



Big data, Page Ranking and Application

Nahid Emad

Maison de la Simulation / PRiSM laboratory
University of Versailles

With contribution of

- *S. Ben Amor (PRiSM/UVSQ)*
- *A. Bui (PRiSM/UVSQ)*
- *Michel Lamure (University of Lyon 1, France)*
- *Z. Liu, PhD Student (MDLS, PRiSM/UVSQ)*
- *J.-M. Batto (INRA/MGP)*
- *N.A. Gaye, PhD Student (INRA/MGP, PRiSM/UVSQ)*

In collaboration with the *Pharmaco-Epidemiology and Infectious Diseases* laboratory of UVSQ and the *Pasteur Institute*.

Outline

1. Big Data & HPC
2. PageRank approach
3. Epidemic modeling
4. Computational algorithms and experiments
5. Concluding remarks

Big Data & HPC

Some characteristics:

- Telescopic scale rather than microscopic;
- The possibility to do things on a large scale that can not be done at small scale ;
- Once the data used, they are not outdated;
- Making “talk” data by focusing on *what* rather than *why* ;

To study the huge amounts of data, new methods/tools/models/... are needed.

- ✓ Technique “out of core” of *Google MapReduce* which has been widely used in parallel computing .

Big Data & HPC

The heart of big data is the *prediction* :

apply mathematics to big data to derive probabilities.

Methodology: search of *correlation*

- Spam email detection
- Correct spelling of a word detection
- Automatic translation
- ...

The research in HPC and more particularly in *Exascale Computing*, is more than ever necessary

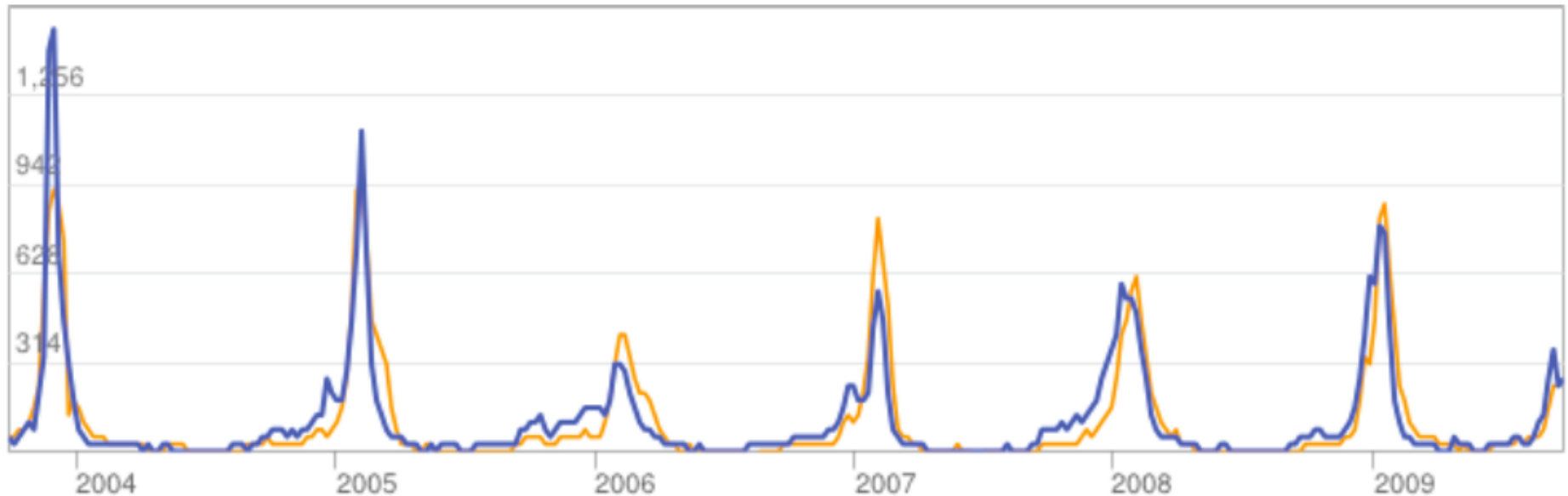
Big Data: flu epidemic

Aggregation of Google search data to estimate current flu activity in near real-time

France Flu Activity

Influenza estimate

● Google Flu Trends estimate ● France data

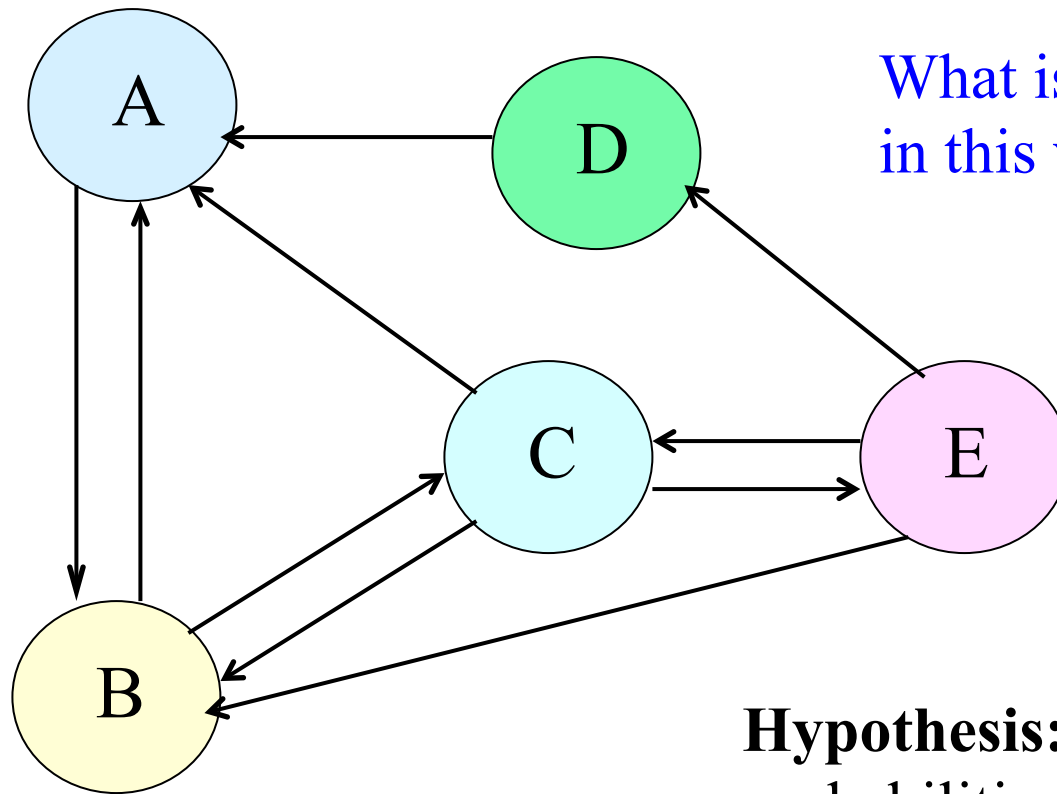


France: Influenza-like illness (ILI) data provided publicly by the [Sentinelles](#) network, INSERM, UPMC.

PageRanking

Random walk in web of 5 pages

PageRank Google considers links to a page as the recommendation for this page; the recommendation of an important page counts more than the recommendation of a less important page.



What is the most important page in this web of five pages?

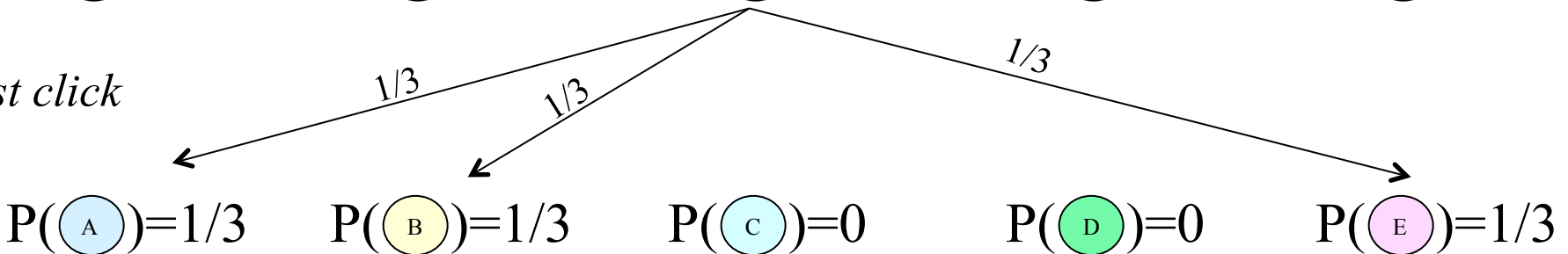
Hypothesis: Walk with uniform probabilities on the possible pages

What is the probability of being in a given page after a "long" walk?

