

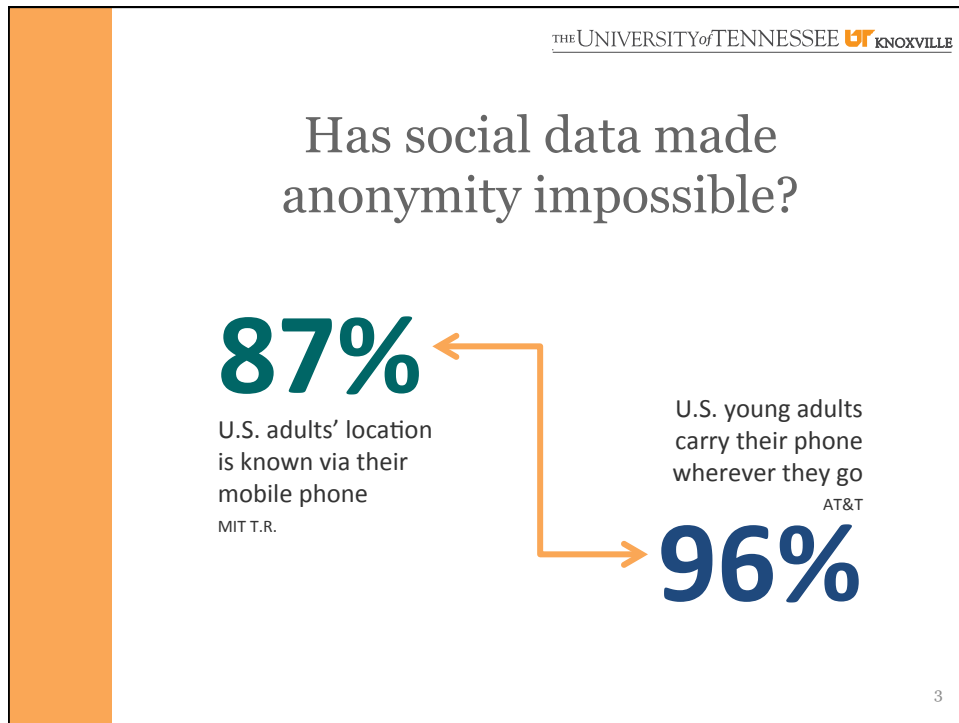
# Social Fingerprinting: Identifying Users of Social Networks by their Data Footprint

The University of Tennessee  
Electrical Engineering and Computer Science  
Dissertation Defense

Denise Koessler Gosnell



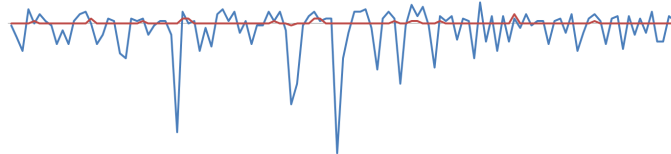
Is a social network  
user's behavior  
unique?



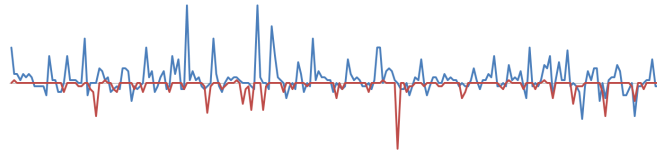
- THE UNIVERSITY of TENNESSEE **UT** KNOXVILLE
- ## Outline:
1. Define and Motivate
  2. Model and Software
  3. Individual Printing Algorithm
  4. Community Printing Algorithm
  5. Conclusions
- 4

# A Social Fingerprint

**Social Network User A**

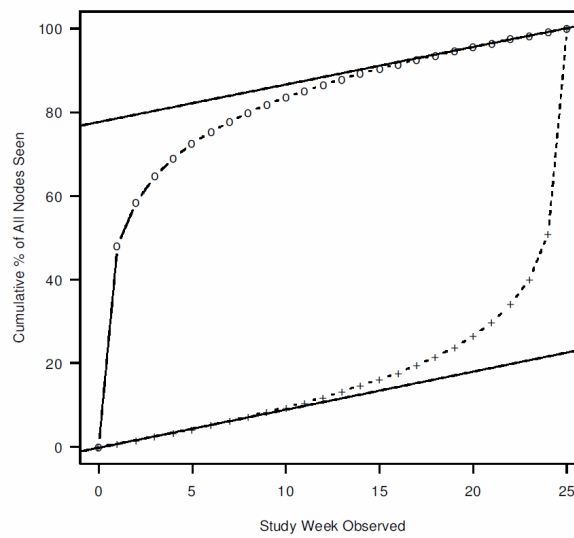


**Social Network User B**



5

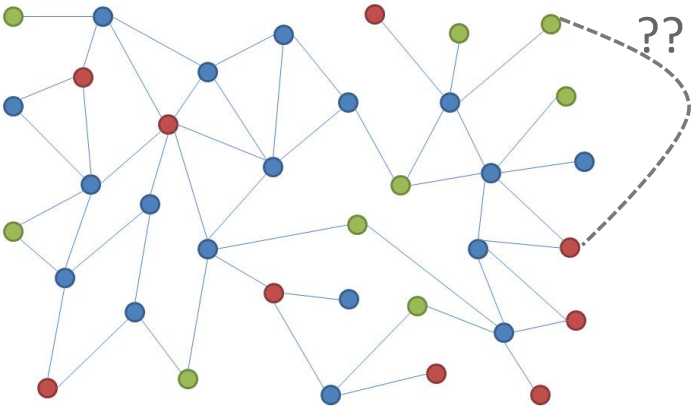
# Why is this hard?



