

# DiSh: Dynamic Shell-Script Distribution

Tammam Mustafa  
MIT

Konstantinos Kallas  
University of Pennsylvania

Pratyush Das  
Purdue University

Nikos Vasilakis  
Brown University

## Abstract

Shell scripting remains prevalent for automation and data-processing tasks, partly due to its dynamic features—*e.g.*, expansion, substitution—and language agnosticism—*i.e.*, the ability to combine third-party commands implemented in any programming language. Unfortunately, these characteristics hinder automated shell-script distribution, often necessary for dealing with large datasets that do not fit on a single computer. This paper introduces DiSH, a system that distributes the execution of dynamic shell scripts operating on distributed filesystems. DiSH is designed as a shim that applies program analyses and transformations to leverage distributed computing, while delegating all execution to the underlying shell available on each computing node. As a result, DiSH does not require modifications to shell scripts and maintains compatibility with existing shells and legacy functionality. We evaluate DiSH against several options available to users today: (i) Bash, a single-node shell-interpreter baseline, (ii) PASH, a state-of-the-art automated-parallelization system, and (iii) Hadoop Streaming, a MapReduce system that supports language-agnostic third-party components. Combined, our results demonstrate that DiSH offers significant performance gains, requires no developer effort, and handles arbitrary dynamic behaviors pervasive in real-world shell scripts.

## 1 Introduction

Unix and Linux shell scripting remains prevalent—8<sup>th</sup> most popular language on GitHub in 2022 [20]—for data processing, system orchestration, and other automation tasks. Part of this prevalence can be attributed to a unique combination of features: (1) powerful and language-agnostic primitives for composing components available in any programming language; (2) dynamic features such as command substitution, variable expansion, and state reflection on the file system; and (3) a wide range of useful components called commands, available in the broader environment and tailored to specific tasks. These features enable the composition of succinct and powerful programs on a single computer (§2).

**Tab. 1: Available options for scaling out shell programs.** Compatibility: support unmodified shell scripts. Granularity: support fine-grained distribution. Expressiveness: support arbitrary dynamic behaviors. Agnosticism: support components in any programming language. Equivalence: behavior equivalence with existing shells.

Approach	Compatibility	Granularity	Expressiveness	Agnosticism	Equivalence	Examples
Distributed Shells	☐	■	■	■	☐	[14, 18, 63]
POSH	■	■	☐	■	☐	[49]
Cluster Comp. Frameworks (CCF)	☐	■	☐	☐	☐	[44, 57, 68, 71]
Language-agnostic CCFs	☐	■	☐	■	■	[25, 30]
Job Scheduling Tools	■	☐	☐	■	■	[19, 29, 58, 69]
Other languages	☐	■	■	☐	☐	[16, 60, 66]
DiSH	■	■	■	■	■	

Unfortunately, these features also hinder automated shell-script scale-out to multiple computers. Such scale-out is often necessary not only to accelerate computations, but also to compute over data that either do not fit on a single computer or are naturally distributed across multiple computers.

**State of the art:** Shell users dealing with large datasets that do not fit on a single computer are left with only a few options (Tab. 1). One option is to use a distributed shell [14, 18, 63]. Distributed shells require rewriting scripts manually and only support a small subset of UNIX features—often with limited, if any, dynamic features and varying support for composition constructs. A recent distributed shell named POSH [49] can handle a subset of shell scripts without rewriting—although that subset is limited to dataflow-only computations and also does not include arbitrary dynamic shell behaviors. In addition, since POSH is a shell reimplement, it is not behaviorally equivalent with existing shells and thus risks breaking ported scripts. A second option is to rewrite (parts of) the script in a cluster-computing framework [11, 44, 68, 71].

These only support pure computations (*e.g.*, batch, stream), require manual rewriting, and only rarely [25, 30] support language-agnostic components. Another option is job scheduling tools [19, 29, 58, 69], but these operate at a coarse granularity and do not leverage parallelism available in individual commands. Yet another option is to rewrite scripts in languages that support distribution [3, 40, 42, 66], foregoing the shell’s succinctness and language agnosticism. To summarize, these options operate on a subset of the shell, require significant manual effort, risk breaking correctness, or—most often—suffer from a combination of these limitations (see §8 for more details).

**Dynamic shell-script distribution:** This paper presents DiSH, a system designed to scale out shell scripts operating on distributed filesystems while maintaining full POSIX compatibility. DiSH satisfies all requirements in Table 1: it operates on existing shell scripts; it distributes scripts at the granularity of individual commands; it handles arbitrary dynamic shell features such as substitution and expansion; it allows the use of commands and utilities of any language; and, most importantly, it is behaviorally equivalent to Bash.

DiSH first instruments the execution of a script to identify regions that may benefit from distribution. At runtime, it compiles these regions to an intermediate representation which it then optimizes to introduce appropriate parallelism, buffering, communication, and coordination. DiSH then executes each compiled region in a distributed fashion using the same shell interpreter, components, and data as the original script.

**Implementation and results:** DiSH is implemented as a shim layer (rather than a shell) that wraps and orchestrates the (completely unmodified) user shell, delegating all execution to the underlying shell available on each computing node. This design hides distribution from the user and avoids modifying the underlying shell interpreter: the user thinks that their original script is being executed (but faster); each underlying shell is given a part of the distributed script to execute. As a result, DiSH achieves a new milestone in automated shell-script distribution: it offers significant performance benefits, it avoids modifications to shell scripts, and it maintains full POSIX compatibility. Additionally, this modular design allows further research and improvements without modifications in the underlying shell.

We characterize DiSH’s performance on a 4-node on-premise cluster and a 20-node cloud deployment using 76 scripts—including ones not trivially expressible in modern distributed computing frameworks, such as scripts with `for` loops, side-effects, and complex third-party components. DiSH surpasses the speedups achieved by production-grade systems on existing benchmarks and extends speedups to new ones: it achieves significant speedups over (1) Bash (avg: 13.6×; max: 136.3×), a single-node shell-interpreter baseline; (2) PASH (avg: 8.9×; max: 108.8×), a shell-script parallelization system; and (3) Hadoop Streaming (avg: 7.2×;

max: 32.3×), a cluster computing framework that supports language-agnostic components and shell scripts. Moreover, whereas Hadoop Streaming does not support 27/76 scripts and requires rewriting 7/76 scripts, DiSH runs all scripts without any modifications; in fact, DiSH is able to execute the entire POSIX shell test suite, only diverging in one error code out of thousands of assertions.

**Paper outline and contributions:** The paper begins with an example and overview (§2) of DiSH’s use and techniques. Sections 3–6 present DiSH’s key components:

- Dynamic orchestration (§3): DiSH parses, pre-processes, expands, and orchestrates its input script to enable dynamic distribution at runtime.
- Compilation (§4): During script execution, DiSH compiles certain regions to an intermediate representation and applies a series of optimizations.
- Distribution (§5): DiSH distributes each region to a set of workers in a way that promotes co-location of processing primitives and the data blocks these operate on.
- Runtime support (§6): DiSH bundles additional runtime primitives supporting correct and efficient communication in the context of distributed shell script execution.

The paper then presents DiSH’s evaluation (§7) and related work (§8), before concluding (§9).

**DiSH limitations:** DiSH currently does not tolerate failures such as worker aborts or network partitions. In such occasions, users are expected to rerun their scripts similar to how they do in non-distributed executions: due to the shell’s dynamic features and its support for third-party components, users often re-run failing scripts from the start. The current DiSH prototype does not implement support for security features such as encryption and containment.

**Availability:** All the work described in this paper has been implemented and incorporated into PASH—an MIT-licensed project—and is available by the Linux Foundation at <https://github.com/binpash/dish>.

## 2 Background, Example, and Overview

DiSH allows everyday shell scripts to reap the benefits of distributed computing: execute on data that do not fit on a single machine, often also speeding up expensive computations.

**Intended use:** DiSH is designed to support a variety of use cases, depending on the details of the distributed environment on which the system is executing. The most common case is one where input data are downloaded and stored in a distributed file system such as HDFS<sup>1</sup> and then processed using

<sup>1</sup>The choice of HDFS is not binding. DiSH could work on top of any distributed file system (*e.g.*, NFS or Alluxio [35]) that exposes the locations of file blocks. To achieve performance benefits due to co-location, there also needs to be available compute on the nodes that host that file system.

various analyses. This is useful for datasets that do not fit on a single computer, that are naturally distributed across multiple computers, or that can be processed faster in a data-parallel fashion. DiSH will distribute the computation appropriately, often running data-parallel instances on multiple machines and multiple processors per machine. DiSH also supports hybrid operation where data resides on both distributed and local file systems; this is useful for computations that contain CPU-intensive stages over datasets that do not necessarily reside on distributed file systems.

**Example script:** Fig. 1 shows a shell script that calculates maximum and average temperatures across the US, on datasets hosted on the National Oceanic and Atmospheric Administration (NOAA). The script is split into three parts: (p. 1) an 11-stage pre-processing pipeline to download data from NOAA and store them on HDFS, with the data range controlled upon invocation via dynamic arguments `$1` and `$2`; (p. 2, 3) two 5-stage pipelines calculating and storing maximum and average temperatures to the local file system.

