

On Monitoring Asynchronous Components, Asynchronously

ANONYMOUS AUTHOR(S)

Runtime Monitoring is becoming an important analysis tool for improving software quality. The prevailing opinion within the software development community is that inline monitoring is preferred over outline monitoring, mainly because it leads to lower runtime overheads. This paper argues that software has evolved enough over the last few years to put this commonly-held view into question. We provide a series of qualitative arguments in favour of outline monitoring in the case of component-based distributed software. We also develop an algorithm for the correct outline monitoring of dynamic decentralised systems. Finally we conduct a rigorous analysis of the overheads induced by both inline and outline monitoring over models of component-based systems, which enables us to assess more precisely the overhead discrepancy induced by the two variants of the runtime analysis technique.

Additional Key Words and Phrases: Asynchronous component systems, Decentralised monitoring, Dynamic reconfiguration

1 INTRODUCTION

Software has changed dramatically over the last decades. The rise of the app economy on mobile devices, the widespread use of streaming services, together with the impending wave of IoT, have fundamentally altered the manner in which software is developed, the tasks it is expected to conduct, and the environments in which it is required to execute. In these cases, software runs autonomously, under constrained resources, and in decentralised fashion. Viewed globally, this software is structured as a collection of encapsulated components that are massively replicated [Joutsutis 2007]; they are expected to run without interruptions for days, months or even years, and scale in response to fluctuating circumstances [Garg 2015]. These components interact with one another via asynchronous messaging [Hohpe and Woolf 2003] (e.g. as microservices [Jamshidi et al. 2018]). Invariably, these components are developed by third parties using different technologies.

This landscape poses new challenges to developers. Software is expected to adhere to stringent requirements (e.g. streaming services need to ensure adequate levels of QoS) and increasingly handles sensitive information (e.g. mobile devices access our financial data, and medical implants regulate insulin levels), raising the stakes of understanding how it really behaves at runtime. At the same time, the behaviour of component-based software has become harder to understand and predict. This is due to a number of reasons. First, the proper functioning of a system does not depend solely on that of its individual components, but also relies on the manner in which they are *integrated* with one another; this information is rarely readily available when components are provided by third parties (e.g. webservices or binary libraries), or when their connections are determined at runtime (e.g. dynamic service discovery). Second, the sheer scale and distribution of said software further complicates the acquisition and comprehension of this information. Third, these systems execute in open environments, where they are subject to malicious attacks from adversaries that are hard to model and anticipate statically.

Traditional verification approaches like model checking—conceived for software developed in monolithic fashion with conventional deployment practices—do not apply (at least, not in their present form). They are also bound to suffer from the usual scalability issues for large code-bases. Popular methods such as testing and mocking are also largely ineffective when debugging open, large-scale distributed systems with a multitude of execution paths [Alshahwan et al. 2019; Alvaro et al. 2016; Arora et al. 2016]. Rather, these factors have increased the need to complement the

2020. 2475-1421/2020/1-ART1 \$15.00

<https://doi.org/>

50 verification analysis carried out at the design and development phases with validation at the *post-*
51 *deployment* phase. For instance, techniques that traditionally scale well, such as type-systems, have
52 evolved to support the integration of dynamic analyses [Siek and Taha 2007; Takikawa et al. 2012].
53 Unfortunately, these technologies are inherently language-specific and, presently, are not mature
54 enough to cope with software developed using multiple programming languages.

55 Interestingly, industry has witnessed a profusion of tools that enable the *observation and mon-*
56 *itoring* of such systems *at runtime*. These technologies are broadly classified under Application
57 Performance Monitoring tools (APMs) [Heger et al. 2017]. They include commercial solutions such
58 as DataDog, Instana and New Relic One, platform-specific frameworks such as inspectIT-Ocelot (a
59 JVM Agent) and WombatOAM (for Erlang/OTP), and open source offerings such as Dapper, and
60 Zipkin. These tools extend traditional profilers to support distributed tracing and telemetry, log
61 aggregation, data storage, processing and presentation, anomaly detection and threshold-violation
62 alerting, root cause isolation, and also automation for runtime system adaptation. APMs are used
63 extensively for maintenance and performance tuning to identify hotspots and reduce bottlenecks;
64 they presently have an edge on static analysis tools for critical-path analysis and unearthing per-
65 formance anti-patterns [Smith and Williams 2001, 2002]. Reported load-time errors and statistics
66 on end-to-end response times are used to improve user experience. The tracing of events such as
67 exceptions and process failures is used for debugging (live or offline), whereas audit trails are used
68 for forensic analysis in the case of security breaches. APMs may also turn program information
69 that used to be ephemeral and uncertain into something that is concrete and analysable through
70 Machine Learning technologies.

71 The verification counterpart to APMs is Runtime Verification (RV) [Bartocci et al. 2018], where
72 executable code is synthesised from formal specifications to observe the behaviour of a running
73 system against said specifications. Although there is a clear case for using RV for decentralised
74 and distributed scenarios [Francalanza et al. 2018; Sánchez et al. 2019] there is one fundamental
75 difference between present-day distributed RV approaches and APMs. Concretely, most APMs
76 operate as *external entities*, running *asynchronously* to the system under scrutiny (SuS) to analyse its
77 behaviour via intermediaries such as log files and data warehouses. By treating the SuS as a black-
78 box, APMs become largely programming-language agnostic. Moreover, by operating externally,
79 APMs provide added assurances that their monitoring does not directly interfere with the execution
80 of the SuS. In contrast, the state-of-the-art in decentralised and distributed RV [Colombo et al. 2009;
81 El-Hokayem and Falcone 2020; Jin et al. 2012; Kim et al. 2001; Reger et al. 2015; Sen et al. 2004,
82 2006] is dominated by tools that still runs synchronously to the SuS, typically using weaving via
83 code injection (inlining). One reason for this is that most efforts are extensions of mature tools that
84 were originally developed for local, single-threaded RV. There, inlining is the preferred method
85 of instrumentation [Bartocci et al. 2018] because it yields lower overheads [Cardoso et al. 2017,
86 2016]; seminal work in security also highlights the advantages inline instrumentation begets when
87 analysing insecure software [Erlingsson 2004; Erlingsson and Schneider 1999]. However, inlining
88 and synchronous instrumentation may not necessarily be the best approach for decentralised and
89 distributed monitoring. For instance, inlining relies on assumptions, such as full access to the SuS
90 source code, that may not always be possible in this setting; inlining is also programming-language
91 dependent and difficult to administer on heterogenous distributed systems; it is also more intrusive
92 and harder to undo once a properly running SuS is attained.

93 Despite the fact that low overheads are a central concern for any monitoring system, this paper
94 contends that the prevailing view about inline and outline monitoring warrants revisiting. To this
95 end, we present a detailed study of asynchronous monitoring, where instead of considering the
96 analysis aspect of the problem (see [Francalanza et al. 2018] for a detailed survey), we focus on the
97 instrumentation part that determines how the runtime analysis hooks onto the running system. To
98

our knowledge, this aspect of RV has seldom been studied in its own right, even though it probably contributes more to runtime overheads than the runtime analysis itself. Concretely, we make the following contributions:

- (1) We detail an algorithm for concurrent asynchronous monitoring that scales in line with a SuS that grows and shrinks. We make minimal assumptions on the operational model to ensure that our algorithm is sufficiently general to be instantiated in a variety of languages and technologies; the algorithm is also agnostic of the runtime analysis carried out, making it applicable for monitoring both functional and non-functional requirements alike (Sec. 3).
- (2) We build models to evaluate outline monitoring *quantitatively*. We use a series of systematic experiments that compare it with inline monitoring, for a selection of typical system loads. Although we confirm that inline monitoring induces lower overhead, we debunk the generally-held assumption that asynchronous monitoring is necessarily infeasible. We are unaware of similar results within the scientific community (Sec. 4).

2 BACKGROUND

Online monitors analyse the execution of the SuS while it is running. The analysis typically moves forwards in time, discarding already processed events of interest to keep monitors as lightweight as possible. Depending on the monitoring application, trace events can be analysed within the system itself at the point they occur (inline monitors), or transmitted to an external entity that performs the analysis without the system (outline monitors). Inline monitors necessarily execute in synchronous fashion with the SuS. Outline monitors typically execute asynchronously in a separate thread, but can also run synchronously when sharing a common thread with the SuS, or in lock-step via a handshake communication protocol.

The nature of overheads. Inlining determines *statically* the points in the system where the events of interest occur, and the monitor instructions are injected accordingly; these segments lay dormant and get activated only when certain execution paths are followed. Outline instrumentation defers this decision until runtime. In order to scale dynamically, it needs to analyse every event generated by the SuS to determine whether to instrument additional monitors; this is clearly more flexible but also more expensive.

Intrusiveness. Code injection via inlining relies on elements of the program structure. In Object Orientated (OO) programs, where the unit of decomposition is the object, code weaving patterns like Aspect-Oriented Programming (AOP) [Kiczales et al. 1997] package *aspect* code into objects that interact with other system objects through method invocation. In concurrent paradigms, systems are structured as *independent* process units such as actors that interact either via asynchronous message passing or via synchronous mechanisms (e.g. channels or locks). In either case, the interaction between processing entities is determined *at runtime* by the scheduler. This complicates the task of inlining monitor code since this code has to account for these interactions; this becomes even harder to manage when inlined monitors themselves interact as well. Outlining naturally keeps monitors and system processes separate, reducing the risks of subtle bugs from occurring when runtime monitors are introduced.

Separation of Concerns. The separation of monitors and system processes as distinct computational units (induced by outlining) adheres better to established software engineering principles. Inlined monitors are sometimes perceived as functionality that can be aspectised in order to organise the system and monitor code at the software design level. This separation however, does not permeate down to the runtime level, since both system and monitor code executes on the *same* thread. A dependency is created between the two, such that if a monitor embedded in one system process

148 crashes, so does the process; the reverse is also possible, and the runtime analysis is lost. Moreover,
 149 inlined monitors are harder to remove or disable in a running system once weaved.

150 *Flexibility.* There are cases where monitors cannot be inlined because code injection is not possible.
 151 In a setting such as ours, certain components might be offered as-a-service or in the form of
 152 a commercial library where code modifications are prohibited due to availability or licensing
 153 agreements. Obfuscated third-party libraries, while possible to reverse engineer [Chen and Chen
 154 2006], are be hard to instrument, and this certainly cannot be accomplished without intimate
 155 knowledge of the decompiled binary instructions. In instances where it can be done, monitor
 156 inlining generally demands a *redployment* of the instrumented system components, which could
 157 be infeasible for long-running systems. Outline monitoring often relies on *tracing* as a mechanism
 158 to acquire runtime information about the SuS. One advantage that many tracing frameworks offer is
 159 the capability of dynamically switching tracing on or off *without* the need to recompile or redeploy
 160 the traced system. Tracing is typically intended for use in production due to the minimal levels
 161 of overhead it induces; this also makes it an invaluable tool when it comes to detect and analyse
 162 problems that occur in real-time. Many programming language frameworks come equipped with
 163 tracing mechanisms that can be configured programmatically (e.g. Erlang). There are also tracing
 164 frameworks such as and LTng [Desnoyers and Dagenais 2006] and DTrace [Cantrill 2006; Cantrill
 165 et al. 2004] that work at the operating system level; DTrace for instance, also supports tracing at
 166 the application (e.g. MySQL, Firefox) and programming language levels (e.g. C, Java, Erlang, etc.).
 167

168 3 DECENTRALISED OUTLINE MONITORING

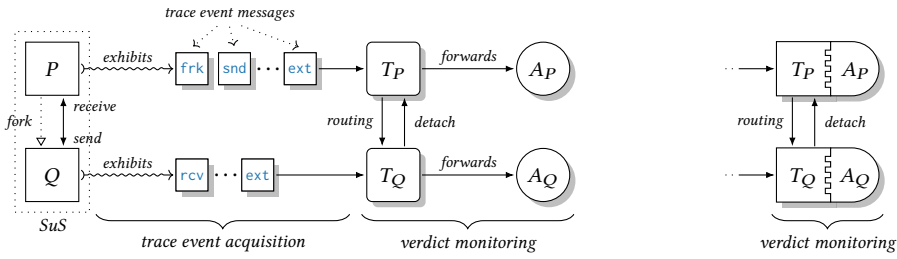
169 We present an outline monitoring algorithm to analyse the behaviour of the SuS by observing its
 170 components in a *decentralised* fashion. Our solution rests on these general assumptions:

- 171 A_1 *No global clock.* System components are not synchronised through a common clock.
- 172 A_2 *System is dynamic.* The number of system components may fluctuate at runtime.
- 173 A_3 *Messages can be reordered.* This does *not* apply for point-to-point communication: successive
 174 messages between the same source and destination are delivered in the sequence issued.
- 175 A_4 *Communication is reliable.* Messages sent are not tampered with, and communication links
 176 never fail (i.e., message delivery is guaranteed and messages duplication does not arise).
- 177 A_5 *Components are reliable.* Components never fail-stop or exhibit Byzantine failures.

178 Our investigation is scoped to execution-monitors/sequence recognisers [Ligatti et al. 2005;
 179 Schneider 2000] where monitors reach irrevocable verdicts after observing a *finite* sequence of
 180 system trace events [Aceto et al. 2019b]. We want our monitors to abide by the following require-
 181 ments:

- 182 R_1 *Monitoring is passive* and only reacts to SuS events.
- 183 R_2 *Monitoring should minimise interference on SuS execution.*
- 184 R_3 *Monitoring is decentralised* without a central coordinating entity.
- 185 R_4 *Monitoring does not miss events or analyse them out of order.*

186 Monitors are instrumented to run in *asynchronous* fashion, in line with assumption A_1 and what is
 187 normally found in distributed setups; although this is outside the scope of our present investigation,
 188 distribution could be obtained by weakening assumptions A_4 and A_5 . Asynchrony may occasionally
 189 affect timely detections. Assumption A_2 and requirements R_2 and R_3 also call for monitoring to *scale*
 190 *dynamically*, continually reconfiguring its choreography in response to certain events exhibited by
 191 the SuS whilst the runtime analysis is in progress. This complicates outline monitoring substantially,
 192 since it must contend with the potential race conditions that may arise. Requirement R_4 addresses
 193



(a) Tracer and analyser organised into separate processes (b) Tracer and analyser as a single process

Fig. 1. Outline verdict monitoring set-up consisting of tracer and analyser roles

problems caused by assumption A_3 . It is vital for execution-monitors, which are usually sensitive to the temporal ordering of the observed events (e.g. RV, root cause analysis, etc.).

3.1 Overview

We proposed *two* outline monitoring set-ups. The choreography in fig. 1a, consisting of independent *tracer* and corresponding *analyser* processes, teases apart the task of trace event routing and monitor reorganisation, performed by tracers, from the task of trace event examination, effected by the analysers. This separation of concerns favours the single responsibility [Agha et al. 1997; Martin 2013] design approach at the expense of introducing an extra process into the monitoring set-up. By contrast, fig. 1b merges the tracing and analysis tasks to forgo this extra process. Our outline approach assumes the existence of a *tracing mechanism* that provides streams of execution events in the form of *messages* for the components of interest in the SuS. The mechanism also allows tracers to control the tracing configuration dynamically at runtime (see discussion in sec. 2). In fig. 1, trace event messages are shown as issuing from processes P and Q and directed to their respective tracers T_P and T_Q; these messages are forwarded to monitors A_P and A_Q for analysis in fig. 1a, or analysed directly as in fig. 1b. The tracing portion of our algorithm relies on these assumptions:

A_6 *Tracers cannot share system processes.* A system process can be traced by (i.e., trace event messages are sent to) *at most one* tracer at any point in time.

A_7 *System processes may share tracers.* A tracer may trace more than one system process.

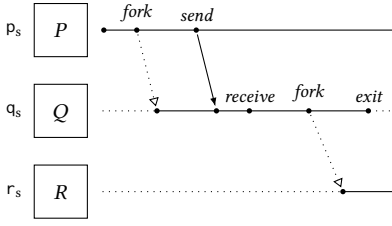
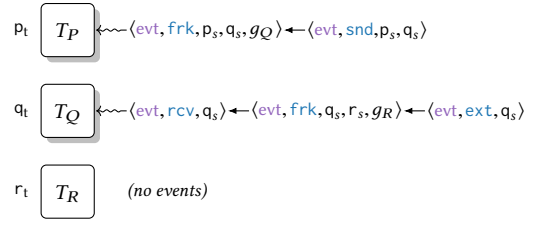
A_8 *System processes inherit tracers.* A system process that is forked by another process that is being traced becomes automatically traced by the *same* tracer.

Assumption A_6 means that for one tracer to start tracing a system process currently being traced, it must first *stop* the active tracer before it can take over and continue tracing this process itself.

3.2 Definitions and Notation

Processes. We assume a denumerable set of process identifiers (PIDs) to uniquely refer to processes. We distinguish between system, tracer and analyser process forms, denoting them respectively by the sets PID_s , PID_t and PID_a , where $p_s \in PID_s$, $p_t \in PID_t$, $p_a \in PID_a$. New processes are created via the function $fork(g)$ that takes the signature of the code to be run by the forked process, $g \in SIG$, returning its *fresh* PID. We refer to the process invoking $fork$ as the *parent*, and to the forked process as the *child*. To create monitor processes, the function $fork$ is overloaded to accept verdict-flagging code, $v \in MON$, and return the corresponding PID p_a ; tracer processes are spawned analogously. Processes communicate with one another through asynchronous messages. Each process is equipped with

246
247
248
249
250
251
252
253

(a) Process creation sequence for P , Q and R (b) Trace event streams for T_P , T_Q and T_R (see tbl. 4)

254
255
256
257
258

Fig. 2. SuS with processes P , Q , and R instrumented with three independent monitors

259
260
261

a *message queue*, K , from where it can read messages *out-of-order* and in a *non-blocking* fashion. Unless stated otherwise, we use the terms *tracer* and *analyser* synonymously since the distinction between the two notions is unimportant for the remainder of this section.

262
263
264

Messages. Messages, $m \in \text{Msg}$, are represented as tuples $\langle q, d_0, d_1, \dots, d_n \rangle$, where q is a *message qualifier* indicating the message type, and $d_{i \in \mathbb{N}}$ are the data elements comprising the message payload. We classify between three messages types, $q \in \{\text{evt}, \text{dte}, \text{rtd}\}$, described thus:

265
266
267

- $q = \text{evt}$: trace events obtained via the tracing mechanism to be analysed;
- $q = \text{dte}$: detach commands that tracers exchange to reorganise the monitoring choreography;
- $q = \text{rtd}$: trace event or command messages that are *routed* between tracers.

