

Programming Elastic Services with AEON and PLASMA

Patrick Eugster

Università della Svizzera italiana (USI)
PL4DS&DDM 2019, Dagstuhl

Work in collaboration with Bo Sang, Srivatsan Ravi, Hui Lu,
Gustavo Petri, Masoud Saeida Ardekani, Pierre-Louis Roman

Elastic Services

- Exploit *distributed* cloud (and edge) datacenter resources
at need
- Benefits service provider
- Benefits cloud provider



Landscape

- Autoscalers
 - Focus mostly on middle tier
 - Straightforward based on containers, VMs etc.
- AWS Lambda et al.
 - Stateless functions
- Specifically designed elastic services
 - E.g. map/reduce, graph processing, key/value storage
- More generic?



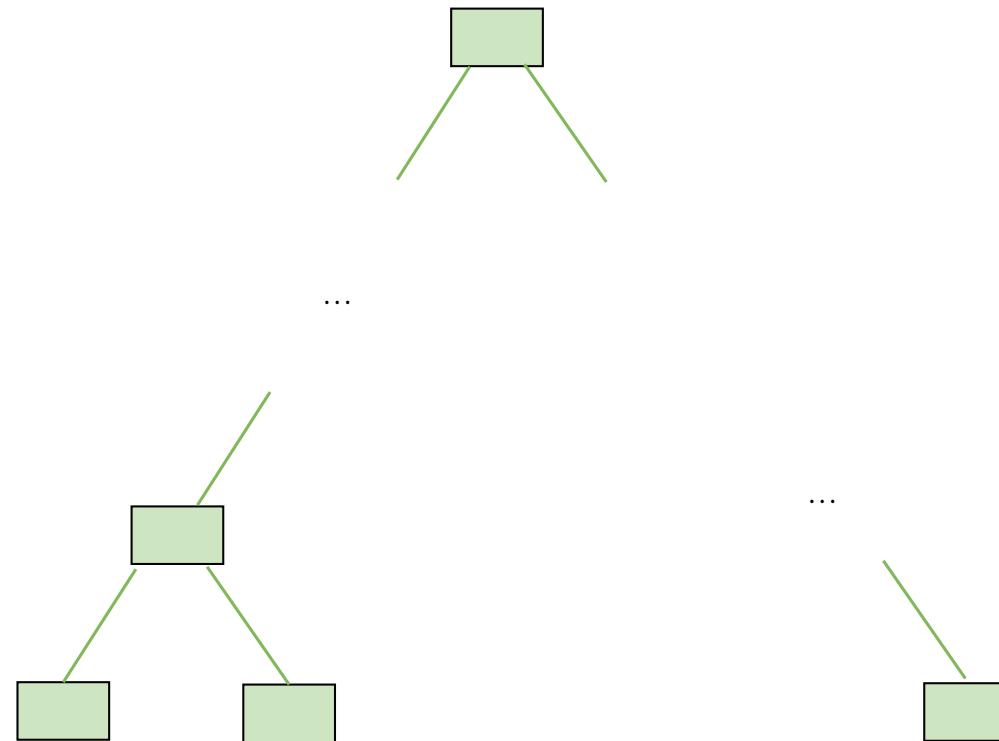
Roadmap

- Motivation
- AEON: Scalable serializable actor programming
- PLAS²MA: Elasticity programming

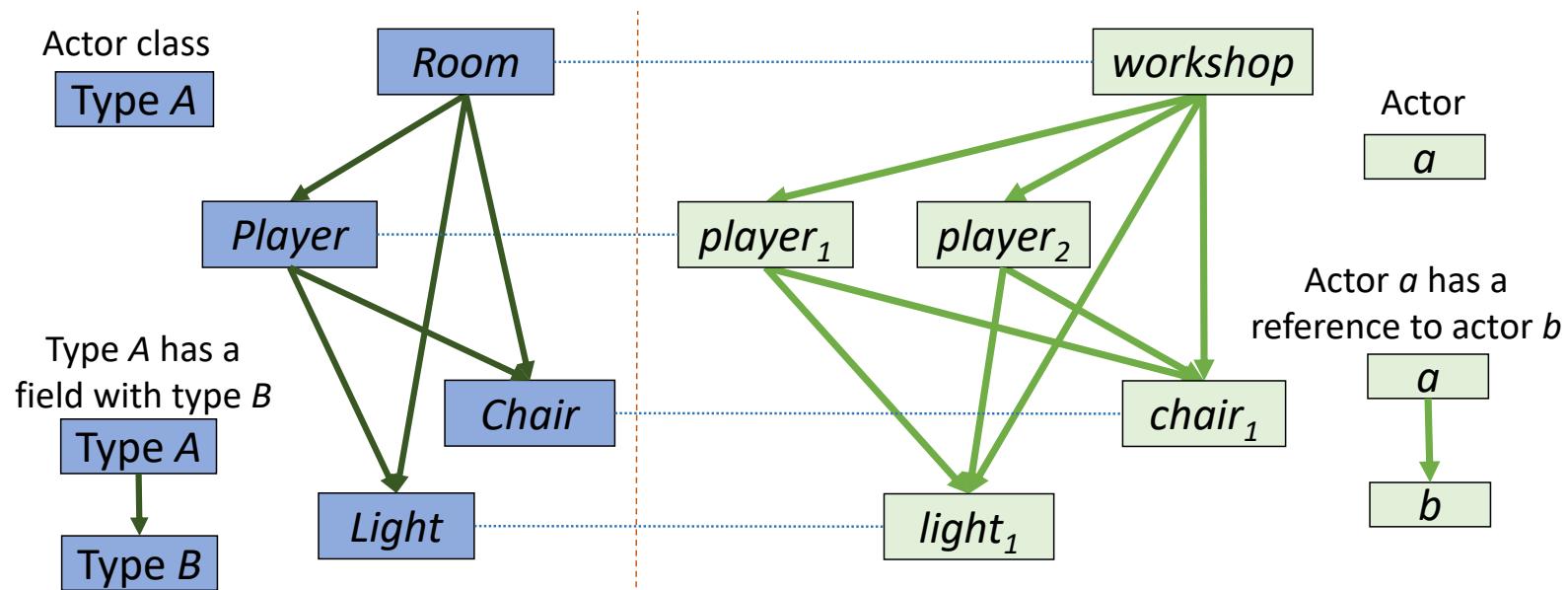
Actors and Events Ownership Network (AEON)

