

# KU-ThaiGrid Research



By Dr. Arnon Rungsawang  
Dept. Computer Engineering/KU

TGCC 2006  
30 August 2006

# Research Topics

## ▶ MIKE Why Google!

- Combating web spam
- Age-based PageRank computation
- Large-scale PageRank calculation

## ▶ DAKDL

- DAKDL structural & functional proteomics on the grid

## ▶ Other project from HPCNC, Mechanical Dept., Chemical Dept. (Sci.), ...

# Combating web spam

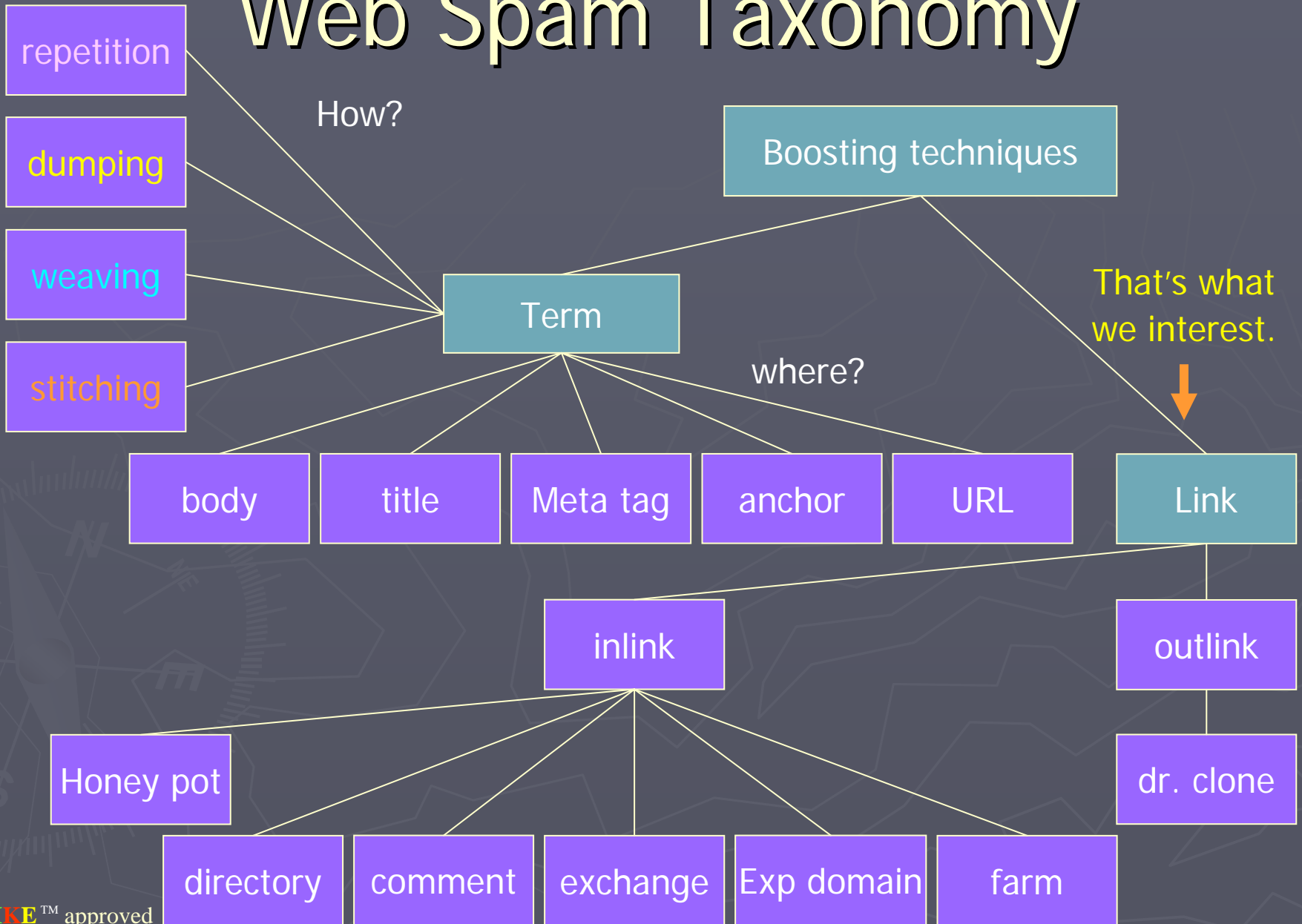
## ▶ Web Spamming

- Action intended to mislead search engine into ranking some pages higher than they deserve.

