

PageRank, alebo ako sa na trhu presadil Google

Michal „Mišof“ Forišek

Department of Theoretical Computer Science
Faculty of Mathematics, Physics and Informatics
Comenius University
Bratislava, Slovakia

7. júna 2017

Vyhľadávanie v bežnom texte

Bežná webstránka má do 10 kB *textu*.

Ctrl-F (find) v prehliadači funguje takmer okamžite.

Kniha *Isaac Asimov: Ja, robot* má 425 kB.

Vyhľadanie všetkých výskytov
slova „robot“: 0.007 s.

Kým máme jedného človeka a málo textu,
všetko funguje ako má.

Vyhľadávanie v bežnom texte

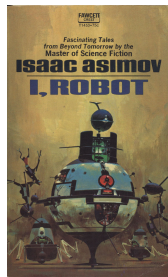
Bežná webstránka má do 10 kB *textu*.

Ctrl-F (find) v prehliadači funguje takmer okamžite.

Kniha *Isaac Asimov: Ja, robot* má 425 kB.

Vyhľadanie všetkých výskytov
slova „robot“: 0.007 s.

Kým máme jedného človeka a málo textu,
všetko funguje ako má.



Vyhľadávanie v bežnom texte

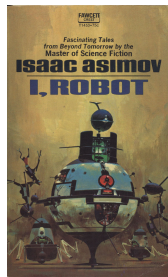
Bežná webstránka má do 10 kB *textu*.

Ctrl-F (find) v prehliadači funguje takmer okamžite.

Kniha *Isaac Asimov: Ja, robot* má 425 kB.

Vyhľadanie všetkých výskytov
slova „robot“: 0.007 s.

Kým máme jedného človeka a málo textu,
všetko funguje ako má.



Množstvo informácie rastie

- Disketa: posledné bežné mali **1.44 MB**
- CD: typicky do **700 MB**
(plné auto kníh)
- DVD: najrozšírenejšie majú **4.7 GB**
- Video v kvalite HD: **28 Mb/s**

Inými slovami: každú sekundu 3.5 megabajtu údajov.
(Každé 4 desatiny sekundy meniť disketu?)

- Bežný pevný disk: rádovo **1 TB**

Množstvo informácie rastie

- Disketa: posledné bežné mali **1.44 MB**
- CD: typicky do **700 MB**
(plné auto kníh)
- DVD: najrozšírenejšie majú **4.7 GB**
- Video v kvalite HD: **28 Mb/s**

Inými slovami: každú sekundu 3.5 megabajtu údajov.
(Každé 4 desatiny sekundy meniť disketu?)

- Bežný pevný disk: rádovo **1 TB**

Množstvo informácie rastie

- Disketa: posledné bežné mali **1.44 MB**
- CD: typicky do **700 MB**
(plné auto kníh)
- DVD: najrozšírenejšie majú **4.7 GB**
- Video v kvalite HD: **28 Mb/s**

Inými slovami: každú sekundu 3.5 megabajtu údajov.
(Každé 4 desatiny sekundy meniť disketu?)

- Bežný pevný disk: rádovo **1 TB**

Množstvo informácie rastie

- Disketa: posledné bežné mali **1.44 MB**
- CD: typicky do **700 MB**
(plné auto kníh)
- DVD: najrozšírenejšie majú **4.7 GB**
- Video v kvalite HD: **28 Mb/s**

Inými slovami: každú sekundu 3.5 megabajtu údajov.
(Každé 4 desatiny sekundy meniť disketu?)

- Bežný pevný disk: rádovo **1 TB**

Množstvo informácie rastie

- Disketa: posledné bežné mali **1.44 MB**
- CD: typicky do **700 MB**
(plné auto kníh)
- DVD: najrozšírenejšie majú **4.7 GB**
- Video v kvalite HD: **28 Mb/s**

Inými slovami: každú sekundu 3.5 megabajtu údajov.
(Každé 4 desatiny sekundy meniť disketu?)

