

Predictable, System-Level Fault Tolerance in a Component-based RTOS

(Dissertation Defense)

Presented by **Jiguo Song**

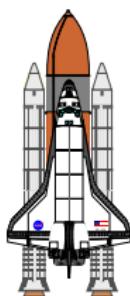
May 16, 2016

Computer Science Department
The George Washington University

Committee: Prof. Gabriel Palmer (Advisor)
Prof. Timothy Wood

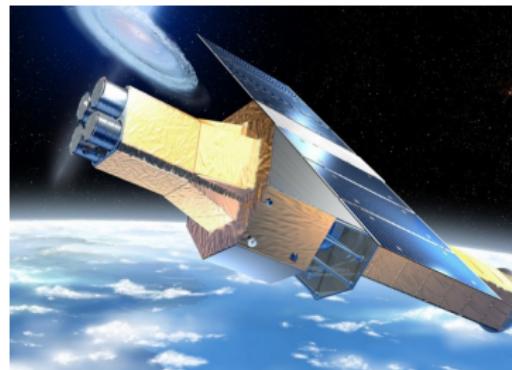
Prof. Rahul Simha
Prof. Howie Huang

Real-Time and Embedded Systems



Consequences of Embedded System Faults

Japanese X-ray astronomy satellite Hitomi lost in 2016



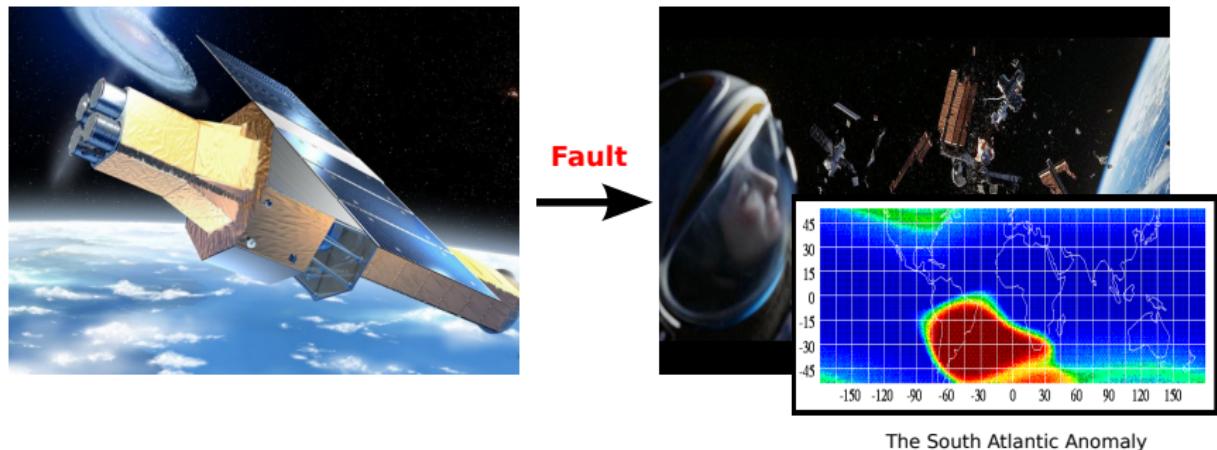
Fault
→



- Financial losses of \$286 million USD
- “*It's a scientific tragedy*” - Richard Mushotzky, UMD

Consequences of Embedded System Faults

Japanese X-ray astronomy satellite Hitomi lost in 2016



- Financial losses of \$286 million USD
- “*It's a scientific tragedy*” - Richard Mushotzky, UMD

Consequences of embedded system faults

Toyota Sudden Unintended Acceleration (SUA) in 2004 – 2010



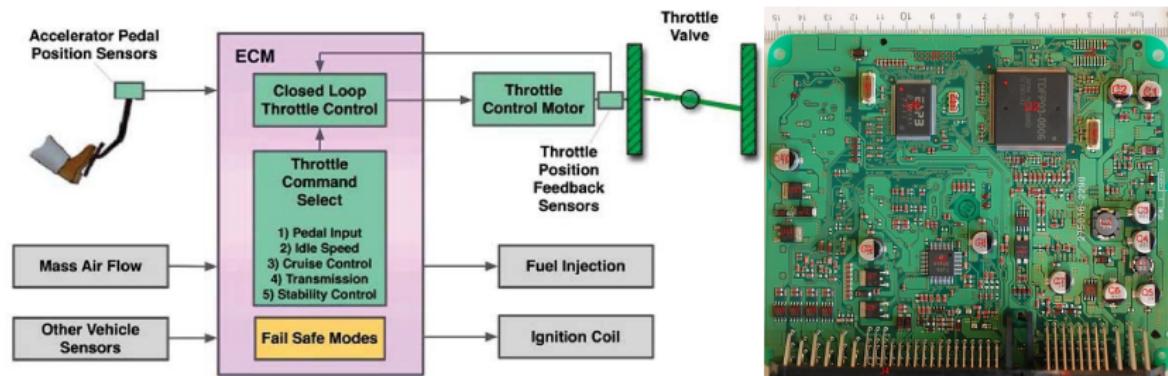
Fault →



- 89 deaths as of May 2010 and nearly 400 U.S. lawsuits
- Recall 10+ million vehicles and pay \$1.2 Billion USD fine

<http://www.cbsnews.com/news/toyota-unintended-acceleration-has-killed-89/>

Sudden Unintended Acceleration

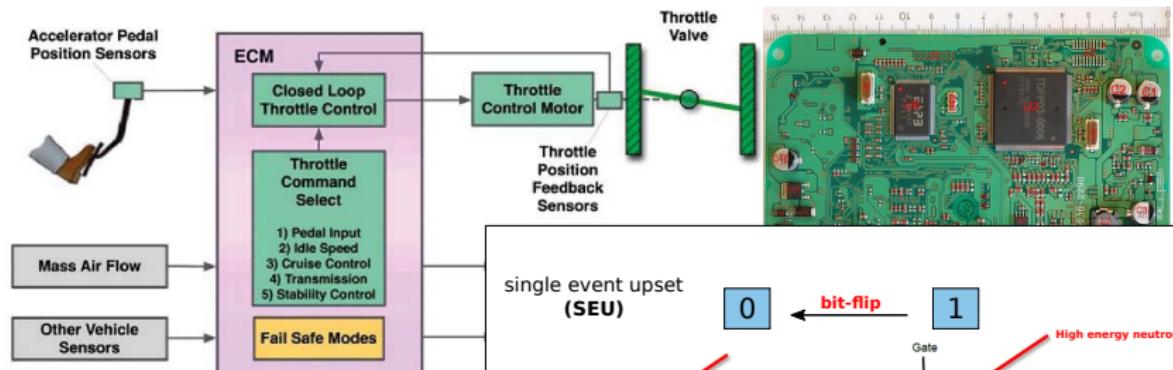


Uncontrolled acceleration in Toyota Camrys

- Electronic Throttle Control System (ETCS)
- OSEK OS, 24 tasks, 280K LOC of C
- bit flip in scheduler data-structures
→ reproducible 30-sec uncontrolled acceleration

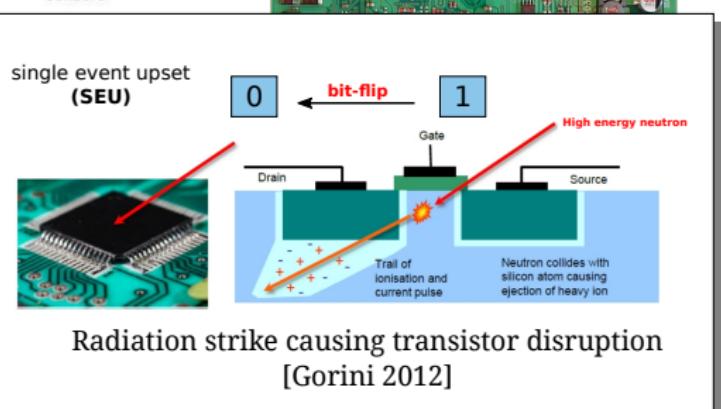
“a bit-flip there, will have the effect of killing one of the tasks”
– Bookout v Toyota

Sudden Unintended Acceleration



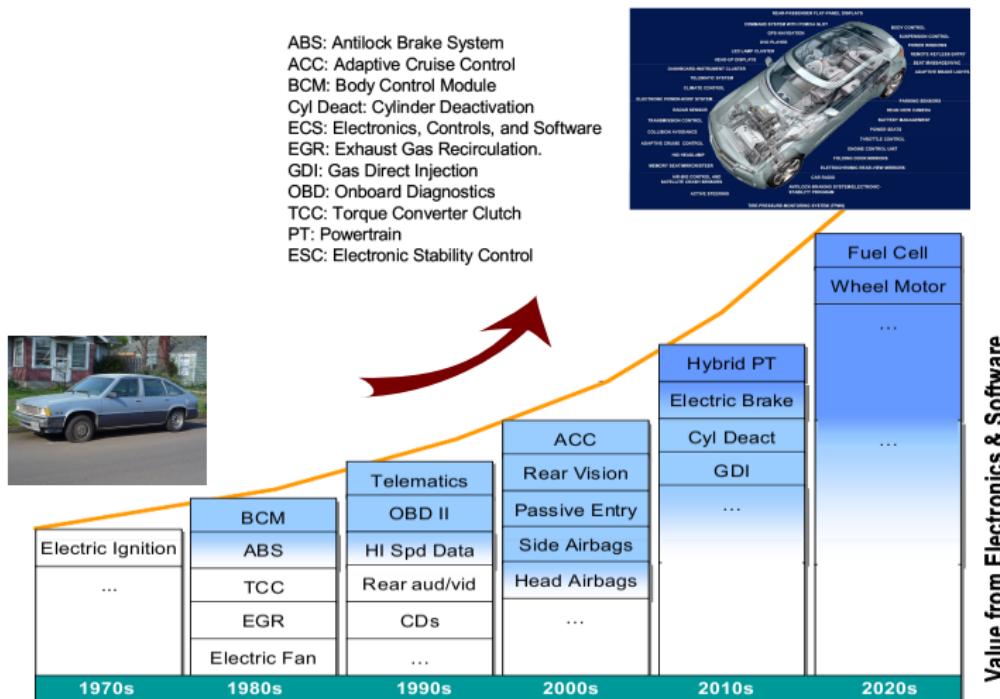
Uncontrolled acceleration in

- Electronic Throttle Con
- OSEK OS, 24 tasks, 28
- bit flip in scheduler data-structures
→ reproducible 30-sec uncontrolled acceleration



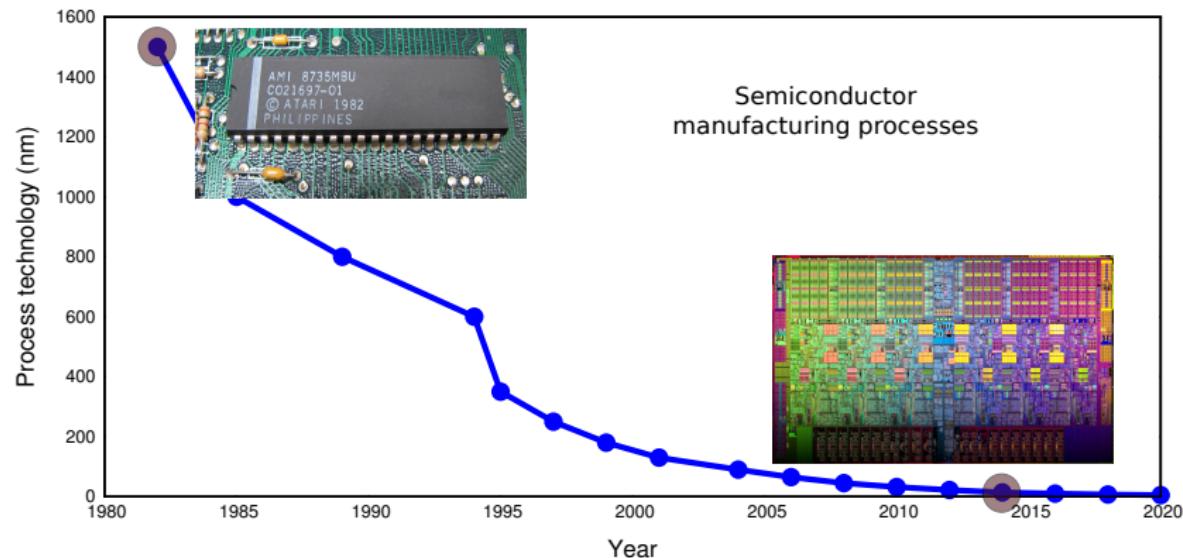
“a bit-flip there, will have the effect of killing one of the tasks”
– Bookout v Toyota

Embedded faults: Bad Now, Worse Tomorrow



- + more functionality
- more complexity → dependability more challenging

Embedded faults: Bad Now, Worse Tomorrow

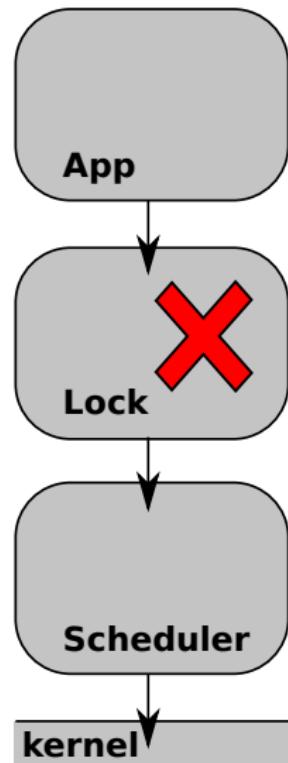


Decreasing process sizes → 5nm

- + faster
- + less power
- + smaller
- increased vulnerability to HW transient faults

Toward Dependable Real-Time Systems

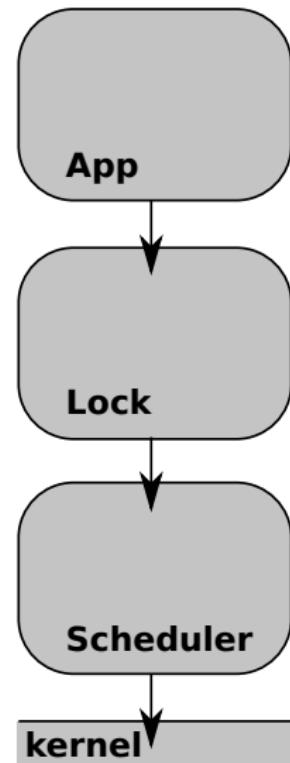
- Keys in dependable system design
 - isolation
 - detection
 - recovery
 - predictability



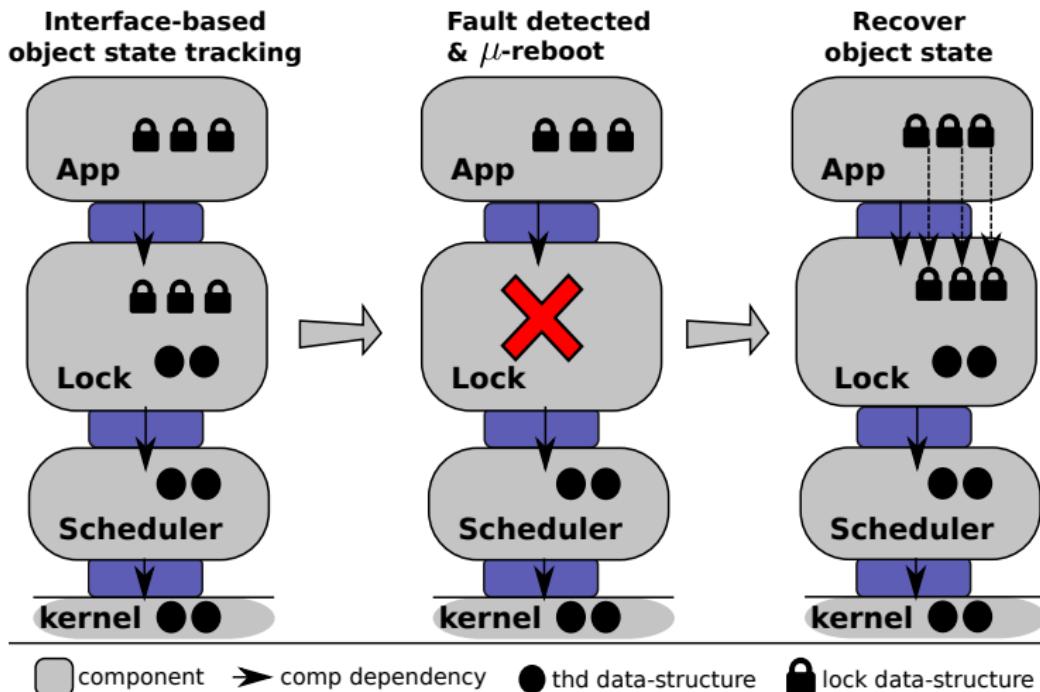
- Keys in dependable system design
 - isolation
 - detection
 - recovery
 - predictability

Goal

- ▶ predictably detect and recover
- ▶ system-level fault



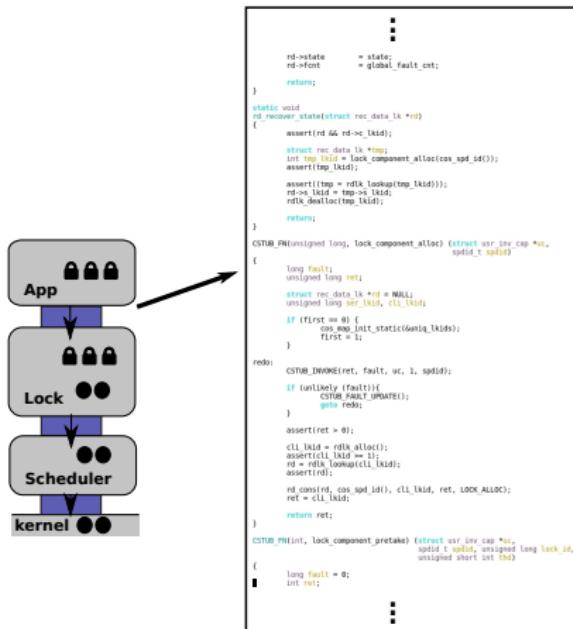
C³ Review – System-Level Fault Recovery



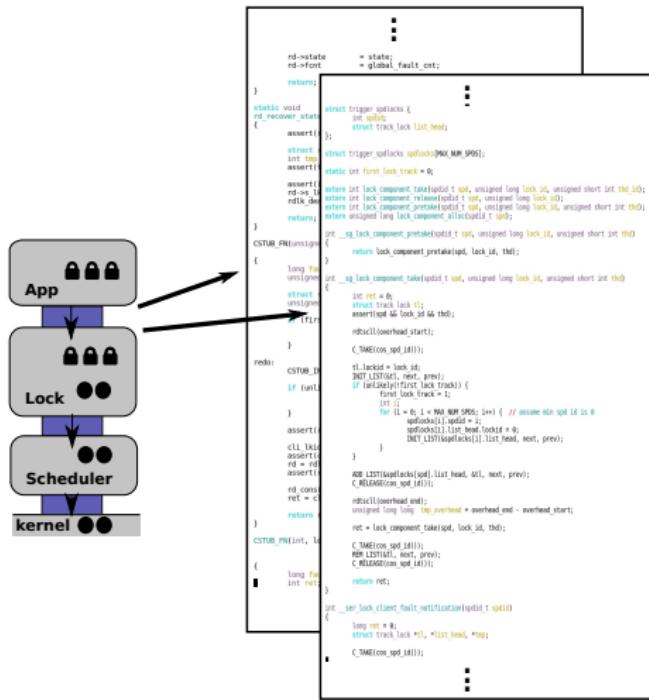
C³ – an interface-driven, predictable system-level fault recovery mechanism

