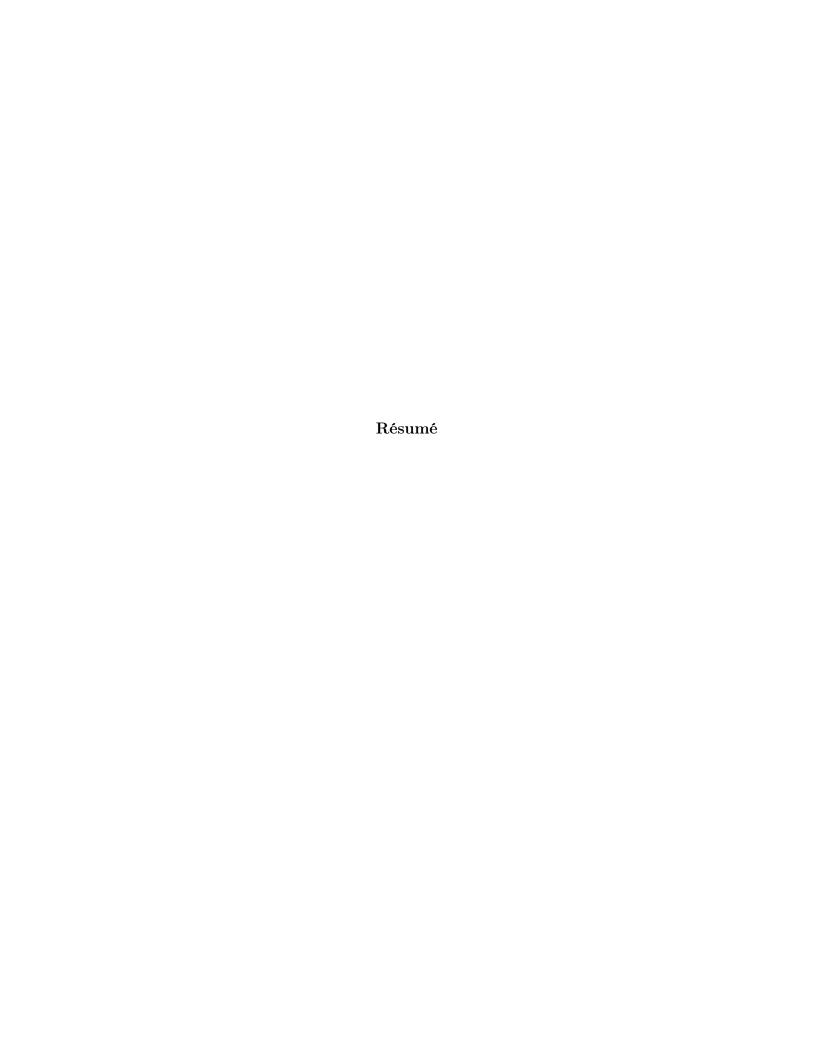
Operating Systems mutations, how and why integrate the user in the digital era?

13 mars 2015



0.1 Introduction

Chapitre 1

Events and Threads

1.1 About concurrent systems

This demonstration focus on application depending on long waiting operations like I/O operations. Particularly, this demonstration focus on real-time web services. It is irrelevant for application heavily relying on CPU operations, like scientific applications.

Threads-based system and event-based system evolved significantly over the last half century. These evolutions were fueled by the long-running debate about which design is better. We try to succinctly and roughly retrace theses evolutions to understand the positions of each community. This demonstration show that thread and events are two faces of the same reality.

Lauer and Needham [Lauer1979] presented an equivalence between Procedure-oriented Systems and a Message-oriented Systems.

Adya et. al. analyzed this debate and presented fives categories through which to present the problem [Adya2002]. These two categories were often associated with thread-based systems and event-based systems. Their advantages and drawbacks were mistaken with those of thread and events. Adya et. al. explain in details two of these categories that are most representative, Task management and Stack management. We paraphrase these explanations.

