

# “Determinación de la eficiencia y Estrategias de Tolerancia a Fallos en Arquitecturas Multiprocesador para aplicaciones de procesamiento de datos”

Jorge R. Osio<sup>1,2</sup>, Diego Montezanti<sup>1,4</sup>, Eduardo Kunysz<sup>1</sup>, Daniel Martin Morales<sup>1,3</sup>

<sup>1</sup>Instituto de Ingeniería y Agronomía - UNAJ

<sup>2</sup>UIDET CeTAD –Fac. de Ingeniería -UNLP

<sup>3</sup>Codiseño HW SW para Aplicaciones en Tiempo Real - UTN - FRLP

<sup>4</sup>Instituto de Investigación en Informática LIDI - Fac. de Informática –UNLP

{josio, dmontezanti, ekunysz, martin.morales}@unaj.edu.ar

## Resumen

Dentro de la línea de investigación que se está desarrollando, existen dos enfoques diferentes. Por un lado se está trabajando sobre la implementación de algoritmos de procesamiento de imágenes sobre dispositivos reconfigurables. El objetivo es utilizar una combinación de diferentes técnicas de concurrencia y paralelismo para tener en cuenta aspectos comunes de dichos algoritmos, y así mejorar la eficiencia de procesamiento sobre las imágenes médicas. Por otra parte, debido a que el procesamiento en paralelo requiere de la implementación de sistemas de múltiples procesadores en dispositivos reconfigurables, se ha incorporado al proyecto una línea en la que se está trabajando en el desarrollo de metodología de tolerancia a fallos transitorios, que son cada vez más frecuentes en las arquitecturas paralelas, y que afectan especialmente a las aplicaciones de cómputo intensivo y ejecuciones de larga duración.

**Palabras clave:** *arquitecturas paralelas, procesamiento de imágenes, checkpoints de*

*capa de sistema, tolerancia a fallos, sistemas multicores, dispositivos reconfigurables.*

## Contexto

Las líneas de Investigación descritas en este trabajo forman parte del Proyecto de Investigación Científico-Tecnológico “Tecnologías de la información y las comunicaciones mediante IoT para la solución de problemas en el medio socio productivo”, que se desarrolla en la Universidad Nacional Arturo Jauretche (UNAJ).

El proyecto cuenta además con financiamiento en el marco del programa “Universidad, Diseño y Desarrollo Productivo” del Ministerio de Educación a través del proyecto “Sistema de eficiencia energética” y con financiamiento de la convocatoria “UNAJ Investiga 2017”.

Parte de las líneas de investigación desarrolladas se encuentran enmarcadas en los convenios de colaboración en Actividades de Investigación firmados por la UNAJ con la UIDET-CeTAD y el Instituto LIDI.

## 1. Introducción

En los últimos años, se ha buscado expandir el concepto del procesamiento paralelo con computadoras basadas en multicores hacia la utilización de plataformas de procesamiento más específicas. Para obtener mayor eficiencia, los fabricantes de computadoras de altas prestaciones, han introducido unidades FPGAs (arreglo de compuertas programables en campo) en su diseño como soporte para el cómputo ([1-3]).

Si bien el estudio de sistemas paralelos con múltiples procesadores, es un campo bien desarrollado, la utilización de múltiples cores en sistemas reconfigurables es un terreno que tiene múltiples posibilidades de exploración [4]. En el presente proyecto se exploran mejoras en la implementación de los algoritmos mediante procesamiento paralelo de SW y Concurrencia en VHDL (lenguaje de descripción de hardware). Adicionalmente, se investiga sobre la tolerancia a fallos de sistemas con gran poder de cómputo y basados en HPC.

### Plataformas FPGAs para procesamiento paralelo

La implementación de paralelismo en plataformas FPGAs consiste en el uso de procesadores embebidos para ejecutar aplicaciones y en la utilización de las características que provee la lógica programable para manejar las porciones de código que se ejecutan concurrentemente [4].

La facilidad de implementar procesadores embebidos en forma rápida, junto con la posibilidad de proveer concurrencia mediante la programación en HW permite combinar las FPGAs con el paralelismo obtenido mediante sistemas multicore para alcanzar la máxima eficiencia.