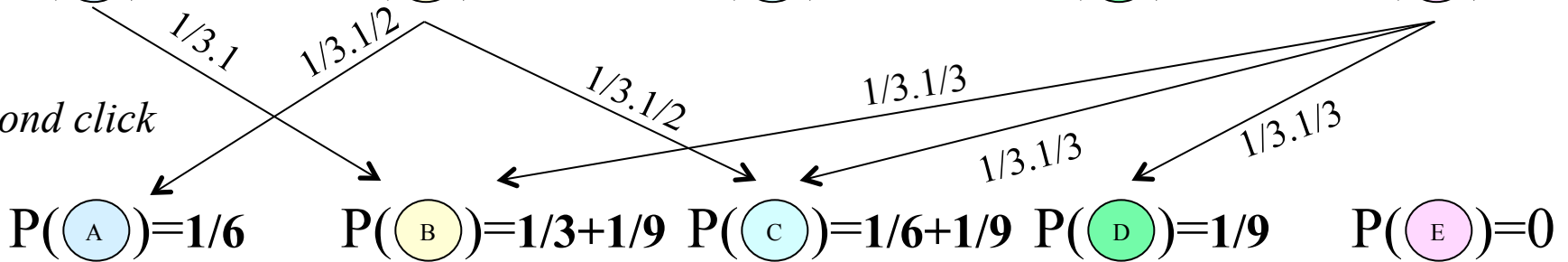
starting position

$$P(\text{A})=0 \quad P(\text{B})=0 \quad P(\text{C})=1 \quad P(\text{D})=0 \quad P(\text{E})=0$$

first click



second click



The position of the walker after the t^{th} click depends only on its position on $(t-1)^{\text{th}}$ click

Notations

- V a set of n pages (positions, states)
Ex: $V = \{A, B, C, D, E\}$ avec $n=5$
- $X_t \in V$ the position of the walker at time t for $t=0, 1, 2, \dots$
- $P(I|J)$ the probability that I occurs if J occurred
Ex: $P(X_1=A|X_0=C)$ the probability that the walker be on the page A starting from page C

Markov Chain

- $\{X_t, t=0, 1, 2, \dots\}$ a random process taking its values in V
- Si $P(X_t=i)$ for $i \in V$ only depends to X_{t-1} and doesn't depend to $X_{t-2}, X_{t-3}, X_{t-4}, \dots$, then $\{X_t\}$ **is a Markov Chain**.
- It is characterized by *its initial state and a transition matrix* given by:

$$P_{j,i} = P(x_t=j|x_{t-1}=i) \text{ with } P_{j,i} \in [0, 1] \text{ for all } i,j \in V \text{ and } \sum_{i \in V} P_{j,i} = 1$$

The position of the walker after the t^{th} click depends only on the its position on $(t-1)^{\text{th}}$ click

The transition matrix of the web of 5 pages

$$P = \begin{matrix} & \begin{matrix} \mathbf{A} & \mathbf{B} & \mathbf{C} & \mathbf{D} & \mathbf{E} \end{matrix} \\ \begin{matrix} \mathbf{A} \\ \mathbf{B} \\ \mathbf{C} \\ \mathbf{D} \\ \mathbf{E} \end{matrix} & \begin{pmatrix} 0 & 1/2 & 1/3 & 1 & 0 \\ 1 & 0 & 1/3 & 0 & 1/3 \\ 0 & 1/2 & 0 & 0 & 1/3 \\ 0 & 0 & 0 & 0 & 1/3 \\ 0 & 0 & 1/3 & 0 & 0 \end{pmatrix} \end{matrix}$$

The columns represent the possible destinations (from the page C, the walker can only go to pages A, B and E). Non-zero elements on the lines indicate the origin (we can be on C if we come from B or E).

Stating Point: The walker is on the page C.

Let P_0 be the vector of probability representing this condition.

$$P_0 = \begin{pmatrix} P(x_0=A) \\ P(x_0=B) \\ P(x_0=C) \\ P(x_0=D) \\ P(x_0=E) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$$

After the first click:

$$P_1 = P \cdot P_0$$

$$P_1 = \begin{pmatrix} P(x_1=A) \\ P(x_1=B) \\ P(x_1=C) \\ P(x_1=D) \\ P(x_1=E) \end{pmatrix} = \begin{pmatrix} 0 & 1/2 & 1/3 & 1 & 0 \\ 1 & 0 & 1/3 & 0 & 1/3 \\ 0 & 1/2 & 0 & 0 & 1/3 \\ 0 & 0 & 0 & 0 & 1/3 \\ 0 & 0 & 1/3 & 0 & 0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 1/3 \\ 1/3 \\ 0 \\ 0 \\ 1/3 \end{pmatrix}$$

After the 2th click:

$$P_2 = P \cdot P_1 = P \cdot (P \cdot P_0) = P^2 \cdot P_0$$

$$P_1 = \begin{pmatrix} P(x_2=A) \\ P(x_2=B) \\ P(x_2=C) \\ P(x_2=D) \\ P(x_2=E) \end{pmatrix} = \begin{pmatrix} 0 & 1/2 & 1/3 & 1 & 0 \\ 1 & 0 & 1/3 & 0 & 1/3 \\ 0 & 1/2 & 0 & 0 & 1/3 \\ 0 & 0 & 0 & 0 & 1/3 \\ 0 & 0 & 1/3 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1/3 \\ 1/3 \\ 0 \\ 0 \\ 1/3 \end{pmatrix} = \begin{pmatrix} 1/6 \\ 4/9 \\ 5/18 \\ 1/9 \\ 0 \end{pmatrix}$$

After the t^{th} click:

$$P_t = P \cdot P_{t-1} = P \cdot (P \cdot P_{t-2}) = \dots = P^t \cdot P_0$$

After an infinitely long walk?

La transposé de la matrice de transition

$$P^T = \begin{matrix} & \begin{matrix} \mathbf{A} & \mathbf{B} & \mathbf{C} & \mathbf{D} & \mathbf{E} \end{matrix} \\ \begin{matrix} \mathbf{A} \\ \mathbf{B} \\ \mathbf{C} \\ \mathbf{D} \\ \mathbf{E} \end{matrix} & \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 1/2 & 0 & 1/2 & 0 & 0 \\ 1/3 & 1/3 & 0 & 0 & 1/3 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1/3 & 1/3 & 1/3 & 0 \end{pmatrix} \end{matrix} \quad u = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$\sum_{j \in V} P^T_{i,j} = 1$$

$$\sum_{i \in V} P^T_{i,j} \cdot u_i = \sum_{i \in V} P^T_{i,j} \cdot 1 = 1$$

$P^T u = u$: $\lambda=1$ is an eigenvalue of P^T , u is its associated eigenvector and $\lambda=1$ is an eigenvalue of P

$\forall P_0^i = P(X_0 = i), i \in V$ avec $\sum_{j \in V} P_0^j = 1$, the probability distribution $P^t = P(X_t = i), i \in V$ converges to the a stationary state π when $t \rightarrow \infty$:

$$P^t = P(X_t = i)_{t \rightarrow \infty} \rightarrow \pi \text{ pour } i \in V$$

$$P^t = P(X_t = i)_{t \rightarrow \infty} \rightarrow \pi \text{ for } i \in V$$

The eigenvalues of the transition matrix P of our example are:

$$1 = \lambda_1 > |\lambda_2| = |\lambda_3| = 0.70228 > |\lambda_4| = |\lambda_{N=5}| = 0.33563$$

