

PL4DS

Patrick Eugster (USI)

Shonan Village

May 2019

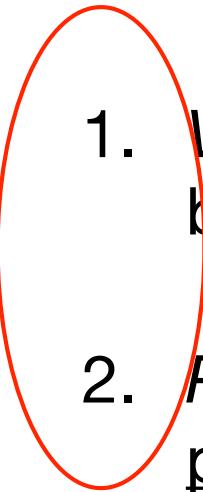
I Had a Dream...

- ... of a PL for all end-to-end DS programming
- But DSs are very “fragmented” — layers, tiers, ...
- Nowadays just develop PL techniques for DSs

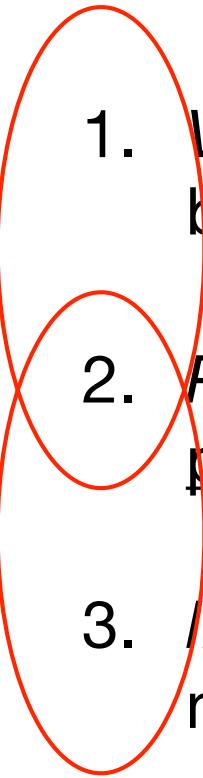
Ongoing Work

1. *Verification* of DSs, e.g., behavioral typing for DSs, SAT-based smart contracts
2. *Programming models* for DSs, e.g., elasticity programming, switch buffer management
3. *Infrastructure* for DSs, e.g., synchronous communication, network flow shape prediction

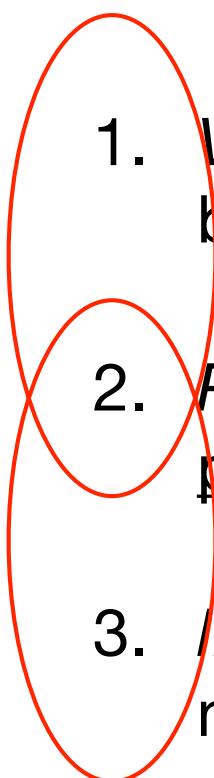
Ongoing Work

- 
1. *Verification* of DSs, e.g., behavioral typing for DSs, SAT-based smart contracts
 2. *Programming models* for DSs, e.g., elasticity programming, switch buffer management
 3. *Infrastructure* for DSs, e.g., synchronous communication, network flow shape prediction

Ongoing Work

- 
1. *Verification* of DSs, e.g., behavioral typing for DSs, SAT-based smart contracts
 2. *Programming models* for DSs, e.g., elasticity programming, switch buffer management
 3. *Infrastructure* for DSs, e.g., synchronous communication, network flow shape prediction

Ongoing Work

- 
1. *Verification* of DSs, e.g., behavioral typing for DSs, SAT-based smart contracts
 2. *Programming models* for DSs, e.g., elasticity programming, switch buffer management
 3. *Infrastructure* for DSs, e.g., synchronous communication, network flow shape prediction

Co-design

1. Behavioral Typing

- Multi-party session types (MPSTs)
 - Targeting *middleware* for networked DSs, e.g., Spark
- Select features
 - E.g., no delegation, shared channels
 - Protocol types

Story in 3 Chapters

- I. App-level failures (“exceptions”) [FORTE’17]
- II. Partial (crash) failures, consensus service [ESOP’18]
- III. Failures *integrated* with events, role sets, subsessions

Story in 3 Chapters

- I. App-level failures (“exceptions”) [FORTE’17]
- II. Partial (crash) failures, consensus service [ESOP’18]
- III. Failures *integrated* with events, role sets, subsessions

G_{top} ::= $\{C_i(\overline{r_R}, \underline{r_R}, R) = G \text{ with } r_R @ r_R.G\}_I$

G ::= spawn $C(\overline{r_R}, \underline{R}, \overline{R}).G$ | $r \rightarrow \mathbb{r}\{l_i : G_i\}_{i \in I}$ | $\mu t.G$ | t | end

