Proceedings of the 2nd Workshop on

Advanced tools, programming languages, and PLatforms for Implementing and Evaluating algorithms for Distributed systems (ApPLIED)

Held in conjunction with DISC 2019

Budapest, Hungary October 14, 2019

Editors: Yanhong Annie Liu, Miguel Matos

Contents

Foreword	3
Organization	4
Sponsors	5
Program	6
Invited Talk Abstracts and Speaker Bios	7
Invited Speakers	7
Building and Testing Byzantine Fault Tolerant State Machines	
Ethan Buchman, Interchain Foundation, Switzerland	8
Algorand: From Theory to Practice	
Jing Chen, Algorand, USA	25
Transactional Data Structure Libraries	
Idit Keidar, Technion, Israel	26
Weak Models for Distributed Computing	
Gadi Taubenfeld, Interdisciplinary Center, Israel	38
Approaches to Data Sharing in Edge FaaS	
Zoltán Turányi, Ericsson, Hungary	46
To build, or Not to Build, That Is the Question	
Nitin Vaidya, Georgetown University, USA	50
Invited Talk Summaries	61
Algorand: from Theory to Practice	
Jing Chen, Algorand Inc., USA	61
Panel and Discussion Summaries	66
Invited panel: Challenges and Current State	66
Discussion session:	
Consensus and Distribute Ledger: Scalability and Assurance $\ . \ . \ .$	66
Invited panel and open discussion: Summary and Directions	66
Short Papers	67
plcli - a Tool for Running Distributed Applications on PlanetLab	
Axel Niklasson, Chalmers University of Technology, Sweden \ldots	67

Foreword

It is our great pleasure to welcome you to the 2019 Workshop on Advanced tools, programming languages, and PLatforms for Implementing and Evaluating algorithms for Distributed systems — ApPLIED 2019. The purpose of this workshop is to bring together designers and practitioners of distributed systems from both academia and industry to share their point of views and experiences on implementing and evaluating distributed algorithms and systems. This second installment of the workshop was co-located with the 33rd International Symposium on Distributed Computing (DISC 2019), and took place on October 14th, 2019, in Budapest, Hungary.

The workshop featured keynote lectures, discussion panels, and presentations of short research papers.

The program included six keynote lectures by Ethan Buchman (Interchain Foundation, Switzerland), Jing Chen (Algorand, Inc.), Idit Keidar (Technion, Israel), Gadi Taubenfeld (Interdisciplinary Center, Israel), Zoltán Turányi (Ericsson, Hungary), Nitin Vaidya (Georgetown University, USA).

There were three panel and discussion sessions by invited panelists, other experts, and all participants, on (1) Challenges and Current State, (2) Consensus and Distributed Ledger: Scalability and Assurance, and (3) Summary and Directions.

Organization

General Chairs

- Chryssis Georgiou, University of Cyprus, Cyprus
- Elad Michael Schiller, Chalmers University of Technology, Sweden

Program Committee Chairs

- Yanhong Annie Liu, Stony Brook University, USA
- Miguel Matos, INESC-ID & IST, University of Lisboa, Portugal

Technical Program Committee

- Amy Babay, University of Pittsburgh, USA
- Ken Birman, Cornell University, USA
- Ioannis Chatzigiannakis, Sapienza University of Rome, Italy
- John Field, Google, New York City, USA
- Seif Haridi, Royal Inst. of Technology and Swedish Inst. of Computer Science, Sweden
- Wolfgang John, Ericsson Research, Sweden
- Marc Shapiro, Sorbonne-Universit-LIP6 and INRIA, France
- Srikumar Venugopal, IBM Research, Ireland

Sponsors

- University of Cyprus, Cyprus
- Chalmers University of Technology, Sweden
- Stony Brook University, USA
- Universidade de Lisboa & Angainor project (PTDC/CCI-COM/31456/2017), Portugal

Program

9:00	Opening and introduction: Chryssis Georgiou
	Morning Sessions: Chair: Miguel Matos
9:15	Invited talk: Idit Keidar Transactional Data Structure Libraries
10:00	Coffee break
10:30	Invited talk: Zoltan Turanyi
	Approaches to Data Sharing in Edge FaaS
11:15	Invited talk: Nitin Vaidya
	To build, or Not to Build, That Is the Question
11:45	Invited panel: Challenges and Current State
	Chryssis Georgiou, Zoltán Turányi, and Nitin Vaidya
12:30	Lunch
	Afternoon Sessions: Chair: Annie Liu
14:00	Invited talk: Ethan Buchman
	Building and Testing Byzantine Fault Tolerant State Machine
14:45	Invited talk: Jing Chen
	Algorand: From Theory to Practice
15:25	Discussion session
	Consensus and Distributed Ledger: Scalability and Assurance
16:00	Coffee break
16:30	Invited talk: Gadi Taubenfeld
	Weak Models for Distributed Computing
17:10	Student experience session: Axel Niklasson
	plcli - a Tool for Running Distributed Applications on PlanetLab
17:20	Invited panel and open discussion: Summary and Directions Ethan Buchman, Jing Chen, and Gadi Taubenfeld
17:50	Closing remarks

Machines

Invited Talk Abstracts and Speaker Bios

Invited Speakers

- Ethan Buchman, Interchain Foundation, Switzerland
- Jing Chen, Algorand, USA
- Idit Keidar, Technion, Israel
- Gadi Taubenfeld, Interdisciplinary Center, Israel
- Zoltán Turányi, Ericsson, Hungary
- Nitin Vaidya, Georgetown University, USA

Building and Testing Byzantine Fault Tolerant State Machines

Ethan Buchman, Interchain Foundation, Switzerland

Abstract: Building and testing fault tolerant state machines, let alone Byzantine Fault Tolerant (BFT) ones, is notoriously tricky business. The handful of industrial solutions, like Apache Zookeeper, CoreOS's etcd, and Hashicrop's Consul, are not BFT, and their state machines are restricted to relatively simple service-discovery oriented key-value stores. Blockchain solutions, on the other hand, add a myriad of complex state machines, including new virtual machine designs, but fail to provide adequate performance and mature development environments.