268
269
270

The meta-variables e , c , and r are reserved to refer to messages of types `evt`, `dte` and `rtd` respectively. We use the suggestive dot notation (\cdot) to access specific data elements through indexable *field names* (e.g. the message qualifier is accessible through $m.\text{type}$). Trace event messages are structured as $\langle q = \text{evt}, d_0 = a, d_1, \dots, d_n \rangle$, where $a \in \text{ACT}$ identifies the kind of action exhibited by the SuS, and d_1, \dots, d_n designate the data particular to the event. For our exposition, we let $\text{ACT} = \{\text{frk}, \text{ext}, \text{snd}, \text{rcv}\}$, respectively denoting the process actions fork (`frk`), exit (`ext`), send denoted via “!” (`snd`) and receive (`rcv`). We abuse the notation and use a in lieu of the full trace event message data (i.e., q and d_1, \dots) when this simplifies the explanation. The data elements particular to the four trace events are accessed using the field names catalogued in tbl. 4 of app. A.

271
272
273

3.3 The Monitoring Approach

274
275
276

We present our outline decentralised monitoring algorithm incrementally, highlighting the issues that arise when the monitoring choreography reorganises itself as the SuS executes. The algorithm covers *both* arrangements outlined in fig. 1. In the pseudocode, we also highlight the technical differences between the two variants, namely: (i) whether trace events are analysed by the separate analyser, as in fig. 1a, or directly by the tracer, as in fig. 1b, and; (ii) depending on the variant, whether or not the separate analyser is created and terminated. The core logic found in each monitor in a choreography is described in lsts. 1–3; auxiliary logic may be found in app. A. Our exposition focusses on the tracer logic, and is agnostic of the analyser code. Each tracer state comprises of three maps: the *routing map*, Π , describing how to re-route events to other tracers, the *instrumentation map*, Φ , describing when new monitors need to be launched, and a *tracked-processes map*, Γ , recording the system processes the tracer is currently monitoring. We detail how these maps are used below.

277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294

Dynamic process creation. To reorganise the monitoring choreography as the SuS executes, tracers are programmed to react to specific events observed in the trace; in our setting, these are fork (`frk`)

295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343

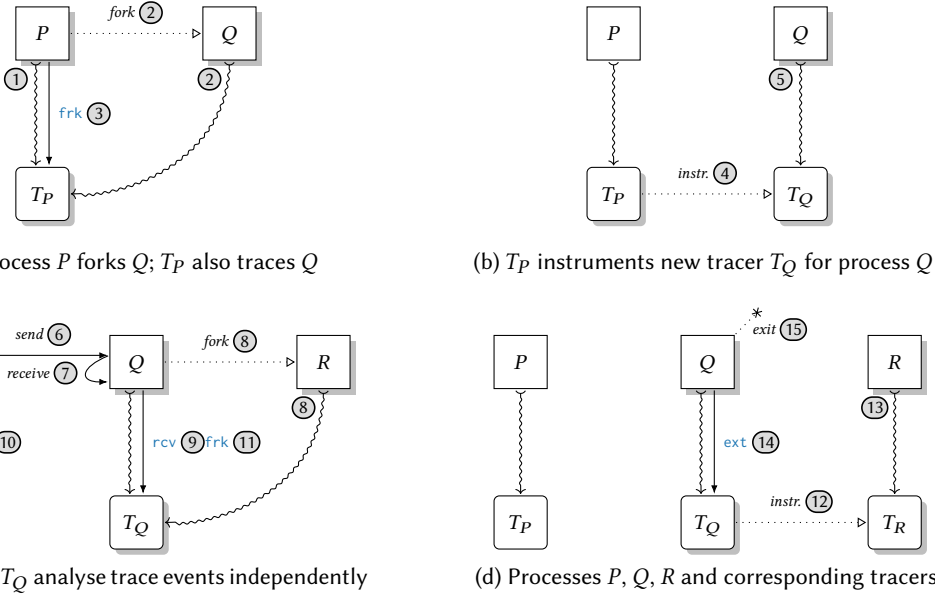


Fig. 3. Outline tracer instrumentation for processes P, Q and Q (analysers omitted)

and exit (`ext`). System processes are typically created in a *hierarchical fashion*, starting from the top-level level process that forks one or more child processes [Armstrong 2007]; we borrow the standard terminology used to describe the relationships between nodes in a tree (*i.e.*, root, ancestor, descendant, *etc.*) when referring to processes. Fig. 2a depicts our running example where the root P forks a child Q and communicates with it; independently Q spawns R and exits. Our example assumes that a dedicated monitor will be assigned per process; our exposition will focus on the tracers, *i.e.*, T_P, T_Q and T_R in this case, where fig. 2b depicts the order of trace events each of these monitors is expected to analyse.

Trace event acquisition. The tracing mechanism alluded to in sec. 3.1 is defined by the operations TRACE, CLEAR and PREEMPT listed in lst. 4 of app. A. TRACE enables a tracer p_t to register its interest in being notified about trace events of a system process p_s . This operation can be undone using CLEAR, which blocks the caller, and returns only when all the trace event messages for p_s that are in the process of being delivered are deposited into the message queue of p_t . PREEMPT combines CLEAR and TRACE, enabling a tracer p_t to take over the tracing of process p_s from another tracer p'_t . Following assumption A_8 , tracing is *inherited* by every child process that a traced system process forks; CLEAR or PREEMPT can therefore be used to alter this arrangement.

Decentralised trace processing. Fig. 3 demonstrates how the process creation sequence of the SuS can be exploited to systematically instrument tracers and evolve the choreography at runtime. The system processes P, Q and R in fig. 3 are created (with PIDs $p_s, q_s,$ and r_s), following the interaction protocol of fig. 2a. A tracer instruments other tracers whenever it encounters fork events in the execution. In fig. 3a, the root tracer T_P analyses the top-level process P , step ①, and instruments a new tracer, T_Q , for process Q when it observes the fork event (`evt, frk, p_s, q_s, g_Q`) exhibited by P in step ③. The field `e.tgt` carried by the fork trace event designates the SuS process that is to be instrumented with the new tracer. Thereafter, T_Q takes over the tracing of process Q by calling

344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392

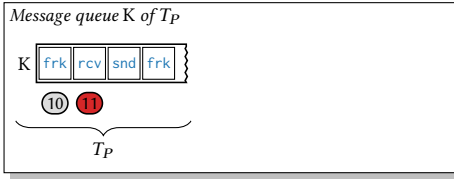
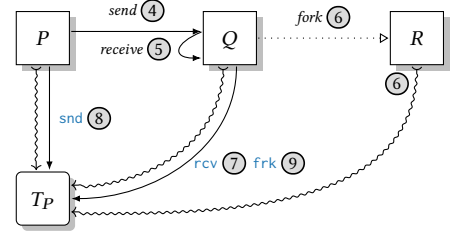
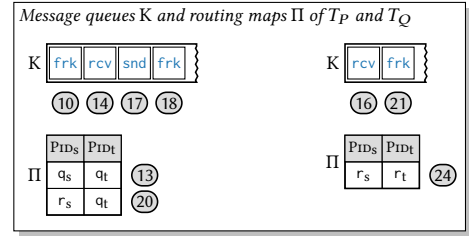
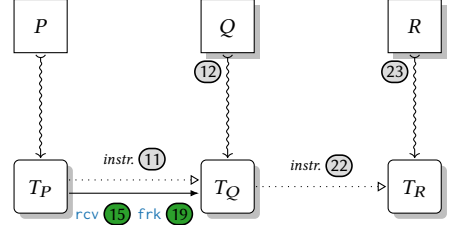
(a) Trace events for P , Q , and R observed by T_P (b) Trace events for Q routed from T_P to T_Q

Fig. 4. Hop-by-hop trace event routing using local tracer routing maps (analysers omitted)

PREEMPT with T_P and $e.src$ to continue tracing Q *independently* of T_P , steps ④ and ⑤ in fig. 3b. The root tracer resumes its own analysis in parallel, where it receives the send event $\langle evt, snd, p_s, q_s \rangle$ in step ⑩ after P issues a message to Q in step ⑥. Subsequent fork events observed by T_P and T_Q are handled in the same manner. Figs. 3c and 3d show how the next tracer, T_R , is instrumented as Q forks its child process R . Recall that prior to the instrumentation of tracers T_Q and T_R , processes Q and R automatically start sharing tracers with their respective parents P and Q when forked, as indicated in steps ② and ⑧.

Trace event routing. Different interleaved executions may still arise for the creation sequence depicted in fig. 2a, due to the asynchrony between the SuS and tracer components. Fig. 4 shows an interleaving alternative to the one captured in figs. 3b–3d. In fig. 4a, the root tracer T_P is slow to handle the fork event exhibited by process Q (step ① in fig. 3a), failing to instrument T_Q promptly. Consequently, in fig. 4a, the trace events due to Q are received by T_P in the sequence indicated by steps ⑦ and ⑨. As a result, the receive event $\langle evt, rcv, q_s \rangle$ is processed by T_P in step ⑩, rather than by the correct tracer T_Q that is eventually instrumented by T_P . This behaviour could derange the runtime analysis, since the events that are expected to be processed by particular analysers unintentionally reach a different monitor.

We address this problem by programming tracers to *filter* the events that are to be analysed locally, and *forward* the rest to other tracers. Fig. 4b shows how the root tracer T_P first instruments T_Q with Q in step ⑪. It subsequently processes the events $\langle evt, rcv, q_s \rangle$ and $\langle evt, frk, q_s, r_s, g_R \rangle$ in steps ⑭ and ⑰, forwarding them to T_Q , steps ⑮ and ⑱. T_Q acts on these events in steps ⑲ and ⑳, where a second tracer, T_R , is instrumented with R . Concurrently, the event $\langle evt, snd, p_s, q_s \rangle$ is *processed locally* by T_P in step ⑰. Trace event routing is accomplished by maintaining a partial map inside tracers, $\Pi: PID_s \rightarrow PID_t$, relating system and tracer PIDs. A tracer queries its instance of the *routing map* Π for every trace event it processes, to determine whether the event should be handled locally or directed elsewhere. The source PID of the event (field $e.src$ in tbl. 4 of app. A) is used to this effect. Trace events are *forwarded* to the tracer with PID p_t when $\Pi(e.src) = p_t$, and


```

393 1 def LOOPo( $\sigma, p_a$ )
394 2   forever do
395   3   # Read routed messages or direct trace events
396   4    $m \leftarrow$  next message from queue K
397   5   if  $m.type = evt$  then
398   6      $\sigma \leftarrow$  HANDLEEVENTo( $\sigma, m, p_a$ )
399   7   else if  $m.type = dtc$  then
400   8     # dtc command received from descendant
401   9     # tracer: route back to sender
402  10     $\sigma \leftarrow$  ROUTEDTC( $\sigma, m$ )
403  11    else if  $m.type = rtd$  then
404  12     $\sigma \leftarrow$  RELAYRTDo( $\sigma, m, p_a$ )
405  13    end if
406  14  end forever
407  15 end def
408
409 13 def HANDLEEVENTo( $\sigma, e, p_a$ )
410 14 if  $e.act = frk$  then
411 15    $\sigma \leftarrow$  HANDLEFORKo( $\sigma, e, p_a$ )
412 16 else if  $e.act = ext$  then
413 17    $\sigma \leftarrow$  HANDLEEXITo( $\sigma, e, p_a$ )
414 18 else if  $e.act \in \{snd, rcv\}$  then
415 19   HANDLECOMMo( $\sigma, e, p_a$ )
416 20 end if
417 21 return  $\sigma$ 
418 22 end def
419
420 23 def HANDLEFORKo( $\sigma, e, p_a$ )
421 24 if ( $p_t \leftarrow \sigma.\Pi(e.src) \neq \perp$ ) then
422 25   ROUTE( $e, p_t$ )
423 26   # New route for events of child process  $e.tgt$ 
424 27   # goes through the same tracer  $p_s$  of its parent
425 28    $\sigma.\Pi \leftarrow \sigma.\Pi \cup \{ \langle e.tgt, p_t \rangle \}$ 
426 29 else
427 30   Monitor  $p_a$  analyses event  $e$ 
428 31    $\sigma \leftarrow$  INSTRUMENTo( $\sigma, e, self()$ )
429 32 end if
430 33 return  $\sigma$ 
431 34 end def
432
433 33 def HANDLEEXITo( $\sigma, e, p_a$ )
434 34 if ( $p_t \leftarrow \sigma.\Pi(e.src) \neq \perp$ ) then
435 35   ROUTE( $e, p_t$ )
436 36 else
437 37   Monitor  $p_a$  analyses event  $e$ 
438 38   # Remove terminated process  $e.src$  from group
439 39    $\sigma.\Gamma \leftarrow \sigma.\Gamma \setminus \{ \langle e.src, \circ \rangle \}$ 
440 40   TRYGC( $\sigma, p_a$ )
441 41 end if
442 42 return  $\sigma$ 
443 43 end def
444
445 43 def HANDLECOMMo( $\sigma, e, p_a$ )
446 44 if ( $p_t \leftarrow \sigma.\Pi(e.src) \neq \perp$ ) then
447 45   ROUTE( $e, p_t$ )
448 46 else
449 47   Monitor  $p_a$  analyses event  $e$ 
450 48 end if
451 49 end def
452
453 Expect:  $\sigma.\Pi(c.tgt) \neq \perp$ 
454 50 def ROUTEDTC( $\sigma, c, p_a$ )
455 51 if ( $p_t \leftarrow \sigma.\Pi(c.tgt) \neq \perp$ ) then
456 52   ROUTE( $e, p_t$ )
457 53    $\sigma.\Pi \leftarrow \sigma.\Pi \setminus \{ \langle c.tgt, p_t \rangle \}$  # Remove route
458 54   TRYGC( $\sigma, p_a$ )
459 55 end if
460 56 return  $\sigma$ 
461 57 end def
462
463 58 def RELAYRTDo( $\sigma, r, p_a$ )
464 59  $m \leftarrow r.emb$ 
465 60 if  $m.type = dtc$  then
466 61    $\sigma \leftarrow$  RELAYDTC( $\sigma, r, p_a$ )
467 62 else if  $m.type = evt$  then
468 63    $\sigma \leftarrow$  RELAYEVT( $\sigma, r$ )
469 64 end if
470 65 return  $\sigma$ 
471 66 end def
472
473 67 def RELAYDTC( $\sigma, r, p_a$ )
474 68  $c \leftarrow r.emb$ 
475 69 if ( $p_t \leftarrow \sigma.\Pi(c.tgt) \neq \perp$ ) then
476 70   RELAY( $r, p_t$ )
477 71    $\sigma.\Pi \leftarrow \sigma.\Pi \setminus \{ \langle c.tgt, p_t \rangle \}$  # Remove route
478 72   TRYGC( $\sigma, p_a$ )
479 73 end if
480 74 return  $\sigma$ 
481 75 end def
482
483 Expect:  $\sigma.\Pi(r.emb.src) \neq \perp$ 
484 76 def RELAYEVT( $\sigma, r$ )
485 77  $e \leftarrow r.emb$ 
486 78 if ( $p_t \leftarrow \sigma.\Pi(e.src) \neq \perp$ ) then
487 79   RELAY( $r, p_a$ )
488 80   # New route for events of child process  $e.tgt$ 
489 81   # goes through the same tracer  $p_s$  of its parent
490 82   if  $e.act = frk$  then
491 83      $\sigma.\Pi \leftarrow \sigma.\Pi \cup \{ \langle e.tgt, p_t \rangle \}$ 
492 84   end if
493 85 end if
494 86 return  $\sigma$ 
495 87 end def

```

Lst. 1. Tracer loop that handles direct events, message routing and relaying

handled by the tracer itself when no such route exists, i.e., $\Pi(e.src) = \perp$. HANDLEFORK, HANDLEEXIT and HANDLECOMM in lst. 1 implement this logic on lines 24, 34 and 44.

A tracer extends its routing map Π whenever it processes a fork event $\langle evt, frk, p_s, p'_s, g \rangle$. It has to consider the following two cases:

Expect: $e.\text{act} = \text{frk}$

```

442 1 def INSTRUMENTo( $\sigma, e, p_t$ )
443 2    $p_s \leftarrow e.\text{tgt}$ 
444 3   if ( $v \leftarrow \sigma.\Phi(e.\text{sig}) \neq \perp$ ) then
445 4      $p'_t \leftarrow \text{fork}(\text{TRACER}(\sigma, v, p_s, p_t))$ 
446 5      $\sigma.\Pi \leftarrow \sigma.\Pi \cup \{ \langle p_s, p'_t \rangle \}$ 
447 6   else
448 7     # In  $\circ$  mode, there is no process  $p_s$  to detach
449 8     # from an ancestor tracer; add  $p_s$  to group
450 9      $\sigma.\Gamma \leftarrow \sigma.\Gamma \cup \{ \langle p_s, \circ \rangle \}$ 
451 10  end if
452 11 return  $\sigma$ 
453 12 end def

```

Expect: $e.\text{act} = \text{frk}$

```

11 def INSTRUMENT•( $\sigma, e, p_t$ )
12    $p_s \leftarrow e.\text{tgt}$ 
13   if ( $v \leftarrow \sigma.\Phi(e.\text{sig}) \neq \perp$ ) then
14      $p'_t \leftarrow \text{fork}(\text{TRACER}(\sigma, v, p_s, p_t))$ 
15      $\sigma.\Pi \leftarrow \sigma.\Pi \cup \{ \langle p_s, p'_t \rangle \}$ 
16   else
17     # Take over  $p_s$  from tracer  $p_t$ ; add  $p_s$  to group
18     DETACH( $p_s, p_t$ )
19      $\sigma.\Gamma \leftarrow \sigma.\Gamma \cup \{ \langle p_s, \bullet \rangle \}$ 
20   end if
21 return  $\sigma$ 
22 end def

```

 Lst. 2. Instrumentation operations for direct (\circ) and priority (\bullet) tracer modes

- C_1 $\Pi(p_s) = \perp$. This means that the tracer needs to adapt the choreography in response to the newly forked process itself. It launches a new (child) tracer $T_{p'}$ with fresh PID p'_t to be instrumented with the forked process p'_s , and extends Π with the mapping $p'_s \mapsto p'_t$; or,
- C_2 $\Pi(p_s) = p'_t$. This means that a route to the tracer with PID p'_t exists for events originating from p_s . Accordingly, the tracer forwards the fork event for p_s to p'_t , and again extends Π with the mapping $p'_s \mapsto p'_t$ (i.e., future events from the new process p'_s will also be forwarded to p'_t).