R ::= W | M | ... $\mathbb{r} ::= r_R$ | R

Story in 3 Chapters

- I. App-level failures (“exceptions”) [FORTE’17]
- II. Partial (crash) failures, consensus service [ESOP’18]
- III. Failures *integrated* with events, role sets, subsessions

Story in 3 Chapters

- I. App-level failures (“exceptions”) [FORTE’17]
- II. Partial (crash) failures, consensus service [ESOP’18]
- III. Failures *integrated* with events, role sets, subsessions

$P ::= (\bar{H}, c : R)$	Event loop
loop	cEnd
$x[y]!l.P$	Send
$x[y]?l.P$	Receive
$x.\text{spawn } C(\bar{y})$	Spawn

...

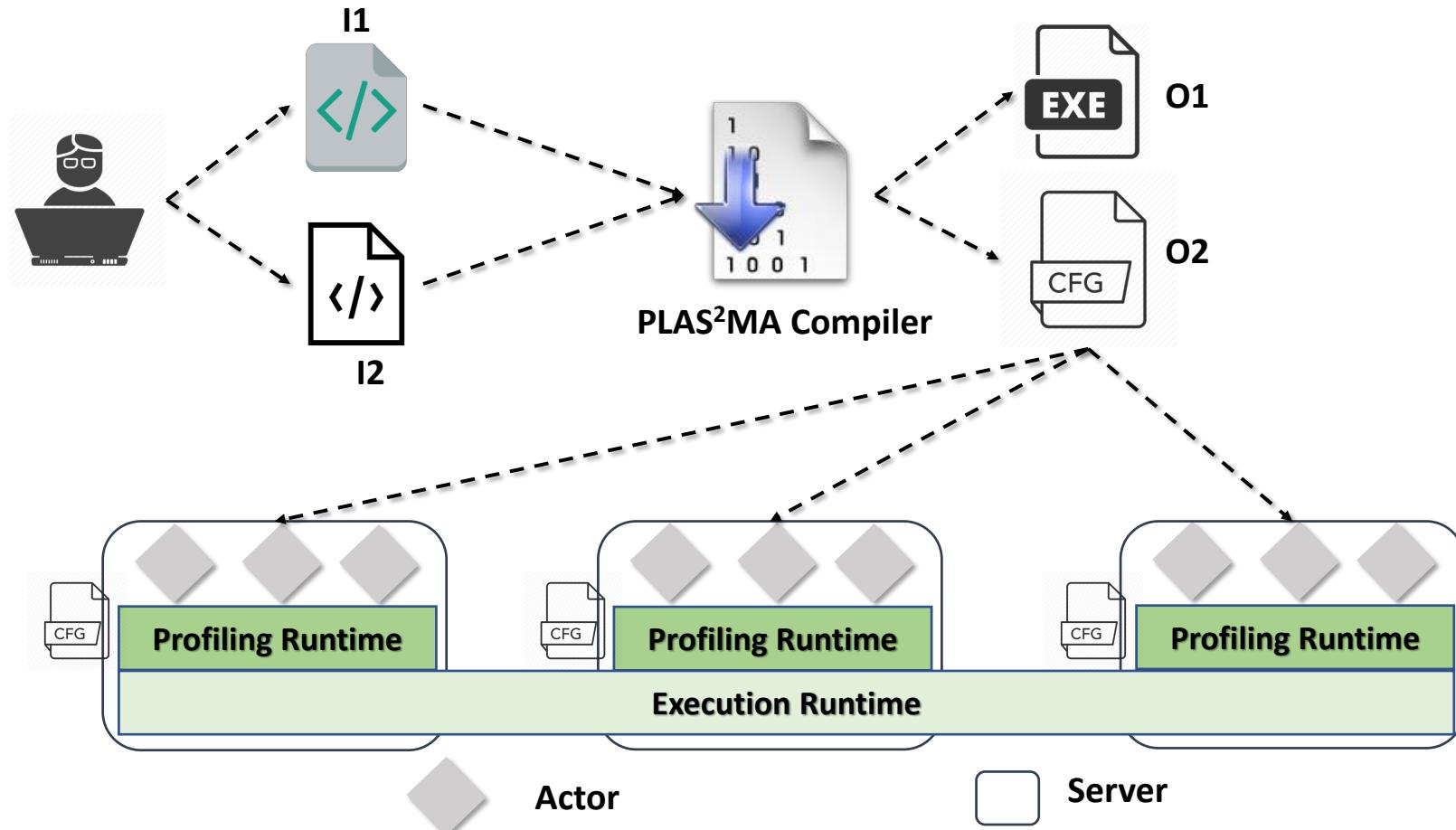
2. Elasticity Programming

