

De 2 à 22 millions d'images; Création, Indexation et Recherche par le contenu avec Piria

contact : patrick.hede@cea.fr

Commissariat à l'Energie Atomique
List

list

cea



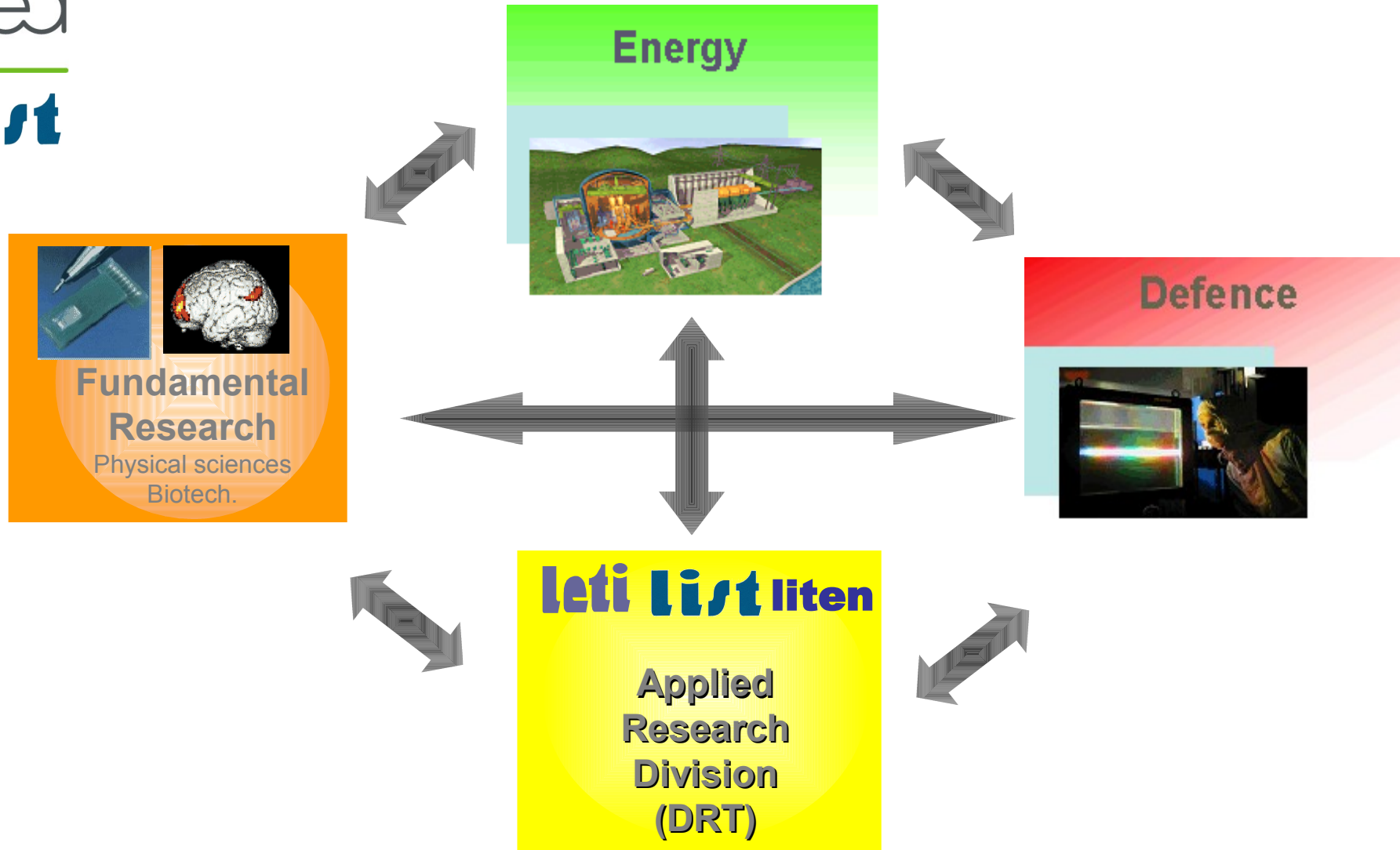


- **Le CEA, le Lic2m**
- **Le moteur image Piria**
- **Le projet Fame II**
 - ✓ **Les partenaires, architecture BULL**
- **Le grand challenge image**
- **Travail réalisé**
 - ✓ **Le corpus (base d'images)**
 - ✓ **Piria MPI**
 - **Indexation**
 - **Recherche**
- **Performances**
- **Conclusions & perspectives**

Le CEA: Commissariat à l'Energie Atomique

cea

list





Création en 2002 (Christian Fluhr)

Héritage: moteur de recherche cross langue Spirit, compétences sur le traitement linguistique et la recherche d'information scientifique et technique.

