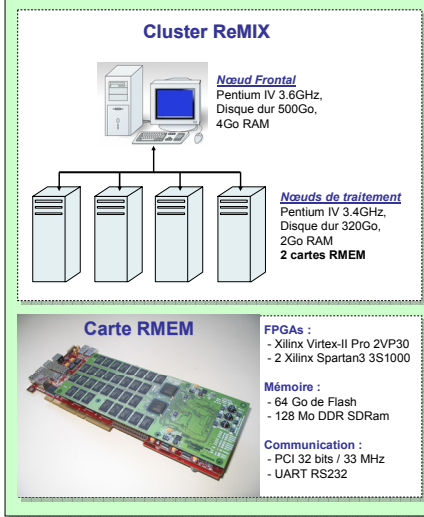


Mémoire Reconfigurable pour l'Indexation de Masses de Données

L. Amsaleg³, F. Charot², S. Derrien², G. Georges¹, D. Lavenier¹, P-F. Marteau⁴, G. Ménier⁴, E. Popovici⁴, F. Raimbault⁴, S. Rubini¹.
Laboratoires : IRISA (Symbiose¹, R2D2², TEXMEX³) - VALORIA (APRIM⁴)

Plate-forme matérielle

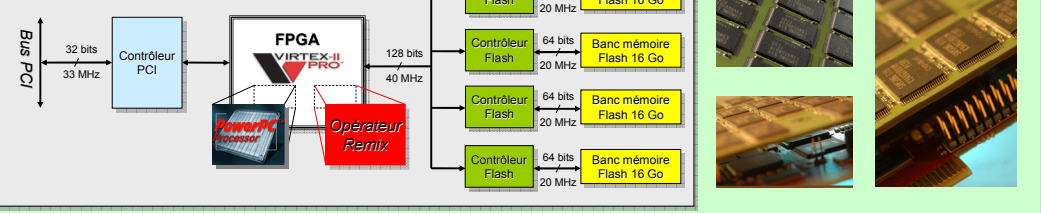


L'ACI ReMIX propose l'élaboration d'une mémoire spécialisée de très grande taille, dans le but d'accélérer la recherche d'informations dans des bases de données indexées. Une architecture matérielle dédiée a été développée. Trois champs disciplinaires représentatifs serviront de support pour valider cette proposition : la génomique, la recherche d'images par le contenu et la recherche documentaire basée sur les textes et leurs structures.

Caractéristiques :

- Cluster de 5 PCs dont 4 nœuds de traitement et un nœud frontal. Chaque nœud de traitement est doté de 2 cartes mémoires reconfigurables (RMem) de 64Go, portant la mémoire d'index globale à 512Go.
- Index distribués sur 8 supports physique indépendants.
- Technologie mémoire « Flash Nand » : latence d'accès en 200 à 500 fois inférieure à un disque dur.
- Bande passante agrégée en sortie de mémoire est d'environ 5Go/s.
- Implantation d'opérateurs matériels reconfigurables, dans un FPGA, permettant un traitement/tri efficace des données en sortie de la mémoire.
- Système de fichiers dédié à la gestion bancs mémoires Flash.
- Environnement de programmation de haut niveau parallélisant les traitements au niveau du cluster.

Architecture de la RMem :



Recherche par le contenu dans les bases de séquences génomiques

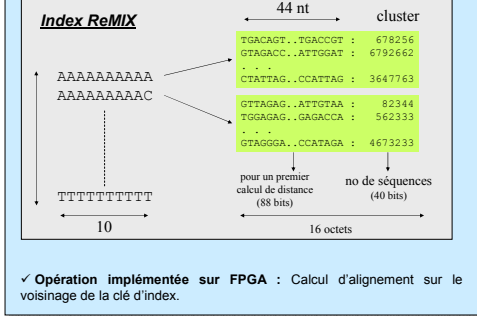
- Objectif :** Indexation du programme BLAST sur ReMIX.
- BLAST (Basic Local Alignment Search Tool) :** Programme le plus utilisé en biologie moléculaire pour la recherche d'alignement entre des séquences génomiques.
- Temps de calcul** proportionnel à la taille des bases de données qui double tous les 14 mois.
- Heuristique :** Dans un alignement les 2 séquences partagent au moins W caractères.

```

Banque
  tggcggatggcgcagaccgagactggcgataacgattga
Requête  gttgagaccgagactggccgga      W=4
    
```

Principe de BLAST : Indexation de la séquence requête, et parcourt de toute la banque génomique.
=> Temps de calcul limité par le débit des données.

Indexation sur ReMIX (iBLAST) : Indexation de la banque et parcourt de la requête.



Recherche par le contenu dans des banques d'images

- Objectif :** Retrouver une image éventuellement modifiée dans une banque d'images protégées par le copyright.
- 94% du chiffre d'affaire généré par le marché de la distribution de photos à usage professionnel (plus de 1600M Euros) est lié au copyright.
=> Importance de la vérification du respect du copyright.
- Recherche par le contenu :** basée sur la présence de similitudes visuelles entre l'image piratée et l'image originale.
- Nécessité d'être robuste** aux altérations éventuellement sévères.
- Une image est définie par des descripteurs locaux de 24 dimensions : 50 à 1000 descripteurs par image.
- Description d'une image (base ou requête) par détection automatique des points d'intérêts (Harris), puis traitement du signal autour de chaque point (convolution, dérivation, mélange spécifique).
- Processus de recherche :**
=> 1 : Description de l'image requête,
=> 2 : Interrogation de la base pour chaque descripteur,
=> 3 : Sélection, liste de plus proches voisins pour chaque descripteur,
=> 4 : Election des meilleures images.