Fig. 4b depicts the routing maps of tracers T_P and T_Q . T_P adds the mapping $q_s \mapsto q_t$ in step 13, after handling the event $\langle \text{evt}, \text{frk}, p_s, q_s, g_Q \rangle$ to instrument T_Q with Q in steps 10 and 11; this is an instance of case C_1 . Lst. 2 describes function INSTRUMENT where, on line 5, the mapping $e.\text{tgt} \mapsto p'_t$ (with $e.\text{tgt} = p'_s$) is added to Π following the creation of tracer p'_t . Step 20 of fig. 4b constitutes an instance of case C_2 . T_P adds the map $r_s \mapsto q_t$ after processing $\langle \text{evt}, \text{frk}, q_s, r_s, g_R \rangle$ for R , step 18. Crucially, T_P does not instrument a new tracer, but simply delegates this task to T_Q by forwarding the event in question. Lines 26 and 81 in lst. 1 (and later line 26 in lst. 3) are manifestations of this, where the mapping $e.\text{tgt} \mapsto p'_t$ is added after the fork event e is routed to the next tracer p'_t .

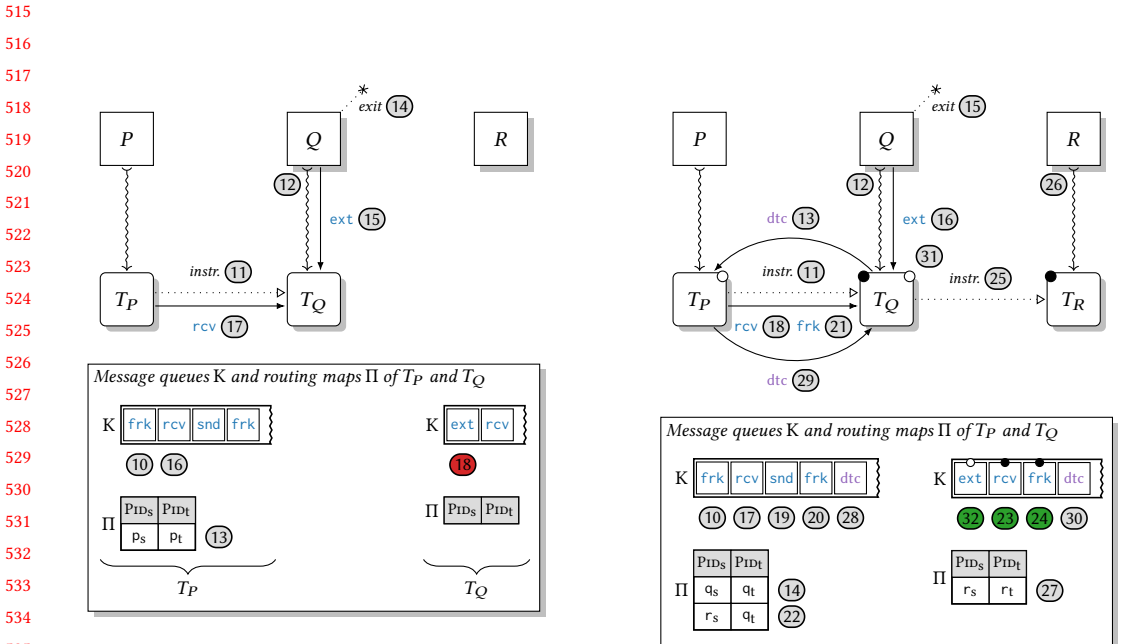
We note that in fig. 4b both mappings inside T_P , created in steps 13 and 20, point to tracer T_Q , and the mapping 24 in T_Q points to T_R . This routing map configuration arises as a result of cases C_1 and C_2 , and implies that any given tracer can *only* forward trace events to adjacent neighbours. For instance, trace events exhibited by R (to be collected by T_P) need to be forwarded twice to reach the intended tracer T_R : from tracer T_P to T_Q , and from T_Q to T_R . This *hop-by-hop routing* [Baker 1995] between tracers forms a connected DAG, and ensures that every message is *eventually delivered* by the tracer choreography. Our algorithm performs routing using two operations, ROUTE and RELAY from lst. 5 in app. A. ROUTE creates a new message, r , with type *rtd* and embeds the message that needs to be routed. Messages routed to a tracer can *either* be analysed or forwarded using RELAY.

Trace event order preservation. Trace event routing does not guarantee that a tracer will receive events in the sequence that should be processed by each monitor, as depicted in fig. 2b, in order to reflect the execution of the SuS shown in fig. 2a. The situation arises when the tracer *simultaneously* actively traces a system component while receiving routed events for that component from another tracer. Fig. 5a highlights the deleterious effect this can have on the runtime analysis should events be deposited out-of-order in the tracer's message queue (assumption A_3). Tracer T_Q takes over tracing process Q from T_P in step 12, and collects the event *ext*, step 15, *before* it receives the routed event *rcv* for Q in step 17. If T_Q naively analyses the events based on their position in the message queue, step 18, it would violate the (correct) order stated in fig. 2b; in fact Q cannot receive a message after exiting. To address this issue, tracers *prioritise* the processing of routed trace event

491 messages. This captures the invariant that out of all events to be analysed by a monitor, routed
 492 events must have temporally preceded all other events. Tracers operate on two levels, *priority* mode
 493 and *direct* mode, denoted by \bullet and \circ in our algorithm. Fig. 5b shows that when in priority mode,
 494 tracer T_Q dequeues the routed events *rcv* and *frk* (labelled with \bullet), and handles them first: *rcv*
 495 is analysed in step 23, whereas *frk* results in the instrumentation of a new tracer T_R in step 25.
 496 Events that should not be analysed by the tracer are forwarded as described earlier in fig. 4. We
 497 note that T_Q can still receive trace events from process Q while this is ongoing, but these events
 498 are only considered once the tracer transitions to direct mode later. Newly-instrumented tracers
 499 default to *priority* mode, so that routed trace events are processed first (see line 5 in lst. 6 of app. A).

500 Lst. 3 shows the priority processing logic, $LOOP_{\bullet}$, where routed trace events are dequeued and
 501 handled (lines 3 and 6). $HANDLEFORK$, $HANDLEEXIT$ and $HANDLECOMM$ for the two tracer modes,
 502 $LOOP_{\circ}$ and $LOOP_{\bullet}$ in lsts. 1 and 3, handle trace events differently. In priority mode, tracers *only*
 503 dequeue routed trace events, and these can be either analysed or relayed (e.g. the branching
 504 statement between lines 24 and 31 in lst. 3). By contrast, tracers in direct mode can relay events
 505 that have been routed their way, *but also* start routing trace events themselves when these are to
 506 be handled by other tracers.

507 *Transitioning safely between tracing modes.* A router tracer is one that currently receives events
 508 from a system process that is configured to be tracked by another tracer; the latter tracer must
 509 be in priority mode. In fig. 4b, T_P is the router tracer for T_Q , since Q (originally set to be traced
 510 by T_Q) shares T_P with process P once forked in fig. 4a, following assumption A_8 . Similarly, T_P is
 511 also the router tracer for T_R . Our tracer choreographies observe the invariant that every tracer in
 512 priority mode has exactly one router tracer. Moreover, if any other tracer along the path between
 513 this tracer and the router tracer is also in priority mode, it must share the same router tracer.
 514



(a) T_Q observes event *ext* before T_P routes *rcv* (b) T_Q processes priority events routed by T_P first

Fig. 5. Trace event order preservation using priority (\bullet) and direct (\circ) tracer modes (analysers omitted)

```

540 1 def LOOP•( $\sigma, p_a$ )
541 2   forever do
542   3   # Trace event messages collected directly are
543   4   # left in the queue to be handled in  $\circ$  mode
544   5    $r \leftarrow$  next rtd message from queue K
545   6    $m \leftarrow r.emb$ 
546   7   if  $m.type = evt$  then
547   8      $\sigma \leftarrow$  HANDLEEVENT•( $\sigma, r, p_a$ )
548   9   else if  $m.type = dtc$  then
549   10    # dtc command routed back from ancestor
550   11     $\sigma \leftarrow$  HANDEDTC( $\sigma, r, p_a$ )
551   12  end if
552   13 end forever
553   14 end def
554
555 12 def HANDLEEVENT•( $\sigma, r, p_a$ )
556 13   $e \leftarrow r.emb$ 
557 14  if  $e.act = frk$  then
558 15     $\sigma \leftarrow$  HANDLEFORK•( $\sigma, r, p_a$ )
559 16  else if  $e.act = ext$  then
560 17     $\sigma \leftarrow$  HANDLEEXIT•( $\sigma, r, p_a$ )
561 18  else if  $e.act \in \{snd, rcv\}$  then
562 19    HANDLECOMM•( $\sigma, r, p_a$ )
563 20  end if
564 21 end def
565
566 22 def HANDLEFORK•( $\sigma, r, p_a$ )
567 23   $e \leftarrow r.emb$ 
568 24  if ( $p_t \leftarrow \sigma.\Pi(e.src)$ )  $\neq \perp$  then
569 25    RELAY( $r, p_t$ )
570 26     $\sigma.\Pi \leftarrow \sigma.\Pi \cup \{ \langle e.tgt, p_t \rangle \}$ 
571 27  else
572 28    # Monitor  $p_a$  analyses event  $e$ 
573 29     $p'_t \leftarrow r.rtr$ 
574 30     $\sigma \leftarrow$  INSTRUMENT•( $\sigma, e, p'_t$ )
575 31  end if
576 32  return  $\sigma$ 
577 33 end def
578
579 34 def HANDLEEXIT•( $\sigma, r, p_a$ )
580 35   $e \leftarrow r.emb$ 
581 36  if ( $p_t \leftarrow \sigma.\Pi(e.src)$ )  $\neq \perp$  then
582 37    RELAY( $r, p_t$ )
583 38  else
584 39    # Monitor  $p_a$  analyses event  $e$ 
585 40    # Remove terminated process  $e.src$  from group
586 41     $\sigma.\Gamma \leftarrow \sigma.\Gamma \setminus \{ \langle e.src, \bullet \rangle \}$ 
587 42    TRYGC( $\sigma, p_a$ )
588 43  end if
589 44  return  $\sigma$ 
590 45 end def
591
592 45 def HANDLECOMM•( $\sigma, r, p_a$ )
593 46   $e \leftarrow r.emb$ 
594 47  if ( $p_t \leftarrow \sigma.\Pi(e.src)$ )  $\neq \perp$  then
595 48    RELAY( $r, p_t$ )
596 49  else
597 50    # Monitor  $p_a$  analyses event  $e$ 
598 51  end if
599 52 end def
600
601 Expect:  $r.emb.iss = self() \vee \sigma.\Pi(r.emb.tgt) \neq \perp$ 
602
603 53 def HANDEDTC( $\sigma, r, p_a$ )
604 54   $c \leftarrow r.emb$ 
605 55  if ( $p_t \leftarrow \sigma.\Pi(c.tgt)$ )  $\neq \perp$  then
606 56    RELAY( $r, p_t$ )
607 57  else
608 58    # Mark process  $c.tgt$  in group as detached
609 59     $\sigma.\Gamma \leftarrow \sigma.\Gamma \setminus \{ \langle c.tgt, \bullet \rangle \}$ 
610 60     $\sigma.\Gamma \leftarrow \sigma.\Gamma \cup \{ \langle c.tgt, \circ \rangle \}$ 
611 61     $\gamma = \{ \langle p_s, d \rangle \mid \langle p_s, d \rangle \in \sigma.\Gamma, d = \bullet \}$ 
612 62    if  $\gamma = \emptyset$  then
613 63      LOOP◦( $\sigma, p_a$ ) # Switch tracer to  $\circ$  mode
614 64    end if
615 65  end if
616 66  return  $\sigma$ 
617 67 end def

```

Lst. 3. Tracer loop that handles priority trace events and message relaying

A tracer in priority mode coordinates with its router tracer to determine whether all of the events for its tracked system processes have been routed. A tracer must effect this procedure for every process it currently tracks, recorded in Γ , before it can safely transition to direct mode, and start operating on the trace events it collects directly. The tracer issues a special detach command message, `dtc`, to notify the router tracer that it is now responsible for tracing a particular system process. Detach commands contain the PIDs of the issuer tracer in priority mode and system process in question, accessed via the fields `iss` and `tgt` respectively, described in tbl. 5 in app. A.

Fig. 5b shows tracer T_Q in priority mode (\bullet) sending the command $\langle dtc, q_t, q_s \rangle$ for Q , step 13, after it starts tracing this process directly, step 12. This transaction is implemented by DETACH in lst. 6 of app. A, where $\langle dtc, p'_t, p_s \rangle$ is sent to the router tracer p_t on line 10, after it invokes PREEMPT. In fig. 5b, `dtc` issued by T_Q follows `rcv` and `frk` in the message queue of tracer T_p . The router tracer processes its message queue sequentially in steps 10, 17, 19, 20 and 28. These trace events are forwarded to neighbouring tracers as necessary in steps 18 and 21 (see lines 3–5 in lst. 1). It also routes the `dtc` command back to the issuer tracer in step 29 where, once handled, marks the system process as *detached* from the router tracer. HANDEDTC in lst. 3 effects this update on the routing

589 map $\Gamma : \text{PID}_s \rightarrow \{o, \bullet\}$ of the issuer tracer on lines 58 and 59. Once all the processes in Γ become
 590 detached, the tracer transitions to direct mode by executing `LOOPo`; this check is performed on lines
 591 60 and 61 in lst. 3. While in priority mode, T_Q handles the prioritised (\bullet) events forwarded by T_P
 592 in the *correct order* stipulated earlier in fig. 2b (steps 23 and 24). This is followed by handling the
 593 command `dtc` in step 30. The transition from priority to direct mode for T_Q in fig. 5b takes place in
 594 step 31. Finally, the trace event `ext` is handled in the correct order in step 32 (as opposed to step
 595 18 in fig. 5a).

596 A detach command c originating at the router tracer may be relayed through multiple intermediate
 597 tracers until it reaches its destination. Every intermediate tracer purges the association $c.tgt \mapsto p_t$
 598 from its routing map Π for some neighbouring tracer PID p_t . This functionality is provided by
 599 `RELAYDTC` and `ROUTEDTC` in lst. 1: despite their similar logic, `ROUTEDTC` is used by the tracer to
 600 commence the routing of detach commands, whereas `RELAYDTC` merely forwards commands to
 601 other tracers. While these steps are not shown in fig. 5b, we briefly remark that tracer T_P would
 602 remove from Π the mapping $q_s \mapsto q_t$, calling `ROUTEDTC` to start routing back the detach command
 603 $\langle \text{dtc}, q_t, q_s \rangle$ it receives from T_Q . In due course, T_P also removes $r_s \mapsto q_t$ for process R once it handles
 604 $\langle \text{dtc}, q_t, r_s \rangle$ sent by tracer T_R . When it receives the routed detach command $\langle \text{rtd}, p_t, \langle \text{dtc}, q_t, r_s \rangle \rangle$ from
 605 T_P , T_Q removes $r_s \mapsto r_t$ from Π and relays it, in turn, to tracer T_R using `RELAYDTC`.

606
 607 *Selective instrumentation.* In practice, one might want to have the flexibility to *group processes* under
 608 a single monitor to analyse them as one *component*. Our algorithm selectively instruments (new)
 609 tracers for particular system processes using the map, $\Phi : \text{SIG} \rightarrow \text{MON}$: it maps the code signatures,
 610 g (of the system process forked), to the monitoring code, v (to be executed by the newly spawned
 611 monitor). `INSTRUMENT` in lst. 2 applies Φ to the code signature, where $e.\text{sig} = g$, in the fork event e
 612 on lines 3 and 13. When $\Phi(g) = \perp$, instrumentation is not performed, and the tracer is automatically
 613 shared by the new process $e.tgt$, according to assumptions A_7 and A_8 .

614
 615 *Garbage collection.* Our outline set-up can shrink in size by discarding tracers that are no longer
 616 needed. A tracer self-terminates after its routing map Π and tracked-processes map Γ become
 617 empty; this check is performed by `TRYGC` in lst. 6 in app. A. The tracer purges process references
 618 from Γ when it handles exit trace events via `HANDLEEXITo` and `HANDLEEXIT\bullet` (lsts. 1 and 3). Note
 619 that, even when $\Gamma = \emptyset$ and the tracer has no processes to analyse, it might still be required to route
 620 trace events to adjacent tracers, *i.e.*, $\Pi \neq \emptyset$. The garbage collection check is therefore performed
 621 each time mappings from Π or Γ are removed on lines 39, 54 and 72 in lst. 1, and line 41 in lst. 3. In
 622 fig. 5b, tracer T_Q would terminate sometime after handling the exit event `ext` for process Q in step
 623 32, once the routed detach command $\langle \text{rtd}, p_t, \langle \text{dtc}, q_t, r_s \rangle \rangle$ it receives from T_P is relayed to tracer T_R .

624 4 EVALUATION

625
 626 We give a comprehensive evaluation to assess quantitative aspects of inline and outline monitoring.
 627 Our evaluation takes the form of a case study that instantiates the monitoring problem from sec. 3 to
 628 a RV setting, where event streams are analysed to reach acceptance/rejection verdicts in connection
 629 to satisfactions/violations of correctness properties [Bartocci et al. 2018; Francalanza et al. 2017].
 630 The set-up follows that of fig. 1, where the analysis components (used in fig. 1 and in inlined
 631 monitors) are synthesised from syntactic descriptions of the properties of interest. Our synthesis
 632 compiles properties down to *automata-based* monitors following [Aceto et al. 2019a]. We evaluate
 633 the different approaches in terms of runtime overheads and, by this, assess their viability. We follow
 634 an approach similar to [Bartocci et al. 2019], and consider the following overhead performance
 635 metrics: (i) mean scheduler utilisation, as a percentage of the total available capacity, (ii) mean
 636 memory consumption, measured in GB, (iii) mean round trip time (RTT), measured in milliseconds

637

638 (ms), and, (iv) mean execution duration, measured in seconds (s). Our measurements are collected
639 *globally* by sampling the runtime environment in which the SuS and monitoring system execute.

640

641

4.1 Scope

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

4.2 Methodology

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

We use Erlang [Armstrong 2007] to implement our evaluation set-up and monitoring algorithms. Erlang adopts the actor model of computation [Agha et al. 1997], implementing them as *lightweight* processes. Actors interact via *asynchronous* messaging, changing their (local) internal state based on messages received. Every actor owns a message queue, called the *mailbox*, where messages can be taken out-of-order. Actors can also *fork* other actors to execute independently in their own process space. Every actor is identified via a PID that is assigned to it when forked. We use the term actor and process interchangeably in the rest of this section.