Héritage complémentaire de compétences sur l'analyse et l'indexation d'image

Effectif:

20 personnes, 13 permanents

Acteur majeur intégré au réseau de recherche française en TAL et en recherche d'information multimédia

CNRS, INRIA, Universités, Ecole de Mines, ENST

Engagement sur projets (MuSCLE, MEDAR, imageval, Fame2)

Transfert de technologie

Partenariat avec partenaires industriels sur projets collaboratifs (Bull, Thales, Alcatel..)

A l'origine de la start-up NewPhenix

Activités de recherche



Image

descripteurs locaux, fusion de descripteurs, indexation, recherche
catégorisation, modèles 3D

Vidéo

structure, résumé

Texte

analyse sémantique
extraction d'information, catégorisation, résumé, filtrage
question réponse

Analyse et indexation document multimédia

moteur de recherche multimédia (fusion de descripteurs, fusion de
résultats de comparaison)
analyse de la structure de documents multimédia

Extensions

constitution de ressources
passage à l'échelle: indexation de gros volume de données
amélioration de la transcription (syntaxe+sémantique)
constitution de ressources (pour la traduction: alignement de texte,
dictionnaires bilingues, pour la sémantique: carte sémantique générale
et métier)



➤ CBIR Piria

- ✓ Programme d'Indexation et de Recherche d'images par Affinité créé en 2002
- ✓ Moteur développé au CEA LIST écrit en C++ utilisant les STL
- ✓ **Caractéristique réponse rapide**,
convertit les signatures numériques des descripteurs en langage naturel, manipule l'image le texte la vidéo ;
bientôt un moteur multimédia.

1 milliard de secondes = ...

Analyse des images: Création de signatures



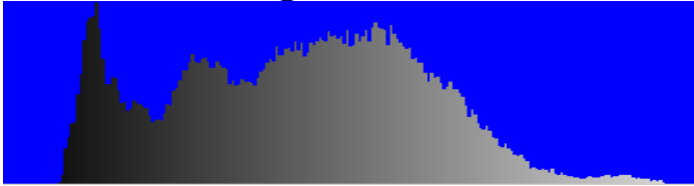
Extraction de caractéristiques

Texture, Couleur, Forme, Points d'intérêt

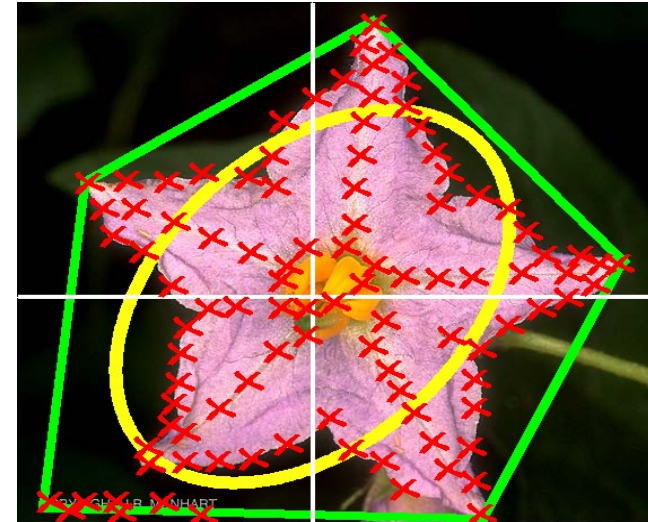
Information globale

Information locale

Ex: Histogramme



Ex :Points singuliers



Avec segmentation systématique

signature de l'image

Moteur indexation/recherche



➤ CBIR Piria

✓ Recherche par le contenu (Hors Google, FlickrR)

PIRIA
LIST/DTSI/SCRI/LIC2M Moteur Images

Base d'images de 100 objets Coil100 de l'Université de Columbia - 7 200 images

MELANGER : Ambiance (Couleurs) Texture (Texture) Silhouette (Forme)