J. Song, J. Wittrock, and G. Palmer, "Predictable, efficient system-level fault tolerance in C³" in RTSS, 2013

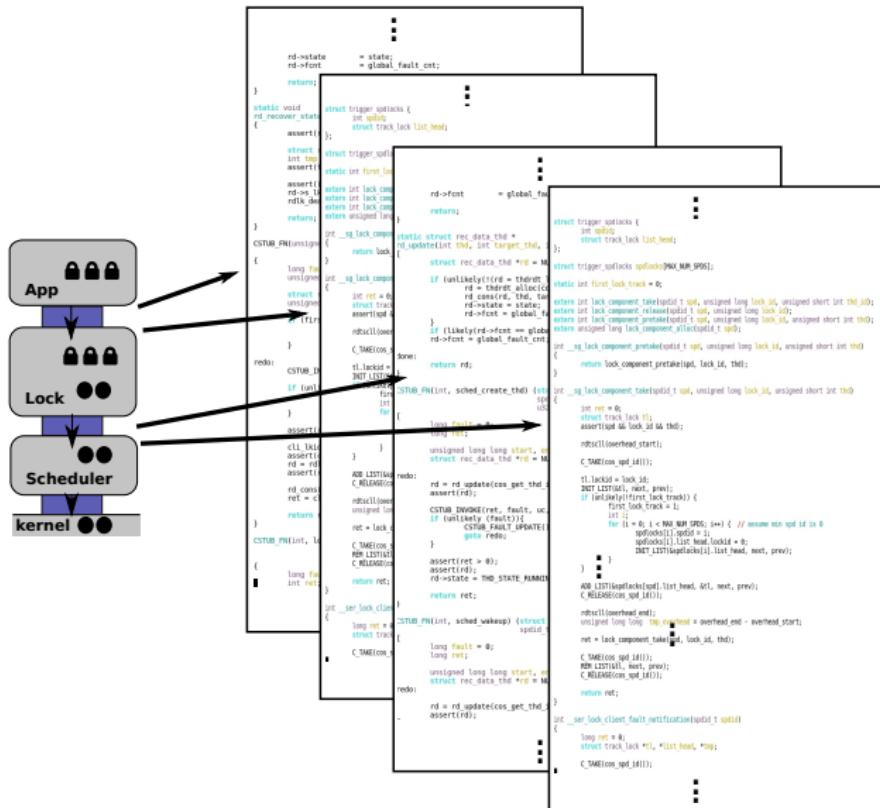
C³ Implementation – Writing Code Manually



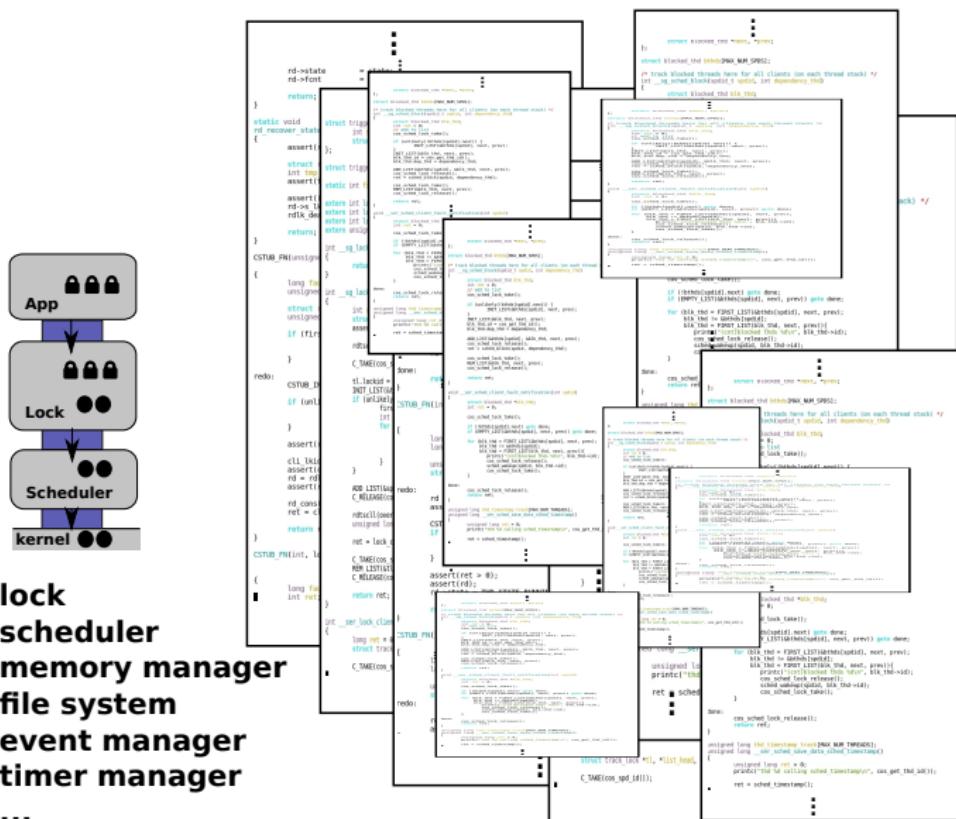
C³ Implementation – Writing Code Manually



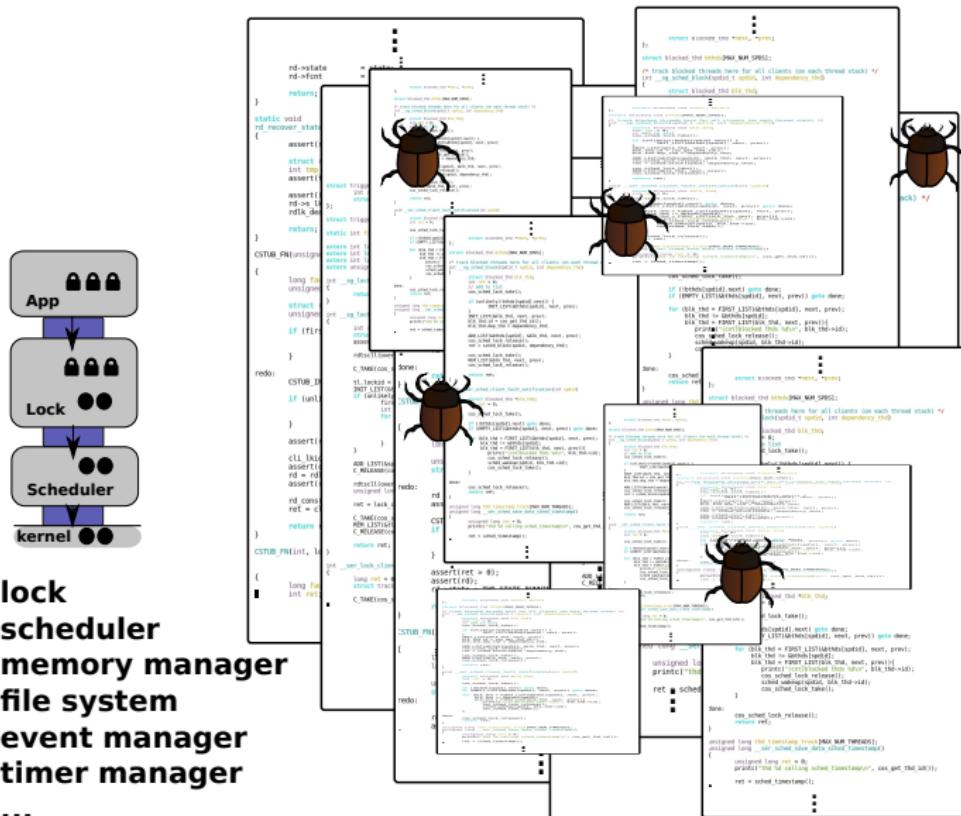
C³ Implementation – Writing Code Manually



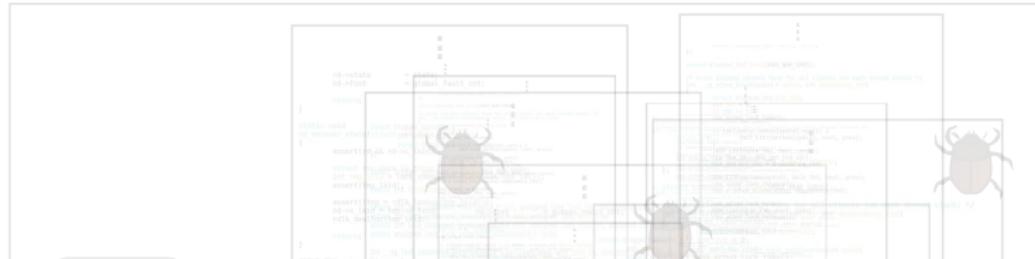
C³ Implementation – Writing Code Manually



C³ Implementation – Writing Code Manually



C³ Implementation – Writing Code Manually



Manually writing is ad-hoc and error-prone

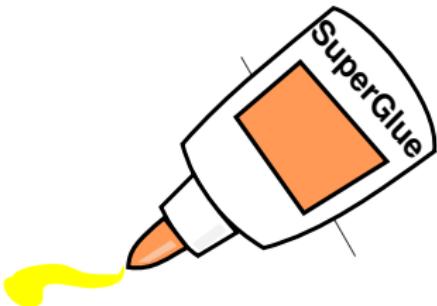
Goal → automatically generate C³-style recovery code

lock
scheduler
memory manager
file system
event manager
timer manager
...



SuperGlue

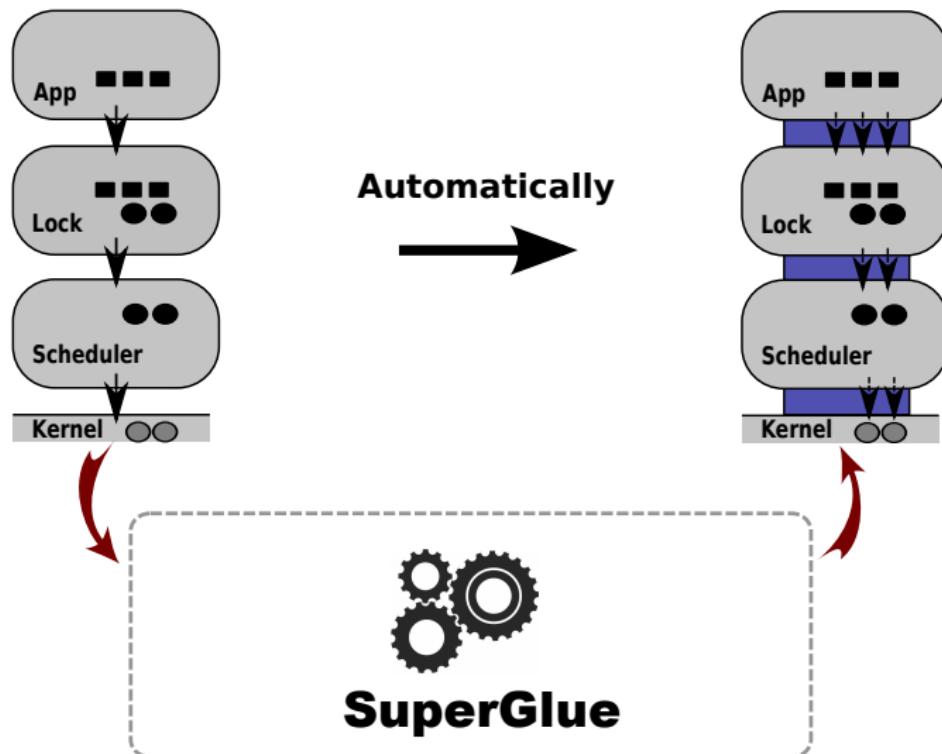
- automatically generate interface-driven recovery code
- for system-level services



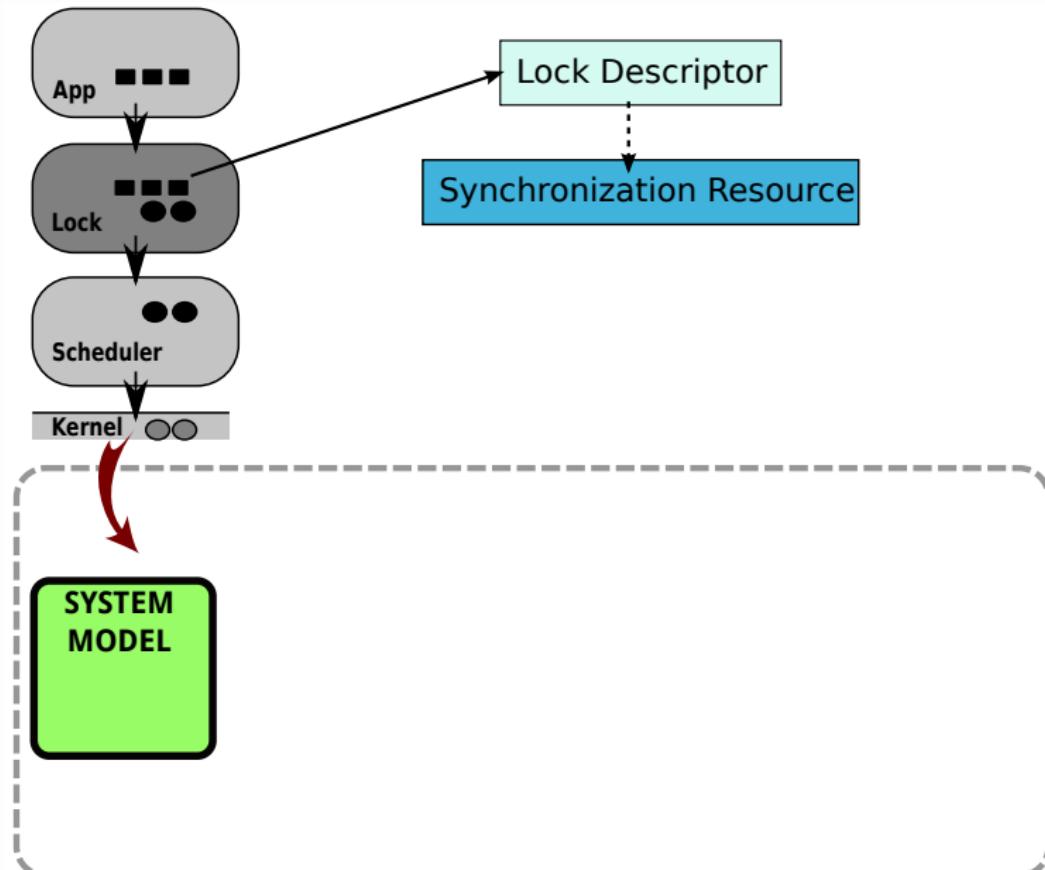
Main ideas

- A system **model** and **interface description language (IDL)**
- An **IDL compiler** synthesizes recovery code from the model

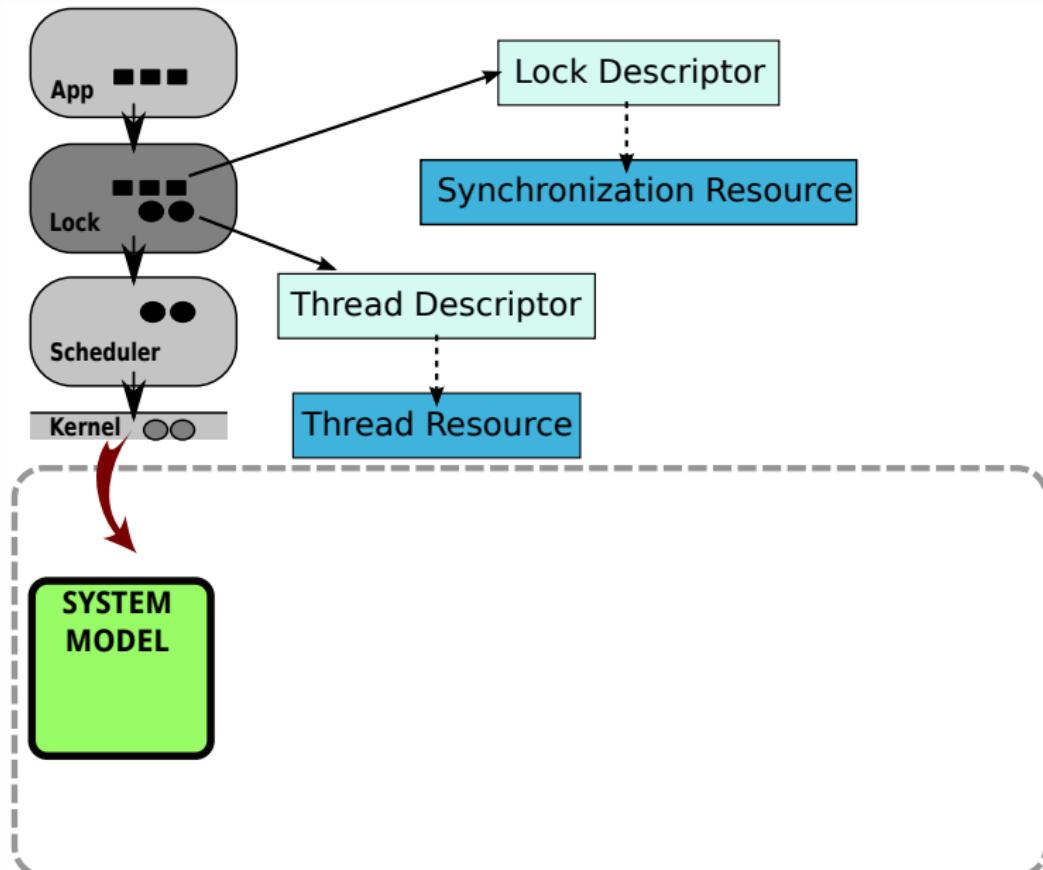
SuperGlue – IDL-based System-Level Fault Tolerance