Task management Consider a task as an encapsulation of part of the logic of a complete application. All the task access the same shared state. The Task management is the strategy chosen to arrange the task executions

in available space and time.

Preemptive task management executes each task concurrently. Their executions interleave on a single core, or overlap on multiple cores. It allows to leverage the parallelism of modern architectures. This parallelism has a cost however, developers are responsible for the synchronization of the shared memory. While accessing a memory cell, it must be locked so that no other task can modify it. Synchronization mechanism impose the developer to be especially aware of race condition, and deadlocks. These synchronization problems make concurrency hard to program with preemptive task management.

The opposite approach, Serial task management, executes each task to completion before starting the next. The exclusivity of execution assures an exclusive access on the memory. Therefore, it removes the need for synchronization mechanism. However, this approach is ill-fitted for modern applications, where concurrency is needed.

A compromise approach, Cooperative task management, allows tasks to yield voluntarily. A task may yield to avoid monopolizing the core for too long. Typically, it yields to avoid waiting on long I/O operations. It merges the concurrency of the preemptive task management, and the exclusive memory access. Thus, it relieves the developer from synchronization problems. But at the cost of dropping parallel execution.

Threads are associated with preemptive task management, and events with Cooperative task management. For this reason, it is commonly believed that synchronization mechanisms make threads hard to program [Ousterhout1996]. While it is really Preemptive task management that is responsible for these synchronization problems [Adya2002].

Stack management Consider a task is composed of several subtasks interleaved with I/O operations. Each I/O operation signal its completion with an event. The task stops at each I/O operation, and must wait the event to continue the execution. The stack management is the strategy chosen to express the sequentiality of the subtasks.

The automatic stack management is what is mostly used in imperative programming. The execution seems to wait the end of the operations to continue with the next instruction. The call stack is kept intact. This is what is commonly called synchronous programming.

In the manual stack management, developers need to manually register

the handlers to continue the execution after the operation. The execution immediately continues with the next instruction, without waiting the completion of the operation. It implies to rip the call stack in two functions; one to initiate the operation, and another to retrieve the result. This is what is commonly called asynchronously programming.

What we argue is that synchronous is good because it is linear, it avoids stack ripping. But asynchronous is good because it allows parallelism by default.

Threads are associated with the automatic stack management, and events with manual stack management. For this reason, it is commonly believed that threads are easier to program. [Thread systems allow programmers to express control flow] [Behren2003]. However, the automatic stack management is not exclusive to threads. Fibers, presented by Adya et. al. is an example of cooperative task management with automatic stack management [Adya2002] . Fibers present the advantage of cooperative task management, without the disadvantage of stack ripping. That is the ease of programming because of the absence of synchronization, without the difficulty of stack ripping.

We argue that the advantages of manual stack management outweigh its drawbacks for web services. Because of the numerous $\rm I/O$ operations, parallelism is

But what is actually highlighted is the automatic state management provided by threads. And with lighter context change, threads are a good choice which provide parallelism.

Historically, events-based system are associated with manual state management, while threads-based systems are associated with automatic state management. Manual state management imposed stack ripping [Adya2002]. With closure, it is not the case anymore. Events now have automatic state management as well [Krohn2007].

Now, there is implementation of thread model with cooperative management, with context-switch overhead improved enough to fill the gap with events model. And there is implementation of event model with automatic state management filling the gap with thread model. In this condition, we ask, what really is the difference between thread and events. We argue there is none. Except the isolation, versus sharing of the memory, which, again is not significant of either. In the first case, the different execution threads exchange messages, while in the second, they use synchronization mechanism to assure invariants in their states

For a single thread of execution, both model could avoid synchronization

through cooperative task management, which assure invariants. Or avoid procedure slicing (if any) using synchronization. These are the two ends of a design spectrum. One end (cooperative task management) fits better for small processing with heavy use of shared resources. While the other end (synchronization) fits better for long processing with small use of shared resources. When one end of the design spectrum is used while the other should be used, one might expect unresponsiveness because of too heavy events, or performance fall due to interlocking.