HDFS is a distributed file system for handling large data sets on commodity hardware. Scripts like the one in Fig. 1 that process files stored in distributed file systems spend most of their execution time moving files across the network. On a 4-node cluster (§7) and 3.6GB of input, running just `hdfs dfs -cat` takes 346s; computing pipeline 2 (maximum temperature) only adds 6s. This phenomenon is due to pipeline parallelism: the execution time of all concurrently executing commands is mostly shadowed by `hdfs dfs -cat`.

**Opportunities for scale-out:** There are ample opportunities for improving the performance of this script. Since all parts contain stages that operate on large datasets, we should be able to execute (at least some of) their stages in a data-parallel fashion. For example, we should parallelize commands that process their input independently, such as `cut` and `grep`, by having them operate in parallel over partial inputs.

Additionally, carefully colocating computation and data should also improve performance. For example, we should schedule the data-parallel execution of the aforementioned `cut` and `grep` instances on machines that store the respective data segments. Directly operating on distributed file segments, rather than gathering and processing data on a subset of the machines, eliminates most data-movement overheads.

Finally, the execution of program fragments that do not depend on each other could become concurrent: since parts 2 and 3 are independent on each other, we should be able to overlap their execution in a task-parallel fashion.

**Key challenges:** Unfortunately, exploiting these opportunities to scale out execution automatically is particularly challenging in the context of the shell. First, exposing opportunities at the level of individual commands such as `cut` and `grep` is challenging—and this is why prior systems often focused on coarser, script-level or job-level granularity [19, 69].

Second, pervasive dynamic features, file-system introspec-

```
NOAA=${NOAA:-http://ndr.md/data/noaa/}
TEMPS=${TEMPS:-/noaa/temps.txt}
hdfs dfs -mkdir /noaa

## Pipeline 1: Download temperature data
##             and store to HDFS
seq $1 $2 | sed "s;^;$NOAA;" |
  sed 's;$/;' | xargs -r -n 1 curl -s | grep gz |
  tr -s ' \n' | cut -d ' ' -f9 |
  sed 's;^\(.*\) \(20[0-9][0-9]\) \.gz;\2/\1\2\2.gz;' |
  sed "s;^;$NOAA;" | xargs -n1 curl -s |
  gunzip | hdfs dfs -put - $TEMPS

## Pipeline 2: Compute maximum temperature
##             over all data
hdfs dfs -cat $TEMPS | cut -c 89-92 | grep -v 999 |
  sort -rn | head -n1 > max.txt

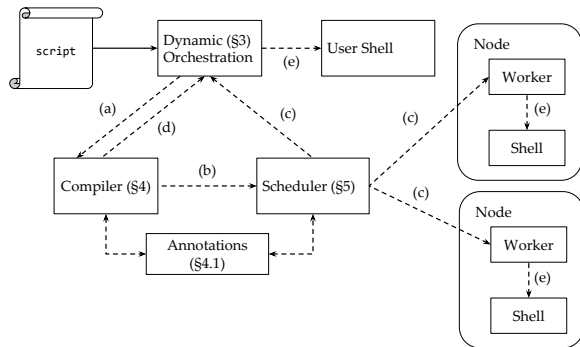
## Pipeline 3: Compute average temperature
##             over all data
hdfs dfs -cat $TEMPS | cut -c 89-92 | grep -v 999 |
  awk "{ t += $1; i++ } END { print t/i }" > avg.txt
```

**Fig. 1: Example script.** Downloading a temperature dataset, storing on a distributed file system, and running analysis to extract statistics.

tion, and other side-effects impede traditional distribution approaches based on static transformation—this is why prior shell-script distribution work [25, 49] focuses on side-effect-free dataflow subsets. These challenges are compounded by the presence of more elaborate control flow such as `for` loops, `break`, and `trap` statements present in ordinary shell scripts.

Third, behavioral equivalence with existing shells is practically unattainable, especially with shell reimplementations—after all, even production-grade shells such as Bash and zsh diverge subtly in their POSIX behavior [23]. A new distributed shell [14, 49] has little hope of *not* breaking some scripts.

**DiSH overview:** To overcome these challenges DiSH (1) extracts details about the behavior of commands through command annotations, (2) deals with dynamic features and side-effects by analyzing scripts at runtime using dynamic orchestration, and (3) achieves behavioral equivalence with Bash by only performing script transformations and delegating execution to the underlying interpreter. DiSH is designed to dynamically orchestrate, compile, schedule, and support the execution of shell scripts (Fig. 2). DiSH’s orchestration (§3) kicks in when a potentially distributable script region is identified, saves a snapshot of the user’s shell environment (variables, configuration) and invokes the DiSH compiler with the candidate region (Fig. 2a). The compiler analyzes this region and if possible, translates it to a dataflow graph—which it then optimizes to introduce parallelism, buffering, *etc.* (§4), finally passing it off to the scheduler (Fig. 2b); or aborts compilation (Fig. 2d) because it cannot guarantee that the region is pure,

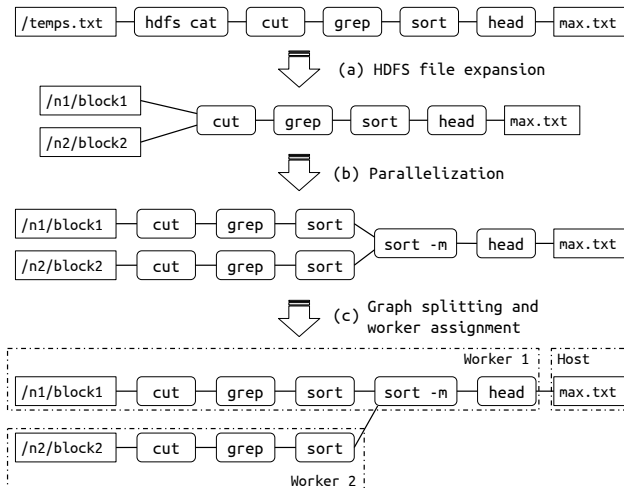


**Fig. 2: DiSH architecture overview.** Steps: (a) compile script region; (b) schedule compiled dataflow; (c) send dataflow subgraphs to workers; (d) compilation failed, fall back to original region; and (e) execute script region (compiled or original).

*i.e.*, side-effect-free. The scheduler (§5) divides the compiled dataflow graph into different subgraphs which it sends to available cluster workers (Fig. 2c). In response to these execution requests, workers apply a second pass of optimizations to better utilize available resources, translate the dataflow graph back to a shell script (Fig. 2e), load the snapshot of the shell environment stored by the orchestrator, and execute the script using the local, unmodified shell interpreter (§6).

**Applying DiSH:** DiSH preprocesses the script in Fig. 1 to identify script regions that could benefit from distribution—in this case, all three pipelines. It then replaces each of these regions with calls to the dynamic orchestrator and attempts to distribute them at runtime. During execution, the orchestrator queries the DiSH compiler to determine whether a region is pure and thus distributable: if the compiler succeeds, it translates the region to a dataflow graph. Since regions contain arbitrary black-box commands, DiSH cannot analyze them directly. Instead, it employs a command specification framework that contains partial specifications of command invocations such as their inputs and outputs. For example, DiSH’s compiler uses these specifications to determine that `hdfs dfs -cat /noaa/temps.txt` reads from the HDFS file `/noaa/temps.txt` and writes to `stdout`. Once a region is in dataflow form, DiSH applies transformations to distribute it.

Fig. 3 shows the distribution stages for pipeline 2 (maximum temperature). DiSH first detects operations on HDFS files (*i.e.*, HDFS `cat`) and expands each distributed file to its segments (datablocks), often stored on different physical machines. Informed by command annotations, DiSH applies parallelization transformations: commands like `cut` and `grep` are parallelizable directly and can be executed on the machine with the raw input datablock. The scheduler then splits the compiled graph into subgraphs and maps them to workers in a data-aware fashion. Finally, each worker translates the graph back to a shell script, adds additional runtime primitives (commands), and executes it locally.



**Fig. 3: DiSH dataflow graph stages.** (a) HDFS files are expanded to sequences of blocks. (b) the graph is parallelized based on the command specifications. (c) the scheduler splits the graph and assigns subgraphs to workers.

The result? DiSH drops the execution of pipeline 2 from 352s to 6s while maintaining full behavioral equivalence and requiring no modifications to the user shell.

### 3 Dynamic Shell Orchestrator

To facilitate adoption, an important desideratum in the design of DiSH is to achieve behavioral equivalence with the underlying shell interpreter. To achieve this, DiSH is not designed to operate as another shell, but rather wraps the user’s existing shell interpreter and the shell interpreters on the worker machines. As a result, DiSH hides parallelization and distribution from both the user and the underlying shells: the user thinks that their original script is being executed—just faster—and each underlying shell simply executes a standard non-distributed shell script. This allows DiSH to achieve exceedingly high compatibility with the underlying shell implementation (§7.3), while also minimizing maintenance costs since updates and modifications on the underlying shell are reflected in DiSH without any change.