- Scalability is prerequisite for elasticity
- Actors to the rescue?
 - Scalable, familiar
 - But... little help for reasoning beyond single messages
 - Isolation
- Scalability vs consistency in distributed systems

Tree Intuition



Actor DAG



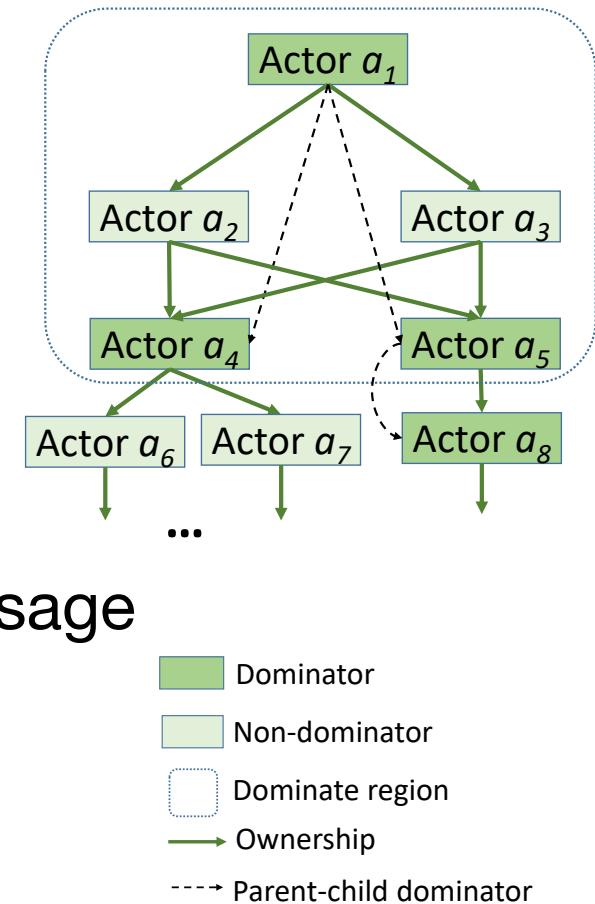
Restricted actor (interaction) graph

DAG Enforcement

- Ownership-style typing
 - Multiple owners
- Type-based program analysis
 - Program variables
- Leveraged for serializability and deadlock freedom

Leveraging the DAG

- Top-down locking
 - Cf. “hand-over-hand” locking
- Dominator
 - “LUB” for all actors involved in a message
 - Dominate region, queues
- Can use also for checkpointing



Programming Model

- Messages tagged at callsite
 - Events (no return) — external vs internal
 - Synchronous methods — default
 - Asynchronous methods (no return)
- Practical extensions
 - Calls outside DAG via yield fields — (sub)events
 - Readonly
 - Direct recursion

Mandatory Code Snippet

```
1  actorclass Room {  
2      // players in this room  
3      vector<Player> players;  
4      vector<Chair> chairs;  
5      vector<Light> lights;  
6      ...  
7      void turnOnLights() {  
8          // turn on lights in parallel  
9          for(int i=0 upto lights.size())  
10             async lights[i].turnOn();  
11     }  
12     void chairColorUpdate(...) {...}  
13 }  
14  
15 actorclass Light {  
16     bool isOn;  
17     ...  
18     void turnOn() { isOn = true; }  
19     void turnOff() { isOn = false; }  
20 }  
21  
22 actorclass Chair {  
23     ...  
24     void paint(...) {...}  
25     void putBag(...) {...}  
26 }  
27 class Role {...}  
28  
29 actorclass Player {  
30     Light light;  
31     Chair chair;  
32     Role role;  
33     Room room;  
34     vector<Chair> chairs;  
35     ...  
36     void cleanLight() {  
37         chair.putBag(...);  
38         light.turnOff();  
39     }  
40     void paintChair() {  
41         light.turnOn();  
42         chair.paint(...);  
43         event room.chairColorUpdate(...);  
44     }  
45     void paintChairs(...) {  
46         for(int i=0 upto chairs.size())  
47             // 3 options  
48             async chairs[i].paint(...);  
49             // chairs[i].paint(...);  
50             // event chairs[i].paint(...);  
51     }  
52 }
```

