

# Deadline scheduling in the Linux kernel

Juri Lelli, Claudio Scordino, Luca Abeni and Dario Faggioli

---

Benno Fünfstück

July 8, 2019

# Introduction

---

## Story: A new coffee machine



### Components:

- Brewing controller
- User interface controller
- Web interface

First approach: run Linux as a task in real-time hypervisor

- examples: RTAI, RTLinux, Xenomai
- maintenance of HAL and microkernel for real-time
- custom tools and API necessary for real-time part

First approach: run Linux as a task in real-time hypervisor

- examples: RTAI, RTLinux, Xenomai
- maintenance of HAL and microkernel for real-time
- custom tools and API necessary for real-time part

Second approach: make the Linux kernel itself suitable for real-time

- PREEMPT\_RT patchset

First approach: run Linux as a task in real-time hypervisor

- examples: RTAI, RTLinux, Xenomai
- maintenance of HAL and microkernel for real-time
- custom tools and API necessary for real-time part

Second approach: make the Linux kernel itself suitable for real-time

- PREEMPT\_RT patchset
- **real-time scheduler**

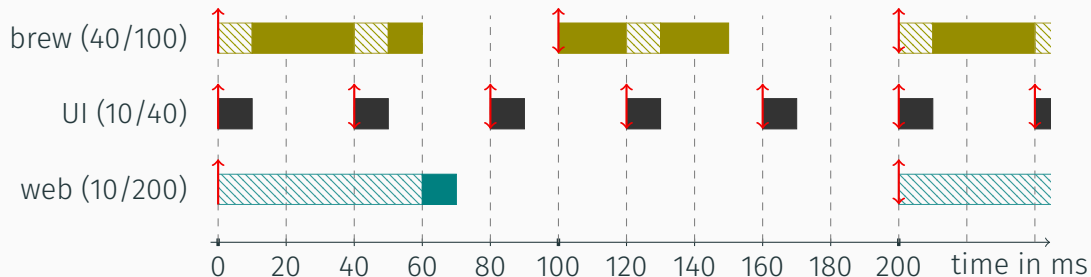
## Design of a realtime scheduler

---

# Real-time Tasks

## Properties of real-time tasks: *runtime*, *deadline* and *period*

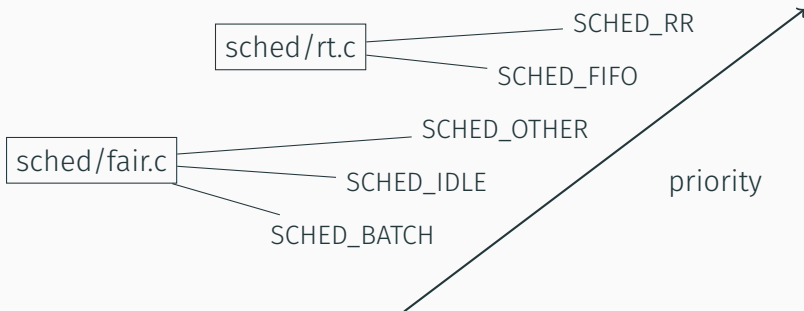
For our coffee machine:





## Related Work

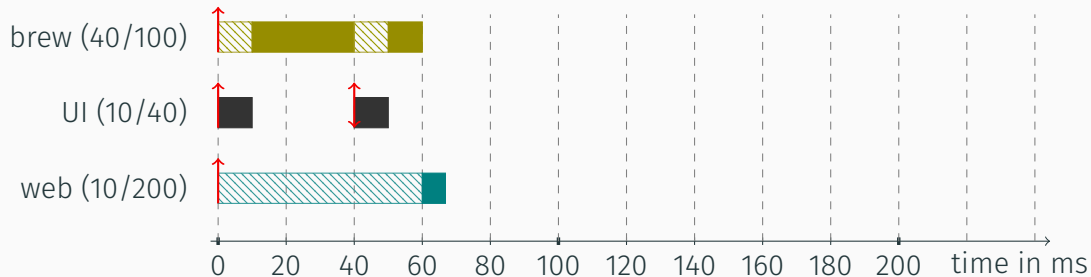
- Linux/RK: resource reservation, CPU (fixed priority) and disk (EDF)
- SCHED\_SPORADIC: priority based, aimed for inclusion in mainline but failed
- OCERA: resource-reservation based scheduler, as loadable kernel module
- LITMUS<sup>RT</sup>: real-time scheduling testbed, not aiming to be production quality
- ExShed: kernel extension to allow scheduler implementation in user space
- Linux modular scheduling framework



