

Control of Computer Systems: From Embedded to the Cloud

Karl-Erik Årzén
Lund University
Lund, Sweden

Cloud
Control

ELLIIT LCCC

etc
Adaptivity & Control of Resources in Embedded Systems

Keynote at CTSE: 1st Workshop on Control Theory for Software Engineering, Bergamo, Aug 31 2015

Content

- Motivation and Background
 - A simple queue length control example
- Resource Management for Multi-Core Embedded Systems
 - The ACTORS Resource Manager
 - Video demo
 - Game-Theory Resource Manager
- The Cloud
 - Problems and Challenges
 - Brownout-inspired resource management for web-service applications

Control of Computer Systems

- Apply control as a technique to manage uncertainty and achieve performance and robustness in computer and communication systems.
- Applications in
 - Internet
 - Servers and data centers, i.e., the cloud
 - Cellular phone systems
 - Embedded systems

Control of Computer Systems

Alternative names:

- Dynamic/adaptive resource management
 - Control as means for managing limited resources
 - Adaptivity from a CS point of view
- Feedback computing/scheduling
- Autonomous/autonomic computing
- Reconfigurable computing

Why?

- System complexity increases

Complexity

6

Why?

- System complexity increases
- Complete information about all use cases and their resource requirements is often not available at design-time
- Green computing → power consumption constraints increasingly important
- Increased hardware density → thermal constraints increasingly important
- Hardware platforms increasingly complex → increasing difficulties in providing good off-line estimates of resource consumption
- Hardware variability increases
- Hardware increasingly often supports adaptivity
- Increased requirements on predictability in the cloud

Control of Computer Systems

- Active research area since around 2000
- However, feedback has been applied in ad hoc ways for long without always understanding that it is control, e.g. TCP/IP
- Control of computing systems can benefit from a lot of the classical control results
 - However, several new challenges

Some Examples

Example 1: A multi-mode embedded system where the resource requirements for all the tasks in all the modes are known at design time

- Use schedulability analysis to ensure that the deadlines are met in all modes and then use a mode-change protocol that ensures that all deadline also are met during the transition between the modes

Example 2: An embedded system with a constant set of hard-RT applications/tasks but where the WCET analysis possible on the selected hardware is too pessimistic and leads to too low resource utilization or where the age- or process-induced variability is too large

- Measure the actual resource consumption and adjust, e.g. the task rates in order ensure that the schedulability condition is fulfilled

Some Examples

Example 3: Open embedded systems where the number of applications and their characteristics change dynamically (e.g. smartphones)

- Measure resource consumption and decide how much resources that should be allocated to each application in order to maximize QoS/QoE while minimizing power consumption and avoiding thermal hotspots

Example 4: A distributed embedded system where one for dependability reasons must be able to ensure system functionality also in case of single-node failures

- Detect node failures and then adapt the task mapping and the schedules so that the system performance is still acceptable

Some Examples

Example 5: An FPGA-based system with multiple modes that is too large to fit in a single FPGA or where the power consumption will be too high

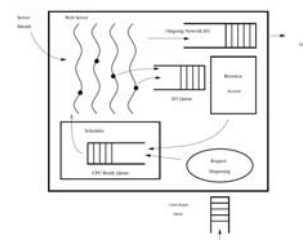
- Use run-time reconfiguration to change the FPGA function dynamically

Example 6: A cloud deployed web-service application where the incoming load varies a lot over time

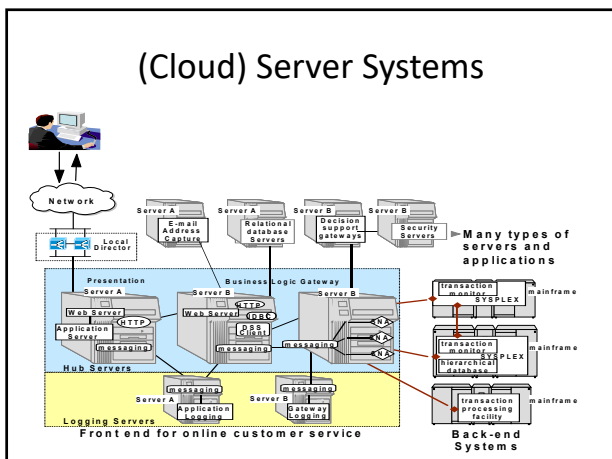
- Dynamically add or remove virtual machines to match the load (elasticity control/auto-scaling)

Computer Internals

- Discrete Event Dynamic System
 - Tasks/requests arrive (queued) and depart (dequeued)
- Execution/service times
- Queuing delays



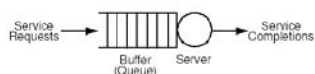
(Cloud) Server Systems



Modeling and Control Formalisms

- Discrete Event Formalisms
 - Automata theory (e.g., Supervisory Control Theory)
 - Petri nets
 - Often problem with scalability
 - Queuing theory
- Continuous-Time Formalisms
 - Liquid ("flow") models + continuous-time control
 - Queues = tanks, computations = flows
 - Average values assuming large number of requests/jobs
 - Sometimes event-driven sampling and control

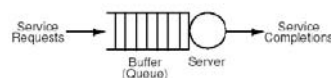
Example: Queuing System



Work requests (customers) arrive and are buffered
 Service level objectives (e.g., response time for request belonging to class X should be less than Y time units)
 Reduce the delay caused by other requests, i.e., adjust the buffer size and redirect or block other requests
 Admission control

Example: Queue Length Control

Assume an M/M/1 queuing system:

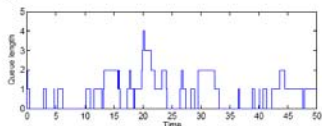


- Random arrivals (requests), Poisson distributed with average λ per second
- Random service times, exponentially distributed with average $1/\mu$
- Queue containing x requests

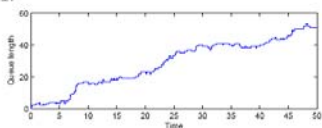
Intuition: $x \rightarrow \infty$ if $\lambda > \mu$