Capsule Semantics

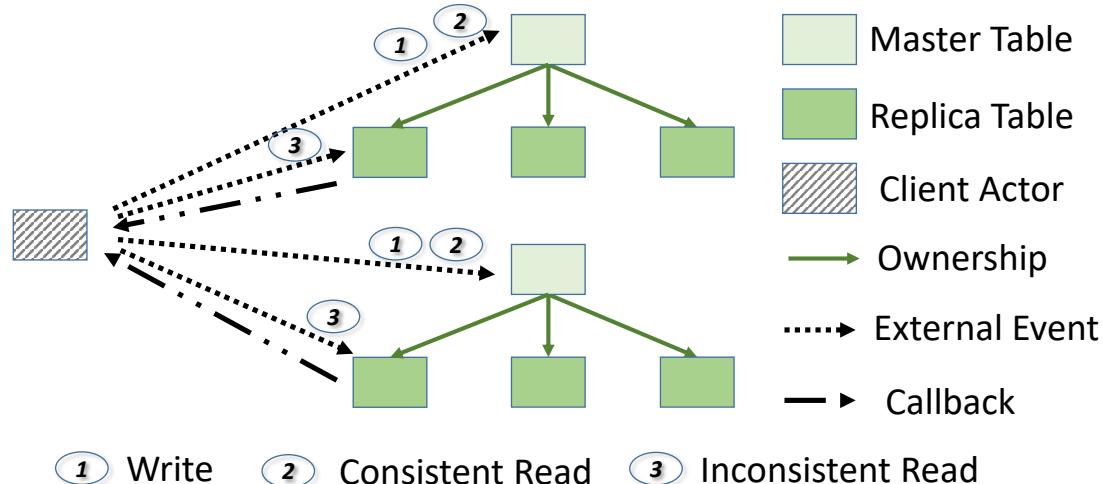
Caller	Call ($dc^?$)	Field ($\text{yield}^?$)	Method ($\text{ro}^?$)	Result	Blocking	Execution
External	event	yield	Any	Callback	Non-blocking	Serialized
Internal non-ro	event	Any	Any	Callback	Non-blocking	Serialized
Internal ro	event	Any	ro	Callback	Non-blocking	Serialized
Internal non-ro	Sync	Non-yield	Any	Return (T)	Blocking	Atomic
Internal non-ro	async	Non-yield	Any	Callback	Non-blocking	Atomic
Internal ro	Sync	Non-yield	ro	Return (T)	Blocking	Atomic
Internal ro	async	Non-yield	ro	Callback	Non-blocking	Atomic

Experiences

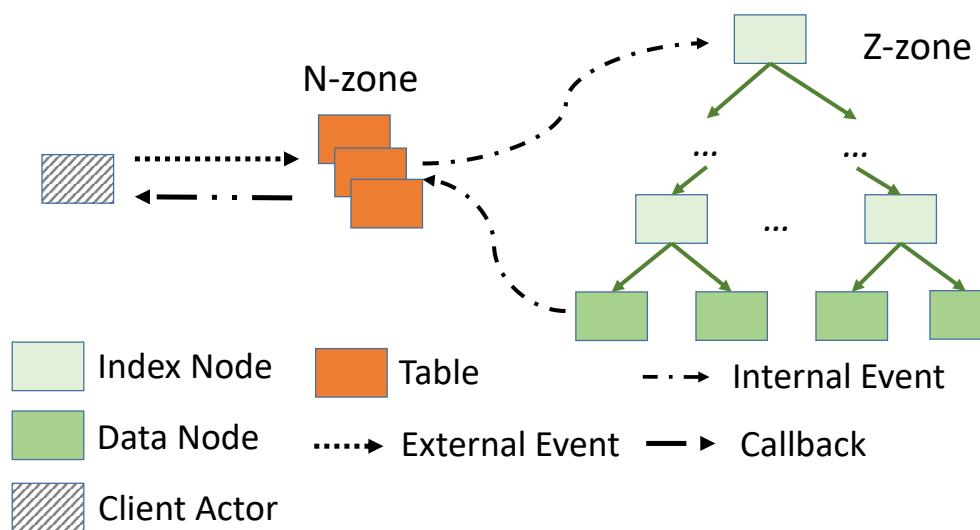
Application	Type	Ownership constraint	Multiple clusters	Uses		Uses int. events	Actor classes	LoC
				ro	async			
Halo Galactic	Game	None	Yes	No	No	Yes	3	313
Halo Presence	Game	None	Yes	No	No	Yes	3	287
Multi-player game	Game	DAG	No	No	No	No	4	564
Chirper	Social net	None	No	No	No	Yes	2	118
Graph engine	Computation	None	No	No	No	Yes	2	469
MapReduce	Computation	None	No	No	No	Yes	3	390
TPC-W-mini	Benchmark	None	No	No	No	Yes	2	236
TPC-C	Benchmark	DAG	No	Yes	Yes	No	7	1798
zExpander	Caching	Tree	Yes	Yes	Yes	Yes	3	506
Memcached	Caching	None	No	Yes	No	No	1	220
BigTable	Store	Tree	No	Yes	Yes	No	3	216
Cassandra	Store	Tree	No	Yes	Yes	No	2	221
Metadata store	Store	Tree	No	Yes	Yes	No	2	552
Silo	Store	DAG	No	Yes	Yes	No	7	1899
Bank account	Finance	Tree	No	No	No	No	4	176
Skip list	Data structure	DAG	No	Yes	Yes	No	2	322
LSM	Data structure	Tree	Yes	Yes	Yes	Yes	2	1587
B+ tree	Data structure	Tree	No	Yes	Yes	No	2	1457
B tree	Data structure	Tree	No	Yes	Yes	No	2	1437
Piazza	Course management	DAG	No	Yes	Yes	No	5	259