# Real-time Tasks

Properties of real-time tasks: *runtime*, *deadline* and *period*

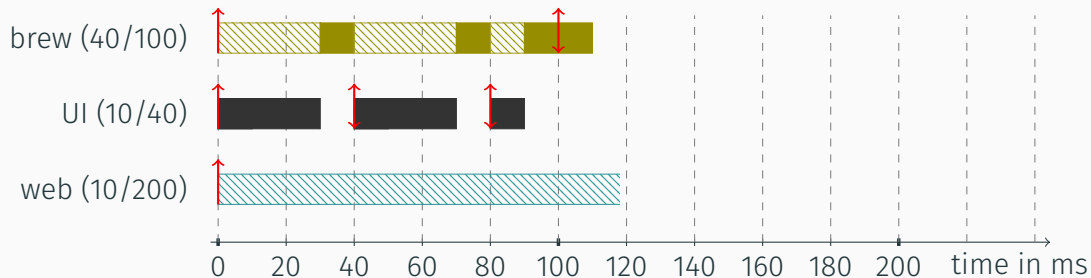
For our coffee machine:



# Real-time Tasks

Properties of real-time tasks: *runtime*, *deadline* and *period*

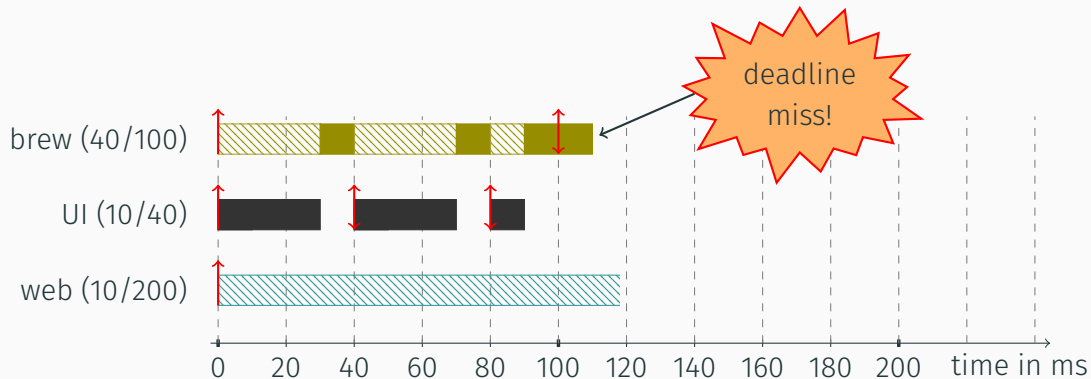
For our coffee machine:



# Real-time Tasks

Properties of real-time tasks: *runtime*, *deadline* and *period*

For our coffee machine:



# Algorithm

For each Task ( $Q_i / T_i$ ), keep track of:

- budget (remaining runtime)  $q_i$  and
- scheduling deadline  $d_i$

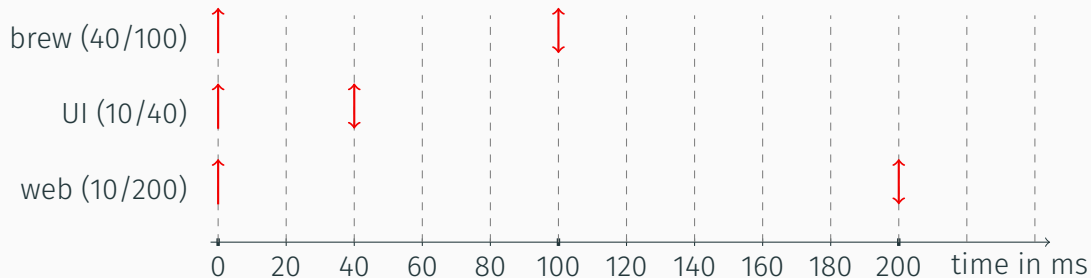
When a task wakes up:

- set new deadline  $d_i = t + T_i$  and recharge budget if  $q_i < (d_i - t_i) \frac{Q_i}{T_i}$
- pick task with earliest deadline and run it
- while running, decrease budget according to runtime
- if budget reaches zero, task is throttled until recharging happens at time  $d_i$

In single-CPU case, all deadlines are hit if sum of  $Q_i/T_i$  is less than 1

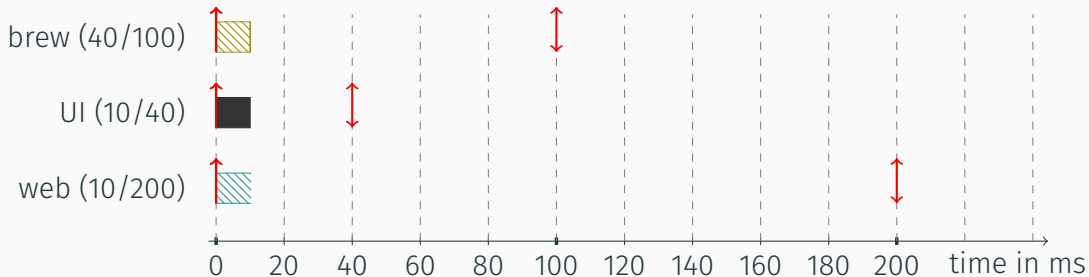
# Example

	brew	UI	web
budget	40	10	10
deadline	100	40	200



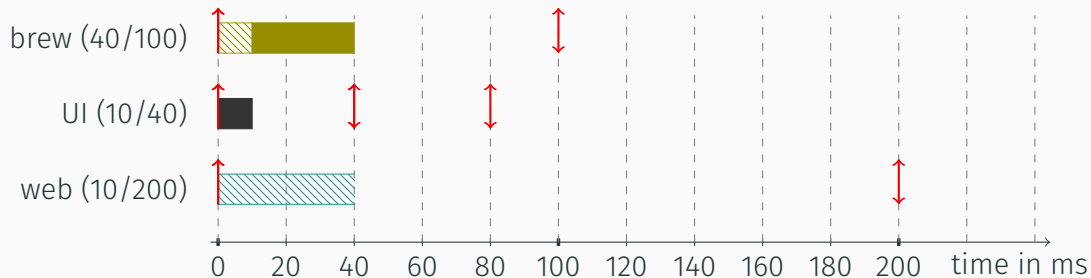
# Example

	brew	UI	web
budget	40	0	10
deadline	100	40	200



# Example

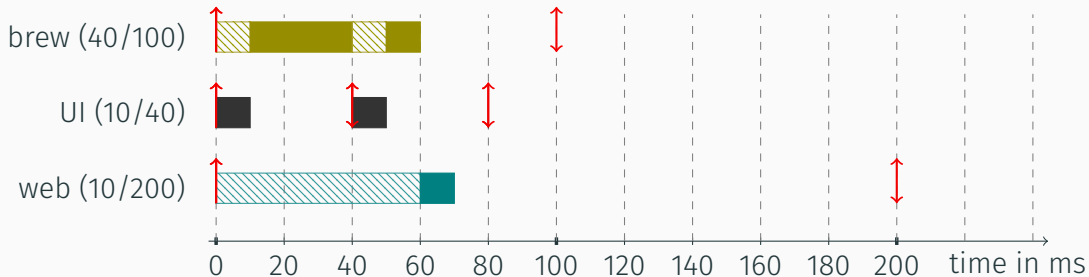
	brew	UI	web
budget	10	10	10
deadline	100	80	200