Here we present Tendermint, a production-grade Byzantine Fault Tolerant State Machine Replication engine written in Go. Tendermint supports BFT replication for state machines written in any language by using a socket protocol to communicate between the state machine and the replication engine, allowing applications to be built and tested in a developer's language of choice. Tendermint is being used on the public Internet today to secure upwards of 1 Billion USD in value, with deployments supporting hundreds of consensus nodes. Here we provide an overview of the Tendermint system and how to build Byzantine Fault Tolerant applications in Go and Javascript.

Biography Ethan is an Internet Biophysicist. He received a B.Sc. in Physical Science and a M.Sc. in Engineering Systems and Computing, both from the University of Guelph. With background in cell biology, neuroscience, mathematics, machine learning, and distributed computing, his goal is to build tools that encourage humans to self-organize into functional systems, much like molecules managed to self-organize into Life. Ethan is a co-founder of the Tendermint and Cosmos projects, and currently serves as Technical Director of the Interchain Foundation, a non-profit with a mission to research, develop, and promote open, decentralized networks. He is focused primarily on building general purpose technology to support interoperable replicated state machine on the public Internet, and in proving their correctness. He also runs CoinCulture CryptoConsulting, which offers in-depth courses on blockchain technology for both non-technical and technical audiences.

















Core Module: Tendermint

• Mature, state-of-the-art BFT consensus & networking protocol

• Vertical scaling solution for your base-layer chain

• Instant confirmation: 1 block finality

• 170+ validators on latest testnet

• Formally proven

Core Module: Governance

• On-chain proposal system:

- Text Updates
- Automatic Parameter Changes

Software Upgrades

• Liquid democracy for bonded stakeholders



vernance

Core Module: Proof of Stake

• Use your own coin for slashing and staking mechanics

• Maximize decentralization with delegation (through skin in the game)

• Mitigate long range attacks with an unbonding period



Core Module: Rewards & Fees

• Use your own coin to pay fees when running a transaction - no lock-in.

• Multi-coin support for improved user experience.



Core Module: IBC

• The Inter-Blockchain Communication (IBC) Protocol enables the Internet of Blockchains

• Move coins & data between different blockchains

• Basis for horizontal scaling and interoperability

• Analogous to TCP/IP for the Internet



Modular Blockchains















Algorand: From Theory to Practice Jing Chen, Algorand, USA

Abstract: Blockchains stand to revolutionize the way a modern society operates. They can secure all kinds of traditional transactions, such as payments, in the exact order in which the transactions occur; and enable totally new transactions, such as cryptocurrencies and smart contracts. They can remove intermediaries and usher in a new paradigm for trust. As currently implemented, however, blockchains scale poorly and cannot achieve their enormous potential. Algorand is the first blockchain that is truly secure, scalable and decentralized. It is permissionless and works in a highly asynchronous environment. It dispenses with proof of work and miners and requires only a negligible amount of computation. Moreover, its transaction history does



not fork, guaranteeing immediate finality of a transaction the moment the transaction enters the blockchain. In this talk, I will introduce Algorands core technology, recent development and roadmap.

Biography: Jing Chen is Chief Scientist and Head of Theory Research at Algorand, and Assistant Professor in the Computer Science Department at Stony Brook University. Her main research interests are distributed ledgers, game theory, and algorithms. Jing received her bachelor and masters degrees in computer science from Tsinghua University, and her PhD in computer science from MIT. She did a one-year postdoc at the Institute for Advanced Study, Princeton. Jing received the NSF CAREER Award in 2016.

Transactional Data Structure Libraries Idit Keidar, Technion, Israel

Abstract: We introduce transactions into libraries of concurrent data structures: such transactions can be used to ensure atomicity of sequences of data structure operations. By focusing on transactional access to a well-defined set of data structure operations, we strike a balance between the ease-ofprogramming of transactions and the efficiency of custom-tailored data struc-We exemplify this concept tures. by designing and implementing a library supporting transactions on any number of maps, sets (implemented as skiplists), queues, stacks, producerconsumer pools, and logs. Our library



offers efficient and scalable transactions, which are an order of magnitude faster than state-of-the-art transactional memory toolkits.

We further introduce nesting into our transactional data structure library. Nested transactions create checkpoints within a longer transaction, so as to limit the scope of abort. We then conduct a case study of pipelined network intrusion detection. In this benchmark, nesting improves throughput by up to 15x. Finally, we discuss cross-library nesting, namely dynamic composition of transactional data structure libraries.

(Based on joint works with Alexander Spiegelman, Gal Assa, Guy Golan-Gueta, and Hagar Meir.)