Entre las ventajas que proveen las FPGAs actuales para el cómputo paralelo se pueden enumerar:

- SoftCores (o procesadores embebidos) que permiten realizar las tareas de administración de datos de baja tasa. Y la posible conexión de IP-Cores

(CoProcesadores) específicos con las siguientes características para realizar el procesamiento duro de señales:

- Capacidad de procesamiento con millones de MACs por segundo (operaciones de Multiplicación / acumulación)
- DSP o Delay Locked Loops (DLL) que permiten la multiplicación o división de la frecuencia de reloj, entre otras tareas.
- Interfaces DRAM / SRAM de alta velocidad y rendimiento.
- Manejo del ancho de banda en señales y buses ahorrando pines I/O.
- Registros de desplazamiento, útiles para buffers de línea o FIFOs.
- RAM distribuida para almacenar coeficientes o pequeñas FIFOs.
- Block RAM con capacidad "true dual-port" para almacenar datos de fotograma, líneas o porciones de imagen, grandes tablas o FIFOs.

Con las mejoras constantes que aporta la evolución de la tecnología sobre las FPGAs pueden lograrse diseños de gran magnitud, a tal punto que la tendencia actual es implementar microprocesadores de propósito general, conjuntamente con todo el hardware de propósito específico que requiere la aplicación, dentro de una FPGA.

### Tolerancia a Fallo en sistemas multicore

El aumento en la escala de integración, con el objetivo de mejorar las prestaciones en los procesadores, y el crecimiento del tamaño de los sistemas de cómputo, han producido que la confiabilidad se haya vuelto un aspecto crítico. En particular, la creciente vulnerabilidad a los fallos transitorios se ha vuelto altamente relevante, a causa de la capacidad de estos fallos de alterar los resultados de las aplicaciones.

El impacto de los fallos transitorios aumenta notoriamente en el contexto del Cómputo de Altas Prestaciones, debido a que el Tiempo Medio Entre Fallos (MTBF) del sistema

disminuye al incrementarse el número de procesadores. En un escenario típico, en el cual cientos o miles de núcleos de procesamiento trabajan en conjunto para ejecutar aplicaciones paralelas, la incidencia de los fallos transitorios crece en el caso de que las ejecuciones tengan una elevada duración, debido a que el tiempo de cómputo y los recursos utilizados desperdiciados resultan mayores [9 - 10]. Por lo tanto, el alto costo que implica volver a lanzar una ejecución desde el comienzo, en caso de que un fallo transitorio produzca la finalización de la aplicación con resultados incorrectos, justifica la necesidad de desarrollar estrategias específicas para mejorar la confiabilidad y robustez en sistemas de HPC. En particular, el foco está puesto en lograr detección y recuperación automática de los fallos silenciosos, que no son detectados por ninguna capa del software del sistema, por lo que, sin producir una finalización abrupta, son capaces de corromper los resultados de la ejecución.

### **Antecedentes del Grupo de Trabajo**

Los antecedentes del grupo de trabajo parten de la investigación acerca de procesamiento de imágenes sobre diferentes arquitecturas paralelas [3-6] y de la investigación de la Tolerancia a fallos en sistemas multicore ([3], [7] y [8]).

Las líneas de investigación se enmarcan fundamentalmente en la búsqueda de técnicas innovadoras de procesamiento paralelo, sobre diferentes arquitecturas, que optimicen las prestaciones en equipos de imágenes médicas (PET, X-Ray, Ecógrafo). Adicionalmente, se pretenden encauzar temas de estudios de posgrado actuales y futuros de docentes y profesionales surgidos de la UNAJ dentro de estas líneas de investigación. El director del proyecto participa en proyectos de investigación desde el año 2005 en la Facultad de Ingeniería de la UNLP, además dirige proyectos de investigación desde el año 2015 en la UNAJ. Actualmente, el grupo de trabajo investiga en los temas: Procesamiento Digital

de Imágenes Médicas sobre plataformas FPGA, Procesamiento de Imágenes en sistemas HPRC y tolerancia a fallos pertenecientes a estudios de Maestrías y Doctorados que se realizan en el marco de acuerdos de colaboración entre el Laboratorio CeTAD, el Instituto Lidi y la UNAJ a través de los respectivos proyectos de investigación.

Como resultado del trabajo realizado recientemente en la UNAJ se han publicado artículos en diferentes congresos nacionales de Ingeniería Informática y Ciencias de la Computación [3 - 7].

## **2. Líneas de Investigación y Desarrollo**

El grupo de investigación que se ha constituido recientemente en la UNAJ es multidisciplinario, y sus miembros cuentan con experiencia en sistemas multiprocesador, procesamiento de imágenes.

