.g., HDD vs. NVMe). To strike a balance, FRACTAL employs a heuristic-driven *dynamic persistence* policy that makes persubgraph decisions based on static cluster profiling (automatically collected by **frac**) and runtime workload characteristics (*e.g.*, commands in subgraphs, their inputs).

Moreover, subgraphs created at HDFS block boundaries that have no downstream dependents, called *singular* subgraphs, never persist outputs because their results will not be reused on recovery. When transformed into DFGs, a shell script may produce both <u>singular</u> subgraphs and non-singular ones. Therefore, FRACTAL's dynamic persistence makes persubgraph decisions: it disables output persistence for *singular* subgraphs and selectively enables it for others based on cost heuristics. This fine-grained policy minimizes unnecessary I/O during normal execution while still caching results to accelerate recovery of expensive upstream tasks.

**Progress Monitor:** The progress monitor maintains all metadata needed for fault recovery: subgraph-to-node assignments, completion events, and inter-subgraph dependencies. Upon completing a send or receive operation, each subgraph emits a 17-byte completion event to the progress monitor, containing its serialized edge ID and a one-byte flag indicating

whether it sent or received. The scheduler then uses these compact messages to determine exactly which subgraphs must be re-executed after a failure. When persistence is enabled, the monitor also tracks file locations so that already-persisted fragments are not re-executed after a failure.

After each subgraph finishes its execution, whether it receives or sends data, it sends a special message to the progress monitor signaling its completion. This metadata is eventually used by the scheduler to decide what needs to be reexecuted. These messages are intentionally small—only 17 bytes—consisting of a serialized ID identifying the communication and a single byte indicating whether the receiver or sender has finished.

There is a specialized component of the progress monitor primarily responsible for discovery between subgraphs that need to communicate, as it is neither possible nor efficient to know every possible endpoint statically. The discovery service acts as a barrier between the sources and targets of communication during runtime, blocking one while it waits for the other until they discover each other and initiate data transfer. Other important responsibilities of the discovery service include tracking the persisted file locations when dynamic switching is enabled. This ensures that already-persisted redundant subgraphs are not re-executed in case of failures. Additionally, the discovery service supports two different kinds of protocols for singular and non-singular subgraphs, as they use different means to communicate.

**Health Monitor**: The health monitor polls the HDFS namenode's JMX endpoint for each data node's lastContact heartbeat timestamp. Nodes whose lastContact exceeds a configurable threshold are flagged offline. This threshold involves a balance between false positives and false negatives. A small threshold may result in frequent false positives, where temporary network slowdowns are mistaken for faults, triggering costly fault recovery mechanisms. Conversely, a high threshold could cause the system to wait unnecessarily for outputs from faulty nodes. Therefore, the threshold is designed to be configurable. The default value of 10 seconds is selected arbitrarily to ensure it does not dominate the execution time during evaluation (§7), while allowing a substantial portion of the execution to complete before initiating fault recovery mechanisms. This liveness information drives the scheduler's fault recovery decisions, triggering the reassignment or re-execution of affected subgraphs on healthy nodes.

Relying on HDFS heartbeats is an intentional choice to avoid scenarios where FRACTAL nodes appear to be available but HDFS nodes are not, or vice versa. This creates an additional dependency between HDFS and FRACTAL, but since any distributed file system must include a heartbeat mechanism, it should be possible to use these heartbeats as an indication of liveliness for FRACTAL nodes.

Executor Runtime: The executor runtime receives serial-

ized subgraphs from the coordinator and deserializes them into shell scripts. These scripts are then staged in a temporary directory whose path, along with some metadata, and enqueued as execution events.

Every 0.1s, the executor runtime performs three actions: (1) reclaims completed tasks—removing them from the active pool and recording timing and debug metadata; (2) applies pending kill requests from the coordinator by dropping targeted events from the queue; and (3) launches queued subgraphs up to the configured concurrency limit by spawning new processes.

Additionally, the executor runtime also manages environment setup and teardown (e.g., terminating remnants of rescheduled subgraphs to avoid duplicated executions), collects timing and diagnostic metadata, and enables controlled fault injection during evaluation.

Command Annotations: Previous systems [32,47,49,62] uses command annotations to identify opportunities in shell parallelization. FRACTAL uses annotations extended from PaSh-JIT and offers a flexible JSON interface that lets developers supply or override annotations for any third-party or black-box command. These annotations enable FRACTAL to distinguish safely re-executable regions (e.g., pure data transformations) from non-re-executable ones (e.g., side-effectful or non-deterministic operations), ensuring only subgraphs containing all safe commands are re-executed on failure. Regions cannot be safely re-executed are offloaded to be part of the main subgraph, which is executed on the client node.

# 5 Optimizations

This section presents targeted optimizations to FRACTAL's critical-path components, reducing control-plane overhead, and address implementation-specific challenges.

Event Driven Architecture: The executor runtime (§4.2)'s event loop is one of FRACTAL's most performance-sensitive components. To eliminate synchronization overhead, the executor runtime relies exclusively on atomic operations—integer assignments and list append/pop—instead of locks. Completion events are kept to 17 bytes each (edge ID plus direction flag) to minimize messaging overhead. The loop polls every 0.1s to balance kill-signal responsiveness against CPU utilization, and its concurrency level, the maximum number of subgraphs launched in parallel, is configurable per node, defaulting to the CPU core count to match hardware capacity.