Implementation. The *inline monitoring* tool we developed for this study assumes access to the SuS source code. It instruments monitoring instructions into the target system via *code injection* by manipulating its parsed abstract syntax tree. The *modified* syntax tree is then compiled into an executable form, and the instrumented instructions perform the runtime analysis in a synchronous manner as the SuS executes. Our implementation of the *outline monitoring* algorithm in sec. 3 maps tracer processes to Erlang actors. Tracers collect the trace events by leveraging the *native* tracing infrastructure provided by the Erlang Virtual Machine (EVM). This infrastructure complies with assumptions A_6 – A_8 . EVM tracing directs trace event messages from system processes to tracer mailboxes acting as the tracer messages queues K of sec. 3.2. The maps Π , Φ and Γ are implemented using Erlang maps for efficient access. We implement the two trace analysis variants of fig. 1. For the arrangement in fig. 1a, the analysis is forked as a separate actor where tracers forward their event messages. Line 4 in lst. 6 of app. A indicates the point at which the actor tasked with the analysis is created whereas line 14 signals said actor to terminate when garbage collection takes place. The analysis is incorporated directly into tracers for the merged monitor case in fig. 1b.

¹We could have employed a peer-to-peer set-up, but this complicates the evaluation considerably.

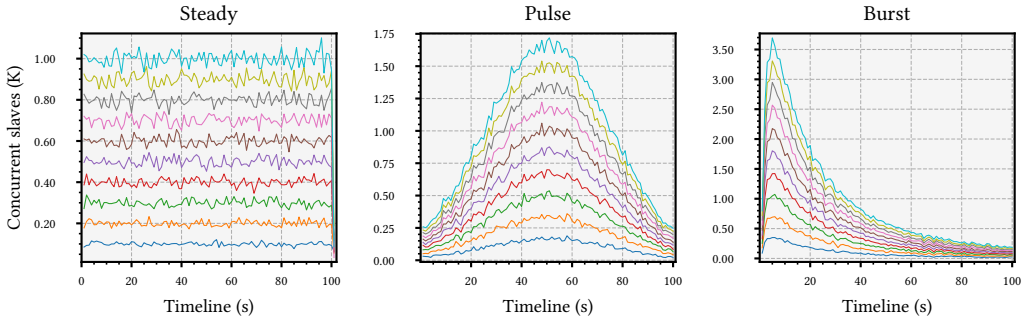


Fig. 6. Steady, Pulse and Burst load distributions with 100 k slaves for the duration of 100 s

The SuS. We opt for a custom-built evaluation platform that emulates models of master-slave systems. The decision *not* to go with off-the-shelf (e.g. web servers, thread pools, etc.) systems stems from three core drawbacks these have, namely: (i) they make it challenging to *precisely control* particular experiment parameters conducive to *repeatable results*, (ii) do not provide hooks that permit *accurate* measurement taking, (iii) often embody *highly-specific* use cases that make it difficult to generalise the findings obtained. Our evaluation platform is *parametrisable* to emulate different system models. The *tasks* farmed out by the master consist of *work requests* that a slave receives, processes and echoes back. A slave is set to terminate once all of its work requests have been handled and acknowledged by the master. The parameter w in our framework regulates the number of work requests that can be batched in one task; the *actual* amount of work requests *per slave* is drawn randomly from a normal distribution with mean $\mu = w$, and a standard deviation $\sigma = w \times 0.02$. This ensures a degree of variability in the amount of messages exchanged between the master and each slave. The speed with which the system reacts to work requests can be controlled via the parameters $\text{Pr}(\text{send})$ and $\text{Pr}(\text{recv})$. To distribute tasks uniformly amongst slaves, the master interleaves the sending and receiving of work requests: $\text{Pr}(\text{send})$ and $\text{Pr}(\text{recv})$ can bias this behaviour. $\text{Pr}(\text{send})$ determines the probability that a work request is assigned by the master to a slave. $\text{Pr}(\text{recv})$ controls the probability that a work request received by the master is handled and acknowledged. Load on the system is induced by the master when it creates slave processes; the total number of slaves that are created during one experiment is set using the parameter n .

Load models. Our system considers three load shapes (fig. 6) that establish how the creation of slaves based on the parameter n is distributed along the *load timeline* t . The load timeline is represented as a sequence of discrete *logical time units* that denote instants at which a new set of slaves is created by the master. *Steady* loads reproduce executions where a system operates under stable conditions. These are modelled on a homogeneous Poisson distribution with *rate* λ , specifying the mean number of slaves that are created at every time instant along the load timeline with duration $t = \lceil n/\lambda \rceil$. *Pulses* emulate scenarios where a system undergoes gradually-increasing load peaks. The pulse load shape is parametrised by t and the spread, s , that controls how slowly or sharply the system load increases as it approaches its peak halfway along t . Pulses follow a normal distribution with $\mu = t/2$ and $\sigma = s$. *Burst* loads capture scenarios where a system is stressed due to instant load spikes: these are based on a log-normal distribution with $\mu = \ln(m^2/\sqrt{p^2 + m^2})$, $\sigma = \sqrt{\ln(1 + p^2/m^2)}$ where $m = t/2$ and p is the pinch controlling the intensity of the initial load burst.

Experiment set-up. To meet the objectives set out in this section, we conduct two case studies where the SuS is configured with $n = 10\text{k}$ for *moderate* loads and $n = 100\text{k}$ *high* loads. The number of

work requests per task is set to $w = 100$. $\Pr(\text{send}) = \Pr(\text{recv}) = 0.9$ fixes the probability of sending and acknowledging work requests: this emulates a system that reacts promptly to load, but at the same time, exhibits slight processing delays that arise in a master-slave architecture. Our chosen parameter values instantiate the SuS to model *realistic* web response time where the request intervals observed at the server follow a Poisson process [Ciemiewicz 2001; Kayser 2017; Liu et al. 2001]. Further detail regarding the validation of this model are given in app. B. For these experiments, the total loading time is set to $t = 100$ s. We use the term *experiment* to denote a series of ten *benchmarks* where the SuS is configured with one particular monitoring set-up (e.g. with outline monitors). Load to the set-up is added *incrementally* at each benchmark until the maximum load is reached, e.g. for the case study with $n = 10$ k slaves, we start with the first benchmark set to $n_1 = 1$ k and progress to $n_{10} = 10$ k in steps of 1 k. We repeated *ten* readings for each experiment, and aggregated the results by computing the weighted mean for the performance metrics mentioned above. Consult app. B for the full list of precautions. The experiments were conducted on an Intel Core i7 M620 64-bit machine with 8GB of memory, running Ubuntu 18.04 and Erlang/OTP 22.2.1.

4.3 Results and Discussion

Our results are reported in tbls. 1 and 2 and figs. 7–10, plotting each performance metrics (y -axis) against the slave processes (x -axis) for every monitoring mode; the *unmonitored* system is inserted as a baseline reference. Fitted data plots corresponding to figs. 7–10 are given in app. C.

Moderate loads. Our first batch of results considers loads that are slightly higher than those employed by the state-of-the-art to evaluate decentralised, concurrent and distributed runtime monitoring, e.g. [Attard and Francalanza 2017; Berkovich et al. 2015; Cassar and Francalanza 2016; Colombo and Falcone 2016; El-Hokayem and Falcone 2017; Francalanza and Seychell 2015; Mostafa and Bonakdarpour 2015; Neykova and Yoshida 2017a,b; Scheffel and Schmitz 2014]; works like [Chen and Rosu 2007, 2009; Reger et al. 2015] consider higher loads, but they evaluate sequential monitoring. Crucially, neither of the aforementioned studies employ different load shapes in their analysis: they either use loads modelled on a Poisson process, *i.e.*, Steady load, or do not specify the load types considered. The SuS set with $n = 10$ k slaves and $w = 100$ work requests per slave generates $\approx n \times w \times (\text{work requests and responses}) = 2$ M messages exchanged between the master and slaves, producing $2M \times (\text{snd and rcv trace events}) = 4$ M trace events. Tbl. 1 reports the percentage overhead at $n = 10$ k. It shows that inline and the two flavours of outline monitoring induce negligible execution slowdown for all three load shapes (e.g. 0.77% maximum for DM under Burst load); the memory consumption overhead behaves similarly. At the Steady load illustrated in fig. 7, memory consumption and RTT (round trip) grow linearly in the number of slave processes. Under the Pulse and Burst loads in tbl. 1, inline monitoring induces negligible scheduler overhead. This is markedly higher for outline monitoring (DS and DM), mostly caused by the dynamic reconfiguration of the

Filler	Steady			Pulse			Burst		
	I	DS	DM	I	DS	DM	I	DS	DM
Scheduler	1.68	42.13	22.60	1.54	35.17	18.12	2.08	38.92	26.03
Memory	0.01	0.10	0.09	0.01	0.08	0.04	0.03	0.10	0.06
RTT	4.37	67.05	58.36	7.72	82.85	60.79	20.17	859.91	666.46
Execution	0.09	0.40	0.28	0.10	0.32	0.22	0.12	1.09	0.77

Inline (I), Decentralised outline separate (DS), Decentralised outline merged (DM)

Tbl. 1. Percentage runtime overhead taken at the maximum load $n = 10$ k

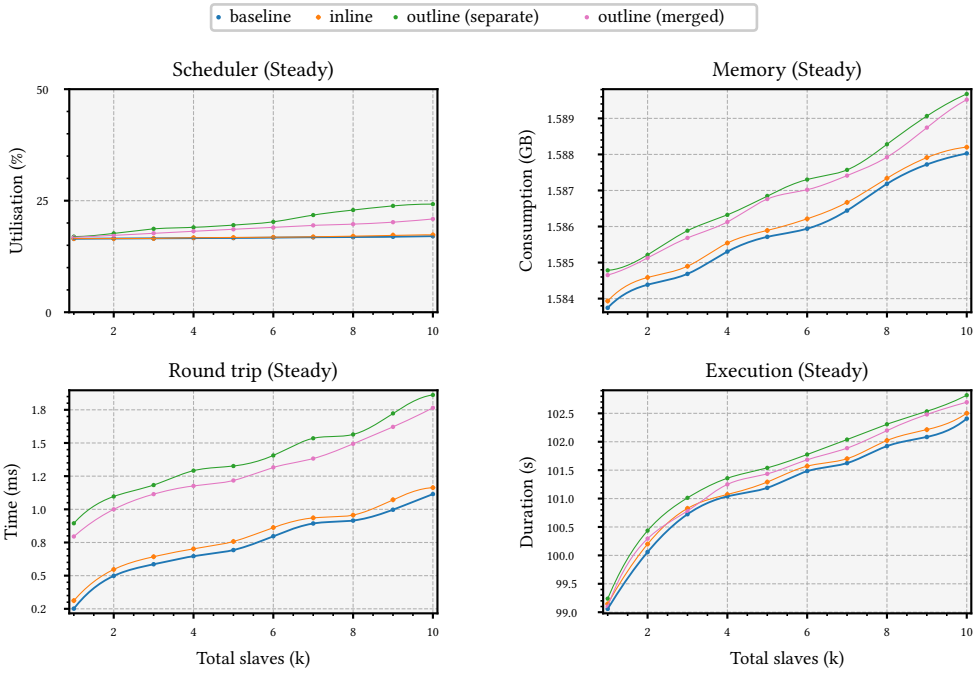


Fig. 7. Mean runtime overhead for monitoring the master and slave processes (10 k slaves)

monitoring choreography. Tbl. 1 also suggests that the RTT is very sensitive to the type of load applied, and it increases for the load shapes Steady, Pulse and Burst respectively. In fact, the latter load shape induces a sharp growth in the RTT for outline monitoring at around 9k ~ 10k slaves, as illustrated in fig. 8. This indicates that specific load shapes prompt very different behaviours from the monitors, and should be taken into account.

Despite the clear discrepancies (percentage-wise) in scheduler and RTT overheads between inline and outline monitoring in tbl. 1, these are *comparable* (value-wise) for the loads that are typically used in other bodies of work, as shown in figs. 7 and 8 (e.g. the worst discrepancy is 11ms for RTT under Burst). Merging the analysis with tracing as in fig. 1b yields improvements, but its effect is negligible. For certain performance metrics, our data plots do not allow us to confidently extrapolate our results. A case in point is the RTT Burst plot for outline monitoring, which raises the question of whether the trend remains consistent when the number of slaves exceeds 10k.

High loads. We increase the number of slaves to $n=100k$ and keep $w=100$, to generate 20 M messages and 40 M trace events. Our aim is to assess how the monitored system performs *under stress*, and whether this reveals aspects that do not emerge at lower loads. Since these loads span a broader range, this also gives us a reasonable level of confidence when extrapolating our observations. Particular, we also include the measurements obtained for the SuS with a *centralised* monitoring set-up for this case study, to better isolate the effects of outline monitoring.

Tbl. 2 confirms that inline monitoring induces lower overheads. However, dissecting these results uncovers a few surprising aspects. For instance, the memory overhead between inline and outline monitoring with the separate analysis is 13.3 % under a Steady load at its highest point of 100 k slaves. This overhead is arguably tolerable for a number of applications. When merging outline tracing with its analysis as in fig. 1b, this discrepancy goes down to respectable 6.8 %. Centralised

834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882

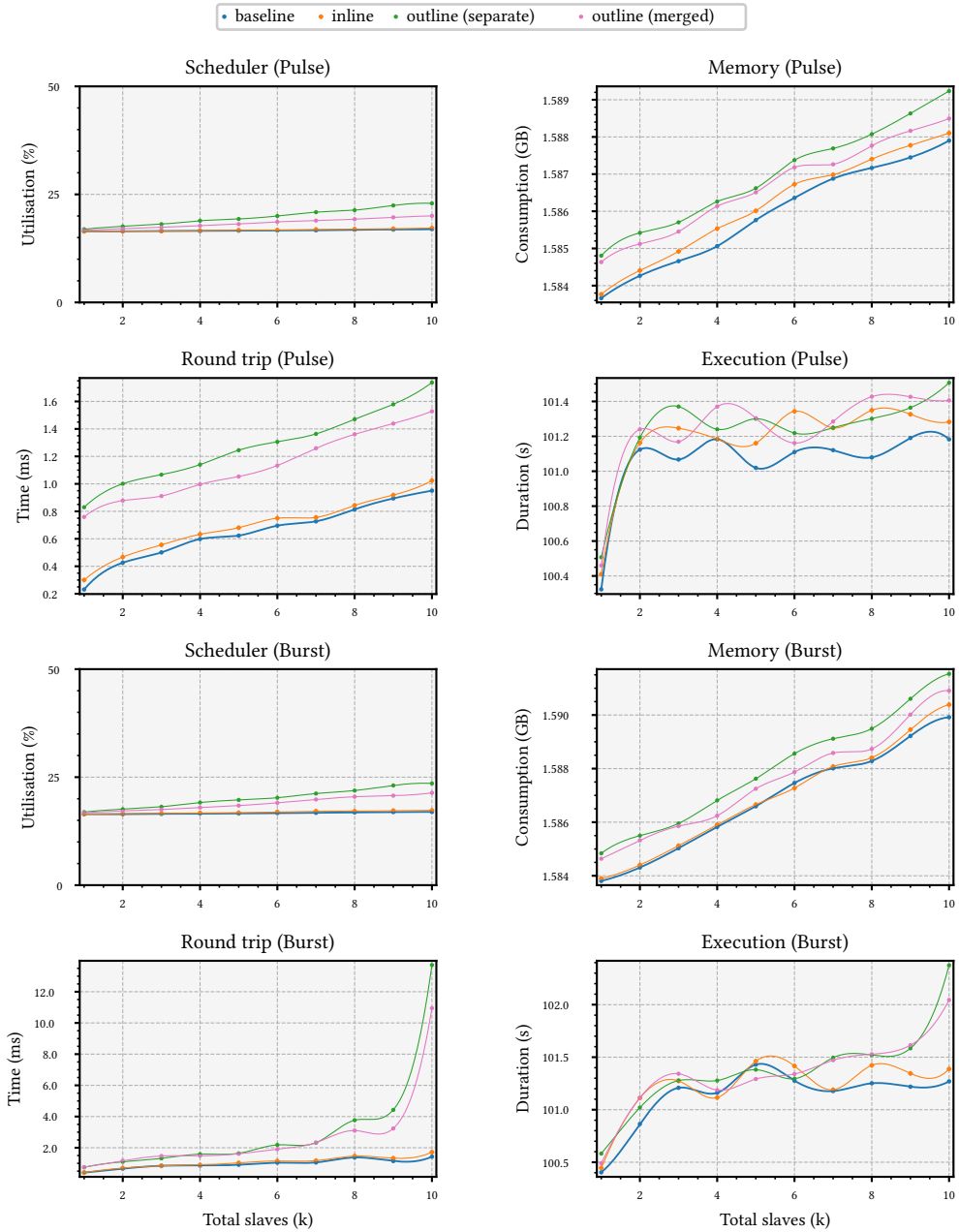


Fig. 8. Mean runtime overhead for monitoring the master and slave processes (10 k slaves, cont.)

outline monitoring further lower this difference to a negligible 0.6 %; this seems to debunk the general assumption that outline monitoring necessarily leads to infeasible overheads. The plots in fig. 9 also show that under Steady loads, the overhead for the memory, RTT and execution duration are comparable to inline monitoring up to the considerable load of around 40 k slaves (*i.e.*, 8 M

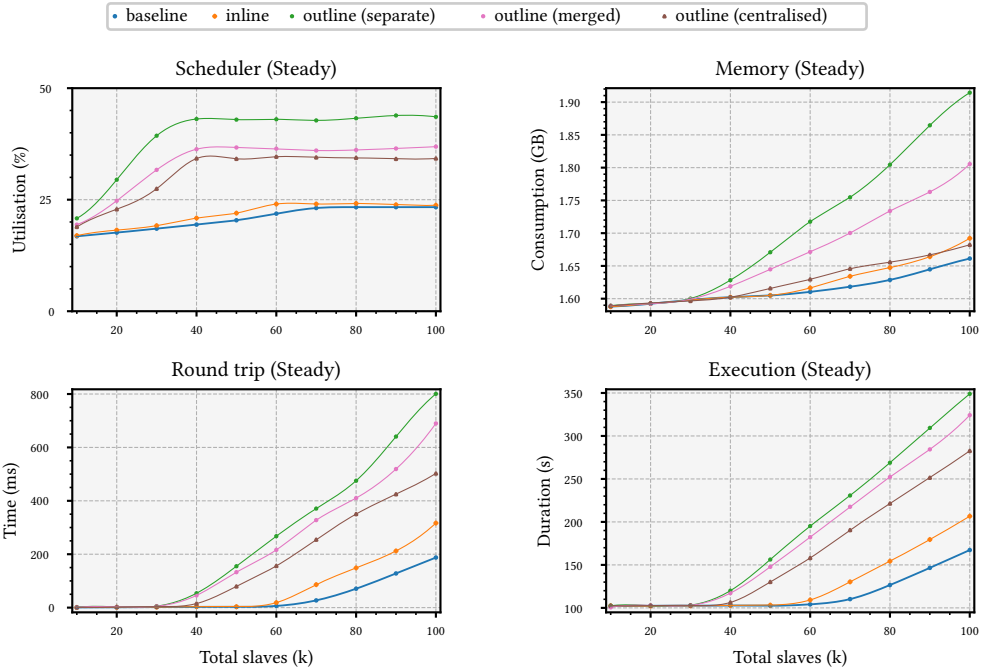


Fig. 9. Mean runtime overhead for monitoring the master and slave processes (100 k slaves)

messages and 16 M trace events). Tbl. 2 indicates that the RTT overhead for outline monitoring *decreases* for the load shapes Steady, Pulse and Burst respectively; this trend is also mirrored in the execution slowdown metric. These results contrast with the ones in tbl. 1, where the overhead for said metrics *gradually increases* under the same load shapes. This suggests that outline monitoring exhibits a degree of robustness at high numbers for loads like Pulse and Burst, whose shapes induce higher stress in the SuS in comparison to consistent loads. In these two instances, outline monitoring pays a price in terms of memory overhead (tbl. 1), although the *maximum* overhead reported in our results, 24.6 %, may be acceptable for many scenarios. Merging the analysis with tracing in outline monitoring *consistently* yields lower overheads when compared to the variant with separate analysis, irrespective of the load shape.