- Bežný pevný disk: rádovo **1 TB**

Množstvo informácie rastie ešte viac

americká Library of Congress: **50 TB** informácií
v tlačenej podobe



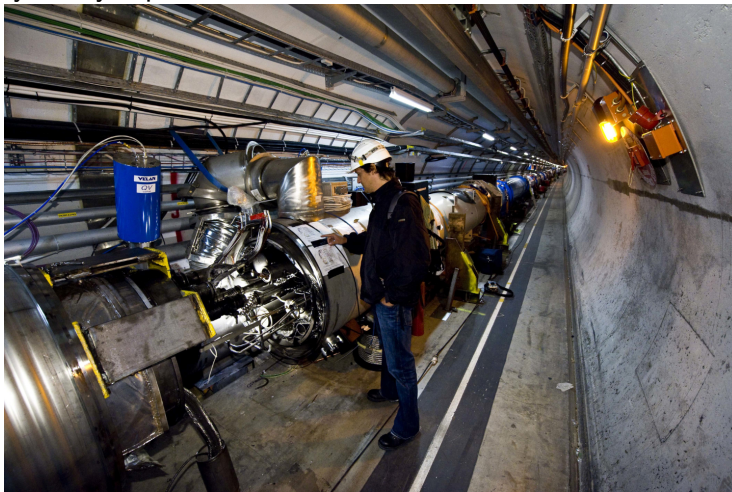
Množstvo informácie rastie ešte viac

špeciálne efekty pre film *Avatar*: **1 PB** = 1000 TB



Množstvo informácie rastie ešte viac

fyzikálny experiment LHC: získame **15 PB** dát za rok



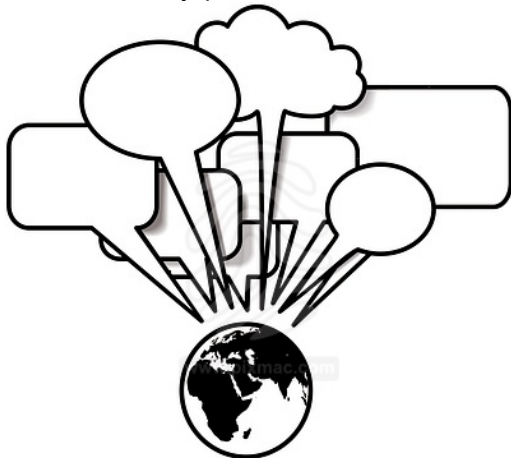
Množstvo informácie rastie ešte viac

Steam: distribuuje vyše **16 PB** dát každý týždeň



Množstvo informácie rastie ešte viac

Všetky slová, ktoré kto kedy povedal: **5 EB = 5000 PB**



Množstvo informácie rastie ešte viac

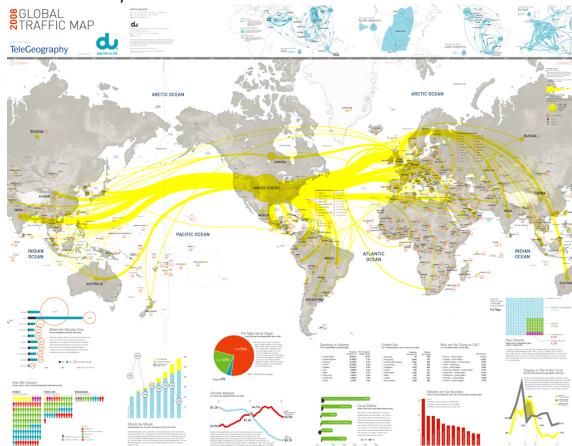
Google: spracuje **24 PB** dát denne

celé dáta používané pri vyhľadávaní majú aspoň stovky PB
úplne všetky uložené dáta sú odhadované na desiatky EB



Množstvo informácie rastie ešte viac

Celkový prenos dát internetom (v 2015):
72.5 EB za mesiac, z toho **3.7 EB** mobilnou sieťou



Web rastie

- 1998: Založené Google, pozná 26 miliónov stránok
- 2015: Google pozná **120 biliónov** (1.2×10^{14}) rôznych URL

Ako dlho by trvalo nájsť výskyty slova „robot“ na celom webe?

rátajme spolu:

100 biliónov stránok \times trebárs 1 kB = 100 PB textu.

to je zhruba štvrt' bilióna Asimovových kníh

\Rightarrow bude nám to trvať cca 50 rokov

Web rastie

- 1998: Založené Google, pozná 26 miliónov stránok
- 2015: Google pozná **120 biliónov** (1.2×10^{14}) rôznych URL

Ako dlho by trvalo nájsť výskyty slova „robot“ na celom webe?

rátajme spolu:

100 biliónov stránok \times trebárs 1 kB = 100 PB textu.

to je zhruba štvrt' bilióna Asimovových kníh

\Rightarrow bude nám to trvať cca 50 rokov

Web rastie

- 1998: Založené Google, pozná 26 miliónov stránok
- 2015: Google pozná **120 biliónov** (1.2×10^{14}) rôznych URL

Ako dlho by trvalo nájsť výskyty slova „robot“ na celom webe?

rátajme spolu:

100 biliónov stránok \times trebárs 1 kB = 100 PB textu.

to je zhruba štvrt' bilióna Asimovových kníh

\Rightarrow bude nám to trvať cca 50 rokov

Web rastie

- 1998: Založené Google, pozná 26 miliónov stránok
- 2015: Google pozná **120 biliónov** (1.2×10^{14}) rôznych URL

Ako dlho by trvalo nájsť výskyty slova „robot“ na celom webe?

rátajme spolu:

100 biliónov stránok \times trebárs 1 kB = 100 PB textu.

to je zhruba štvrt' bilióna Asimovových kníh

\Rightarrow bude nám to trvať cca 50 rokov

Web rastie ešte viac

Realita: **0.25 sekundy**

Roky čakať na výsledok nechceme!

Riešenie: paralelizácia

„nech si Google nakúpi viac počítačov :-“

Na čo zabúdame?

Na ostatných ľuďoch!

Denne Google spracuje **3.5 miliardy** vyhľadávaní.

To by bolo cca 200 miliárd rokov výpočtu.

Web rastie ešte viac

Realita: **0.25 sekundy**

Roky čakať na výsledok nechceme!

Riešenie: paralelizácia

„nech si Google nakúpi viac počítačov :-“

Na čo zabúdame?

Na ostatných ľuďoch!

Denne Google spracuje **3.5 miliardy** vyhľadávaní.

To by bolo cca 200 miliárd rokov výpočtu.

Web rastie ešte viac

Realita: **0.25 sekundy**

Roky čakať na výsledok nechceme!

Riešenie: paralelizácia

„nech si Google nakúpi viac počítačov :-“

Na čo zabúdame?