Queue Length Control: Simulation

$\lambda = 0.5, \mu = 1:$



$\lambda = 2.0, \mu = 1:$



Queue Length Control: Model

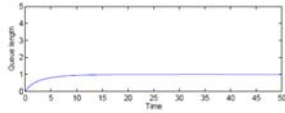
Approximate the system with a nonlinear flow model (Tipper's model from queuing theory)

The expectation of the future queue length x is given by

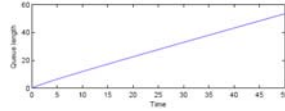
$$\dot{x} = \lambda - \mu \frac{x}{x+1}$$

Queue Length Control: Model

$\lambda = 0.5, \mu = 1:$



$\lambda = 2.0, \mu = 1:$



Queue Length Control: Control Signal

Control the queue length by only admitting a fraction u (between 0 and 1) of the requests

$$\dot{x} = \lambda u - \mu \frac{x}{x+1}$$

Admission control

Queue Length Control: Linearization

Linearize around $x = x^o$

Let $y = x - x^o$

$$\dot{y} = \lambda y - \mu \frac{1}{(x^o + 1)^2} y = \lambda u - \mu \alpha y$$

Queue Length Control: P-control

$$u = K(r - y)$$

$$\dot{y} = \lambda K(r - y) - \mu \alpha y$$

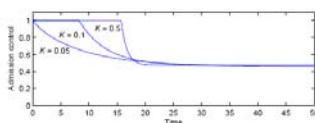
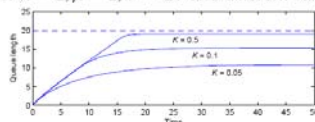
$$(s + \lambda K + \mu \alpha)Y(s) = \lambda K R(s)$$

$$G_{cl}(s) = \frac{\lambda K}{s + \lambda K + \mu \alpha}$$

With K the closed loop pole can be placed arbitrarily

Queue Length Control: P-Control

Simulations for $\lambda = 2, \mu = 1, x^o = 20$ and different values of K



Queue Length Control: PI-control

$$G_P(s) = \frac{\lambda}{s + \mu \alpha}$$

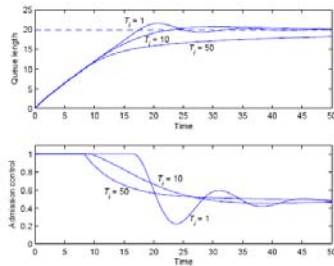
$$G_R(s) = K(1 + \frac{1}{sT_i})$$

$$G_{cl}(s) = \frac{G_P G_R}{1 + G_P G_R} = \frac{\lambda K(s + \frac{1}{T_i})}{s(s + \mu \alpha) + \lambda K(s + \frac{1}{T_i})}$$

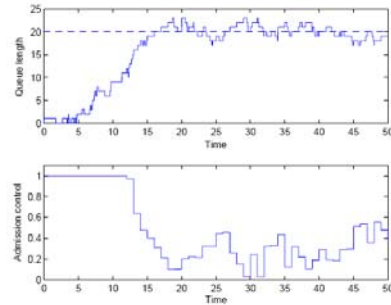
With K and T_i the closed loop poles can be placed arbitrarily

Queue Length Control: PI-control

Simulations for $\lambda = 2, \mu = 1, x^0 = 20, K = 0.1$ and different values of T_i

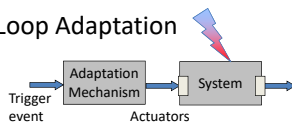


PI-control on Real Queue



Adaptation Mechanisms

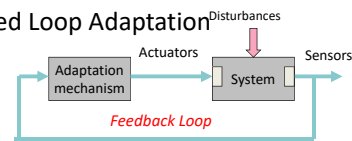
- Open Loop Adaptation



- Feedforward
- Assumes perfect information (model) of the system
- Assumes that there are no external disturbances

Adaptation Mechanisms

- Closed Loop Adaptation



- Feedback
- Adaptation Mechanism == Controller
- Requires sensors
- May cause unstabilities

Adaptation Formulations

- Often formulated as an optimization-problem or as a control-problem
- Optimization Formulations:

$$\begin{aligned} &\text{maximize/minimize } \textit{resource-consumption objective} \\ &\text{s.t. } \textit{performance constraint} \end{aligned}$$

$$\begin{aligned} &\text{maximize/minimize } \textit{performance objective} \\ &\text{s.t. } \textit{resource consumption constraint} \end{aligned}$$

- Performed off-line, online when some change has occurred or periodically, off-line+on-line, ...
- ILP, Bin-packing, MILP, QP, NLP (B&B, GA, CP ...)
- Centralized or distributed

Adaptation Formulations

- Control Formulations:

- System modelled as (linear) dynamic system
- Classical linear control design techniques

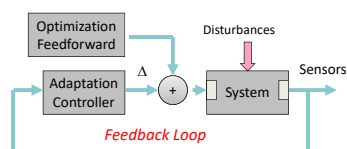
- PID
- LQG
-

$$u(t) = K(e(t) + \frac{1}{T_i} \int e(s) ds + T_D \frac{de(t)}{dt})$$

- Designed to obtain a stable closed loop system with desired dynamic performance

Adaptation Formulations

- Combined Optimization and Control Formulations:
 - Model-Predictive Control (MPC)
 - Optimization problem solved each sample
 - Only the first control signal is used (receding horizon principle)
 - Feedforward + feedback structures



Actuators

- Change the applications / threads
 - For example:
 - Accept or reject decision
 - Change the rates of periodic processes
 - Change between alternative versions (service/quality levels)
 - Anytime formulations
 - Requires support from the applications
- Change the mapping of the application onto the execution platform
 - Priority
 - Schedule
 - Processor allocation

Actuators

- Change the execution platform
 - Number of processors (virtual or physical)
 - DPM techniques
 - Speed of processors
 - DVFS
 - Change the bandwidth of the VM or of the bandwidth server
 - Functionality (hardware-based systems)
 - Micro-code in soft-cores
 - FPGA netlist

Sensors

- What we can (or would we like to) measure?
 - Application performance
 - Obtained QoS
 - Throughput
 - Latency
 - OS / CPU level
 - CPU cycles / task
 - CPU utilization
 - Deadline miss ratio
 - Power and temperature
 - Power consumption for each unit
 - Temperature of each heat source (core, coprocessor, memory controller, ...)