- How to develop elastic software systems?
 - Middle tier
 - Scalability prerequisite
 - AEON (Atomic Events via Ownership Network)
[Middleware'16]
 - Fine-grained elasticity, e.g., with actors?

Programmable Elasticity for Stateful Serverless Computing Apps (PLAS²MA)

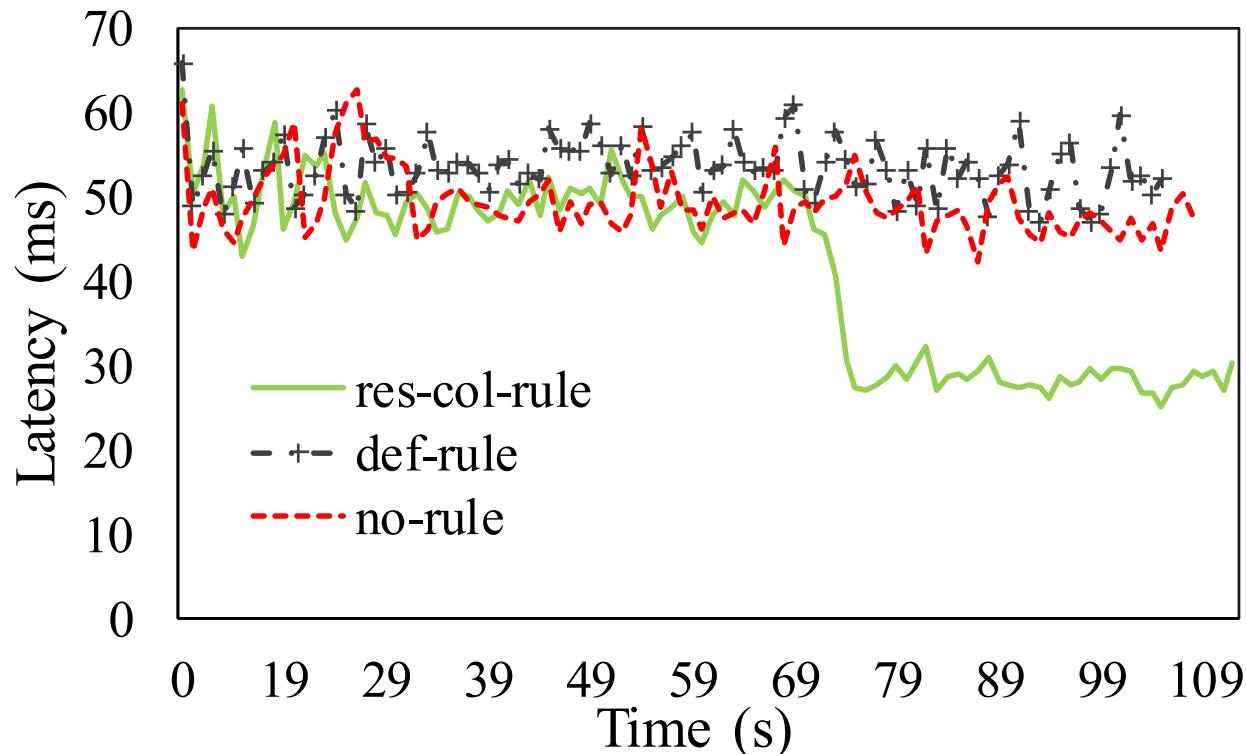
<i>Policy</i>	$pol ::= \overline{rul}$	
<i>Rule</i>	$rul ::= cond \Rightarrow \overline{beh};$	
<i>Actor</i>	$actor ::= atype(var) \mid atype \mid var$	
<i>Actor type</i>	$atype ::= fname \mid \text{any}$	
<i>Condition</i>	$cond ::= cond \text{ or } cond \mid cond \text{ and } cond$ $\text{true} \mid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