Fig. 4 shows an overview of the structure of DiSH’s dynamic orchestration. To achieve dynamic shell script orchestration without any shell-interpreter modification, DiSH opts for a light-weight script instrumentation pre-processing step: it instruments *potentially* distributable regions with invocations to the orchestration engine. It chooses regions with the goal of maximizing distribution benefits: intuitively, it focuses on commands and pipelines rather than control-flow statements and variable assignments. However, the choice of these region boundaries is not binding—the preprocessor just needs to be precise enough to determine potential regions, but DiSH will eventually decide whether or not (and if yes, how) to distribute

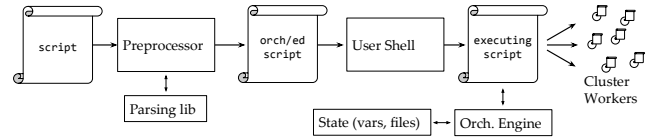
a candidate region at runtime. The preprocessor first parses the original script, it then replaces the relevant program regions with orchestration prefixes, and then un-parses (emits) it back as an instrumented script that is given for execution to the user’s shell interpreter.

The instrumented script then makes calls to the orchestration engine. The orchestration engine is itself a shell script coordinating with the compiler and worker manager and attempting to distribute the upcoming region (see §4 and 5 for details). If it succeeds, it runs the distributed version of the region. If it aborts, it just falls back to the original region, executing it normally. Reasons for aborting include the region being side-effectful, *e.g.*, modifying some environment variable, or lacking relevant command annotations.

**Preprocessor:** The preprocessor searches for maximal potentially distributable regions by processing the AST bottom-up, combining distributable subtrees when they are composed using constructs that do not introduce scheduling constraints (*e.g.*, `&`, `|`). When a region cannot outgrow a certain subtree, DiSH replaces it with a call to the orchestration engine. If the region is successfully compiled (at runtime), DiSH translates it to a dataflow representation—a convenient and well-studied model amenable to transformation-based optimizations [26]. At a later point, DiSH running on each node translates the instrumented AST resulting from the compilation back to shell code and passes it to the underlying shell for execution.

**Parsing library:** DiSH invokes parsing and unparsing routines frequently, and therefore needs them to be very efficient. To that end, it uses an internal Python implementation [32] of POSIX-shell-script parsing and unparsing based on `libdash` [22, 23]. The DiSH parser contains several optimizations such as caching, inlining, and careful array appending to achieve improved performance.

**Orchestration engine:** DiSH’s orchestration engine is designed to maintain the original script behavior and minimize runtime overhead—as it is invoked multiple times per script. The engine is a reflective shell script: it coordinates transparently with the compiler to determine whether or not to parallelize a script by inspecting the state of the shell and that of the broader system. DiSH constantly switches between two execution modes when executing scripts: (1) conventional shell mode, where scripts execute in the original shell context, and (2) DiSH mode, where the runtime reflects on shell state and invokes the compiler to determine whether to execute the original or an optimized version of the target region. To switch from shell mode to DiSH mode, the engine saves the state of the user’s shell; to switch back, it restores the state of the user’s shell. The state of a shell is quite complex: apart from saving and restoring variables, DiSH must account for various shell flags along with other internal shell state (*e.g.*, the previous exit status, working directory). During an invocation, the engine first switches to DiSH mode, communicates with the compiler and scheduler to determine whether a region can



**Fig. 4: Dynamic orchestration overview.** DiSH instruments scripts with calls to the orchestration engine, which passes program fragments to the worker manager at run-time.

be safely distributed, and it then switches back to shell mode to execute the original or distributed version of the script.

**Environment sharing:** The distributed version of the script region might execute on a different shell (or even machine). Therefore, a challenge that DiSH needs to address is to make sure that all regions execute in the correct environment—including access to the latest variable values and function definitions. To achieve that the engine takes a snapshot of the environment right before execution. It then transfers the snapshot to the distributed workers, which they load before executing the incoming script fragment. This is safe to do since successful distribution of a region implies that it is pure (and therefore does not affect the environment), and thus the snapshot will be valid until the region finishes execution.

**String expansion:** To correctly determine if a script region is safe to distribute, the compiler needs to expand all strings in that region. Since DiSH performs compilation and distribution of each script region at runtime, right before execution, it has access to all the latest variables and system state to fully expand all strings in the region. DiSH only implements a common and safe subset of all available expansions, and avoids implementing side-effectful expansions that have the risk of affecting the environment (*e.g.*,  `${x=foo}`: set `x` to `foo` if `x` is unset). Note that DiSH keeps expansion local: it does not expand regions succeeding the target region, as these might depend on the execution of the target region.

## 4 Compiler

This section describes the compiler of DiSH, which builds on the PASH parallelizing compiler [64]. The compiler is given the AST of an input script fragment and information about the commands in that fragment (§4.1). It then attempts to transform it to a dataflow graph (§4.2), an intermediate representation amenable to parallelizing transformations. If the compiler succeeds in transforming a script region to a parallel dataflow graph, that graph is then passed to the scheduler which then decides how to map subgraph components to the available worker nodes. As the compiler operates at runtime in a just-in-time fashion, it exploits ample opportunities for parallelization even across subgraphs (§4.3).

## 4.1 Command Annotations

DISH needs to support analyses and transformations over third-party commands, without access to their source code. To achieve this, DISH uses annotations à la PASH [64] and POSH [49], capturing information about a command invocation’s parallelizability class, inputs, and outputs. Command annotations act as an intermediate layer that provides restricted but sufficient information about the behavior of a command to analysis and transformation systems like DISH. They also enable reuse, as they are not tied to a particular analysis and can thus be reused by different tools. For this work, DISH reuses the set of annotations developed by the authors of PASH [64] extended annotations for commands that appear in the evaluation of DISH (§7).

A command annotation in DISH encodes information at the level of individual command invocations, *i.e.*, precise instantiations of a command’s flags, options, and arguments. Among other information, annotations determine how a command invocation affects its environment, and specifically whether it is pure, *i.e.*, whether it only affects its environment by writing and reading to and from a well-defined set of files—information which DISH uses when translating commands to and from dataflow nodes (§4.2). For example, the annotation for `grep` can be used to extract that the script fragment `grep -f dict.txt src.txt > out.txt` contains two input files `dict.txt` and `src.txt` and one output file `out.txt`. This knowledge of input and output files is used by DISH to enable location-aware distribution, by scheduling the computation on nodes that contain relevant data blocks. Additionally, annotations describe parallelization opportunities—*e.g.*, that `grep "pattern" src.txt` processes each line of `src.txt` independently and thus can be parallelized at a line boundary.

## 4.2 Dataflow Model

The core of DISH’s compiler is an order-aware dataflow model that captures pure shell script regions that read from a well-defined set of input files and write to a well-defined set of output files—*i.e.*, they do not modify their environment in any other way. This model is expressive enough to capture a shell subset used pervasively in data processing scripts [26].

In this model, nodes represent commands and edges represent files, pipes, named FIFOs, and file descriptors. The model is order-aware in the sense that it keeps information about the order in which nodes read from their inputs, which is important for the script’s semantics. For example, `grep "pattern" in1.txt - in2.txt` first reads from `in1.txt`, then from its standard input, and then from `in2.txt`. This order awareness allows DISH to perform transformations that optimize execution of a script—*e.g.*, by exposing parallelism—but preserve its original behavior.

**Translation workflow:** Given an AST representation of an input script region, the compiler uses annotations to deduce

whether commands are pure *i.e.*, they only affect their environment through a well-defined set of output files, and attempts to transform them to dataflow nodes. If all commands in the region are pure the compiler transforms the region to a dataflow graph. It then applies transformations (described below), optimizing the graph to expose parallelism and improve the script’s performance. Finally, it serializes the graph back to a (now optimized) shell script, by translating every node back to a command and connecting them all together with appropriate channels (*e.g.*, FIFOs, RFIFOs, redirections).

**Transformations:** DISH’s transformations enable data-parallel execution by replicating nodes in the graph and adding appropriate split and merge nodes around them. They apply a pass over the graph to remove pairs of inverse nodes—*i.e.*, pairs of nodes whose semantic effects cancel out but whose performance effects are additive—for example, a concatenation-style merge followed by a linear split. For commutative commands, *i.e.*, commands that produce the same output regardless of their input-line order, DISH applies transformations that pack and unpack metadata across the graph—achieving better performance by avoiding unnecessary blocking and buffering. Finally, to improve the flow of data across the graph, DISH applies additional transformations that inject hybrid memory-disk buffer nodes in points in the graph that are likely to become bottlenecks.

**Remote file resources and HDFS files:** To support scripts that perform data analysis on a combination of HDFS and local files, DISH extends the dataflow model with remote-file resources (RFRs) that encode file blocks in different nodes. RFRs usually represent blocks of files that are partitioned and replicated in HDFS, and contain information about the location of the data in the distributed environment. This information could contain multiple locations to support replication, and is used by the scheduler to assign script fragments to different workers. When the DISH compiler comes across an HDFS file path, it queries HDFS to determine the locations of its file blocks and then expands that file to a sequence of RFRs, each of which represents a block.

## 4.3 Dynamic Dependency Untangling (DDU)

Scripts often contain regions that are independent, *i.e.*, they have different (file) working sets. Independent regions could potentially run in parallel, better utilizing computational resources and improving the execution times of the scripts in which they belong. However, inferring independence statically and ahead of time is challenging as shell scripts make extensive use of dynamic features. Figure 5 shows an example script that contains independent fragments but also features dynamic behavior. This script iterates over all files in an HDFS directory, compresses them using `gzip`, and finally stores them as independent files.