**Buffered I/O**: To mark the end of an remote pipe stream, the writer appends an 8-byte EOF token. However, detecting and removing this sentinel on-the-fly is challenging because the reader cannot buffer the entire stream or perform full-stream scans. To address these challenges, the reader employs a buffered I/O strategy with several optimizations. It first allocates a configurable buffer, typically 4096 bytes, and

ensures that at least 8 bytes are initially read—this is always possible because the presence of the EOF token guarantees at least 8 bytes of data. After this initial setup, the reader enters a loop that (1) performs another read to fill the buffer following the initial 8-byte segment; (2) sends the buffer's contents, except for the final 8 bytes, downstream; (3) checks whether these last 8 bytes match the EOF token and, if they do, stops reading; and (4) moves the last 8 bytes to the start of the buffer, ready for the next iteration. This approach reduces overhead to at most an 8-byte copy for each iteration without generating unnecessary garbage, offering a significant improvement over simpler, linear parsing methods.

**Batched Scheduling**: If a script's input is relatively large or consists of many smaller files, FRACTAL may generate an excessive number of subgraphs to schedule, track, and execute. In such cases, distributing subgraphs can become more time-consuming than executing them. To address this issue, FRACTAL collects and batches all subgraphs with identical targets into a single request and sends these batches asynchronously to all cluster members. This kind of batching becomes increasingly important as the cluster size grows.

## 6 Fault Injection

To aid parameter selection and recovery characterization, FRACTAL's fault-injection subsystem, available as a command-line tool called **frac**, allows injecting runtime fail-stop and fail-restart faults in large-scale distributed deployments. The **frac** subsystem is agnostic to deployment and component internals, and has been used to inject faults to FRACTAL, DISH, and AHS across a variety of environments. Contrary to manual killing, **frac** offers automation, operates at byte-level and millisecond-level precision, can be driven by key events, allows automated restarts—and, by operating at the process-tree level, offers significant performance improvements over complete node shutdown, accelerating parameter selection and recovery characterization.

**Hard faults**: Manually shutting down compute nodes running the executor processes, termed *hard fault*, ensures they end up offline by issuing commands to the host environment. Unfortunately, hard faults are hard to automate at large-scale experiments. Existing VM shutdown tools are not ideal because the aforementioned experiments require non-graceful shutdowns, and verifying that nodes have truly come back up requires custom Docker-level health checks (e.g., polling until each block reaches its target replication factor).

Hard faults additionally do not support fine-grained control over the timing of fault events, crucial for precisely characterizing a systems' recovery behavior. Precision is particularly important for systems that prioritize minimizing runtime over load balance. Contrary to these systems, ones that prioritize balancing load over minimizing latency make roughly the same progress across all participating nodes—for exam-

ple, AHS' mapper and reducer executions are load-balanced across the cluster and thus a fault injection will hit a nodes at roughly the same execution point as other nodes. But other systems such as FRACTAL see imbalanced progression across regular and merger nodes, thus requiring and benefiting from improved precising in fault injection.

**Soft faults**: The **frac** tool supports two modes of *soft faults*. A data-plane mode injects into the data stream a special faultsequence token that uniquely matches a wrapper in one of the nodes. The token flows through the entire DFG, propagated downstream by command wrappers, and is duplicated by DFG splitters—i.e., commands that split the input data to identical subgraphs that implement parallel execution. Most wrappers propagate the token upon receipt, i.e., they do not feed it into the command they wrap but propagate it to the output stream—except for the one wrapping the ingress edge of the DFG subgraph targeted by the fault token. Upon receipt, the target wrapper kills all processes in the subgraph. Dataplane soft faults offer fine-grained byte-level precision for determining the exact point at which to inject a node fault, when the precise fault conditions hinge on specific elements of the data stream.

A control-plane mode sends a special token directly to the node responsible either at a specific time point or by the trigger of a specific event. Additional automation collects baseline execution times about the various jobs on each node. In a subsequent run, the coordinator injects the fault at a configurable time or percentage of a node's fault-free execution time. Focusing on the completion percentage of individual nodes is important in cases where the execution is not balanced—for example, 50% end-to-end job completion does not translate to 50% of AHS's map or FRACTAL's regular node execution, as these nodes typically consume a minority of the runtime. Once it receives a message from the coordinator, the fault terminates its HDFS datanode process and kills the corresponding process. Control-plane soft faults offer coarser-grained precision for determining the poitn at which to inject a fault, often incorporating higher-level goals such as completion percentages.

### 7 Evaluation

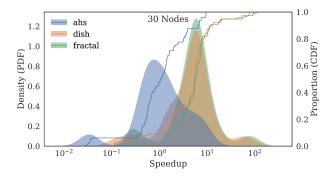
This section characterizes FRACTAL's performance under fail-free and fault-induced execution.

**Baselines**: We compare FRACTAL with (1) Bash [59], a standard single-node shell interpreter; (2) Apache Hadoop Streaming AHS [25], a production-grade fault-tolerant distributed-computing system that supports black-box Unix commands; and (3) DISH [47], a state-of-the-art fault-intolerant shell-script distribution system.