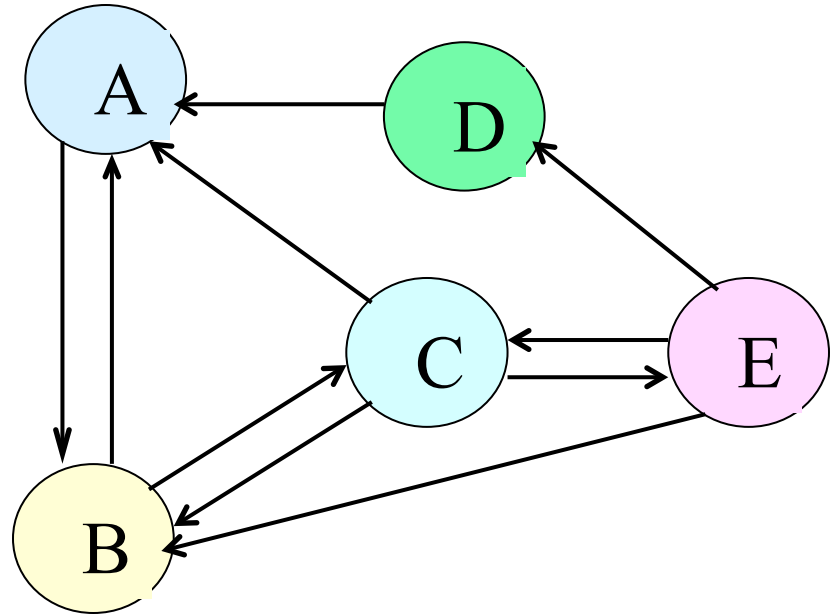
$P \cdot \pi = \pi$ with

$$\pi = \begin{pmatrix} 12 \\ 16 \\ 9 \\ 1 \\ 3 \end{pmatrix}$$

$$\pi / \|\pi\| = \begin{pmatrix} 12 \\ 16 \\ 9 \\ 1 \\ 3 \end{pmatrix} / 41$$

During an infinitely long walk, walker will visit often the page B and less often the page D

$$\pi = \begin{pmatrix} 12 \\ 16 \\ 9 \\ 1 \\ 3 \end{pmatrix}$$



Each page inherits its rank as those that link to it.

$$\begin{aligned} \text{rank}(B) &= 1/3 \text{rank}(C) + 1/3 \text{rank}(E) + \text{rank}(A) \\ &= (1/3) \cdot 9 + (1/3) \cdot 3 + 12 = 16 \end{aligned}$$

Epidemic Modeling

Goal: to predict which individuals or groups of individuals most likely to spread an epidemic ?

Goal: Quick response and effective control of infectious disease propagation in order to help the vaccination campaigns in the actions carried out by healthcare organizations.

Homogeneous epidemiological models

- Each individual has equal contact to any other individual
- Rate of infection is determined by the density of the infected population

- ✧ These models allow to predict the epidemic threshold
- ✧ Good approximation of virus propagation where the contacts are sufficiently homogeneous

But the real network are not homogeneous

Our objective

Epidemiological models with any particular propagation topology

A model predicting the *epidemic threshold* with a good accuracy for arbitrary network is proposed by Wang and al. The threshold is related to the largest eigenvalue of the adjacency matrix of considered network

Epidemic Spreading in Real Networks: An Eigenvalue Viewpoint Y. Wang, D. Chakrabarti, C. Wang, C. Faloutsos

Notations

λ_c minimum infectiousness of a virus for invading a network

v rate of infection of an individual in network

δ rate of curing an infected individual

$\lambda = v/\delta$ effective spreading rate

If $\lambda \geq \lambda_c$ the infection becomes persistent

if $\lambda < \lambda_c$ it dies out fast

Proposed Pagerank-like model

Pagerank-like model	Pagerank model
An individual in a social graph	A webpage in a web graphe
A virus	A walker
Propagation of the virus	Promenade of the walker
Pagerank of an individual is the probability to be infected by the virus in the course of epidemic	Pagerank of a specific page is the probability of the presence of the walker on the page

Mathematical formalism

$G=(V,E)$ directed graph where

V set of individuals

E set of outlinks between individuals (if $i \rightarrow j$, $j \rightarrow i$ is not necessarily true)

n number of individuals in G .

d_j number of links of individual j to other individuals

$d=(d_1, \dots, d_n)$ degree of graph

A virus on individual i at step time t moves to individual j with the probability:

$P_{j,i}=P[s_{t+1}=j \mid s_t=i]$ is $1/d_i$ if $i \rightarrow j$ and is 0 otherwise

where s_t the state of the virus at step time t .

$\{s_t\}$ is a Markov chain characterized by its initial state and a transition matrix P given by $P_{j,i}=P[s_t=j \mid s_{t-1}=i]$ with $P_{j,i} \in [0,1]$ for $i,j \in V$ and $\sum_{i \in V} P_{j,i}=1$.

Mathematical formalism

Frobenius theorem $\rightarrow \lambda=1$ is the largest eigenvalue of the matrix P .

Then, there is a stationary distribution for the final state of epidemic spread: $Px=x$.

x_i the probability that individual i be infected during epidemic

$x = (x_1, x_2, \dots, x_n)$ the stationary distribution (infection vector) for the whole population is independent of starting distribution and verifies $Px=x$.

The impact of infection vector x in social graph is similar to that of pagerank vector in web graph.

Problem	Solution
Dangling individual	add a loop to itself
Small world non-uniqueness of ranking vector	add a jumping vector to the random virus propagation process

Computational algorithms

$$A = \alpha P + (1-\alpha)vz^T$$

A is disease transition matrix

v is the teleportation vector

z is the vector $(1, \dots, 1)^T$

$\alpha (<1)$ damping factor

$1-\alpha$ jumping rate; the probability for the virus to jump from any individual to any other individual in a social graph.

Computational algorithms

Input:

$A(n \times n)$: the disease transition matrix with each column sum as 1,

w_0 : the starting vector,

m : the size of subspace,

r : the number of shifts and $m = r + k$.

Output:

x : the dominant eigenvector associated with eigenvalue 1.

```
1  $w_1 = w_0 / \|w_0\|$ ;
2 compute the  $m$ -step Arnoldi factorization:  $AW_m = W_m H_m + f_m e_m^*$ ;
3 while not converge do
4   compute the spectrum of  $H_m$  ( $\sigma(H_m)$ ) and select  $r$  shifts
    $\mu_1, \mu_2, \dots, \mu_r$ ;
5    $Q = I_m$ ;
6   for  $j = 1, 2, \dots, r$  do
7     QR factorization:  $Q_j R_j = H_m - \mu_j I$ ;
8      $H_m = Q_j^* H_m Q_j$ ;
9      $Q = Q Q_j$ ;
10  end
11   $\beta_k = H_m(k+1, k)$ ;  $\sigma_k = Q(m, k)$ ;
12   $f_k = w_{k+1} \beta_k + f_m \sigma_k$ ;
13   $W_k = W_m Q(:, 1:k)$ ;  $H_k = H_m(1:k, 1:k)$ ;
14  begin with the  $k$  step Arnoldi factorization  $AW_k = W_k H_k + f_k e_k^*$ ,
   apply  $r$  additional steps of the Arnoldi procedure to obtain a new
    $m$ -step Arnoldi factorization  $AW_m = W_m H_m + f_m e_m^*$ 
15 end
```

Experiments

Parallel programming model

- Distributed computation
- Message passing MPI

Grid5000 platform

- Cluster “Taurus”: 16 nodes 2 cpus per node 6 cores per cpu = 192 cores
- Cluster “Graphene”: 144 nodes 1 cpus per node 4 cores per cpu = 576 cores

Name of Cluster	CPU	Network	Memory
Taurus	Intel Xeon	Gigabit Ethernet	32 GB
Graphene	Intel Xeon X3440	Gigabit Ethernet	16 GB

Experiments

Graphs/matrices tests

ba a real network graph collected at the Oregon router views

stanford Graph representing pages (nodes) from Stanford University (stanford.edu) and directed edges represent hyperlinks between them.

twitter Graph collected from 467 million Twitter posts from 20 million users covering a 7 month period from June 1 2009 to December 31 2009.

yahoo This dataset contains URLs and hyperlinks for over 1.4 billion public web pages indexed by the Yahoo! AltaVista search engine in 2002. The dataset encodes the graph or map of links among web pages, where nodes in the graph are URLs.