Models

- It is unrealistic to assume sensors for everything
- Must be combined with realistic models that allow us to estimate entities which we cannot measure
- Dynamic calibration using sensor readings (Kalman filter / dynamic observers)
- Power models:
 - Dynamic and static power consumption
- Temperature models
 - Heat transfer between cores
 - Active cooling
 - Multi-tier 3D chips
- Interplay between power and temperature models
 - Temperature dependent leakage power
- Model parameters through system identification

Problems of Feedback

Feedback can introduce new problems:

- The feedback mechanism itself consumes resources
- Harder to provide formal guarantees about the system → not suitable for safety-critical hard real-time application, or?

What about Safety-Critical Systems?

- In many cases control systems
- Due to the feedback errors in the space domain are natural
- Control system designed using
 - Numerous approximations
 - Model reduction, linearization,
 - Verified through extensive simulations
 - Large safety margins when selecting, e.g., sampling periods
- **Why is it then so unthinkable to use feedback also at the implementation level?**



Problems of Feedback

Feedback can introduce new problems:

- The feedback mechanism itself consumes resources
- Harder to provide formal guarantees about the system → not suitable for safety-critical hard real-time application, or?
- Adds to the complexity
- May complicate the design process (modeling, V&V, ...)
- Requires tuning
- Sensors and actuators are necessary
- Models are necessary
 - Of the system
 - Of the feedback mechanism itself
- Feedback may cause instability

Golomb on Modeling

- "Mathematical Models: Uses and Limitations" – Simulation, Apr 70
- Don't apply a model until you understand the simplifying assumptions on which it is based and can test their applicability.
- Distinguish at all times between the model and the real world. **You will never strike oil by drilling through the map!**
- The purpose of notation and terminology should be to enhance insight and facilitate computation – not to impress or confuse the uninitiated



Solomon Wolf Golomb
(1932) mathematician and engineer and a professor of electrical engineering at the University of Southern California.

Problems of Feedback

Feedback can introduce new problems:

- The feedback mechanism itself consumes resources
- Harder to provide formal guarantees about the system → not suitable for safety-critical hard real-time application, or?
- Adds to the complexity
- May complicate the design process (modeling, V&V, ...)
- Requires tuning
- Sensors and actuators are necessary
- Models are necessary
 - Of the system
 - Of the feedback mechanism itself
- Feedback may cause instability
- Feedback may introduce measurement noise
 - Only when you measure physical entities!

Content

- Motivation and Background
 - A simple queue length control example
- **Resource Management for Multi-Core Embedded Systems**
 - **The ACTORS Resource Manager**
 - Video demo
 - Game-Theory Resource Manager
- The Cloud
 - Problems and Challenges
 - Brownout-inspired resource management for web-service applications

ACTORS

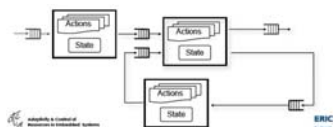


- Adaptivity and Control of Resources in Embedded Systems
- EU FP7 STREP project
 - 2008-2011
 - Coordinated by Ericsson (Johan Eker)
 - Lund University, TU Kaiserslautern, Scuola Superiore Sant'Anna di Pisa, EPFL, AKAtch, Evidence
- Media applications (soft real-time) for smart phones
- Control applications



ACTORS: Key Ingredients

1. Data-Flow Programming
 - CAL Actor Language



2. Adaptive Resource Management of service-level aware applications
 - Soft real-time media applications
 - Control applications

Service Level-Aware Applications

- Application knob
 - Decides the QoS achieved and the amount resources required
 - High SL → high QoS & high resource usage
 - Low SL → low QoS & low resource usage
- Discrete
 - "application modes"
- Continuous
 - e.g., sampling rate in a controller

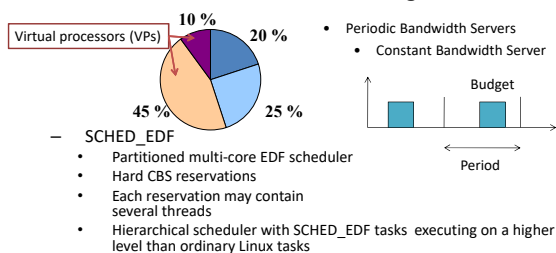
Service level example



- SL: Resolution and/or frame rate of a video stream
- QoS and required CPU for encoding and decoding depends on the SL

ACTORS: Key Ingredients

3. Reservation-Based CPU Scheduling

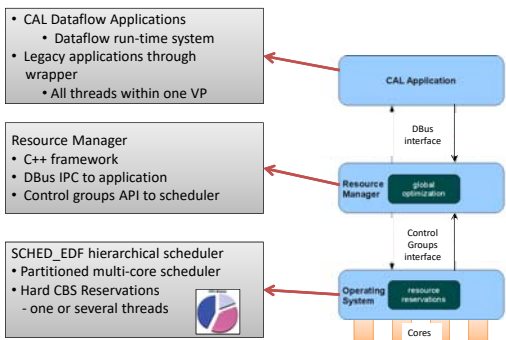


- SCHED_EDF
 - Partitioned multi-core EDF scheduler
 - Hard CBS reservations
 - Each reservation may contain several threads
 - Hierarchical scheduler with SCHED_EDF tasks executing on a higher level than ordinary Linux tasks

4. Multicore Linux Platforms

- ARM 11, x86

Overview



Static Information

From applications to RM at registration:

- Service Level Table

Service Level	QoS	BW Requirement	BW distribution	Timing Granularity
0	100	240	60-60-60-60	20 ms
1	75	180	45-45-45-45	20 ms
2	40	120	30-30-30-30	20 ms

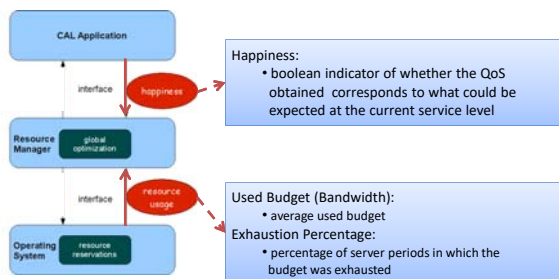
- Thread IDs and how they should be grouped

From system administrator to RM at startup:

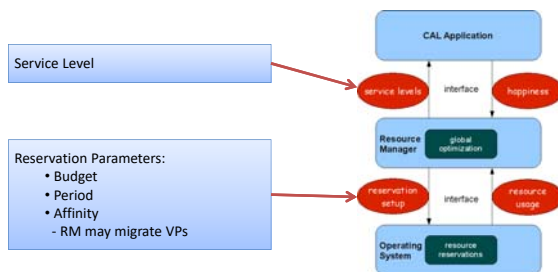
- Application Importance

Appl.	Importance
Appl 1	10
Appl 2	20
Appl 3	100
Default	10

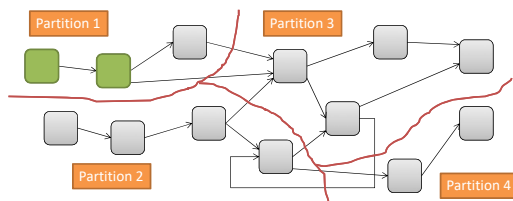
Dynamic Inputs



Outputs

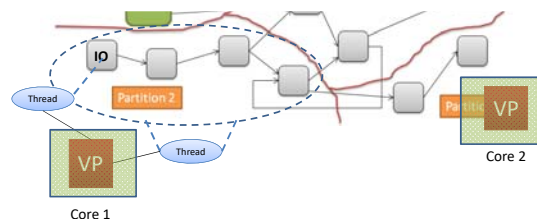


Dataflow Analysis



- Static partitioning
- Automatic analysis for SDF/CSDF actors
- Automatic merging of SDF/CSDF actors to improve run-time performance

CAL Run-Time System

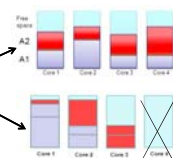


- One thread per partition executing its actors using round-robin
- One thread per "system actor" (IO, time)
- The threads from the same partition are executed by the same virtual processor
- If possible the VPs are mapped to different physical cores in order to enable parallel execution

Resource Manager Functionality

- Assign service levels
 - When applications register or unregister
 - Formulated as a ILP problem
 - Importance as weight
 - glpk solver
- Mapping & bandwidth distribution
 - Map reservations to cores
 - Distribute the total BW to the reservations
 - Two Approaches:
 - Spread out the VPs and balance the load
 - Pack the VPs in as few cores as possible
 - Allow turning off unused cores
 - Bin packing

$$\begin{aligned} \max \sum_{i=1}^n w_i q_i x_i \\ \sum_{i=1}^n \alpha_i x_i \leq C \\ \forall i, \sum x_i = 1 \end{aligned}$$



- At most 90% of each core is used for SCHED_EDF tasks

Resource Manager Functionality

- Separate service level assignment and BW Distribution
 - The best service level assignment may lead to an unfeasible BW distribution
 - Approach 1:
 - New SL assignment that generates the next best solution
 - New BW Distribution
 - Approach 2:
 - Compress the individual VPs

Resource Manager Tasks

- Bandwidth adaptation
 - Adjust the servers bandwidth dynamically based on measured resource usage and obtained happiness

Changes what is meant by sufficiently close based on EP:

Changes the AB so that the UB lies sufficiently close:

Resource Manager Tasks

- Multiple bandwidth adaptation strategies
- Strategy 1:
 - A VP may never consume more bandwidth than what was originally assigned to it
 - BW controller may reduce the BW if not used
- Strategy 2:
 - A VP may use more resources than originally assigned to it
 - If there are free resources available, or
 - If there are VPs of less important applications that use more BW than originally allocated to them
 - In the latter case the less important applications are compressed
 - All applications are always guaranteed to obtain their originally assigned values (can never be compressed beyond that)

RM Support Software

- GUI
 - VP to core assignment
 - AB, UB, and EP
 - Service Level Table
 - Event history
 - Itself an application under the control of the RM
- Load Generator
 - Generates artificial load for testing
- Application Wrapper
 - Wrapper for non-Actors aware applications

MPEG-4 Video Decoding Example

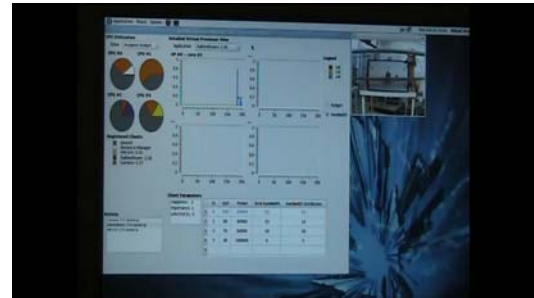
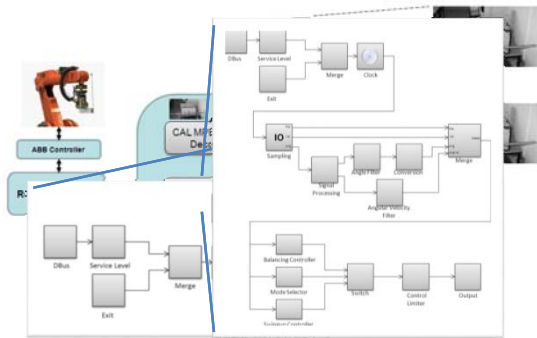
- MPEG-4 SP decoder implemented in CAL
- Connected to an Axis network camera
- Service level changes results in commands from the decoder to the camera to reduce the frame rate and/or resolution