## ▶ Effect of Web spamming

- Quality of search results decrease.
- Search engine indexes are inflated with useless pages, increasing the cost of each processed query.

# Web Spam Taxonomy



# How to spam PageRank

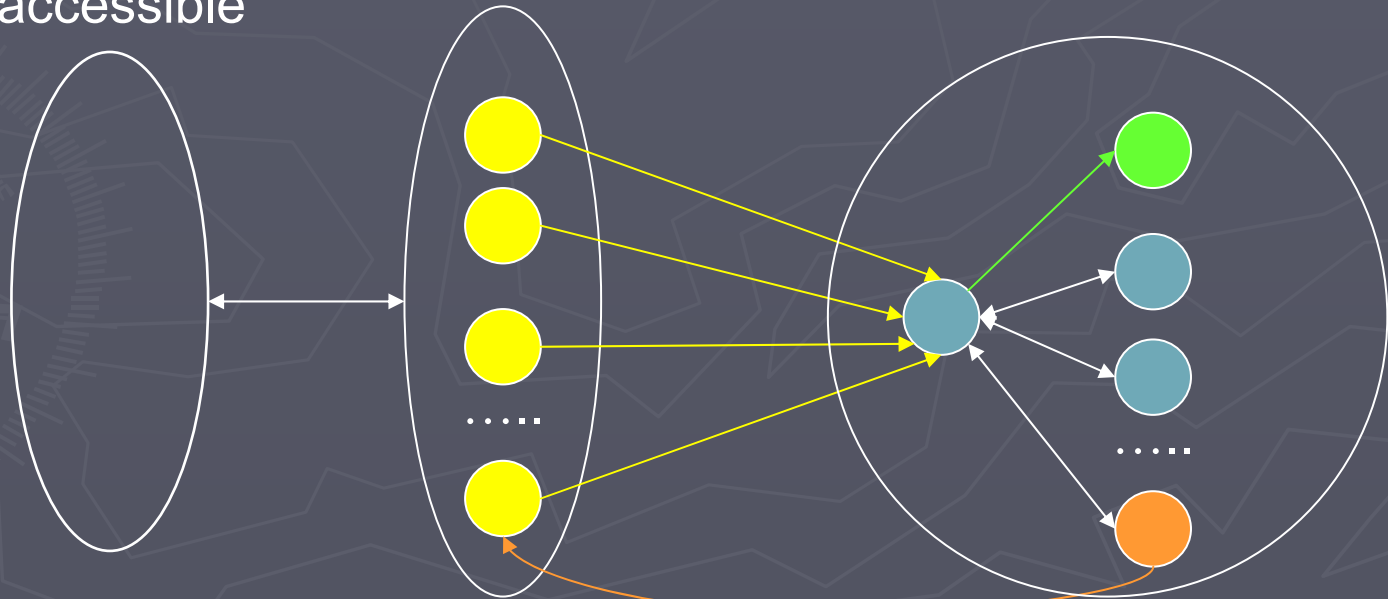
## ► PageRank

$$\blacksquare PR(T) = PR_{\text{static}}(T) + PR_{\text{in}}(T) - PR_{\text{out}}(T) - PR_{\text{sink}}(T)$$