Examples

- Cassandra

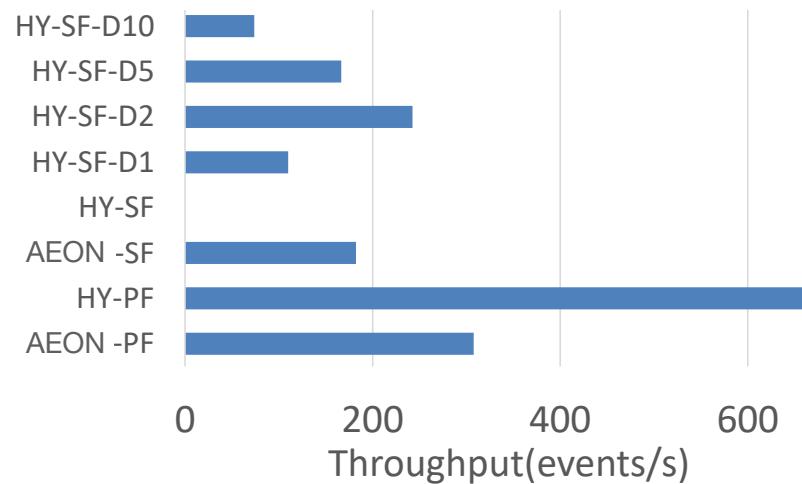


- zExpander

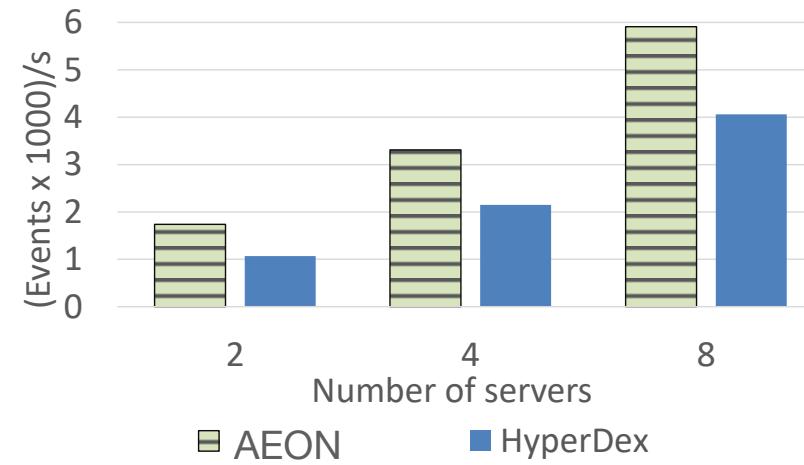


Performance

- AEON vs HyperDex Warp
 - Metadata server for Warp transactional file system



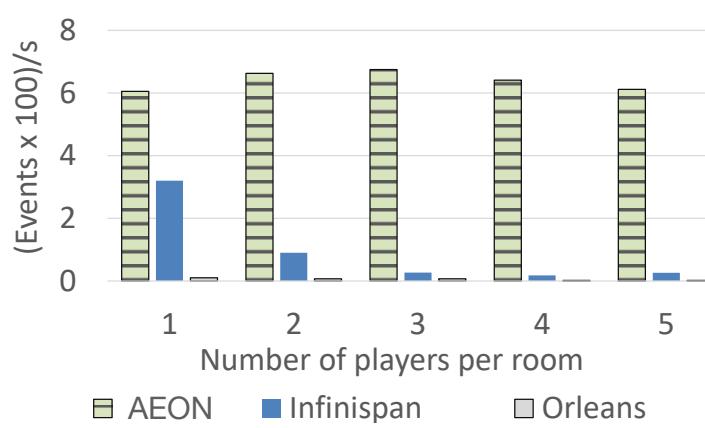
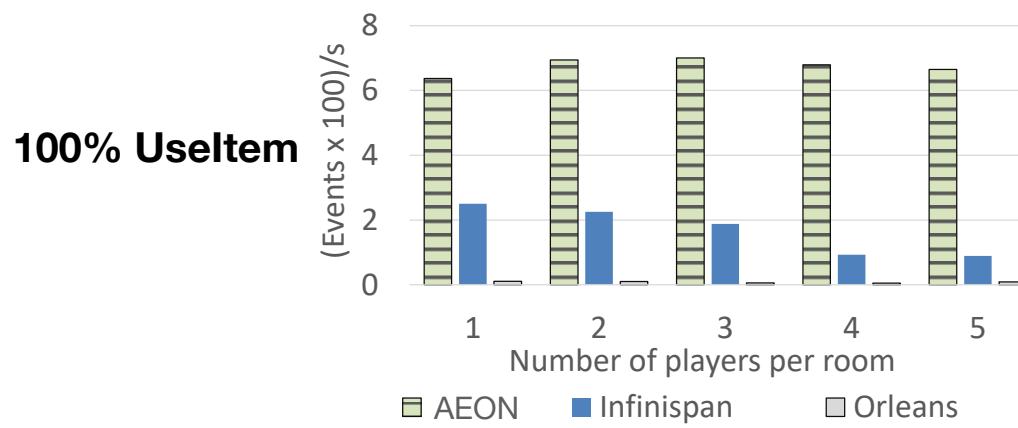
File creation