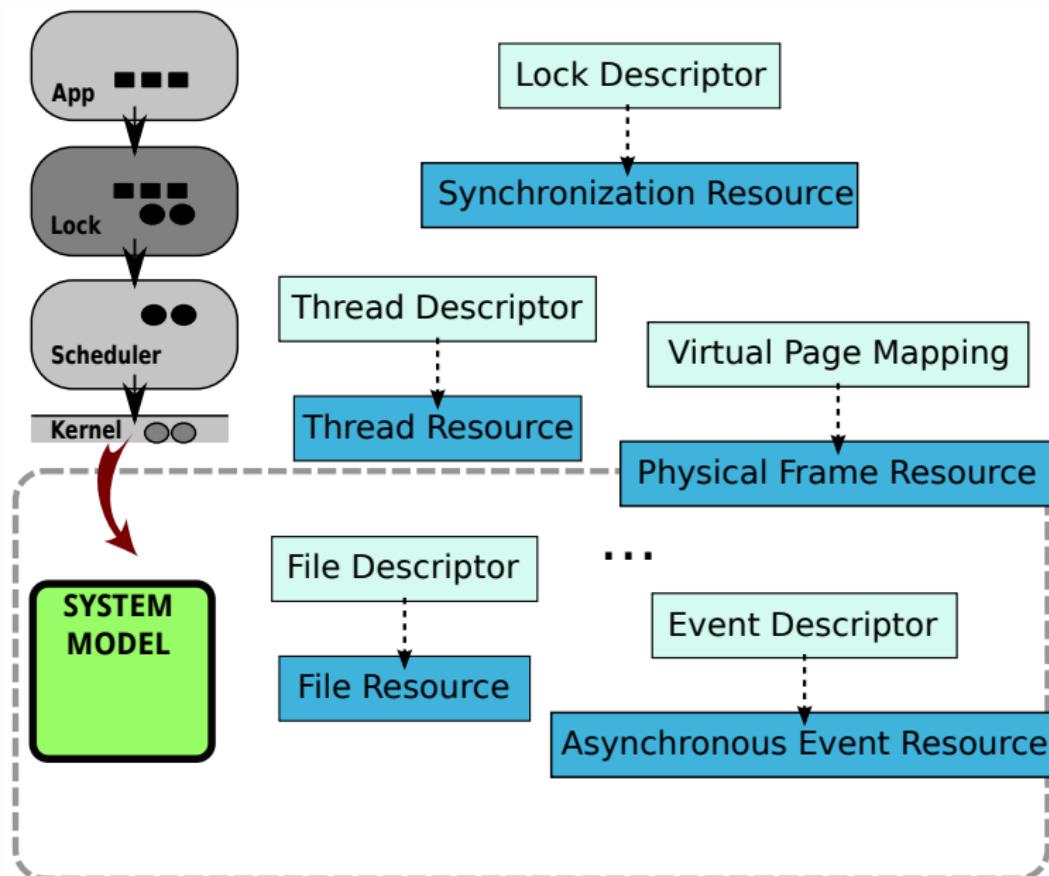
Descriptor-Resource Model



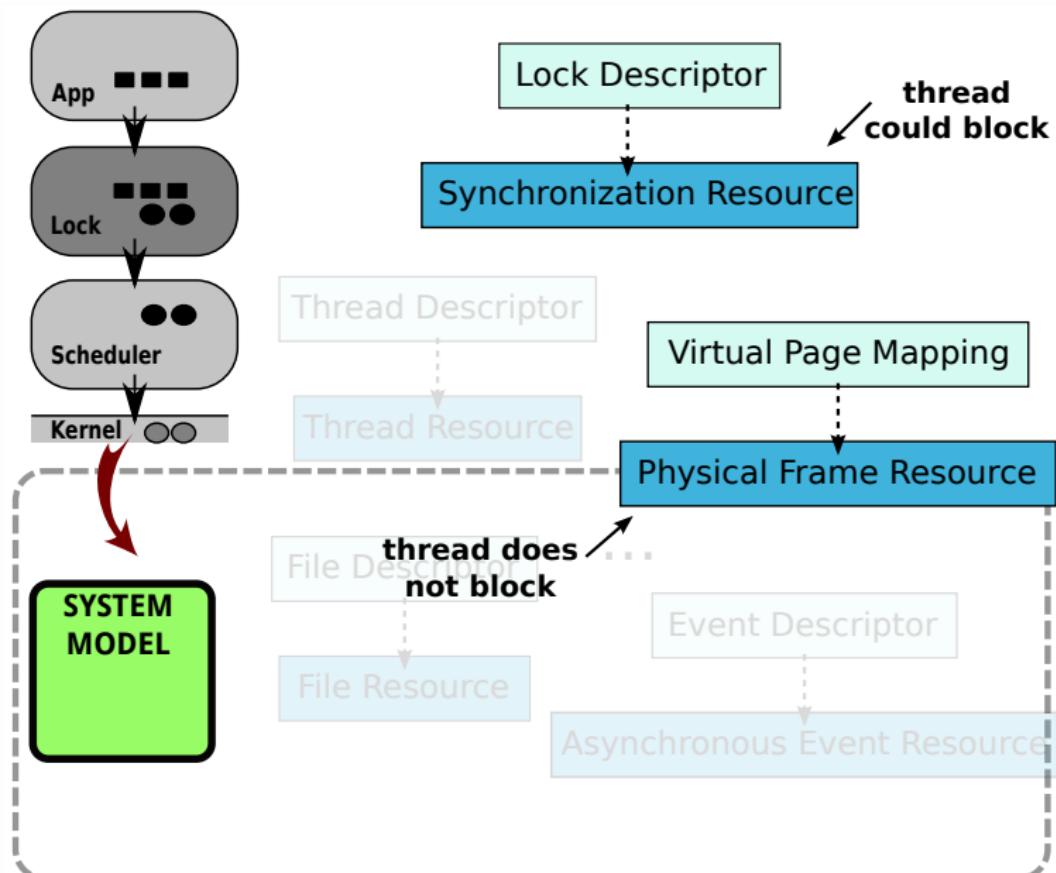
Descriptor-Resource Model



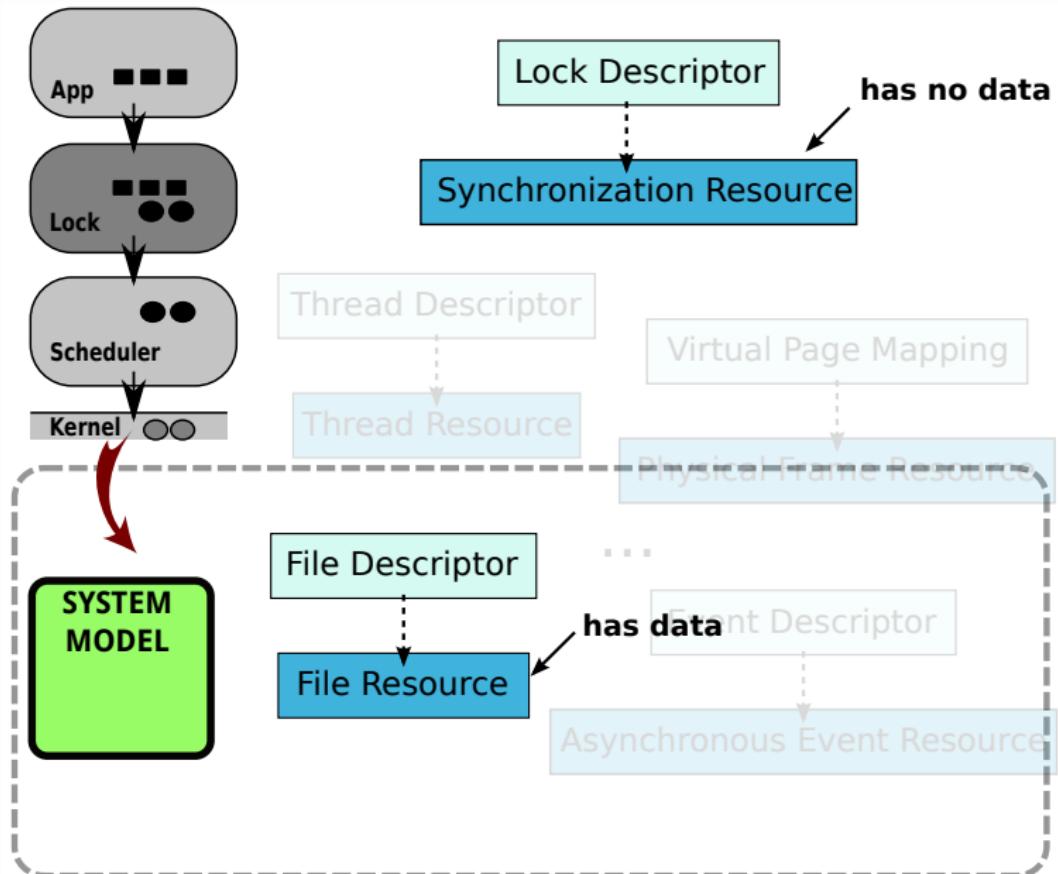
Descriptor-Resource Model



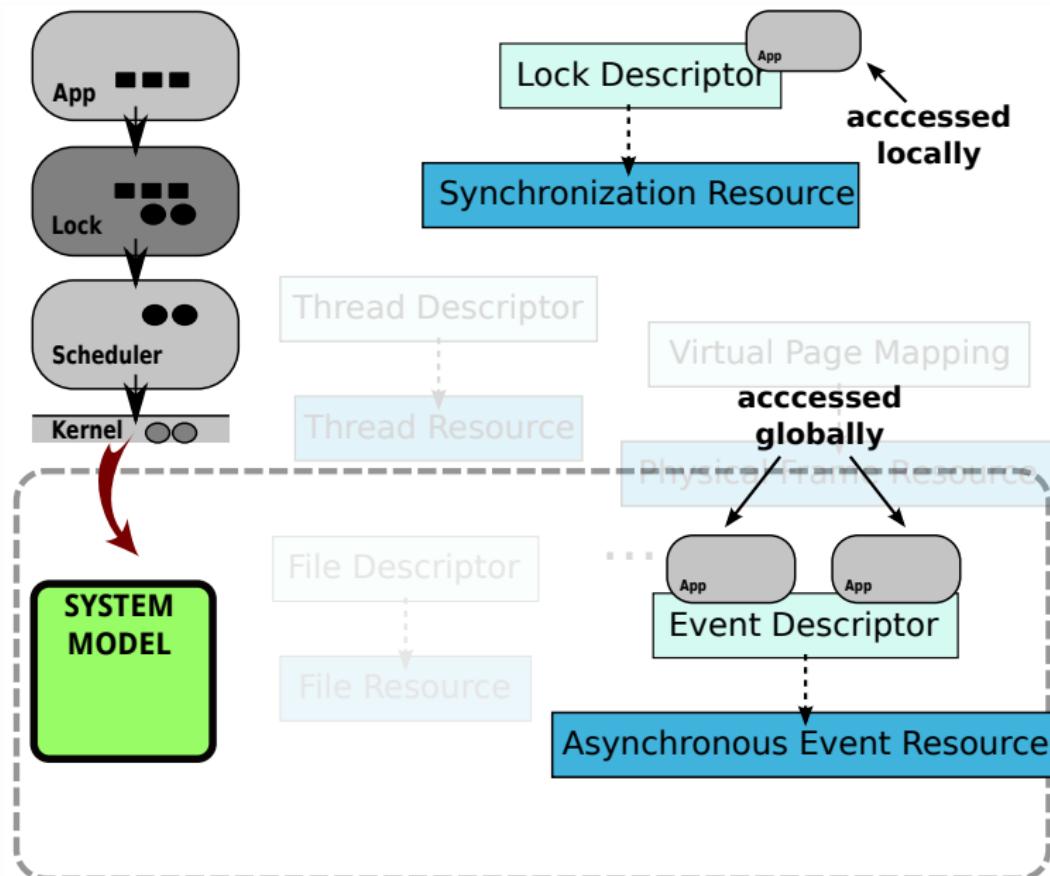
Descriptor-Resource Model



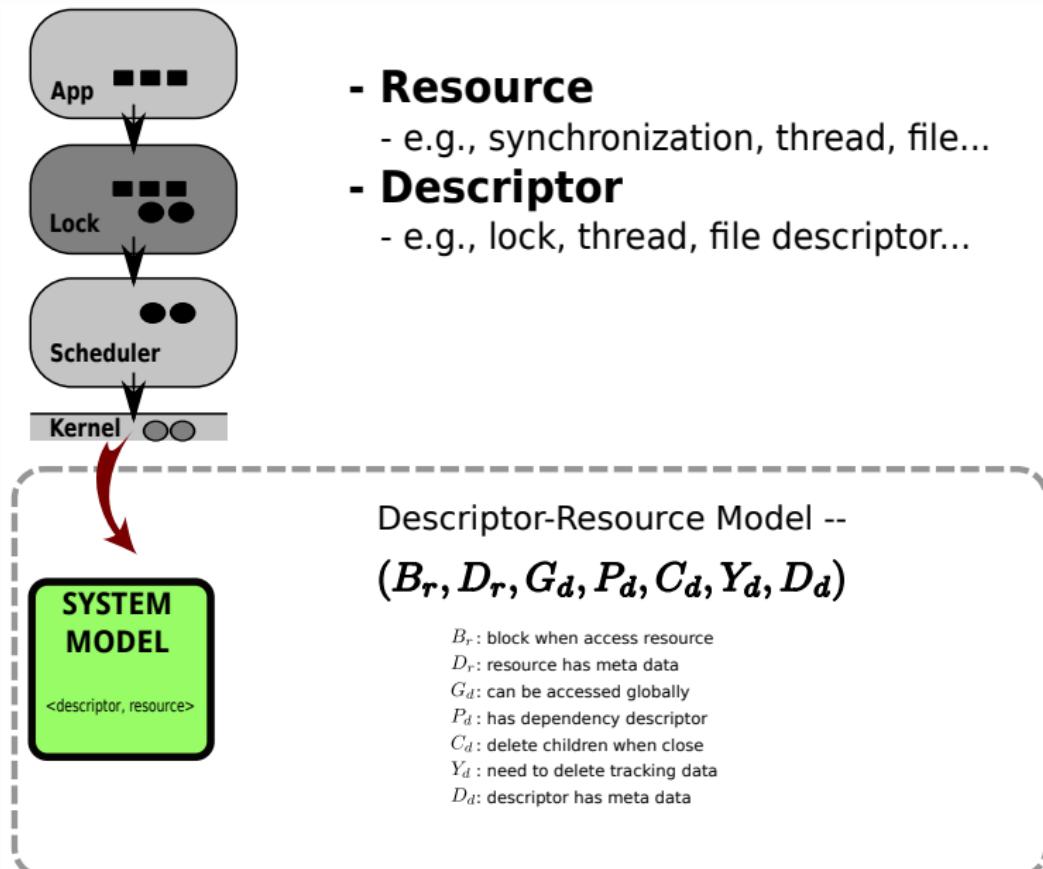
Descriptor-Resource Model



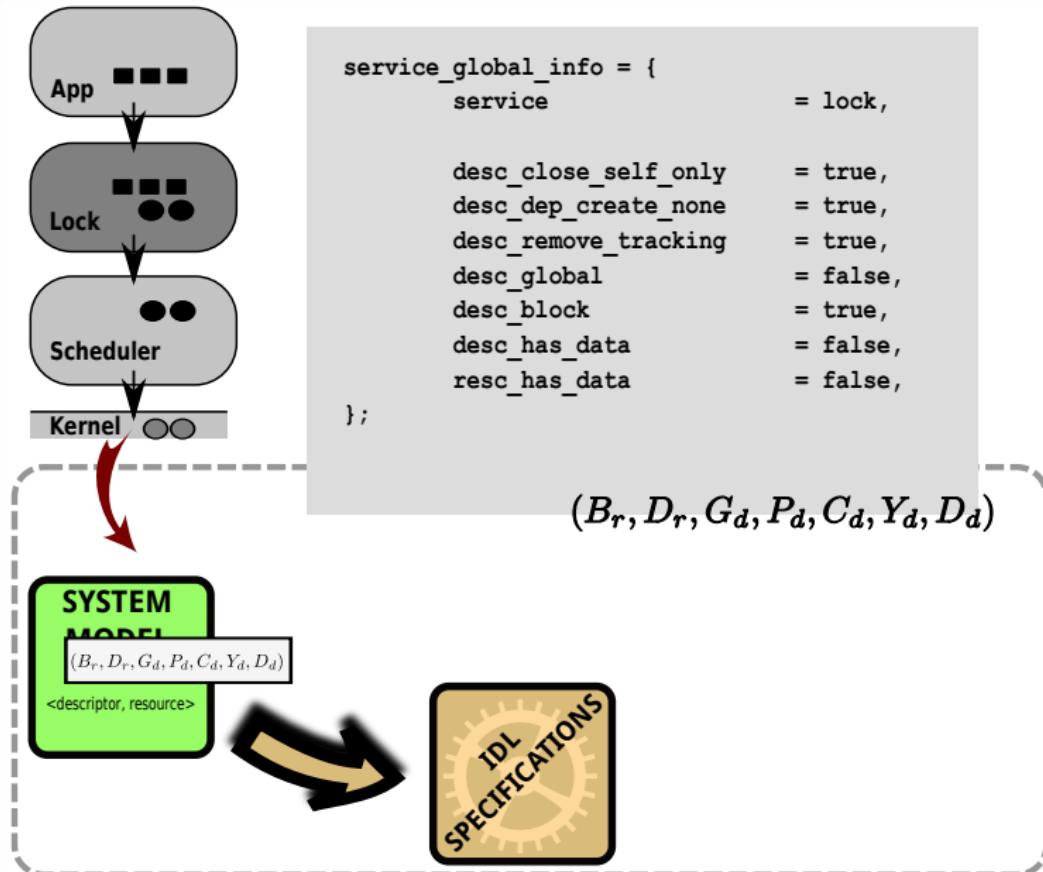
Descriptor-Resource Model



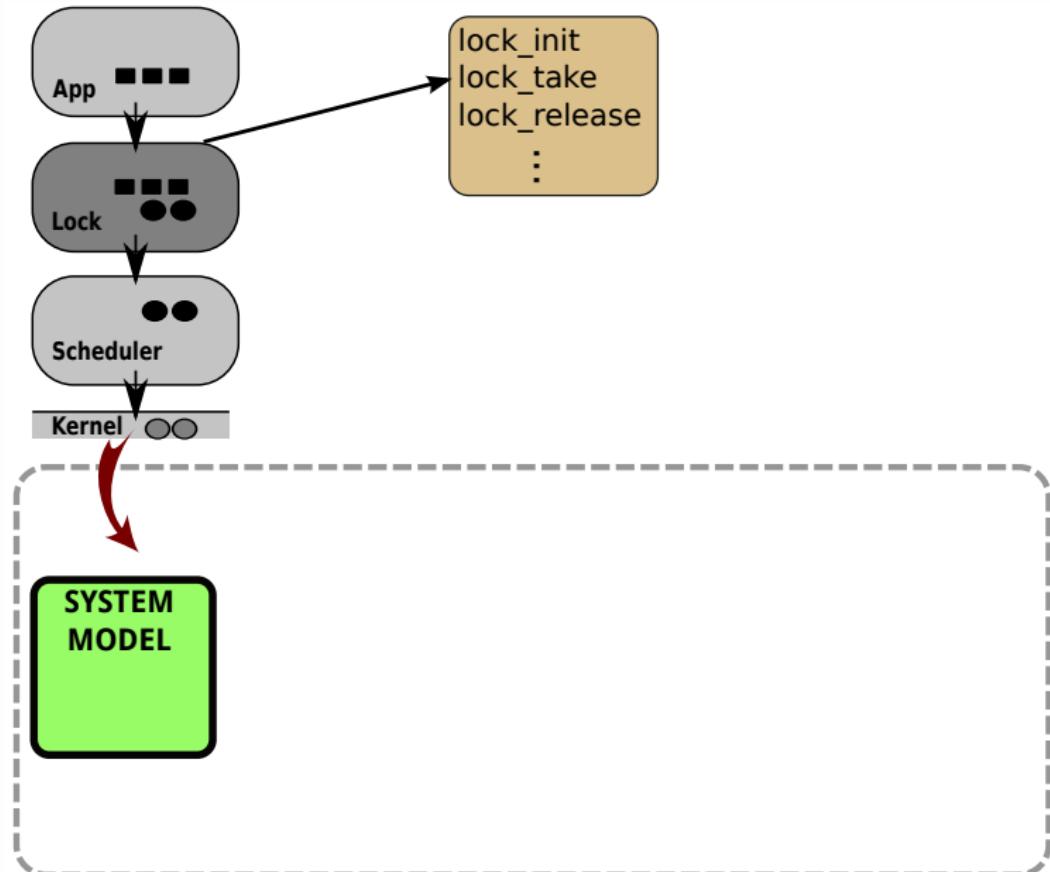
Descriptor-Resource Model



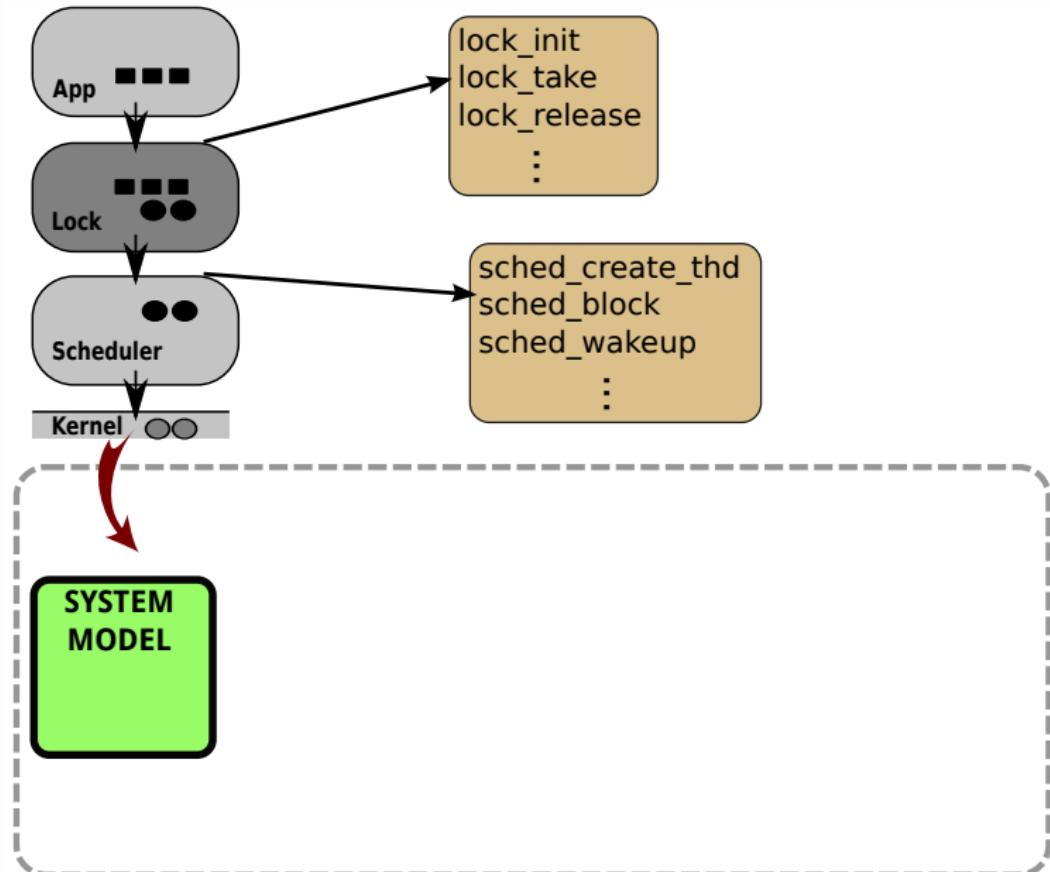
Descriptor-Resource Model → IDL



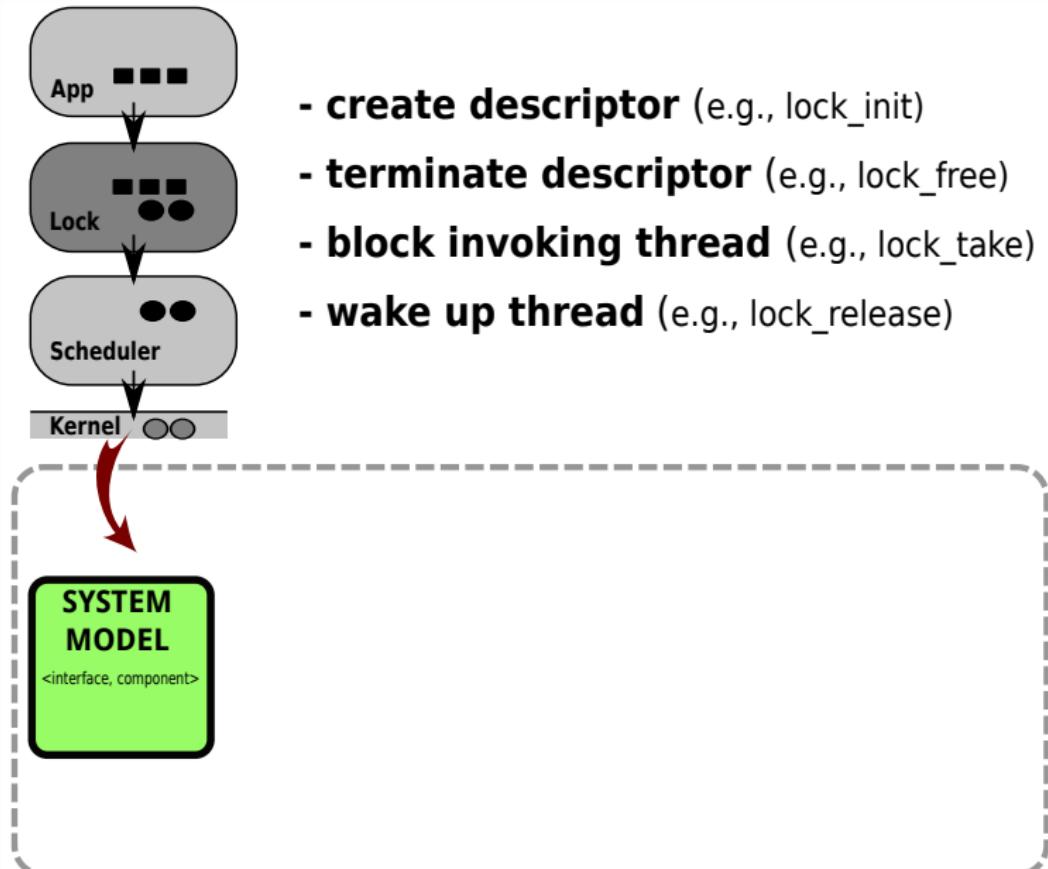