Na ostatných ľuďoch!

Denne Google spracuje **3.5 miliardy** vyhľadávaní.

To by bolo cca 200 miliárd rokov výpočtu.

Potrebujeme lepšie algoritmy

Prvý nápad:

pre každé slovo budeme mať zoznam stránok, ktoré ho obsahujú.

Na jednej strane žiadna výhra (veľa pamäte; čo zložitejšie otázky?).
Na druhej strane ani toto nestačí.

Potrebujeme lepšie algoritmy

Prvý nápad:

pre každé slovo budeme mať zoznam stránok, ktoré ho obsahujú.

Na jednej strane žiadna výhra (veľa pamäte; čo zložitejšie otázky?).

Na druhej strane ani toto nestačí.

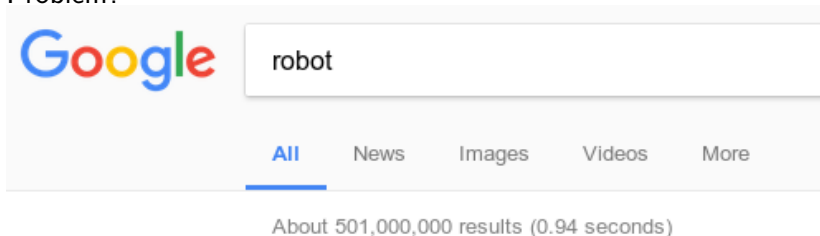
Potrebuje lepšie algoritmy

Prvý nápad:

pre každé slovo budeme mať zoznam stránok, ktoré ho obsahujú.

Na jednej strane žiadna výhra (veľa pamäte; čo zložitejšie otázky?).
Na druhej strane ani toto nestačí.

Problém?



Ktoré z nich ukázať užívateľovi?

Ako spoznať, ktoré sú relevantné?

PageRank

Nový algoritmus pri zrode Googlu.

Autori Larry Page (podľa neho sa volá) a Sergej Brin.

Vyhodnocuje, ktoré stránky sú dôležité a ktoré nie.

Myšlienka: dôležité stránky sú tie, ktoré ľudia navštevujú.

Problém: ale my nevieme, ktoré to sú. . .

Riešenie: oni nám to povedia – spravia linky!

PageRank

Nový algoritmus pri zrode Googlu.

Autori Larry Page (podľa neho sa volá) a Sergej Brin.

Vyhodnocuje, ktoré stránky sú dôležité a ktoré nie.

Myšlienka: dôležité stránky sú tie, ktoré ľudia navštevujú.

Problém: ale my nevieme, ktoré to sú. . .

Riešenie: oni nám to povedia – spravia linky!

PageRank

Nový algoritmus pri zrode Googlu.

Autori Larry Page (podľa neho sa volá) a Sergej Brin.

Vyhodnocuje, ktoré stránky sú dôležité a ktoré nie.

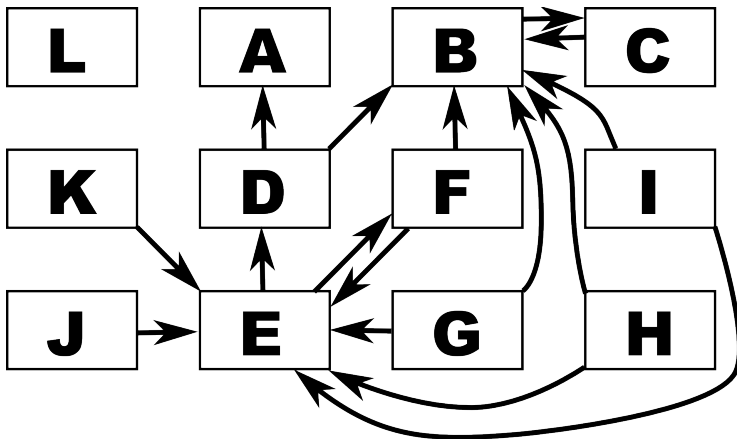
Myšlienka: dôležité stránky sú tie, ktoré ľudia navštevujú.

Problém: ale my nevieme, ktoré to sú. . .

Riešenie: oni nám to povedia – spravia linky!

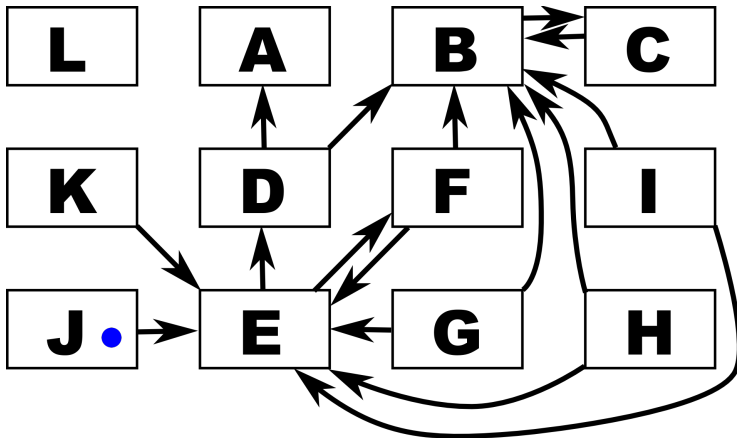
Základná myšlienka PageRanku

Človek náhodne kliká po linkách.