MPEG-4 Video Decoding Demonstrator

Feedback Control Demonstrator

- Industrial robot balancing an inverted pendulum
 - Pendulum controller in CAL
 - Service level changes correspond to changes in sampling period
- Ball and Beam Processes
 - Controller in CAL
 - Service level changes correspond to changes in sampling period
- MPEG-4 Video Decoder
 - Service level changes correspond to changes in resolution/frame rate

Feedback Control Demonstrator



Drawbacks

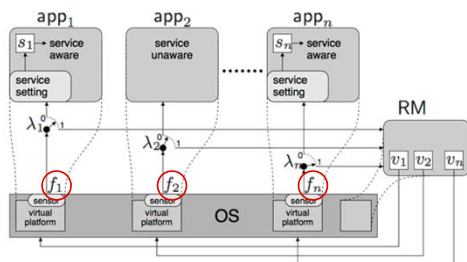
- ILP does not scale well
 - Requires a lot of information from the applications
 - More natural if the applications select their service levels and the resource manager adjust the vp size
- The game-theory inspired approach

Content

- Motivation and Background
 - A simple queue length control example
- Resource Management for Multi-Core Embedded Systems
 - The ACTORS Resource Manager
 - Video demo
 - Game-Theory Resource Manager
- The Cloud
 - Problems and Challenges
 - Brownout-inspired resource management for web-service applications

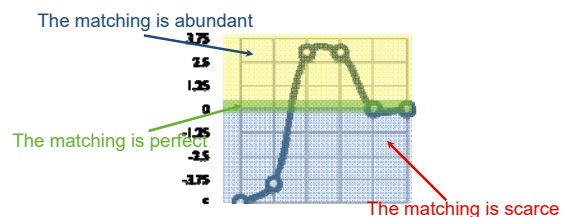
Towards decentralization

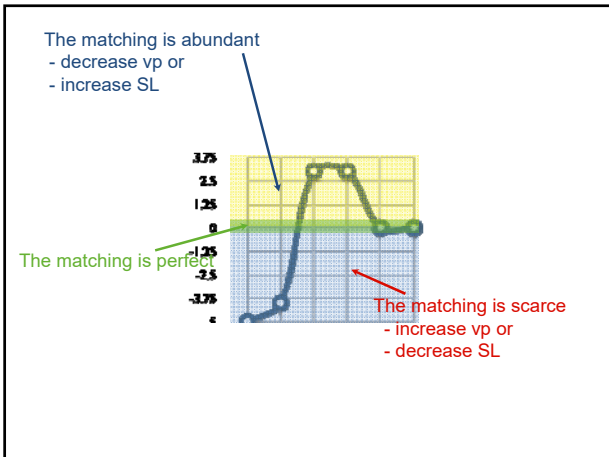
- The resource manager allocates resources and applications choose their service levels based on current performance



Matching function

Measures how well the resources assigned to the application (vp size) matches the resource requirements (SL) of the applications





Application weights

- $\lambda_i \in [0..1]$
- $\lambda_i = 0$ means that only the application should adjust
- $\lambda_i = 1$ means that only the resource manager should adjust

Matching Function

- Application executes a series of jobs
 - Execution time $C_i = a_i s_i$
 - Response time for each job $R_i = \frac{C_i}{v_i} = \frac{a_i s_i}{v_i}$
 - Want this to be equal to the deadline D_i
 - Hence

$$f_i = \frac{D_i}{R_i} - 1 = \frac{D_i v_i}{a_i s_i} - 1 = \beta_i \frac{v_i}{s_i} - 1$$
 - Can be measured from job start and stop times
 - Depends on service level and virtual processor speed

Resource Manager

At each step:

- Measure f_i for all the applications
 - The applications report the start and stop of each job by writing to shared memory
 - RM reads from shared memory and calculates the response time
- Updates the virtual processors:

$$v_i(k+1) = v_i + \varepsilon_{RM}(k) \left(-\lambda_i f_i(k) + \sum_{j=1}^n \lambda_j f_j(k) v_i(k) \right)$$

$$\varepsilon_{RM}(k) = \frac{1}{k+1}$$

Service Level Adjustment

- Should set s_i so that f_i becomes 0
- Naive approach: $s_i(k+1) = \beta_i v_i(k+1)$
 - Assumes knowledge of β_i
- Instead estimate β_i

$$\beta_i = (1 + f_i(k)) \frac{s_i(k)}{v_i(k)}$$
- Gives

$$s_i(k+1) = (1 + f_i(k)) \frac{v_i(k+1)}{v_i(k)} s_i(k)$$
 - Continuous service levels
 - Robustified

Game Theory

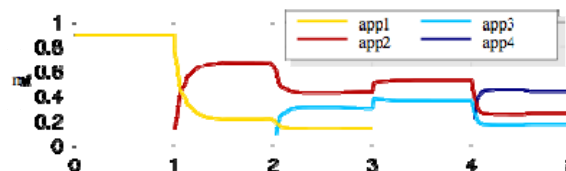
- Provides convergence results and results about stationary values
- For example
 - A stationary point satisfies the following condition
 - The matching function is zero
 - OR
 - The service level is the smallest possible AND the matching function is negative

Implementation

- Using SCHED_DEADLINE
 - New EDF+CBS scheduling policy in Linux
 - Developed by Evidence within ACTORS

Evaluation

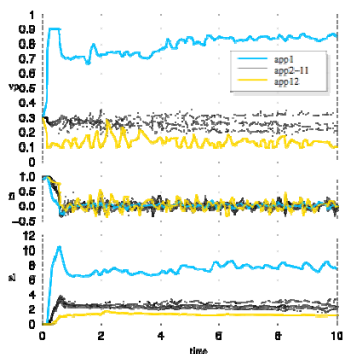
- Convergence of virtual processors
- Four applications
 $\lambda_1 = 0.1, \lambda_2 = 0.3, \lambda_3 = 0.2$ and $\lambda_4 = 0.5$