6

THE UNIVERSITY of TENNESSEE **UT** KNOXVILLE

## Why is this hard?



7

THE UNIVERSITY of TENNESSEE **UT** KNOXVILLE

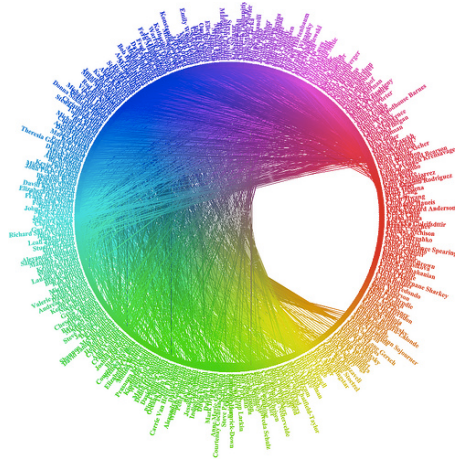
## Previous Work

ACADEMIA	INDUSTRY
<ul style="list-style-type: none"> <li>▪ Bounds of Human Privacy, MIT</li> </ul>	<ul style="list-style-type: none"> <li>▪ IBM Watson SMSim</li> </ul>
<ul style="list-style-type: none"> <li>▪ SNAP Project, Stanford</li> </ul>	<ul style="list-style-type: none"> <li>▪ Topology Statistics</li> </ul>
<ul style="list-style-type: none"> <li>▪ Graph Modeling Libraries</li> </ul>	<ul style="list-style-type: none"> <li>▪ Private Companies</li> </ul>
<ul style="list-style-type: none"> <li>▪ Simulated topologies and statistics</li> </ul>	<ul style="list-style-type: none"> <li>▪ ...?</li> </ul>

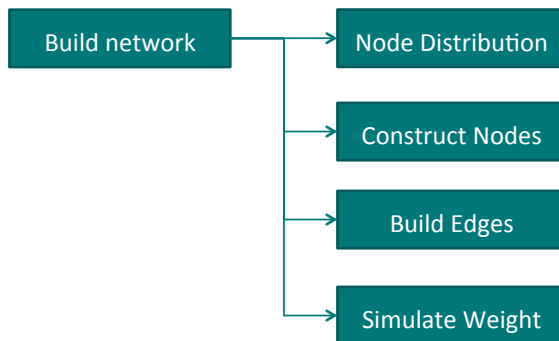
8

# The Model

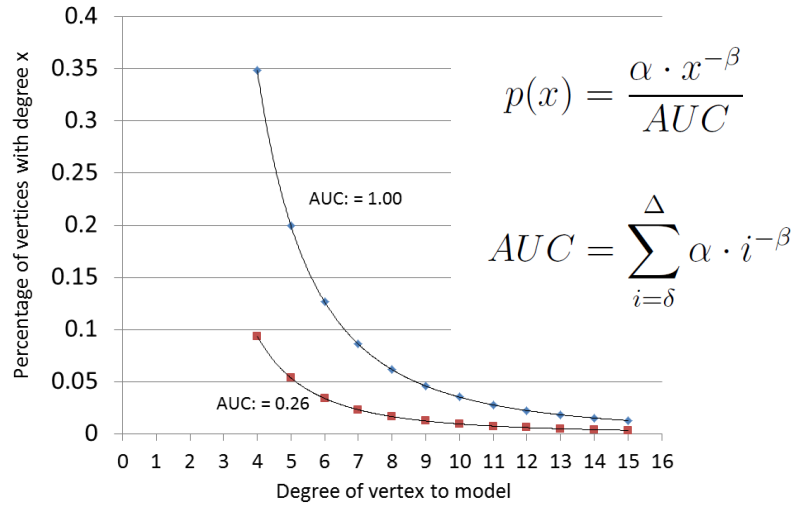
- Connectivity?
- Time?
- Multiple Edges?
- Empirical Data vs. Assumptions?



# The Model



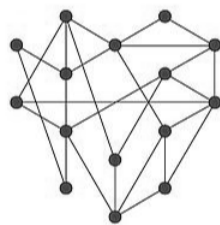
# Scale-free Distribution



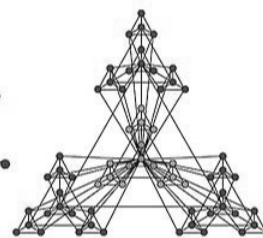
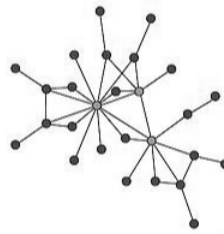
11

# The Base Model

## Scale-Free Graph



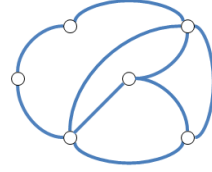
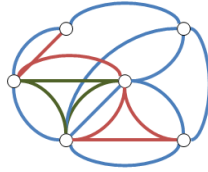
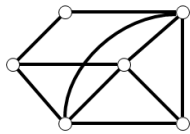
Random Graph