Component Interface Function



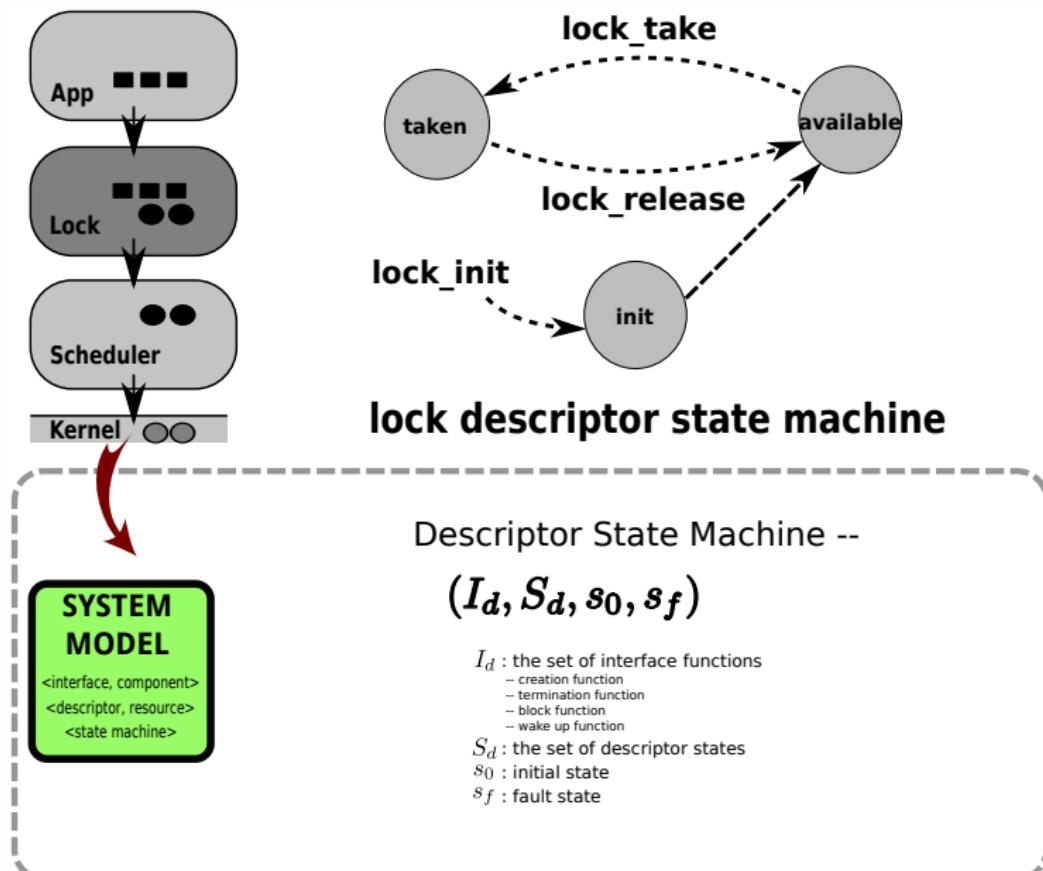
Component Interface Function



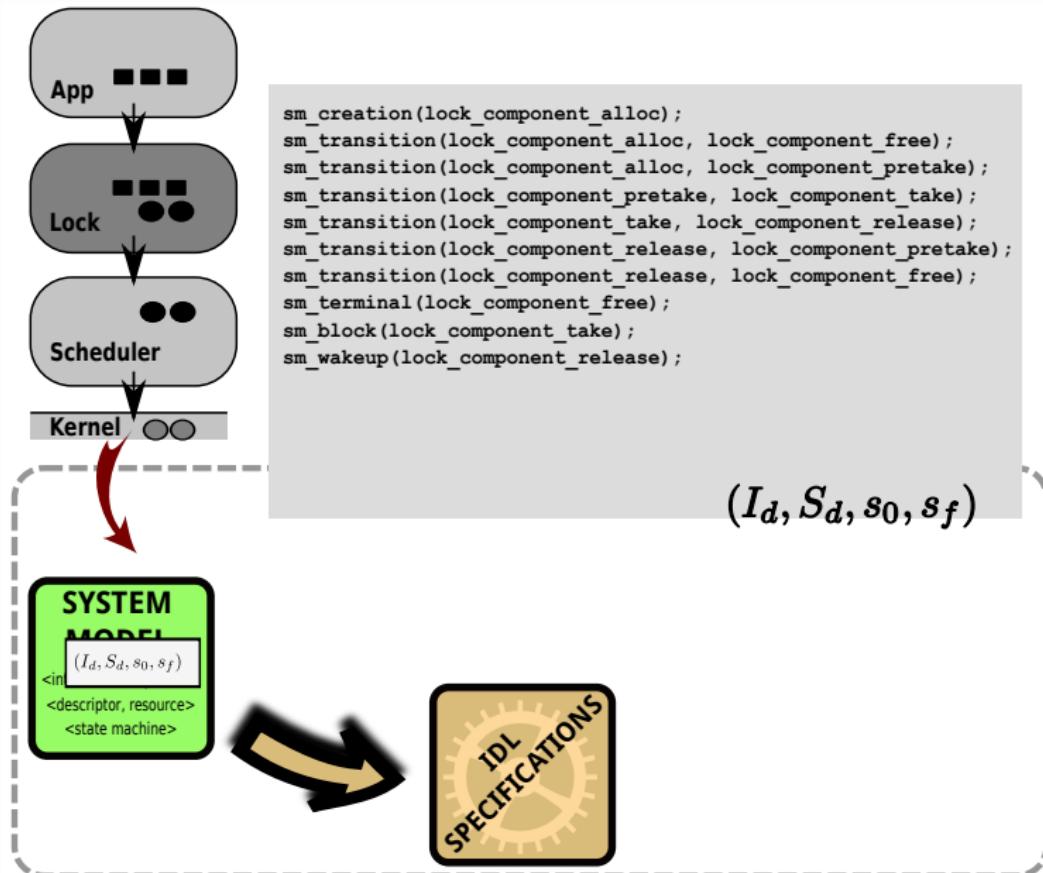
Component Interface Function



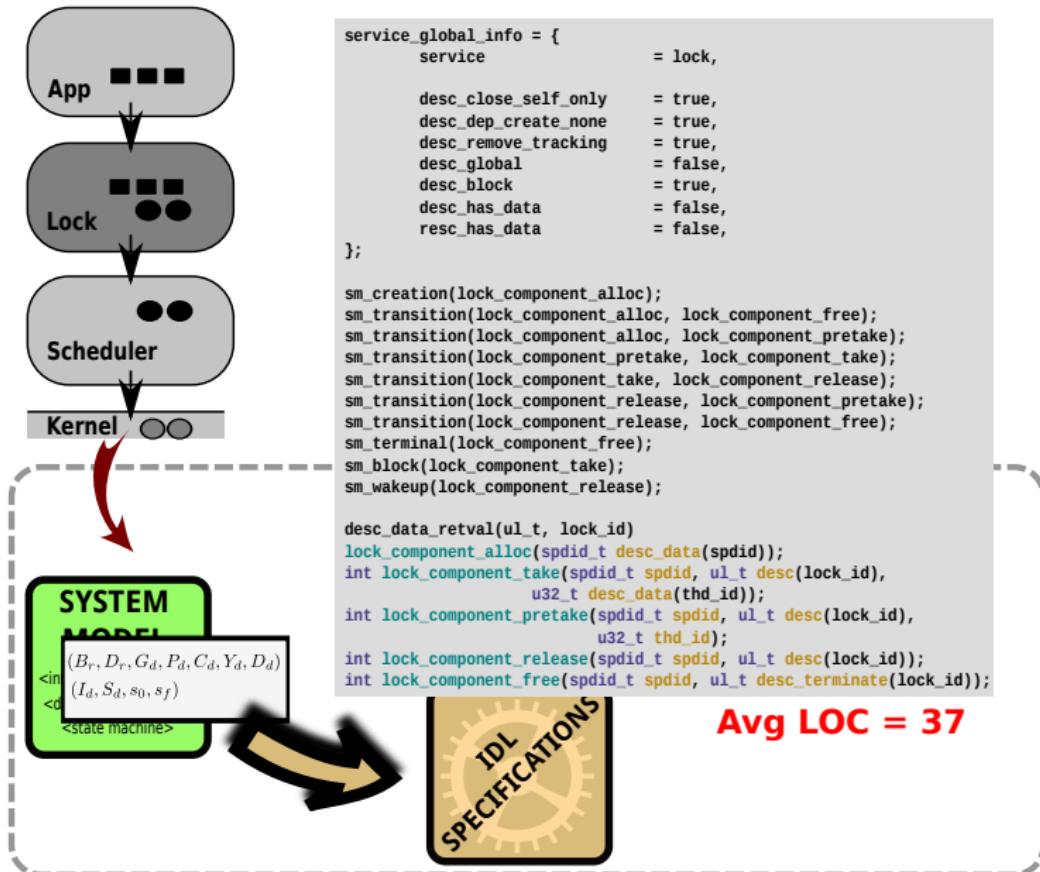
Descriptor State Machine



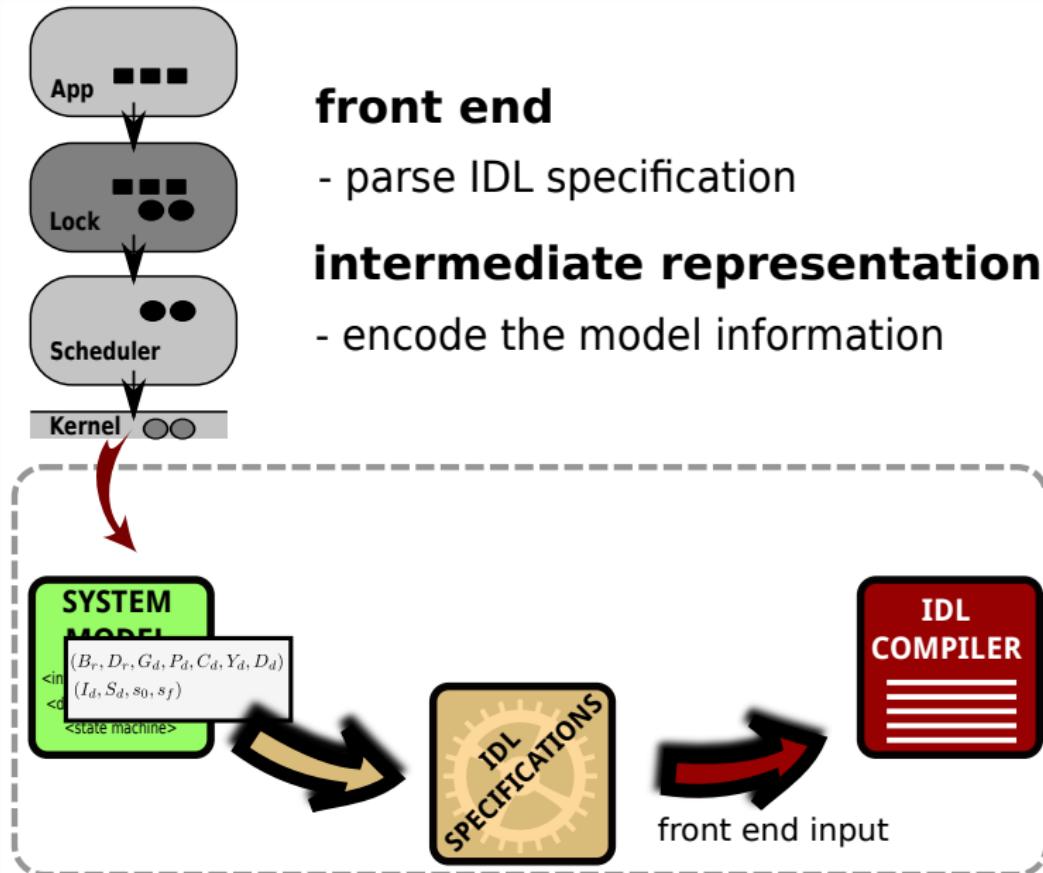
Descriptor State Machine → IDL



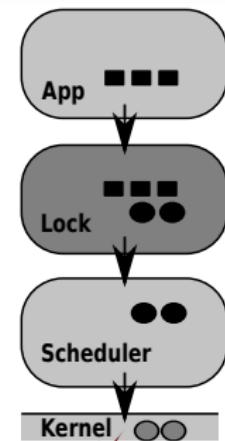
IDL-based Specification



Toward Generating Recovery Code

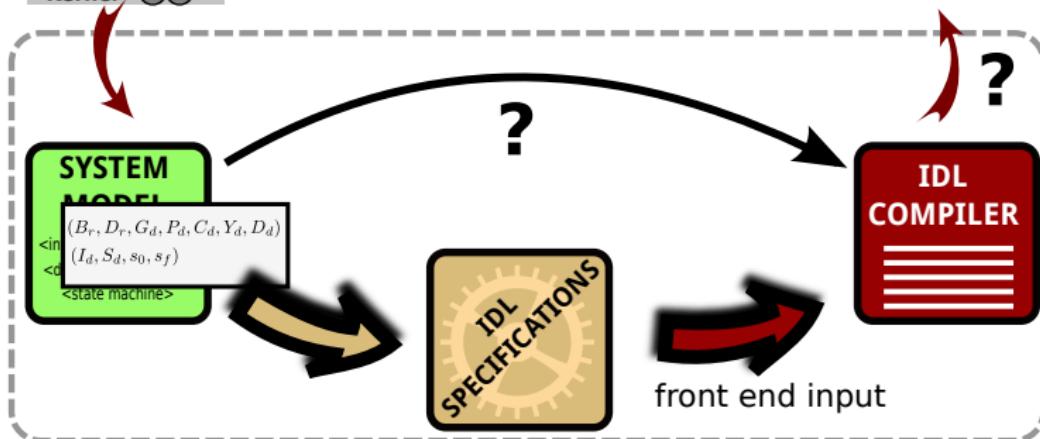


Toward Generating Recovery Code

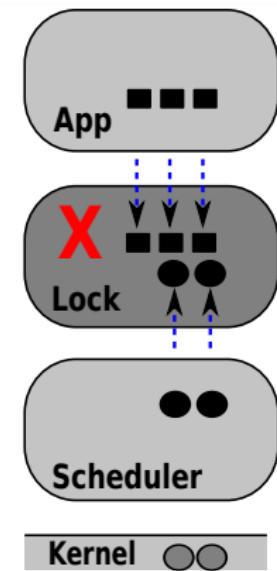


back end

- what recovery mechanism should be used?
- how to synthesize the code?