- Resources assigned proportionally to the normalized weights

Evaluation

- Multicore
- 12 identical apps
- Different weights
 - 1 – high
 - 12 – low
 - 2-11 – medium
- Resources proportional to the service levels



Content

- Motivation and Background
 - A simple queue length control example
- Resource Management for Multi-Core Embedded Systems
 - The ACTORS Resource Manager
 - Video demo
 - Game-Theory Resource Manager
- The Cloud
 - Problems and Challenges
 - Brownout-inspired resource management for web-service applications



- As soon as you use any networked computing unit (laptop, smart phone, sensor device, ...) the likelihood that the computations will be performed in a data center somewhere in the cloud, rather than locally, is very large
- News, mail, photos, tickets, books, clothes, maps, social media, television, music, banking, administrative systems, technical software,

The Datacenter



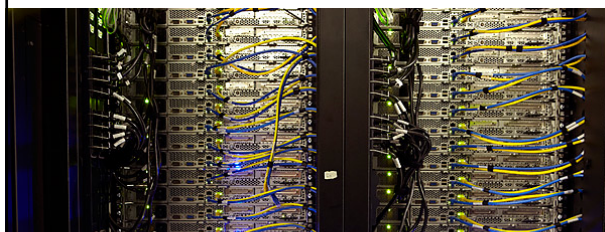
What's inside?



Rack or blade-mounted servers

A modern Amazon data center consists of 50-80.000 physical servers
Each server typically has 4-8 cores

What's inside?



Networking

What's inside?



Power supplies

What's inside?



Cooling

What runs in the cloud?

- Everything that can execute on physical hardware can in principle do so on virtual hardware in the cloud
 - But constraints on latency, throughput, IO, ...
- Two major classes:
 - Web service applications:
 - E-Commerce
 - Massively parallel applications
 - Map/Reduce programming model (HADOOP), MPI
 - Google, Facebook, Twitter,



Problems with the Cloud

- Low level of determinism and predictability
 - No hard performance guarantees

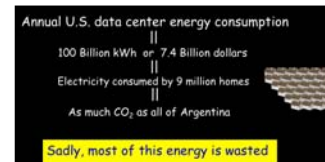
Service Cor

AWS will use commercially reasonable efforts to make Amazon EC2 and Amazon EBS each available with a Monthly Uptime Percentage (defined below) of at least 99.95%, in each case during any monthly billing cycle (the "Service Commitment"). In the event Amazon EC2 or Amazon EBS does not meet the Service Commitment, you will be eligible to receive a Service Credit as described below.

Problems with the Cloud

- Low level of determinism and predictability
 - No hard performance guarantees
- Data centers consume a lot of energy

Cloud Energy Consumption

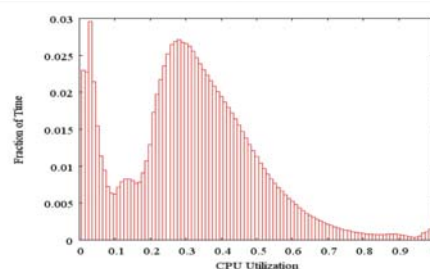


Facebook in Luleå, Sweden will consume 1 TWh/year – approx 40,000 family houses

Problems with the Cloud

- Low level of determinism and predictability
 - No hard performance guarantees
- Data centers consume a lot of energy
- Data centers have low utilization

Datacenter Utilization



5000 google servers over 6 months
Source: "The Datacenter as a computer", Barroso et al

Problems with the Cloud

- Low level of determinism and predictability
 - No hard performance guarantees
- Data centers consume a lot of energy
- Data centers have low utilization

➔ Room for better resource management techniques

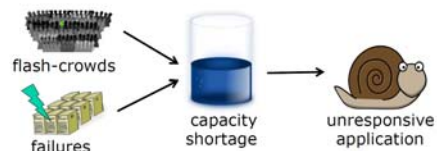
Content

- Motivation and Background
 - A simple queue length control example
- Resource Management for Multi-Core Embedded Systems
 - The ACTORS Resource Manager
 - Video demo
 - Game-Theory Resource Manager
- The Cloud
 - Problems and Challenges
 - **Brownout-inspired resource management for web-service applications**

From Embedded to the Cloud

- Apply the game-theoretic resource manager to cloud applications
- Focus on state-free, request-based applications

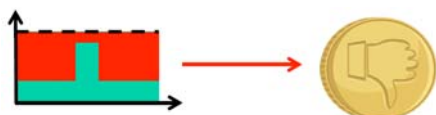
Problem: Unexpected Events



- 25% of end-users leave if load time > 4s*
- 1% reduced sale per 100ms load time*
- 20% reduced income if 0.5s longer load time**

* Amazon ** Google

Standard Practice

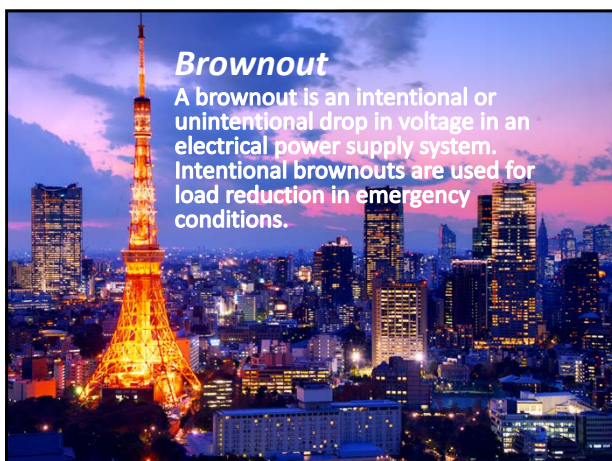


- Overprovisioning
 - Economically impractical

Brownouts



- We borrow from the concept of **brownouts** in power grids
- A brownout-compliant application can degrade user performance when needed to face unexpected conditions