Scalability is achieved through parallelism, which is itself achieved in our case (web servers) through cluster of commodity machines.

With distribution, this design spectrum gets a better contrast.

The synchronization of distributed, shared resources is limited through the CAP theorem [Gilbert2002a]. Partition tolerance is a requirement of a distributed system. One needs to choose good latency (availability) or consistency. The CAP theorem is generalized into a broader theroem about [Gilbert2012]

The isolation of resources implies to split the architecture in different stages, like Ninja [Gribble2001], SEDA [Welsh2000], or Flash [Pai1999]. This splitting is difficult for the developer. The splitting which is good for the machine, is not the same as the one good for the design in modules.

The two ends of this design spectrum presented map directly onto the two kinds of parallelism advocated for scalability. That is pipeline parallelism, and data parallelism.

Pipeline parallelism is good for data locality, and important throughput. But each stage adds an overhead in latency.

Data parallelism is good for latency, because one request is processed from beginning to the end without waiting in queues. But it implies that the different machines share a common database. Which is a shared resource, and is limited by the CAP theorem.

Both parallelism have advantages and drawbacks, and both could be combined, like in the SEDA architecture. Ultimately, it would be possible to design a design spectrum to choose which kind of parallelism for a set of requirements. But we leave this for future works.

Splitting an architecture in stages is a difficult process, which prevent future code refactoring, and module modifications. We argue that the design for the technical architecture, and the design for the human minds should not be the same. Threads belong to the mental model, the design granularity Events belong to the execution model, the architecture granularity It is a mistake to attempt high concurrency without help from the compiler [Behren2003]. Through compilation, we want to transform an event-loop based program (cooperative task management, no synchronization) into a pipeline parallelism distributed system. So, basically, we argue that it is possible to distribute one loop event onto multiple execution core.

/! WARNING The paper Why events are a bad idea states that : the control flow patterns used by these applications fell into three simple categories : call/return, parallel calls, and pipelines. Indeed, it is no coincidence that common event patterns map cleanly onto the call/return mechanism of threads. Robust systems need acknowledgements for error handling, for storage deallocation, and for cleanup; thus, they need a "return" even in the event model. » Why is it completly false? It is crucial to find an answer. Moreover, Ayda et. al. state that : For the classes of applications we reference here [file servers and web servers], processing is often partitioned into stages. Other system designers advocated non-threaded programming models because they observe that for a certain class of high-performance systems [...] substantial performance improvements can be optained by reducing context switching and carefully implementing application-specific cacheconscious task scheduling.

The paper Why events are a bad idea states that: One could argue that instead of switching to thread systems, we should build tools or languages that address the problems with event systems (i.e., reply matching, live state management, and shared state management). However, such tools would effectively duplicate the syntax and run-time behavior of threads. » Well, yes ... With the exception of the stack junction. The paper on Duality had it right, their graph is correct, but for threads, it cannot be distributed because of stacks, while for events, it can.

Software evolution substantially magnifies the problem of function ripping: when a function evolves from being compute-only to potentially yielding, all functions, along every path from the function whose concurrency semantics have changed to the root of the call graph may potentially have to be ripped in two. (More precisely, all functions up a branch of the call graph will have to be ripped until a function is encountered that already makes its call in continuation-passing form.) We call this phenomenon "stack ripping" and see it as the primary drawback to manual stack management. Note that, as with all global evolutions, functions on the call graph may be maintained by different parties, making the change difficult. » Stack ripping is what I am talking about. While the stack are joined, it is not possible to distribute.

If they say that stack ripping is necessary, that means it is not possible to encapsulate asynchronous function into synchronous function.

1.2 Related Works

1.3 Conclusion