**Benchmarks**: We used five sets of real-world benchmarks (Tab. 2), totaling 77 scripts and 547 lines of shell code (LoC) excluding empty lines and comments. The Clas-

**Tab. 2: Benchmark summary.** Summary of all the benchmarks used to evaluate FRACTAL and their characteristics.

Benchmark	Scripts	LoC	AHS	Input
Classics	10	103	/	3 GB
Unix50	34	34	1	10 GB
NLP	22	280	×	10 GB
Analytics	5	62	/	33.4 GB
Automation	6	68	×	2.1-30 GB



**Fig. 4:** Fault-free performance summary. Summary of AHS, DISH, and FRACTAL 30-node fault-free speedups over Bash across all benchmarks: FRACTAL is comparable to DISH and significantly faster than AHS.

sics [1,2,31,42,58] and Unix50 [3,36] benchmarks comprise scripts that extensively invoke UNIX and Linux built-in commands. The NLP [6] benchmarks features scripts from a tutorial focused on developing natural language processing programs using UNIX and Linux utilities. The Analytics [61,64] benchmarks features data-processing scripts, including actual telemetry data from mass-transit schedules during a large metropolitan area's COVID-19 response and multi-year temperature data across the US. The Automation [49, 51, 53] benchmarks features scripts for processing, transforming, and compressing video and audio files, typical system administration and network traffic analyses over log files, and aliases for encrypting and compressing files.

Hardware & software setup: Experiments were conducted on two clusters: (1) 30 x Cloudlab m510 nodes, each with 8-core Intel Xeon D-1548 CPU at 2.0 GHz, 64GB RAM, 256 GB NVMe, and 10-Gb connection; (2) 4 x on-premise Raspberry Pi-5 nodes, each with a 4-core Arm Cortex A76 CPU at 2.4 GHz, 8GB RAM, 1TB SSD, and 1-Gb connection.

To improve reproducibility and ease deployment, we use Ubuntu 22.04-based Docker Swarm images on both 4 and 30 node clusters. We used Bash v5.1.16, Apache Hadoop v3.4.0, Python v3.10.12, and Go v1.22.2.

FRACTAL is developed on top of the PaSh-JIT compiler [32], which includes a Python re-implementation of libdash [22] a POSIX-compliant shell-script parser. FRAC-

**Tab. 3: Fault-free performance comparison highlights.** Average, minimum, and maximum speedups over Bash for FRACTAL, DISH, and AHS across all benchmarks.

		4 Node		30 Node		
System	Avg	Min	Max	Avg	Min	Max
FRACTAL	5.93	0.28	18.55	9.64	0.22	107.8
DISH	5.88	0.15	19.04	8.20	0.10	78.35
AHS	1.27	0.01	6.94	1.99	0.02	9.48

TAL adds 2K lines of Python (scheduler, monitors, executor runtime), 1.1K lines of Go (remote pipe and services), and 0.73K lines of shell script. An additional 389 lines of Python and 4.1K lines of shell scripts comprise the frac tool.

**Correctness:** Apart from careful engineering and many unit tests, the results of over 100 repetitions across several dozen distributed deployments and fault scenarios, over 70 scripts, and over 200GB of inputs are identical to those of the sequential script execution, offering significant confidence about FRACTAL's correct execution and recovery.

#### 7.1 Fault-Free Execution

This section characterizes the speedup of FRACTAL over Bash against DISH and AHS (Fig. 4).

**Experiments**: We execute all benchmarks on Bash, DISH, AHS, and FRACTAL on both clusters without injecting any faults. While Bash, DISH, and FRACTAL execute all shell scripts without modifications, AHS requires modifications. Not all shell scripts are expressible in AHS; those that are (Tab. 2, col. AHS) are used to compare AHS with FRACTAL.

**Results**: Fig. 5 shows the speedup of the three systems over Bash on the two clusters (key comparisons in Tab. 3). On the 30-node cluster, FRACTAL achieves an average speedup of  $9.64\times$  (max:  $107.84\times$ ) compared to  $8.2\times$  (max:  $78.35\times$ ) for DISH and  $1.99\times$  (max:  $9.48\times$ ) for AHS. On the 4-node cluster, FRACTAL achieves an average speedup of  $5.93\times$  (max:  $18.55\times$ ), compared to DISH's  $5.88\times$  (max:  $19.04\times$ ) and AHS's  $1.27\times$  (max:  $6.94\times$ ).

Excluding the Unix50 and NLP benchmarks, which are not well-suited for scaling across large clusters, FRACTAL achieves an average speedup of  $4.66\times$  on a 4-node cluster and  $21.90\times$  on a 30-node cluster.

Section 3's log-analysis script (Fig. 2), part of Automation, processes 30GB in 2140s on Pi-5 and 1524s on m510. FRACTAL brings it to 436s  $(4.90\times)$  on the 4-node Pi-5 cluster and 484s  $(3.15\times)$  on the 30-node Cloudlab cluster.

**Discussion:** FRACTAL is almost always faster than Bash, but the exact speedup achieved depends largely on the parallelization characteristics of each script. Scripts whose regular subgraphs consist of filters (e.g., grep) or folds (e.g., wc) perform better, as they reduce the runtime fraction used for I/O

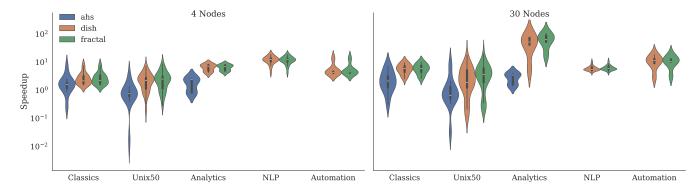
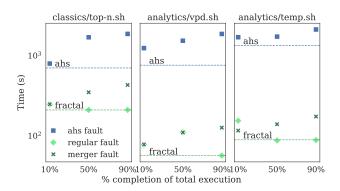


Fig. 5: Fault-free performance comparison (*Cf.*§7.1). Comparison between fault-free execution speedups of AHS, DISH, and FRACTAL, relative to single-node Bash, on the 4-node Pi-5 cluster (left) and the 30-node Cloudlab cluster (right).