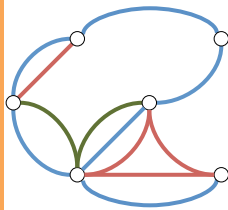
Hierarchical Graph

12

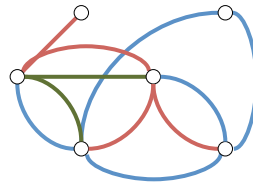
# Multiple Edges



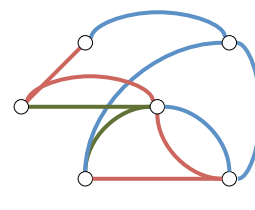
# Dynamic Edges



$G_1$

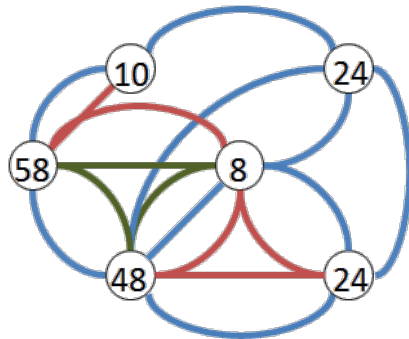


$G_2$



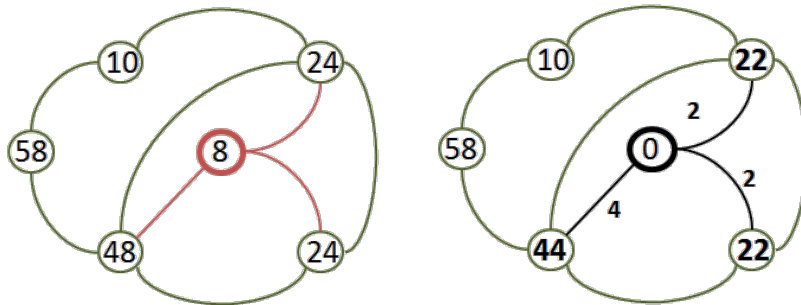
$G_3$

# Simulating Edge Weights



15

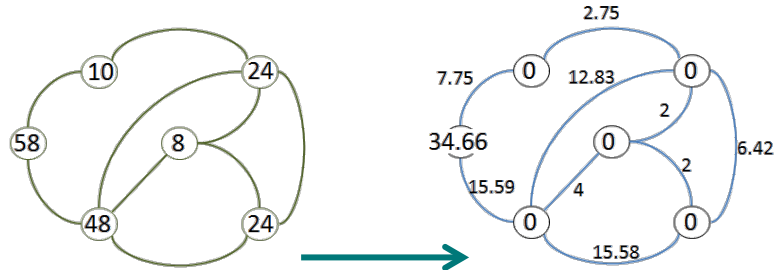
# Maximal Diffusion Algorithm



16



## Maximal Diffusion Algorithm



Capacity constraints: 
$$\sum_{v_j \in N(v_i)} w(e(i, j)) \leq c_{i,b}$$

Flow conservation: 
$$r_{i,b} + \sum_{v_j \in N(v_i)} w(e(i, j)) = c_{i,b}$$

17

## Maximal Diffusion Algorithm

**Result:** A fractional estimation of an optimally weighted graph via the diffusion of node capacity.

**Input:** time step  $t$  and edge type  $b$ ;

Sort all nodes such that  $c_{i,b}^t \leq c_{i+1,b}^t$ ;

for each  $v_i \in V(G)_{\text{sorted}}$  do

    Initialize:  $r_{i,b}^t = c_{i,b}^t$ ;

    Determine total community capacity  $C$  of  $v_i$ ;

$$C = \sum_{v_j \in N(v_i)} c_{j,b}^t;$$

    where  $N(v_i) = \{v_j \in V(G^t) | e(i, j)_b^t \in E(G^t)\}$ ;

    for each  $v_j \in N(v_i)$  do

$$w = c_{i,b}^t \cdot \frac{c_{j,b}^t}{C};$$

$$e(i, j)_b^t = w;$$

$$c_{j,b}^t = c_{j,b}^t - w;$$

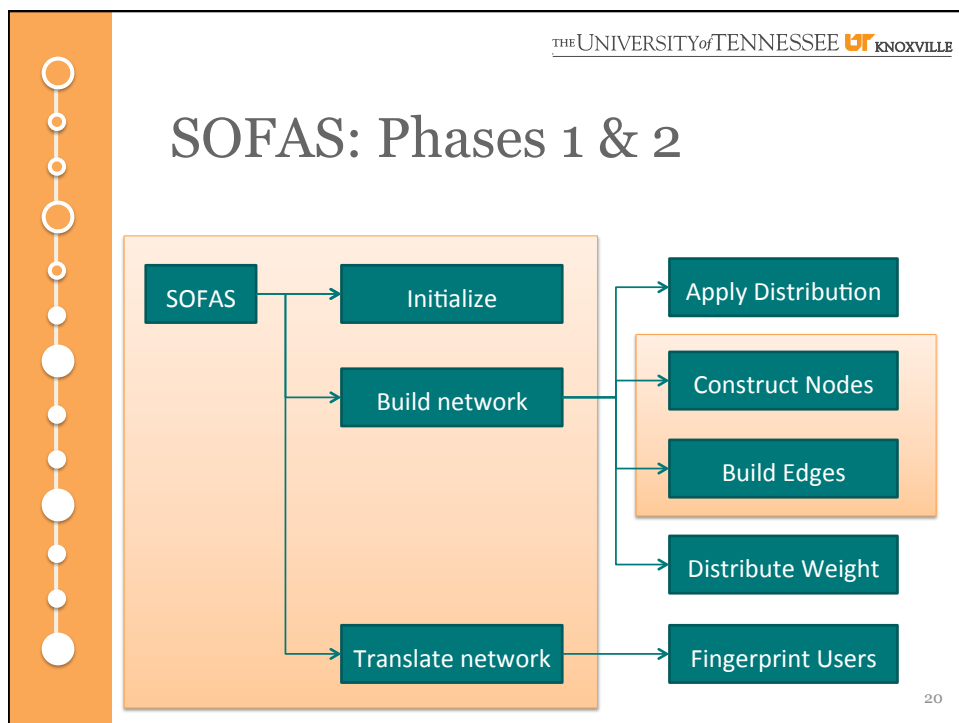
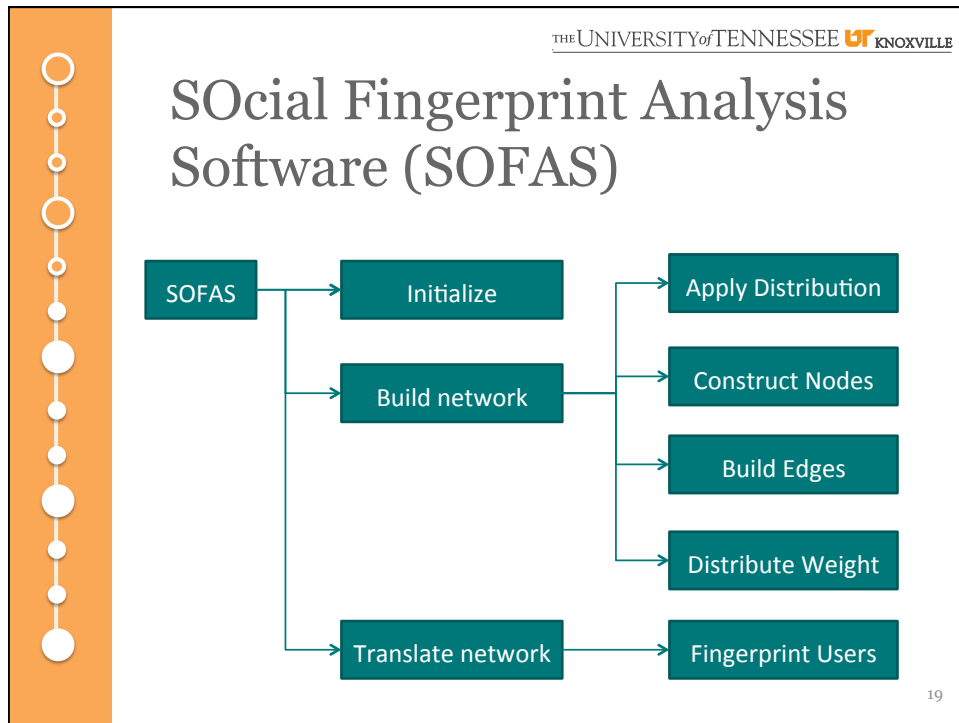
$$r_{i,b}^t = r_{i,b}^t - w;$$