MODE : Affichage des Réponses

contact
La base d'images du Professeur J.WANG
La base d'images de l'Université de Columbia

CEA LIST - Démonstration du moteur de recherche image développé par le LIST: PiRIA - Mozilla Firefox

http://www-list.cea.fr/fr/programmes/systemes_interactifs/labo Lic2m/piria/w3/pirianet.php?bdi=coil-100&cide=tle

CEA LIST/DTSI/SCRI/LIC2M Moteur Images

Base d'images de 100 objets Coil100 de l'Université de Columbia - 7 200 images

MELANGER : Ambiance (Couleurs) Texture (Texture) Silhouette (Forme)

MODE : Affichage des Réponses

contact
La base d'images du Professeur J.WANG
La base d'images de l'Université de Columbia

Rechercher : millo Surligner tout Respecter la casse

Terminé



Pourquoi indexer?



list

- **En aout 2005 Google recence 2,1 milliards d'images sur le web sans doute environ 5 milliards en Février 2008**
- **En Octobre 2006 FlickrR 3000 images/jour , 4000 en Février 2008**
- **Ina, Nasa, les particuliers sont des gros générateurs de contenu numériques. Mais très peu, moins de 10% du web annotent.**
- **Un disque dur de 1Tera Octets pour moins de 500€ en grands magasins**

volume





- **Problème : On recherche une certaine image(ou document:texte multilingue, vidéo, sons)**

Solution : Sans doute une version proche existe sur internet!

Savoir que l'on dispose d'une information sans pouvoir y accéder revient à ne pas disposer de cette information.

Sous Problème : comment la trouver?

- **Solution utiliser un moteur de recherche**

Fame2

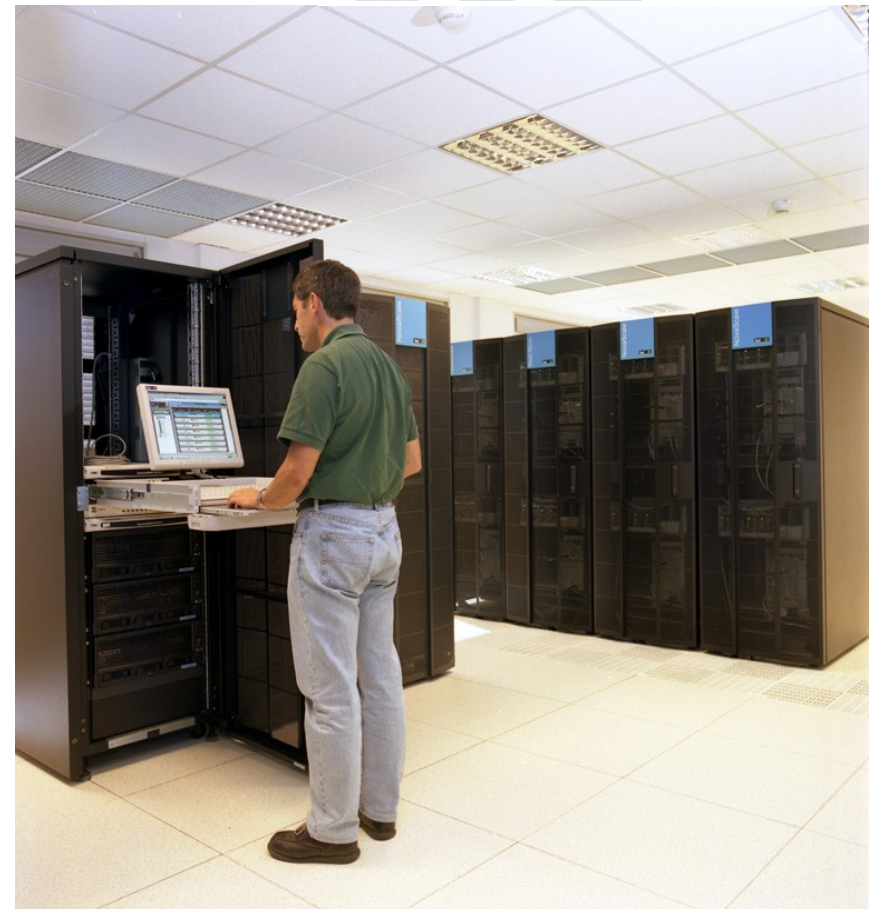
cea

list

- 88 cœurs de calcul hiérarchie de parallélismes intra et inter socket et inter modules
- Accélérateurs spécifiques
- 50 Terabytes de disques