Interface-Driven Recovery Mechanisms

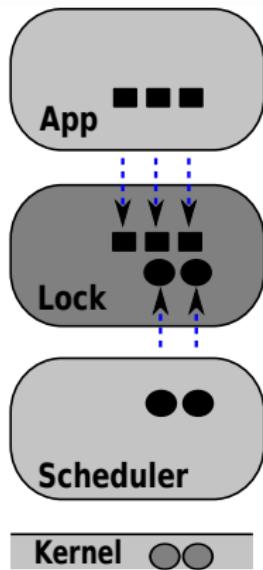


basic recovery

-- through component operation (**always**)

basic
recovery

Interface-Driven Recovery Mechanisms



basic recovery

-- through component operation (**always**)

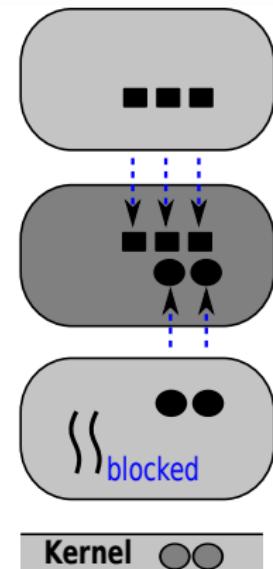
timing of recovery - on-demand

-- recover only when accessed (**always**)

basic
recovery

timing
of recovery

Interface-Driven Recovery Mechanisms



basic recovery

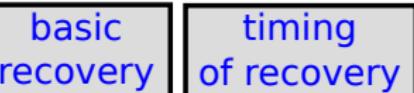
-- through component operation (**always**)

timing of recovery - on-demand

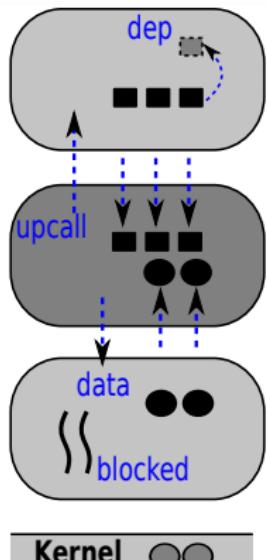
-- recover only when accessed (**always**)

timing of recovery - eager

-- eagerly wake up blocking threads (**B_r**)



Interface-Driven Recovery Mechanisms



basic recovery

- through component operation (**always**)

timing of recovery - on-demand

- recover only when accessed (**always**)

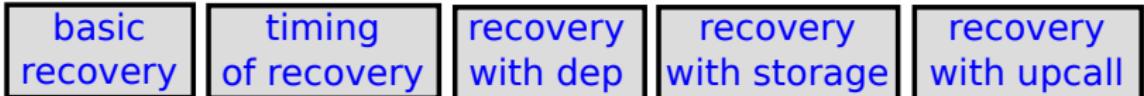
timing of recovery - eager

- eagerly wake up blocking threads (**B_r**)

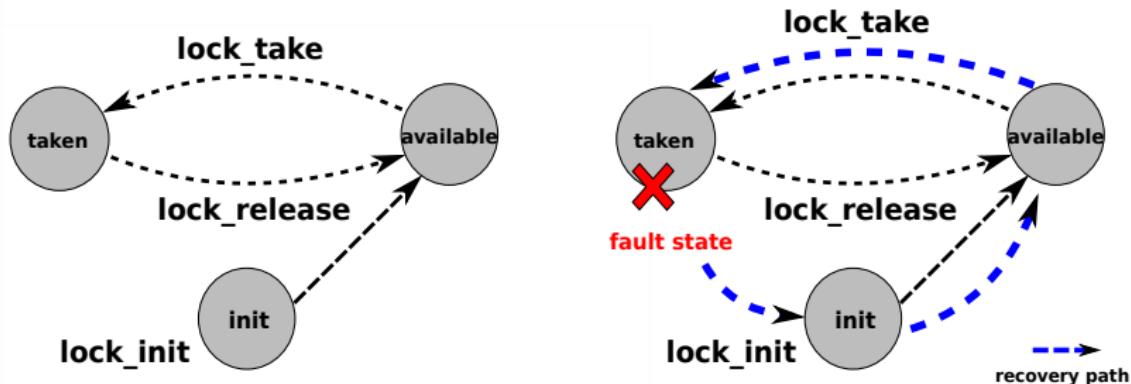
recovery with dependency

- require to reconstruct parent (**P_d**)
- require to reconstruct children (**C_d**)

:

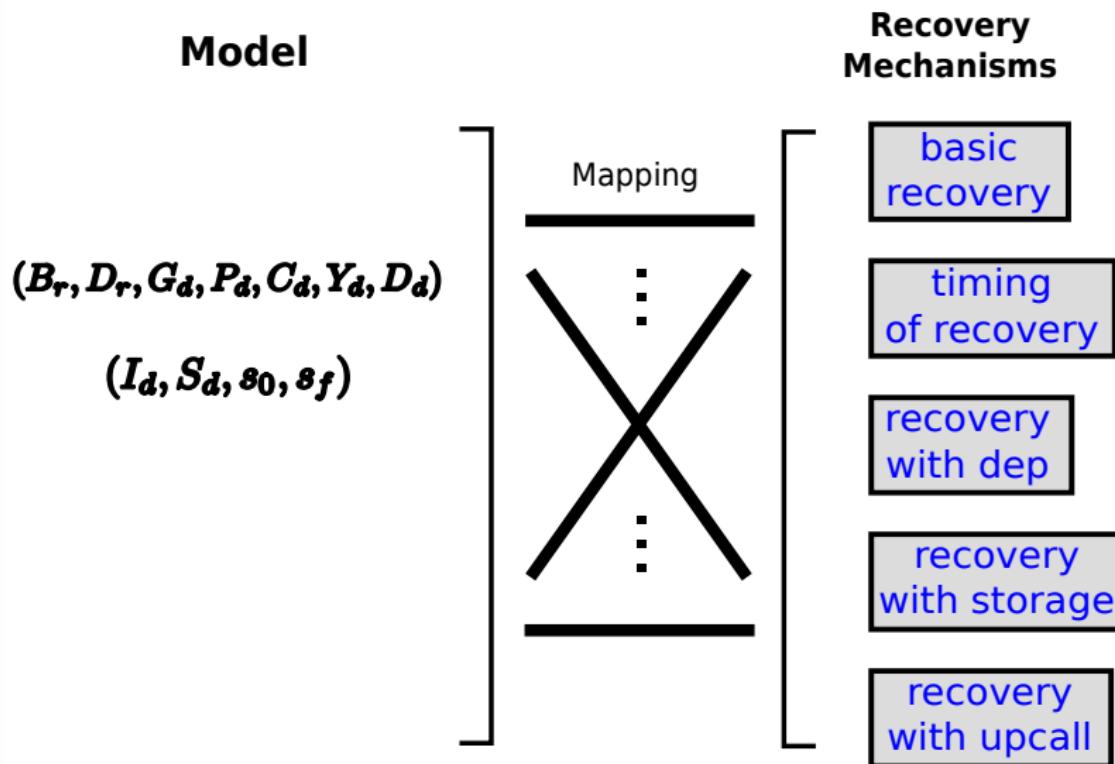


Interface-Driven Recovery Mechanisms



Descriptor State Machine

Model → Recovery Mechanisms



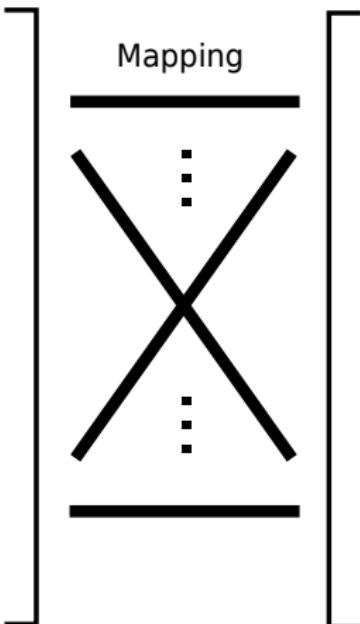
Synthesize the Recovery Code

Predicates

```

true
create_fn ∧ ¬G_d
wakeup_fn ∧ B_r
block_fn ∧ B_r
terminate_fn ∧ C_d
G_d ∧ P_d ∧ C_d ∧ D_d
D_r ∧ ¬P_d ∧ Y_d

```



Templates

Synthesize the Recovery Code

Predicates

Templates

Mapping

```
/* predicate: true */
CSTUB_FN(IDL_fntype, IDL_fname) (IDL_parsdecl) {
    long fault = 0;
    int ret = 0;
redo:
    cli_if_desc_update_IDL_fname(IDL_params);

    ret = cli_if_invoke_IDL_fname(IDL_params);
    if (fault){
        CSTUBFAULT_UPDATE();
        if (cli_if_desc_update_post_fault_IDL_fname()) goto redo;
    }
    ret = cli_if_track_IDL_fname(ret, IDL_params);
    return ret;
}
```

```
/* predicate: true */
CSTUB_FN(IDL_fntype, IDL_fname) (IDL_parsdecl) {
    long fault = 0;
    int ret = 0;
redo:
    cli_if_desc_update_IDL_fname(IDL_params);

    ret = cli_if_invoke_IDL_fname(IDL_params);
    if (fault){
        CSTUBFAULT_UPDATE();
        if (cli_if_desc_update_post_fault_IDL_fname()) goto redo;
    }
    ret = cli_if_track_IDL_fname(ret, IDL_params);
    return ret;
}
```

true

create_fn \wedge

wakeup_fn \wedge

block_fn \wedge I

terminate_

G_d \wedge P_d \wedge C

D_r \wedge \neg P_d \wedge

=

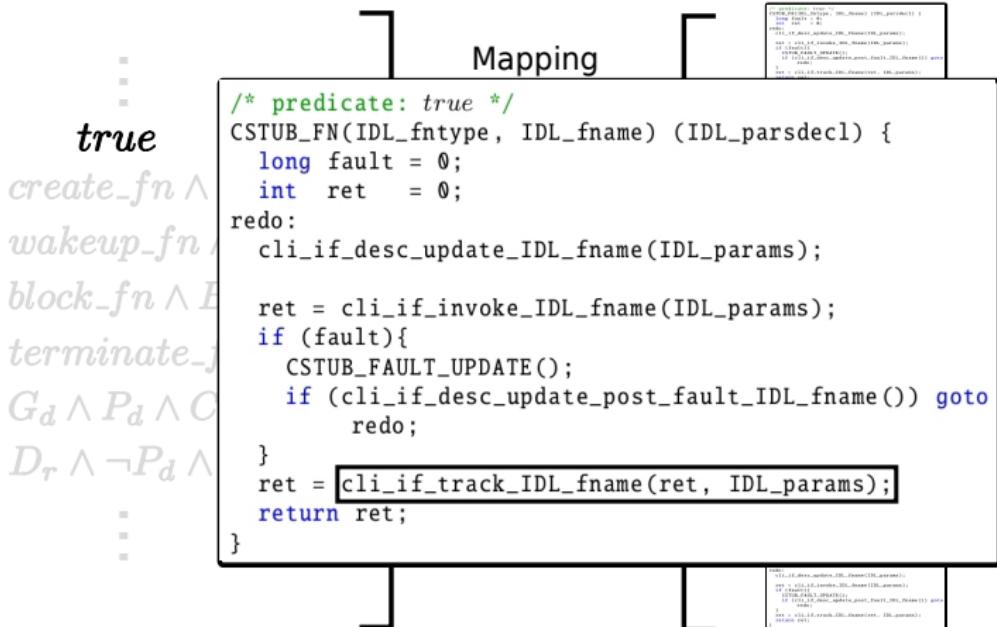
=

=

Synthesize the Recovery Code

Predicates

Templates



Synthesize the Recovery Code

Predicates

=
:
:
true
 $create_fn \wedge \neg G_d$
 $wakeup_fn \wedge \neg G_d$
 $block_fn \wedge \neg G_d$
 $terminate \wedge \neg G_d$
 $G_d \wedge P_d \wedge$
 $D_r \wedge \neg P_d$

Templates

```
/* predicate: f ∈ Icreatedr ∧ ¬Gdr */  
static inline int cli_if_track_IDL_fname(int ret,  
                                       IDL_parsdecl) {  
    if (ret == -EINVAL) return ret;  
    struct desc_track *desc = call_desc_alloc();  
    if (!desc) return -ENOMEM;  
    call_desc_track(desc, ret, IDL_params);  
  
    return desc->IDL_id;  
}
```

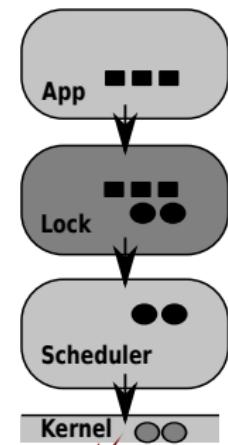
```
/* predicate: f ∈ Iwakeupdr ∧ ¬Gdr */  
static inline int cli_if_wakeup_IDL_fname(int ret,  
                                         IDL_parsdecl) {  
    if (ret == -EINVAL) return ret;  
    struct desc_track *desc = call_desc_alloc();  
    if (!desc) return -ENOMEM;  
    call_desc_track(desc, ret, IDL_params);  
  
    return desc->IDL_id;
```

```
/* predicate: f ∈ Iblockdr ∧ ¬Gdr */  
static inline int cli_if_block_IDL_fname(int ret,  
                                       IDL_parsdecl) {  
    if (ret == -EINVAL) return ret;  
    struct desc_track *desc = call_desc_alloc();  
    if (!desc) return -ENOMEM;  
    call_desc_track(desc, ret, IDL_params);  
  
    return desc->IDL_id;
```