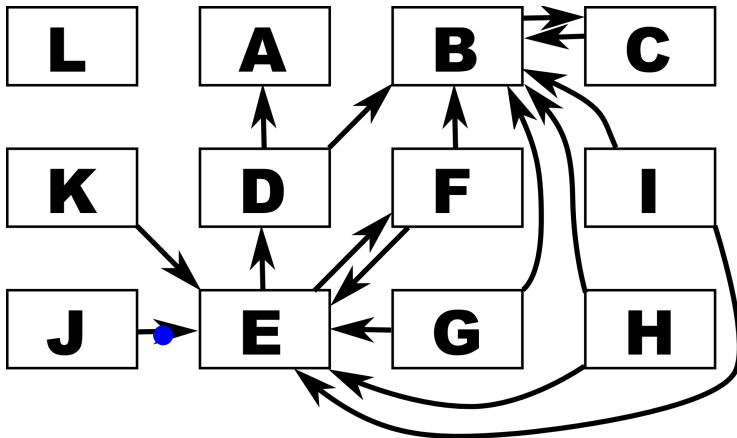
Základná myšlienka PageRanku

Človek náhodne kliká po linkách.



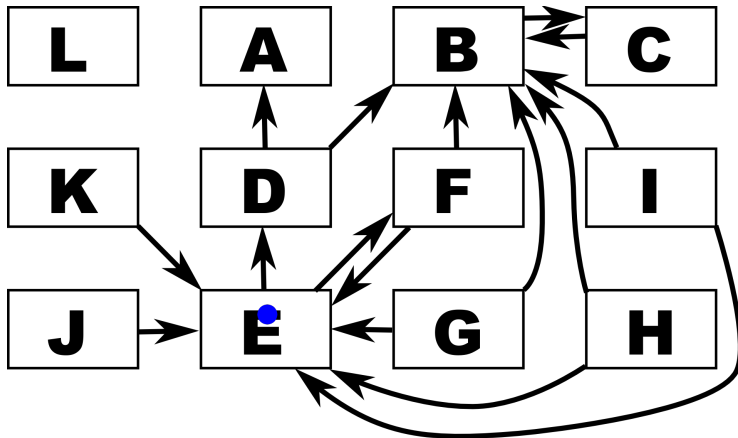
Základná myšlienka PageRanku

Človek náhodne kliká po linkách.



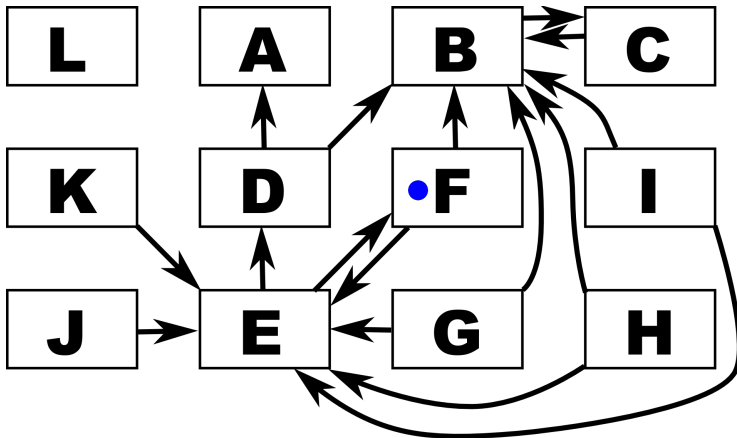
Základná myšlienka PageRanku

Človek náhodne kliká po linkách.



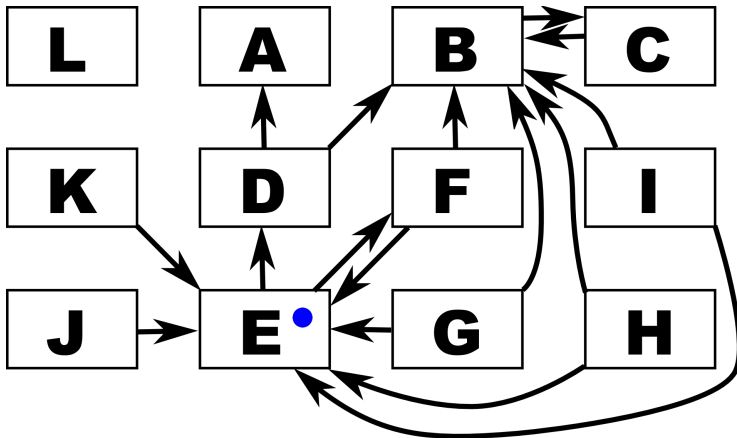
Základná myšlienka PageRanku

Človek náhodne kliká po linkách.



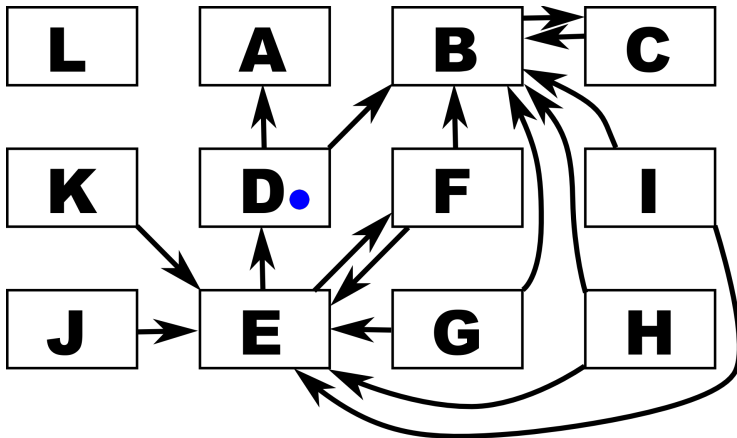
Základná myšlienka PageRanku

Človek náhodne kliká po linkách.