Biography: Idit Keidar received her B.Sc. (summa cum laude), M.Sc. (summa cum laude), and Ph.D. from the Hebrew University of Jerusalem in 1992, 1994, and 1998, respectively. She was a Postdoctoral Fellow at MITs Laboratory for Computer Science. She is currently a Professor at the Technions Viterbi Faculty of Electrical Engineering, where she holds the Lord Leonard Wolfson Academic Chair. She serves as the Head of the Technion Rothschild Scholars Program for Excellence, and also heads the EE Facultys EMET Excellence Program. Her research interests are in fault-tolerant distributed and concurrent algorithms and systems, theory and practice. Recently, shes mostly interested in distributed storage and concurrent data structures and transactions.



Concurrent Data Structure Libraries (CDSLs) Are They Really Thread Safe?? balance ←map.get (key=Yoni) Each operation executes atomically balance <- map.get(key=Yoni) Custom-tailored implementation preserves semantics newBalance ← balance + deposit newBalance ← balance + deposit map.set(key=Yoni, newBalance) balance←map.get (key=Yoni) Ann map.set(key=Yoni, newBalance) repQ.enq("Yoni's balance=new") newBalance ← balance + deposit repQ.eng("Yoni's balance=new") map.set(key=Yoni, newBalance) · Oops! Atomic operations are not enough repQ.enq("Yoni's balance=new") Software Transactional Memory (STM) CDSL vs STM Why? • TL2, TinySTM, SwissTM,... Transactions (TXs) = atomic sections including multiple DS CDSL STI operations Performance 1 **Exploit DS semantics** √ Begin_TX x Арр balance←map.get(key=Yoni) Used in practice 1 £ newBalance ← balance + deposit Generality × map.set(key=Yoni, newBalance) Composability x repQ.enq("Yoni's balance=new") End_TX STM Overhead Explained Agenda • Common STM solution: Motivation • Each object has a version • Concurrent Data Structure Libraries (CDSLs) vs Transactional Memory • TX reads a "consistent snapshot" • Introducing: Transactional Data Structure Libraries (TDSL) · Validates versions have not changed by commit time • Example TDSL algorithm Otherwise aborts TX updates occur during commit Skiplist Additional objects • Sources of overhead: • Fast abort-free singletons • Global version clock (GVC) ⇒ contention • Read- and write-sets \Rightarrow tracking and validation overhead Library composition & nesting • Conflicts \Rightarrow aborts • Evaluation







<section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></section-header>	<section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></section-header>
Logs • Optimistic read(pos), hasNext(pos) • If hasNext(pos) returns false - track in read-set "last = pos" • No other tracking • Pessimistic append(val) • Lock the entire log • Track appends in write-set until commit time	 Producer-Consumer Pools Fine-grain: multiple slots Each with its own lock Pessimistic produce/consume Lock only the affected slot Track locally until commit Optimistic consume from earlier produce Produce and consume cancel each other out
<section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><section-header><section-header><section-header><section-header><section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></section-header></section-header></section-header></section-header></section-header></list-item></list-item></list-item></list-item></list-item></list-item></list-item></section-header>	<section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></section-header>

<section-header><section-header><section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></section-header></section-header></section-header>	<section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></section-header>
<section-header><list-item><list-item><list-item><list-item><list-item><list-item><table-container></table-container></list-item></list-item></list-item></list-item></list-item></list-item></section-header>	Transaction Adjustments • Transactions validate that no node in read-set has • a version equal to the transaction's snapshot version & • a true singleton bit • If yes, abort and increment GVC before retrying
<section-header><section-header><section-header><section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></section-header></section-header></section-header></section-header>	<section-header><section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></section-header></section-header>








- New concept for concurrent programing
 Composable DSs supporting TX as well as fast singletons
- Support for nesting
- A prototype library with a host of DSs
 - Map (based on skiplist), queue, stack, log, pool
- We hope that the community will adopt this concept • Build and use more such libraries

Weak Models for Distributed Computing Gadi Taubenfeld, Interdisciplinary Center, Israel

Abstract: The talk has three parts. In the first part, I will show how to model the process of genome-wide epigenetic modifications, which allows cells to utilize their DNA, as an anonymous shared memory system. This is done by formulating a particular consensus problem and presenting algorithms for solving the problem. In the second part, I will discuss results for anonymous shared memory systems which are composed of shared objects for which there is no a priori agreement between the processes on the names of the objects. In the third part, I will motivate and explore the new notion of weak failures, which should be viewed as fractions of traditional failures.



Biography: Gadi Taubenfeld is a professor and past dean of the School of Computer Science at the Interdisciplinary Center in Herzliya, Israel. He is an established authority in the area of concurrent and distributed computing and has published widely in leading journals and conferences. He authored the books "Synchronization Algorithms and Concurrent Programming" and "Distributed Computing Pearls". His primary research interests are in concurrent and distributed computing. Gadi was the head of the computer science division at Israel's Open University; member of technical staff at AT&T Bell Laboratories; consultant to AT&T Labs - Research; and a research scientist and lecturer at Yale University. Gadi served as the program committee chair of PODC 2013 and DISC 2008 and holds a Ph.D. in Computer Science from the Technion - Israel Institute of Technology.