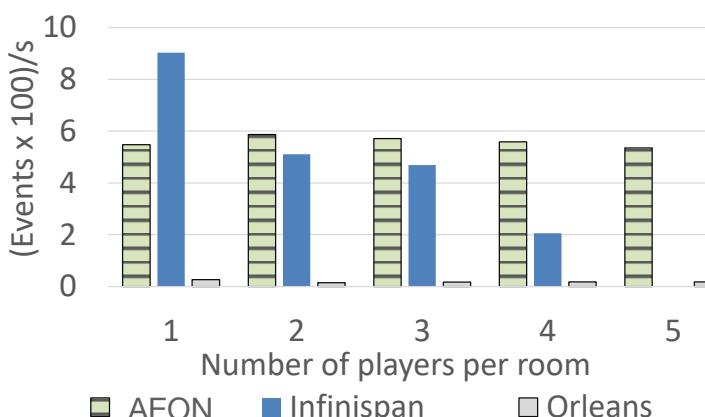
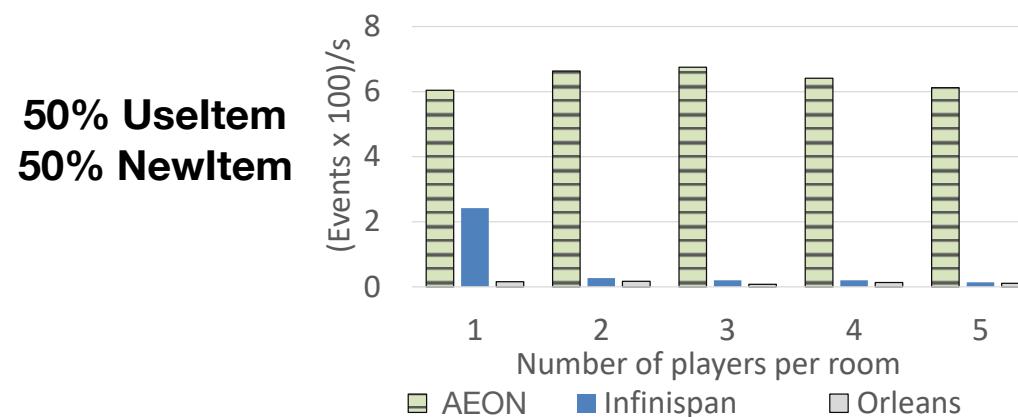
Open private inode

Performance (2)

- Game application



**80% UseItem
20% NewItem**



100% NewItem

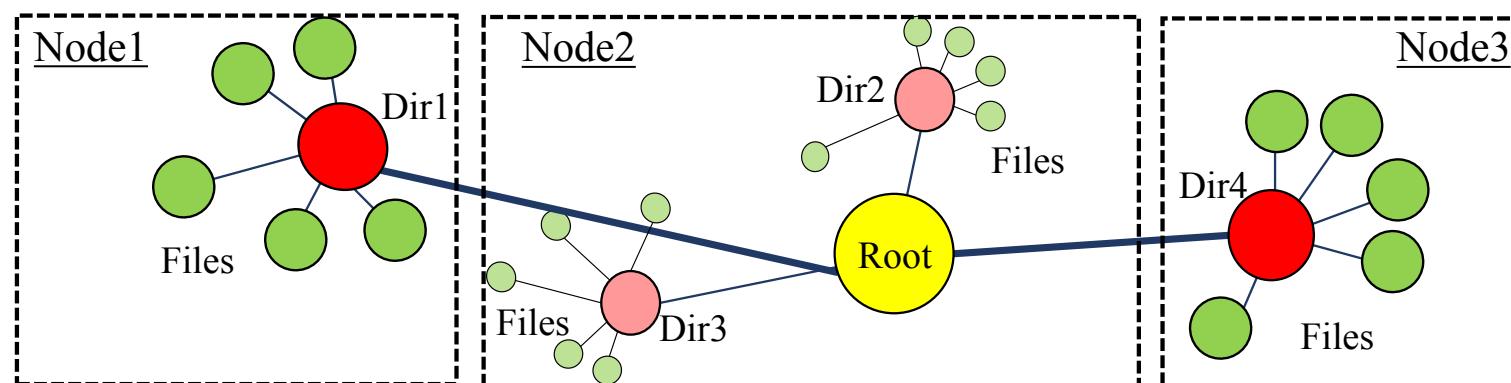
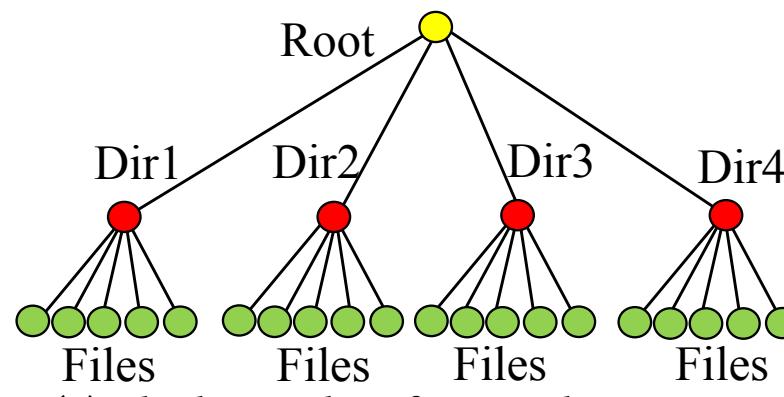
Programmable Elasticity for Stateful Serverless Computing Applications (PLAS²MA)

Scalable distributed programming model

Migration

- How to know when/what to migrate?
 - Automated approaches for VMs etc.