BULL



- **Le grand challenge image**

- **Etat de l'art :Le système Cortina (Univ. Santa Barbara)**
 - ✓ **Indexation : 11 millions en ligne**
 - ✓ **Recherche : 15 secondes**

- **Piria**
 - ✓ **Indexation : la plus grosse base 50 000 images**
 - ✓ **Recherche : 5s**

- **But: faire mieux!**



- **Source du Moteur sous Win32 .net**
- **10 descripteurs globaux, 3 locaux**
- **Tests sur des bases de 1000 a 50 000 images**
- **http://www-list.cea.fr/fr/programmes/systemes_interactifs/labo_lic2m/piria/w3/pirianet.php**



- **Difficulté N°1 : comment trouver plusieurs millions d'images © ?**

- **Multiplication des données**
 - ✓ **CLIC (CEA Lic2m Image Collection)**
 - Transformations géométriques et chromatiques 1million d'images a partir d'un noyau de 15 200.

- **Fame II**
 - ✓ **Corpus multimédia multilingue: wikipédia fr, gb,**



➤ **Fame II**

- ✓ **Corpus multimédia multilingue: wikipédia fr, gb aspiré**
 - Un peu moins d'un million d'images
- ✓ **Filtrage des ©, récupération des images**
- ✓ **Mise à dimension 320x200 maxi et conversion au format JPEG**
- ✓ **En conservant l'arborescence limitant 10 000 images maxi par répertoire**
- ✓ **Application de 25 transformations sur toutes les images**
- ✓ **Obtention de 22 millions d'images**



- ✓ **Choix d'un descripteur de Piria existant**
- ✓ **Amélioration de ce descripteur**
 - Réduction de la dimension de sa signature par 2
 - Accélération du code de calcul par optimisation
 - Élimination de fuite mémoires non détectées avant le passage à l'échelle.
- ✓ **Conception et Codage d'une méthode permettant l'indexation de plusieurs répertoires**
- ✓ **Écriture d'un wrapper permettant via les MPI l'exécution simultanée de plusieurs instances du moteur Piria**
- ✓ **Monté en volume par itération successives**

- ✓ **Moins de 100 heures de calcul**



➤ Couleur Contour BIC

- ✓ *Ref: A compact and efficient image retrieval approach based on border/interior pixel classification Stehling, Nascimento, Falcão CIKM 2002*



list

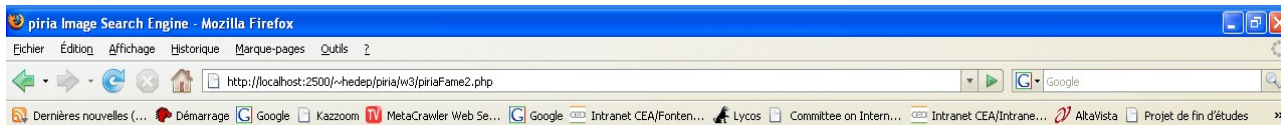
- ❖ Portage Linux P4 (Woodcrest), I64 (Montécio)
- ❖ Aspiration des images de wikipédia (français, anglais, commons)
 - ❖ Mise en forme (jpeg, resize)
- ❖ Création d'une base de 22 millions par transformations, duplication
- ❖ Amélioration d'un descripteur (réduction de la taille de sa signature par 2)

- ❖ Écriture d'un wrapper MPI Bull v 1.
- ❖ Architecture logicielle 'scalable'
 - ❖ En nombre de cœur de 1 a 'n' (test: 1, 16, 48 et 80 coeurs)
 - ❖ Arborescence d'image quelconque (lustre plus de 700 répertoires)
- ❖ Indexation de 10 000, 60 000, 101 000, 700 000, et
22Millions d'images 3To de données
- ❖ Tests de pertinence et de temps de réponse
- ❖ Conception et écriture d'une Interface php (Avec Bull serveur web)
pour l'interrogation et l'affichage.

Démonstrateur



Soyez indulgent...