Approaches to Data Sharing in Edge FaaS Zoltán Turányi, Ericsson, Hungary

Abstract: Function-as-a-service systems are stateless by nature - every function starts with no memory of what happened before. Such context is typically placed in an external database, where functions can read and update as side effect. The time to execute functions significantly depends on access time to this database - accessing data in a different location is much slower than locally. It is thus beneficial to co-locate data and execution. This is especially relevant for the Edge scenario, where remote means geographically remote. In this talk we overview some approaches, how such placement can be combined with methods to ensure data consistency among various locations. Our key goals are high performance and ease of use.



Biography: Zoltn Richard Turnyi is currently an expert of 5G Network Architectures within Ericsson Research. He currently works on a Function-as-a-Service concept built on a fast, distributed, in-memory key-value store. His background is in IP networking and mobile core with recent addition of Software Defined Networking and Network Function Virtualization. He is the author of more than 50 patent applications and works with Ericsson Research for more than 20 years. He holds an M.Sc. in Computer Science from the Technical University of Budapest. He is a scout master for 25 years.







To build, or Not to Build, That Is the Question Nitin Vaidya, Georgetown University, USA

Abstract: The talk will explore necessity and sufficiency of experimental evaluations, using examples of past projects, and make suggestions for future directions in experimental research in distributed computing.

Biography: Nitin Vaidya is the Robert L. McDevitt, K.S.G., K.C.H.S. and Catherine H. McDevitt L.C.H.S. Chair of Computer Science at Georgetown University. He received Ph.D. from the University of Massachusetts at Amherst. He previously served as a Professor and Associate Head in Electrical and Computer Engineering at the University of Illinois at Urbana-



Champaign. He has co-authored papers that received awards at several conferences, including 2015 SSS, 2007 ACM MobiHoc and 1998 ACM MobiCom. He is a fellow of the IEEE. He has served as the Chair of the Steering Committee for the ACM PODC conference, as the Editor-in-Chief for the IEEE Transactions on Mobile Computing, and as the Editor-in-Chief for ACM SIGMOBILE publication MC2R.



















How to unwind this clock?	<section-header><section-header><section-header><section-header><section-header><section-header><image/></section-header></section-header></section-header></section-header></section-header></section-header>
Minimalism	Litmus Test
 Often less is more Don't build just because you can There may be better things to do with your time and resources 	 Would you be willing to publicly post the <i>exact problem statement</i>? before developing the solution If not, find something better to do
57	
17 5 59	
Thanks!	
disc.georgetown.domains	

Invited Talk Summaries

Algorand: from Theory to Practice Jing Chen, Algorand Inc., USA

Algorand: from Theory to Practice

Jing Chen Algorand Inc., Boston, MA 02116 jing@algorand.com

Abstract

A summary of my talk at the ApPLIED Workshop at DISC 2019.

Introduction

Blockchains stand to revolutionize the way a modern society operates. They can secure all kinds of traditional transactions, such as payments, in the exact order in which the transactions occur; and enable totally new transactions, such as cryptocurrencies and smart contracts. They can remove intermediaries and usher in a new paradigm for trust. As currently implemented, however, blockchains scale poorly and cannot achieve their enormous potential. Algorand is the first blockchain that is truly secure, scalable and decentralized. It is permissionless and works in a highly asynchronous environment. It dispenses with "proof of work" and "miners" and requires only a negligible amount of computation. Moreover, its transaction history does not "fork", guaranteeing immediate finality of a transaction the moment the transaction enters the blockchain. In this talk, I will briefly introduce Algorand's core technology, recent development and roadmap. The readers may refer to [7, 8, 4, 6] for more details.

Underlying the Algorand blockchain is a new Byzantine Agreement protocol that is highly efficient. It works under an adversarial model where the adversary can dynamically corrupt any user at any time, control the actions of a corrupted user, and perfectly coordinate the actions of all corrupted users. Even with such a strong adversary, the protocol achieves asynchronous safety and guarantees that the Algorand blockchain doesn't fork even when the underlying propagation network is partitioned. Accordingly, any transaction that appears on the blockchain is immediately final and can be relied upon, without the need of waiting for more blocks being added after it. On the other hand, the protocol achieves liveness as long as messages propagated by honest users are received by other honest users within a known time bound.

Another important idea that makes the Algorand blockchain scalable is cryptographic self-selection. Indeed, having millions of users participate in the Byzantine Agreement in order to select the next block is unrealistic. One possibility is to *publicly* select, at random, a subset of users to form a *committee* and participate on behalf of everybody. However, once the identities of the committee members become public, the adversary can corrupt them so that they behave maliciously. Instead, the Algorand blockchain has each user self-select herself into a committee. Thanks to cryptographic primitives such as unique signatures, cryptographic hash functions and verifiable random functions, a user can privately generate a unique "lottery ticket", which she can use to prove her membership in the committee if she is selected, but cannot cheat and convince others to accept her proof otherwise. When participating in the Byzantine agreement, a committee member propagates her winning ticket together with her proposal or voting message. No other communication is needed to find out who is selected. In this way, the adversary learns the fact that a particular user is selected only after the user has sent out her message in the protocol instead of before, and corrupting the user doesn't let the adversary control the message being sent out.

One more idea is needed here. If a selected committee participates in multiple steps of the Byzantine agreement, the adversary can learn the identities of its members and corrupt all of them after seeing their first messages, so that they behave maliciously in remaining steps. The Algorand blockchain is immune to this problem because its Byzantine agreement has an important property referred to as *user replaceability*: the protocol doesn't rely on users keeping private states, and the message that a user should send in a step can be determined solely based on messages that have been propagated to him/her in previous each steps. As such, the protocol has a committee randomly and independently selected for every step. Corrupting the committee members for a step does not give the adversary more power than random corruption in terms of controlling committee members for future steps.