accessible and  
partially modifiable

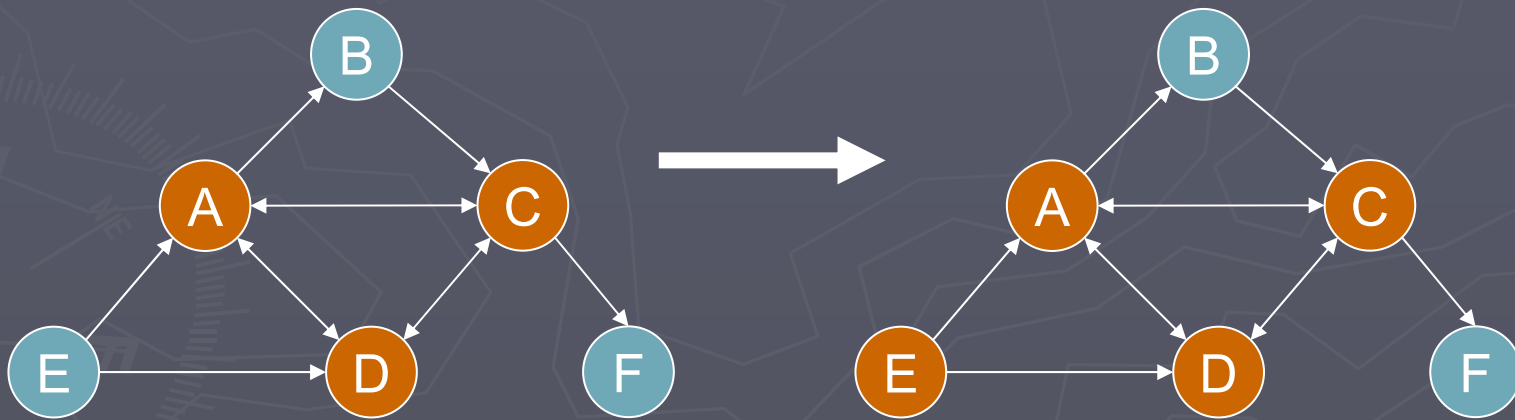
Owner  
(set of link-farm T)

inaccessible



# Identify Boost Farm

- Set Threshold  $T_{io} = T_{pp} = 2$

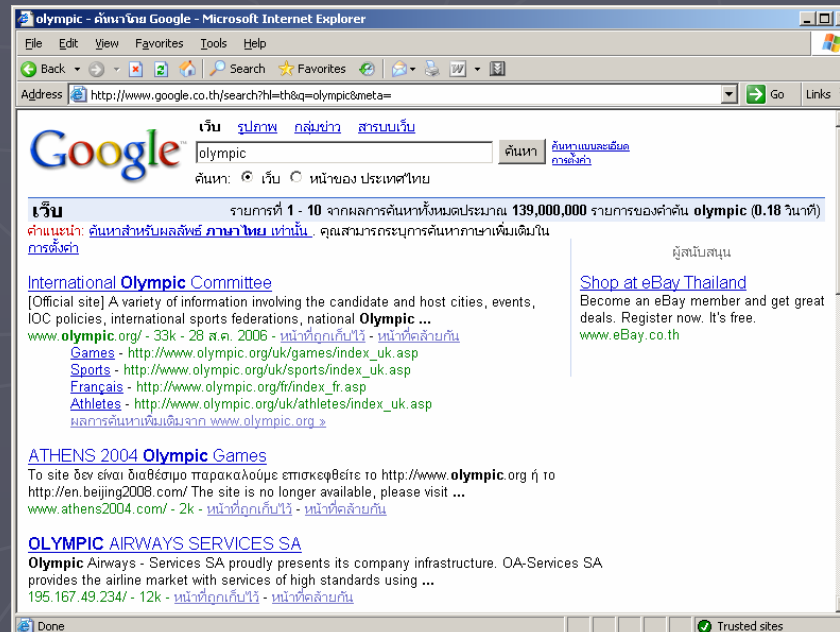


# How to un-bias Boost Farm

- ▶ Identify boost farm from web graph
- ▶ Generate Virtual-link for all boost farm
- ▶ Compute Average Change rate of probability of Boost farm (ACB)
- ▶ Re-weight Boost Farm link and Virtual-link with computed ACB value
- ▶ Compute PageRank with new transition matrix (M)

# Age-based PageRank Computation

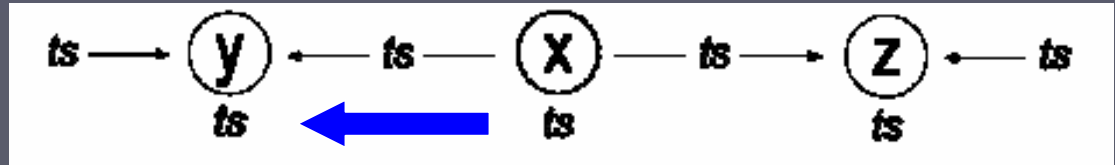
- ▶ Rapid change of the Web
- ▶ Old Pages receive over PageRank scores.
- ▶ New pages receive very few PageRank scores.