Mapping



Generate Recovery Code



front end

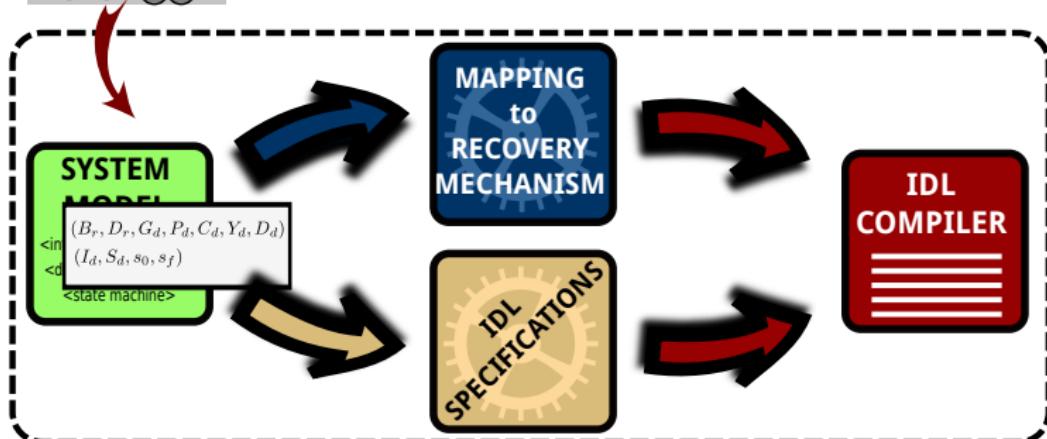
- parse IDL-based specification

intermediate representation

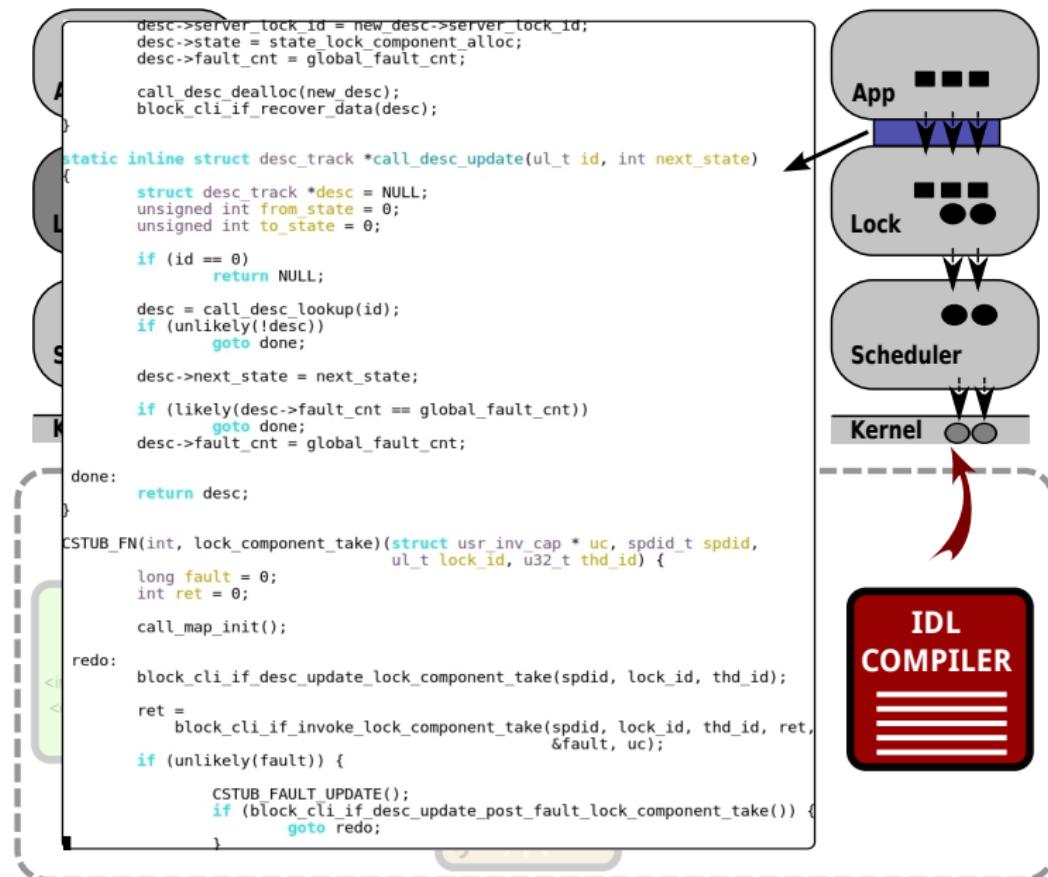
- encode the model information

back end

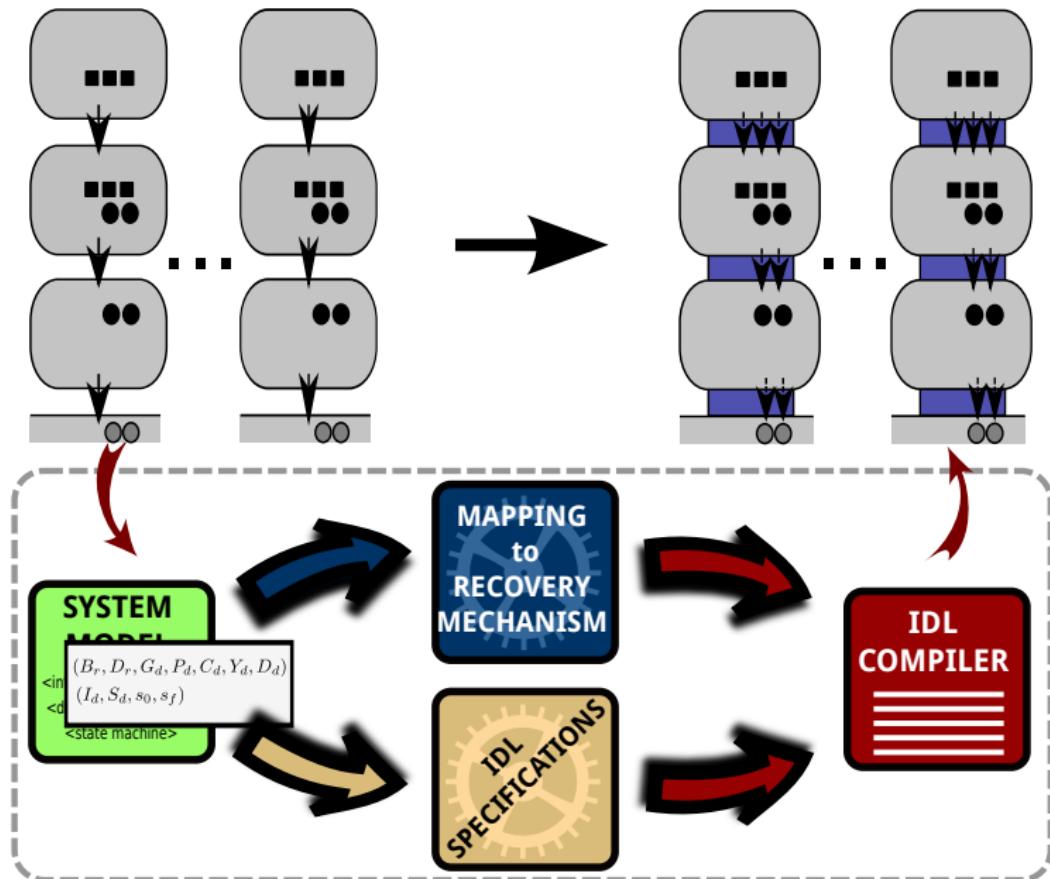
- evaluate predicates
- generate the code from templates



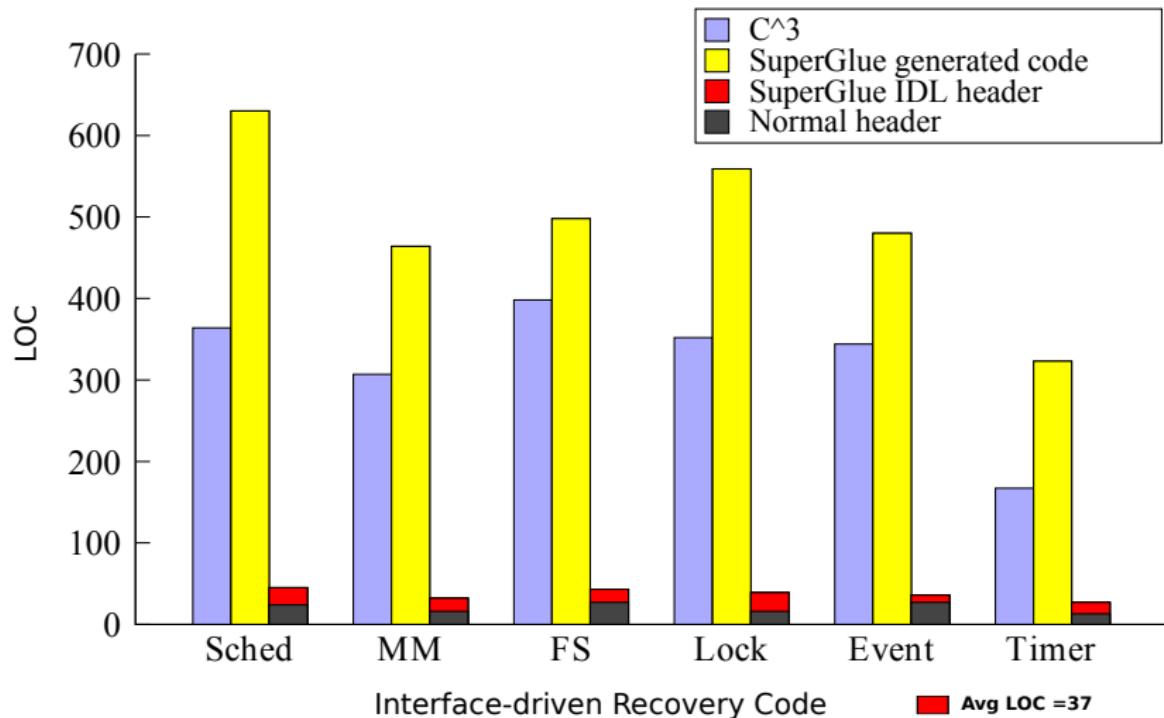
Generate Recovery Code



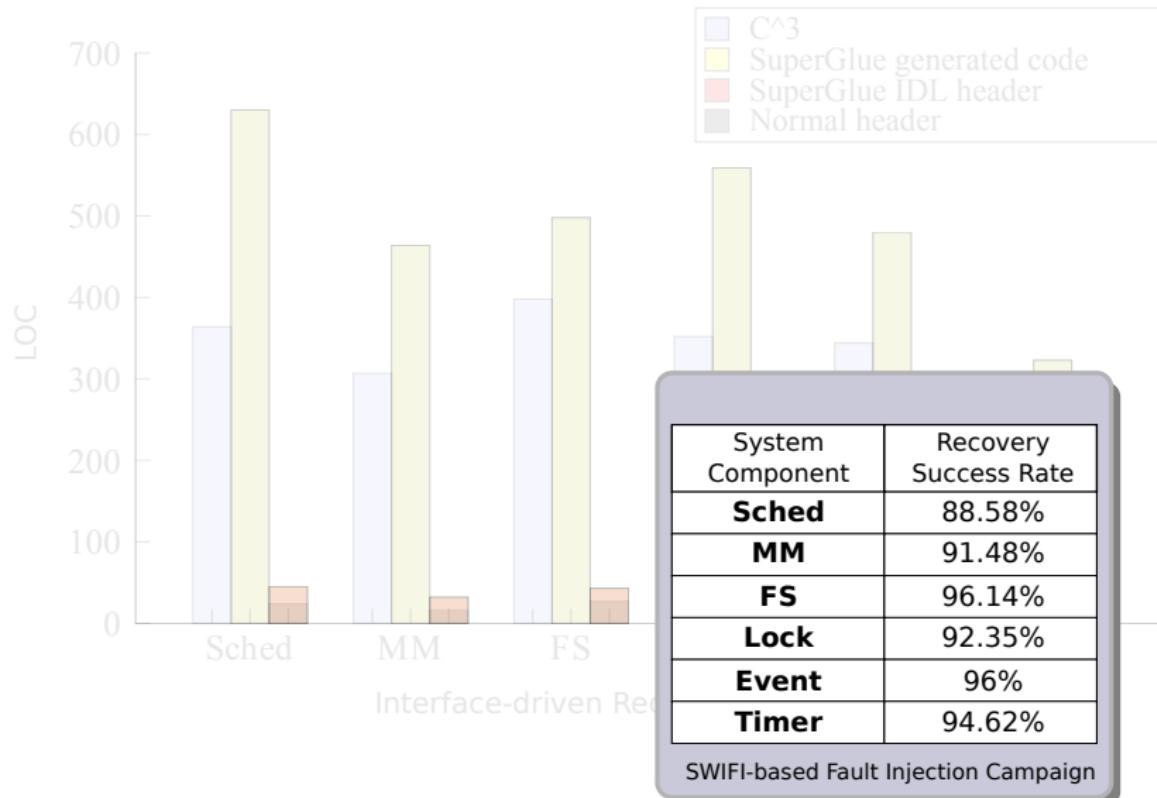
Generate Recovery Code



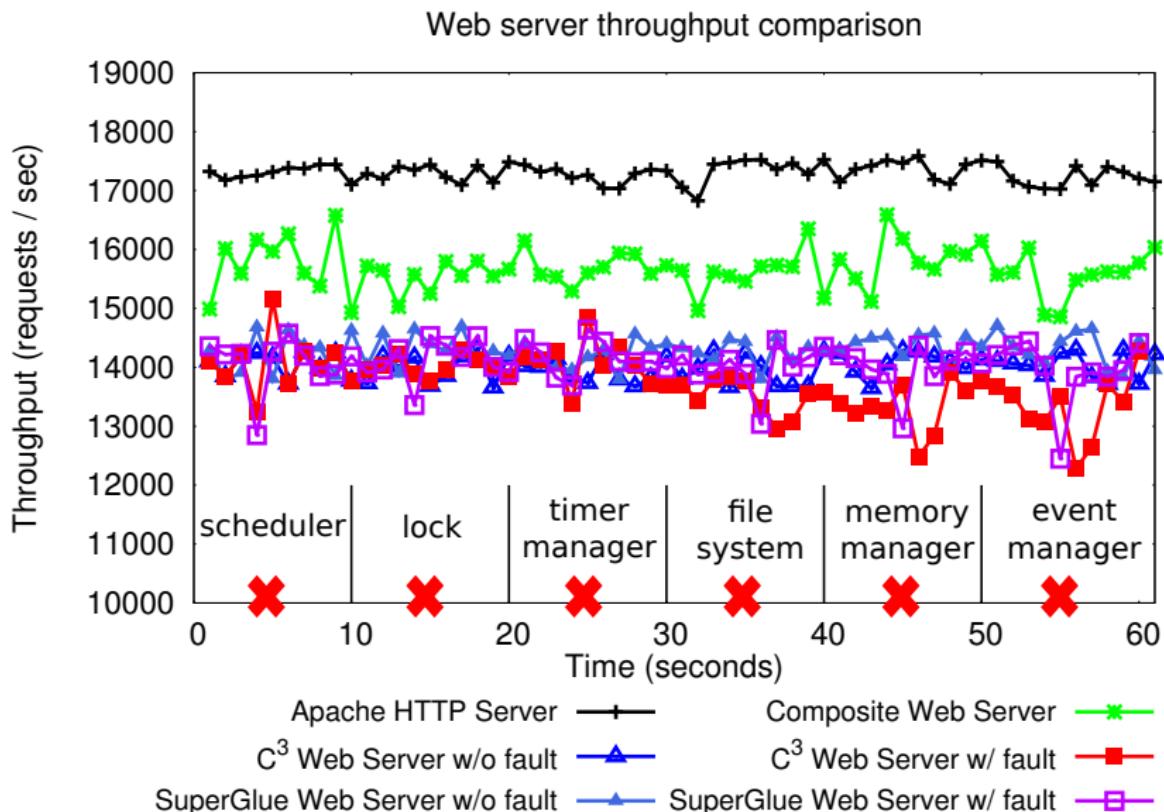
Code Generation Result



Fault Injection Result



Web Server Evaluation with Injected Faults



Conclusion on SuperGlue

SuperGlue – code generation for system-level fault tolerance

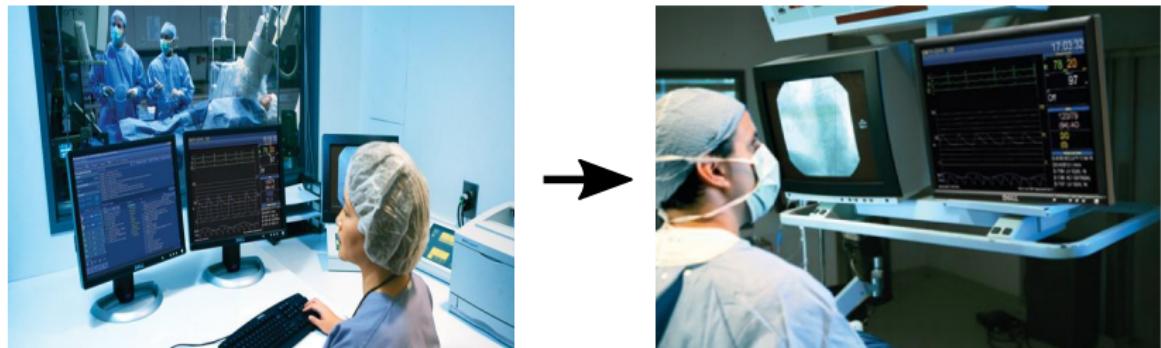
- descriptor-resource **model** and descriptor **state machine**
- **IDL-based** declarative specifications
- **compiler** for synthesizing C³-style recovery code

Publication

J. Song, G. Bloom, and G. Parmer, "SuperGlue: IDL-Based, System-Level Fault Tolerance for Embedded Systems" in DSN, 2016

Threats and Vulnerabilities in Cyber-Physical Systems

Medical equipment (Merge Hemo) crashed in Feb 2016



- Equipment went black and manually rebooted
- 5 minutes delay

Threats and Vulnerabilities in Cyber-Physical Systems

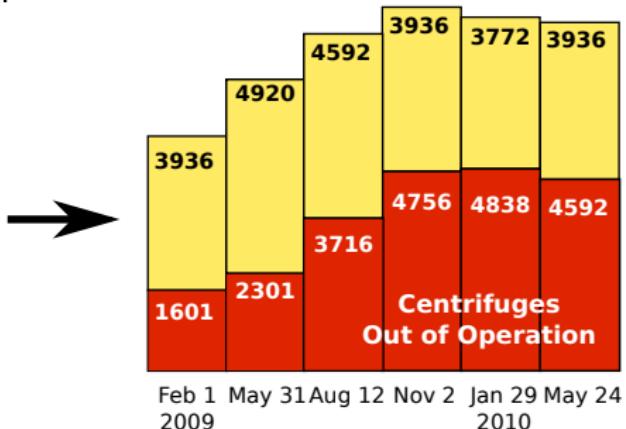
Medical equipment (Merge Hemo) crashed in Feb 2016



- Equipment went black and manually rebooted
- 5 minutes delay

Threats and Vulnerabilities in Cyber-Physical Systems

The World's First Cyber Weapon: Stuxnet

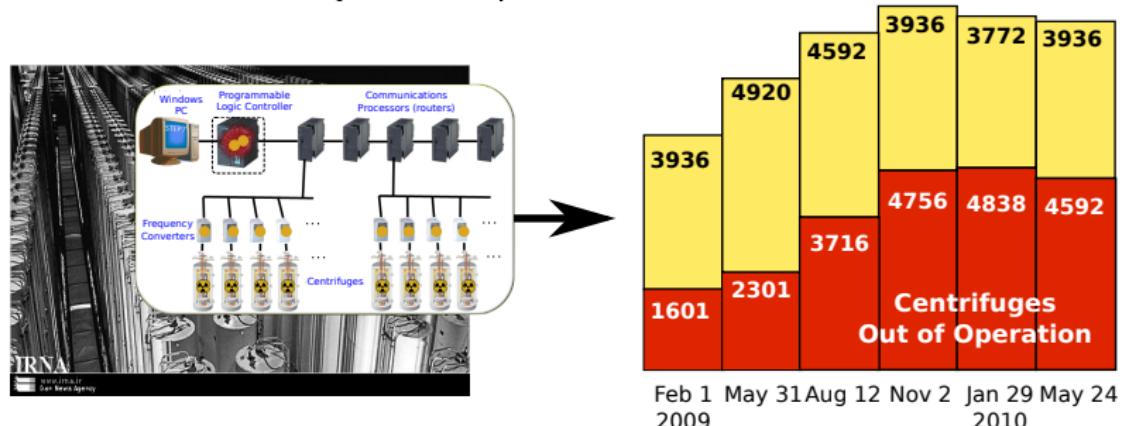


- Physical infrastructure destruction
- Could raise nuclear tensions between the nations

Adapted from "The 1-hour Guide to Stuxnet" presentation by Carey Nachenberg, Symantec Fellow

Threats and Vulnerabilities in Cyber-Physical Systems

The World's First Cyber Weapon: Stuxnet



- Physical infrastructure destruction
- Could raise nuclear tensions between the nations

Adapted from "The 1-hour Guide to Stuxnet" presentation by Carey Nachenberg, Symantec Fellow

CAML– System-Level Anomaly Detection Infrastructure

CAML (*an ongoing collaborative project*)

- system-level anomaly detection
- predictably
- using machine learning techniques

(by Charles River Analytics Inc)



Main ideas

- Extend the system events logging infrastructure from C'Mon
- Pre-process logged events
- Use probabilistic techniques to detect the anomaly

1 Tracking

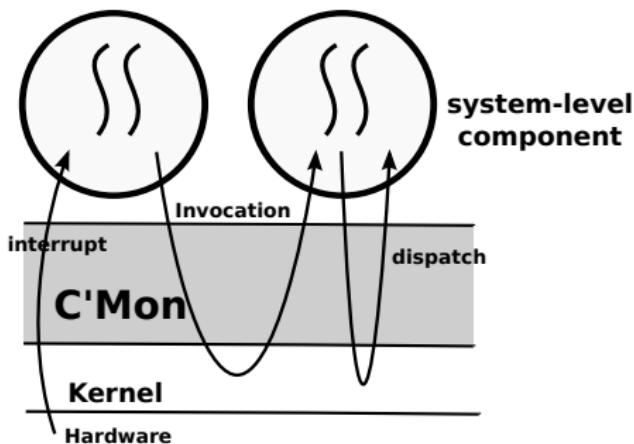
- system communication events
- system timing information