In order to deal with Sybil attacks, the selection probability is the same for every token: that is, in effect, tokens are selected at random, and users that own the selected tokens participate in the Byzantine agreement. The users do not need to delegate their right of participation to a small group of super nodes, neither do they need to lock up their tokens for a long time in order to participate in the consensus protocol. Indeed, the approach here is a pure form of proof-of-stake.

Many other ideas have been introduced in the Algorand blockchain, but rather than covering them all today, I'd like to report on some recent developments on actually implementing the blockchain and putting it to work. The Algorand MainNet launched in mid-June 2019, and has been running smoothly since. The TestNet has been running since April, 2019; it runs the same version of the blockchain as the MainNet, so that developers can test their software (such as wallets) before running them on the MainNet. The code has been audited by third parties, and the entire code repository is open-sourced and available at [3]. Various tools for developers, such as SDKs for multiple programming languages and tutorials, can be found at [2]. We continue to enlarge the tool set, and have also launched a bug bounty program [1]. In addition to a pen-and-paper analysis, with collaborators at Runtime Verification, we have begun formal verification of the consensus protocol using the Coq proof assistant. Our model explicitly incorporates timing issues and adversarial actions, reflecting a more realistic environment that may be faced by a public blockchain. We have proved asynchronous safety of the protocol under this model, and a paper reporting on the progress is available at [5].

Since the MainNet launch, we continue to develop new technology to enable important applications on the blockchain. For example, in forthcoming versions, users will be able to issue their own fungible tokens directly in layer-1 of the blockchain. Moreover, users will be able to clear multiple transfers, among arbitrary sets of users and for arbitrary sets of layer-1 currencies, in a single transaction: that is, in a truly atomic way without relying on devices such as hashed time-locks. And there is still more to come. Interested readers can find the most up-to-date information at https://www.algorand.com/. Stay tuned!

References

- [1] https://bugcrowd.com/algorand.
- [2] https://developer.algorand.org/.
- [3] https://github.com/algorand/go-algorand.

- [4] Algorand blockchain features. https://github.com/ algorandfoundation/specs/blob/master/overview/Algorand_ v1_spec-2.pdf, 2019.
- [5] M. A. Alturki, J. Chen, V. Luchangco, B. Moore, K. Palmskog, L. Peña, and G. Roşu. Towards a verified model of the Algorand consensus protocol in Coq. In *FMBC'19: Workshop on Formal Methods for Blockchains, 3rd Formal Methods World Congress*, 2019.
- [6] J. Chen, S. Gorbunov, S. Micali, and G. Vlachos. Algorand Agreement: Super fast and partition resilient Byzantine agreement. Cryptology ePrint Archive, Report 2018/377, 2018. https://eprint.iacr.org/2018/377.
- [7] J. Chen and S. Micali. Algorand: A secure and efficient distributed ledger. *Theoretical Computer Science*, 777:155–183, 2019.
- [8] Y. Gilad, R. Hemo, S. Micali, G. Vlachos, and N. Zeldovich. Algorand: Scaling Byzantine agreements for cryptocurrencies. In SOSP, pages 51– 68, 2017.

Panel and Discussion Summaries

These sessions were well attended and had live and involved discussions. During each session, discussion questions are addressed by each panelist or expert and also discussed among all participants.

Various aspects of the current state were discussed, from tools for simulation to work on verification. The conclusion is that, overall, a lot of progress is being made, but significant future work is needed.

Invited panel: Challenges and Current State

Panelists: Chryssis Georgiou, Zoltán Turányi, and Nitin Vaidya

Chair: Miguel Matos

Discussion questions:

- 1. On the processes and practices from designing an algorithm to developing a prototype implementation suitable for running on real systems. What are the current practices and best current practices?
- 2. On languages, libraries, platforms, and tools for distributed algorithms and systems, addressing correctness, performance, ease of use, and result reproducibility. What are the current ones and best current ones?

Discussion session: Consensus and Distribute Ledger: Scalability and Assurance

Chair: Annie Liu

Discussion questions:

- 1. What is the current state and problems for scalability?
- 2. What is the current state and problems for assurance?

Invited panel and open discussion: Summary and Directions

Panelists: Ethan Buchman, Jing Chen, and Gadi Taubenfeld

Chair: Annie Liu

Discussion questions:

- 1. What are the learned lessons on the design and evaluation of distributed algorithms and systems? Include both positive and negative lessons.
- 2. What are future directions? Consider both short term and longer terms.

Short Papers

plcli - a Tool for Running Distributed Applications on PlanetLab

Axel Niklasson, Chalmers University of Technology, Sweden

plcli - a Tool for Running Distributed Applications on PlanetLab

Axel Niklasson axelni@student.chalmers.se Department of Computer Science and Engineering Chalmers University of Technology, Sweden

Abstract

Implementing and testing distributed applications on platforms with many servers such as PlanetLab is made easier through the help of tooling. Various solutions exist, but most tools are either outdated or lack sufficient documentation. In this paper, a tool referred to as *plcli* is introduced which, among other things, enables engineers working with PlanetLab to deploy and run distributed applications on PlanetLab servers. The intended functionality of *plcli* has been summarised in three use cases; running distributed experiments, distributed debugging and node health monitoring. Furthermore, a pilot implementation written in the programming language Go is offered with this paper that implements support for experiment deployment. To see how *plcli* performs, an evaluation has been carried out that has shown that the deployment time is nearly kept constant as the number of application instances and servers are scaled up to as much as 120 instances on fifteen servers. For example, a 400% increase in the number of servers only resulted in a 7% increase in deployment time.