Metadata Server Revisited



2-Level Programming

1. Application programming (e.g., AEON)

2. Elasticity programming

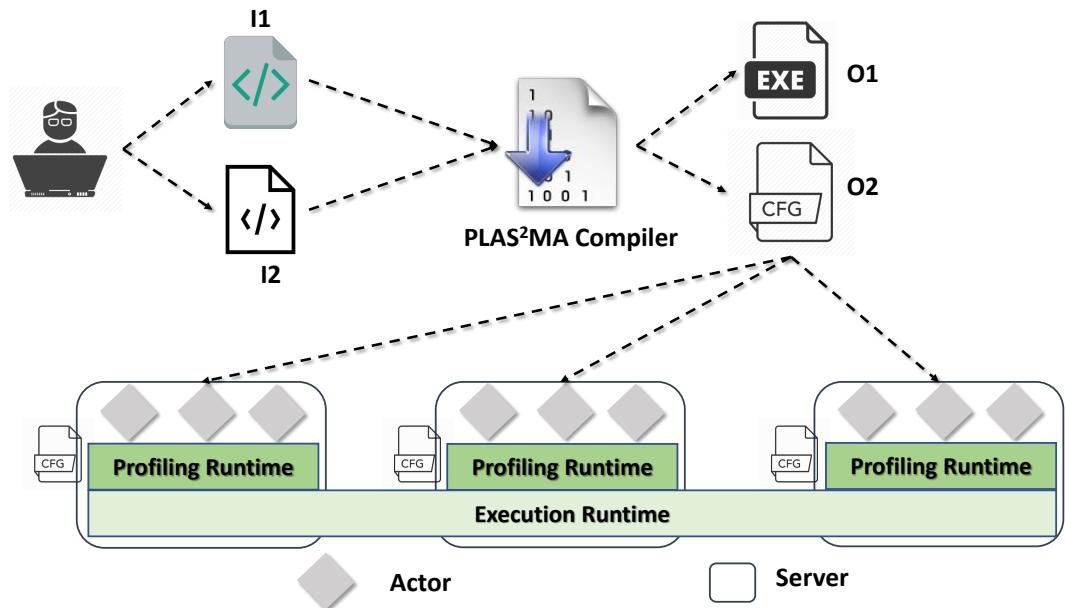
- Assumes actor types

- Elasticity rules

[R-R] Resource elasticity rules

- Actor (e.g., CPU usage) and server resources (e.g., network usage)

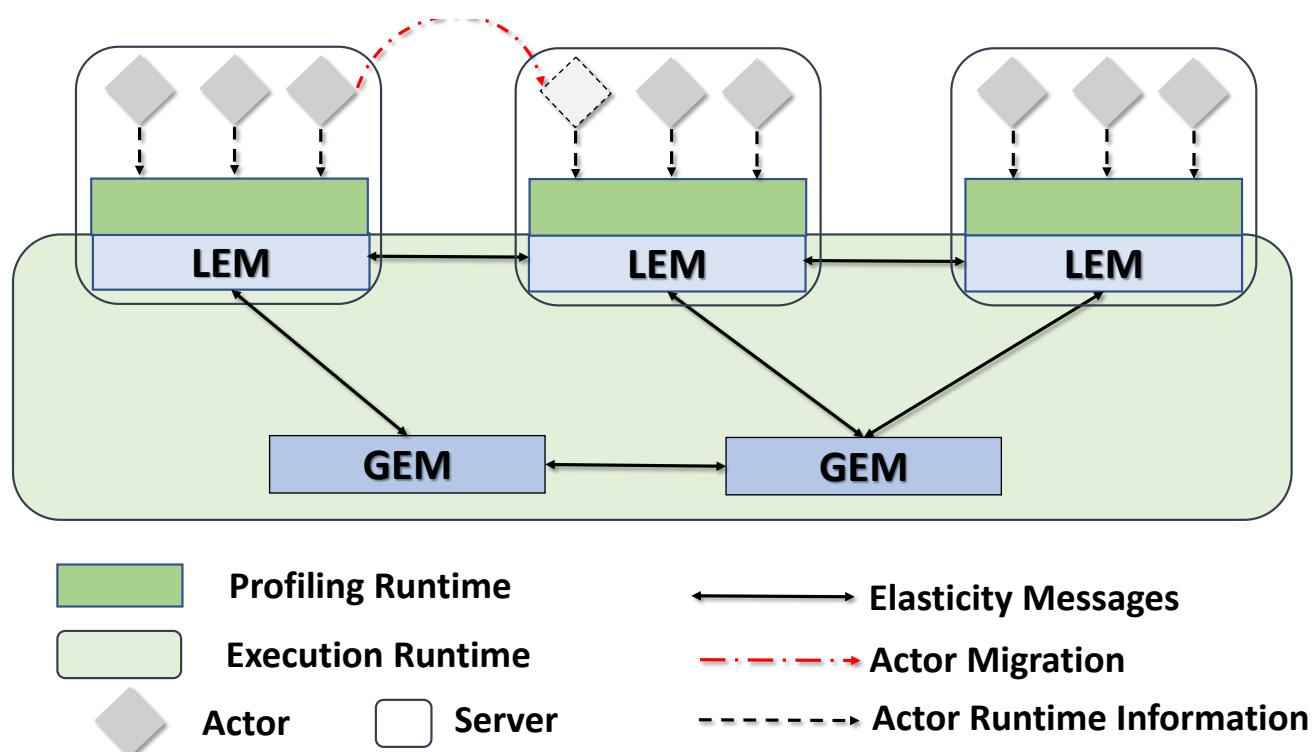
[R-I] (Actor) interaction elasticity rules (e.g., message rate)