We draw attention to the charts in figs. 9 and 10, where the memory consumption plot for centralised outline monitoring crosses over that of inline monitoring. This behavior emerges

	Steady				Pulse				Burst			
	I	DS	DM	C	I	DS	DM	C	I	DS	DM	C
Scheduler	1.8	86.8	58.1	46.8	2.9	85.5	55.6	47.5	3.1	84.4	57.9	51.6
Memory	1.9	15.2	8.7	1.3	2.9	18.1	11.8	2.2	3.1	24.6	15.4	1.9
RTT	68.9	326.9	267.9	167.9	72.7	257.8	238.7	189.0	28.4	120.6	114.3	100.1
Execution	23.5	108.6	93.8	68.9	24.5	101.8	93.5	72.8	15.7	82.0	77.5	69.3

Inline (I), Decentralised outline separate (DS), Decentralised outline merged (DM), Centralised outline (C)

Tbl. 2. Percentage runtime overhead taken at the maximum load $n = 100k$

932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980

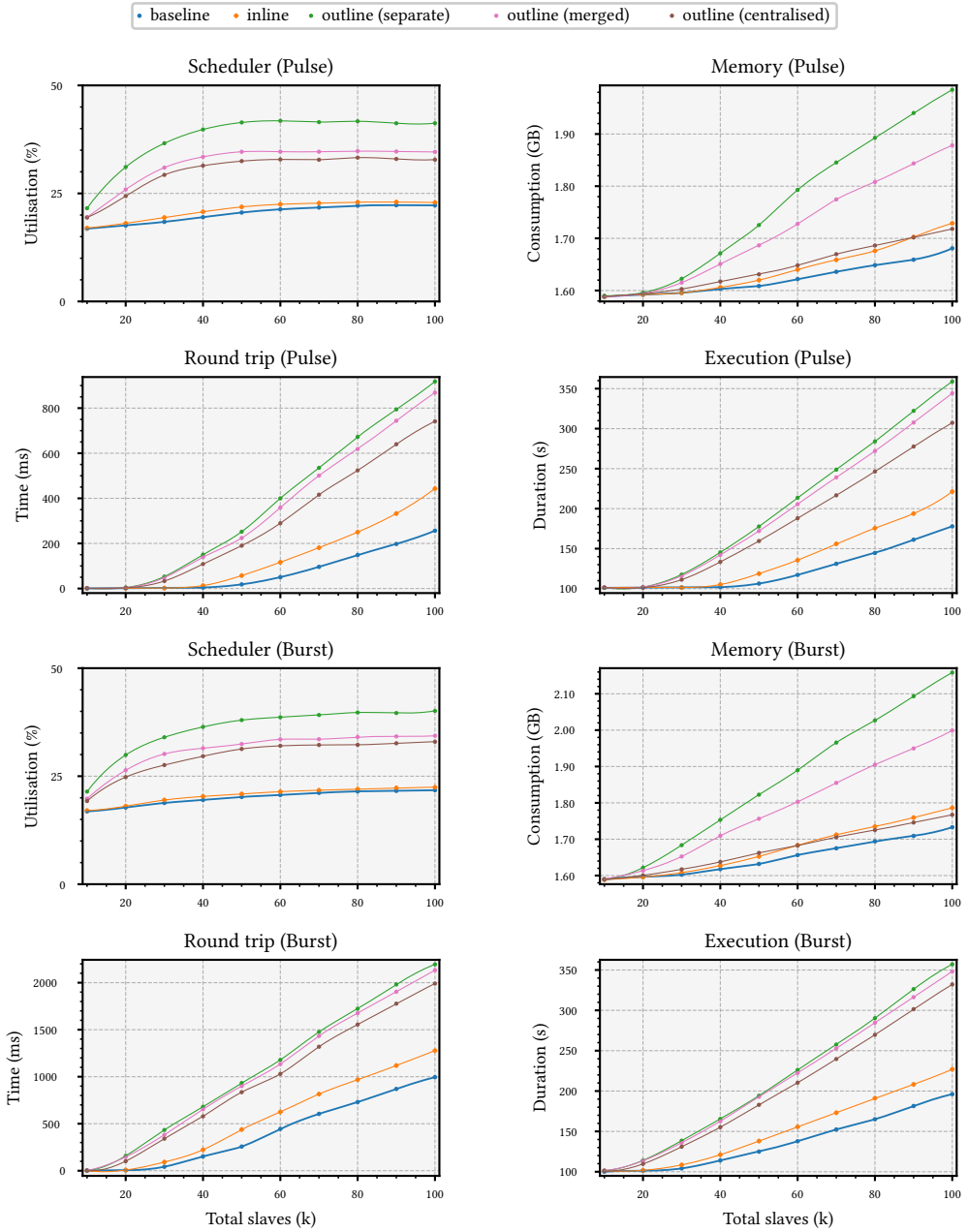


Fig. 10. Mean runtime overhead for monitoring the master and slave processes (100 k slaves, cont.)

because the former method consumes less memory than inline monitoring *on average*, but it then executes for a longer period of time. Figs. 9 and 10 illustrate the mean measurements obtained per experiment; a different depiction that shows the *total memory consumed* during the experiment can be found in app. C. From the figures reported in tbl. 2, one could even make a case that for

	Steady				Pulse				Burst			
	DS _m	DM _m	OS _s	DM _s	DS _m	DM _m	DS _s	DM _s	DS _m	DM _m	DS _s	DM _s
Scheduler	20.4	10.1	0.001	0.001	21.4	14.4	0.001	0.000	23.2	15.4	0.001	0.000
Memory	1.4	0.7	0.000	0.000	2.0	1.5	0.000	0.000	1.4	0.9	0.000	0.000
RTT	194.6	134.4	0.001	0.001	200.5	185.4	0.001	0.001	87.2	72.7	0.000	0.000
Execution	74.9	61.1	0.000	0.000	79.2	73.7	0.000	0.000	60.7	50.3	0.000	0.000

Decentralised outline separate on master (DS_m), Decentralised outline merged on master (DM_m)

Decentralised outline separate per slave (DS_s), Decentralised outline merged per slave (DM_s)

Tbl. 3. Percentage amortised runtime overhead on each slave taken at the maximum load $n = 100k$

settings where memory is limited but execution time is not, *outlined* centralised monitoring is more appropriate than *inlined* decentralised monitoring.

The memory consumption, RTT, and execution duration plots in figs. 9 and 10 exhibit a *linear growth* beyond specific x -axis thresholds. This contrasts with the plots in figs. 7 and 8 for $n = 10k$, where different trends may be observed: the execution duration plot under a Steady load shape grows (negative) quadratically in fig. 7, but follows a cubic trend in the case of Pulse and Burst loads in fig. 8; a similar effect is obtained in the RTT for Burst (consult the fitted data plots in app. C). These differences in runtime behaviour underscore the value of performing tests using reasonably high loads, as this increases the chances of *observing likely trends*. For instance, the empirical evidence obtained for moderate loads could mislead one to assert that outline monitoring scales very poorly in the case of RTT under moderately-sized Burst loads.

Estimated overhead on slaves. Our results thus far present an *overall view* of the overhead induced by runtime monitoring. In certain cases however, this measure is too coarse since we would be interested more in quantifying the overhead incurred at *each slave*: this bears particular relevance to distributed setting where the processing capability of the system is spread over heterogenous machines (e.g. in an IoT set-up deployed on edge nodes [Shi et al. 2016] with limited computing power, understanding slave overheads is essential). Our experiment set-up does not allow us to directly measure the overhead at each slave since our measurements are collected globally. Instead, Tbl. 3 shows the percentage overhead for decentralised outline monitoring induced on the master process *only*, together with the estimated *overhead apportioned* over each slave. The overhead incurred by the SuS when monitoring the master process, DS_m and DM_m in tbl. 3, is obtained by setting up the experiment with $n = 100k$ and $w = 100$ as before. We can then approximate the *combined overhead* induced by the slaves by subtracting DS_m and DM_m from the *total overhead* obtained when the master and slaves are monitored together (DS and DM in tbl. 2). An *apportioned* overhead per slave can therefore be obtained by dividing this combined overhead by the number of slaves, i.e., 100 k, to give DS_s and DM_s in tbl. 3. Figs. 18 and 19 in app. C show the gap in overhead between the SuS fixed with one monitor on the master process and to the fully-monitored system. The estimated figures in tbl. 3 clearly indicate that the two flavours of decentralised runtime monitoring from fig. 1 induce nominal overhead per slave for the load shapes we consider.

5 CONCLUSION

We provide a detailed study of asynchronous outline monitoring, an alternative to inline monitoring that is often discarded due to its high overheads. Our study makes a case that there are instances where outlining is the only available solution for analysing a system at runtime and that the overheads are tolerable in certain scenarios. To the best of our knowledge, the algorithm presented

1030 in sec. 3 differs from the state-of-the art in three fundamental ways: (i) it asynchronously gathers
1031 events from the SuS, (ii) effects the analysis using outline monitors, and, (iii) dynamically scales
1032 the runtime set-up as the SuS grows and shrinks. Our experiments in sec. 4 give scenarios that
1033 indicate the uses-cases where outline monitoring is best applied. They establish a pessimistic point
1034 of departure for outline monitoring: (i) RV monitoring typically does not exclude any events but
1035 less stringent analyses can resort to sampling to lower overheads [Sigelman et al. 2010]; (ii) RV
1036 analysis was carried out on every slave until termination, but RV verdicts may be reached early in
1037 the execution; We anticipate that less demanding settings such as general APM tools will lower
1038 further the overhead discrepancies reported.

1039

1040

1041 5.1 Related Work

1042 *Decentralised monitoring in RV.* Following standard texts on distributed computing [Buyya et al.
1043 2011; Coulouris et al. 2005; Tarkoma 2010], distributed systems are *necessarily* concurrent, but can
1044 be either centralised or decentralised, although the opposite *does not* always hold. Specifically,
1045 assumption A_1 in sec. 3 renders a system concurrent whereas lifting assumptions A_3 and A_4 changes
1046 this to distributed. Systems that rely on a global clock are neither concurrent nor distributed.
1047 For example, works such as [Colombo and Falcone 2016], while decentralised, do not qualify as
1048 distributed solutions; system components operate in synchronous rounds whereby a *unique* global
1049 trace can be reconstructed by combining the different traces collected at each component. Their
1050 approach does not address challenges such as message reordering (assumption A_3 in sec. 3).

1051 Code injection is used in a number of tools targeting concurrent and distributed component
1052 systems. For example, [Sen et al. 2006] study decentralised monitors that are attached to differ-
1053 ent threads to collect and process trace events locally. In an earlier work [Sen et al. 2004], this
1054 investigation is conducted in a distributed setting using decentralised monitors that are weaved
1055 into components of the SuS. The authors focus on the efficiency of monitor communication but
1056 do not study nor quantify the overhead induced by runtime monitoring. Minimising overhead is
1057 also the focus of [Mostafa and Bonakdarpour 2015]. In this setting, the SuS consists of distributed
1058 asynchronous processes that communicate together via message-passing primitives over reliable
1059 channels. Similar to ours, their monitoring algorithm does not rely on a global notion of timing, and
1060 does not assume failing system components. The work by [Basin et al. 2015] is one of the few that
1061 considers distributed system monitoring where components and network links may fail. Despite
1062 the absence of a global clock, their monitoring algorithm is based on the timed asynchronous model
1063 for distributed systems [Cristian and Fetzer 1999] that assumes *highly-synchronised* physical clocks
1064 across nodes. In a different manner, [Bonakdarpour et al. 2016; Fraigniaud et al. 2014] address the
1065 problem of when the monitors themselves crash. Failure is an aspect that we do not presently
1066 address (see assumptions A_4 and A_5). The tools in [Basin et al. 2015; Cassar and Francalanza 2016;
1067 Jin et al. 2012] weave special instructions to enable the system to externalise its monitors, similar to
1068 fig. 1a. Crucially, inlining spares their algorithms from having to deal with reordered trace events.

1069 Tools such as [Attard and Francalanza 2017; Neykova and Yoshida 2017a] target Erlang. In [Neykova
1070 and Yoshida 2017a], the authors propose a method that statically analyses the program commu-
1071 nication flow, specified in terms of a multiparty protocol. Monitors attached to system processes
1072 check that the messages received coincide with the projected type, and in the case of failure, the
1073 associated processes are restarted. The authors show that their recovery algorithm induces less
1074 communication overhead, and improves upon the static process structure recovery mechanisms of-
1075 fered by the Erlang/OTP platform. Similarly, [Attard and Francalanza 2017] focus on decentralised
1076 outline monitoring in a concurrent setting. By contrast to [Neykova and Yoshida 2017a], they
1077 leverage the native tracing infrastructure offered by the EVM.

1078

1079 We remark that the works above all rely on bespoke evaluation platforms, making it hard to
1080 reproduce or directly compare their empirical results to ours. They either use loads modelled on
1081 a Poisson process, *i.e.*, Steady load, or fail altogether to specify the load types considered. Our
1082 empirical study has shown that different load shapes are indeed relevant.

1083
1084 *Evaluation and benchmarking tools.* Savina [Imam and Sarkar 2014] addresses the lack of a common
1085 benchmarking tool for actor-based systems. In the spirit of DaCapo [Blackburn et al. 2006], it
1086 provides a suite of diverse benchmarks that represent compute (as opposed to IO) intensive appli-
1087 cations. These range from micro-benchmarks to classical concurrency problems and parallelism
1088 benchmarks. Similar our evaluation framework, Savina includes instantiations of the master-slave
1089 set-up (Trapezoidal Approximation and Precise Pi Computation), that are configurable with the
1090 number of slaves. However, our implementation accommodates more parameters: (i) the number
1091 of work requests per slave, (ii) the probability of allocating and acknowledging work requests,
1092 and, (iii) the type of load shape. In addition, we also support dynamic set-ups, as opposed to the
1093 static ones currently included in this suite. Savina measures the mean execution duration; we also
1094 collect the scheduler utilisation, memory consumption and RTT between the master and slaves. The
1095 requirement for a dynamic set-up, together with the performance metrics outlined in sec. 4 made
1096 Savina unsuitable to our use-case. Presently, Savina targets JVM actor-based languages. Like our
1097 implementation, Savina does not yet include benchmarks based on the peer-to-peer architecture.

1098 Themis [El-Hokayem and Falcone 2017] is a tool that aims to facilitate the design and analysis
1099 of decentralised monitoring algorithms. It supports static set-ups where the number of system
1100 components and corresponding monitors is known and remains fixed at runtime. Unlike Savina
1101 or our benchmarking tool, Themis processes only pre-recorded tracers supplied via text files,
1102 making it incompatible with online monitoring. In [El-Hokayem and Falcone 2017], the authors
1103 claim that these trace files may be obtained from instrumented programs. Monitor and monitor
1104 communication in Themis is simulated via method calls that deposit messages inside blocking
1105 queues linked to each monitor. Like Savina, Themis is developed for Java applications.

1106 The Behaviour, Interaction, Priority (BIP) framework models heterogenous real-time component
1107 systems. In BIP [Basu et al. 2006; El-Hokayem et al. 2018], the interaction between components is
1108 specified using syntactic descriptions that are parsed by a Java front-end and translated to C++
1109 code. The automata-based operational model of BIP is implemented into their back-end platform
1110 that executes the generated code. BIP supports synchronous and asynchronous components that
1111 may run on the same or separate threads. While the back-end implementation relies on POSIX
1112 threads [Butenhof 1997] for easy integration with C++, this also limits the scalability of BIP when
1113 many asynchronous components are used, since each pthread takes kernel resources from the
1114 system. By contrast, the green processes used by Erlang allows our evaluation platform to scale
1115 considerably while incurring manageable overhead. Erlang process scheduling is performed by the
1116 EVM, making these much more lightweight when compared to pthreads. BIP is principally built as
1117 a flexible modelling tool that is inapplicable to our benchmarking requirements from sec. 4.

1118 Kollaps [Gouveia et al. 2020] emulates distributed network conditions from an application-level
1119 perspective that considers the observable end-to-end properties such as latency or packet loss.
1120 The tool simplifies the network view by abstracting over the state of physical network appliances
1121 that sit in between nodes of the distributed application. Kollaps is fully-decentralised and agnostic
1122 of the application language and transport protocol. The authors show that Kollaps can closely
1123 model realistic network conditions. We plan to integrate Kollaps into our evaluation framework
1124 when extending it to account for further experiment variables such as packet loss and node failure
1125 (assumptions A_4 and A_5).