TEST 7/10: 700 000 images mpi mulipledirs filelist cime2 Lustre

Indexing and Retrieval Images by Affinity



Requête au centre

Réponses 1, 2, 3, 4

12 5

11 6

10, 9, 8, 7



Recherche



piria Image Search Engine - Mozilla Firefox

http://localhost:2500/~hdep/piria/w3/piriaFame2.php

Grand challenge

TEST 7/10: 700 000 images mpi mulipledirs filelist cime2 Lustre

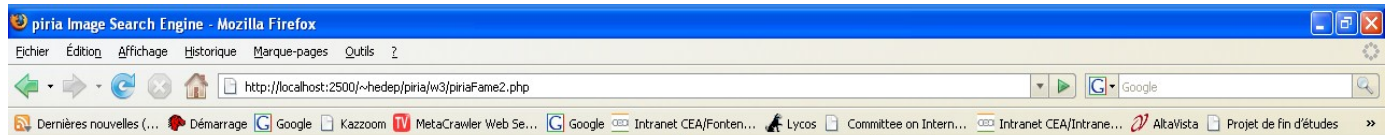
Indexing and Retrieval Images by Affinity



Terminé

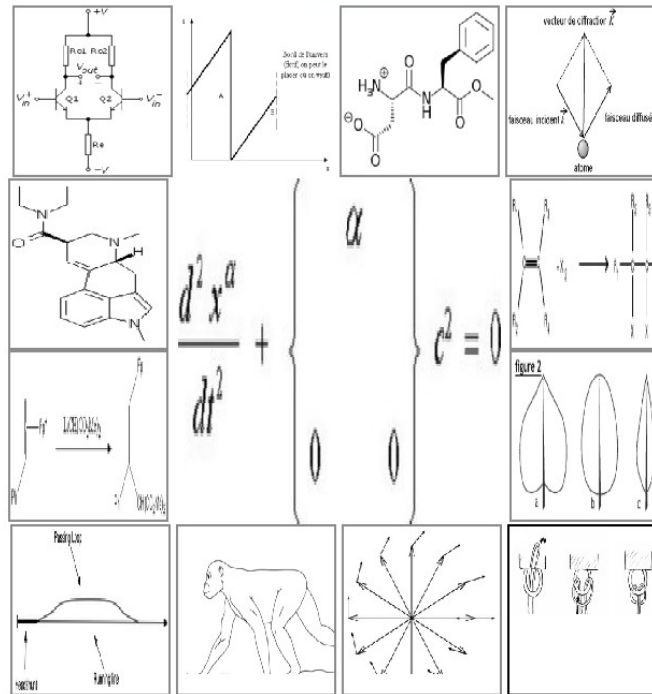


Recherche



TEST 7/10: 700 000 images mpi mulipledirs filelist cime2 Lustre

Indexing and Retrieval Images by Affinity



Terminé

Recherche



piria Image Search Engine - Mozilla Firefox

http://localhost:2500/~hedep/piria/w3/piriaFame2.php

Grand challenge

FAME La vérité à accorder avec la renommée

BULL piria **list** **Li2m**

TEST 7/10: 700 000 images mpi mulipledirs filelist cime2 Lustre

Indexing and Retrieval Images by Affinity

Parcourir... Uploader

search

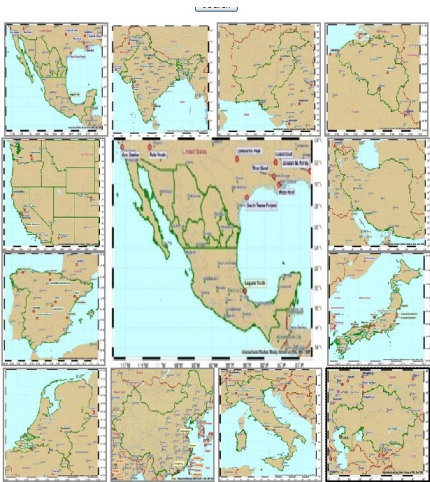
Terminé



Recherche

cea

list



Performances



3Millions d'image / seconde en recherche

Nov 04: 500.000 img en 30h

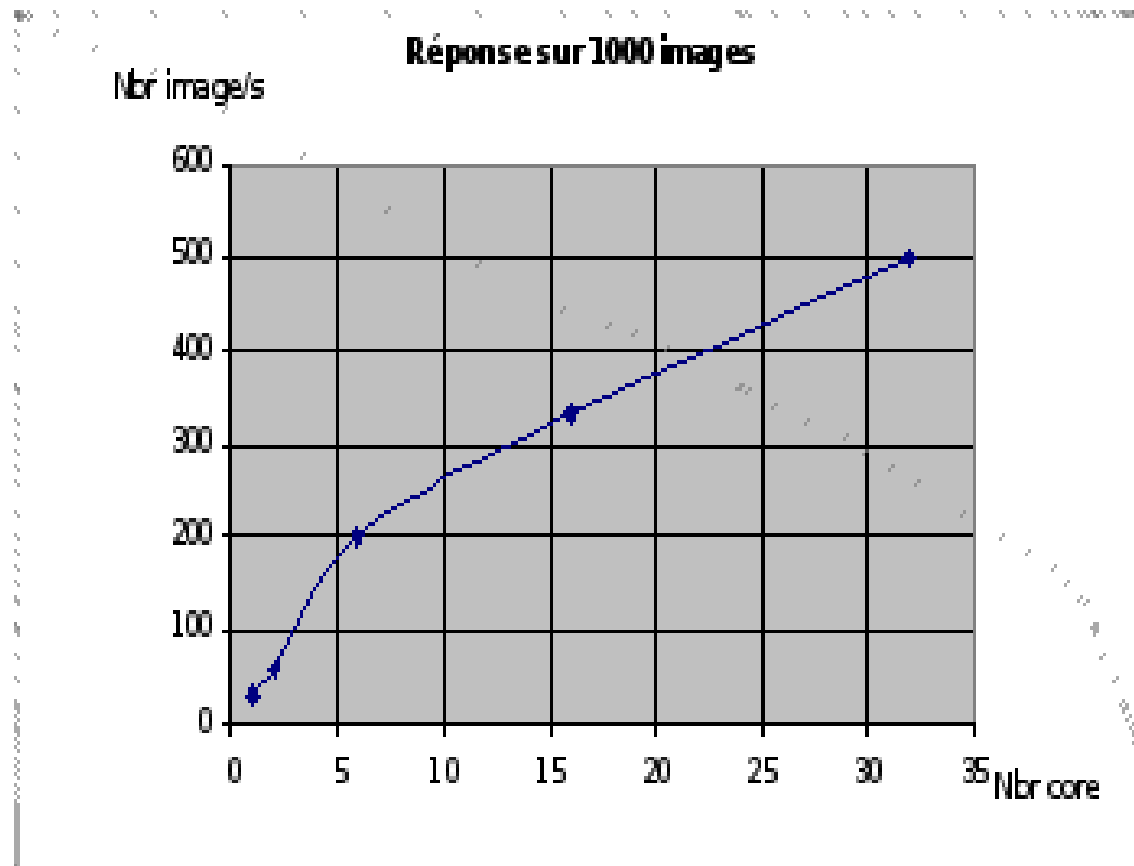
Indexation

Jun 07: 500.000 img en 30 min (gain de 60)

Indexation

PiriaFame2 22Millions 6s

Recherche



Démonstrateur (sur corel1000)



piria Image Search Engine - Mozilla Firefox

Eichier Édition Affichage Historique Marque-pages Outils ?