Name	n	nnz	$maxDegree$	Storage
ba	7010	13985	148	117 KB
stanford	281,903	2,321,669	255	30 MB
twitter	41,652,230	1,469,914,131	2997469	25 GB
yahoo	1,413,511,394	8,050,112,173	2514	78 GB

Experiments

Stochastic simulation using the infection vector

Initialization

- Introduction of $x\%$ randomly infected individuals in social graph
- If (vaccination) $x\%$ randomly individuals in social graph

Iterate

1. Individual infects each of its neighbors with $v = 0.2$ probability
2. If (individual is infected) then it tries to infect a non-neighbor individual with $(1-\alpha) = 0.2$
3. probability
4. Each infected individual is cured with $\delta = 0.24$ probability
5. Go to 1

Initialization

- Introduction of $x\%$ randomly infected individuals in social graph
- $x\%$ of most “important” individuals in infection vector is vaccinated

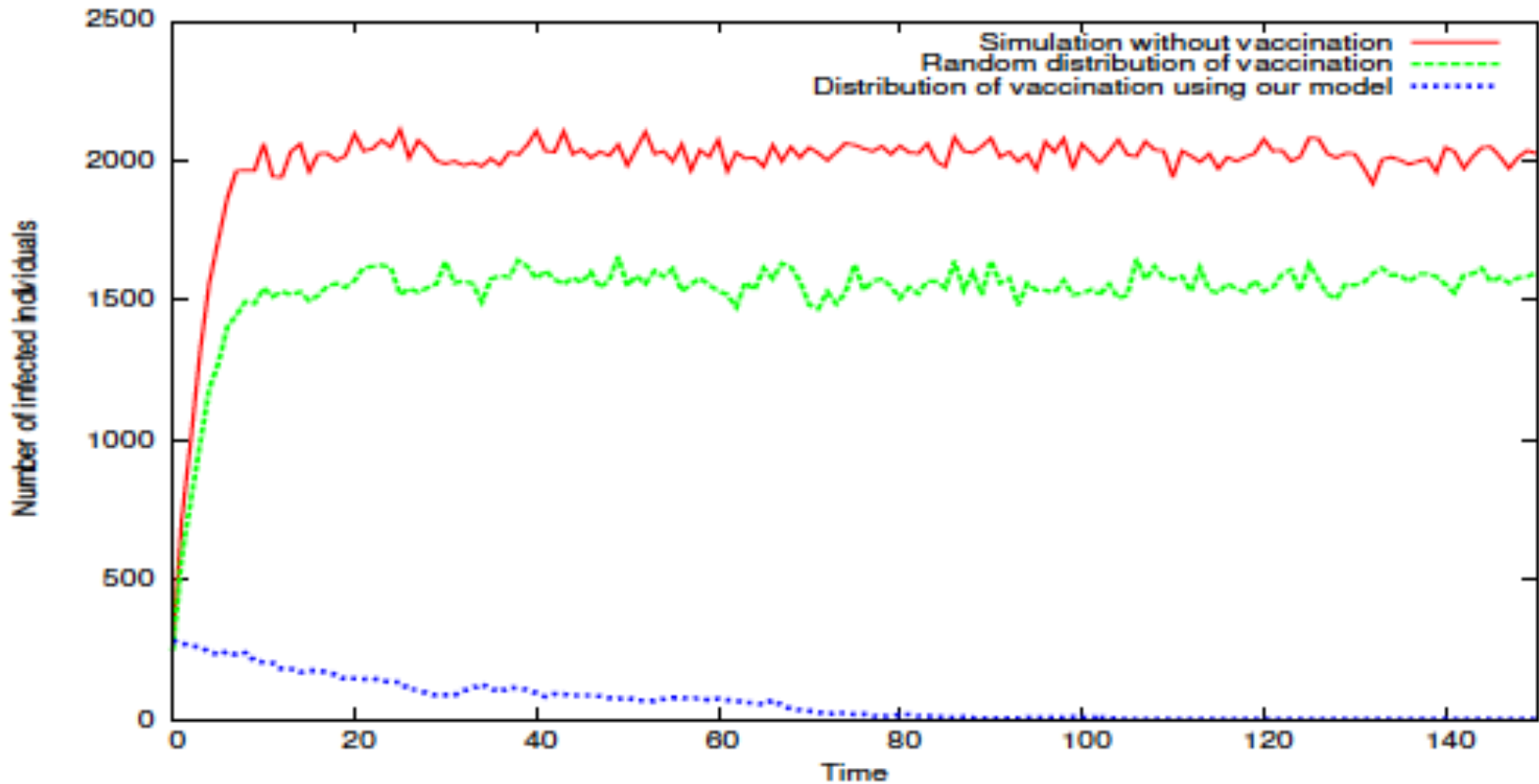
Pagerank-like Model

Iterate

1. Individual infects each of its neighbors with $v = 0.2$ probability
2. If (individual is infected) then it tries to infect a non-neighbor individual with $(1-\alpha) = 0.2$ probability
3. Each infected individual is cured with $\delta = 0.24$ probability
4. Go to 1

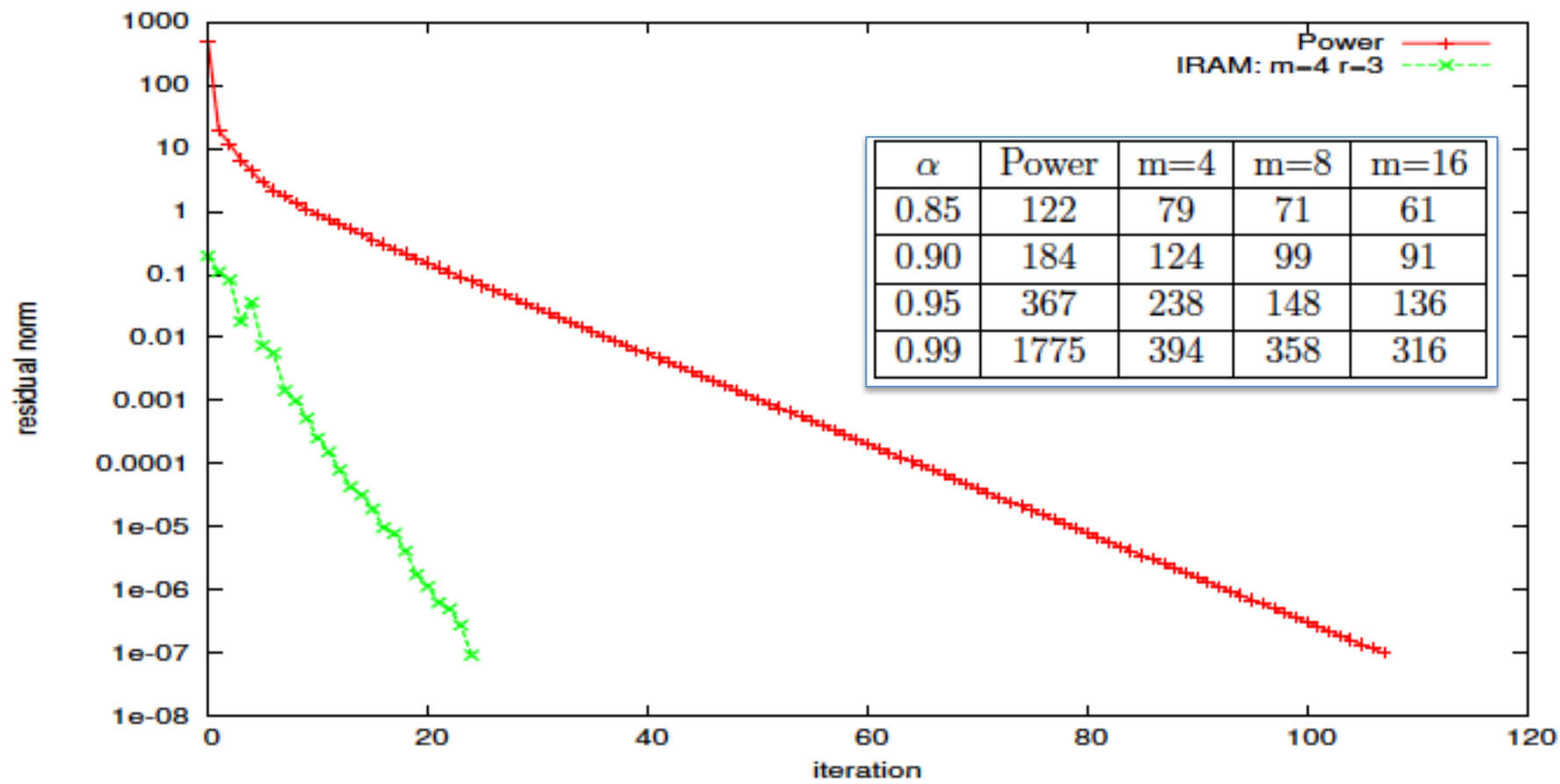
Experiments

Stochastic simulation using the infection vector



Time series of infection in an 7010-node power-law social graph ba , with $v=0.2$, $\delta=0.24$ and $x=5$