REFERENCES

- 1128
1129 Luca Aceto, Antonis Achilleos, Adrian Francalanza, Anna Ingólfssdóttir, and Karoliina Lehtinen. 2019a. Adventures in
1130 monitorability: from branching to linear time and back again. *Proc. ACM Program. Lang.* 3, POPL (2019), 52:1–52:29.
1131 <https://doi.org/10.1145/3290365>
- 1132 Luca Aceto, Antonis Achilleos, Adrian Francalanza, Anna Ingólfssdóttir, and Karoliina Lehtinen. 2019b. An Operational
1133 Guide to Monitorability. In *SEFM (LNCS)*, Vol. 11724. Springer, 433–453. https://doi.org/10.1007/978-3-030-30446-1_23
- 1134 Gul Agha, Ian A. Mason, Scott F. Smith, and Carolyn L. Talcott. 1997. A Foundation for Actor Computation. *JFP* 7, 1 (1997),
1–72.
- 1135 Nadia Alshahwan, Andrea Ciancone, Mark Harman, Yue Jia, Ke Mao, Alexandru Marginean, Alexander Mols, Hila Peleg,
1136 Federica Sarro, and Ilya Zorin. 2019. Some Challenges for Software Testing Research (Invited Talk Paper). In *ISSTA*.
ACM, 1–3.
- 1137 Peter Alvaro, Kolton Andrus, Chris Sanden, Casey Rosenthal, Ali Basiri, and Lorin Hochstein. 2016. Automating Failure
1138 Testing Research at Internet Scale. In *SoCC*. ACM, 17–28.
- 1139 Joe Armstrong. 2007. *Programming Erlang: Software for a Concurrent World* (first ed.). Pragmatic Bookshelf.
- 1140 Vinay Arora, Rajesh Kumar Bhatia, and Maninder Singh. 2016. A systematic review of approaches for testing concurrent
1141 programs. *CCPE* 28, 5 (2016), 1572–1611.
- 1142 Duncan Paul Attard and Adrian Francalanza. 2017. Trace Partitioning and Local Monitoring for Asynchronous Components.
1143 In *SEFM (LNCS)*, Vol. 10469. Springer, 219–235.
- 1144 Fred Baker. 1995. Requirements for IPv4 Routers. <https://www.ietf.org/rfc/rfc1812.txt>
- 1145 Ezio Bartocci, Yliès Falcone, Borzoo Bonakdarpour, Christian Colombo, Normann Decker, Klaus Havelund, Yogi Joshi, Felix
1146 Klaedtke, Reed Milewicz, Giles Reger, Grigore Rosu, Julien Signoles, Daniel Thoma, Eugen Zalinescu, and Yi Zhang. 2019.
1147 First international Competition on Runtime Verification: rules, benchmarks, tools, and final results of CRV 2014. *STTT*
21, 1 (2019), 31–70. <https://doi.org/10.1007/s10009-017-0454-5>
- 1148 Ezio Bartocci, Yliès Falcone, Adrian Francalanza, and Giles Reger. 2018. Introduction to Runtime Verification. In *Lectures on*
1149 *Runtime Verification*. LNCS, Vol. 10457. Springer, 1–33.
- 1150 David A. Basin, Felix Klaedtke, and Eugen Zalinescu. 2015. Failure-Aware Runtime Verification of Distributed Systems. In
1151 *FSTTCS (LIPIcs)*, Vol. 45. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 590–603.
- 1152 Ananda Basu, Marius Bozga, and Joseph Sifakis. 2006. Modeling Heterogeneous Real-time Components in BIP. In *SEFM*.
1153 IEEE Computer Society, 3–12.
- 1154 Shay Berkovich, Borzoo Bonakdarpour, and Sebastian Fischmeister. 2015. Runtime Verification with Minimal Intrusion
1155 through Parallelism. *FMSD* 46, 3 (2015), 317–348.
- 1156 Stephen M. Blackburn, Robin Garner, Chris Hoffmann, Asjad M. Khan, Kathryn S. McKinley, Rotem Bentzur, Amer Diwan,
1157 Daniel Feinberg, Daniel Frampton, Samuel Z. Guyer, Martin Hirzel, Antony L. Hosking, Maria Jump, Han Bok Lee,
1158 J. Eliot B. Moss, Aashish Phansalkar, Darko Stefanovic, Thomas VanDrunen, Daniel von Dincklage, and Ben Wiedermann.
2006. The DaCapo Benchmarks: Java Benchmarking Development and Analysis. In *OOPSLA*. ACM, 169–190.
- 1159 Borzoo Bonakdarpour, Pierre Fraigniaud, Sergio Rajsbaum, David A. Rosenblueth, and Corentin Travers. 2016. Decentralized
Asynchronous Crash-Resilient Runtime Verification. In *CONCUR (LIPIcs)*, Vol. 59. Schloss Dagstuhl - Leibniz-Zentrum
fuer Informatik, 16:1–16:15.
- 1160 David R. Butenhof. 1997. *Programming with POSIX threads* (first ed.). Addison-Wesley.
- 1161 Rajkumar Buyya, James Broberg, and Andrzej Goscinski. 2011. *Cloud Computing: Principles and Paradigms (Wiley Series on*
1162 *Parallel and Distributed Computing)*. Wiley.
- 1163 Bryan Cantrill. 2006. Hidden in Plain Sight. *ACM Queue* 4, 1 (2006), 26–36.
- 1164 Bryan Cantrill, Michael W. Shapiro, and Adam H. Leventhal. 2004. Dynamic Instrumentation of Production Systems. In
1165 *USENIX Annual Technical Conference, General Track*. USENIX, 15–28.
- 1166 João M.P. Cardoso, José Gabriel F. Coutinho, and Pedro C. Diniz. 2017. Chapter 5 - Source code transformations and
1167 optimizations. In *Embedded Computing for High Performance*, João M.P. Cardoso, José Gabriel F. Coutinho, and Pedro C.
Diniz (Eds.). Morgan Kaufmann, Boston, 137 – 183. <https://doi.org/10.1016/B978-0-12-804189-5.00005-3>
- 1168 João M. P. Cardoso, José Gabriel F. Coutinho, Tiago Carvalho, Pedro C. Diniz, Zlatko Petrov, Wayne Luk, and Fernando M.
1169 Gonçalves. 2016. Performance-driven instrumentation and mapping strategies using the LARA aspect-oriented program-
1170 ming approach. *SPE* 46, 2 (2016), 251–287.
- 1171 Ian Cassar and Adrian Francalanza. 2016. On Implementing a Monitor-Oriented Programming Framework for Actor Systems.
1172 In *IFM (LNCS)*, Vol. 9681. Springer, 176–192.
- 1173 Francesco Cesarini and Simon Thompson. 2009. *Erlang Programming: A Concurrent Approach to Software Development* (first
1174 ed.). O'Reilly Media.
- 1175 Feng Chen and Grigore Rosu. 2007. MOP: An Efficient and Generic Runtime Verification Framework. In *OOPSLA*. ACM,
569–588.
- 1176 Feng Chen and Grigore Rosu. 2009. Parametric Trace Slicing and Monitoring. In *TACAS (LNCS)*, Vol. 5505. Springer, 246–261.

- 1177 Kung Chen and Ju-Bing Chen. 2006. On Instrumenting Obfuscated Java Bytecode with Aspects. In *SESSICSE*. ACM, 19–26.
- 1178 David M. Ciemiewicz. 2001. What Do You Mean? - Revisiting Statistics for Web Response Time Measurements. In *CMG*.
1179 Computer Measurement Group, 385–396.
- 1180 Christian Colombo and Yliès Falcone. 2016. Organising LTL Monitors over Distributed Systems with a Global Clock. *FMSD*
1181 49, 1-2 (2016), 109–158.
- 1182 Christian Colombo, Gordon J. Pace, and Gerardo Schneider. 2009. LARVA—Safer Monitoring of Real-Time Java Programs
(Tool Paper). In *SEFM*. IEEE Computer Society, 33–37.
- 1183 George Coulouris, Jean Dollimore, and Tim Kindberg. 2005. *Distributed Systems: Concepts and Design* (fourth ed.). Addison
1184 Wesley.
- 1185 Flaviu Cristian and Christof Fetzer. 1999. The Timed Asynchronous Distributed System Model. *IEEE Trans. Parallel Distrib.*
1186 *Syst.* 10, 6 (1999), 642–657.
- 1187 Mathieu Desnoyers and Michel R. Dagenais. 2006. *The LTTng tracer : A low impact performance and behavior monitor for*
1188 *GNU / Linux*. Technical Report.
- 1189 Antoine El-Hokayem and Yliès Falcone. 2017. Monitoring Decentralized Specifications. In *ISSTA*. ACM, 125–135.
- 1190 Antoine El-Hokayem and Yliès Falcone. 2017. THEMIS: A Tool for Decentralized Monitoring Algorithms. In *ISSTA*. ACM,
1191 372–375.
- 1192 Antoine El-Hokayem and Yliès Falcone. 2020. On the Monitoring of Decentralized Specifications: Semantics, Properties,
1193 Analysis, and Simulation. *ACM Trans. Softw. Eng. Methodol.* 29, 1, Article 1 (Jan. 2020), 57 pages. [https://doi.org/10.1145/](https://doi.org/10.1145/3355181)
1194 [3355181](https://doi.org/10.1145/3355181)
- 1195 Antoine El-Hokayem, Yliès Falcone, and Mohamad Jaber. 2018. Modularizing behavioral and architectural crosscutting
1196 concerns in formal component-based systems - Application to the Behavior Interaction Priority framework. *JLAMP* 99
(2018), 143–177. <https://doi.org/10.1016/j.jlamp.2018.05.005>
- 1197 Úlfar Erlingsson. 2004. *The Inlined Reference Monitor Approach to Security Policy Enforcement*. Ph.D. Dissertation. Cornell
1198 University.
- 1199 Úlfar Erlingsson and Fred B. Schneider. 1999. SASI Enforcement of Security Policies: A Retrospective. In *NSPW*. ACM,
1200 87–95.
- 1201 Pierre Fraigniaud, Sergio Rajsbaum, and Corentin Travers. 2014. On the Number of Opinions Needed for Fault-Tolerant
1202 Run-Time Monitoring in Distributed Systems. In *RV (Lecture Notes in Computer Science)*, Vol. 8734. Springer, 92–107.
- 1203 Adrian Francalanza, Luca Aceto, Antonis Achilleos, Duncan Paul Attard, Ian Cassar, Dario Della Monica, and Anna
1204 Ingólfssdóttir. 2017. A Foundation for Runtime Monitoring. In *Runtime Verification (RV) (LNCS)*, Vol. 10548. Springer,
1205 8–29. https://doi.org/10.1007/978-3-319-67531-2_2
- 1206 Adrian Francalanza, Jorge A. Pérez, and César Sánchez. 2018. Runtime Verification for Decentralised and Distributed
1207 Systems. In *Lectures on Runtime Verification*. LNCS, Vol. 10457. Springer, 176–210.
- 1208 Adrian Francalanza and Aldrin Seychell. 2015. Synthesising Correct Concurrent Runtime Monitors. *FMSD* 46, 3 (2015),
1209 226–261.
- 1210 Vijay K. Garg. 2015. *Elements of Distributed Computing* (first ed.). Wiley India.
- 1211 Paulo Gouveia, João Neves, Carlos Segarra, Luca Liechti, Shady Issa, Valerio Schiavoni, and Miguel Matos. 2020. Kollaps:
1212 Decentralized and Dynamic Topology Emulation. In *EuroSys*. ACM, 23:1–23:16.
- 1213 Duncan A. Grove and Paul D. Coddington. 2005. Analytical Models of Probability Distributions for MPI Point-to-Point
1214 Communication Times on Distributed Memory Parallel Computers. In *ICA3PP (LNCS)*, Vol. 3719. Springer, 406–415.
- 1215 Christoph Heger, André van Hoorn, Mario Mann, and Dusan Okanovic. 2017. Application Performance Management: State
1216 of the Art and Challenges for the Future. In *ICPE*. ACM, 429–432.
- 1217 Gregor Hohpe and Bobby Woolf. 2003. *Enterprise Integration Patterns: Designing, Building, and Deploying Messaging Solutions*
1218 (first ed.). Addison-Wesley Professional.
- 1219 Shams Mahmood Imam and Vivek Sarkar. 2014. Savina - An Actor Benchmark Suite: Enabling Empirical Evaluation of
1220 Actor Libraries. In *AGERE!@SPLASH*. ACM, 67–80.
- 1221 Pooyan Jamshidi, Claus Pahl, Nabor C. Mendonça, James Lewis, and Stefan Tilkov. 2018. Microservices: The Journey So Far
1222 and Challenges Ahead. *IEEE Software* 35, 3 (2018), 24–35.
- 1223 Dongyun Jin, Patrick O’Neil Meredith, Choonghwan Lee, and Grigore Rosu. 2012. JavaMOP: Efficient Parametric Runtime
1224 Monitoring Framework. In *ICSE*. IEEE Computer Society, 1427–1430.
- 1225 Nicolai M. Josuttis. 2007. *SOA in Practice: The Art of Distributed System Design: Theory in Practice* (first ed.). O’Reilly Media.
- Bill Kayser. 2017. What Is the Expected Distribution of Website Response Times? [https://blog.newrelic.com/engineering/](https://blog.newrelic.com/engineering/expected-distributions-website-response-times)
[expected-distributions-website-response-times](https://blog.newrelic.com/engineering/expected-distributions-website-response-times)
- Gregor Kiczales, John Lamping, Anurag Mendhekar, Chris Maeda, Cristina Videira Lopes, Jean-Marc Loingtier, and John
Irwin. 1997. Aspect-Oriented Programming. In *ECOOP (LNCS)*, Vol. 1241. Springer, 220–242.
- Moonjoo Kim, Sampath Kannan, Insup Lee, Oleg Sokolsky, and Mahesh Viswanathan. 2001. Java-MaC: a Run-time Assurance
Tool for Java Programs. *Electr. Notes Theor. Comput. Sci.* 55, 2 (2001), 218–235.

- 1226 Jay Ligatti, Lujó Bauer, and David Walker. 2005. Edit automata: enforcement mechanisms for run-time security policies.
1227 *International Journal of Information Security (IJIS)* 4, 1-2 (2005), 2–16. <https://doi.org/10.1007/s10207-004-0046-8>
- 1228 Zhen Liu, Nicolas Niclausse, and César Jalpa-Villanueva. 2001. Traffic Model and Performance Evaluation of Web Servers.
1229 *Perform. Evaluation* 46, 2-3 (2001), 77–100.
- 1230 Robert Martin. 2013. *Agile Software Development, Principles, Patterns, and Practices* (first ed.). Pearson.
- 1231 Menna Mostafa and Borzoo Bonakdarpour. 2015. Decentralized Runtime Verification of LTL Specifications in Distributed
1232 Systems. In *IPDPS*. IEEE Computer Society, 494–503.
- 1233 Romyana Neykova and Nobuko Yoshida. 2017a. Let It Recover: Multiparty Protocol-Induced Recovery. In *CC*. ACM, 98–108.
- 1234 Romyana Neykova and Nobuko Yoshida. 2017b. Multiparty Session Actors. *Logical Methods in Computer Science* 13, 1
1235 (2017).
- 1236 Giles Reger, Helena Cuenca Cruz, and David E. Rydeheard. 2015. MarQ: Monitoring at Runtime with QE. In *TACAS (LNCS)*,
1237 Vol. 9035. Springer, 596–610.
- 1238 Richard J. Rossi. 2018. *Mathematical Statistics: An Introduction to Likelihood Based Inference*. Wiley.
- 1239 César Sánchez, Gerardo Schneider, Wolfgang Ahrendt, Ezio Bartocci, Domenico Bianculli, Christian Colombo, Yliès Falcone,
1240 Adrian Francalanza, Srdan Krstic, João M. Lourenço, Dejan Nickovic, Gordon J. Pace, José Rufino, Julien Signoles, Dmitriy
1241 Traytel, and Alexander Weiss. 2019. A survey of challenges for runtime verification from advanced application domains
1242 (beyond software). *FMSD* 54, 3 (2019), 279–335.
- 1243 Torben Scheffel and Malte Schmitz. 2014. Three-Valued Asynchronous Distributed Runtime Verification. In *MEMOCODE*.
1244 IEEE, 52–61.
- 1245 Fred B. Schneider. 2000. Enforceable Security Policies. *Transactions on Information and System Security (TISSEC)* 3, 1 (2000),
1246 30–50. <https://doi.org/10.1145/353323.353382>
- 1247 Koushik Sen, Abhay Vardhan, Gul Agha, and Grigore Rosu. 2004. Efficient Decentralized Monitoring of Safety in Distributed
1248 Systems. In *ICSE*. IEEE Computer Society, 418–427.
- 1249 Koushik Sen, Abhay Vardhan, Gul Agha, and Grigore Rosu. 2006. Decentralized Runtime Analysis of Multithreaded
1250 Applications. In *IPDPS*. IEEE.
- 1251 Weisong Shi, Jie Cao, Quan Zhang, Youhuizi Li, and Lanyu Xu. 2016. Edge Computing: Vision and Challenges. *IEEE Internet*
1252 *of Things Journal* 3, 5 (2016), 637–646.
- 1253 Jeremy G. Siek and Walid Taha. 2007. Gradual Typing for Objects. In *ECOOP (LNCS)*, Vol. 4609. Springer, 2–27. https://doi.org/10.1007/978-3-540-73589-2_2
- 1254 Benjamin H. Sigelman, Luiz André Barroso, Mike Burrows, Pat Stephenson, Manoj Plakal, Donald Beaver, Saul Jaspan, and
1255 Chandan Shanbhag. 2010. *Dapper, a Large-Scale Distributed Systems Tracing Infrastructure*. Technical Report. Google, Inc.
1256 <https://research.google.com/archive/papers/dapper-2010-1.pdf>
- 1257 Connie U. Smith and Lloyd G. Williams. 2001. Software Performance AntiPatterns; Common Performance Problems and
1258 their Solutions. In *CMG§*. Computer Measurement Group, 797–806.
- 1259 Connie U. Smith and Lloyd G. Williams. 2002. New Software Performance AntiPatterns: More Ways to Shoot Yourself in
1260 the Foot. In *CMG*. Computer Measurement Group, 667–674.
- 1261 Asumu Takikawa, T. Stephen Strickland, Christos Dimoulas, Sam Tobin-Hochstadt, and Matthias Felleisen. 2012. Gradual
1262 typing for first-class classes. In *OOPSLA*, Gary T. Leavens and Matthew B. Dwyer (Eds.). ACM, 793–810. <https://doi.org/10.1145/2384616.2384674>
- 1263 Sasu Tarkoma. 2010. *Overlay Networks: Toward Information Networking* (first ed.). Auerbach Publications.

1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274

Event	Action (<i>e.act</i>)	Field names	Description
fork	frk	src	PID p_s of the (parent) process invoking fork(g)
		tgt	PID p_s of the forked (child) process
		sig	Code signature g run by the forked process
exit	ext	src	PID p_s of the terminated process
send	snd	src	PID p_s of the sender process
		tgt	PID p_s of the recipient process
receive	rcv	src	PID p_s of the recipient process

Tbl. 4. Trace event messages data field names

A DECENTRALISED OUTLINE MONITORING ALGORITHM AUXILLARY CODE

Field access notation for tracer messages. Just as the message qualifier is accessible through the field name $m.type$, so are the data elements of the respective trace event types `frk`, `ext`, `snd`, and `rcv`. These are catalogued in tbl. 4.

Trace event acquisition. The tracing mechanism is defined by the operations TRACE, CLEAR and PREEMPT listed in lst. 4. TRACE enables a tracer p_t to register its interest in being notified about trace events of a system process p_s ; this operation can be undone using CLEAR. CLEAR blocks the caller, and returns only once all the trace event messages for p_s that are in the process of being delivered are deposited into the message queue of p_t . PREEMPT combines CLEAR and TRACE, enabling a tracer p_t to take over the tracing of process p_s from another tracer p'_t . The preemption instructions on lines 19–20 are ideally executed *atomically* to prevent potential trace event loss that could occur when switching tracers. This guarantee however, depends exclusively on the implementation of the underlying tracing mechanism. We recall that, following assumption A_8 , tracing is *inherited* by every child process that a traced system process forks; CLEAR or PREEMPT can then be used to alter this arrangement.