**Fig. 6: Recovery comparison** (*Cf.*§7.2). Comparison between the FRACTAL and AHS recovery times for 3 representative scripts (left, mid, right), with faults introduced at 10%, 50%, and 90% of the execution—and without faults (dashed lines).

or the sequential merger. Conversely, scripts that do not filter as much or spend more time merging results experience lower speedups. In the limit, short-running scripts such as Unix50's 4.sh and 34.sh experience slowdowns, as their runtime is dominated by near-instant heads—but still remain within 1s.

FRACTAL performs better than AHS due to its whole-program optimizations, exploiting more opportunities for parallelism: AHS programs often contain multiple map and reduce stages, thus leaving pipeline parallelism, data parallelism, and task parallelism due to DLOpt unexploited.

FRACTAL at times performs better than DISH, as the benefits from its optimizations offset the (insignificant, in fault-free execution) costs to support fault tolerance. FRACTAL's asynchronous batching results in significant benefits as the deployment grows, which explains the more pronounced differences between FRACTAL and DISH on the 30-node cluster. FRACTAL's persistent subgraph outputs, aimed at avoiding re-computation under faults, result in additional benefits by allowing downstream components to access them: unlike DISH's TCP communication, which requires incremental—

and often blocking—generation due to buffer size constraints, FRACTAL allow persistence-enabled subgraphs to operate with effectively unbounded buffers. This in turn allows FRACTAL to pre-compute and store larger execution units, offering a significant advantage in scenarios with complex, interdependent subgraphs where DISH's buffer constraints lead to bottlenecks. And FRACTAL fault-tolerance overheads in fault-free execution paths are minimal—*e.g.*, each executor in Classics adds 136B over the network and writes 1MB to disk, imperceptible overheads relative to DISH's fault-intolerant execution.

### 7.2 Performance of Fault Recovery

We characterize FRACTAL's fault recovery with one experiment comparing FRACTAL with AHS failing at various stages and another characterizing FRACTAL under more scenarios.

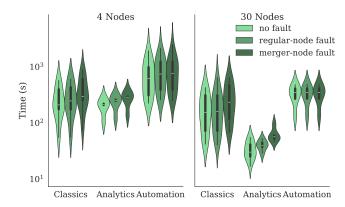
**Experiments**: The first experiment assesses the time required for FRACTAL and AHS to recover and successfully complete the job on the 4-node deployment. We introduce faults at approximately 10%, 50%, and 90% of the baseline execution time: at 10%, AHS executes mappers and FRACTAL executes regular subgraphs; at 90%, AHS executes reducers and FRACTAL executes merger subgraphs. Faults are introduced to an arbitrary AHS and both a base-case regular node and a worst-case merger node (in separate runs). Combining all these configurations with hard and soft faults to confirm **frac** across systems results in prohibitive manual effort, thus for this experiment we focus on three of the collected benchmarks.<sup>1</sup>

The second experiment zooms into FRACTAL's fault recovery under different failure scenarios on both clusters. It employs soft faults using **frac**, and all benchmarks except NLP and Unix50 as they contain many short-running scripts.

<sup>&</sup>lt;sup>1</sup>Total: 3 completion percents  $\times$  3 system configs (AHS, regular, merger)  $\times$  2 failure modes  $\times$  5 repetitions  $\times$  3 benchmarks = 270 experiments (about a week of manual effort) instead of 6,930 experiments.

**Tab. 4:** FRACTAL's speedup over AHS for different failure conditions and recovery scenarios. Format: avg (min-max).

	Regular Recovery	Merger Recovery
Fail at 10%	$7.8 \times (3.2 - 16.0 \times)$	8.5× (3.2–15.9×)
Fail at 50%	$12.1 \times (8.0 - 19.7 \times)$	$8.3 \times (4.8 - 13.9 \times)$
Fail at 90%	$16.4 \times (8.9 - 32.9 \times)$	$8.0 \times (4.3 - 14.3 \times)$



**Fig. 7: Recovery comparison (soft faults) (***Cf.***§7.2).** FRACTAL execution times for three benchmark sets (Classics, Analytics, Automation) with no faults, regular faults, and merger faults on a 4-node Pi-5 cluster (left) and a 30-node Cloudlab cluster (right).

**Results**: Fig. 6 summarizes the first experiment. The xaxis shows different completion percentages and the y-axis shows the time it takes AHS and FRACTAL (both regular and merger nodes) to recover. For context, dashed lines (constant across the x-axis) represent the fault-free executions for AHS (avg: 937.6s) and FRACTAL (avg: 118.9s). Under regular faults, it takes FRACTAL 160s (134.5% vs. fault-free), 136.5s (114.8%), and 118.8s (100.1%) to recover for each execution percentile; under merger faults, these become 147.7s (124.2% vs. fault-free), 200.2s (168.3%), and 244.4s (205.6%) respectively. For the same percentiles, it takes AHS 1,248.8s (133.2% vs. AHS fault-free, 780.9% vs. regular, 845.5% vs. merger), 1,655.8s (176.6% vs. AHS fault-free, 1,213.2% vs. regular, 727.1% vs. merger), and 1,953.4s (208.3%, 1,644.3%, 799.1%). Tab. 4 summarizes the comparison between AHS and FRACTAL.