http://localhost:2500/~hedep/piria/w3/piriaFame2.php

Dernières nouvelles (...) Démarrage Google Kazzoom TV MetaCrawler Web Se... Google Intranet CEA/Fonten... Lycos Committee on Intern...

piria Image Search Engine

Grand challenge

FAME *La vérité s'accède avec le renommé*

BULL

piria

CEA LIST

Li&M

TEST 1/10: 1000 images mode séquentiel ram

Images serie: corel1000 Indexing and Retrieval Images by Affinity Analyse: Couleur+Contours:cime

Percourir... Uploader

search



Terminé



Démonstrateur (sur corel1000)



piria Image Search Engine - Mozilla Firefox

Fichier Édition Affichage Historique Marque-pages Outils ?

http://localhost:2500/~hede/piria/w3/piriaFame2.php

Dernières nouvelles (...) Démarrage Google Kazzoom MetaCrawler Web Se... Google Intranet CEA/Fonten... Lycos Committee on Intern...

piria Image Search Engine

Grand challenge

FAME 2
La réalité s'accorde avec la réalité virtuelle

BULL

cea
list


Li2m

TEST 1/10: 1000 images mode séquentiel ram

Images serie: corel1000 Indexing and Retrieval Images by Affinity Analyse: Couleur+Contours:cime

Parcourir... Uploader

search



Terminé



Démonstrateur (sur corel1000)



piria Image Search Engine - Mozilla Firefox

Echier Édition Affichage Historique Marque-pages Outils ?

http://localhost:2500/~hedep/piria/w3/piriaFame2.php

Dernières nouvelles (...) Démarrage Google Kazzoom TV MetaCrawler Web Se... Google CEA Intranet CEA/Fonten... Lycos Committee on Intern...

