**Carnegie Mellon**

# Connection subgraphs

*Christos Faloutsos* (CMU & IBM)
*Kevin McCurley* (IBM)
*Andrew Tomkins* (IBM)

---

**IBM**

**Carnegie Mellon**

## Outline

- Introduction / Motivation
- Survey
- Proposed Method
- Algorithms
- Experiments
- Conclusions

---

**IBM**

**Carnegie Mellon**

## Introduction

- What are the best path**s** between 'Kidman' and 'Diaz'?

---

**IBM**

**Carnegie Mellon**

## Problem definition

- Given a graph, and two nodes $s$ and $t$, and a 'budget' $b$ of nodes
- Find the best $b$ nodes that capture the relationship between $s$ and $t$

---

**IBM**

**Carnegie Mellon**

## Problem definition

- Given a graph, and two nodes $s$ and $t$, and a 'budget' $b$ of nodes
- Find the best $b$ nodes that capture the relationship between $s$ and $t$

---

**IBM**

**Carnegie Mellon**

## Problem definition

- Part 1: How to quantify the goodness?
- Part 2: How to pick 'best few' nodes?
- Part 3: Scalability: large graphs (10**7 nodes)

## Survey

- Graph Partitioning
  - [Karypis+Kumar]; [Newman+];
  - [Virtanen]; …
- Communities
  - [Flake+]; [Tomkins, Kleinberg+]
- External distances [Palmer+]

---

## Outline

- Introduction / Motivation
- Survey
➡ • Proposed Method
- Algorithms
- Experiments
- Conclusions

---

## Proposed method

- part 1: measuring goodness:
  - electricity
- part 2: finding good paths
  - dynamic programming
- part 3: scalability
  - heuristics

---

## Electricity

- Why not shortest path?

---

## Electricity

- Why not shortest path?
- Why not net. flow?

---

## Electricity

- Why not shortest path?
- Why not net. flow?
- Why not plain 'voltages'?

+1V     0V

2

**Carnegie Mellon**

## Electricity

- Why not shortest path?
- Why not net. flow?
- Why not plain 'voltages'?

**+1V**      **0V**

**+0.5V**

---

**Carnegie Mellon**

## Electricity, cont'd

- Proposed method: voltages **with** universal sink:
  - ~ 'tax collector'
- goodness of a path:
- its electric current[*]!   **+1V**      **0V**

s      t

f

**0V**

...

---

**Carnegie Mellon**

## Outline

- Introduction / Motivation
- Survey
- Proposed Method
➡ - Algorithms
- Experiments
- Conclusions

---

**Carnegie Mellon**

## Electricity – Algorithm

- Voltages/Amperages can be computed easily ( $O(E)$ )
- without universal sink:

  $v(i) = \Sigma \, [v(j) * C(i,j) \, / \, C(i,*) \,]$

  $i \neq source, sink$

  $v(source)=1; \; v(sink)=0$

---

**Carnegie Mellon**

## Electricity – Algorithm

**With** universal sink:

$v(i) = \mathbf{1/(1+a)} \, \Sigma \, [v(j) * C(i,j) \, / \, C(i,*) \,]$

*(~ insensitive to a (=1))*

---

**Carnegie Mellon**
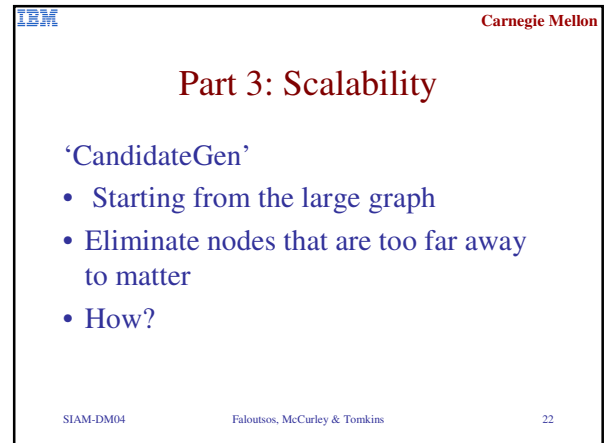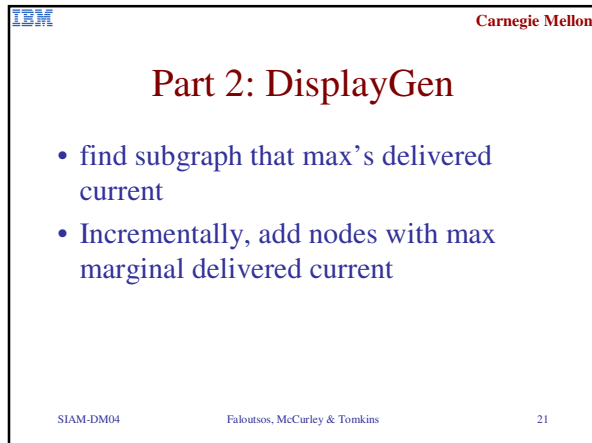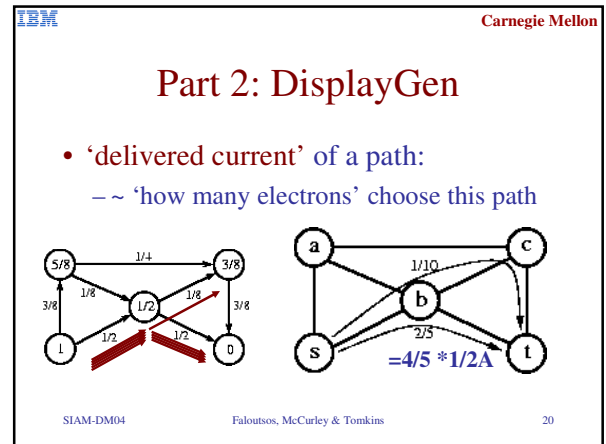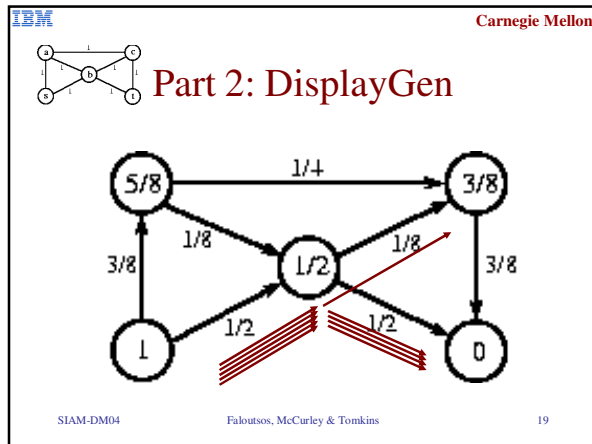
## Part 2: DisplayGen