Trace routing and relaying. Our algorithm performs routing using two operations, ROUTE and RELAY in lst. 5. ROUTE creates a *new* message, r , with type `rtd`, that embeds trace events or `drc` commands

```

1309 1 def TRACE( $p_s, p_t$ )
1310 2 if  $p_s$  is not traced then
1311 3 Set the tracer for  $p_s$  to  $p_t$ ;  $p_t$  will trace new
1312 4 descendant processes  $p_{s_1}, p_{s_2}, \dots$  forked
1313 5 by  $p_s$  automatically (assumption  $A_8$ )
1314 6 while  $p_s$ 's tracer is set do
1315 7    $s \leftarrow$  read next event for  $p_s$  from
1316 8   trace event source
1317 9    $e \leftarrow$  encode  $s$  as a message
1318 10   $p_t ! e$ 
1319 11 end while
1320 12 end if
1321 13 end def
1322
1323 14 def CLEAR( $p_s, p_t$ )
1324 15 if  $p_s$  is traced then
1325 16 Clear the tracer  $p_t$  for  $p_s$ ;  $p_t$  still traces
1326 17 the descendant processes  $p_{s_1}, p_{s_2}, \dots$  of  $p_s$ 
1327 18 Block until the trace events for  $p_s$  that are in
1328 19 transit are delivered to  $p_t$ 
1329 20 end if
1330 21 end def
1331
1332 Expect:  $p_s$ 's tracer is set
1333 22 def PREEMPT( $p_s, p_t$ )
1334 23  $p'_t \leftarrow p_s$ 's tracer
1335 24 CLEAR( $p_s, p'_t$ )
1336 25 TRACE( $p_s, p_t$ )
1337 26 end def

```

Lst. 4. Trace event acquisition, clear, and preemption operations offered by the tracing mechanism

```

1324 Expect:  $m.type = evt \vee m.type = dtc$ 
1325 1 def ROUTE( $m, p_t$ )
1326 2  $p_t ! \langle rtd, self(), m \rangle$ 
1327 3 end def
1328
1329 Expect:  $m.type = rtd$ 
1330 4 def RELAY( $m, p_t$ )
1331 5  $p_t ! m$ 
1332 6 end def

```

Lst. 5. Message routing and relaying operations

```

1333 1 def TRACER( $\sigma, v, p_s, p_t$ )
1334 2 # Executable monitor map  $\sigma, \Phi$  is copied to the
1335 3 # new (child) tracer state  $\sigma'$ ; component group  $\Gamma$ 
1336 4 # is initialised with the process being traced,  $p_s$ 
1337 5  $\sigma' \leftarrow \langle \Pi \leftarrow \emptyset, \sigma, \Phi, \Gamma \leftarrow \{ \langle p_s, \bullet \rangle \} \rangle$ 
1338 6 DETACH( $p_s, p_t$ )
1339 7  $p_a \leftarrow \text{fork}(v)$  executable monitor
1340 8 # New tracer is started in  $\bullet$  mode to process
1341 9 # routed events before locally-traced ones
1342 10 LOOP $\bullet$ ( $\sigma', p_a$ )
1343 11 end def
1344
1345 7 def DETACH( $p_s, p_t$ )
1346 8  $p_t' \leftarrow \text{self}()$ 
1347 9 PREEMPT( $p_s, p_t'$ )
1348 10  $p_t ! \langle dtc, p_t', p_s \rangle$ 
1349 11 end def
1350
1351 12 def TRYGC( $\sigma, p_a$ )
1352 13 if  $\sigma.\Gamma = \emptyset \wedge \sigma.\Pi = \emptyset$  then
1353 14   Signal monitor  $p_a$  to terminate
1354 15   Terminate tracer
1355 16 end if
1356 17 end def

```

Lst. 6. Operations used by the (◦) and priority (●) tracer loops

that need to be routed. The PID of the tracer process invoking ROUTE is included into a routed message: we refer to this process as the *router tracer* that it is responsible for injecting the message into the tracer choreography. This PID is retrievable using the field $m.rtr$ (see tbl. 5), and enables other tracers to identify the tracer that initiated the message dispatch. Routed messages can *only* be handled by tracers or forwarded using RELAY.

Starting the system. START in lst. 7 launches the SuS and monitoring system in tandem. The operation accepts the code signature g , as the entry point of the SuS, together with the map of executable monitors, Φ . As a safeguard that prevents the initial loss of trace events, the SuS is launched in a paused state (line 7.2) to permit the root tracer to start tracing the top-level system process. ROOT resumes the system (7.8), and begins its trace inspection in *direct* mode, as shown in line 7.10.

Message	Type ($m.type$)	Field name	Description
routed	rtd	rtr	PID p_t of the (ancestor) tracer that starts routing the message
		emb	The embedded trace event e or command c
detach command	dtd	tgt	PID p_s of the system process that is, from this point, traced by the new tracer
		iss	PID p_t of the new tracer issuing the detach command to the <i>router tracer</i>

Tbl. 5. Routed messages and detach command data field names

```

1373 1 def START( $g, \Phi$ )
1374   # Pausing allows root tracer to be set
1375   # up; no initial message loss
1376 2  $p_s \leftarrow \text{fork}(g)$  in paused mode
1377 3  $p_t \leftarrow \text{fork}(\text{ROOT}(p_s, \Phi))$ 
1378 4 return  $\langle p_s, p_t \rangle$ 
1379 5 end def
1380
1381
1382
1383
1384 6 def ROOT( $p_s, \Phi$ )
1385 7 TRACE( $p_s, \text{self}()$ )
1386 8 Resume system  $p_s$ 
1387 9  $\sigma \leftarrow \langle \Pi \leftarrow \emptyset, \Phi, \Gamma \leftarrow \{ \langle p_s, \circ \rangle \} \rangle$ 
1388   # Root tracer has no monitor
1389 10 LOOP $_c(\sigma, \perp)$ 
1390 11 end def

```

Lst. 7. System starting operation and root tracer

B EXPERIMENT SET-UP AND EVALUATION

Inline monitoring implementation. We synthesise *automata-based* monitors from high-level correctness specifications. These monitors are encoded as executable functions that can be represented as an AST. Fig. 11 outlines how our monitors are inlined in the SuS. In step ①, the Erlang source code of the system is pared into the corresponding AST, step ② The Erlang compilation process contains a *parse transform* phase step ③ provides a hook that allows for the AST to be post-processed [Cesarini and Thompson 2009]. We leverage this mechanism through our custom-built weaver, step ④, that injects into the AST of the SuS the AST of the monitor in step ⑤. It performs two types of code transformations:

- C_1 *Monitor bootstrapping.* The function encoding the synthesised monitor is stored in the process dictionary (a key-value map) of the monitored system process to make it globally accessible from within said process;
- C_2 *Instrumentation points.* The AST of the system is instrumented with calls at the points of interest: these calls constitute the trace event actions that are to be analysed.

The instrumented calls in transformation C_2 retrieve the monitor function stored the process dictionary in transformation C_1 , and apply it to the trace event in question. This function application on the event returns the *monitor continuation* that is used to replace the current monitor in the process dictionary. Our two-step weaving process produces the instrumented code in step ⑥ which can be subsequently compiled by the Erlang compiler into the application binary. We note that the same monitor ASTs synthesised for use in inline monitoring are used by our outline monitoring algorithm as well.

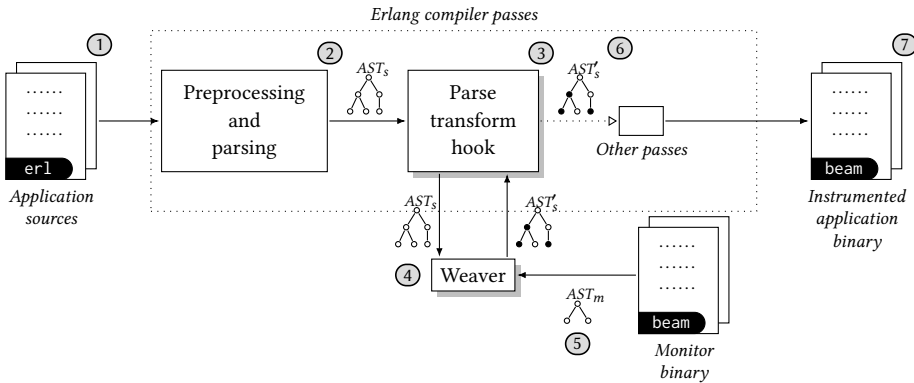


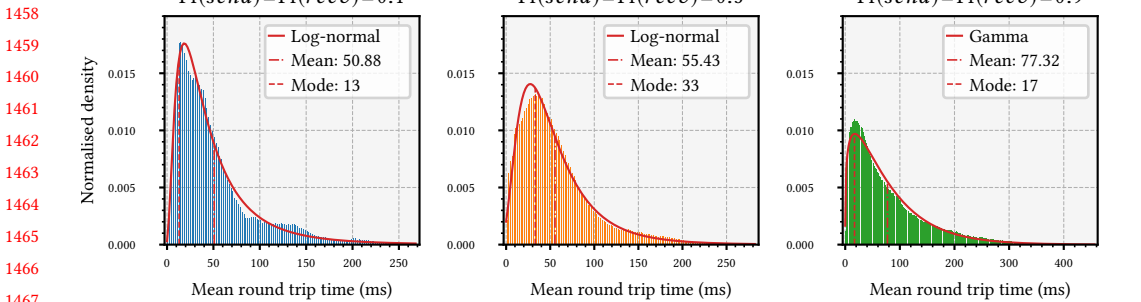
Fig. 11. Instrumentation pipeline for inline monitors

1422 *Validation of system model parameters.* Our SuS, when configured with steady load, models web
 1423 server traffic where the requests observed at the server are known to follow a Poisson process. The
 1424 probability distribution of the RTT of web application requests is generally right-skewed, and can
 1425 be approximated to a log-normal [Ciemiewicz 2001; Grove and Coddington 2005; Liu et al. 2001] or
 1426 Erlang (a special case of gamma) distribution [Kayser 2017]. We conduct three experiments using
 1427 steady loads fixed with $n = 10k$ and $w = 100$. $\Pr(\text{send}) = \Pr(\text{recv})$ are varied through 0.1, 0.5 and 0.9
 1428 to establish whether the RTT for our chosen set-ups resembles the aforementioned probability
 1429 distributions. Our results, summarised in fig. 12, were obtained as follows. The parameters for a
 1430 series of candidate probability distributions (e.g. normal, log-normal, gamma, etc.) were estimated
 1431 using Maximum Likelihood Estimation [Rossi 2018] on the RTT obtained from *each* experiment.
 1432 We then performed goodness-of-fit tests on these parametrised distributions, selecting the most
 1433 appropriate RTT fit for each of the three experiments. Our goodness-of-fit measure was derived
 1434 using the Kolmogorov-Smirnov test. The fitted distributions in fig. 12 indicate that the RTT of our
 1435 SuS confirms the findings reported in [Ciemiewicz 2001; Grove and Coddington 2005; Kayser 2017],
 1436 which show that web response times follow log-normal or Erlang distributions.

1437 *Experiment Precautions.* Further to the set-up parameters discussed in sec. 4.2, the following pre-
 1438 cautions were also taken:
 1439

- 1440 P_1 *Ten repeated readings.* The number of repeated readings to take was determined empirically
 1441 based on the coefficient of variation, $CV = \frac{\sigma}{\bar{x}} \times 100$, that was calculated for experiments with
 1442 different repetitions.
- 1443 P_2 $\Pr(\text{send}) = \Pr(\text{recv}) = 0.9$. Lower values of $\Pr(\text{send})$ and $\Pr(\text{recv})$ detract from the veracity of
 1444 the experiments because slaves become frequently idle.
- 1445 P_3 *Scheduler utilisation.* Sampled every 500 ms asynchronously, not to affect the SuS. Samples
 1446 were obtained using EVM function calls to get the most accurate reading. We did not measure
 1447 the CPU at the OS-level, because the EVM keeps scheduler threads momentarily spinning to
 1448 remain reactive, and this inflates the utilisation metric. This EVM feature could have been
 1449 switched off, but we decided to use the default settings and instead, measure the utilisation
 1450 internally.
- 1451 P_4 *Memory consumption.* Sampled every 500 ms asynchronously, not to affect the SuS. Samples
 1452 were obtained using EVM function calls to get the most accurate reading.
- 1453 P_5 *Mean RTT.* Sampled every 10% out of the total number of messages exchanged between
 1454 master and each slave. The sampling window of 10% was determined empirically via a series
 1455 of tests. The RTT is calculated as a running mean of *each* sample taken; the overall drift w.r.t.

1456
 1457



1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 Fig. 12. Fitted probability distributions on mean RTT for steady loads of 10k slaves

1471 the mean calculated over all samples was $\approx \pm 1.4\%$. We gathered the RTT thus to avoid as
1472 much as possible perturbations in the SuS that would arise due to data collection.

1473 P₆ *Weighted mean*. We aggregated the sampling records collected from repetitions of the same
1474 experiment using the weighted mean to account for the differing number of records counts
1475 that were obtained at each run.

1476 P₇ *Randomisation seed*. We fixed the randomisation seed to ensure experiment repeatability.
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519

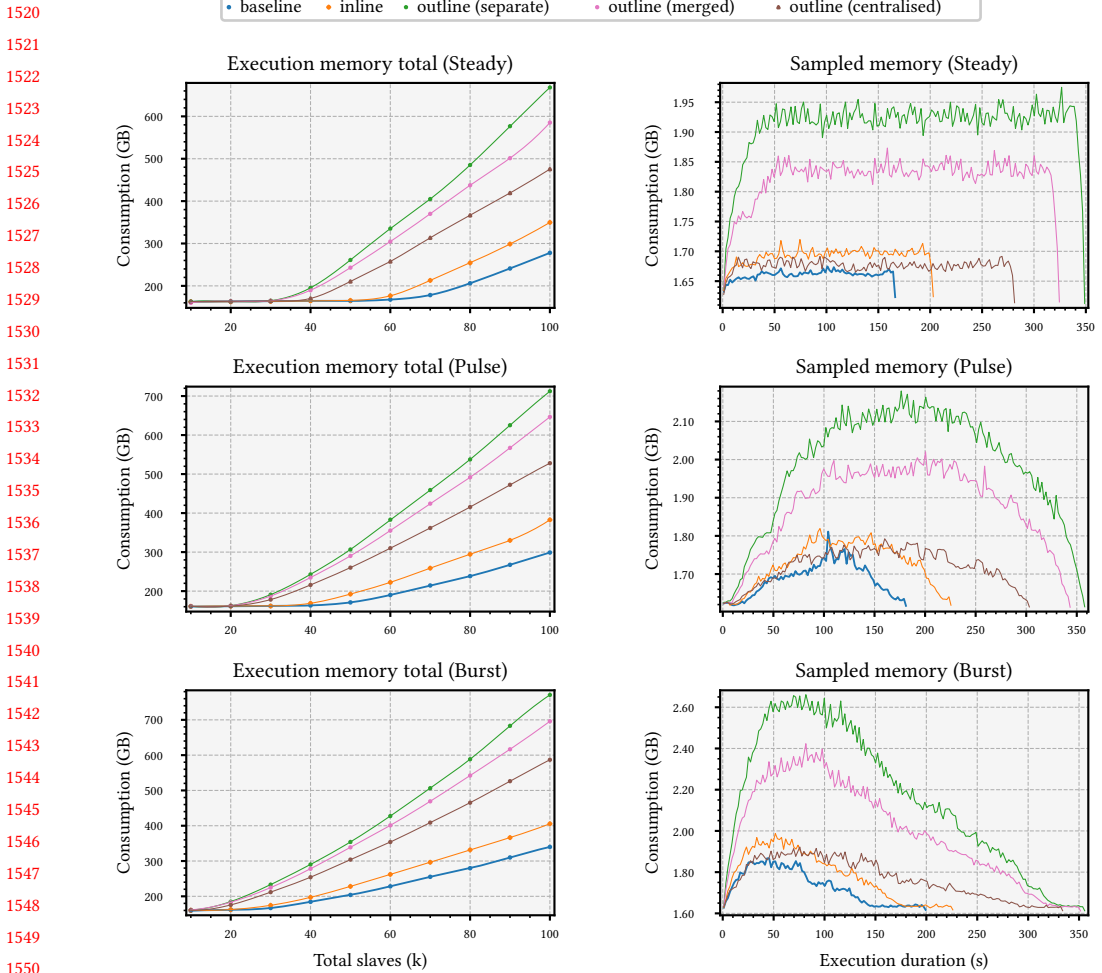


Fig. 13. Total and sampled memory consumption over entire execution duration with 100 k slaves

C SUPPORTING DATA PLOTS

The plots in figs. 14–17 have been fitted with linear, quadratic and cubic polynomials where the R^2 is above 0.96.