Základná myšlienka PageRanku

Človek náhodne kliká po linkách.



Nie je náhodné klikanie blbosť?

Človek predsa nekliká po linkách *náhodne!*

Jeden človek nie – ale všetci dokopy áno!

Presnejší model správania sa ľudí:

- 85% prípadov: klikne na linku
- 15% prípadov: odíde na nesúvisiacu stránku

Nie je náhodné klikanie blbosť?

Človek predsa nekliká po linkách *náhodne!*

Jeden človek nie – ale všetci dokopy áno!

Presnejší model správania sa ľudí:

- 85% prípadov: klikne na linku
- 15% prípadov: odíde na nesúvisiacu stránku

Nie je náhodné klikanie blbosť?

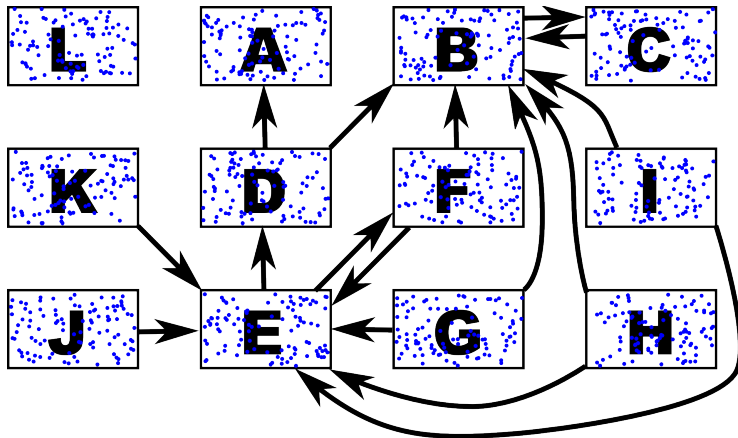
Človek predsa nekliká po linkách *náhodne!*

Jeden človek nie – ale všetci dokopy áno!

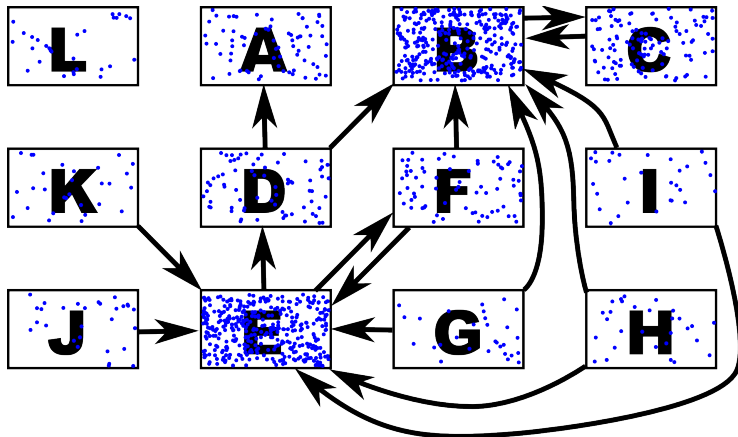
Presnejší model správania sa ľudí:

- 85% prípadov: klikne na linku
- 15% prípadov: odíde na nesúvisiacu stránku

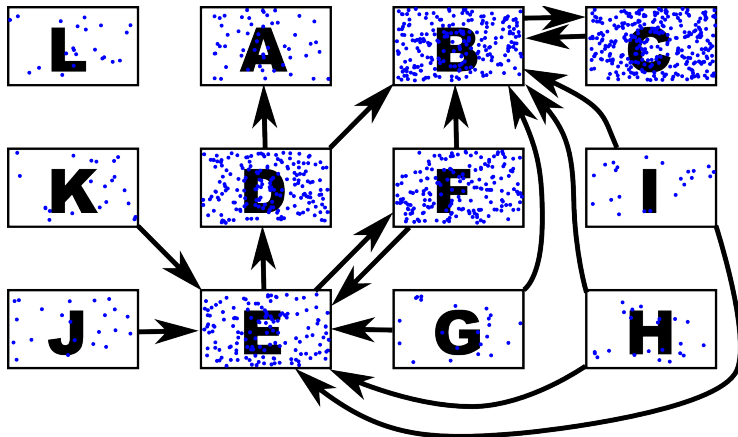
Simulácia pre veľa ľudí



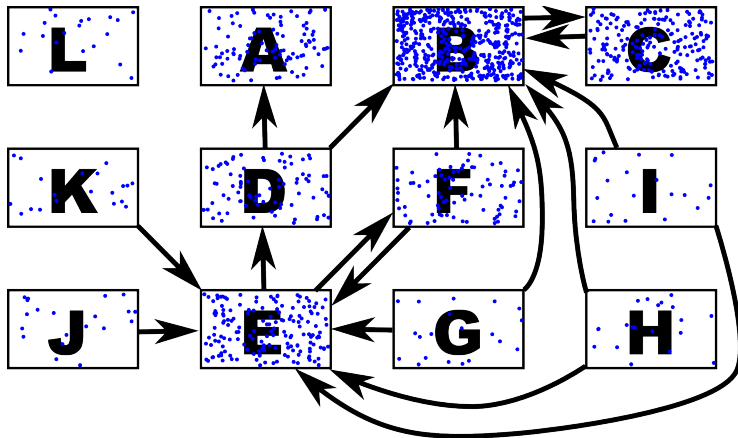
Simulácia pre veľa ľudí



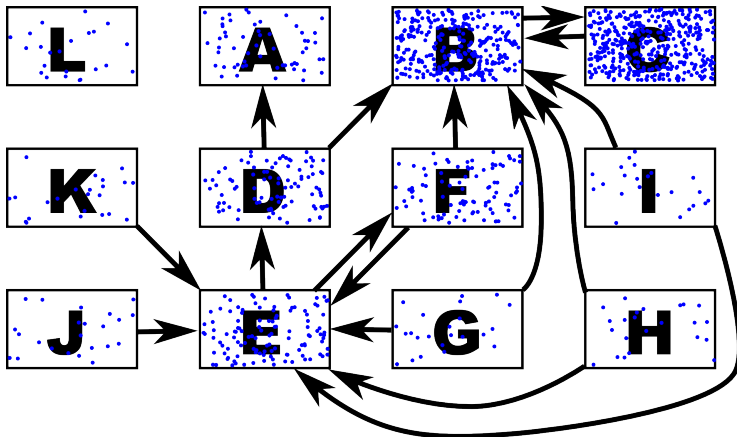
Simulácia pre veľa ľudí



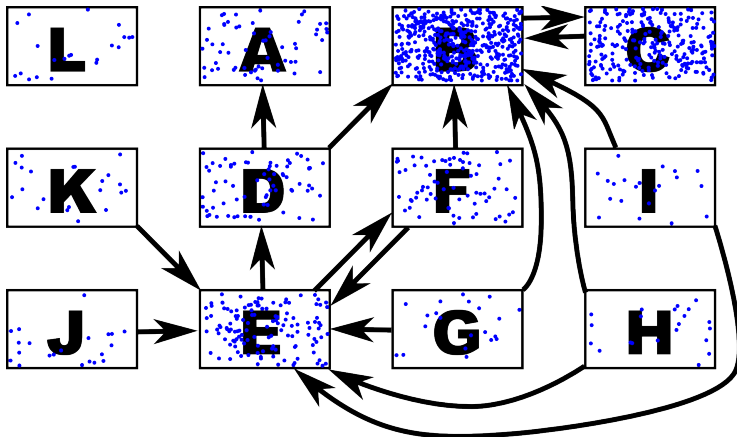
Simulácia pre veľa ľudí



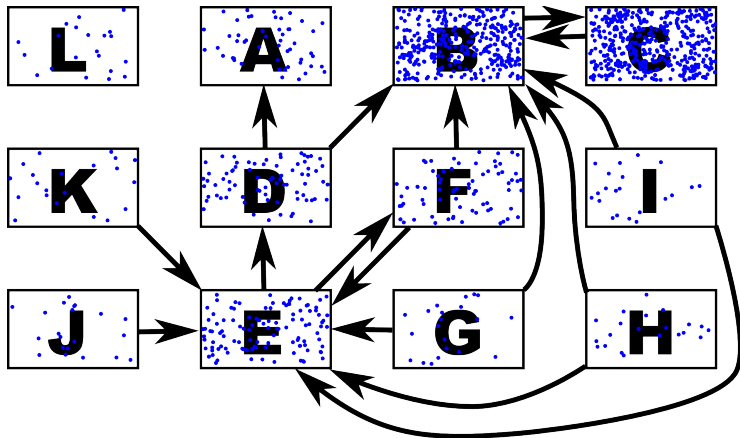
Simulácia pre veľa ľudí



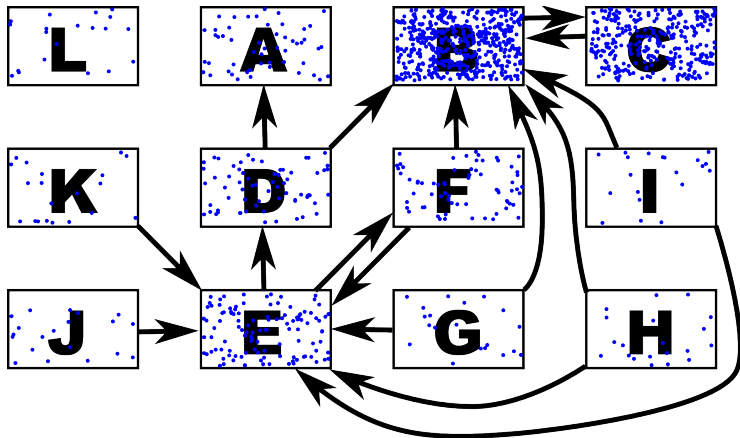