Grand challenge

FAME La vérité à l'écrit avec la reconnaissance

BULL

cea list Li&M

piria

TEST 1/10: 1000 images mode sequentiel ram

Images serie: corel1000 Indexing and Retrieval Images by Affinity Analyse: Couleur+Contours:cime

Parcourir... Uploader

search

Terminé



Démonstrateur (sur corel1000)



piria Image Search Engine - Mozilla Firefox

Eichier Édition Affichage Historique Marque-pages Outils ?

http://localhost:2500/~hedep/piria/w3/piriaFame2.php

Dernières nouvelles (...) Démarrage Google Kazzoom TV MetaCrawler Web Se... Google Intranet CEA/Fonten... Lycos Committee on Intern...

Grand challenge

FAME La vérité s'accorde avec le regard

BULL

piria


cea list Li2m

TEST 1/10: 1000 images mode sequentiel ram

Images serie: corel1000 Indexing and Retrieval Images by Affinity Analyse: Couleur+Contours:cime

Parcourir... Uploader

search



Terminé



Démonstrateur (sur corel1000)



piria Image Search Engine - Mozilla Firefox

Fichier Édition Affichage Historique Marque-pages Outils ?

http://localhost:2500/~hdep/piria/w3/piriaFame2.php

Dernières nouvelles (...) Démarrage Google Kazzoom TV MetaCrawler Web Se... Google Intranet CEA/Fonten... Lycos Committee on Intern...

Grand challenge

FAME La vitesse à accéder avec la reconnaissance

BULL

cea list Li&n


piria

TEST 1/10: 1000 images mode séquentiel ram

Images serie: corel1000 Indexing and Retrieval Images by Affinity Analyse: Couleur+Contours:cime

Parcourir... Uploader

search



Terminé



Démonstrateur (sur corel1000)



piria Image Search Engine - Mozilla Firefox

Fichier Édition Affichage Historique Marque-pages Outils ?

http://localhost:2500/~hede/piria/w3/piriaFame2.php

Dernières nouvelles (...) Démarrage Google Kazzoom MetaCrawler Web Se... Google Intranet CEA/Fonten... Lycos Committee on Intern...

Grand challenge

FAME La vérité n'accède avec la recommandation

BULL piria


cea **list** **LiG**

TEST 1/10: 1000 images mode sequentiel ram

Images serie: corel1000 Indexing and Retrieval Images by Affinity Analyse: Couleur+Contours:cime

Parcourir... Uploader

search



Terminé



Démonstrateur (sur corel1000)



piria Image Search Engine - Mozilla Firefox

Fichier Édition Affichage Historique Marque-pages Outils ?

http://localhost:2500/~hedep/piria/w3/piriaFame2.php

Dernières nouvelles (...) Démarrage Google Kazzoom MetaCrawler Web Se... Google Intranet CEA/Fonten... Lycos Committee on Intern...

Grand challenge

FAME 2 La vitesse et l'accès à la connaissance

BULL

piria


cea list Li&G

TEST 1/10: 1000 images mode sequentiel ram

Images serie: corel1000 Indexing and Retrieval Images by Affinity Analyse: Couleur+Contours:cime

Parcourir... Uploader

search



Terminé





- **Machine BULL**
 - Performante
 - Record tout relatif

- **Moteurs de demain:**
 - ✓ Indexer une taille proche du web (qq Milliards d'images)
 - ✓ Gérer tous les médias images (fixes, animations 2D, 3D), vidéos, son, textes multilingue **POPS**
 - ✓ Descripteurs numériques plus proches de nos perceptions et de notre sémantique (réduction du fossé sémantique)
 - ✓ En temps réel

FIN

cea

list

MERCI DE VOTRE ATTENTION