Fig. 7 summarizes the second experiment, showing execution times for fault-free, regular recovery, and merger recovery. Benchmarks with fewer parallel pipelines such as Classics and Analytics take 20.3–32.1% longer to recover from merger (209.4–344.8s) than regular node faults (150.4–277.6s). For other benchmarks, there is not a significant difference between the recovery of different nodes (335.3–845.0s for regular vs. 335.7–856.5s for merger).

**Discussion**: Overall, the first experiment shows that FRACTAL recovers at a fraction (6.08–12.8%) of AHS's recovery time. Most benefits are due to its selective re-execution of

affected-only portions, while still enjoying all parallelization benefits (§7.1) during re-execution.

Failures that occur later generally result in longer recovery times (Fig. 6), as they require more upstream subgraphs to be re-executed—except for FRACTAL's regular failures, whose recovery does not interfere with end-to-end performance due to their completion. This non-interference of regular nodes will be important in practical deployments, as we have found that the vast majority of nodes in a distributed execution are regular nodes—thus have significantly higher probability to be affected by a fault in the underlying infrastructure.

As this experiment compares *hard* faults introduced manually and *soft* faults injected by **frac**, it confirms that the two modes result in identical executions across 270 experiments—but **frac** completes experiments at about 2–5% of the hardfault time, and without the mental overhead of keeping manual track of various experiment timepoints.

Diving into various types of recovery (Fig. 7) indicates that the pipeline-to-node ratio of pipelines is correlated with merger -to- regular node recovery performance. This observation is intuitive for two reasons. First, benchmarks that rely on a single pipeline—such as Classics and Analytics—experience longer recovery times from merger faults than from regular faults, since regular faults involve re-executing fewer subgraphs. Second, benchmarks with many pipelines (e.g., Automation) are indifferent to merger and regular recovery times: having a larger number of merger subgraphs distributes the workload evenly, effectively making every node a merger node and neutralizing the impact type.

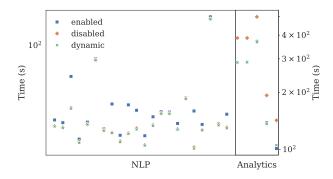
# 7.3 Microbenchmark: Dynamic Persistence

This experiment zooms into the benefits of dynamically persisting subgraph outputs, which accelerate fault recovery but add overhead to fault-free execution.

**Experiments**: We explore output-persistence trade-offs using two benchmark sets under different fault conditions and with both persistence options (enabled and disabled): NLP (no faults) features many parallel pipelines, each using a small input (<128MB); Analytics (merger faults) features long-running regular (upstream) subgraphs.

**Results**: Fig. 8 compares the runtime between three outputpersistence options. For fault-free NLP (resp. faulty Analytics) suite, enabling (resp. disabling) persistence results to 21.0% (resp. 38.7%) overhead on average.

**Discussion**: Enabling output persistence for short-running workloads introduces significant overheads and for relatively small benefits—as regular and merger subgraphs tend to be co-located, thus becoming unavailable upon faults and necessitating re-execution irrespective of persistence. Failed long-running scripts see significant benefits from output persistence, as they avoid significant upstream re-execution. FRACTAL first-order workload heuristics—*e.g.*, size of DFG graphs,



**Fig. 8: Microbenchmark: dynamic persistence** (*Cf.*§7.3). Faultfree NLP benchmark (left) and fault-injected Analytics benchmark (right) with dynamic persistence enabled, disabled, and set dynamically by FRACTAL's heuristics.

input sizes—decide whether to enable persistence for the vast majority of these benchmarks.

#### 8 Related Work

FRACTAL is related to a large body of work in distributed shells and shell-related utilities, distributed computing frameworks, and language-based distributed systems.

**Distributed shells and utilities**: Several command-line job-scheduling tools allow distributing workloads on Unix systems—*e.g.*, qsub for the Sun Grid Engine [19] and parallel for GNU Parallel [56]—but their invocation requires careful manual orchestration and does not come with fault tolerance built-in. Slurm [65], a workload manager for distributing batch jobs across computing clusters, provides periodic check-pointing for later resumption using DMTCP [33]. This mechanism focuses only on recovering Slurm-visible state, does not support complex commands, and fails to account for thorny shell semantics such as append (§3).

Shells like Rc [12], Dgsh [54], and gsh [41] offer scalable, often non-linear and acyclic, extensions to Unix pipelines but require manual rewriting and do not tolerate failures.

Recent systems offering automated parallelization and distribution of shell programs such as PaSh [32,62], POSH [49], and DISH [47] optimize the execution of shell programs by offloading and scaling out distribution automatically. Similar to FRACTAL, they employ a series of techniques for automatically compiling scripts to an internal representation which they then parallelize and distribute, but different from FRACTAL tolerate no failures.

**Distributed computing frameworks**: FRACTAL combines elements from distributed batching and streaming systems [10, 43, 45, 46, 48, 55, 63, 66]. These systems offer the ability to tolerate failures, often by tracking lineage similarly to FRACTAL [43, 66], but require their users to (re-)implement

their programs using the abstractions these systems provide. These systems do not support the black-box nature, runtime expansion behaviors, and arbitrary side effects pervasive in the commands typically present in shell programs.