Determining independence statically in this script would

```

for item in $(hdfs dfs -ls -C ${IN});
do
    output_name=$(basename $item).zip
    hdfs dfs -cat $item |
        gzip -c > $OUT/$output_name
done

```

**Fig. 5: Example of independent regions.** This shell script compresses all files in a directory—but each iteration results in an independent body region that can be executed in parallel.

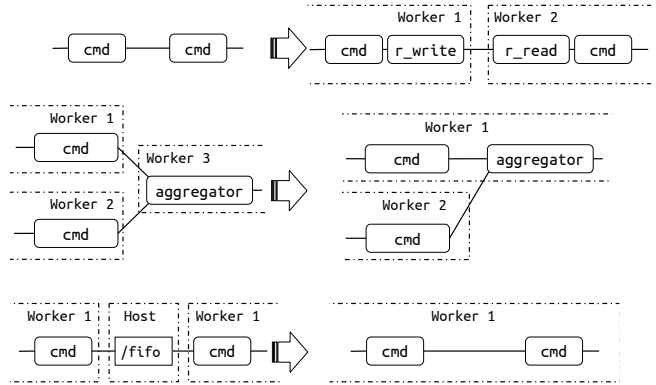
require inferring values of environment variables (like `IN` and `OUT`) and the state of the file system, *e.g.*, `hdfs dfs -ls`. DiSH’s dynamic orchestration (§3) circumvents this challenge by making distribution decisions during the execution of the script when environment variables and the file system state are known. DiSH further exploits this by discovering independent dataflow regions at runtime and executing them in parallel—even if they were not parallel in the original script.

When DiSH successfully compiles a dataflow region (at runtime), it knows that the region is pure and therefore can determine the region’s inputs and outputs—and it does so for free, without additional analysis or inference stages. DiSH then uses this information to check for read-write or write-write dependency conflicts with regions that are running concurrently. If none is found, DiSH passes the region to the scheduler, which orchestrates distributed execution, and then immediately continues the execution of the script until it reaches the next dataflow region. Whenever the compilation of a dataflow region fails, DiSH cannot safely detect the input and output information of this region—and thus it needs to wait until every previous region is done executing to ensure that no dependency will be violated.

Since DDU is done at runtime it is both sound, *i.e.*, it does not execute dependent fragments concurrently, and precise, *i.e.*, it offers significant benefits due to improved parallelism and resource utilization—especially for scripts that do not contain highly data-parallelizable commands, such as the commands in the aforementioned compression script (Fig. 5). Compared to analyses over static languages, DDU cannot identify global optimizations such as reordering the final command in the script to run first. This lack of optimality is not specific to DDU, but applies to any shell script analysis; in fact, as far as we know there is no sound and precise static analysis for shell scripts.

## 5 Distributed Scheduling

This section describes how DiSH’s scheduler distributes a compiled script to a set of workers. The scheduler is given a dataflow graph that is already parallelized and has HDFS files expanded to sequences of remote file resources (RFRs) representing their blocks. The task of the scheduler is then to distribute this graph with the goal of optimizing performance



**Fig. 6: (Top)** Remote writes and reads added during distributed scheduling. **(Mid)** Worker-first aggregation. **(Bot)** Named FIFO teleportation.

by both utilizing available resources and moving computation close to the data. Currently the scheduler knows about the workers in the cluster ahead of time using a configuration file.

The scheduler makes a decision on how to split the graph based on a policy that optimizes performance through collocation of data blocks and the commands that execution over them. The scheduler processes the top-level dataflow graph to generate a set of subgraphs, one for each worker and one for the host machine executing the script. It then replaces edges corresponding to communication channels (*e.g.*, FIFOs, pipes) at the boundaries of each subgraph with remote channels—adding a remote write node on the sender side and a remote read node on the receiver side (see Fig. 6, Top). It also inserts remote reads for subgraph nodes that access files stored on remote workers. The final generated subgraph represents the script fragment that is passed for execution to the user shell running on each worker: the compiled script handles all the redirection to and from local files and the standard input, output, and error streams to and from the worker.

**Data-aware scheduling policy:** The highest performance overhead when executing distributed shell scripts is networked data movement across workers. DiSH addresses this overhead by introducing a greedy scheduling policy that allocates subgraphs in a way that attempts to minimize data movement across workers. If a data file (or block) is available on a worker, then DiSH maps the maximal dataflow subgraph that starts from that file to that worker—*i.e.*, scheduling as much of the processing as possible on the worker. The scheduler also tracks the amount of work that each worker currently has scheduled, which can vary due to dynamic dependency untangling (§4.3): if a data file is replicated across multiple workers, DiSH chooses the worker with the least amount of pending work to execute that subgraph.

**Worker-first aggregation:** The distributed dataflow graphs that DiSH executes often contain aggregation (*i.e.*, merge)

nodes, similarly to the reduce stages in Hadoop Streaming. Regardless of the worker on which the aggregation is performed, data from different workers will need to be combined onto a single worker and thus these dataflow nodes will necessarily result in data movement. DiSH prioritizes performing aggregation on one of the participating workers, because workers already contain a subset of the data used in the aggregation (see Fig. 6, Mid). This optimization is particularly beneficial for scripts that filter and aggregate data, often containing commands such as `grep` and `uniq`, because any filtering stages prior to aggregation result in reduced data transfer.

It is worth noting that, absent additional information about commands [49], the location of aggregators involves challenging trade-offs not addressable with a single optimization policy. For scripts that include aggregators that do not reduce data sizes, DiSH’s worker-first aggregation optimization risks transferring more data. As DiSH’s evaluation confirms (§7), however, worker-first aggregation results in performance benefits for most scripts.

**Delegated script concretization:** DiSH’s scheduler sends workers dataflow subgraphs, encoded in DiSH’s intermediate representation, instead of concrete shell scripts ready for execution. Each dataflow subgraph contains holes that workers are expected to fill in, based on the specifics of their local environment. This choice simplifies DiSH’s distributed execution, as the scheduler does not need to have up-to-date information about several worker details such as the temporary directories they use. Additionally, this choice enables better resource utilization in a heterogeneous environments with different worker capabilities: a worker can apply another optimization pass to the dataflow subgraph it receives to better manage and utilize its resources.

**Named FIFO teleportation:** Scripts often use named FIFOs to share data between concurrently executing processes. Named FIFOs introduce a performance challenge, because they are local files that reside on the host machine where the script was executed. Therefore, by default, all data that would normally go through named FIFOs in the original execution would now have to go back and forth between workers and the machine for which the script was developed. DiSH addresses this challenge by observing that named FIFOs are ephemeral, *i.e.*, they maintain no data after the execution of a dataflow region. Based on this observation, DiSH migrates named FIFOs to workers closer to the data, eventually deleting the migrated versions after the dataflow region has finished executing (see Fig. 6, Bot). This transformation, termed FIFO teleportation, improves performance by avoiding unnecessary data movement in scripts that use FIFOs.

## 6 Runtime Support

DiSH has to address several runtime challenges: communication among workers, identification of HDFS data block lo-

cations, and correctness in view of HDFS blocks split independently of newlines—an assumption necessary for several dataflow transformations. This section describes several components of DiSH’s runtime that address the above challenges.

**Remote FIFO channel:** As described earlier (§5), connections between dataflow nodes are instantiated using UNIX FIFOs in a single-machine setting. Unfortunately, FIFOs do not support networked operation and thus cannot cross worker boundaries. To address this challenge, DiSH introduces a remote FIFO primitive (RFIFO) that is implemented in Go and uses socket-based communication. RFIFOs are intended to operate identically to FIFOs, *i.e.*, implement the semantics of dataflow graph edges, but with support for operation over the network. They have a unique identifier and two ends—a read end and a write end.

Since shell streams are lazy, *i.e.*, a producer blocks until its consumer requests input, the network link is often not fully utilized, lowering throughput and risking introducing significant latency. To avoid these throughput and latency challenges, DiSH adds two buffer nodes to the dataflow graph: one before the write end of the RFIFO, to allow uninterrupted access to data, and one after the read end of the RFIFO, to force the read to request data. This lazy-to-strict optimization maintains correctness and improves performance in most cases; in rare cases, it may lead to unnecessary data transfer between nodes—*e.g.*, when there is a `head` command right after the read end of an RFIFO.

**Port discovery service:** As transformations and optimizations are applied during the execution of a script—contrary to most other distributed environments—DiSH’s scheduler cannot statically predict which ports will be available at runtime for RFIFOs at each worker: different scripts and script fragments running concurrently during a single execution may collide on port usage. To address that, each DiSH worker implements a port discovery service (PDS) that can be accessed by remote FIFOs to (1) advertise their port, and (2) discover the port that their other end uses. The discovery service is implemented in Go with gRPC [61] and supports a few remote procedure calls (RPCs), central among which are a `put` call for advertisement and a `get` call for discovering the port of a remote end. RFIFOs are extended with gRPC clients to advertise ports among local PDS or identify the ports corresponding to their other end by querying the PDS of the respective worker. By deferring port selection until runtime execution, DiSH’s port discovery service facilitates loose subgraph coupling and simplifies remote subgraph execution on multiple workers.