1 Introduction

When implementing and testing applications that are meant to run as distributed systems, the need for tooling to deploy these applications arise. In order to understand how the application in question performs as a real system, it must be deployed as such; while simulating a physically distributed system using tools such as NS-3 [1] [2] is a great alternative during the development phase, testing must be performed in a real-world environment as well. Testing applications for production usage often requires ten or more servers to be part of the deployments, which emphasises the importance of automating the process of deployment. Manually connecting to servers or residing to ad-hoc implemented scripts is not a scalable solution, which is what this kind of tooling aims to provide.

This paper focuses on applications and experiments run on the PlanetLab EU platform [3]. PlanetLab is a platform used for deploying, running and accessing distributed applications in a planetary-scale system [4] [5] [6]. More than 300 universities and research institutes, referred to as *sites*, are providing servers to the network. These servers are called *nodes* and are what makes out the computing capacity of the enormous cluster that is PlanetLab. Users are assigned a *slice*, which essentially is a distributed virtualised virtual machine, that they can use to deploy services to. Since nodes are provided by the different sites, no guarantees on homogeneity on the machines can be given and downtime is to be expected. The varying quality of nodes makes PlanetLab a suitable platform for testing systems in a real-world setting. However, users of this platform need to make sure that desired nodes are actually functioning as intended, which can at times be tedious work.

1.1 Related work

There are public user tools listed on the PlanetLab website that may be used to decrease workload and enable easier usage when working with PlanetLab [7]. The listed tools offer various functionalities such as slice management, package management and others. Most of the tools (8/12) are not available anymore and the list is gravely outdated, but two tools that are still accessible and exhibit similar functionality as *plcli* are *PLDeploy* and *pssh*. PLDeploy functions as a "utility to deploy, configure and control PlanetLab services" which is rather similar to Use Case 1, presented in Section 1.2. However, when deploying services using PLDeploy, it needs to be done in a very special fashion by constructing and attaching what is called *cogs* that are used to deploy services and pulling the results back. There is not much information related to this tool and its documentation has not been updated in the last decade, making a more in-depth comparison hard to perform. pssh on the other hand provides a parallel version of OpenSSH and related tools [8] and its main features of providing parallel execution of commands over ssh is also present in *plcli*. For example, when users want to run a command on several nodes, this is done in a concurrent fashion without the users having to consider it.

1.2 Use Cases

The first use case, **Use Case 1**, targets support for deploying an application to a given number of physical PlanetLab nodes, launching a specified amount of instances on each node and upon termination, gathering of log files. As outlined in Section 2, engineers using *plcli* are able to quickly deploy their experiments through adding a configuration file to their application repository. Use Case 1 is considered the main feature of *plcli* and is the only one implemented in the preliminary version of *plcli*.

Apart from running experiments, being able to carry out distributed debugging is also an attractive feature and is expressed as **Use Case 2**. Finding out what is happening in a distributed system is a very complex task, as clearly outlined by Joyce *et al.* [9] and support for distributed debugging could increase engineering productivity a great deal as well as aiding in finding out why certain problems occur in production systems. For example, an approach similar to D3S presented by Liu *et al.* in [10] or the solution based on the MINHA platform presented by Jorge *et al.* [11] could be taken to implement support for this use case.

The duration of experiments might range from a few minutes to hours or even days and it is important that these experiments are performed on healthy nodes. Features such as healthy node discovery and automatically fixing some of the more simple problems of nodes (for example problems that may be resolved through a simple reboot) could be added to *plcli* to provide a richer toolset when working with PlanetLab nodes, which is referred to as **Use Case 3**. This could be taken even further by considering the PlanetLab platform as a system in itself and deploy planning agents that make decisions based on repair plans and carry out node repairs, which is an approach heavily based on the one presented by Dashofy *et al.* [12].

1.3 Contribution

In this paper a tool for deploying and running distributed applications using the PlanetLab platform known as *plcli* has been introduced, along with three main use cases representing the three main features of the tool. A preliminary implementation of *plcli* is bundled with this paper which provides functionality for Use Case 1. An evaluation of the performance of the tool with respect to Use Case 1 has been carried as well, which is presented and discussed. An implementation satisfying Use Case 1 is offered alongside this paper [13].

2 System Design

plcli is a command-line interface written in the programming language Go, which is a programming language introduced by Google in 2009, designed for fast compilation of source code and easier programming [14] [15] [16]. The Go programming language provides functionality for implementing efficient applications with scalable concurrency mechanisms known as *goroutines*. Goroutines are essentially lightweight threads associated with less overhead and the Go runtime is very efficient in the handling of these goroutines. As shown by Togashi *et al.*, Go outperforms for example Java when it comes to concurrency handling [17]. Mechanisms for efficient concurrency handling are naturally of high interest when developing a tool such as *plcli* that must be able to communicate with tens, or even hundreds, of nodes efficiently. These goroutines are further referred to as *workers* and are used whenever there are I/O-bound operations that need to be executed concurrently, such as calling the PlanetLab API or executing commands on nodes over SSH.

plcli is a preliminary implementation of a full-fledged tool intended to provide support for all three use cases outlined in Section 1.2. Through the implementation of support for Use Case 1, various features have been added to *plcli* such as deployment of applications, file transfer to nodes as well as concurrent command execution. An extensive list of available commands at the time of writing can be found in the README in the project repository [13].