    end

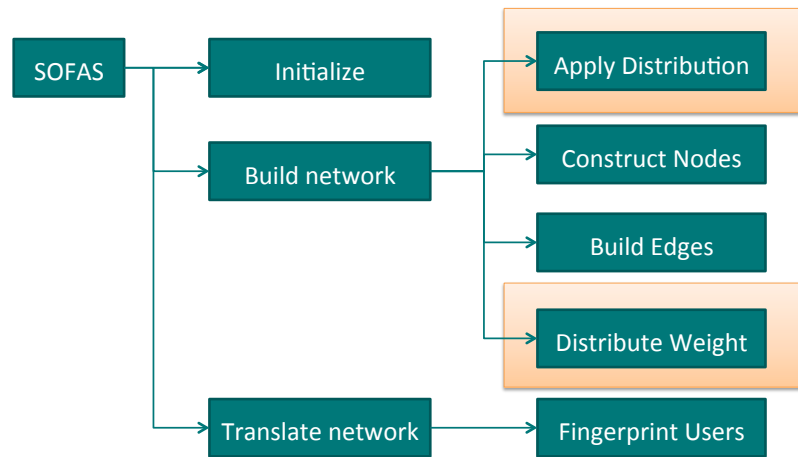
    Sort all nodes starting from  $v_i$  such that  $c_{i,b}^t \leq c_{i+1,b}^t$ ;