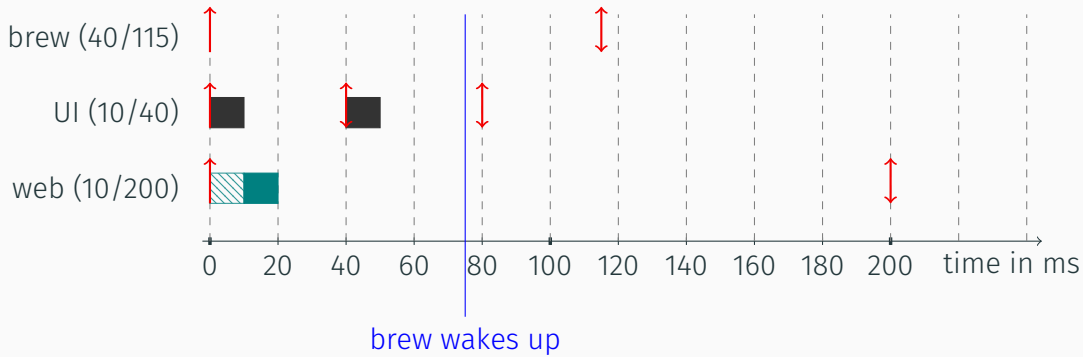
# Example

	brew	UI	web
budget	0	0	0
deadline	100	80	200



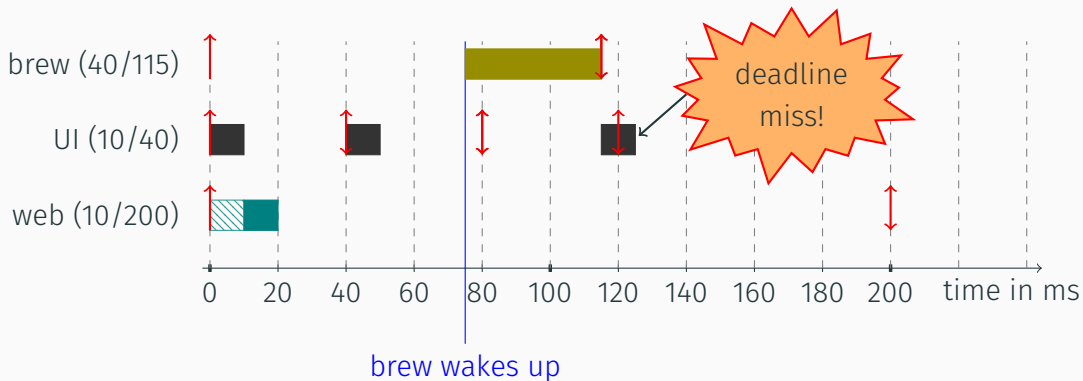
## Example 2

	brew	UI	web
budget	40	0	0
deadline	115	80	200



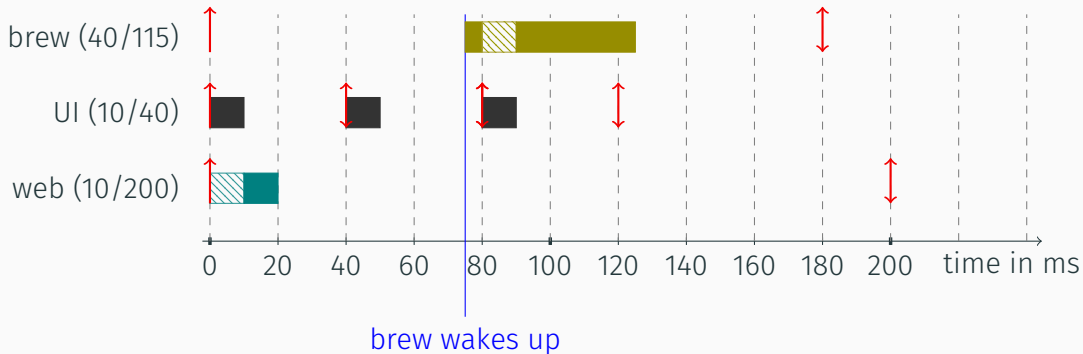
## Example 2

	brew	UI	web
budget	0	10	0
deadline	115	120	200



## Example 2

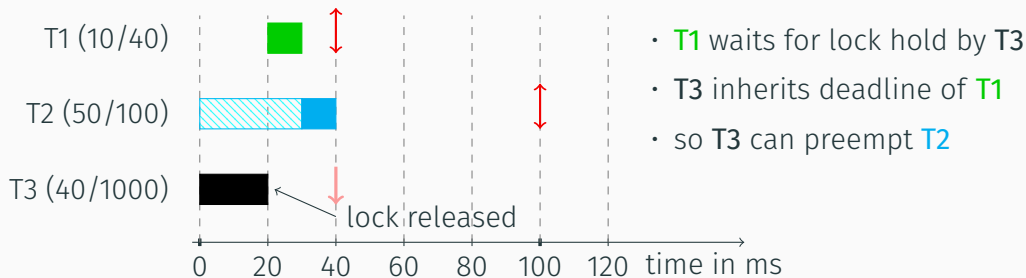
	brew	UI	web
budget	40	10	0
deadline	180	120	200



## API, shared resources, multicore

User level API: two new syscalls `sched_setattr` and `sched_getattr`

Shared resources: inherit deadline of blocked task