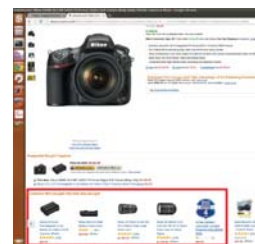
Brownouts



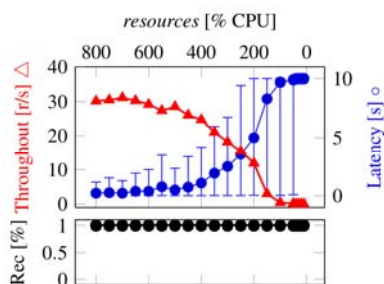
- We borrow from the concept of **brownouts** in power grids
- A brownout-compliant application can degrade user performance when needed to face unexpected conditions
- We assume that the reply to a request consists of a mandatory and an optional part
- During overload and/or lack of resources the percentage of requests that also receives the optional part can be decreased

Brownout Examples

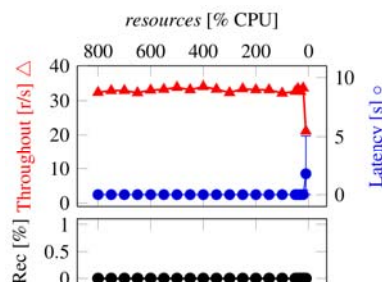
- E-commerce systems with recommendations
- Content adaptation in web server applications
-



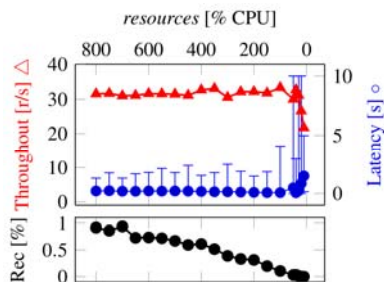
Always Service Optional Part



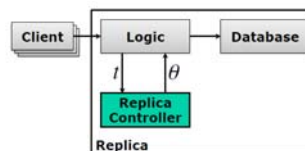
Never Service Optional Part



With Brownout



Brownout Control Loop

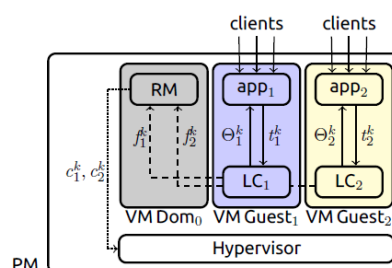


- Measured variable: Response time (average or 95% percentile)
- Control signal: Dimmer value (percentage of requests for which also the optional parts are calculated)
- Setpoint: 2 seconds

Brownout Controllers

- Adaptive PI and/or PID controllers
- Adaptive deadbeat controller
- Feedforward + feedback controller
 - Assumes knowledge of average service time for mandatory and optional parts

Brownout Resource Management



Resource Management Details

- Applications sends value of **matching function**

$$f_i^k = 1 - t_i^k / \bar{t}_i$$

- Resource manager computes size of virtual machine

$$c_i^{k+1} = c_i^k - \epsilon_{rm} \left(f_i^k - c_i^k \cdot \sum_p f_p^k \right)$$

- Proven to be fair and converge using game theory

Evaluation

- Applications
 - RUBiS: eBay-like prototype auction website
 - Added a recommender
 - RUBBoS: Slashdot-like bulletin board website
 - Added a recommender
 - Marked comments as optional
- Effort in lines of code:

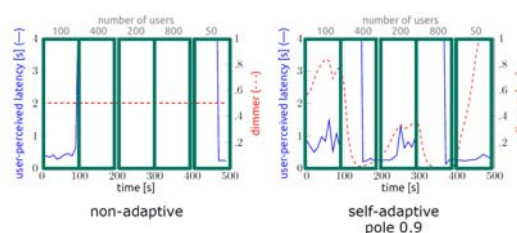
Modification	RUBiS	RUBBoS
Recommender	37	22
Dimmer	3	6
Reporting response time to controller	5	5
Controller	120	120
<i>Total</i>	<i>165</i>	<i>153</i>

Evaluation

- Hardware
 - One physical machine
 - Two AMD Opteron 6272 processors
 - 16 cores / processor
- Hypervisor
 - Xen
 - Each VM contains Apache web server, PHP interpreter, MySQL server, and brownout controller
- Client load generator
 - Custom built – httpmon
 - Closed loop model
 - Clients with think time

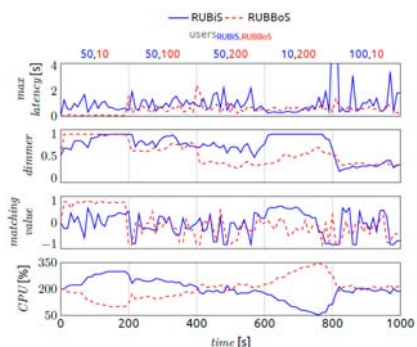
Experimental Results

- RUBiS flash crowds



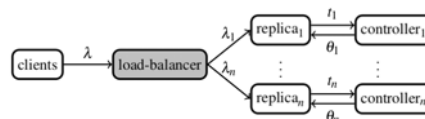
Experimental Results

Four cores used



Load Balancing

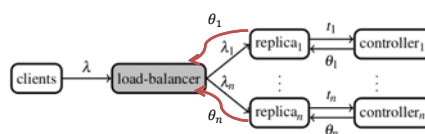
- Multiple replicas
 - Support scaling
 - Robust towards faults
- Load balancing



Existing Load Balancers

- **Dynamic** load balancers often measure response times and decide based on that
 - FRF, FRF-EWMA, 2RC, Predictive ...
- Does not work well in the presence of brownout control
 - the replica controllers keep the response times close to the setpoint