PLAS²MA Toolchain



Meta-Data Server

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

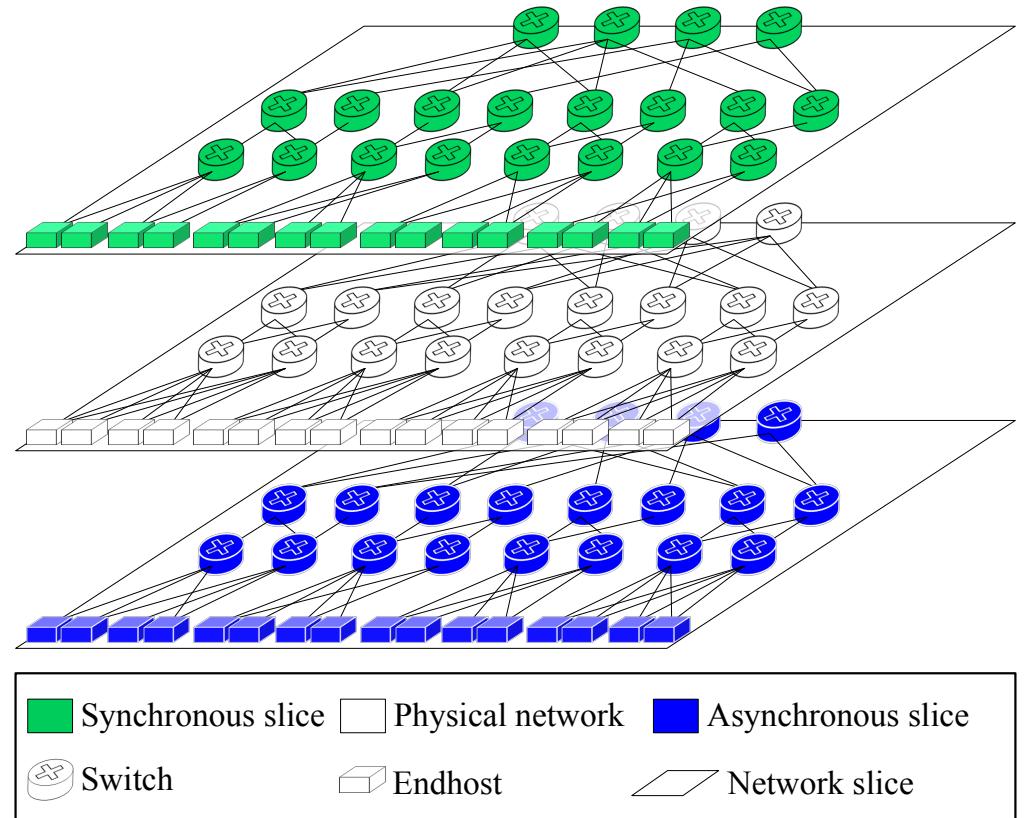


3. Synchronous Communication

- Distributed systems suffer from asynchrony
 - Failure or delay?
- Proposals for overall latency & jitter reduction in datacenter
 - Still “high” latency and/or jitter
 - Not needed for all traffic
- Several proposals for “INP” support for distributed protocols
 - Which one(s)? Ramifications?