end

18

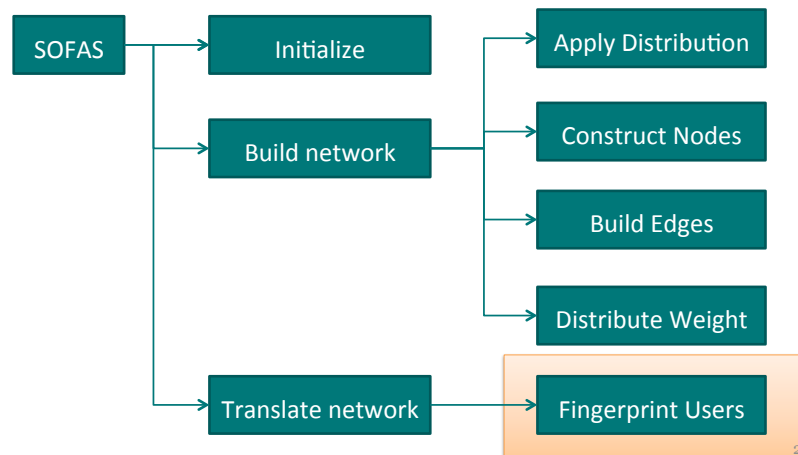


## SOFAS: Phases 3 and 4



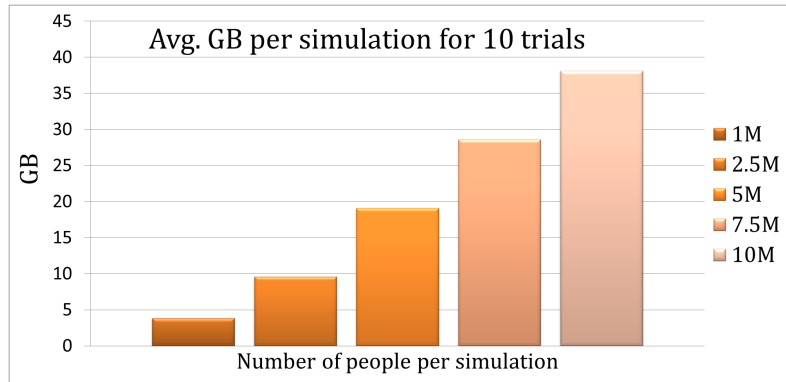
21

## SOFAS: Phase 5



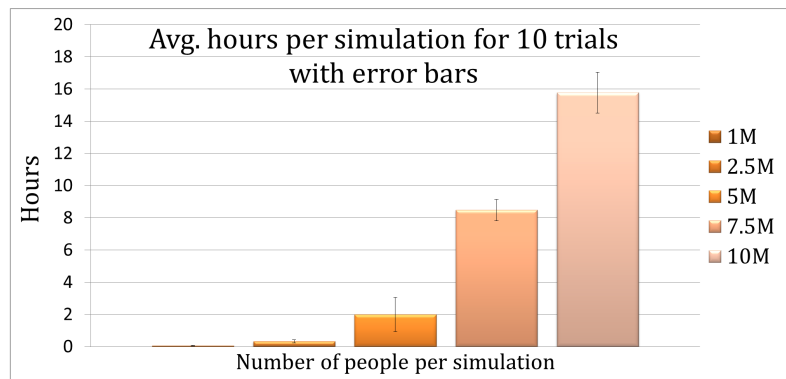
22

# SOFAS: Memory Performance



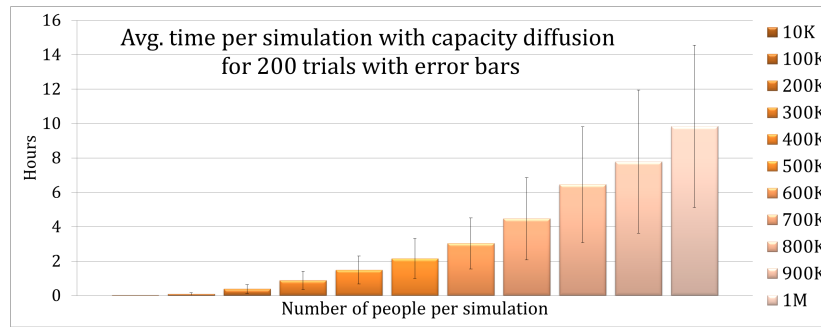
23

# SOFAS: Time



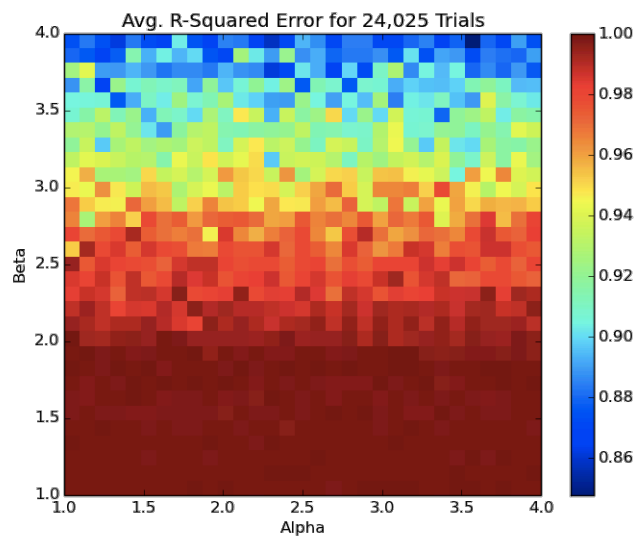
24

## SOFAS: Diffusion Time



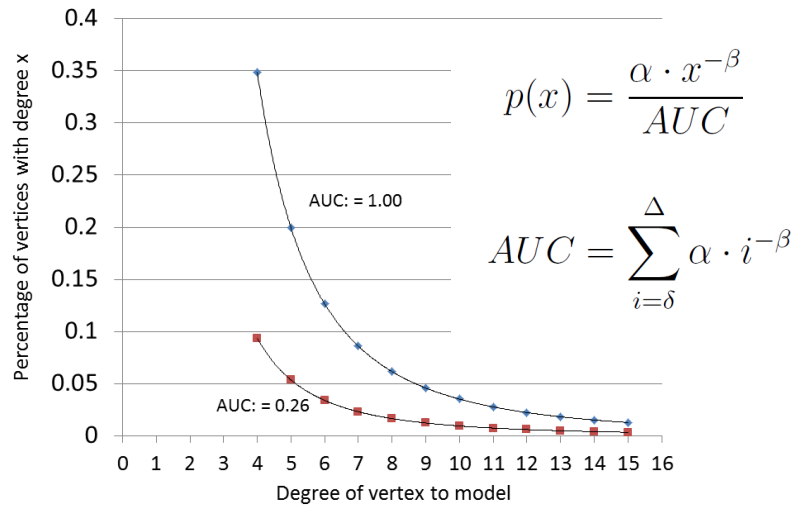