**HDFS data retrieval:** During transformations, the DiSH compiler (§4) needs to retrieve information about HDFS paths to expand them into block sequences. This expansion happens on a critical runtime path and thus needs to be efficient. A prior implementation of DiSH invoked this expansion on every HDFS path using a shell command—by wrapping `fsck`, a command offered by HDFS API for querying the health of



the disk in the cluster, returning information about a file and its partitioning into blocks. This implementation ended up incurring significant latency (> 1s), and thus DiSH switched to the web API reducing expansion to sub-10ms latency.

**Enforcing logical block boundaries:** A key challenge when processing separate file blocks in HDFS is the mismatch between compiler assumptions about the block shape and how blocks are actually stored in HDFS: HDFS blocks might not be split on newline boundaries, but the parallelizing transformations performed by the DiSH compiler (§4) assume that all blocks are logically separated by newlines. This assumption is crucial and depends on the way commands process their input, *e.g.*, `sort` processes its input line by line, and therefore would require a significantly more complex parallelization transformation if its input could be split at arbitrary points. Developing complex custom parallelization transformations for each command would be infeasible in practice due to the sheer number of available commands and would not allow DiSH to reuse the parallelization transformations developed for PASH [64].

Instead of relaxing the compiler assumption, DiSH addresses the mismatch by ensuring it holds during script execution using additional runtime support. DiSH implements a distributed file reader (DFR) primitive that runs as a service on every worker. The DFR service ensures that parallel dataflow nodes only process batches that are split in newline boundaries, independent of how the actual physical blocks are split—providing the illusion of a logical block that ends at a newline to its consumer. Given a distributed file path, DFR reads the local file or block from the worker’s disk going beyond the first newline character in its block. If the block is not terminated with a new line, then the DFR communicates with the reader of the next block (and potentially any readers after that), returning a complete logical block to its consumer. When a compiled dataflow graph is translated back to a script, DiSH prefixes file paths with a command invoking a DFR client that communicates with the relevant DFR service to retrieve the relevant logical block. Both service and client are implemented in Go, communicating using gRPC and protobufs [21].

## 7 Evaluation

We are interested in evaluating two aspects of DiSH: (1) its performance, and (2) its compatibility with Bash.

**Experiments:** We perform four experiments using several real-world shell scripts taken from a variety of sources (Tab. 2). The first two experiments focus on the performance gains (§7.1) achieved by DiSH’s distribution on (1) a 4-node on-premise cluster, and (2) a 20-node cloud deployment—both over a variety of benchmarks and workloads. We compare DiSH’s performance against (1) GNU Bash [50], the *de facto* sequential shell-script execution environment; (2)

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

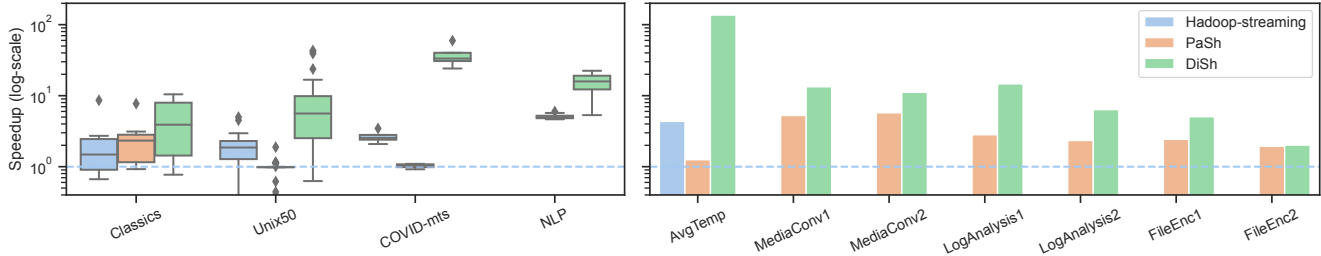
Benchmark	Scripts	Pure HS	LOC	Input	Source
1 Classics	10	7/10	123	3G	[5, 6, 31, 39, 59]
2 Unix50	34	30/34	142	21G	[7, 34]
3 COVID-mts	4	4/4	79	3.4G	[62]
4 NLP	21	-	306	120 books	[9]
5 AvgTemp	1	1/1	31	3.6G	[68]
6 MediaConv	2	-	35	0.8 & 0.4G	[49, 56]
7 LogAnalysis	2	-	63	0.7 & 1.3G	[49, 56]
8 FileEnc	2	-	44	1.3G	[41]

Apache Hadoop Streaming [25] (AHS), a production-grade distributed data-processing framework that supports language-agnostic executables; and (3) in the case of the 4-node setup, PASH [32, 64], a shell-script parallelization system from the Linux Foundation. PASH’s parallelism benefits make it a likely alternative to DiSH for smaller clusters, where DiSH’s anticipated benefits of distribution might be smaller, but this likelihood diminishes as the size of the cluster grows.

The last two experiments evaluate DiSH’s dynamic dependency untangling (§7.2) and DiSH’s correctness (§7.3), *i.e.*, its compatibility with respect to Bash across all scripts and the POSIX shell test suite.

**Benchmarks:** We use 8 sets of real-world benchmarks, totaling 76 shell scripts and 823 LoC. Classics and Unix50 contain classic and recent (c. 2019) scripts that make heavy use of UNIX and Linux built-in commands. COVID-mts contains four scripts used to analyze real telemetry data from mass-transit schedules during a large metropolitan area’s COVID-19 response. NLP contains several scripts from UNIX-for-poets, a tutorial for developing programs for natural-language processing out of UNIX and Linux utilities. AvgTemp contains a large script downloading and processing multi-year temperature data across the US. MediaConv contains two scripts that process, transform, and compress video and audio files. LogAnalysis contains two scripts that apply typical system-administration and network-traffic analyses over log files. Finally, FileEnc contains aliases that encrypt and compress files.

**Baselines and implementations:** Bash, PASH, and DiSH executed every shell script completely unmodified. Apache Hadoop Streaming (AHS) posed significant expressiveness limitations. Only 42 scripts in Classics, Unix50, COVID-mts, and AvgTemp out of 76 scripts can be implemented natively (Tab. 2, col. Pure HS). Another 7 scripts required manual porting by splitting them into mappers, reducers, and additional components: These components were not available natively by AHS—for example, components for reading from two pipelines for `diff.sh` and for sorting after the reducer for `bigrams.sh` (both in Classics). During porting, we put significant care to avoid limiting AHS’s parallelism: we modified 3 AHS scripts in Classics to help HS introduce additional



**Fig. 7: DiSH performance on a 4-node cluster.** DiSH speedup (vs. PASH and Hadoop Streaming whenever possible) over Bash for Tab. 2 rows 1–4 (left, box) and 5–8 (right, bar) (Cf. §7.1). (Log y-axis; higher is better.)

parallelism—for example, we manually expanded `tr -cs` into `tr -c | grep -v` (both stateless). None of the scripts in NLP, MediaConv, or LogAnalysis can be implemented in AHS as they perform processing in loops, the iterations of which depend on the files in a statically indeterminable directory (see Fig. 5) and are therefore not expressible in AHS. We attempted to replace the body of the loop with an AHS invocation but the startup overhead ended up dwarfing the execution time by a factor of ten on average.

**Hardware & software setup:** The 4-node cluster consists of four 6-core Intel(R) Core(TM) i7-10710U CPU nodes each with 64GBs of RAM, located in the same room and connected with an average bandwidth of 90.8 Mbits/sec. The 20-node deployment consists of x1170 Cloudlab [15] nodes, each equipped with 10 × Intel Core E5-2640 2.4 GHz CPUs and 8GB of memory. Single-machine shells (Bash & PASH) were evaluated on a machine with 20 × 2.80GHz Intel(R) Core(TM) i9-10900 CPUs and 32GB of memory.

For ease of deployment and reproducibility, we used Docker swarm to deploy (1) HDFS, and (2) the DiSH runtime. The containers were created using the standard Ubuntu 18.04 image. We use Bash v.5.0.3, PASH v.6e2ecba, and HDFS/Hadoop streaming version 3.2.2. We explicitly disabled checksum verification from HDFS in all configurations, scripts, and measurements. All scripts were executed completely unmodified, using environment variables, loops, and other shell constructs. To minimize statistical non-determinism we repeated the experiments several times noticing imperceptible variance.

The DiSH implementation comprises 6784 lines of Python (preprocessor, compilation server, expansion, compiler, and parser), 1011 lines of shell code (JIT engine and various utilities), and 1174 lines of C (commutativity primitives, and other runtime components). All counts include only semantically meaningful lines of code.

## 7.1 Performance

How does DiSH’s distributed perform on small on-premise clusters and multi-node cloud deployments, and how does it compare to state-of-the-art systems?

**Tab. 3: DiSH performance in 20-node cloud deployment.** DiSH speedup over Hadoop Streaming for scripts that AHS supports.