### **Temas de Estudio e Investigación**

- Implementación de un sistema multiprocesador en Dispositivos Lógicos Programables (FPGAs).
- Análisis y determinación del desempeño logrado en el procesamiento de imágenes mediante la combinación de cómputo paralelo y concurrencia.
- Diseño e implementación de una estrategia que permite al sistema detectar y recuperarse automáticamente de los errores producidos por fallos transitorios. La implementación está basada en replicación de software y es totalmente distribuida; está diseñada para operar en entornos de clusters de multicores donde se ejecutan aplicaciones paralelas científicas de paso de mensajes.

## **3. Resultados Obtenidos/Esperados**

### **Investigación experimental**

Hasta el momento se han obtenido resultados satisfactorios en relación los objetivos principales:

- En cuando al procesamiento de imágenes sobre dispositivos FPGAs, se han implementados varios algoritmos en sistemas con dos microblaze, en donde el acceso a los datos se realizó mediante memoria compartida DDR2-SDRAM de alta velocidad y con múltiples puertos. El procesador master es el encargado de coordinar la lectura y procesamiento de los datos mediante pasaje de mensajes. Respecto a los resultados, el tiempo de procesamiento se disminuye casi al 50% al utilizar dos procesadores. Se han obtenido tiempos muy buenos en los filtros de mediana y de media (4 ms por fila), pero este se duplica al ejecutar algoritmos con mucho cálculo matemático como el algoritmo de sobel (7,5 ms por fila).

- En cuanto a la tolerancia a los fallos transitorios, se ha diseñado e implementado una estrategia de detección basada en la replicación de procesos y el monitoreo de las comunicaciones, y un mecanismo de recuperación basado en el almacenamiento de un conjunto de checkpoints distribuidos de capa de sistema. La recuperación es automática y se realiza regresando atrás al checkpoint correspondiente a la cantidad de veces que se detecta un fallo determinado. Esta estrategia se encuentra en proceso de validación experimental. Para ello, se ha desarrollado un conjunto de casos de prueba. Además, se ha diseñado una estrategia de recuperación alternativa, para el caso en el que se cuente con checkpointing de capa de aplicación. En este caso, es posible almacenar sólo el último checkpoint luego de validarlo.

Para el año en curso, se esperan alcanzar importantes resultados en el área de cómputo paralelo y concurrencia posibilitados por los sistemas basados en FPGAs. Con esta arquitectura, constituida por varios procesadores implementados en una misma FPGA se espera mejorar el desempeño logrado con el sistema multiprocesador basado en master-slave y memoria compartida.

Se pretende medir la eficiencia de ejecución lograda con el sistema multiprocesador mediante los siguientes algoritmos basados en operadores de ventana:

- Filtro de media
- Algoritmo de mediana, erosión y dilatación
- Algoritmos de detección de borde (Sobel)

En la línea de tolerancia a fallos, se ha diseñado e implementado una metodología distribuida basada en replicación de software, diseñada específicamente para aplicaciones paralelas científicas de paso de mensajes, capaz de detectar los fallos transitorios que producirían resultados incorrectos y recuperar de manera automática las ejecuciones [8]. En este tipo de aplicaciones, cuyos procesos cooperan para obtener un resultado, la mayor parte de los datos relevantes son comunicados entre tareas. Por lo tanto, la estrategia de detección consiste en la validación de los contenidos de los mensajes a enviar y de los resultados finales. Esta solución representa un término medio entre un alto nivel de cobertura frente a fallos y la introducción de un bajo *overhead* temporal; además, no se realiza trabajo para detectar fallos que no afectan a los resultados. Por otra parte, de esta forma, la corrupción de los datos que utiliza un proceso de la aplicación se mantiene aislada en el contexto de ese proceso, evitándose la propagación a otros procesos. Así, no solamente se mejora la confiabilidad del sistema, sino que también disminuye la latencia de detección, y por lo tanto el tiempo luego del cual se puede relanzar la aplicación, lo cual es especialmente útil en ejecuciones prolongadas, y beneficia al sistema aprovechando la redundancia intrínseca presente en la arquitectura multicore.