25

## SOFAS: Distribution Error

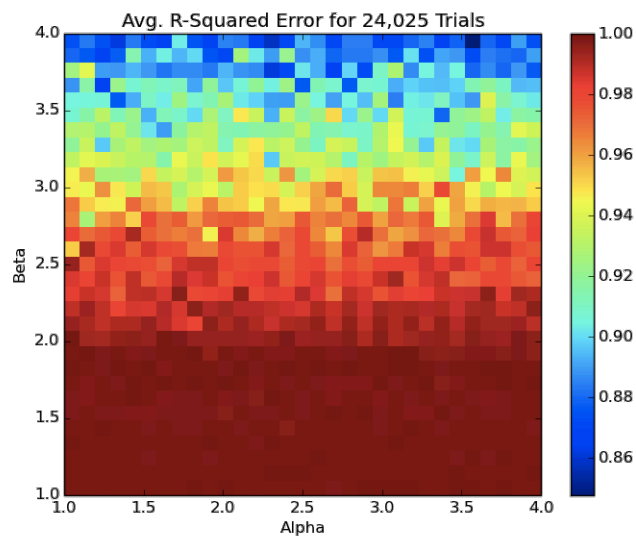


26

## Scale-free Distribution

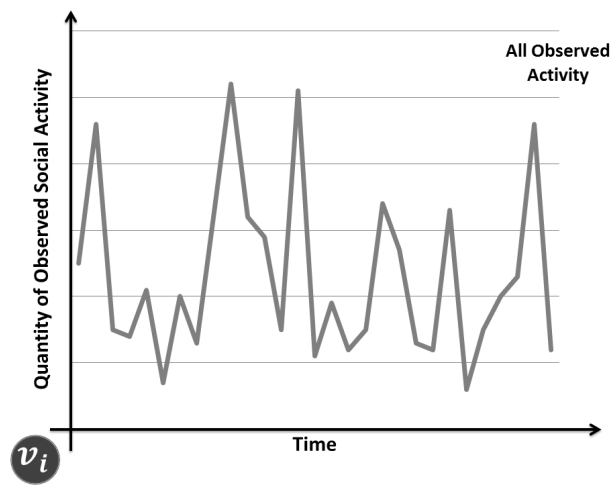


## SOFAS: Distribution Error

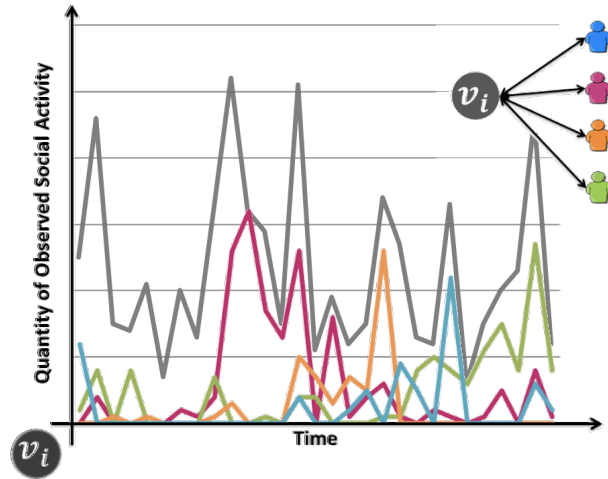


# Social Fingerprinting: Identifying the Individual

## The Theory

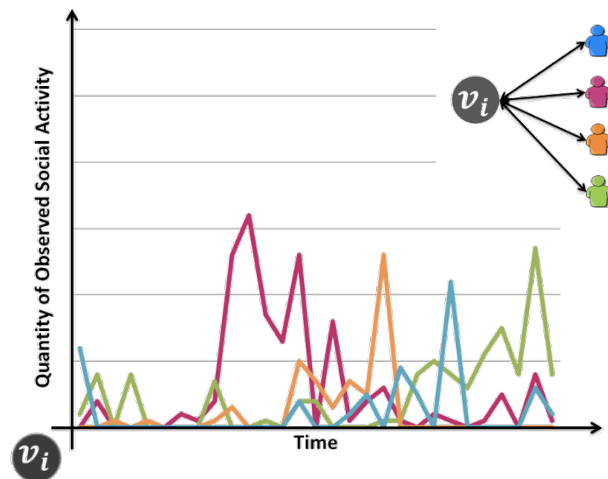


# The Theory



31

# The Theory



32