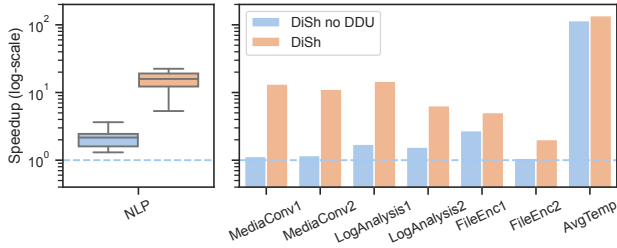
DiSH speedup over AHS						
Benchmark	Avg	Min	25th	50th	75th	Max
Classics	2.74×	0.92×	2.41×	2.60×	2.85×	6.55×
Unix50	6.64×	0.91×	2.85×	5.38×	10.4×	16.9×
COVID-mts	10.4×	6.64×	8.91×	9.27×	-	16.8×
AvgTemp	7.85×	-	-	-	-	-

**Results:** Fig. 7 (note the log y-axis) shows the performance of DiSH, PASH, and AHS on a 4-node on-premise cluster across all benchmarks of Tab. 2. Box plots (left) show result quartiles for multi-benchmark suites (Tab. 2, rows 1–4) and bars (right) show results for individual scripts (Tab. 2, rows 5–8). Across all benchmarks, DiSH achieves an average speedup of 13.6× (vs. 2.55× for PASH and 2.1× for AHS) and a maximum speedup of 136.3× (vs. 7.8× for PASH and 8.6× for Hadoop Streaming). The average execution time of all scripts on Bash is 299s, ranging from 1s for `34.sh` in Unix50 to 2840s for `nfa-regex.sh` in Classics. DiSH is only slower than Bash (737s vs 568s) in the case of `diff.sh` from Classics, for which AHS is even slower (766s). DiSH achieves a performance comparable to Bash (1-2s) in `4.sh` and `34.sh` from Unix50, because both perform a short-running head.

Tab. 3 shows the speedup of DiSH over AHS on a 20-node Cloudlab deployment across all scripts implementable with AHS (Classics, Unix50, COVID-mts, AvgTemp). Across all benchmarks, DiSH achieves an average speedup of 6.17× and a maximum speedup of 16.95× over AHS. DiSH is slower than AHS only for three scripts: `nfa-regex.sh` from Classics (0.92×), `29.sh` and `30.sh` from Unix50 (0.91× and 0.94×).

Across all scripts in both deployments, DiSH’s overheads (startup cost, dynamic orchestration, preprocessing, compilation, scheduling) take less than 1 second.

**Discussion:** DiSH is faster than Bash, PASH, and AHS across Tab. 2’s suites (rows 1–4) with respect to average, and across all of Tab. 2 individual benchmarks (rows 5–8)—often by a significant margin (e.g., 134× for AvgTemp against PASH).



**Fig. 8: Dynamic dependency untangling.** D1SH speedup over Bash when toggling DDU (higher is better).

D1SH’s (and PASH’s) speedup over Bash is due to parallelism. D1SH’s speedup over PASH is due to D1SH’ co-location of data and computation: PASH cannot offload computation and thus first gathers all data onto a single machine—a time-consuming stage—and then starts processing in parallel. D1SH is slower than Bash only for `diff.sh`, because (1) it is not highly parallelizable and (2) it performs no filtering, *i.e.*, its output is the same size as its input. In contrast to Bash, which simply fetches all data and processes it locally, D1SH tries to allocate most commands on the workers, but this leads to increased data movement since moving data between workers does not avoid sending the whole output to the client.

D1SH’s speedup over AHS is due to a few different reasons. One reason is the increased expressiveness of D1SH’s dataflow model: D1SH accepts and parallelizes complete scripts, discovering more opportunities for parallelism. Many of the AHS scripts are broken into multiple map and reduce stages, often leaving pipeline parallelism and data parallelism unexploited. Another reason is D1SH’s dynamic independence discovery, which allows for additional parallelism and better utilization of resources—in ways that AHS does not support; we zoom into these benefits below (§7.2). In the Cloudlab deployment, D1SH is (marginally) slower than AHS in only two cases: (1) a script that is embarrassingly parallel and thus implementable in AHS using only a single mapper (`nfa-regex.sh`), and (2) two scripts in Unix50 that see slightly more benefits from our manual, hand-optimized AHS rewrite than they do from D1SH’s automated distribution.

We found porting scripts to AHS a serious challenge. Many scripts required significant manual effort, resulted in multiple error-and-fix cycles, and led to script size increases. To overcome AHS’s expressiveness limitations, we had to modify a few scripts in unintuitive ways—often combining plain Bash scripts with AHS mappers and reducers. These modifications made scripts significantly more complex and compounded the effort to test and maintain them. Instead, D1SH distributed scripts successfully without any such challenges.

## 7.2 Dynamic Dependency Untangling

What is the speedup due to dynamic dependency untangling?

**Results:** Figure 8 shows D1SH’s speedup over Bash with and without dynamic dependency untangling (DDU, § 4.3). It excludes scripts that contain a single dataflow, for which DDU is not applicable. D1SH’s average speedup over D1SH-w/o-DDU is 6.9×, ranging between 1.2–13.9×.

**Discussion:** Enabling DDU improves performance significantly across all relevant scripts, by running independent dataflow regions in parallel. This allows D1SH to expose parallelism not just within data pipelines but across them, improving utilization. DDU also improves the distributed execution of scripts that operate on many files, many or most of which are small enough to fit on a single HDFS block.

DDU is the main reason why D1SH gets an edge over Bash on scripts that (1) have implicit independences that are not highly parallelizable, and (2) operate on small data that incur imperceptible data-movement costs. Examples of such scripts include `MediaConv1` and `FileEnc2`.

## 7.3 Correctness

What is D1SH’s output compatibility with respect to Bash?

**Results:** To check the correctness of D1SH across all benchmarks, we check that its stdout and exit status are equivalent to the ones produced by Bash. Across all benchmarks, totaling over 650 millions lines (18GB) of output, D1SH produces the same output and exit status as Bash.

We additionally execute the complete POSIX shell-test suite to evaluate D1SH’s compatibility with Bash. Out of all relevant tests, D1SH diverges from Bash in two cases and only with respect to the exit status it returns: both exit with an error, but Bash returns 1 whereas D1SH returns 127, which is outside of the POSIX mandated exit status range between 1–125. The reason is that D1SH always invokes the underlying Bash interpreter using the `-c` flag to set the `$0` variable, and Bash (contrary to most other shells, *e.g.*, `dash`, `ksh`, `mksh`, `sash`, `Smooch`, `yash`, `zsh`) exits with 127 in particular failing cases when called with `-c`.

**Discussion:** All benchmarks in Tab. 2 were executed with D1SH repeatedly. After hundreds of runs over several weeks, we observed dozens of different execution orders. Comparing the output on every run provides significant confidence about the correctness of the resulting distributed execution. The POSIX test suite mostly evaluates the correctness of dynamic orchestration (§3), as it does not feature many opportunities for parallelization and features no opportunities for distribution.

## 8 Related Work

D1SH is related to a large body of prior work.

**Distributed data processing:** Several environments assist in the development of distributed software systems: distributed computing frameworks [11, 44, 45, 57, 71] and domain-specific

languages [3, 8, 13, 40, 42] simplify the development of distributed systems that fall under certain computational classes such as batch processing, stream processing, *etc.* These systems deal with many of the challenges of distribution, but require developers to (re)write their computations manually in models that differ significantly from UNIX shell programming.

Hadoop Streaming and Dryad Nebula are abstractions that allow using third-party language-agnostic components similar to the UNIX shell, atop cluster-computing engines (Hadoop and Dryad, respectively). Both require their users to understand and rewrite their shell scripts using the abstractions provided by each framework. DiSH can operate on arbitrary shell scripts automatically, without requiring any manual effort from its users.

**Distributed shells and tools:** Several packages expose commands for specifying parallelism and distribution on modern UNIXes—*e.g.*, `qsub` [19], SLURM [69], calls to GNU `parallel` [58]. Different from DiSH, their effectiveness is predicated upon explicit and careful invocation and is limited to embarrassingly parallel (and short) programs. Often, these commands provide options to support an array of special sub-cases—a stark contradiction to the celebrated UNIX philosophy. For example, `parallel` contains flags such as `--skip-first-line`, `-trim`, and `--xargs`, that a UNIX user can achieve using `head`, `sed`, and `xargs`; it also includes other programs with complex semantics, such as the ability transfer files between computers, separate text files, and parse CSV. DiSH embraces the UNIX philosophy, attempting to rewrite shell programs to leverage distributed infrastructure.

Several shells [14, 38, 56] add primitives for non-linear pipe topologies—some of which target distribution. Here too, however, developers are expected to manually rewrite scripts to exploit these new primitives.

POSH [49] is a recent shell for scripts operating on NFS-stored data. It brings pipeline components closer to the data on which they operate, but operates only on shell pipelines that are fully expanded—*i.e.*, ones that do not use dynamic features. DiSH operates on shell scripts that use (1) any POSIX composition primitive, and (2) the full set of dynamic features present in the UNIX shell.

**Distributed operating systems:** There is a long history of networked and distributed operating systems [4, 12, 43, 46, 48, 51–53, 67]. These systems offer abstractions that (1) are similar, but not identical, to the ones offered by UNIX, (2) operate at a lower level of abstraction (*e.g.*, that of system calls, rather than shell primitives), and (3) often aim at simply hiding the network rather than offering scalability benefits. Instead of implementing full-fledged distributed operating system, DiSH shows that a thin but sophisticated rewriting-based shim can operate on completely unmodified programs, avoid requiring any user input, and achieve significant speedups by executing fragments in parallel across nodes.