Simulácia pre veľa ľudí



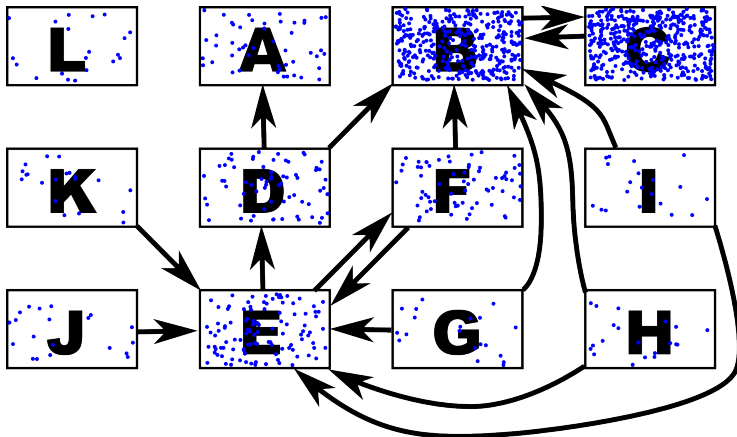
Simulácia pre veľa ľudí



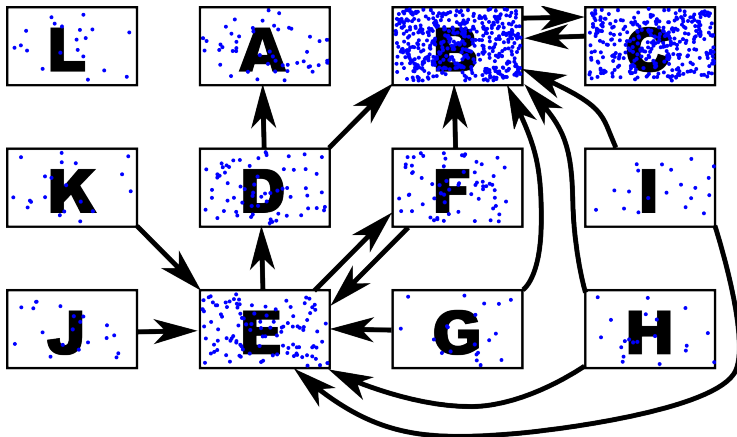
Simulácia pre veľa ľudí



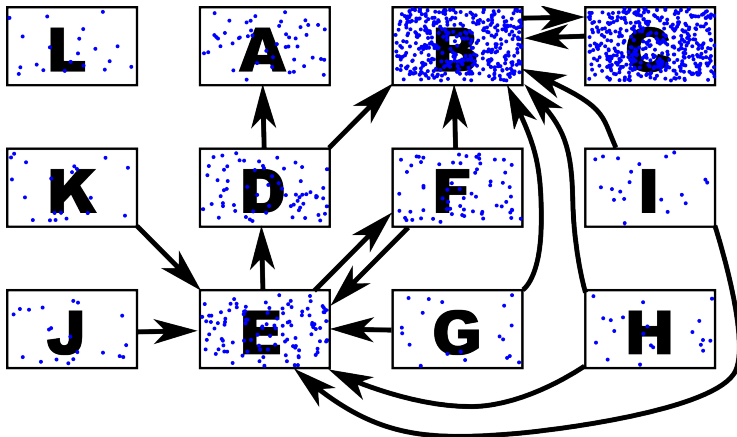
Simulácia pre veľa ľudí



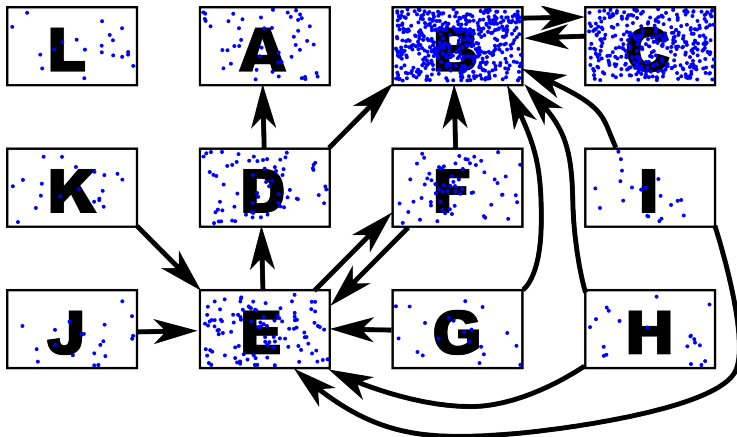
Simulácia pre veľa ľudí



Simulácia pre veľa ľudí

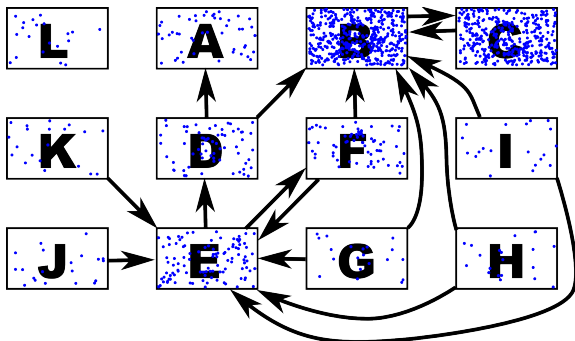


Simulácia pre veľa ľudí



Sústava lineárnych rovníc

Ako zistiť, k čomu to celé smeruje?