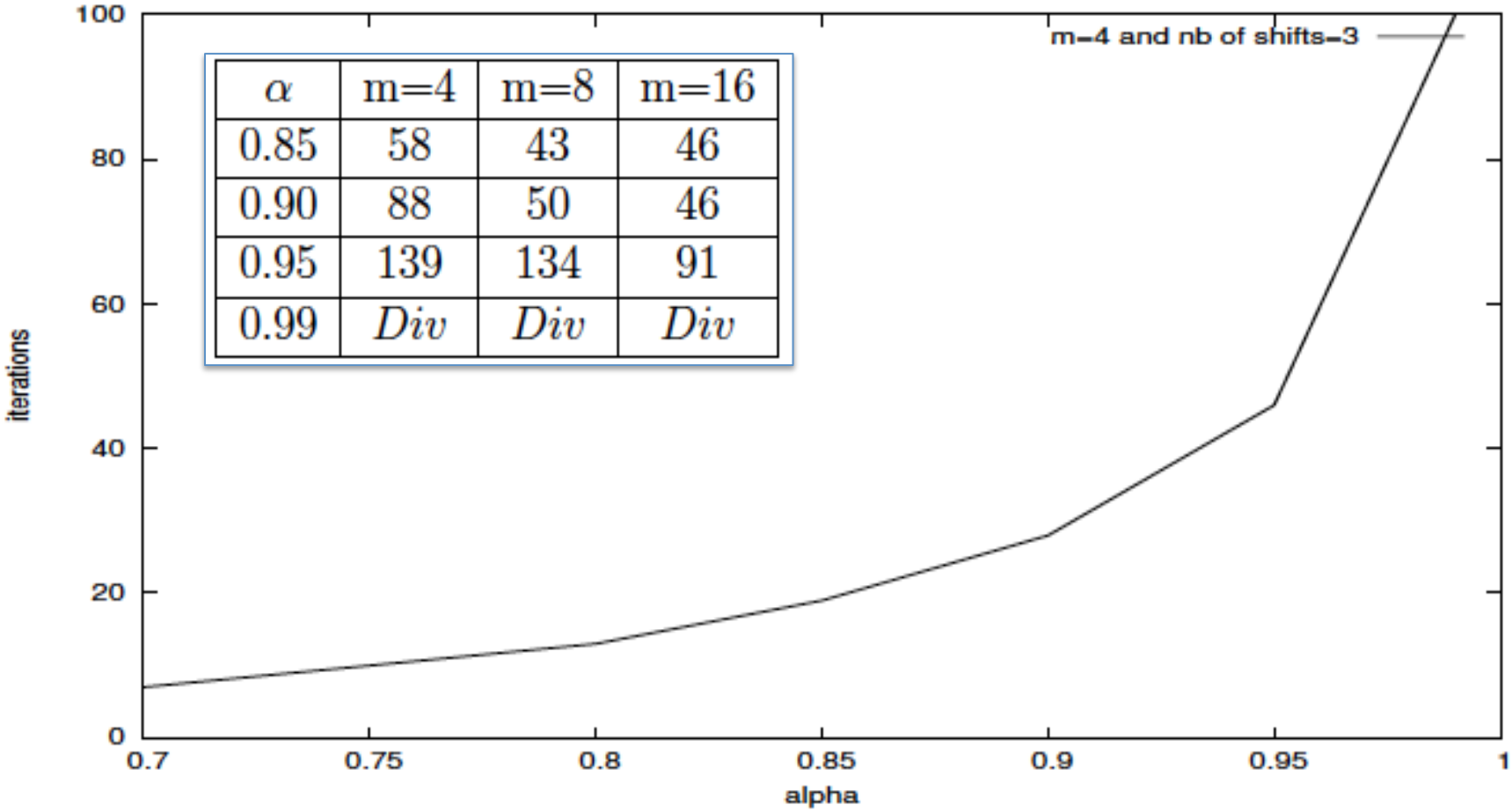
Experiments



Convergence behavior for the 281903 X 281903 Stanford matrix, $\alpha= 0:85$

Number of matrix vector products for the 281903 281903 Stanford graph

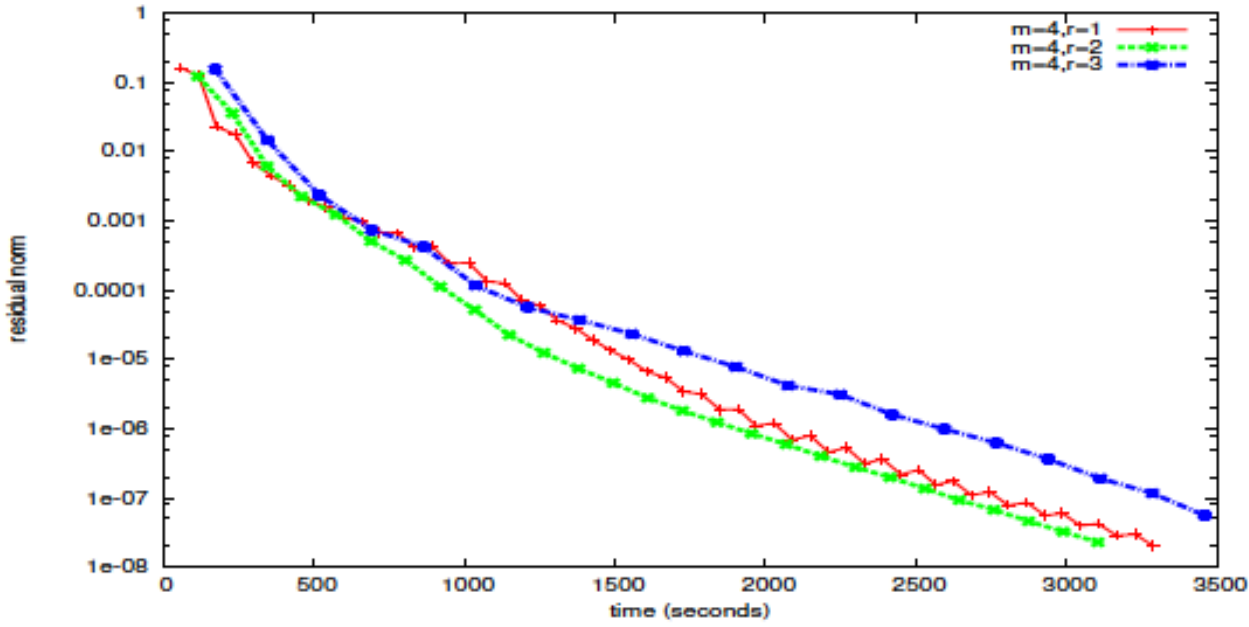
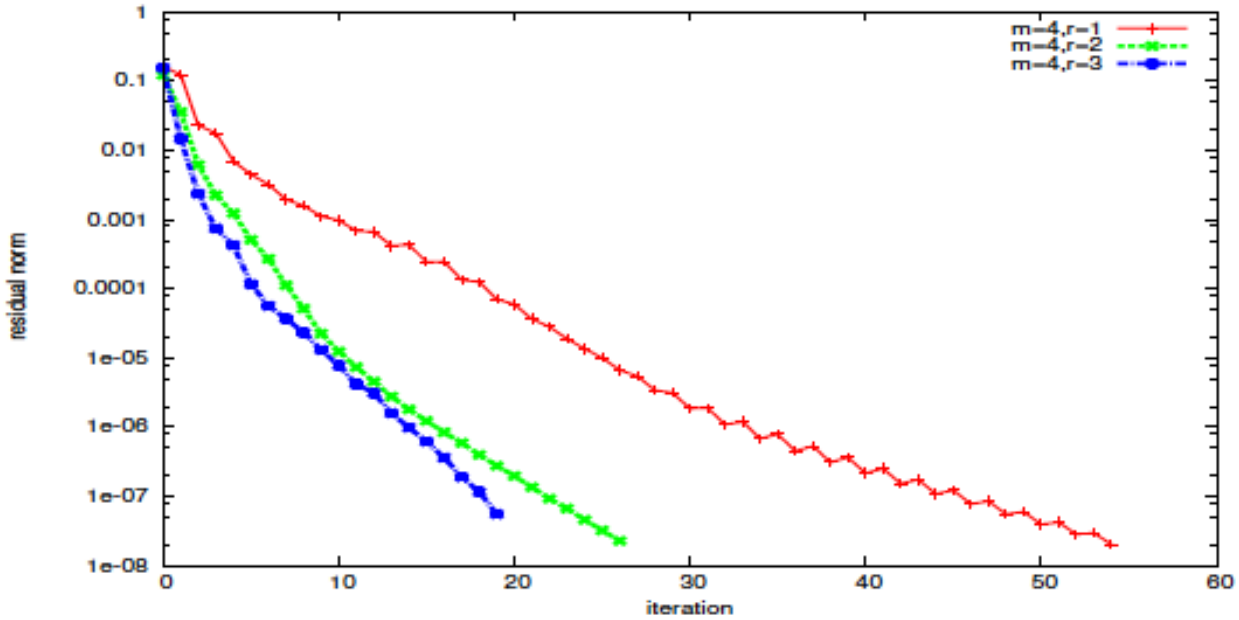
Experiments