Elasticity Programming

<i>Policy</i>	$pol ::= \overline{rul}$	
<i>Rule</i>	$rul ::= cond \Rightarrow \overline{beh};$	
Actor	$actor ::= atype(var) \mid atype \mid var$	
<i>Actor type</i>	$atype ::= fname \mid \textbf{any}$	
Condition	$cond ::= cond \textbf{ or } cond \mid cond \textbf{ and } cond$ true $feat.stat \: comp \: val$ $actor \: \textbf{in} \: \textcolor{red}{ref}(actor.pname)$	[F-IA]
<i>Feature</i>	$feat ::= entity.res$ $cllr.\textcolor{red}{call}(actor.fname)$	[F-IA]
<i>Entity</i>	$entity ::= actor$ server	[F-RA] [F-RS]
<i>Caller</i>	$cllr ::= \textbf{client} \mid actor$	
<i>Statistic</i>	$stat ::= \textcolor{orange}{count} \mid \textcolor{orange}{size} \mid \textcolor{orange}{perc}$	
<i>Resource</i>	$res ::= \textcolor{blue}{cpu} \mid \textcolor{blue}{mem} \mid \textcolor{blue}{net}$	
<i>Comparison</i>	$comp ::= < \mid > \mid >= \mid <=$	
Behavior	$beh ::= \textcolor{red}{balance}(\{\overline{atype}\}, res)$ reserve ($actor, res$) colocate ($actor, actor$) separate ($actor, actor$) pin ($actor$)	[R-R] [R-R] [R-I] [R-I] [R-I]
<i>Value</i>	$val \in \mathbb{N} \cup \mathbb{R}$	

