# Studying and removing performance bottlenecks in fault-tolerant communication primitives.

## Context and Objectives

This research work focuses on identifying and eliminating bottlenecks in fault tolerance primitives for distributed systems. The applications behind this research, which use these primitives, include key-value stores such as memcached and consensus services such as etcd and zookeeper. These primitives are essential because they enable an application running on multiple nodes to continue operating despite the outright failure of certain processes or components, or even despite the malicious behavior or compromission of certain nodes in the system. In this context, ensuring fault tolerance while maintaining good performance in such distributed environments remains complex.

Typically, in a distributed application, several fault tolerance primitives are used -namely broadcast, consensus and state machine replication- where nodes frequently exchange messages in communication "rounds". In these primitives, the duplication of message content within and between communication rounds is significant. Message duplication poses several problems: increased latency, as the message must be delivered at the application level, reduced throughput and rapid bandwidth saturation, as network components (switch and NIC) are busy processing the message, and memory saturation (RAM or persistent storage).

This work aims to characterize this message duplication, and then propose and evaluate solutions to detect, avoid or minimize the duplication of these messages.

## Methodology:

#1 - The first phase of this work consists in characterizing the duplication rate. We will experimentally study several fault-tolerant primitives in a cluster of nodes, for example, a consensus primitive with a key-value store. We will vary the parameters that could impact traffic and duplication rate, and will measure duplication at both system and application level.

#2 - The second phase of this work will build on the characterization obtained above, we will design and implement a system that detect and avoid duplication as early as possible in the journey of a packet between its arrival on the network card, its delivery and its use by the application. We will consider the use of kernel bypass techniques to enable the application and the fault-tolerant primitive to deduplicate messages. This will involve studying both the different approaches to deduplicating these messages (passing duplicated data by reference between the network card and the application) and the extent to which this will be possible (the whole/partial content of each message).

#3 - The third phase of this work will consist in evaluating the approaches proposed in step #2 on a set of fault-tolerant primitives, in order to experimentally quantify the impact of the proposed approaches on performance.

**Supervisors (LaBRI – University of Bordeaux) :**
Joachim Bruneau-Queyreix – joachim.bruneau-queyreix@u-bordeaux.fr
Laurent Réveillère – laurent.reveillere@u-bordeaux.fr