**Annotation-based transformations:** Recent systems [47,

65, 70] lower the developer effort of scaling out program components by performing program transformations based on user-provided annotations. These systems operate in single-language environments, offering declarative DSLs for tuning the semantics of the resulting distributed program. DiSH uses a similar approach, leveraging command annotations from prior projects [49, 64], but operates on-the-fly—within an environment that makes extensive use of dynamic features and that allows combining components from multiple languages.

PASH-JIT [32] parallelizes scripts by dynamically interposing between a shell script and the underlying shell interpreter. This kind of interposition offers significant performance benefits without jeopardizing correctness, *i.e.*, maintains compatibility with the underlying shell interpreter. DiSH uses similar insights and interposition architecture, but operates on a distributed multi-node setting and addresses challenges that are specific to this setting—such as integration with a distributed file system and distributed environment passing.

**Cloud build systems:** Several cloud build systems [1, 2, 17, 27] distribute and parallelize the execution of large builds by constructing dependency graphs using dependency information explicitly specified by their users. Contrary to these systems, DiSH operates on general shell programs without exploiting domain-specific information—*e.g.*, build dependencies—and by taking a just-in-time approach that resolves dependencies during the execution of the script.

**Correct distribution of dataflow graphs:** The DFG is a prevalent model in several areas of data processing, including batch- and stream-processing. Systems implementing DFGs often perform optimizations that are correct given subtle assumptions on the dataflow nodes that do not always hold, introducing erroneous behaviors. Recent work [28, 33, 37, 54] attempts to address this issue by performing optimizations only in cases where correctness is preserved, or by testing that applied optimizations preserve the original behavior. DiSH uses its dynamic orchestration to achieve compatibility with the underlying shell and then achieves correct distribution on a per-region level by building on prior work on provably correct transformations for order-aware dataflow graphs [26]. Similarly to other automated shell script transformation works [49, 64], DiSH’s correctness is predicated upon the correctness of the annotations describing commands.

**Resurgence of shell research:** Recent shell research [10, 23, 24, 32, 36, 41, 49, 55, 56, 64] highlights renewed interest in shell scripting both as a vehicle for impactful research and as a target worthy of scientific attention. We see DiSH as a natural continuation of the insights and research behind recent systems [24, 26, 32, 49, 64], allowing other researchers to leverage DiSH’s POSIX-compliant high-performance dynamic distribution in their future work.

## 9 Discussion

**Programmability:** An important consideration with any automated system is how it affects programmability, and specifically the ability to debug a misbehaving program or to test a program for correctness. DiSH does not negatively affect the developer experience compared to a shell: a developer can use a combination of the many existing tools and commands—*e.g.*, `head` and `grep`—as they would normally do to inspect their script’s output and determine what is wrong. When it comes to shell scripts intended for distributed environments, DiSH in fact improves developer experience: a developer may use the same set of commands for local or distributed interactions—*e.g.*, to inspect and project parts of a file, regardless of whether that file is stored in HDFS or the local system. Furthermore, developers using DiSH can reap the scalability benefits of distribution in analyzing or testing scripts by automatically scaling out load to multiple computers.

**Command annotations:** DiSH’s transformations depend on the existence of command annotations. To maintain soundness, DiSH will avoid compiling and distributing a script region if some command lacks annotations. To increase the distributability of their scripts, DiSH users could opt for more constrained commands—*e.g.*, `cut` instead of `awk` for data projection—and thus enjoy tighter annotations and more applicable optimization transformations. The correctness of applying DiSH’s transformations depends on the correctness of these annotations, and thus annotations are currently expected to be authored by command developers or other experts (but not script developers). Developing automation for testing or synthesizing correct annotations is an interesting avenue for future research that would benefit several systems that use them—*e.g.*, DiSH, PASH [64], and POSH [49].

**Fault tolerance:** DiSH does not tolerate failures such as worker aborts or network partitions (§1). In such cases users are expected to rerun their scripts as shell users normally do in the non-distributed case. Achieving fault tolerance in the context of general shell scripts is in fact particularly challenging due to the prevalence of black-box components that may perform arbitrary side-effects. A fault-tolerant version of DiSH should be able to track all these side-effects and re-execute them appropriately when a script fails. This is in contrast to constrained cluster computing frameworks such as MapReduce and Spark that have precise information about the inputs and outputs of purely functional program components enabling simplified re-execution of dependency graphs (lineage) in the presence of failures. DiSH’s design however combined with incremental script execution [10] creates an opportunity for addressing this challenge with a hybrid approach: employ conventional fault tolerance approaches for script fragments with annotation information, and instrument the rest of the script to capture and replay its side-effects appropriately in cases of failure.

**Conclusion:** DiSH is the first system able to distribute unmodified shell scripts that use (1) any POSIX composition primitive, (2) the full set of dynamic features present in the UNIX shell, and (3) distributed file systems such as HDFS. DiSH uses a dynamic orchestration approach that instruments a given script and dynamically distributes it at runtime to then execute it using the underlying shell interpreter. As a result, DiSH avoids modifications to shell scripts and maintains compatibility with existing shells and legacy functionality. Evaluated against several alternatives available to users today, DiSH offers significant speedups, requires no developer effort, and handles arbitrary dynamic behaviors pervasive in shell scripts. DiSH is open-source software, available by the Linux Foundation.

## Acknowledgments

We would like to thank Ayush Bhardwaj, Felix Stutz, Hannah Gross, Lily Tsai, Malte Schwarzkopf, Michael Greenberg, Neil Ramaswamy, the Brown Systems Group, the NSDI 2023 reviewers, and our shepherd, Rebecca Isaacs, for discussions and feedback on the paper. This material is based upon work supported by DARPA contract no. HR00112020013 and no. HR001120C0191, and NSF award CCF 2124184.

## References

- [1] Bazel dynamic execution. <https://bazel.build/remote/dynamic>, 2022. [Online; accessed Feb 1, 2022].
- [2] Google cloud build. <https://cloud.google.com/build/docs/overview>, 2022. [Online; accessed Feb 1, 2022].
- [3] Peter Alvaro, Neil Conway, Joseph M Hellerstein, and William R Marczak. Consistency analysis in bloom: a calm and collected approach. In *CIDR*, pages 249–260, 2011.
- [4] Amnon Barak and Oren La’adan. The mosix multi-computer operating system for high performance cluster computing. *Future Generation Computer Systems*, 13(4):361–372, 1998.
- [5] Jon Bentley. Programming pearls: A spelling checker. *Commun. ACM*, 28(5):456–462, May 1985.
- [6] Jon Bentley, Don Knuth, and Doug McIlroy. Programming pearls: A literate program. *Commun. ACM*, 29(6):471–483, June 1986.
- [7] Pawan Bhandari. Solutions to unixgame.io, 2020. Accessed: 2020-04-14.

- [8] Martin Biely, Pamela Delgado, Zarko Milosevic, and Andre Schiper. Distal: A framework for implementing fault-tolerant distributed algorithms. In *Proceedings of the 2013 43rd Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, DSN '13, pages 1–8, Washington, DC, USA, 2013. IEEE Computer Society.
- [9] Kenneth Ward Church. Unix™ for poets. *Notes of a course from the European Summer School on Language and Speech Communication, Corpus Based Methods*, 1994.
- [10] Charlie Curtsinger and Daniel W Barowy. Riker: Always-Correct and fast incremental builds from simple specifications. In *2022 USENIX Annual Technical Conference (USENIX ATC 22)*, pages 885–898, 2022.
- [11] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: Simplified data processing on large clusters. *Commun. ACM*, 51(1):107–113, January 2008.
- [12] Sean Dorward, Rob Pike, David L Presotto, Dennis Ritchie, Howard Trickey, and Phil Winterbottom. Inferno. In *Proceedings IEEE COMPCON 97. Digest of Papers*, pages 241–244. IEEE, 1997.
- [13] Cezara Drăgoi, Thomas A Henzinger, and Damien Zufferey. Psync: a partially synchronous language for fault-tolerant distributed algorithms. *ACM SIGPLAN Notices*, 51(1):400–415, 2016.
- [14] Tom Duff. Rc—a shell for plan 9 and unix systems. *AUUGN*, 12(1):75, 1990.
- [15] Dmitry Duplyakin, Robert Ricci, Aleksander Maricq, Gary Wong, Jonathon Duerig, Eric Eide, Leigh Stoller, Mike Hibler, David Johnson, Kirk Webb, Aditya Akella, Kuangching Wang, Glenn Ricart, Larry Landweber, Chip Elliott, Michael Zink, Emmanuel Cecchet, Snigdhaswin Kar, and Prabodh Mishra. The design and operation of CloudLab. In *Proceedings of the USENIX Annual Technical Conference (ATC)*, pages 1–14, July 2019.
- [16] Jeff Epstein, Andrew P. Black, and Simon Peyton-Jones. Towards haskell in the cloud. In *Proceedings of the 4th ACM Symposium on Haskell*, Haskell '11, pages 118–129, New York, NY, USA, 2011. ACM.
- [17] Hamed Esfahani, Jonas Fietz, Qi Ke, Alexei Kolomiets, Erica Lan, Erik Mavrinac, Wolfram Schulte, Newton Sanches, and Srikanth Kandula. Cloudbuild: Microsoft's distributed and caching build service. In *SEIP*. IEEE - Institute of Electrical and Electronics Engineers, June 2016.
- [18] Jim Garlick. pdsh. <https://github.com/chaos/pdsh>, 2022. [Online; accessed September 15, 2022].
- [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] Inc. GitHub. The 2021 state of the octoverse: Top languages over the years. <https://octoverse.github.com/#top-languages-over-the-years>, 2021. [Online; accessed June 1, 2022].
- [21] Google. Protocol Buffers, 2022. Accessed: 2022-06-01.
- [22] Michael Greenberg. libdash. <https://github.com/mgree/libdash>, 2019. [Online; accessed December 6, 2021].
- [23] Michael Greenberg and Austin J. Blatt. Executable formal semantics for the posix shell: Smoosh: the symbolic, mechanized, observable, operational shell. *Proc. ACM Program. Lang.*, 4(POPL):43:1–43:30, January 2020.
- [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*, HotOS '21, page 104–111, New York, NY, USA, 2021. Association for Computing Machinery.
- [25] Hadoop. Hadoop streaming. <https://hadoop.apache.org/docs/r1.2.1/streaming.html>, 2022. [Online; accessed September 15, 2022].
- [26] Shivam Handa, Konstantinos Kallas, Nikos Vasilakis, and Martin C. Rinard. An order-aware dataflow model for parallel unix pipelines. *Proc. ACM Program. Lang.*, 5(ICFP), aug 2021.
- [27] Jason Hickey and Aleksey Nogin. Omake: Designing a scalable build process. In Luciano Baresi and Reiko Heckel, editors, *Fundamental Approaches to Software Engineering*, pages 63–78, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg.
- [28] Martin Hirzel, Robert Soulé, Scott Schneider, Buğra Gedik, and Robert Grimm. A catalog of stream processing optimizations. *ACM Computing Surveys (CSUR)*, 46(4):46:1–46:34, March 2014.
- [29] Lluís Batlle i Rossell. *tsp(1) Linux User's Manual*. <https://vicerveza.homeunix.net/viric/soft/ts/>, 2016.
- [30] 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.