Since the main functionality of the current implementation is to deploy code, a more indepth explanation of how a deployment is performed is provided. *plcli* makes use of what is called a configuration file which is required to be placed in the root of the public repository of applications that should be deployed using *plcli*. This file contains information about environment variables, how the application is prepared for launch and how to launch an instance. *plcli* downloads this file from the application repository and performs the needed steps in order to prepare the PlanetLab nodes for deployment of the given application. The structure of this file is not finalised at the time of writing and omitted for brevity. However, an example can be found in the GitHub repository for the demo application [18].

plcli aims to reduce the time and effort needed to deploy and run experiments on PlanetLab and consequently, it is highly dependant on the performance of the PlanetLab API. In order to retrieve information about what nodes are available and decide which nodes to use in a deployment, the API is queried for up to date information. However, the API is only used internally by *plcli* in an effort to enable users and researchers to focus on what experiments to run rather than how these experiments actually are run.

3 Evaluation Plan

In order to investigate the performance of *plcli*, the time taken to perform a full deployment - the *deployment time* - is evaluated in three different experiments. All three experiments

presented below utilise a basic demo application written in Python 3.7, which can be found on GitHub [18] and were run on a MacBook Pro 15" with a 2,2 GHz Intel Core i7 processor and 16 GB of RAM.

In the first experiment, one application instance was deployed to a varying amount of physical nodes with one worker allocated per node. This was done to investigate changes in deployment time as the deployment grows with respect to the number of nodes. Furthermore, experiments with one instance deployed to a fixed amount of nodes with a varying amount of workers were also conducted, which aided in investigating the performance gain of using workers. Lastly, due to severe problems with finding more than fifteen suitable nodes for the demo application, a varying amount of application instances were launched on the same set of physical nodes to try and investigate the performance at scale.

4 Evaluation Results

The first experiment measuring the performance of *plcli* when deploying to an increasing amount of physical nodes exhibits a nearly constant deployment time, as can be seen in Figure 1. There is a small increase of the trendline as the number of nodes increase, but it is very small and indicates that the overhead introduced by adding more nodes than fifteen will be minimal. For example, the time taken to deploy to three nodes is around fifteen seconds while deploying to fifteen nodes takes around sixteen seconds; that is a 400% increase in the number of nodes with merely 7% increase in deployment time which indicates promising performance when scaling. These results are expected, since one worker is allocated per instance and consequently a lot of work can be carried out in an efficient fashion.



Figure 1: Deploying one instance to an increasing number of physical nodes with one worker per instance.

The results from the second experiment, investigating the effect of workers with respect to the deployment time, can be seen in Figure 2. As can be seen, the deployment time is nearly cut in half as the number of workers are doubled. The reason for the time being a bit unevenly reduced as the number of workers are doubled is most likely due to network latency and other operations that are subject to variation in execution time. The reason for not deploying using more than twelve workers is that since twelve physical nodes are used, more than twelve workers would not make any sense since there would not be any work for the additional workers. These results are expected, since in theory, a 100% increase in workers should yield a 50% decrease in execution time since there are twice as many workers available to perform the deployments.



Figure 2: Deploying one instance to twelve nodes with an increasing number of workers.

Results from the third experiment, examining the change in deployment time when increasing the amount of instances on a constant amount of nodes, can be seen in Figure 3. It shows the increase of deployment time as the number of instances is scaled up to as much as 120 instances launched on fifteen physical nodes. Much similar to Figure 1, the deployment time is slowly growing but almost kept constant as the number of instances is multiplied by 8x. Due to one worker being used per instance launch, an almost constant deployment time is to be expected.



Figure 3: Scaling deployments virtually with one worker per instance to fifteen physical nodes.

5 Conclusion

A first implementation of *plcli* supporting Use Case 1 has been implemented and shown to provide a nearly constant deployment time for up to as much as 120 application instances on 15 physical nodes on the PlanetLab platform. The approach of using workers has been shown to be directly linked to the deployment time and shown to be very beneficial when performing concurrent deployments. Furthermore, *plcli* is bundled with many useful features for engineers using PlanetLab to run applications. When it comes to further work, *plcli* would benefit greatly from implementation of Use Case 2 and Use Case 3, as well as a more rigorous evaluation with respect to scalability (i.e. deploying to more physical nodes). This preliminary implementation is to be considered as a pilot and a foundation for continued work.
References