Multicore: partitioned (via cpuset) and global supported

# Evaluation

---

## Synthetic workload (partitioned)

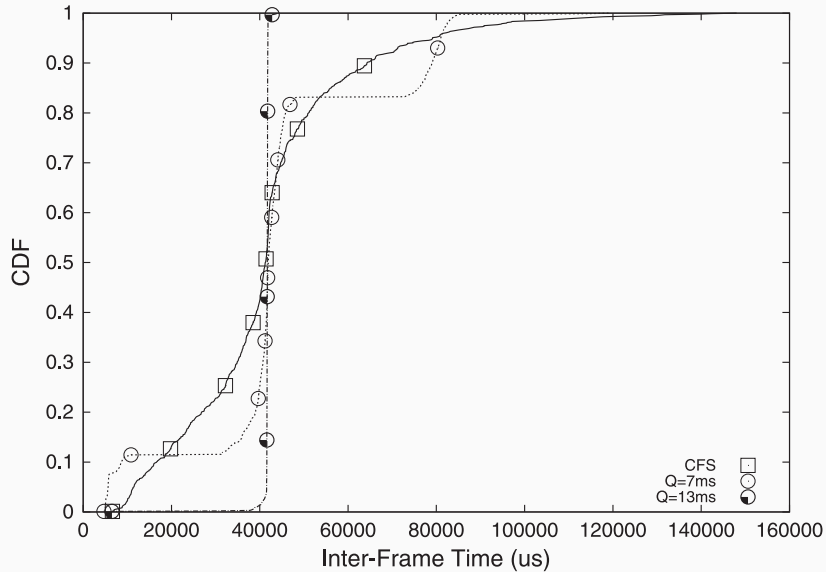
U(%)	SCHED_DEADLINE(%)	SCHED_FIFO(%)	SCHED_OTHER(%)
60	0	0	0.58
70	0	0	1.87
80	0	0.003	6.03
90	0	0.38	10.20

## Synthetic workload (global, six cores)

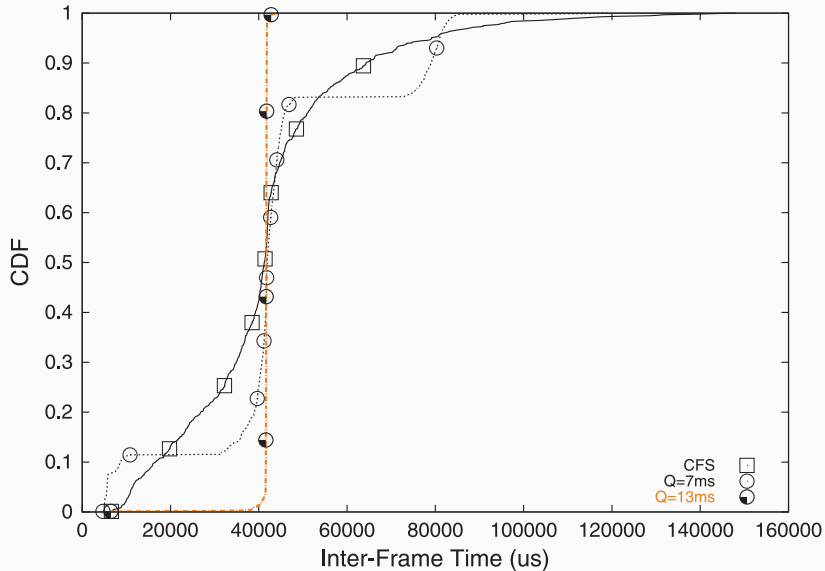
U(%)	SCHED_DEADLINE(%)	SCHED_FIFO(%)	SCHED_OTHER(%)
500	0.027	0.001	3.303
510	0.023	0.002	4.310
520	0.051	0.011	4.992
530	0.099	0.023	6.046
540	0.138	0.230	7.093
550	0.239	0.271	8.097
560	0.289	0.380	9.977
570	0.351	0.640	11.554
580	0.618	1.380	15.384
590	1.295	2.535	19.774