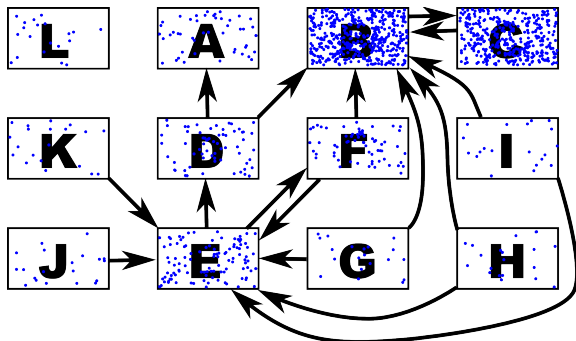


Rovnice!

$$\text{rank}_E = 0.15 \cdot \frac{1}{12} + 0.85 \cdot \left(\frac{\text{rank}_F}{2} + \frac{\text{rank}_G}{2} + \dots + \text{rank}_K \right)$$

Sústava lineárnych rovníc

Ako zistiť, k čomu to celé smeruje?



Rovnice!

$$\text{rank}_E = 0.15 \cdot \frac{1}{12} + 0.85 \cdot \left(\frac{\text{rank}_F}{2} + \frac{\text{rank}_G}{2} + \dots + \text{rank}_K \right)$$

Sústava lineárnych rovníc

Gaussova eliminačná metóda:

$$\begin{array}{rclcl} x + 2y - z & = & -1 \\ -x - 3y + 5z & = & 10 \\ 3x + y + 4z & = & 16 \end{array}$$

Sústava lineárnych rovníc

Gaussova eliminačná metóda:

$$\begin{array}{rclcl} x + & 2y - & z = & -1 \\ & -y + & 4z = & 9 \\ & -5y + & 7z = & 19 \end{array}$$

Sústava lineárnych rovníc

Gaussova eliminačná metóda:

$$\begin{array}{rclcl} x + 2y - z & = & -1 \\ -y + 4z & = & 9 \\ 13z & = & 26 \end{array}$$

a postupným dosadzovaním máme $z = 2$, $y = -1$, $x = 3$

Sústava lineárnych rovníc

Gaussova eliminačná metóda:

$$\begin{array}{rclcl} x + 2y - z & = & -1 \\ -y + 4z & = & 9 \\ 13z & = & 26 \end{array}$$

a postupným dosadzovaním máme $z = 2$, $y = -1$, $x = 3$

Sústava lineárnych rovníc

Gaussova eliminačná metóda je algoritmus.

n premenných \Rightarrow rádovo n^3 operácií

Problém:

miliardy webstránok \Rightarrow miliardy premenných

miliardy premenných \Rightarrow miliardy miliárd miliárd operácií

\Rightarrow stotisíc miliárd rokov výpočtu?!

A samotnú sústavu rovníc tvorí miliarda miliárd čísel

\Rightarrow exobajty pamäte?!

Sústava lineárnych rovníc

Gaussova eliminačná metóda je algoritmus.

n premenných \Rightarrow rádovo n^3 operácií

Problém:

miliardy webstránok \Rightarrow miliardy premenných

miliardy premenných \Rightarrow miliardy miliárd miliárd operácií

\Rightarrow stotisíc miliárd rokov výpočtu?!

A samotnú sústavu rovníc tvorí miliarda miliárd čísel

\Rightarrow exobajty pamäte?!

Sústava lineárnych rovníc

Gaussova eliminačná metóda je algoritmus.

n premenných \Rightarrow rádovo n^3 operácií

Problém:

miliardy webstránok \Rightarrow miliardy premenných

miliardy premenných \Rightarrow miliardy miliárd miliárd operácií

\Rightarrow stotisíc miliárd rokov výpočtu?!

A samotnú sústavu rovníc tvorí miliarda miliárd čísel

\Rightarrow exobajty pamäte?!

Sústava lineárnych rovníc

Gaussova eliminačná metóda je algoritmus.

n premenných \Rightarrow rádovo n^3 operácií

Problém:

miliardy webstránok \Rightarrow miliardy premenných

miliardy premenných \Rightarrow miliardy miliárd miliárd operácií

\Rightarrow stotisíc miliárd rokov výpočtu?!

A samotnú sústavu rovníc tvorí miliarda miliárd čísel

\Rightarrow exobajty pamäte?!

Sústava lineárnych rovníc

Riešenie problému: nové algoritmy.

Šikovne využijeme to, že naša sústava rovníc je „riedka“:
Každá webstránka ukazuje len na pár iných.

Prototyp Googlu: jeden počítač za niekoľko hodín spočítal
PageRank pre „celý internet“ (vtedy pár miliónov stránok)

Súčasnosť: PageRank je jedným z mnohých faktorov stále
používaných Googlom.

Sústava lineárnych rovníc

Riešenie problému: nové algoritmy.

Šikovne využijeme to, že naša sústava rovníc je „riedka“:
Každá webstránka ukazuje len na pár iných.

Prototyp Googlu: jeden počítač za niekoľko hodín spočítal
PageRank pre „celý internet“ (vtedy pár miliónov stránok)

Súčasnoscť: PageRank je jedným z mnohých faktorov stále
používaných Googľom.

Sústava lineárnych rovníc

Riešenie problému: nové algoritmy.

Šikovne využijeme to, že naša sústava rovníc je „riedka“:
Každá webstránka ukazuje len na pár iných.

Prototyp Googlu: jeden počítač za niekoľko hodín spočítal
PageRank pre „celý internet“ (vtedy pár miliónov stránok)

Súčasnosť: PageRank je jedným z mnohých faktorov stále
používaných Googlom.

Koniec

Ďakujem za pozornosť!