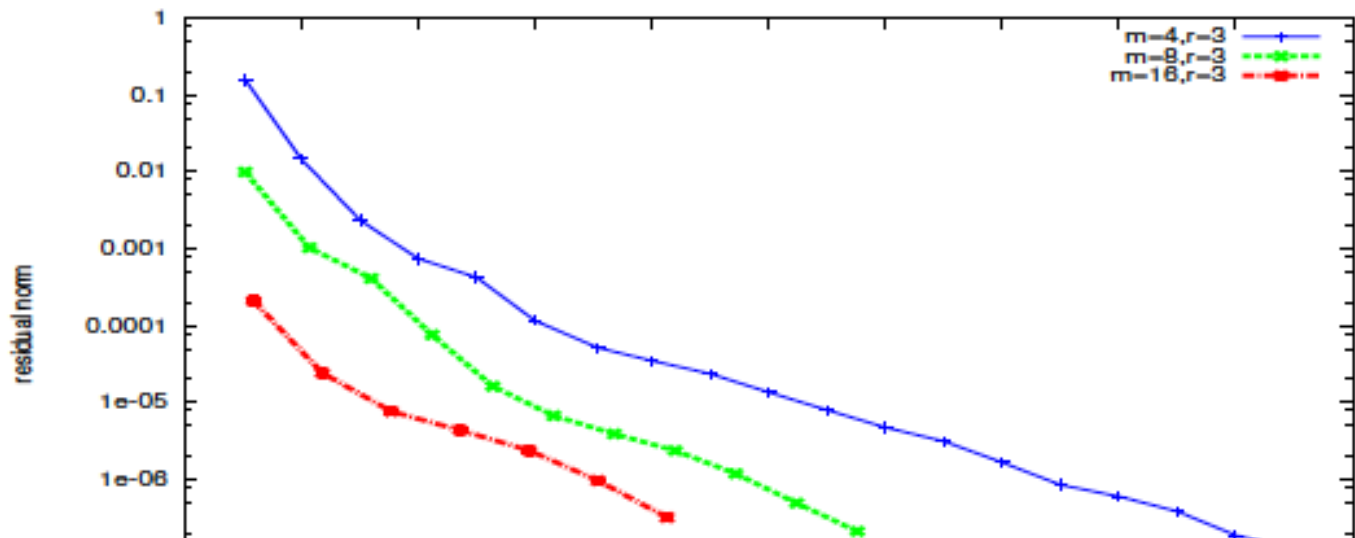
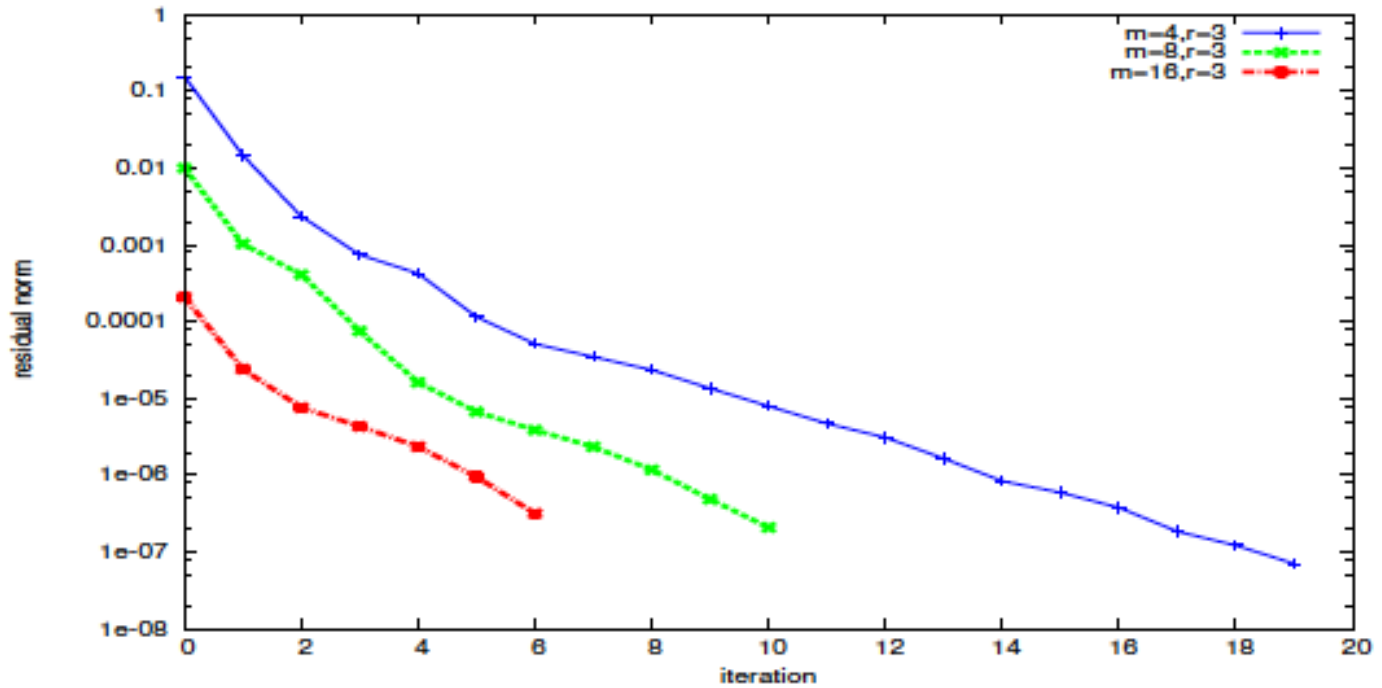
Number of iterations as α grows for the 41652230 X 41652230 twitter graph.

Experiments

Convergence with different number of shifts on twitter graph, where $\alpha = 0.85$.



Convergence with different size of subspace on twitter graph where $\alpha = 0.85$.



Big Data: microbiota (J.-M. BATTO - INRA MGP)



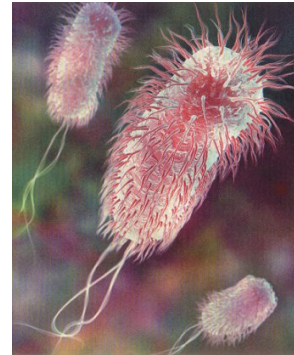
virus



protists



archae



bacteria



fungus

- 2 kg – more bacteria than human cells ($60 \cdot 10^{16}$)
- An unknown organ: intestinal microbiota
- Amount of sequence generated has increased 10^9 times in 20 years.