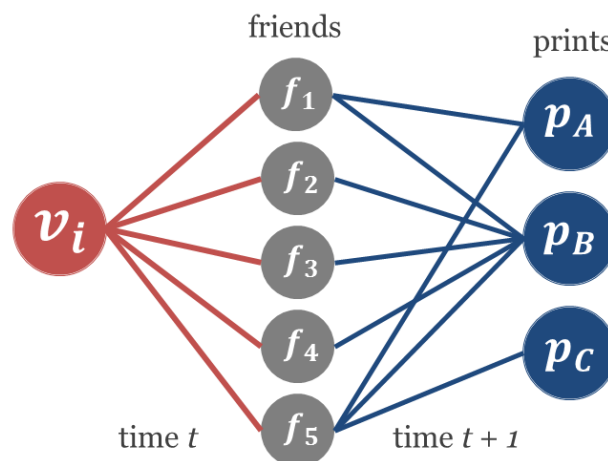


## Social Fingerprint Procedure:

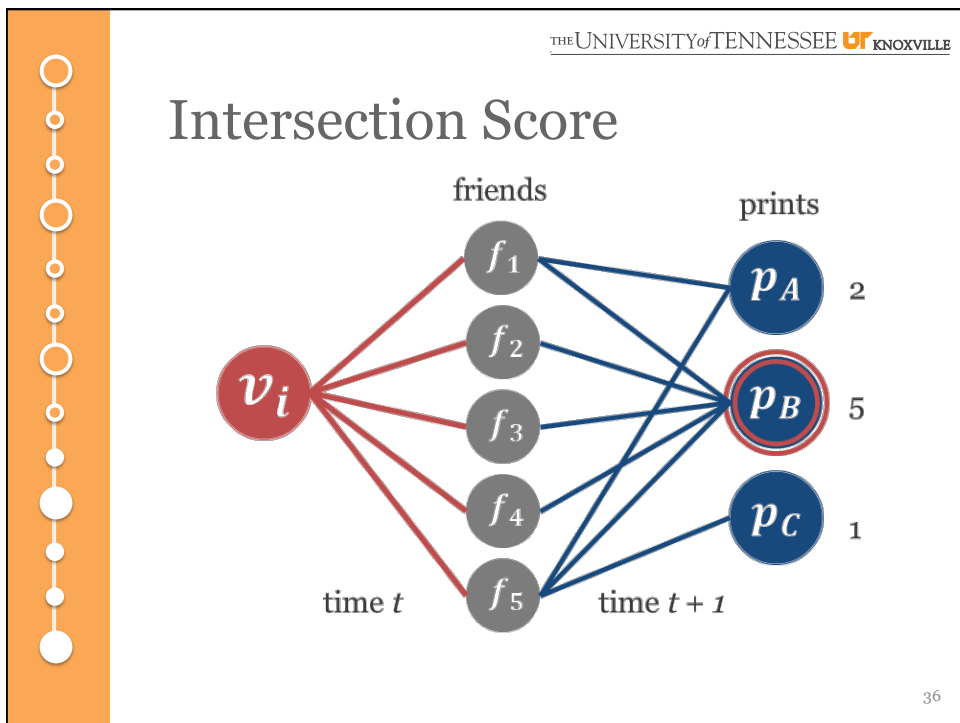
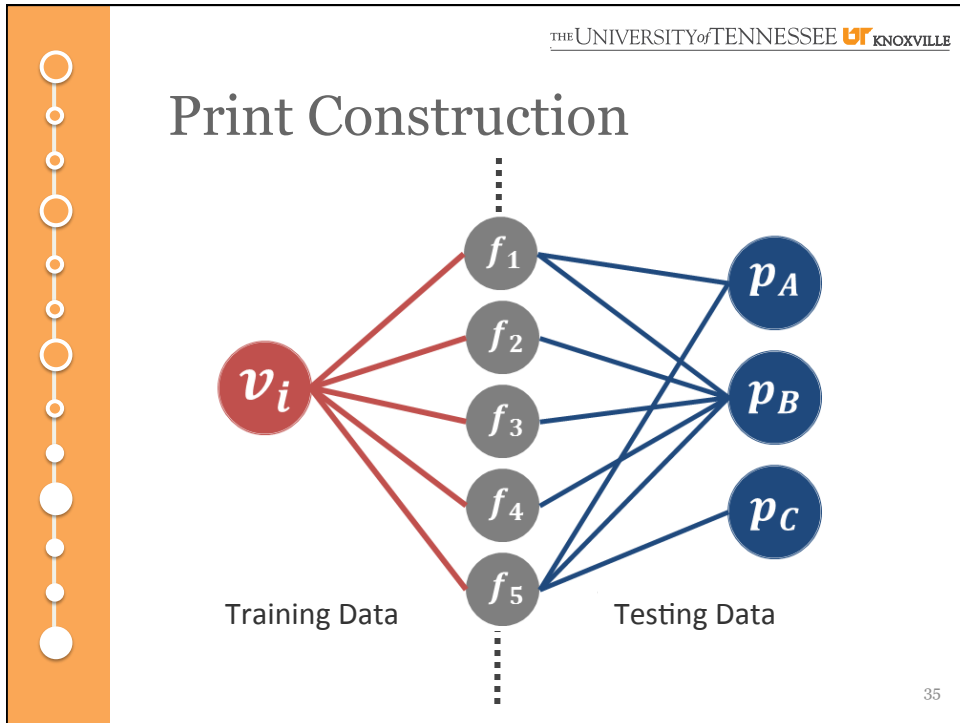
1. For each person, construct a neighborhood graph
2. Given the neighborhood graph, determine candidate prints
3. Rank the prints
4. Evaluate

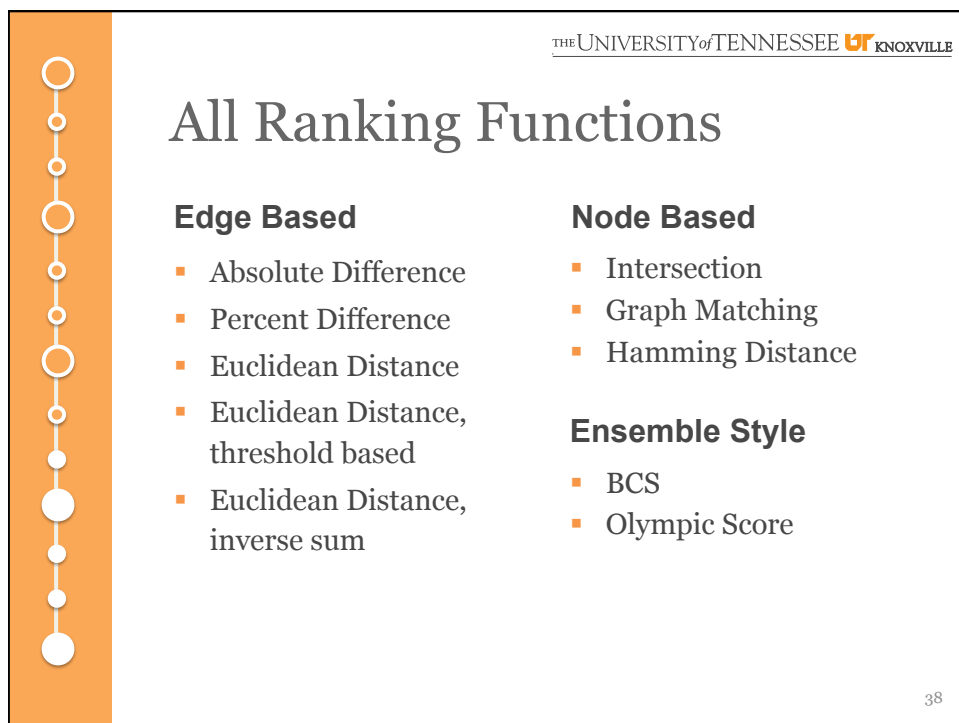
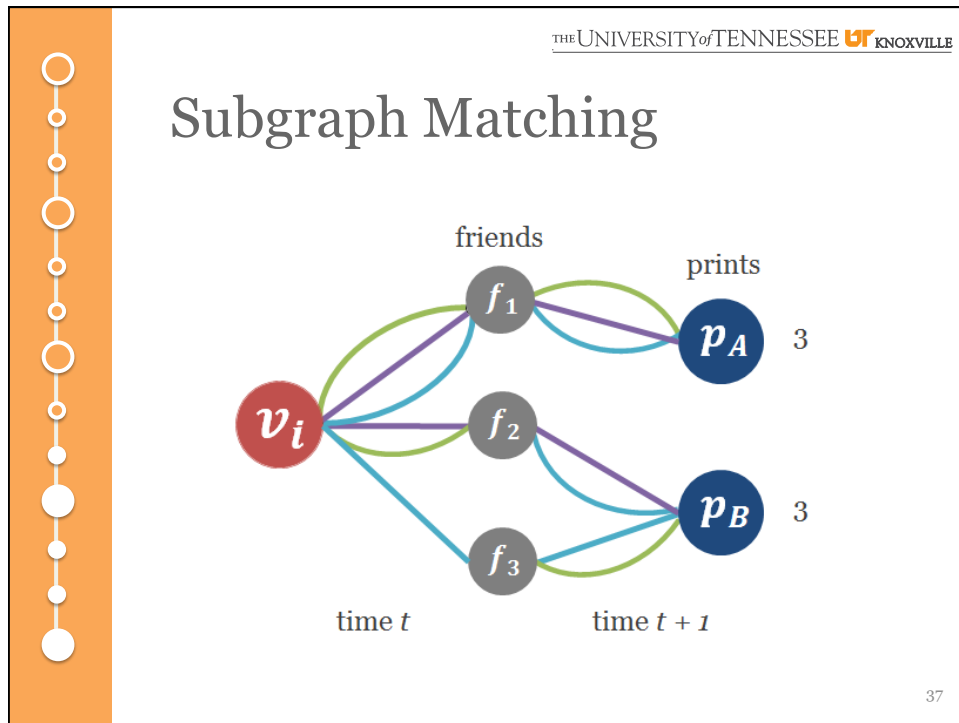
33