- [1] ns-3 a discrete-event network simulator for internet systems. URL: https://www.nsnam.org/, Accessed: 2019-07-12.
- [2] Teerawat Issariyakul and Ekram Hossain. *Introduction to Network Simulator 2 (NS2)*, pages 1–18. Springer US, Boston, MA, 2009.
- [3] Planetlabeurope. URL: https://www.planet-lab.eu/, Accessed: 2019-07-12.
- [4] Brent Chun, David Culler, Timothy Roscoe, Andy Bavier, Larry Peterson, Mike Wawrzoniak, and Mic Bowman. Planetlab: An overlay testbed for broad-coverage services. SIGCOMM Comput. Commun. Rev., 33(3):3–12, July 2003.
- [5] Larry Peterson and Timothy Roscoe. The design principles of planetlab. *SIGOPS Oper. Syst. Rev.*, 40(1):11–16, January 2006.
- [6] Larry Peterson, Andy Bavier, Marc E. Fiuczynski, and Steve Muir. Experiences building planetlab. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation*, OSDI '06, pages 351–366, Berkeley, CA, USA, 2006. USENIX Association.
- [7] User tools. URL: https://www.planet-lab.org/tools, Accessed: 2019-09-14.
- [8] parallel-ssh. URL: https://code.google.com/archive/p/parallel-ssh/, Accessed: 2019-07-03.
- [9] Jeffrey Joyce, Greg Lomow, Konrad Slind, and Brian Unger. Monitoring distributed systems. *ACM Trans. Comput. Syst.*, 5(2):121–150, March 1987.
- [10] Xuezheng Liu, Zhenyu Guo, Xi Wang, Feibo Chen, Xiaochen Lian, Jian Tang, Ming Wu, M. Kaashoek, and Zheng Zhang. D3s: Debugging deployed distributed systems. In *Proc. NSDI*, pages 423–437, 01 2008.
- [11] Tiago Jorge, Francisco Maia, Miguel Matos, José Pereira, and Rui Oliveira. Practical evaluation of large scale applications. In Alysson Bessani and Sara Bouchenak, editors, *Distributed Applications and Interoperable Systems*, pages 124–137, Cham, 2015. Springer International Publishing.
- [12] Eric M. Dashofy, André van der Hoek, and Richard N. Taylor. Towards architecturebased self-healing systems. In *Proceedings of the First Workshop on Self-healing Systems*, WOSS '02, pages 21–26, New York, NY, USA, 2002. ACM.
- [13] plcli github repository. URL: https://github.com/axelniklasson/plcli, Accessed: 2019-09-14.
- [14] The go programming language. URL: https://golang.org/, Accessed: 2019-07-09.
- [15] Rob Pike. The go programming language. Talk given at Google's Tech Talks, 2009.
- [16] Alan AA Donovan and Brian W Kernighan. *The Go programming language*. Addison-Wesley Professional, 2015.

- [17] N. Togashi and V. Klyuev. Concurrency in go and java: Performance analysis. In 2014 4th IEEE International Conference on Information Science and Technology, pages 213–216, April 2014.
- [18] plcli demo app github repository. URL: https://github.com/axelniklasson/plcli-demoapp, Accessed: 2019-09-14.

A Existing PlanetLab tools

Tables 1 and 2 list all the tools listed on the PlanetLab website¹ as of July 2, 2019.

Name	Brief description	State
Plush	Users describe experiments or computation in	Can't access
	XML, and Plush uses it to locate, contact,	webpage.
	and prepare resources. It includes a Neb-	
	ula GUI that allows users to build, visualize	
	and run their applications without using the	
	command-line interface.	
PIMan	PlanetLab Experiment Manager is designed	Can access
	to simplify the deployment, execution and	webpage,
	monitoring of your PlanetLab experiment.	but all links
	The application presents a simple GUI to per-	are broken.
	form common tasks.	
Stork	A software installation utility akin to yum	Broken link.
	and apt available for both users of PlanetLab	
	and for home use. It includes a Stock Slice	
	Manager GUI that simplifies package man-	
	agement and Stork installation on your Plan-	
	etLab slices.	
pShell	A Linux shell like interface providing a few	Broken link.
	basic commands to interact with a Planetlab	
	slice, works as a command center at the local	
	machine and interact with slice nodes.	
AppManager	PlanetLab Application Manager is designed	Broken link.
	to help deploy, monitor, and run applications	
	on PlanetLab. The package gives you the abil-	
	ity to centrally manage, install, upgrade, start,	
	stop, and monitor of applications on a Planet-	
	Lab slice.	

Table 1: Tools listed on the PlanetLab website and their state as of July 2, 2019 (Table 1/2)

¹https://www.planet-lab.org/tools

Name	Brief description	State
Emulab	A network testbed, giving researchers a wide	Accessible,
	range of environments in which to develop,	but different
	debug, and evaluate their systems.	purpose than
		plcli.
plDist	A tool for parallel distribution of files to Plan-	Broken link.
	etlab nodes using BitTorrent or rsync.	
Nixes	A set of bash scripts to install, maintain, con-	Broken link.
	trol and monitor applications on PlanetLab.	
PLDeploy	PlanetLab Slice Deploy Toolkit is a set of	Accessible,
	scripts to help users manage their slices.	comparison
		with <i>plcli</i>
		can be found
		in Section
		1.1.
pssh	Provides the parallel versions of the openssh	Accessible,
	tools. It can be used to control large collec-	comparison
	tions of nodes in the wide-area network.	with <i>plcli</i>
		can be found
		in Section
		1.1.
vxargs	Inspired by xargs and pssh, it provides the	Broken link.
	parallel versions of any arbitrary command,	
	including ssh, rsync, scp, wget, curl, etc.	
PlanetLab	A link emulator for PlanetLab that can be con-	Accessible,
broadband	figured with few important measured char-	but different
link emula-	acteristics of broadband links, such as their	purpose than
tor	asymmetric link bandwidths and queue sizes.	plcli.

 Table 2: Tools listed on the PlanetLab website and their state as of July 2, 2019 (Table 2/2)