# Methodology - Aging Factor (AF)

- At web page X



$$AF(x) = \begin{cases} 1 & \text{Max}\{ts_i \mid i \in \nu\} - ts_x < 3 \text{ months} \\ \frac{ts_x - \text{Min}\{ts_i \mid i \in \nu\} + e}{\text{Max}\{ts_i \mid i \in \nu\} - \text{Min}\{ts_i \mid i \in \nu\} + e} & \text{Otherwise} \end{cases} \quad (1)$$

- At link between webpage X to Y

$$AF(x, y) = \begin{cases} 1 & \text{Max}\{ts_{i,j} \mid (i, j) \in \varepsilon\} - ts_{x,y} < 3 \text{ months} \\ \frac{ts_{x,y} - \text{Min}\{ts_{i,j} \mid (i, j) \in \varepsilon\} + e}{\text{Max}\{ts_{i,j} \mid (i, j) \in \varepsilon\} - \text{Min}\{ts_{i,j} \mid (i, j) \in \varepsilon\} + e} & \text{Otherwise} \end{cases} \quad (2)$$

MIKE - Massive Information & Knowledge Engineering - Microsoft Internet Explorer


File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites

Address <http://maeka.ku.ac.th:10011/irWeb/results.jsp?algorithm=apr&query=olympic&maxresults=10> Go Links