Brownout-Aware Load Balancing



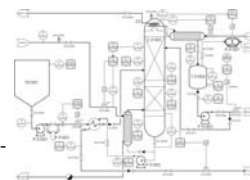
- New load balancing schemes have been developed that make use of the dimmer values
- Piggyback the dimmer values to the replies going back to the load balancer

Challenges

- Modeling of resources and of workload
 - Very difficult to predict peaks
- Scalability
 - What work for 500 servers will most likely not work for 50 000 servers
- Everything will break down all the time
 - Netflix chaos monkey
- Difficult to isolate the phenomenon under study
- Only industry have real data centers
 - Very industry-driven area
 - Google, Microsoft,

Process Automation Comparison

- Compare with process industry in the 1940-50s
 - Manual control
- Next step PI(D) +
 - Feedforward, cascade, ratio selector logic, split-range, mid-range,
 - Controller patterns
 - Flow control, temperature control,
- Now
 - MPC control + optimization-based long-term planning
 - But still PI(D) at the lowest level

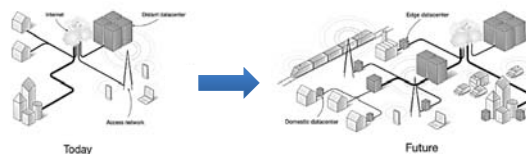


Process Automation Comparison

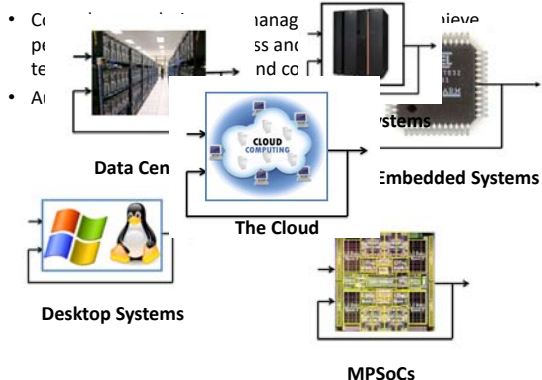
- Lessons:
 - Keep it simple ("KISS")
 - Start with classical, e.g. PID, control techniques
 - Add optimization control on top of this if required
- Differences:
 - Flexibility and generality
 - Behavior not governed by laws of nature
- Maybe one have to constrain the flexibility in order to achieve determinism

The Future Edge/Fog Cloud

- 5G opens up for new application classes
 - Mission-critical applications, e.g., closed loop control
 - Requirements on low latency and high availability → predictability
- Network and cloud convergence
 - Boundary between the network and datacenters disappears
 - Telco cloud, mobile cloud, distributed cloud, edge cloud, fog cloud,
 - Computations to be dynamically deployed in all types of nodes, incl base stations and remote data centers



Feedback Computing



Further Information

- ACTORS Resource Manager
 - E. Bini, G. Buttazzo, J. Eker, S. Schorr, R. Guerra, G. Fohler, K-E. Årzén, V. Romero, C. Scordino: "Resource Management on Multi-core Systems: the ACTORS approach", *IEEE Micro*, 31:3, pp. 72-81, May 2011
- Game-Theoretical Resource Manager
 - Georgios Chasparis, Martina Maggio, Karl-Erik Årzén, Enrico Bini: "Distributed Management of CPU Resources for Time-Sensitive Applications". In 2013 *American Control Conference*, Washington DC, USA, June 2013
 - Martina Maggio, Enrico Bini, Georgios Chasparis, Karl-Erik Årzén: "A Game-Theoretic Resource Manager for RT Applications". In 25th *Euromicro Conference on Real-Time Systems, ECRTS13*, Paris, France, July 2013
 - Georgios Chasparis, Martina Maggio, Enrico Bini, Karl-Erik Årzén: "Design and Implementation of Distributed Resource Management for Time Sensitive Applications", Accepted for publication in *Automatica*

Further Information

Cloud Brownout

- Cristian Klein, Martina Maggio, Karl-Erik Årzén, Francisco Hernández-Rodríguez: "Introducing Service-level Awareness in the Cloud". In 2013 *ACM Symposium on Cloud Computing*, Santa Clara, CA, October 2013
- Cristian Klein, Martina Maggio, Karl-Erik Årzén, Francisco Hernández-Rodríguez: "Brownout: Building more Robust Cloud Applications", 36th *International Conference on Software Engineering (ICSE)*, Hyderabad, India, 2014
- Martina Maggio, Cristian Klein, Karl-Erik Årzén "Control strategies for predictable brownouts in cloud computing", *IFAC World Congress*, Cape Town, South Africa, August 2014
- Jonas Dürango, Manfred Dellkrantz, Martina Maggio, Cristian Klein, Alessandro Vittorio Papadopoulos, Francisco Hernández-Rodríguez, Erik Elmroth, Karl-Erik Årzén, "Control-theoretical load-balancing for cloud applications with brownout", *CDC* 2014
- Cristian Klein, Alessandro Vittorio Papadopoulos, Manfred Dellkrantz, Jonas Dürango, Martina Maggio, Karl-Erik Årzén, Francisco Hernández-Rodríguez, Erik Elmroth, "Improving Cloud Service Resilience using Brownout-Aware Load-Balancing", In 33rd *IEEE Symposium on Reliable Distributed Systems (SRDS)*, 2014

Contributors

ACTORS:

- Johan Eker (Ericsson), Giorgio Buttazzo (SSSA), Enrico Bini (SSSA), Claudio Scordino (Evidence), Gerhard Fohler (TUKL), Stefan Schorr (TUKL), Vanessa Romero Segovia (LU),

Game-Theoretical RM:

- Martina Maggio (LU), Enrico Bini (LU/SSSA), Georgios Chasparis (LU/Software Competence Center Hagenberg)

Cloud Brownout:

- Martina Maggio, Alessandro Papadopoulos, Manfred Dellkrantz, Jonas Dürango (LU), Cristian Klein (Umeå Univ), ...

Questions?