2 Detection

- temporal constraints violation

3 User-level

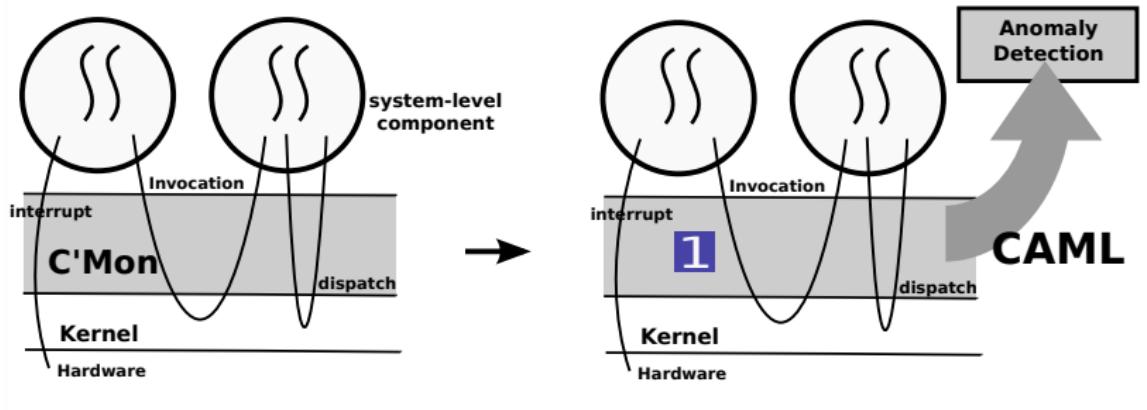
- no kernel modification



Publication

J. Song, and G. Parmer, "C'Mon: a Predictable Monitoring Infrastructure for System-Level Latent Fault Detection and Recovery", in RTAS 2015

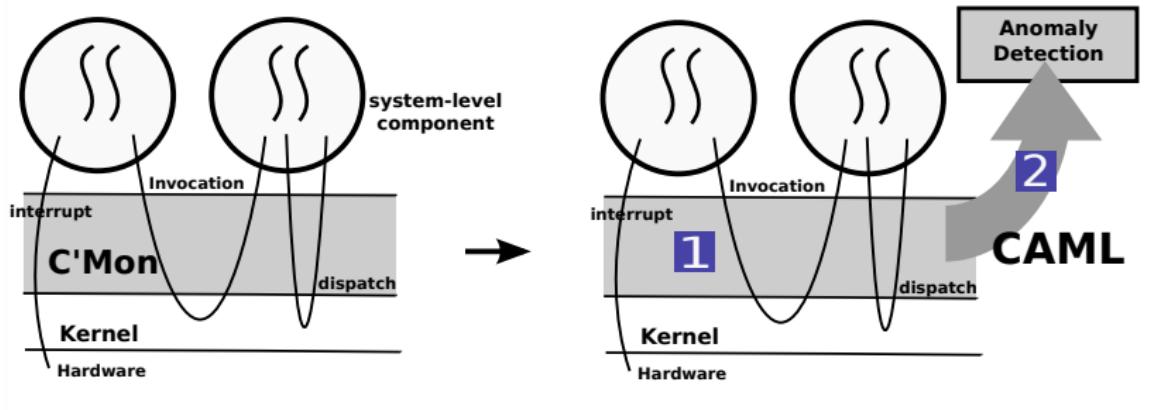
CAML - System-Level Anomaly Detection



Idea →

- 1 logging events (same as in C'Mon)

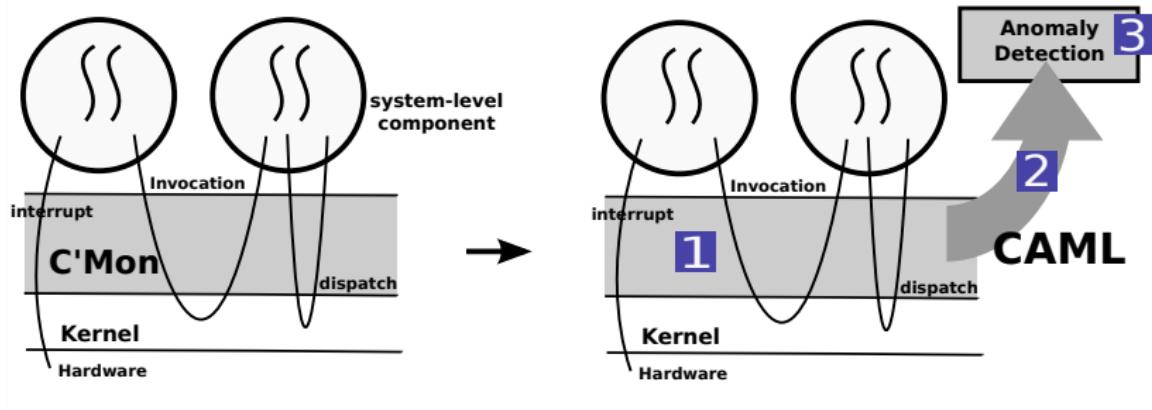
CAML - System-Level Anomaly Detection



Idea →

- 1 logging events (same as in C'Mon)
- 2 pre-processing logged events

CAML - System-Level Anomaly Detection

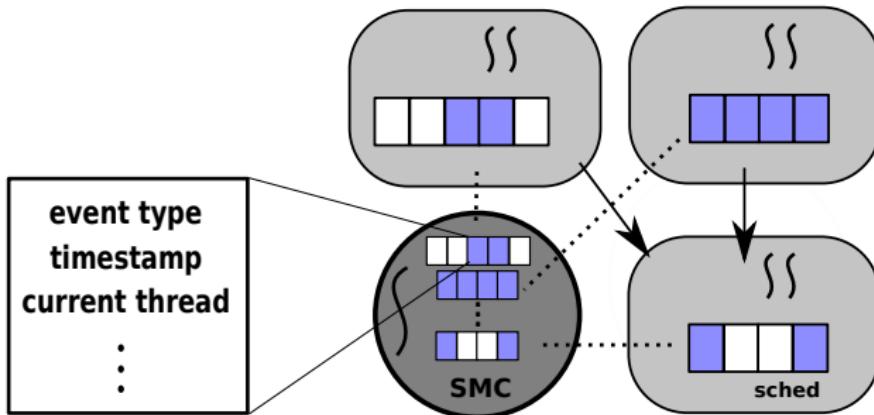


Idea →

- 1 logging events (same as in C'Mon)
- 2 pre-processing logged events
- 3 detecting anomalies using machine learning techniques

SMC – System Monitor Component

- SMC component
 - log events

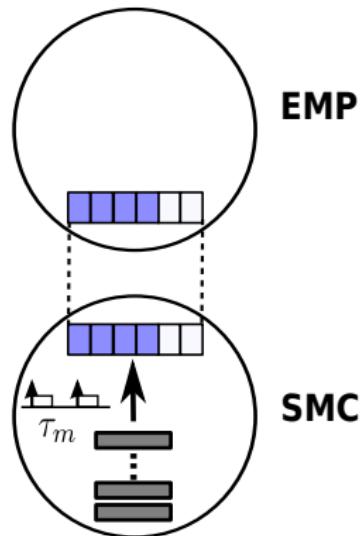


- ▶ *wait-free multi-producer single consumer buffer (same as in C'Mon)*

SMC – System Monitor Component

- SMC component
 - log events
 - copy events into the shared buffer between SMC and EMP

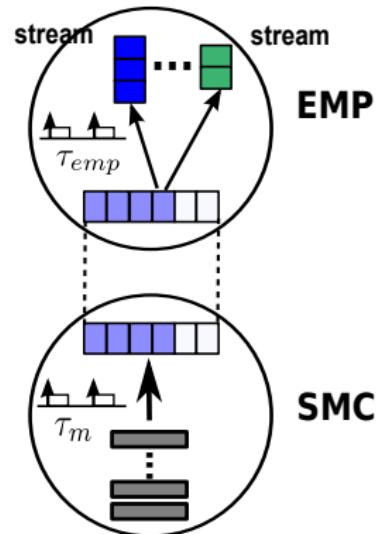
► assume shared buffer is large enough



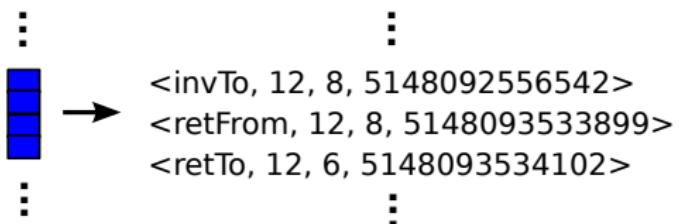
EMP – Event Multiplexing Component

- SMC component
 - log events
 - copy events into the shared buffer between SMC and EMP
- EMP component
 - multiplex logged events into streams

► pre-processing execution time

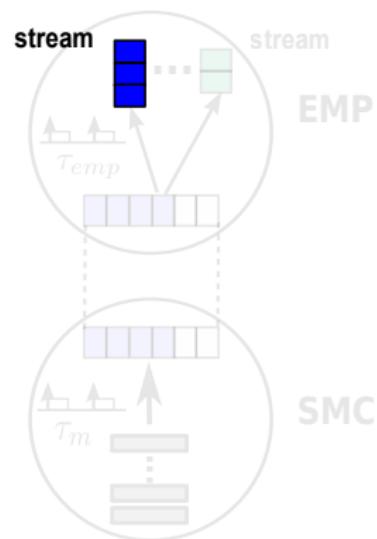


Stream

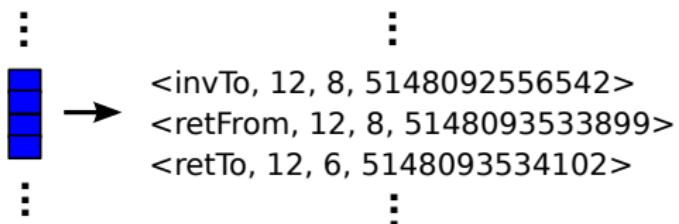


Component invocations of a thread

- <invFrom, thread, component, time>
- <invTo, thread, component, time>
- <retFrom, thread, component, time>
- <retTo, thread, component, time>

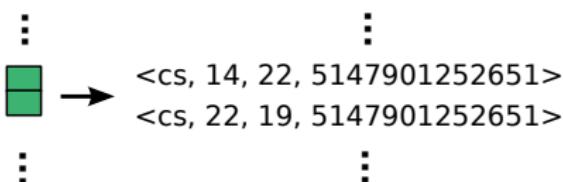


Stream



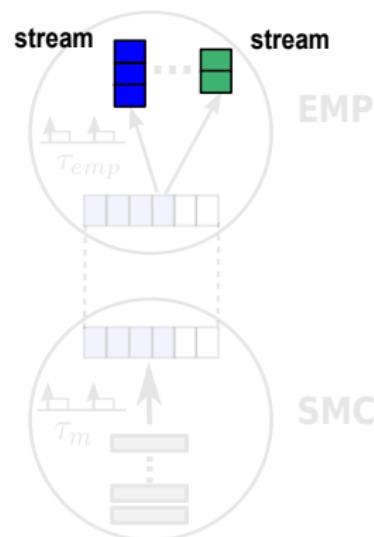
Component invocations of a thread

- <invFrom, thread, component, time>
- <invTo, thread, component, time>
- <retFrom, thread, component, time>
- <retTo, thread, component, time>



thread dispatching

- <cs, from_thread, to_thread, time>



Per-thread stream

- sequence of component invocation/return of a thread
- number of times context switch to/from a thread

:

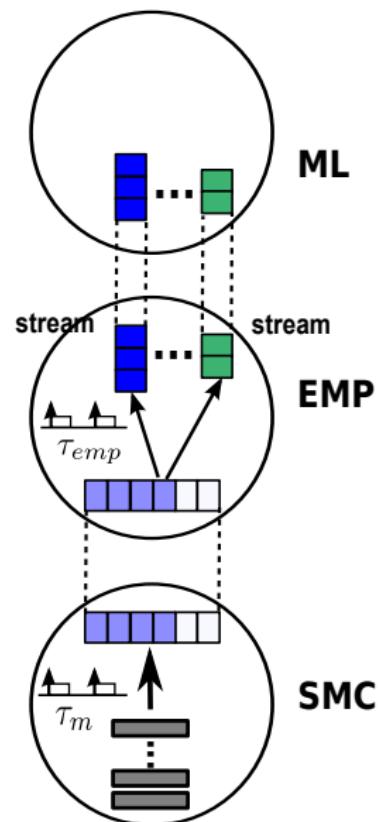
Per-component stream

- invocations/returns of a component by any thread
- number of times any thread is interrupted in a component

:

EMP – Event Multiplexing Component

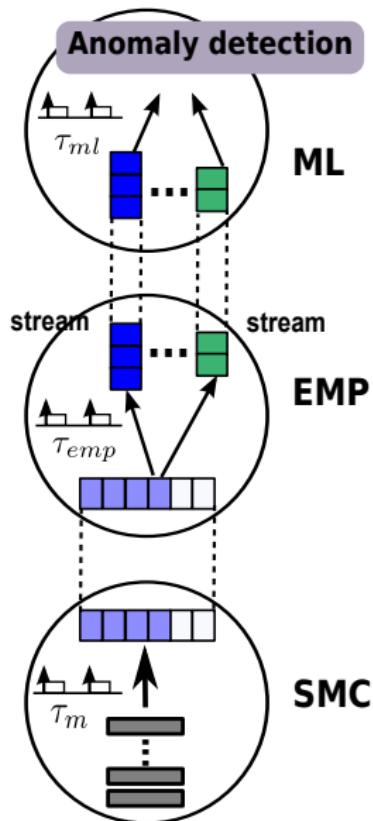
- SMC component
 - log events
 - copy events into the shared buffer between SMC and EMP
- EMP component
 - multiplex logged events into streams
 - copy streams into the shared buffers between EMP and ML



► stream → another layer of abstraction about the system execution

ML – Anomaly Detection Component

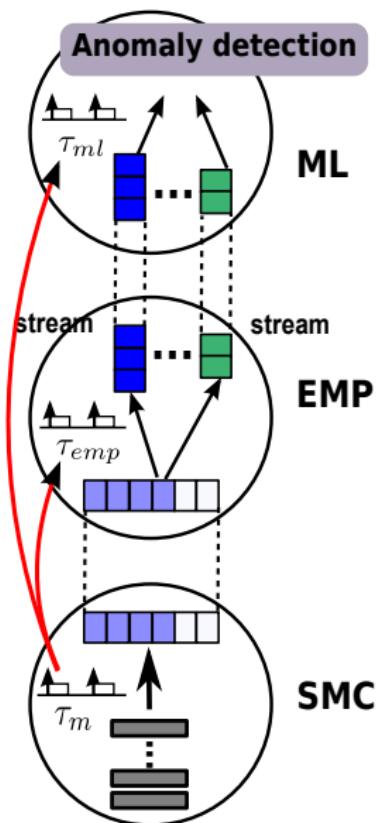
- SMC component
 - log events
 - copy events into the shared buffer between SMC and EMP
- EMP component
 - multiplex logged events into streams
 - copy streams into the shared buffers between EMP and ML
- ML component
 - run machine learning algorithms to detect the anomaly



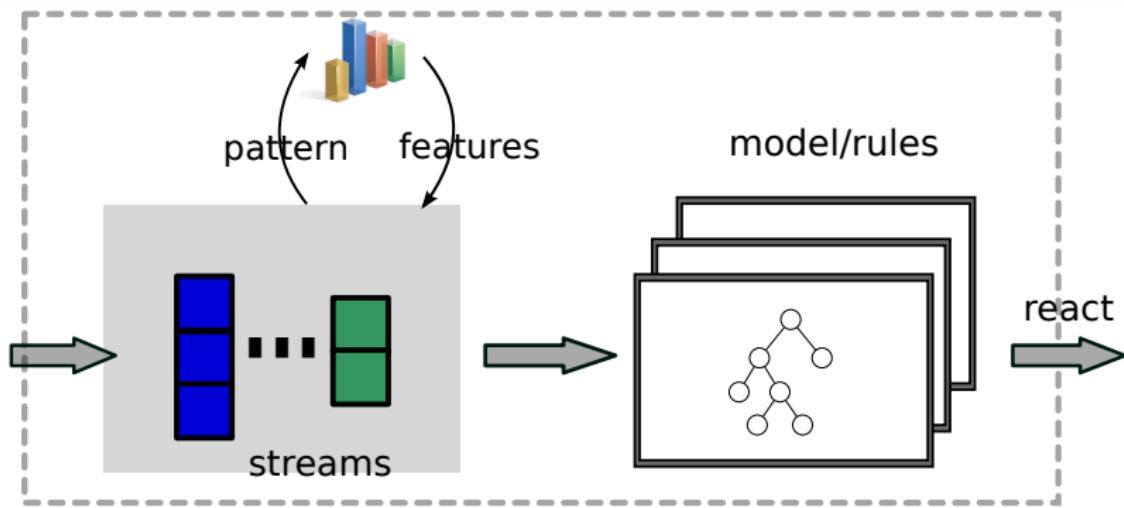
ML – Anomaly Detection Component

- SMC component
 - log events
 - copy events into the shared buffer between SMC and EMP
- EMP component
 - multiplex logged events into streams
 - copy streams into the shared buffers between EMP and ML
- ML component
 - run machine learning algorithms to detect the anomaly

► needs to ensure tasks τ_{emp} and τ_{ml} are activated at right time



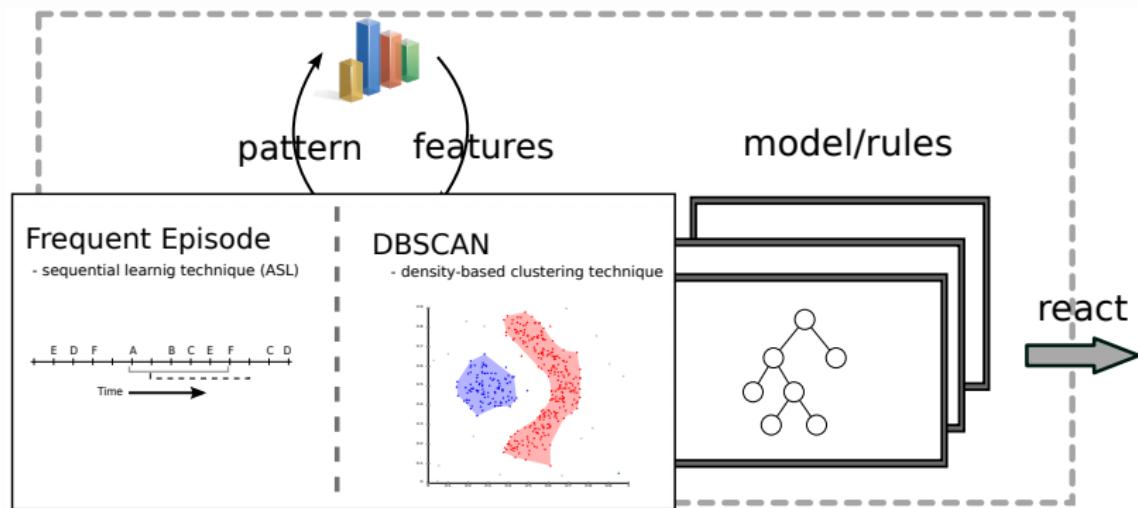
Machine Learning based Anomaly Detection



by Charles River Analytics Inc

- 1 build rules
- 2 detect the violation of the rules (e.g., outlier)
- 3 determine the likelihood of the fault (*in progress*)