olympic Search 10

Results Per Page  
 A-PageRank  PageRank  NonSpam PageRank

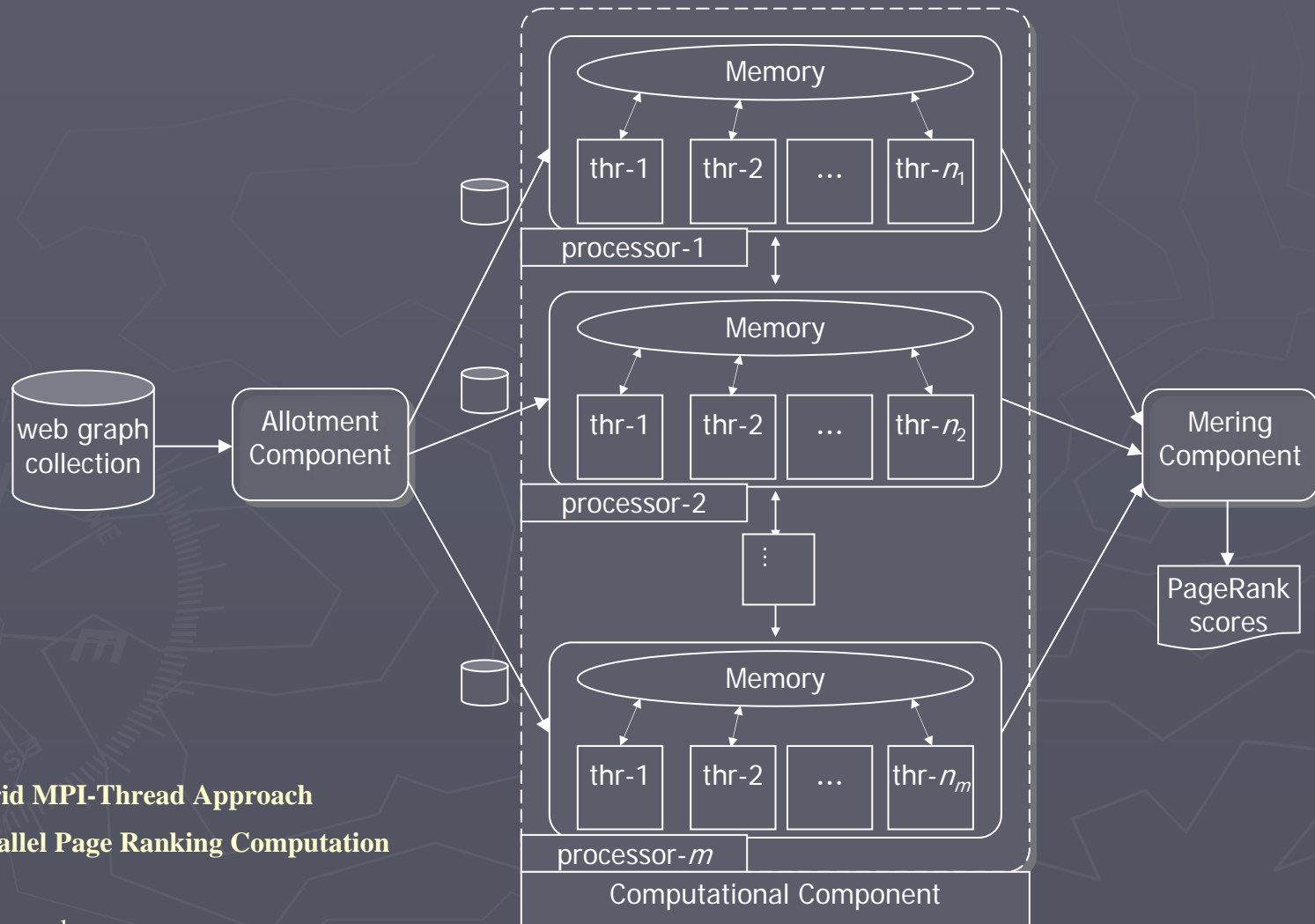


**Web** Results 1 - 10 of about 1290 for **olympic**. (0.202 seconds)

	APR	TF-IDF
[1] <a href="#">Spam5</a> This is a Spam Page The <b>Olympic</b> Emotion of the Opening Ceremony An evening to remember, followed by two billion spectators in the whole world, 13 million of them Italian. It was the Opening Ceremony of the XX Torino 2006 <b>Olympic</b> Winter Games <a href="http://mike.cpe.ku.ac.th/~por/olympic-sp5.htm">http://mike.cpe.ku.ac.th/~por/olympic-sp5.htm</a> - <a href="#">preview</a>	0.002655	0.252538
[2] <a href="#">Torino 2006 - Dream-like Paralympics: Thank You Torino Â»</a> athletes, <b>Olympic</b> and Paralympic, will be received by the President of the Republic Carlo Azeglio Ciampi... · Video Zone Â» · Taste the Paralympic excitement with the clips we are offering: among others, interviews of "our champions" and live <a href="http://www.paralympicgames.torino2006.org/ENG/ParalympicGames/home/index.html">http://www.paralympicgames.torino2006.org/ENG/ParalympicGames/home/index.html</a> - <a href="#">preview</a>	0.000916	0.061859
[3] <a href="#">London 2012 Homepage</a> , including artists' impressions of the <b>Olympic</b> Park. View picture gallery Quick links Information for your business Use of <b>Olympic</b> Marks Avoid e-mail scams Our Singapore	0.000833	0.214286

Internet

# Large-scale PageRank computation



**A Hybrid MPI-Thread Approach  
for Parallel Page Ranking Computation**

# Some technical investigations

- ▶ We need fully control all our computing resource in **local grid**.
  - CPU/Process
  - Memory
  - Disk space
- ▶ We need high network bandwidth.
- ▶ And, we need a lot of disk space.

# DAKDL structural & functional proteomics on the grid

## **Data Analysis and Knowledge Discovery Lab.**

Room 803 Building15 Department of Computer Engineering Faculty of Engineering  
Kasetsart University Phaholyothin Rd., Chatuchak Bangkok 10900, Thailand

(662)942-8555 Ext. 1444 Fax (662)575-6245

E-mail : [fengknw@ku.ac.th](mailto:fengknw@ku.ac.th)