Hadoop Streaming [25] and Dryad Nebula [28] are unusual in that they allow the use of black-box components such as Unix commands. However, they do not target the semantics of the shell and thus require users to manually port their shell programs—often facing the difficulty or inability to express entire classes of shell programs as FRACTAL's evaluation confirms. FRACTAL provides automated scale-out of unmodified shell scripts, supports the shell's dynamic-expansion features, and offers efficient fault tolerance by tracking lineage.

**Other cloud offerings**: Prior work on VM- and container-level replication has used check-pointing [8,9,13,34,40] to tolerate failures. Contrary to FRACTAL, these approaches leverage logging, require infrastructure support, and lack a logical understanding of the workload.

Serverless platforms have started introducing stateful operations [30,39,67] and thus fault tolerance, through a combination of logging and check-pointing. Different from FRACTAL, these systems do not support shell scripts or arbitrary black box commands—users need to (re)write their scripts in the abstractions provided by these systems to see benefits.

The gg system [16] supports scaling black box commands to serverless functions. In contrast to FRACTAL, gg does not support pipeline parallelism and attempts full executions via a re-try mechanism when faults occur.

#### 9 Conclusion

Transparent fault tolerance is a *sine qua non* for scalable shell-script distribution: without it, unmodified POSIX scripts cannot reliably handle the black-box binaries, ad-hoc pipes, non-idempotent side effects, and dynamic control flow that characterize real-world workflows.

FRACTAL is the first system that offers fault-tolerant shell-script distribution by separating recoverable from side-effectful regions. It performs lightweight instrumentation to record byte-level progress and enforce exactly-once semantics. By employing precise dependency and progress tracking at the subgraph level, it offer sound and efficient fault recovery.

Acknowledgements: We thank the NSDI'26 reviewers for their feedback; our shepherd, Eric Eide, for his guidance; the NSDI'26 Artifact reviewers for their time; and the Brown CS2952R (Fall'24) participants for their input on several iterations of this paper. This material is based upon research supported by NSF awards CNS-2247687 and CNS-2312346; DARPA contract no. HR001124C0486; a Fall'24 Amazon Research Award; a Google ML-and-Systems Junior Faculty award; a seed grant from Brown University's Data Science Institute; and a BrownCS Faculty Innovation Award.

### References

- [1] Jon Bentley. Programming pearls: a spelling checker. *Commun. ACM*, 28(5):456–462, may 1985.
- [2] Jon Bentley, Don Knuth, and Doug McIlroy. Programming pearls: a literate program. *Commun. ACM*, 29(6):471–483, jun 1986.
- [3] Pawan Bhandari. Solutions to unixgame.io, 2020. Accessed: 2020-04-14.
- [4] Paris Carbone, Stephan Ewen, Gyula Fóra, Seif Haridi, Stefan Richter, and Kostas Tzoumas. State management in apache flink®: consistent stateful distributed stream processing. *Proc. VLDB Endow.*, 10(12):1718–1729, aug 2017.
- [5] Raul Castro Fernandez, Matteo Migliavacca, Evangelia Kalyvianaki, and Peter Pietzuch. Integrating scale out and fault tolerance in stream processing using operator state management. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, SIGMOD '13, pages 725–736, New York, NY, USA, 2013. Association for Computing Machinery.
- [6] Kenneth Ward Church. Unix<sup>™</sup> for poets, 1994. Notes of a course from the European Summer School on Language and Speech Communication, Corpus Based Methods.
- [7] CRIU community. Checkpoint/restart in userspace (criu). https://criu.org/, 2019. Accessed: April 2025.
- [8] Heming Cui, Rui Gu, Cheng Liu, Tianyu Chen, and Junfeng Yang. Paxos made transparent. In *Proceedings* of the 25th Symposium on Operating Systems Principles, pages 105–120, 2015.
- [9] Brendan Cully, Geoffrey Lefebvre, Dutch Meyer, Mike Feeley, Norm Hutchinson, and Andrew Warfield. Remus: High availability via asynchronous virtual machine replication. In *Proceedings of the 5th USENIX symposium on networked systems design and implementation*, pages 161–174. San Francisco, 2008.
- [10] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, 2008.
- [11] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: a flexible data processing tool. *Commun. ACM*, 53(1):72–77, jan 2010.
- [12] Tom Duff. Rc—a shell for plan 9 and unix systems. *AUUGN*, 12(1):75, 1990.