Runtime System

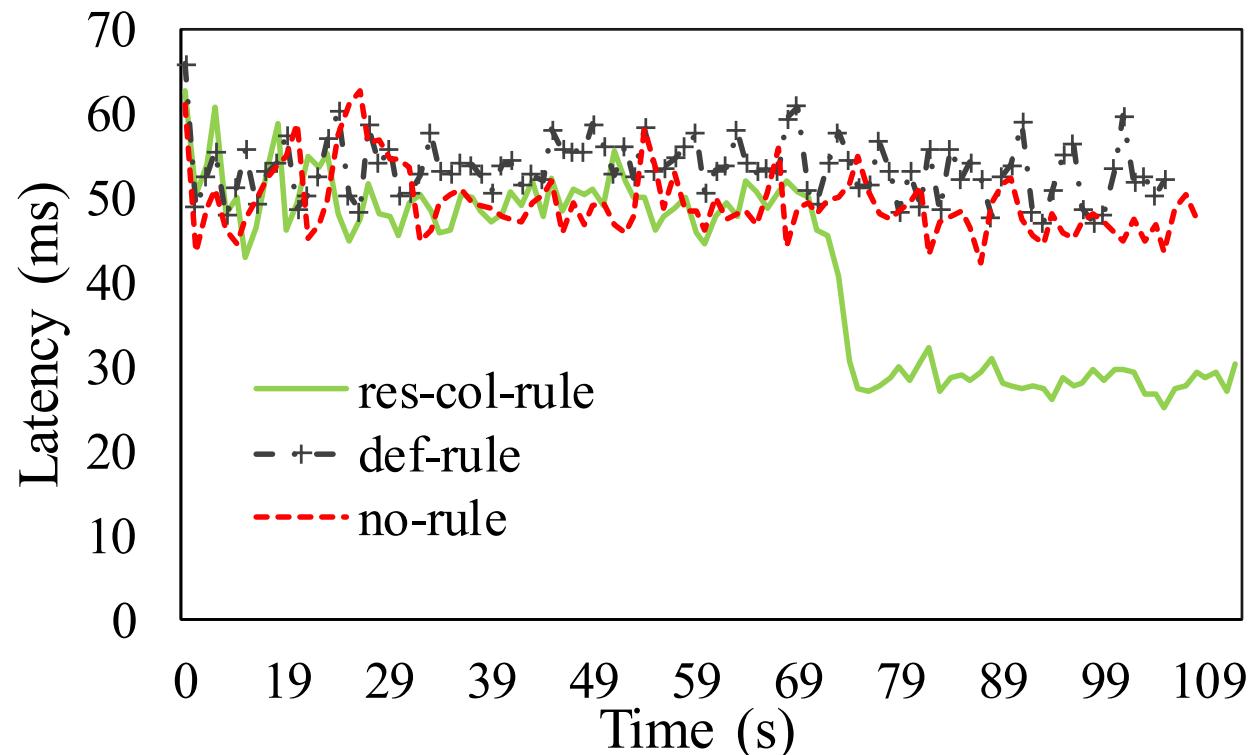


Experiences

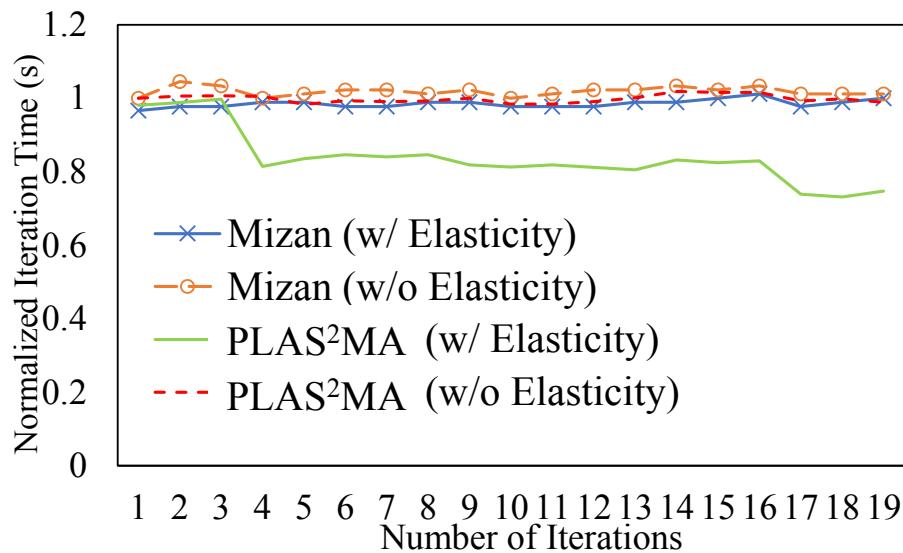
Application	LoC	Elasticity rules (and number of rules)
Metadata Server	253	1. Colocate Folder with Files in it (since they are accessed together)
PageRank	465	1. Balance CPU workload
Halo	287	1. Balance CPU workload of Routers
Presence Service		2. Colocate Session with Players in it
Media Service	756	1. Balance network workload for FrontEnds 2. Provide VideoStream with enough CPU 3. Colocate linked VideoStream and UserInfo 4. Avoid migrating MovieReview 5. Balance CPU workload of ReviewChecker
B+ tree	1457	1. Colocate parent-child inner nodes 2. Put leaf nodes on separate servers
Piccolo	-	1. Balance CPU workload for Workers 2. Colocate Worker and Table that Worker reads data from
E-Store	645	1. Put hot Keys on idle servers 2. Colocate parent-child Keys
zExpander	-	1. Put leaf nodes on idle servers
Cassandra	-	1. Put table replicas on different servers

Metadata Server

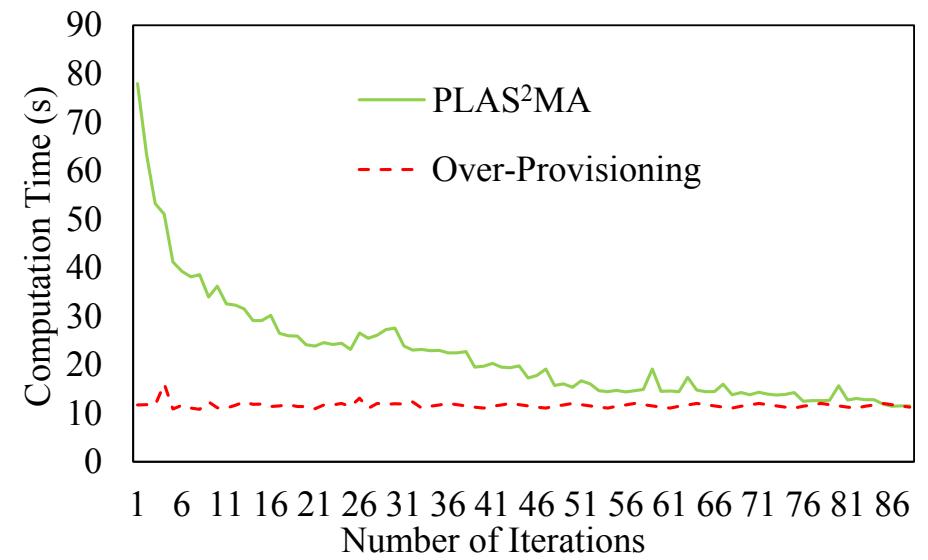
```
server.cpu.perc > 80 and
client.call(Folder(fo).open).perc > 40 and
File(fi) in ref(fo.files) =>
    reserve(fo, cpu); colocate(fo, fi);
```



PageRank



vs Mizan



vs over-provisioning

Conclusions

- AEON: combines scalability and serializability in a well-known programming model with pragmatic tradeoffs
- PLAS²MA: fine-grained elasticity with few programmer instructions
- Ongoing
 - Extensions, e.g., futures
 - Livelocks with ownership updates
 - High availability