Given the voltages and amperages

- Which *b* nodes to keep?
- (and how to spot them quickly?)

a      c

b

s      t

**Carnegie Mellon**

# Part 2: DisplayGen

---

**Carnegie Mellon**

# Part 2: DisplayGen

- 'delivered current' of a path:
  - ~ 'how many electrons' choose this path



=4/5 *1/2A

---

**Carnegie Mellon**

# Part 2: DisplayGen

- find subgraph that max's delivered current
- Incrementally, add nodes with max marginal delivered current

---

**Carnegie Mellon**

# Part 3: Scalability

'CandidateGen'
- Starting from the large graph
- Eliminate nodes that are too far away to matter
- How?

---

**Carnegie Mellon**

# Part 3: Scalability

- By successive, careful expansions



source    s          t    sink

---

**Carnegie Mellon**

# Part 3: Scalability

4

**Carnegie Mellon**

# Part 3: Scalability

---

**Carnegie Mellon**

# Part 3: Scalability

---

**Carnegie Mellon**

# Pseudo-code

Until (*stoppingCriterion*)
   use *pickHeuristic( )* to pick a node *n*
   expand node *n*

---

**Carnegie Mellon**

# Pseudo-code

*pickHeuristic( )* favors
- Nearby nodes with
- Strong connections to source or sink and with
- Small degree

---

**Carnegie Mellon**

# Outline

- Introduction / Motivation
- Survey
- Proposed Method
- Algorithms
- ➡ Experiments
- Conclusions

---

**Carnegie Mellon**

# Experiments

- on large real graph
  - ~15M nodes, ~100M edges, weighted
  - 'who co-appears with whom' (from 500M web pages)
- Q1: Quality of 'voltage' approach?
- Q2: Speed/accuracy trade-off?

**Slide 31:**

# Q1: Quality

- Actors (A); Computer-Scientists (CS)
- Kidman-Diaz (A-A)
- Negreponte-Palmisano (CS-CS)
- Turing-Stone (CS-A)

**Slide 32:**

# (A-A) Kidman-Diaz

- What are the best paths between 'Kidman' and 'Diaz'?

Nicole Kidman — Angelina Jolie — Britney Spears — Carmen Electra — Cameron Diaz

Strong, direct link

**Slide 33:**

# CS-CS: Negreponte - Palmisano

NN                                                                 SP

- Mainly: CEOs of major Computer companies (Dell, Gates, Fiorina, ++)

**Slide 34:**

# CS-CS: Negreponte - Palmisano

NN                                                                 SP

Esther Dyson          Louis Gerstner

**Slide 35:**

# CS-A: Turing - Stone

Turing                              Anderson

Alan Turing — Kate Winslet — Harry Potter — Gillian Anderson — Sharon Stone

Stone

**Slide 36:**

# Outline

- Introduction / Motivation
- ...
- Experiments
  - Q1: quality
  → - Q2: speed/accuracy trade-off
- Conclusions

# Speed/Accuracy Trade-off

delivered
current



Kleinberg-Newell
Rivest-Hoffman
Turing-Stone
Kidman-Diaz

number of nodes kept ('*b*')

---

# Speed/accuracy trade-off

- 80/20-like rule: the first few nodes/paths
- contribute the vast majority of 'delivered current'
- Thus: CandidateGen makes sense

---

# Conclusions

- Defined the problem
- Part 1: Electricity-based method to measure quality
- Part 2: Dynamic programming to spot best paths ('DisplayGen')
- Part 3: Scalability with good accuracy ('CandidateGen')
- Operational system