Big Data: microbiota

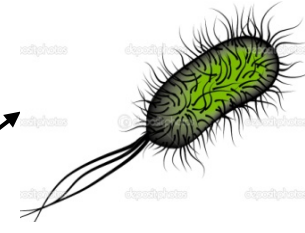
DNA preparation -> Get Sequences -> Compare to reference -> Counting & analyzing

← 100 – 1000 individuals →

↑ Up to 10 millions of genes ↓

id_fragment_MH0001	MH0002	MH0003	MH0004	MH0005	MH0006	MH0007	MH0008	
1	1.0246523439	8.3868647739	1.4874450848	2.3749503763	2.7242547919	4.7625759493	3.3844890635	1.5393658587
12	1.34675319526	1.9313563828	1.65852581486	9.8972700447	1.3267630323	1.3677934976	8.3035359532	6.5820901674
57	2.2892071460	4.3815305244	7.7707332906	2.56497818006	1.5525290094	7.7887477214	1.3471502008	1.6832954357
80	1.25263789532	1.89434148235	2.2545038192	5.3278079500	4.3803850688	5.41654766881	8.2863054205	3.5226732177
129	9.5879880329	4.56817276121	1.35029507561	2.7355209738	3.54721346105	6.761076985	1.6093633649	4.6580802879
144	1.1326224535	5.8291878391	1.5753395670	3.6259632993	3.44260910451	3.5020693235	7.4943601707	2.4709660277
162	6.7023224214	6.7852874733	2.888125396	1.19565205209	8.0145374559	1.1726322944	2.6077647166	2.81531793494
193	4.610749598	0	2.291400069	2.950443761	0	3.1899019391	4.8551623093	1.4218449525
226	1.0065218623	4.5445646333	4.0295384274	1.24537907951	4.4732302079	7.4001496204	5.8218633303	5.5600905259
233	1.0614220078	1.0694423241E	1.37940180995	4.4404063267	2.9665750993	3.4094820984	3.7364968910	7.1377548690
234	1.26274096216	9.3064209779	0	6.0509943978	0	4.9064710222	4.4552937081	0
241	1.0279681344	2.14481232084	3.6982056910	3.91828783581	4.10500698961	5.9720328632	1.06473125057	9.0769585235
305	3.81804695210	2.68218422471	2.66726611017	1.5225369727	3.3943423342	2.7671053495	2.3009686632	6.0857795514
336	4.6774055863	4.31859180624	9.573986878	3.1958774118	7.7577180236	5.91096043165	5.39418291742	1.30789541656
378	9.0895751438	9.5325014259	0	0	0	1.17675524517	9.1589784995	4.3740510020
448	8.2490245187	4.2948632798	2.8563949293	2.0956654007	2.2898678445	8.1295252468	1.30674803134	1.9707262756
451	1.19537420813	5.1089223328	4.5303982976	4.0134307940	3.6206851656	9.2709854413	6.2621724625	1.87541010285
548	1.2977239878	8.0829179307	2.5297423705	3.921906220	3.5969168228	1.2560645402	2.6496412253	4.3634085422
549	4.3167827556	0	3.8085132391	1.2554151013	8.4549969568	3.1811857486	5.772018228	3.9414525513
550	6.7728832890	1.7765162938	0	0	0	1.31582631960	0	0
604	1.69439579176	2.3531803258	0	0	0	2.2734712014	6.2914030326	0
667	4.132636003	5.3217054326	6.5639148629	2.9750722455	2.9143772567	2.7047541014	3.69824849002	4.6192629521
688	3.2878074218	6.4506344601	1.68240147942	2.07966215221	1.4939720306	3.1156798198	1.64081313900	2.61619640687
694	9.711189421830	9.5732293741	2.548334173	2.8875701207	2.3572168988	8.51417030331	2.930076177	1.75819837142
781	5.6827274172	0	1.18689967384	1.7168540157	1.4492040331	9.2528426127	1.15741888833	3.6849881729
793	5.1887202104	6.3853511006	8.6688087482	3.7794077387	5.0123491804	1.57156237954	2.1057595399	5.8415664848
854	6.9878954869	4.8205933143	3.05621432771	1.2592898394	1.18733888236	3.6589409122	1.0827759394	6.3257889453
888	1.10126367464	4.2537844423	2.40041846961	2.8730249126	2.5375844187	8.43367133167	2.9637285793	1.5023537973
903	9.862772130	7.1226494299	0	0	0	1.79765326195	4.9599683898	2.35137513442
952	6.8983070536	7.2376399715	0	0	0	2.3151082061	1.8494669396	0
994	1.0974066210	4.15779317515	1.701676128124	4.0690905774	2.9592138239	6.3171490044	2.3047347697	5.8702484807
1009	7.0509488786	0	2.5408702695	3.1408387343	1.8203850871	1.2733803092	5.0161307695	0
1123	1.10839138175	4.9787587065	2.2074821959	4.6615781882	3.430426223	9.2191652965	2.320332708	3.9979383363
1124	9.125412593	3.9005245355	0	0	0	1.06888823194	0	1.8738429065
1126	8.3226082131	2.4841052444	2.2028053288	2.1934954932	1.8745867269	1.1243996992	2.3867877022	8.3589856014
1130	1.7287429374	0	0	0	0	4.4938541991	1.3450971541	0
1190	9.910739477	3.8007231867	8.4324745894	3.979636407	2.879032478	7.9433265623	7.7561897302	1.1844662099
1191	7.7878789820	8.9040495801	1.5897322236	1.8014189182	2.847001964	5.0902761291	1.725801785	8.2264261966
1214	1.719428247	9.9829596949	4.4262414401	1.1424442652	5.867827302	1.5404522895	6.714244749	0
1289	9.6459769146	6.2112753796	1.270963744	2.9677382203	2.3044244616	4.520046849	1.8895383661	8.7894910066
1300	1.80861620076	4.2871544844	2.2296586371	5.7582204363	3.3410478724	5.0282196828	1.21104878329	4.8952937814
1470	7.0505867594	2.171231991471	0	0	0	9.9163074853	1.0686049941	6.9630060630
1539	5.8386924905	3.3692461836	0	0	0	8.3184422503	6.4744503186	3.0920018704
1590	1.1049126649	2.2461641290	7.9672345922	4.3497690804	2.7987426743	8.76209295037	1.786830871	1.6490675042
1798	2.2967527147	1.129557754187	1.7328735238	2.72194451576	4.8246629900	9.7826719194	1.9709514505	3.3093295324
1803	4.8897307755	2.3193850631	8.2714726673	1.15026859614	1.2394797783	3.2241408626	7.6179093248	4.2800979392
1806	2.3764502768	3.4283557760	0	0	0	3.3957519334	7.2468408829	0
1834	5.6440649075	2.6587248975	3.3294184213	6.8811847960	8.7233060745	5.2513703458	1.0643937598	0
1872	1.0741136439	7.4895225961	3.5878979127	7.8636270808	2.5420959331	8.5016678905	1.0752721454	9.4803219957
1877	8.6847871282	5.3503805265	7.90750630161	2.3505789329	2.0914809142	5.0637297088	1.32165399614	4.6762996371
1880	4.2330520556	3.382325846	0	0	0	0	1.4081929443	0
1883	5.3574565079	2.90071039481	1.28611706845	1.70136944044	4.0573595970	2.2804229575	3.66716191258	6.9652097368
1884	3.5621072053	0	0	0	0	3.0885302836	0	0
1900	6.0560452767	6.2761432746	0	0	0	1.11232196085	1.0821847794	0
1942	4.7673207617	1.4223940720	0	0	0	9.337325194	1.171042841071	1.36717761583
1958	4.3503406516	1.5283205683	0	0	0	5.3351177979	9.93183940571	4.0655537021