Opération implémentée sur FPGA : Calcul de distance euclidienne entre descripteurs multidimensionnels.

Recherche approchée d'information dans des bases de documents semi-structurés

- Objectif :** Indexation et interrogation de bases de documents semi-structurés de taille de l'ordre de 50 à 100 Go.
- Documents** représentés sous la forme d'un arbre DOM, assimilable à un ensemble de chemins (P²).

Arbre DOM (Document Object Model)

```

graph TD
  doc --> title
  doc --> author
  doc --> text
  chapter1[chapter num="1"] --- Omphalos
  chapter2[chapter num="2"] --- Alice
  
```

Type de recherche : Hors contexte, en contexte ou sur la structure des documents.

- Requêtes complexes** fragmentées en requêtes élémentaires réparties sur les nœuds ReMIX.
- Opérations assembleuses** et création du résultat réalisés sur l'hôte.

Opération à implémenter sur FPGA : Calcul de distance entre deux chemins (Distance d'édition de Levenshtein).

Publications

Conférences internationales
E. Popovici, G. Ménier, P-F. Marteau, SIRIUS: A Lightweight XML Indexing and Approximate Search System at INEX 2005, Advances in XML Information Retrieval and Evaluation: Fourth Workshop of the Initiative for the Evaluation of XML Retrieval (INEX 2005), Duisburg, Germany, Lecture Notes in Computer Science, vol.3977, 2006.
D. Lavenier, L. Xinchun, G. Georges, Seed-based Genomic Sequence Comparison using a FPGA/FLASH Accelerator, *International IEEE Conference on Field Programmable Technology (FPT)*, Bangkok, Thailand, 2006
Van Hoa Nguyen, D. Lavenier, Recherche dans les banques d'ADN par indexation parallèle, *4th International Conference on research, innovation & vision for the future*, Ho Chi Minh Ville, Vietnam, 2006
E. Popovici, G. Ménier, P-F. Marteau, Information Retrieval of Sequential Data in Heterogeneous XML Databases, *Adaptive Multimedia Retrieval: User, Context, and Feedback: Third International Workshop*, LNCS 3877, AMR 2005, Glasgow, UK, 2005.
G. Ménier, P-F. Marteau, A. Azarian, PARTAGE: Software Prototype for Dynamic Management of Document and Data, 18th International Conference on Software and Systems Engineering and their Applications, Paris, 2005.
S.-A. Berrani, L. Amsaleg, P. Gros, Robust Content-Based Image Searches for Copyright Protection. Proc. of the ACM International Workshop on Multimedia Databases, pages 70-77, New Orleans, Louisiana, USA, November 2003.
S.-A. Berrani, L. Amsaleg, P. Gros, Approximate Searches: k-Neighbors + Precision. Proc. of the 12th ACM International Conference on Information and Knowledge Management, pages 24-31, New Orleans, Louisiana, USA, November 2003.

Conférences nationales
E. Popovici, G. Ménier, P-F. Marteau, Recherche approchée d'information dans une base de documents semi-structurés, 3ème Conférence en Recherche d'Informations et Applications (CORIA'06), Lyon, France, 2006.
E. Popovici, G. Ménier, P-F. Marteau, Recherche approchée d'information dans une base de documents semi-structurés : une application ReMIX, Panorama des Recherches Initiatives en STIC (PaRiSTIC), LaBRI, Bordeaux, France, 2005.
G. Georges, S. Derrien, S. Rubini, F. Raimbault, L. Amsaleg, D. Lavenier, ReMIX : une architecture pour la recherche dans les masses, *Sympa 2006, Symposium en Architecture de Machines*, Perpignan, France, 2006
N. Ben Zaccour, R. Bouville, D. Lavenier, M. Gauthier, Y. Leloir, Utilisation de l'indexation de séquences et du calcul thermodynamique pour optimiser la spécificité des oligonucléotides, *JOBIM 2005, Journées Ouvertes Biologie, Informatique et Mathématique*, Lyon, juillet 2005.
E. Popovici, G. Ménier, P-F. Marteau, D-Pilgrim : une base relationnelle/documentaire mobile, 15ème journée francophone d'ingénierie des connaissances (IC'2004), Lyon, France, 2004

Contacts :

Dominique LAVENIER
IRISA
Campus de Beaulieu
35042 Rennes Cedex
Tel : (33) 2 99 84 72 17
Email : lavenier@irisa.fr

Gilles GEORGES
IRISA
Campus de Beaulieu
35042 Rennes Cedex
Tel : (33) 2 99 84 73 21
Email : georges@irisa.fr

<http://www.irisa.fr/remix>

Action Incitative Masse de Données

ministère de l'Éducation nationale et de la Recherche
ministère délégué recherche et nouvelles technologies

MINISTÈRE DE LA DÉFENSE
DGA

INRIA