Total memory consumed. Fig. 13 shows the total memory consumed and sampled memory during the experiment runs conducted under Steady, Pulse and Burst loads for the case study with $n = 100k$ slaves. Note that unlike in figs. 9 and 10, the y -axis is labelled in GB. The total memory consumed plotted on the left in fig. 13 corresponds to the *area* under the sampled memory plots on the right. Decentralised outline monitoring consumes the most memory, while the centralised version falls midway between decentralised outline (separate) and inline monitoring. The sampled memory plots reveal that centralised outline monitoring consumes less memory than inline monitoring on average, but does so for a longer time period. This is especially noticeable in the Steady and Pulse plots, suggesting the memory overhead in centralised outline monitoring is induced in a more

1569 consistent manner. Our sampled memory plots reflect the shapes of the loads applied, although
1570 these extend for a longer duration that goes beyond the original loading time of $t = 100\text{s}$ (see fig. 6).
1571

1572

1573

1574

1575

1576

1577

1578

1579

1580

1581

1582

1583

1584

1585

1586

1587

1588

1589

1590

1591

1592

1593

1594

1595

1596

1597

1598

1599

1600

1601

1602

1603

1604

1605

1606

1607

1608

1609

1610

1611

1612

1613

1614

1615

1616

1617

1618 Moderate loads (fitted data plots).

1619

1620

1621

1622

1623

1624

1625

1626

1627

1628

1629

1630

1631

1632

1633

1634

1635

1636

1637

1638

1639

1640

1641

1642

1643

1644

1645

1646

1647

1648

1649

1650

1651

1652

1653

1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666

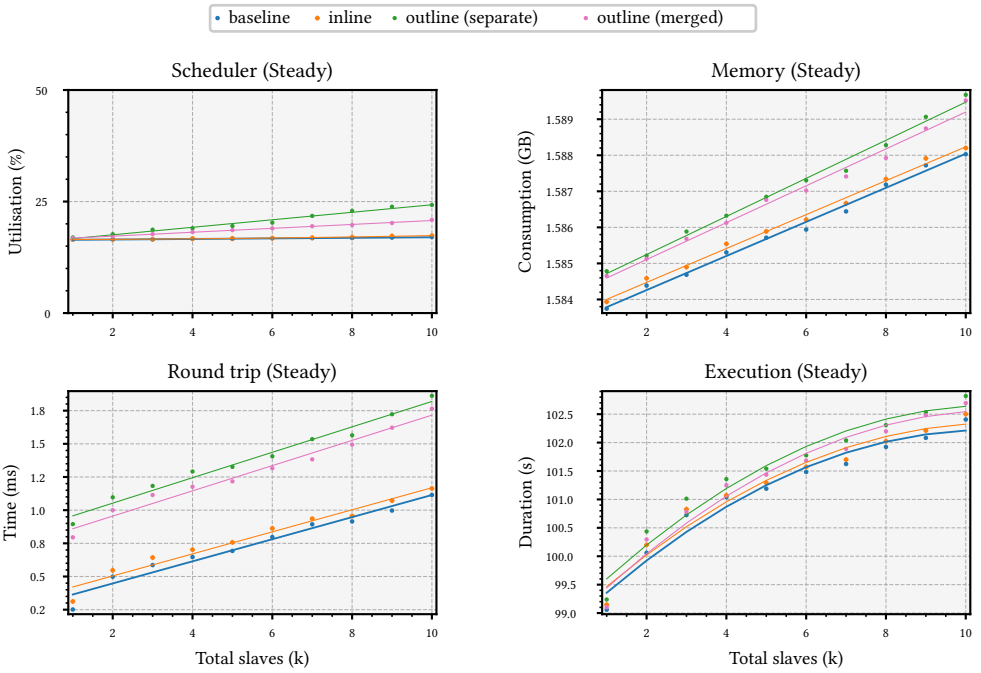


Fig. 14. Mean runtime overhead for monitoring the master and slave processes 10 k slaves

1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715

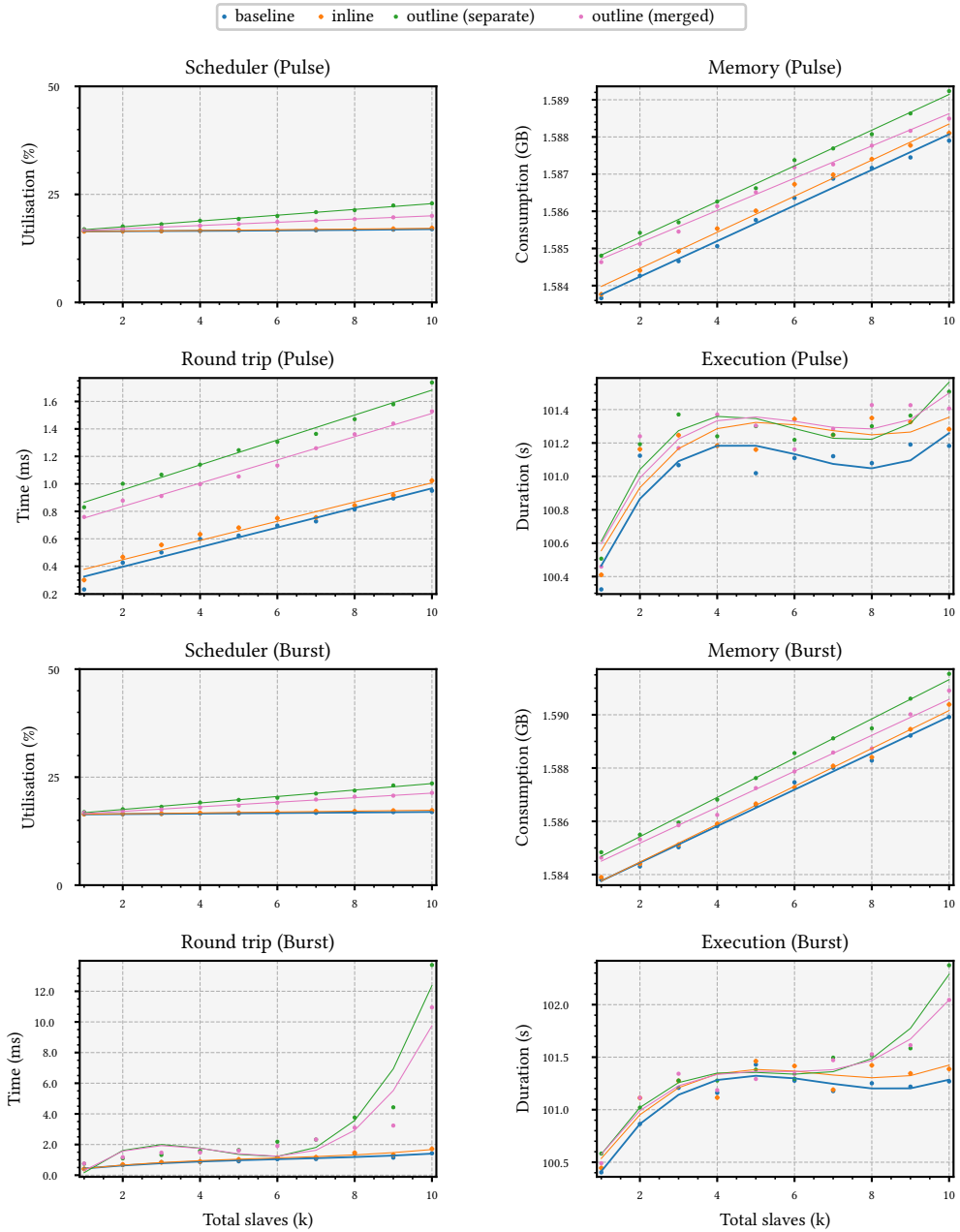


Fig. 15. Mean runtime overhead for monitoring the master and slave processes 10 k slaves (cont.)

1716 High loads (fitted data plots).

1717

1718

1719

1720

1721

1722

1723

1724

1725

1726

1727

1728

1729

1730

1731

1732

1733

1734

1735

1736

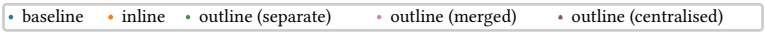
1737

1738

1739

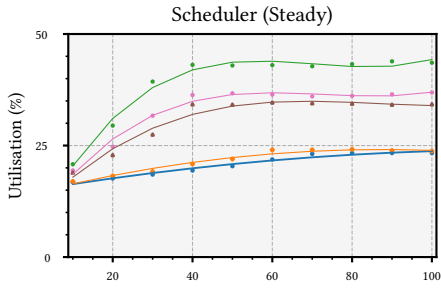
1740

1741



1742

1743



1744

1745

1746

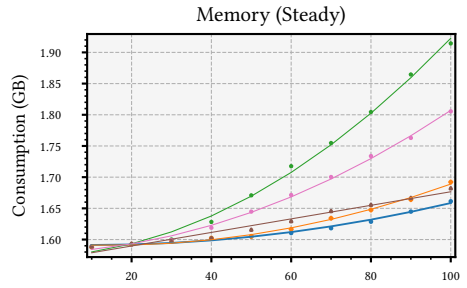
1747

1748

1749

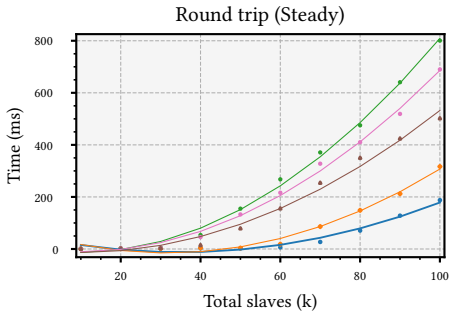
1750

1751



1752

1753



1754

1755

1756

1757

1758

1759

1760

1761

1762

1763

1764

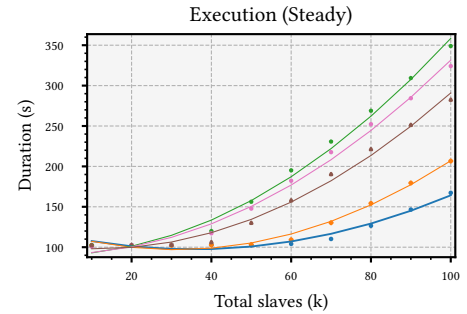


Fig. 16. Mean runtime overhead for monitoring the master and slave processes (100 k slaves)

1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813

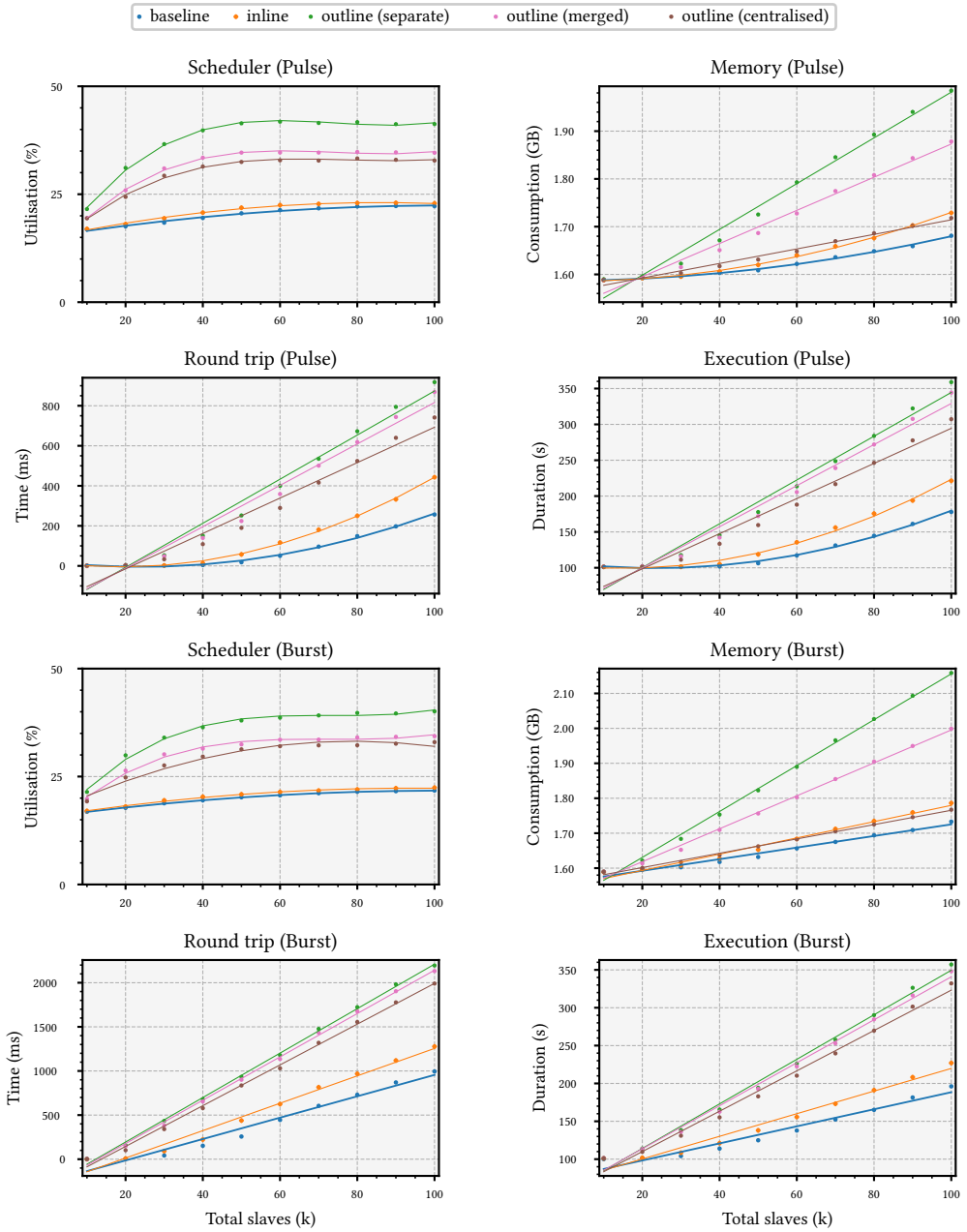


Fig. 17. Mean runtime overhead for monitoring the master and slave processes (100 k slaves, cont.)

1814 *Estimated overhead on slaves.*

1815

1816

1817

1818

1819

1820

1821

1822

1823

1824

1825

1826

1827

1828

1829

1830

1831

1832

1833

1834

1835

1836

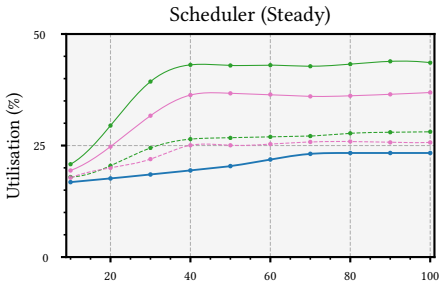
1837

1838

1839 • baseline • separate • merged • separate (master only, dashed) • merged (master only, dashed)

1840

1841



1842

1843

1844

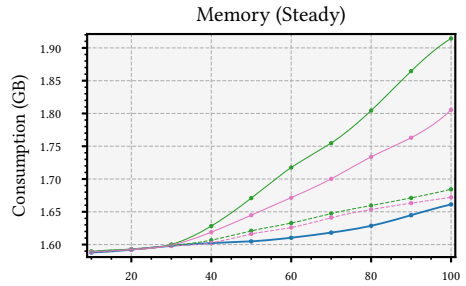
1845

1846

1847

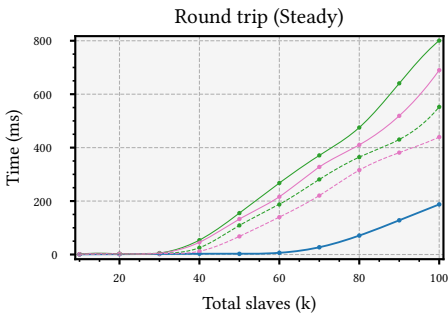
1848

1849



1850

1851



1852

1853

1854

1855

1856

1857

1858

1859

1860

1861

1862

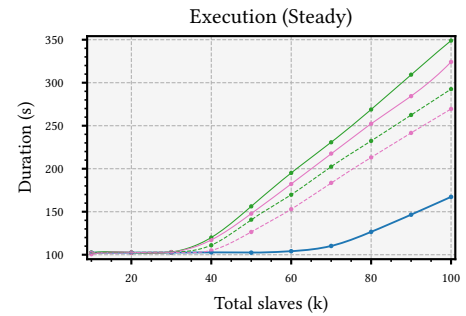


Fig. 18. Mean runtime overhead for monitoring the master process only (100 k slaves)

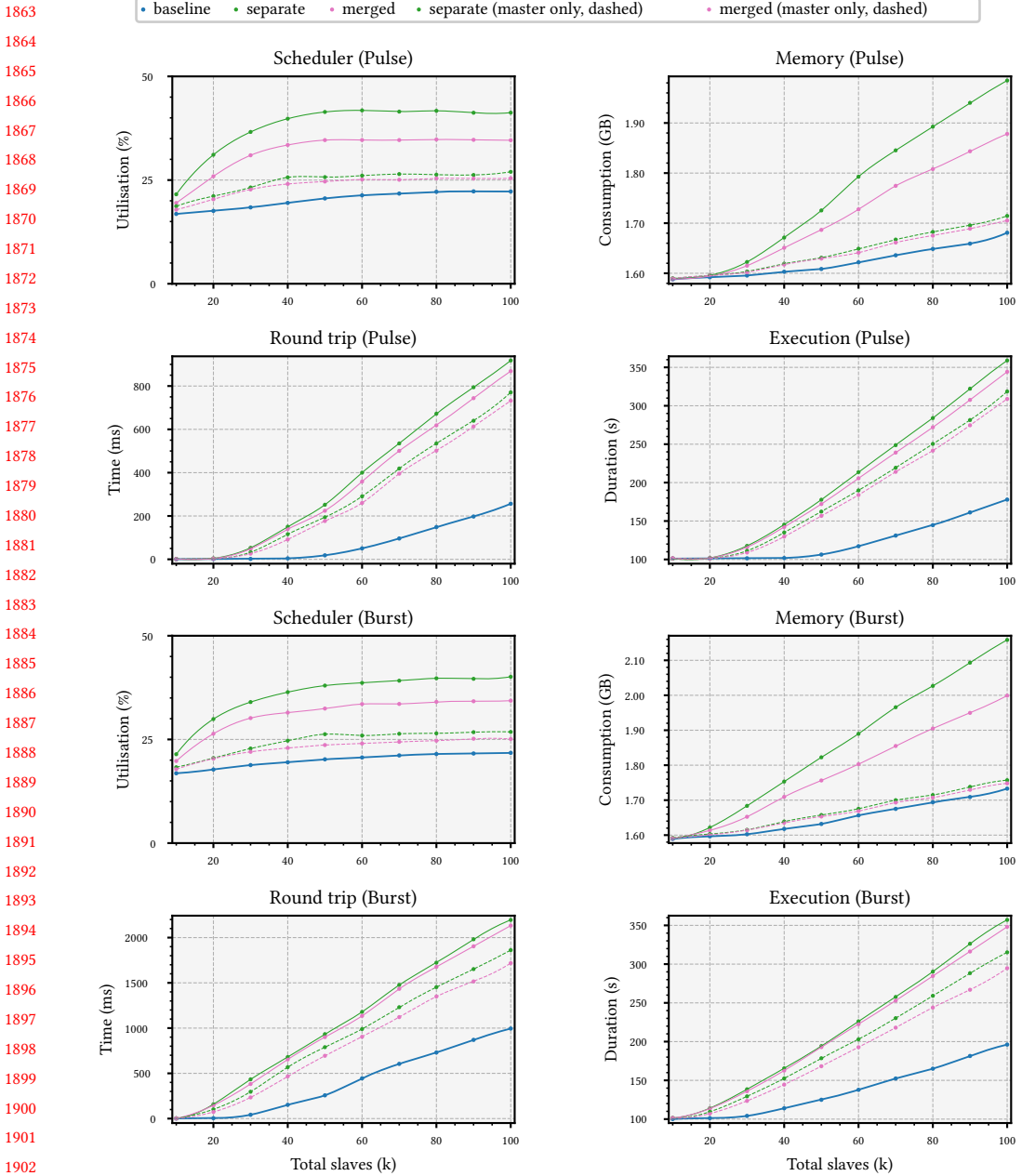


Fig. 19. Mean runtime overhead for monitoring the master process only (100 k slaves, cont.)

1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911