bacteria = ~3000 genes



parasites = ~6000 genes



Virus = ~50 genes



Big Data: microbiota

Samples

genes

id_fragment_	MH0001	MH0002	MH0003	MH0004	MH0005	MH0006	MH0007	MH0008	
1.1024529429	0.385647739	1.0381647739	1.074450086	2.374903793	2.2242547919	1.7629759493	3.2044806235	1.9330595897	
12.1340703358	1.0310766289	1.0305010488	1.0305010488	1.0305010488	1.0305010488	1.0305010488	1.0305010488	1.0305010488	
57.7.282027440	4.3893002441	7.770732306	2.5614978006	1.0525020094	7.7681747214	1.3475020061	1.0083954207	0	
60.12526278920	1.0344444232	2.2540038922	5.7207897900	4.3800305868	6.4947966891	8.2863004202	3.5282702771	0	
120.2857879278	1.0681727193	1.9503399791	2.2902299791	5.4212344362	6.7678921869	1.0036335443	6.9506002079	0	
144.11528224529	5.8281878739	1.5753395670	3.6259832933	3.4420081049	3.0020082036	7.4343007702	2.4709660277	0	
180.6702324241	1.7902874723	2.0802205086	1.1865030209	0.0495374058	1.0258530946	2.0577479162	2.1037078946	0	
183.4.6178749636	0	0.2384030661	2.9504449760	0	0	0.3080838931	4.6558323693	1.4821849528	
220.10086029623	0.5449442323	4.0293584274	1.2457780789	4.4723202079	7.4801846204	5.8239632330	5.9600800289	0	
223.100842007018	1.0084423218	1.2704802262	2.060709093	3.4004002004	3.7284980289	1.0084423218	6.1377784900	0	
234.126274095218	0.9364200779	0	0.0509342978	0	0.4004792022	4.4582937900	0	0	
241.120786051844	1.4443202004	0.0302003007	0.3020782007	4.0030003000	0.3720202000	1.0470200000	0	0	
260.33804696307	0.0021842471	0.0072060707	1.0020089727	3.2043832462	0.2000860289	2.0000860289	6.0057789504	0	
336.4.6174055683	3.0953908264	0.5733888789	3.0558774740	7.7577380236	5.8090804393	5.3044829742	1.0078945468	0	
370.30055750430	0.5025004293	0	0	0	0.1767954657	0.890394995	4.3740050020	0	
448.8.249024587	4.2948632798	2.8583849293	2.0586654007	2.2898784445	0.326525468	1.0674800314	1.9707262756	0	
491.19557420010	0.1089223208	4.5203882970	0.013037940	3.6206889589	0.2703894433	2.6357246225	1.6754000289	0	
548.15257219878	0.0020010000	0.2527432000	0.3201890200	0.2568196200	1.0504846000	0.8484422000	4.3640000000	0	
549.4.387827956	0	0.8085122391	1.2854181043	8.4548968568	3.8118874881	5.7727082828	3.9448258133	0	
590.6.7228202890	1.7708182838	0	0	0	0.1398231880	0	0	0	
604.15849297978	2.253803296	0	0	0	0.227412016	6.2940003228	0	0	
667.4.3812826020	0.921094032	0.563849828	2.9790722495	2.043772967	2.7047541046	3.9382814000	4.1826239220	0	
680.328707478	0.4006440018	1.0304047942	2.0798629232	1.8329720206	0.3186303899	1.8400310000	2.0180004603	0	
694.9.7111842830	0.979293741	2.5483434173	2.0879705207	2.3927268988	0.5147030331	2.3300076177	1.7919897742	0	
781.5.6227274720	0	1.1903967384	1.7769409797	1.4402040331	0.252427627	1.9748095003	3.0484980729	0	
793.5.937230234	0.6363051006	0.4460007462	3.7784877207	0.0024890384	1.9795227094	2.0407390000	0.4496448448	0	
854.6.3878945489	0.025032143	3.0562142271	1.2532098334	1.1872388823	3.858409322	1.9327775934	6.3257800453	0	
888.19028307444	4.257384423	2.4004848486	2.9730249302	2.5370444801	1.4330732963	2.5623205703	1.90239537929	0	
903.6.962772150	1.024494299	0	0	0.1797652015	0.953863088	2.951370442	0	0	
952.6.880207056	1.237639791	0	0	0.2353802061	1.8444849398	0	0	0	
994.10574602004	0.677930705	1.7086780381	4.060090774	2.2602302009	0.311483044	2.3047347630	5.9702448407	0	
1009.7.0509489786	0	0.5480702895	3.1402873429	2.6202650878	1.2723802002	5.5061007685	0	0	
1023.11802818975	0.978760705	2.0748218959	4.688781882	2.4304262210	0.218622805	2.332022708	3.9979383263	0	
1024.9.92543955	0.3006249295	0	0.1068882984	0	0.88893924	1.873843000	0	0	
1026.8.324808231	2.484052444	2.002805328	2.1834854930	1.8749872839	1.3243960026	2.3867877002	8.3588856014	0	
1028.1237429374	0	0	4.4834841891	1.450597000	0	0.5220046701	0.8181870089	0	
1090.8.9892739477	3.803723967	0.4324745684	3.577868442	2.6676324673	7.0433252623	7.7641037982	1.13446222097	0	
1091.7.7670679920	0.9640495060	1.097922206	1.8014993822	2.647003964	0.0902761293	1.722607982	8.226426189	0	
1191.7.738320287	1.902008184	4.8282440012	1.8144420002	0.8078707002	1.044020089	6.742417749	0	0	
1389.9.6459788145	6.2127572902	1.27805637442	2.9677802302	2.3044244616	6.723046648	1.6389503668	8.7684900066	0	
1390.10008620074	4.287584484	1.2389885371	0.786204343	3.340478724	0.0262982829	1.2184878828	4.8990239798	0	
1470.7.0950967804	4.17123889181	0	0.9389074823	1.0388648416	1.8603006206	7.3818418271	9.3818187842	0	
1850.5.3388242495	3.382424930	0	0	0.3384425003	6.4744503380	3.9300019704	0	0	
1900.1048182648	2.448484200	0	0	1.7822245822	4.240700000	2.3785142543	3.7620302007	1.8600500716	0
1798.7.236752147	1.2955775487	1.7238728238	2.7214445574	4.8244629300	3.762879394	1.9709545685	3.3092985244	0	
1800.4.889700755	2.3193850629	0.2714726878	1.9502858614	1.0294797703	3.224408828	7.873003244	4.2000876232	0	
1900.2.378460278	0.420805780	0	0	0	0.338578024	7.248400000	0	0	
1834.5.6440630475	2.6587248075	3.9291842103	6.8811847860	8.7233000745	6.251070458	1.0443037590	0	0	
1872.1074018434	0.4896220005	3.0878079207	7.863627088	2.8420900201	0.5068780800	1.0752714844	5.4803108007	0	
1877.8.644817202	1.2500260205	2.3607000306	2.0018400944	0.0000000000	0.621707306	1.9218629844	4.4762396327	0	
1880.4.2330920956	3.983295846	0	0	0	0	1.40919234430	0	0	
1883.5.387465019	3.9007300481	1.0611705461	1.7010844044	1.0573959703	2.0842228627	3.6178032502	6.0852097368	0	
1884.3.8626072093	0	0	0.3088530286	0	0	1.10232186895	1.0828477946	0	
1800.6.0980452787	6.276432746	0	0	0	0.5276782807	1.7405408424	0	0	
1842.4.972307071	1.6230407026	0	0	0.3337325184	1.1784248071	1.8217789623	0	0	
1858.4.3503406816	1.029205883	0	0	0	0.3391179795	9.83188240571	4.0655037021	0	

Matrix 10^6 genes
by 800 samples

genes

genes

Correlation matrix

Counting matrix

MetaProf →

energy efficiency multiplied by 4.7 with the GPU implementation

Principal Coordinates Analysis applied on the matrices of distances between samples, concentrating the major variations in the samples in a small space implies using many linear algebra techniques.

Numerical methods/algorithms & HPC techniques have to be defined/adapted to increase data-scalability.

Concluding remarks and future work

- Conventional means of investigation are essential;
- Our predictions provide complementary solutions ;
- The virus/individual characteristics have to be integrated
- The impact of social graph structure on propagation of virus have to be extended

For efficient computation solver, many problems arise:

- Methods / algorithms
- Data Processing
- Programming models for Exascale computing (graph computation, PGAS ...)
- ...