Machine Learning based Anomaly Detection



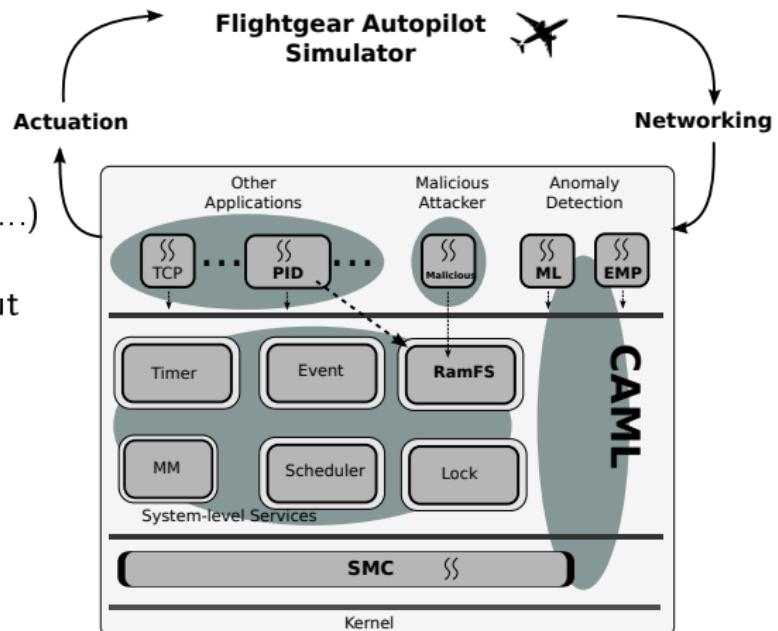
by Charles River Analytics Inc

- 1 build rules
- 2 detect the violation of the rules (e.g., outlier)
- 3 determine the likelihood of the fault (*in progress*)

Evaluation – Flightgear Autopilot Simulator

- Normally

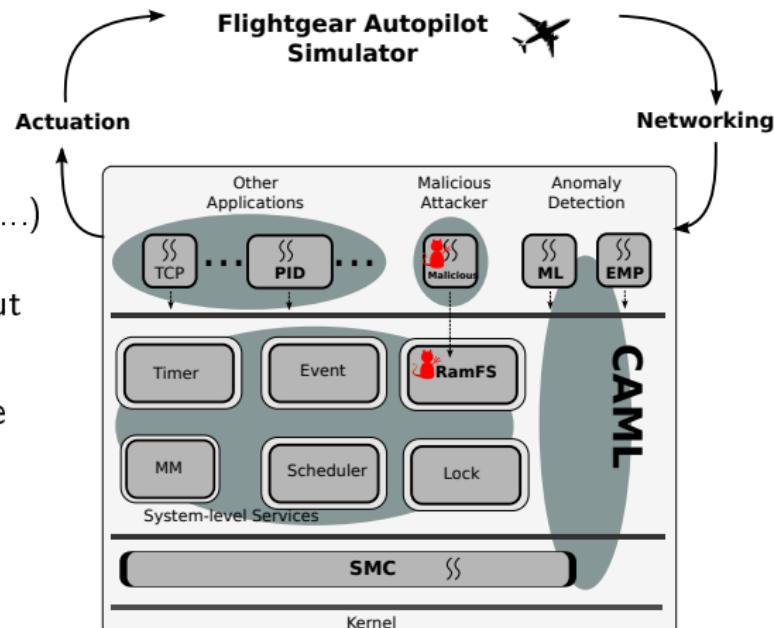
- 1 sensor input (position, velocities, acceleration...)
- 2 PID computation
- 3 corrective action output



Evaluation – Flightgear Autopilot Simulator

- Normally

- 1 sensor input (position, velocities, acceleration...)
- 2 PID computation
- 3 corrective action output



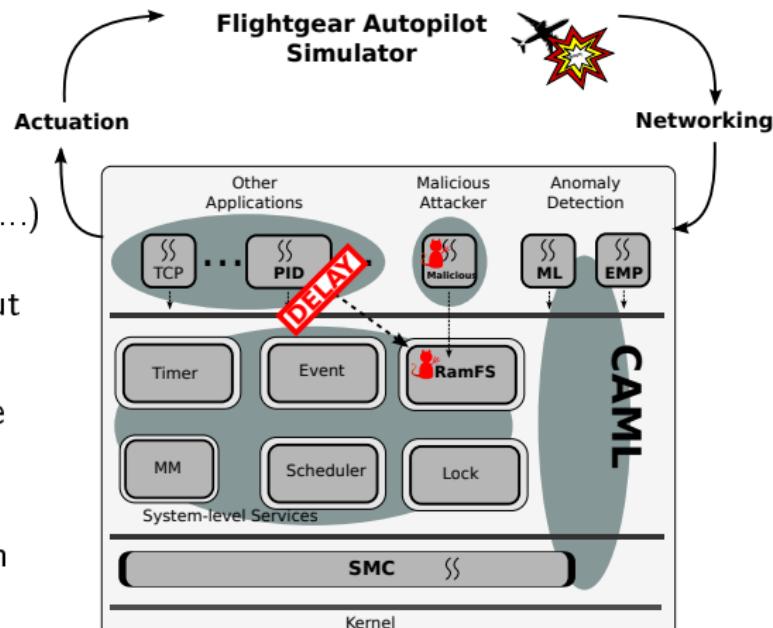
- Malicious user or software bugs

- hog the FS

Evaluation – Flightgear Autopilot Simulator

- Normally

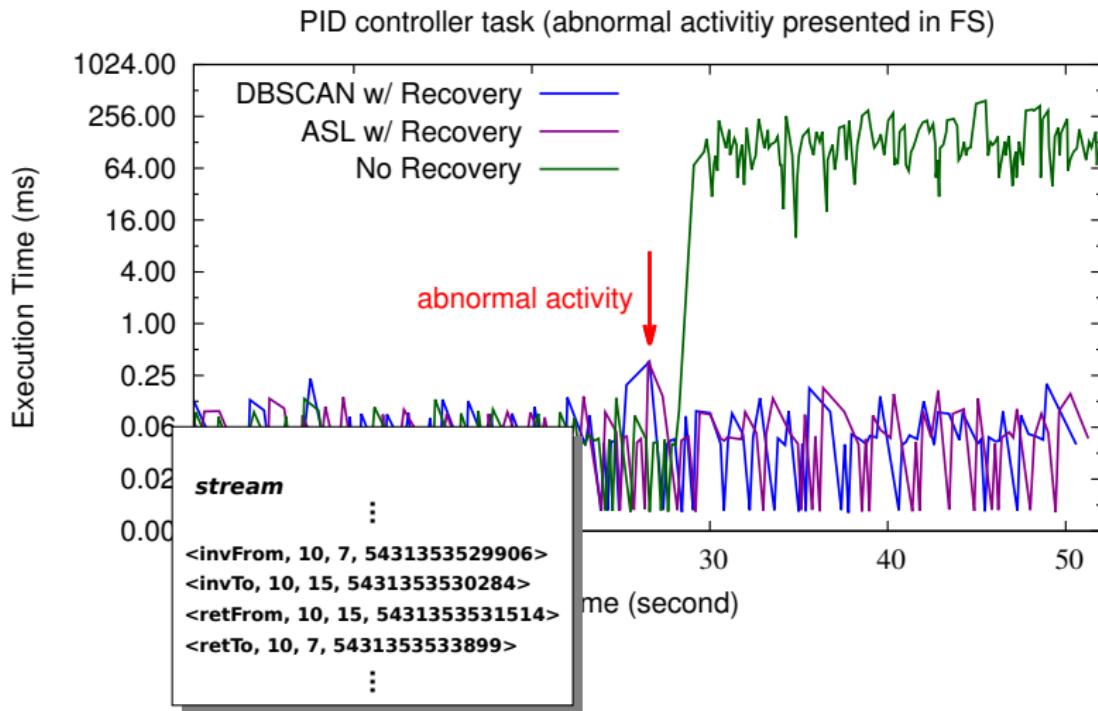
- 1 sensor input (position, velocities, acceleration...)
- 2 PID computation
- 3 corrective action output



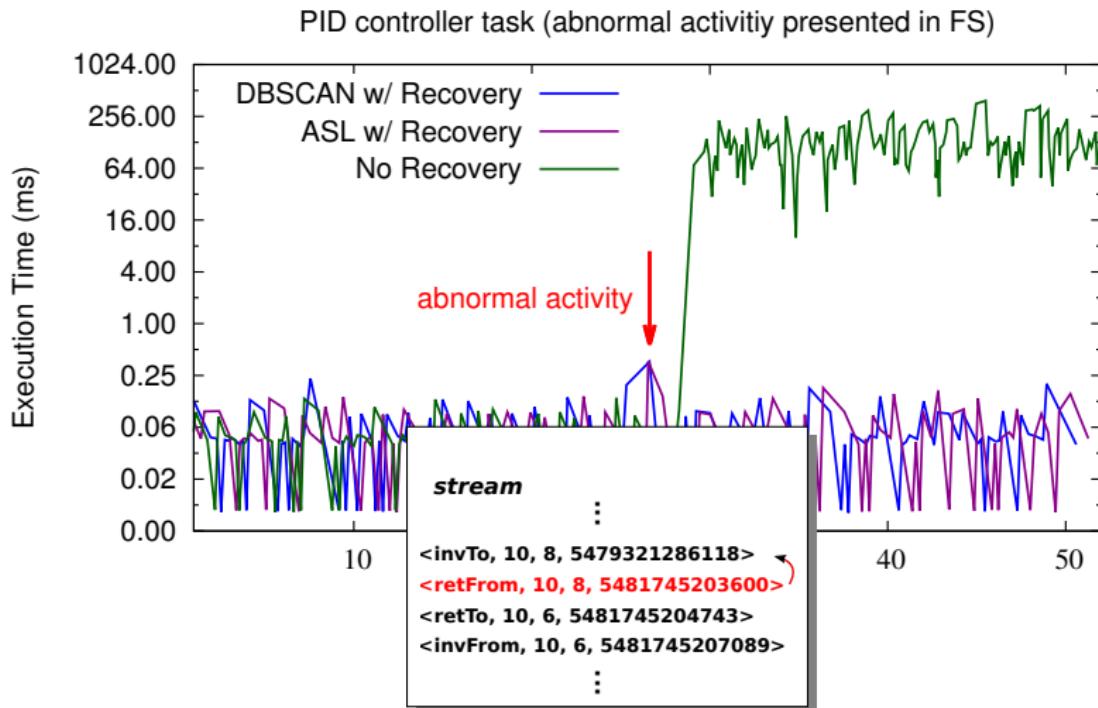
- Malicious user or software bugs

- hog the FS
- delay PID computation

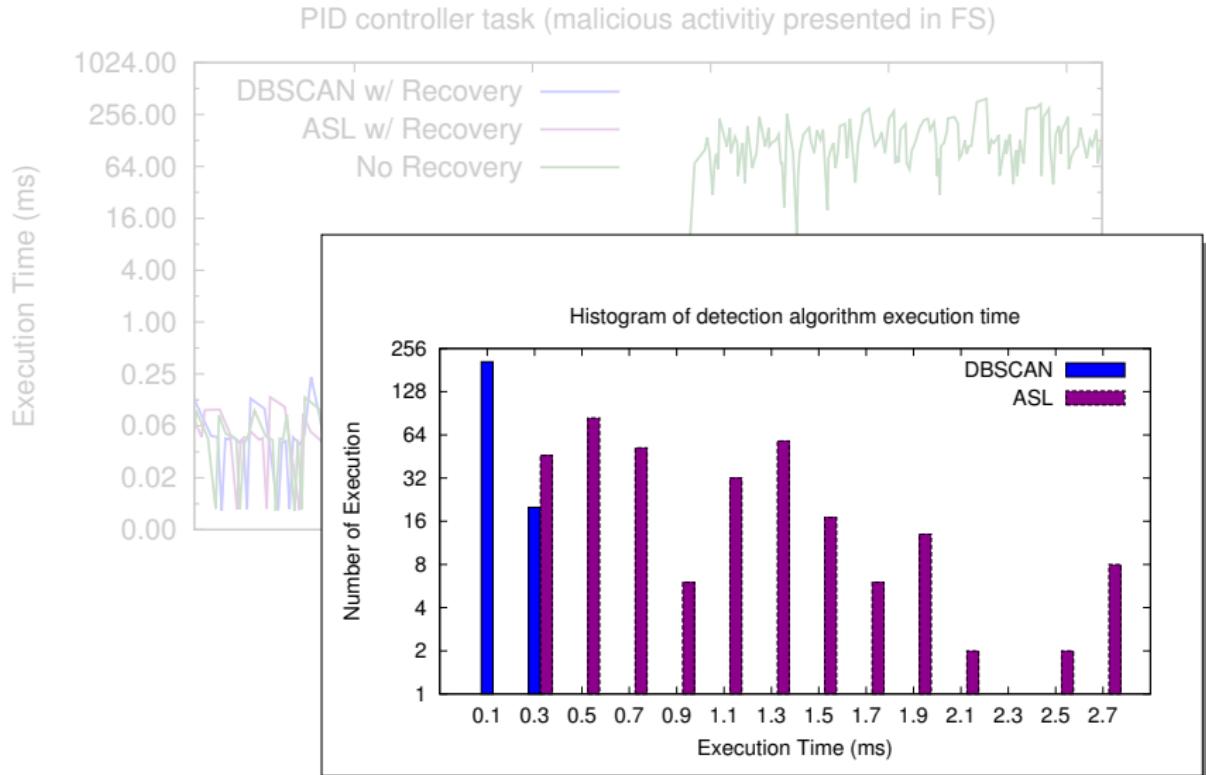
Anomaly Detection Experiment



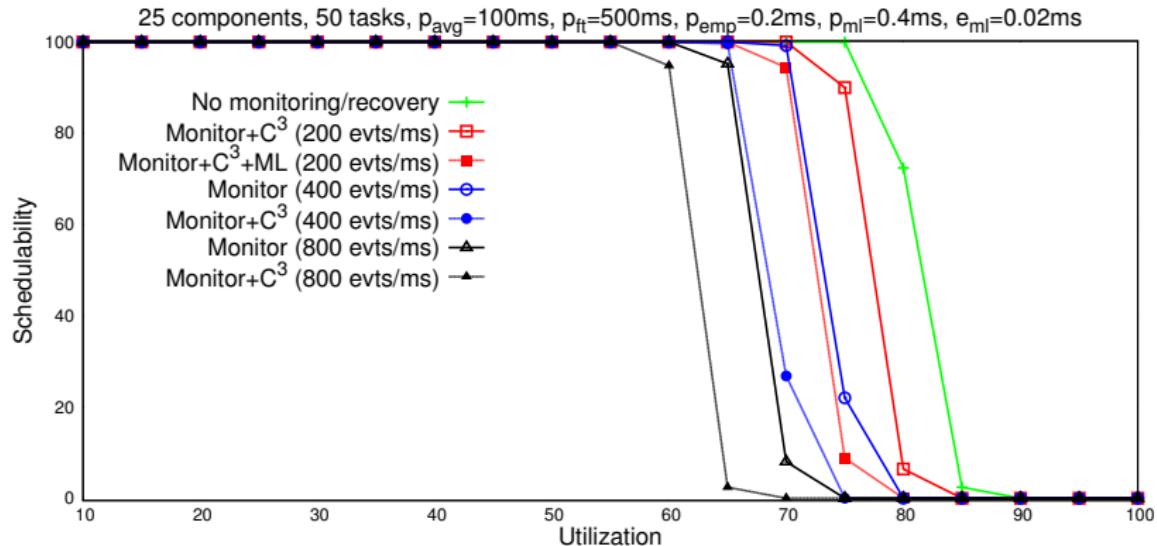
Anomaly Detection Experiment



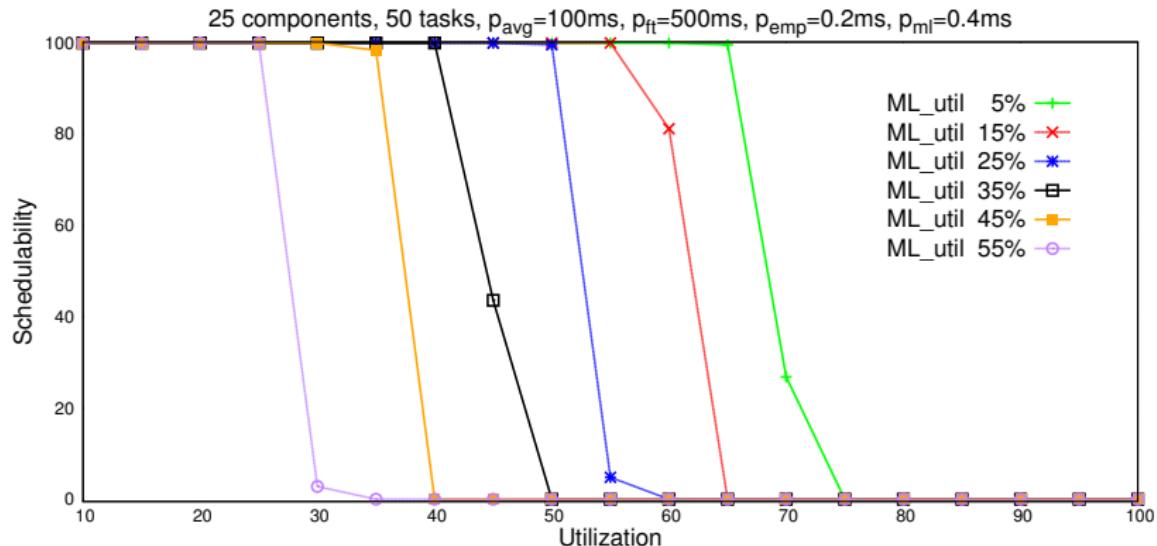
Anomaly Detection Experiment



Response Time Analysis (RTA) Result



RTA Result – Machine Learning Task Utilization



Conclusion on CAML

CAML – an ongoing project that

- aims at detecting **system-level anomaly predictably**
 - *allows the recovery of suspicious component*
- investigates the feasibility using **machine learning** techniques for anomaly detection
 - *collaborate with Charles River Analytics Inc*

Issues that need to be further addressed include

- determine the likelihood of fault occurrence and localize the suspicious service from the rule violations
- study the impact of detection false positives and false negatives on the system schedulability

Publications

- J. Song, G. Bloom and G. Parmer, "SuperGlue: IDL-Based, System-Level Fault Tolerance for Embedded Systems", the 46th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN), 2016
- J. Song, and G. Parmer, "C'Mon: a Predictable Monitoring Infrastructure for System-Level Latent Fault Detection and Recovery", in the 20th RTAS, 2015
- J. Song, J. Wittrock, and G. Parmer, "Predictable, efficient system-level fault tolerance in C³", in the 34th RTSS, 2013
- J. Song, Q. Wang, and G. Parmer, "The state of composite", in the Workshop on Operating Systems Platforms for Embedded Real-Time Applications (OSPERT), 2013
- J. Song and G. Parmer, "Toward predictable, efficient, system-level tolerance of transient faults", in the 5th Workshop on Adaptive and Reconfigurable Embedded Systems (APRES), 2013
- E. Armbrust, J. Song, G. Bloom, and G. Parmer, "On Spatial Isolation for Mixed Criticality, Embedded Systems", in the Workshop on Mixed Criticality, 2014
- Q. Wang, J. Song, G. Parmer, G. Venkataramani, and A. Sweeney, "Increasing memory utilization with transient memory scheduling", in the 33rd RTSS, 2012
- Q. Wang, J. Song, and G. Parmer, "Stack management for hard real-time computation in a component-based OS", in the 32nd RTSS, 2011
- G. Parmer and J. Song, "Customizable and predictable synchronization in a component-based OS", in the International Conference on Embedded Systems and Applications (ESA), 2010

Thanks

? || /* */

composite.seas.gwu.edu

THE GEORGE
WASHINGTON
UNIVERSITY

WASHINGTON, DC