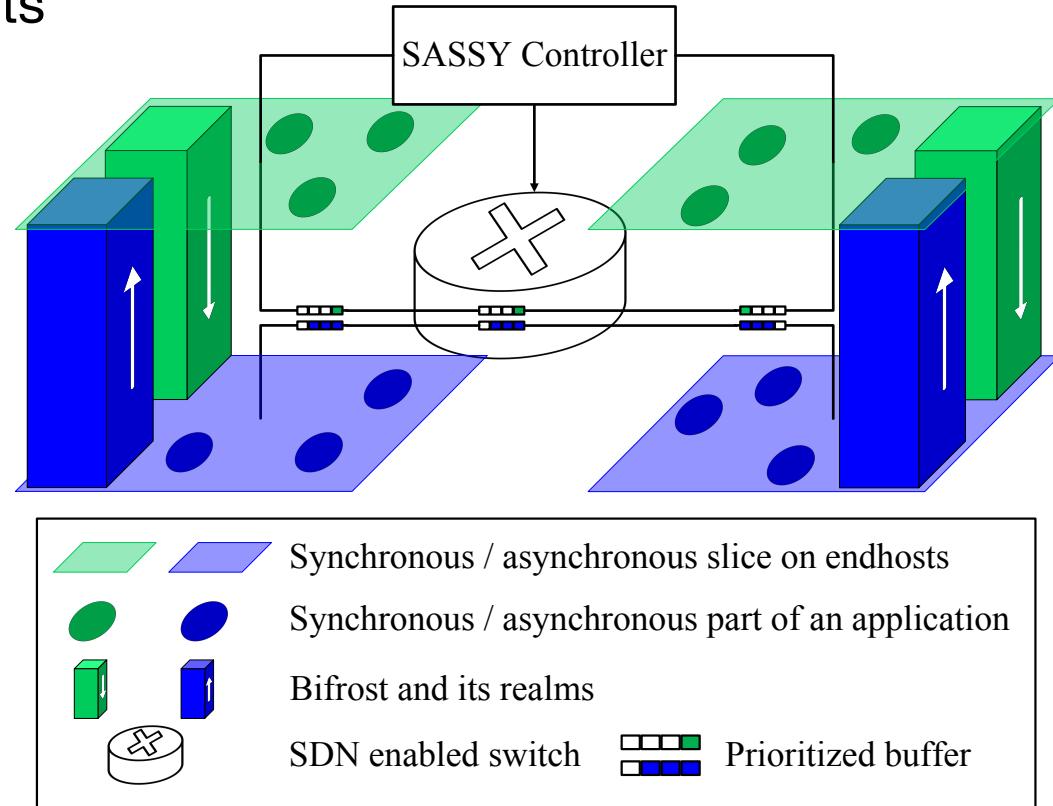
Sliced Network for Asynchronous/Synchronous Communication (SASSY)

- 2 slices
 - 1. Latency & deadline-centric “control” traffic
 - No penalty for not using
 - 2. “Regular” bandwidth-centric traffic
- Commodity HW and SW (Linux)
 - Extensions for SmartNIC



Co-Existing Slices

- API for 1. slice traffic requests
- Switches
 - 2 priorities
 - Cf. buffering
- Links
 - Traffic engineering
- Endhosts
 - Reserve core, mask interrupts, etc.



Synchronous Slice Performance Overview

- SAP datacenter, 17 Intel XEON 52 cores, hyperthreading, 128 GB
- Arista 7280CR-48 switches

Approach	Latency (μs)	Jitter (ns)	FFT
SASSY	5,137	144,052	8,1
SASSY_{sNIC}	4,089	31,995	4,652
SASSY_0	5,253	162,478	10,581
QJump	7726,405	7721311,38	4275821,503
DPDK	3831,707	395873,272	33128,168

(Via 1 switch)

Open Questions

- What abstractions are appropriate for what level/layer?
 - How to deal with “layering”?
 - Do we need more “OSI layers”?
 - Breaking up the “application layer”?