- [31] Dan Jurafsky. Unix for poets, 2017.
- [32] Konstantinos Kallas, Tammam Mustafa, Jan Bielak, Dimitris Karnikis, Thurston H.Y. Dang, Michael Greenberg, and Nikos Vasilakis. Practically correct, Just-in-Time shell script parallelization. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*, pages 769–785, Carlsbad, CA, July 2022. USENIX Association.
- [33] Konstantinos Kallas, Filip Niksic, Caleb Stanford, and Rajeev Alur. Diffstream: Differential output testing for stream processing programs. *Proceedings of the ACM on Programming Languages*, 4(OOPSLA):1–29, 2020.
- [34] Nokia Bell Labs. The unix game—solve puzzles using unix pipes, 2019. Accessed: 2020-03-05.
- [35] Haoyuan Li. *Alluxio: A virtual distributed file system*. University of California, Berkeley, 2018.
- [36] Aurèle Mahéo, Pierre Sutra, and Tristan Tarrant. The serverless shell. In *Proceedings of the 22nd International Middleware Conference: Industrial Track*, pages 9–15, 2021.
- [37] Konstantinos Mamouras, Caleb Stanford, Rajeev Alur, Zachary G. Ives, and Val Tannen. Data-trace types for distributed stream processing systems. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2019*, pages 670–685, New York, NY, USA, 2019. ACM.
- [38] Chris McDonald and Trevor I Dix. Support for graphs of processes in a command interpreter. *Software: Practice and Experience*, 18(10):1011–1016, 1988.
- [39] Malcolm D McIlroy, Elliot N Pinson, and Berkley A Tague. Unix time-sharing system: Foreword. *Bell System Technical Journal*, 57(6):1899–1904, 1978.
- [40] Christopher Meiklejohn and Peter Van Roy. Lasp: a language for distributed, eventually consistent computations with crdts. In *Proceedings of the First Workshop on Principles and Practice of Consistency for Distributed Data*, page 7. ACM, 2015.
- [41] Jürgen Cito Michael Schröder. An empirical investigation of command-line customization. *arXiv preprint arXiv:2012.10206*, 2020.
- [42] Adrian Mizzi, Joshua Ellul, and Gordon Pace. D’artagnan: An embedded dsl framework for distributed embedded systems. In *Proceedings of the Real World Domain Specific Languages Workshop 2018*, pages 1–9, 2018.
- [43] Sape J Mullender, Guido Van Rossum, AS Tanenbaum, Robbert Van Renesse, and Hans Van Staveren. Amoeba: A distributed operating system for the 1990s. *Computer*, 23(5):44–53, 1990.
- [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. ACM.
- [45] Derek G. Murray, Malte Schwarzkopf, Christopher Snowton, Steven Smith, Anil Madhavapeddy, and Steven Hand. Ciel: A universal execution engine for distributed data-flow computing. In *Proceedings of the 8th USENIX Conference on Networked Systems Design and Implementation, NSDI’11*, pages 113–126, Berkeley, CA, USA, 2011. USENIX Association.
- [46] John K Ousterhout, Andrew R. Cherenson, Fred Douglass, Michael N. Nelson, and Brent B. Welch. The sprite network operating system. *Computer*, 21(2):23–36, 1988.
- [47] Shoumik Palkar and Matei Zaharia. Optimizing data-intensive computations in existing libraries with split annotations. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles, SOSP ’19*, pages 291–305, New York, NY, USA, 2019. ACM.
- [48] Rob Pike, Dave Presotto, Ken Thompson, Howard Trickey, et al. Plan 9 from Bell Labs. In *Proceedings of the summer 1990 UKUUG Conference*, pages 1–9, 1990.
- [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] Chet Ramey. Bash reference manual. *Network Theory Limited*, 15, 1998.
- [51] Richard F Rashid and George G Robertson. *Accent: A communication oriented network operating system kernel*, volume 15. ACM, 1981.
- [52] Marc Rozier, Vadim Abrossimov, François Armand, Ivan Boule, Michel Gien, Marc Guillemont, Frédéric Herrmann, Claude Kaiser, Sylvain Langlois, Pierre Léonard, et al. Overview of the chorus distributed operating system. In *Workshop on Micro-Kernels and Other Kernel Architectures*, pages 39–70. Seattle WA (USA), 1992.
- [53] Jan Sacha, Jeff Napper, Sape Mullender, and Jim McKie. Osprey: Operating system for predictable clouds. In *IEEE/IFIP International Conference on Dependable Systems and Networks Workshops (DSN 2012)*, pages 1–6. IEEE, 2012.

- [54] Scott Schneider, Martin Hirzel, Buğra Gedik, and Kun-Lung Wu. Safe data parallelism for general streaming. *IEEE Transactions on Computers*, 64(2):504–517, Feb 2015.
- [55] Jiasi Shen, Martin Rinard, and Nikos Vasilakis. Automatic synthesis of parallel unix commands and pipelines with kumquat. corr abs/2012.15443 (2021). *arXiv preprint arXiv:2012.15443*, 2021.
- [56] Diomidis Spinellis and Marios Fragkoulis. Extending unix pipelines to dags. *IEEE Transactions on Computers*, 66(9):1547–1561, 2017.
- [57] 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.
- [58] Ole Tange. Gnu parallel—the command-line power tool. *login: The USENIX Magazine*, 36(1):42–47, Feb 2011.
- [59] Dave Taylor. *Wicked Cool Shell Scripts: 101 Scripts for Linux, Mac OS X, and Unix Systems*. No Starch Press, 2004.
- [60] Elixir Core Team. Elixir. <https://elixir-lang.org/>.
- [61] The gRPC Authors. *grpc*, 2018. Accessed: 2019-04-16.
- [62] Eleftheria Tsaliki and Diomidis Spinellis. The real statistics of buses in athens. <https://bit.ly/3s112R5>, 2021.
- [63] Junichi Uekawa. dsh. <https://www.netfort.gr.jp/~dancer/software/dsh.html.en>, 2022. [Online; accessed September 15, 2022].
- [64] 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, EuroSys '21*, page 49–66, New York, NY, USA, 2021. Association for Computing Machinery.
- [65] Nikos Vasilakis, Ben Karel, Yash Palkhiwala, John Sonchack, André DeHon, and Jonathan M. Smith. Ignis: Scaling distribution-oblivious systems with light-touch distribution. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2019*, pages 1010–1026, New York, NY, USA, 2019. ACM.
- [66] Robert Virding, Claes Wikström, and Mike Williams. *Concurrent Programming in ERLANG (2Nd Ed.)*. Prentice Hall International (UK) Ltd., Hertfordshire, UK, UK, 1996.
- [67] Bruce Walker, Gerald Popek, Robert English, Charles Kline, and Greg Thiel. The locus distributed operating system. *ACM SIGOPS Operating Systems Review*, 17(5):49–70, 1983.
- [68] Tom White. *Hadoop: The Definitive Guide*. O’Reilly Media, Inc., 4th edition, 2015.
- [69] 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.
- [70] Gina Yuan, Shoumik Palkar, Deepak Narayanan, and Matei Zaharia. Offload annotations: Bringing heterogeneous computing to existing libraries and workloads. In *2020 USENIX Annual Technical Conference (USENIX ATC 20)*, pages 293–306, 2020.
- [71] 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 in-memory cluster computing. In *Proceedings of the 9th USENIX Conference on Networked Systems Design and Implementation, NSDI’12*, pages 2–2, Berkeley, CA, USA, 2012. USENIX Association.