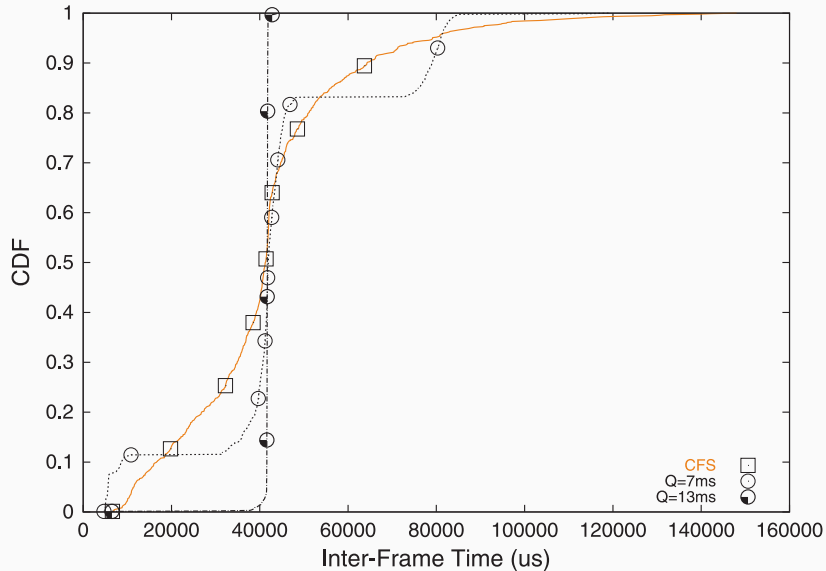
## MPlayer on loaded system



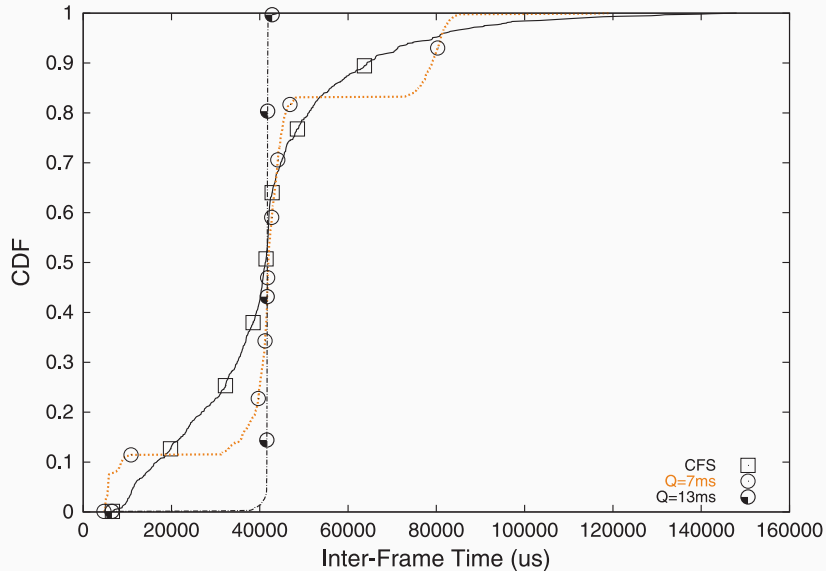
## MPlayer on loaded system



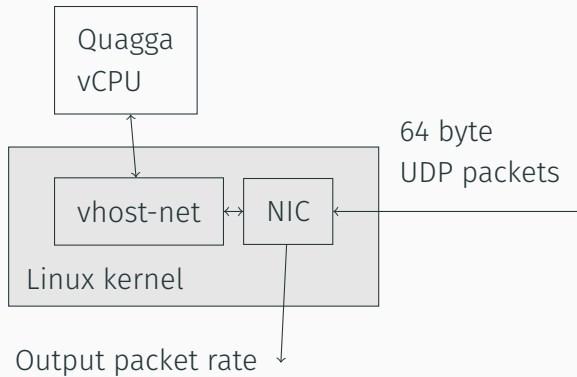
## MPlayer on loaded system



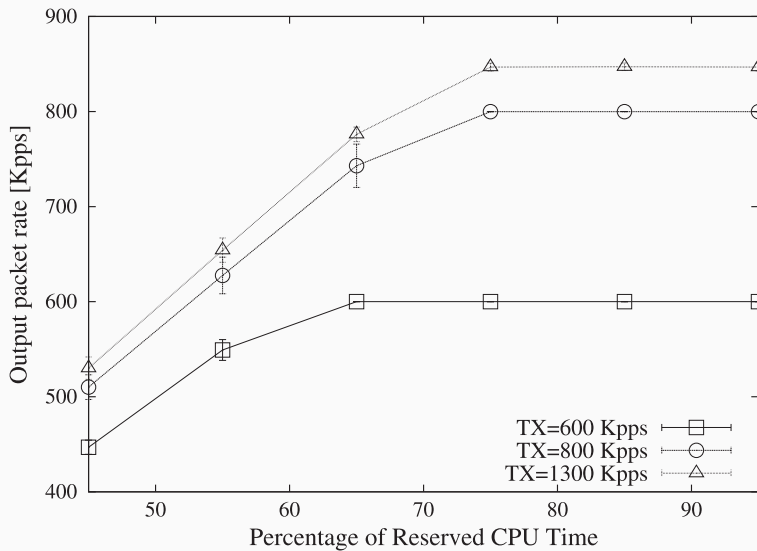
## MPlayer on loaded system



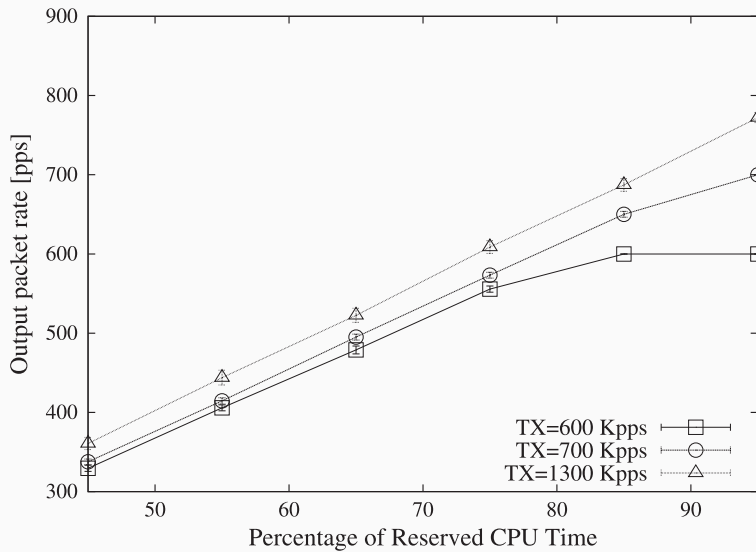
# Virtual router



## Limit vcpu thread



## Limit vhost-net thread



Still some work left:

- M-BWI: Multiprocessor Bandwidth Inheritance
- Power aware algorithms (example: GRUB-PA)
- Support cgroups interface

Not rocket science, but solid implementation of proven concepts (EDF, CBS)

- ... ready for production use
- ... upstream in the Linux kernel
- ... with simple to use API

We can run our coffee machine on stock linux! :)