- [13] George W Dunlap, Dominic G Lucchetti, Michael A Fetterman, and Peter M Chen. Execution replay of multiprocessor virtual machines. In *Proceedings of the* fourth ACM SIGPLAN/SIGOPS international conference on Virtual execution environments, pages 121–130, 2008.
- [14] Johan Eveleens and Chris Verhoef. The rise and fall of the chaos report figures. *IEEE software*, 27(1):30–36, 2009.
- [15] Bent Flyvbjerg and Alexander Budzier. Why your it project might be riskier than you think. *arXiv* preprint *arXiv*:1304.0265, 2013.
- [16] Sadjad Fouladi, Francisco Romero, Dan Iter, Qian Li, Shuvo Chatterjee, Christos Kozyrakis, Matei Zaharia, and Keith Winstein. From laptop to lambda: Outsourcing everyday jobs to thousands of transient functional containers. In 2019 USENIX annual technical conference (USENIX ATC 19), pages 475–488, 2019.
- [17] Aeleen Frisch. Essential system administration: Tools and techniques for linux and unix administration. "O'Reilly Media, Inc.", 2002.
- [18] Ishaan Gandhi and Anshula Gandhi. Lightening the cognitive load of shell programming. PLATEAU 2020, 2020.
- [19] Wolfgang Gentzsch. Sun grid engine: Towards creating a compute power grid. In *Proceedings First IEEE/ACM International Symposium on Cluster Computing and the Grid*, pages 35–36. IEEE, 2001.
- [20] GitHub. The state of the octoverse 2024: The most popular programming languages, 2024. Accessed: 2024-10-31.
- [21] Michael Greenberg. Word expansion supports posix shell interactivity. In *Companion Proceedings of the 2nd International Conference on the Art, Science, and Engineering of Programming*, pages 153–160, 2018.
- [22] Michael Greenberg. libdash. https://github.com/m gree/libdash, 2019. [Online; accessed November 22, 2024].
- [23] Michael Greenberg and Austin J Blatt. Executable formal semantics for the posix shell. *Proceedings of the ACM on Programming Languages*, 4(POPL):1–30, 2019.
- [24] Michael Greenberg, Konstantinos Kallas, and Nikos Vasilakis. Unix shell programming: the next 50 years. In *Proceedings of the Workshop on Hot Topics in Operating Systems*, pages 104–111, 2021.

- [25] Hadoop. Hadoop streaming. https://hadoop.apache.org/docs/r3.4.0/hadoop-streaming/HadoopStreaming.html, 2024. [Online; accessed June 13, 2024].
- [26] Saurav Haloi. *Apache zookeeper essentials*. Packt Publishing Ltd, 2015.
- [27] Jordan Henkel, Christian Bird, Shuvendu K Lahiri, and Thomas Reps. Learning from, understanding, and supporting devops artifacts for docker. In *Proceedings of the ACM/IEEE 42nd international conference on software engineering*, pages 38–49, 2020.
- [28] Michael Isard, Mihai Budiu, Yuan Yu, Andrew Birrell, and Dennis Fetterly. Dryad: distributed data-parallel programs from sequential building blocks. In Proceedings of the 2nd ACM SIGOPS/EuroSys European conference on computer systems 2007, pages 59–72, 2007.
- [29] Jeroen Janssens. Data science at the command line: Facing the future with time-tested tools. "O'Reilly Media, Inc.", 2014.
- [30] Zhipeng Jia and Emmett Witchel. Boki: Stateful serverless computing with shared logs. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*, pages 691–707, 2021.
- [31] Dan Jurafsky. Unix for poets, 2017. Accessed: 2024-09-16.
- [32] Konstantinos Kallas, Tammam Mustafa, Jan Bielak, Dimitris Karnikis, Thurston H.Y. Dang, Michael Greenberg, and Nikos Vasilakis. Practically correct, just-intime shell script parallelization. In 16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22), pages 1–18. USENIX Association, July 2022.
- [33] Gene Kapadia, Jason Ansel, Kapil Arya, Charles Guo, Daniel Maze, Mihir Modi, Cameron Musco, Alexey Lory, and Gene Cooperman. Dmtcp: Distributed multithreaded checkpointing. https://dmtcp.sourceforge.io/, 2024. Accessed: 2024-11-29.
- [34] Manos Kapritsos, Yang Wang, Vivien Quema, Allen Clement, Lorenzo Alvisi, and Mike Dahlin. All about eve:{Execute-Verify} replication for {Multi-Core} servers. In 10th USENIX Symposium on Operating Systems Design and Implementation (OSDI 12), pages 237–250, 2012.
- [35] Jay Kreps, Neha Narkhede, Jun Rao, et al. Kafka: A distributed messaging system for log processing. In *Proceedings of the NetDB*, volume 11, pages 1–7. Athens, Greece, 2011.

- [36] Nokia Bell Labs. The unix game—solve puzzles using unix pipes, 2019. Accessed: 2020-03-05.
- [37] Haoyuan Li. *Alluxio: A virtual distributed file system.* University of California, Berkeley, 2018.
- [38] Georgios Liargkovas, Konstantinos Kallas, Michael Greenberg, and Nikos Vasilakis. Executing shell scripts in the wrong order, correctly. In *Proceedings of the 19th Workshop on Hot Topics in Operating Systems*, pages 103–109, 2023.
- [39] David H Liu, Amit Levy, Shadi Noghabi, and Sebastian Burckhardt. Doing more with less: Orchestrating serverless applications without an orchestrator. In 20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23), pages 1505–1519, 2023.
- [40] Haikun Liu, Hai Jin, Xiaofei Liao, Liting Hu, and Chen Yu. Live migration of virtual machine based on full system trace and replay. In *Proceedings of the 18th ACM international symposium on High performance distributed computing*, pages 101–110, 2009.
- [41] Chris McDonald and Trevor I Dix. Support for graphs of processes in a command interpreter. *Software: Practice and Experience*, 18(10):1011–1016, 1988.
- [42] Malcolm D. McIlroy, Elliot N. Pinson, and Berkley A. Tague. Unix time-sharing system: Foreword. *Bell System Technical Journal*, 57(6):1899–1904, 1978.
- [43] Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard Liaw, Eric Liang, Melih Elibol, Zongheng Yang, William Paul, Michael I Jordan, et al. Ray: A distributed framework for emerging {AI} applications. In 13th USENIX symposium on operating systems design and implementation (OSDI 18), pages 561–577, 2018.
- [44] Derek G. Murray, Frank McSherry, Rebecca Isaacs, Michael Isard, Paul Barham, and Martín Abadi. Naiad: a timely dataflow system. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, SOSP '13, pages 439–455, New York, NY, USA, 2013. Association for Computing Machinery.
- [45] Derek G Murray, Frank McSherry, Rebecca Isaacs, Michael Isard, Paul Barham, and Martín Abadi. Naiad: a timely dataflow system. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, pages 439–455, 2013.
- [46] Derek G Murray, Malte Schwarzkopf, Christopher Smowton, Steven Smith, Anil Madhavapeddy, and Steven Hand. {CIEL}: A universal execution engine for distributed {Data-Flow} computing. In 8th USENIX Symposium on Networked Systems Design and Implementation (NSDI 11), 2011.