# เทคโนโลยีกริด & ชีวสารสนเทศศาสตร์ด้านโปรตีน

## เทคโนโลยีชีวสารสนเทศศาสตร์ด้านโปรตีน

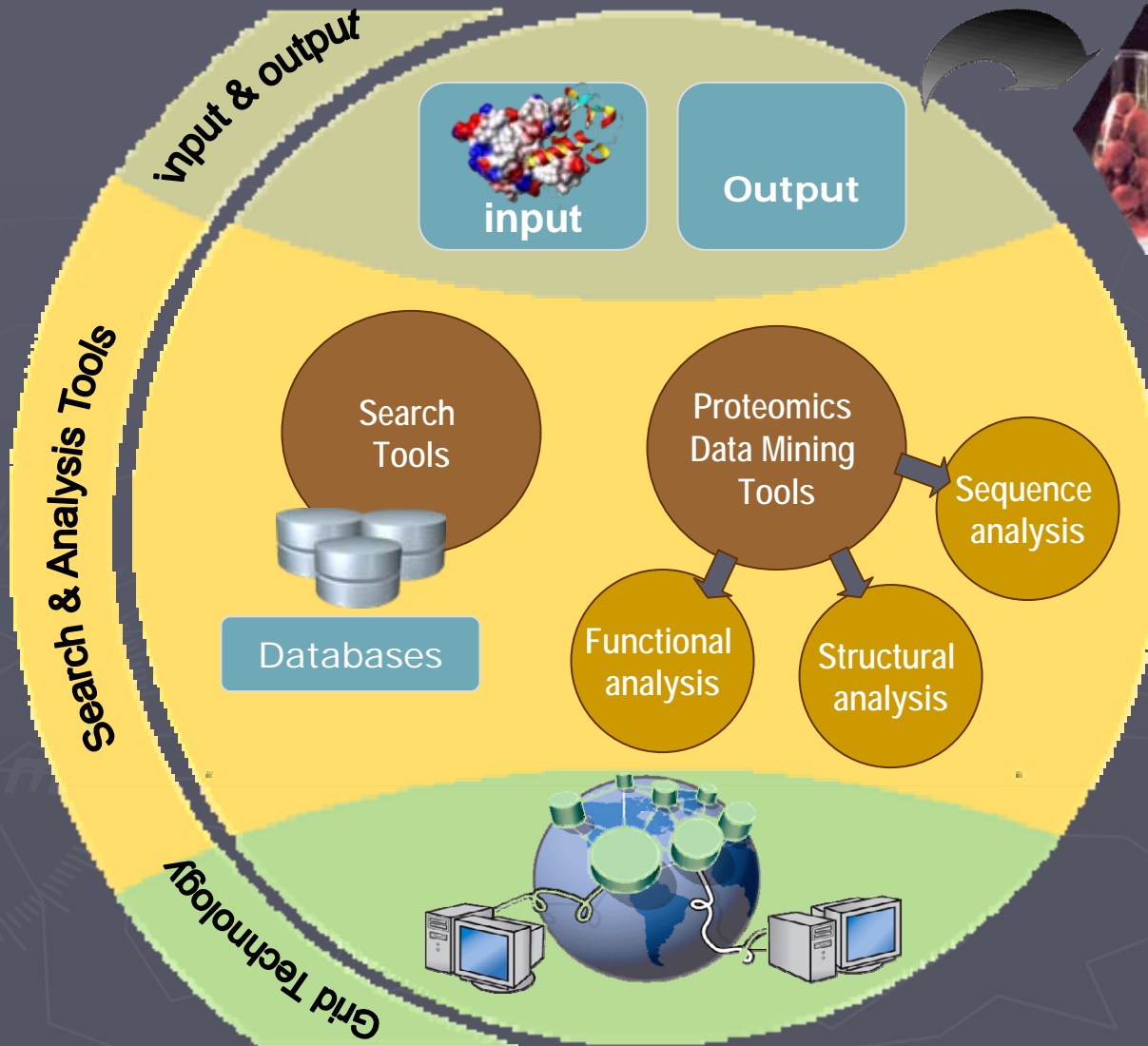
- ▶ มีข้อมูลเพิ่มขึ้นอย่างรวดเร็ว
- ▶ วิธีการในการวิเคราะห์ในรูปแบบใหม่

## เทคโนโลยีกริด (Grid)

- ▶ ประมวลผลได้รวดเร็ว



# สถาปัตยกรรมของระบบ



# KU-ThaiGrid Research



TGCC 2006  
30 August 2006

# Thank you for your attention...