Para recuperar al sistema de los efectos del error, la propuesta se basa en restaurar la aplicación a un estado seguro previo a su ocurrencia. Para lograr este objetivo, se ha integrado la detección con un mecanismo basado en múltiples *checkpoints* coordinados en capa de sistema, construidos con la librería DMTCP (que proporcionan cobertura en el caso de que un *checkpoint* resulte afectado por un fallo, imposibilitando la recuperación desde él),

o la utilización de un único *checkpoint* no coordinado de capa de aplicación (construido ad-hoc, basándose en el conocimiento de la aplicación), que puede ser verificado para asegurar la integridad de sus datos.

Para una aplicación de prueba, se ha desarrollado un modelo formal de verificación, que permite constatar la eficacia del mecanismo de recuperación. Se encuentra en fase de implementación una experimentación basada en casos de inyección controlada de fallos que valide las predicciones dadas por el modelo mencionado. Se busca además caracterizar el comportamiento temporal de cada nivel de solución (sólo detección, recuperación basada en múltiples *checkpoints* de capa de sistema o recuperación basada en un único *checkpoint* de capa de aplicación) y verificar el comportamiento pronosticado.

El uso de estas estrategias posibilitaría prescindir de la utilización de redundancia triple con votación para detectar y recuperar de fallos transitorios, proveyendo opciones que proporcionan flexibilidad para adaptarse a los requerimientos o limitaciones del sistema. Además, como estos fallos no requieren reconfiguración del sistema, la recuperación puede realizarse mediante re-ejecución en el mismo core en el que ocurrió el fallo.

#### 4. Formación de Recursos Humanos

Dentro de la temática de la línea de I+D, todos los miembros del proyecto participan en el dictado de asignaturas de la carrera de Ingeniería Informática de la UNAJ dentro del Área Arquitectura, Sistemas Operativos y Redes.

En este proyecto existe cooperación a nivel nacional. Hay un Doctor, un investigador realizando Doctorado y dos realizando Maestrías en temas relacionados con simulación de sistemas multiprocesador, tolerancia a fallos, sistemas embebidos y sistemas multicore en HPRC.

Adicionalmente, se cuenta con la colaboración de estudiantes avanzados.

#### 5. Bibliografía

- [1] O. Mencer, K. Tsoi, S. Cramer, T. Todman, W. Luk, Ming Wong and P. Leong, "CUBE: a 512-FPGA Cluster", Dept. of Computing, Imperial College London, Dept. of Computer Science and Engineering The Chinese University of Hong Kong. (2009)
- [2] Keith Underwood, "FPGAs vs. CPUs: Trends in Peak Floating-Point Performance", Sandia National Laboratories. (2011)
- [3] J. Osio, D. Montezanti, E. Kunysz, Morales M., "Análisis de eficiencia y tolerancia a fallo en Arquitecturas Multiprocesador para aplicaciones de procesamiento de datos", UNNE, Corrientes, WICC 2018.
- [4] J. Osio, J. Salvatore, E. Kunysz, V. Guarepi, M. Morales, "Análisis de Eficiencia en Arquitecturas Multiprocesador para Aplicaciones de Transmisión y Procesamiento de Datos", UNER, Ciudad de Concordia, WICC 2016.
- [5] J. Osio, D. Montezanti, M. Morales, "Análisis de Eficiencia en Sistemas Paralelos", Ushuahia, Tierra del Fuego, WICC 2014
- [6] E. Kunysz, J. Rapallini, J. Osio, "Sistema de cómputo reconfigurable de alta performance (Proyecto HPRC)", 3ras Jornadas ITE - 2015 - Facultad de Ingeniería - UNLP
- [7] D. Montezanti, A. De Giusti, M. Naiouf, J. Villamayor, D. Rexachs, E. Luque, "A Methodology for Soft Errors Detection and Automatic Recovery", in Proceedings of the 15th International Conference on High Performance Computing & Simulation (HPCS). ISBN: 978-1-5386-3250-5/17. IEEE, 2017, pp. 434
- [8] F. Cappello, A. Geist, W. Gropp, S. Kale, B. Kramer, and M. Snir, "Toward exascale resilience: 2014 update," Supercomputing frontiers and innovations, vol. 1, no. 1, pp. 5–28, 2014.
- [9] Grama A, Gupta A, Karypis G, Kumar V. "Introduction to parallel computing". Pearson Addison Wesley, 2003.
- [10] F. Cappello, A. Geist, W. Gropp, S. Kale, B. Kramer, and M. Snir, "Toward exascale resilience: 2014 update," Supercomputing frontiers and innovations, vol. 1, no. 1, pp. 5–28, 2014.