- [47] Tammam Mustafa, Konstantinos Kallas, Pratyush Das, and Nikos Vasilakis. DiSh: Dynamic Shell-Script distribution. In 20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23), pages 341–356, Boston, MA, April 2023. USENIX Association
- [48] Russell Power and Jinyang Li. Piccolo: Building fast, distributed programs with partitioned tables. In 9th USENIX Symposium on Operating Systems Design and Implementation (OSDI 10), 2010.
- [49] Deepti Raghavan, Sadjad Fouladi, Philip Levis, and Matei Zaharia. POSH: A data-aware shell. In 2020 USENIX Annual Technical Conference (USENIX ATC 20), pages 617–631, 2020.
- [50] Arnold Robbins and Nelson HF Beebe. Classic Shell Scripting: Hidden Commands that Unlock the Power of Unix. "O'Reilly Media, Inc.", 2005.
- [51] Michael Schröder and Jürgen Cito. An empirical investigation of command-line customization. *Empirical Software Engineering*, 27(2), December 2021.
- [52] Konstantin Shvachko, Hairong Kuang, Sanjay Radia, and Robert Chansler. The hadoop distributed file system. In 2010 IEEE 26th Symposium on Mass Storage Systems and Technologies (MSST), pages 1–10, 2010.
- [53] Diomidis Spinellis and Marios Fragkoulis. Extending unix pipelines to dags. *IEEE Transactions on Computers*, 66(9):1547–1561, 2017.
- [54] Diomidis Spinellis and Marios Fragkoulis. Extending unix pipelines to dags. *IEEE Transactions on Computers*, 66(9):1547–1561, 2017.
- [55] Craig A Stewart, Timothy M Cockerill, Ian Foster, David Hancock, Nirav Merchant, Edwin Skidmore, Daniel Stanzione, James Taylor, Steven Tuecke, George Turner, et al. Jetstream: a self-provisioned, scalable science and engineering cloud environment. In Proceedings of the 2015 XSEDE Conference: Scientific Advancements Enabled by Enhanced Cyberinfrastructure, pages 1–8, 2015.
- [56] Ole Tange. Gnu parallel-the command-line power tool. *Usenix Mag*, 36(1):42, 2011.
- [57] Dave Taylor. Wicked Cool Shell Scripts: 101 Scripts for Linux, Mac OS X, and Unix Systems. No Starch Press, 2004.
- [58] Dave Taylor. Wicked Cool Shell Scripts: 101 Scripts for Linux, Mac OS X, and Unix Systems. No Starch Press, USA, 2004.

- [59] The Free Software Foundation. Bash shell, 2009. [Online; accessed 30-October-2024].
- [60] Ankit Toshniwal, Siddarth Taneja, Amit Shukla, Karthik Ramasamy, Jignesh M Patel, Sanjeev Kulkarni, Jason Jackson, Krishna Gade, Maosong Fu, Jake Donham, et al. Storm@ twitter. In Proceedings of the 2014 ACM SIGMOD international conference on Management of data, pages 147–156, 2014.
- [61] Eleftheria Tsaliki and Diomidis Spinellis. The real statistics of buses in athens, 2021.
- [62] Nikos Vasilakis, Konstantinos Kallas, Konstantinos Mamouras, Achilles Benetopoulos, and Lazar Cvetković. Pash: Light-touch data-parallel shell processing. In *Proceedings of the Sixteenth European Conference on Computer Systems*, pages 49–66, New York, NY, USA, 2021. Association for Computing Machinery.
- [63] Tom White. *Hadoop: The definitive guide*. "O'Reilly Media, Inc.", 2012.
- [64] Tom White. *Hadoop: The Definitive Guide*. O'Reilly Media, Inc., 4th edition, 2015.
- [65] Andy B Yoo, Morris A Jette, and Mark Grondona. Slurm: Simple linux utility for resource management. In *Workshop on job scheduling strategies for parallel processing*, pages 44–60. Springer, 2003.
- [66] Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, and Ion Stoica. Resilient distributed datasets: a fault-tolerant abstraction for inmemory cluster computing. In Proceedings of the 9th USENIX Conference on Networked Systems Design and Implementation, NSDI'12, page 2, USA, 2012. USENIX Association.
- [67] Haoran Zhang, Adney Cardoza, Peter Baile Chen, Se-bastian Angel, and Vincent Liu. Fault-tolerant and transactional stateful serverless workflows. In 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20), pages 1187–1204, 2020.