## Print Construction



34



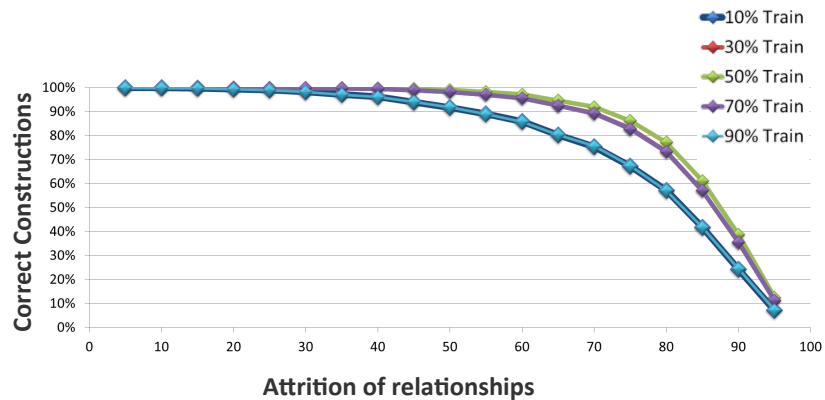


## Graph Prints: Test Cases

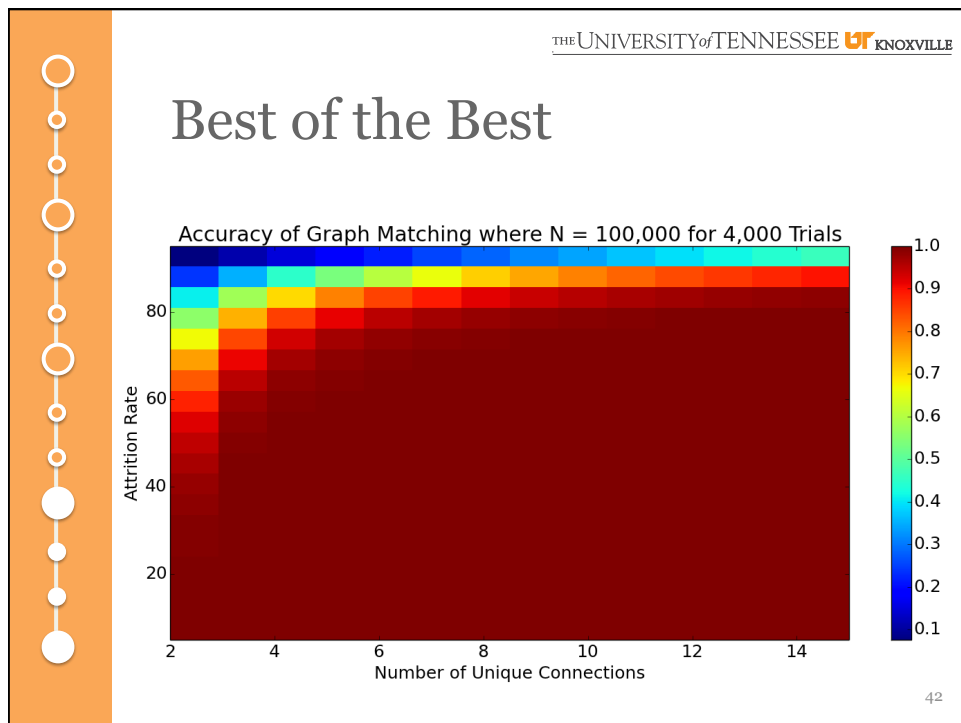
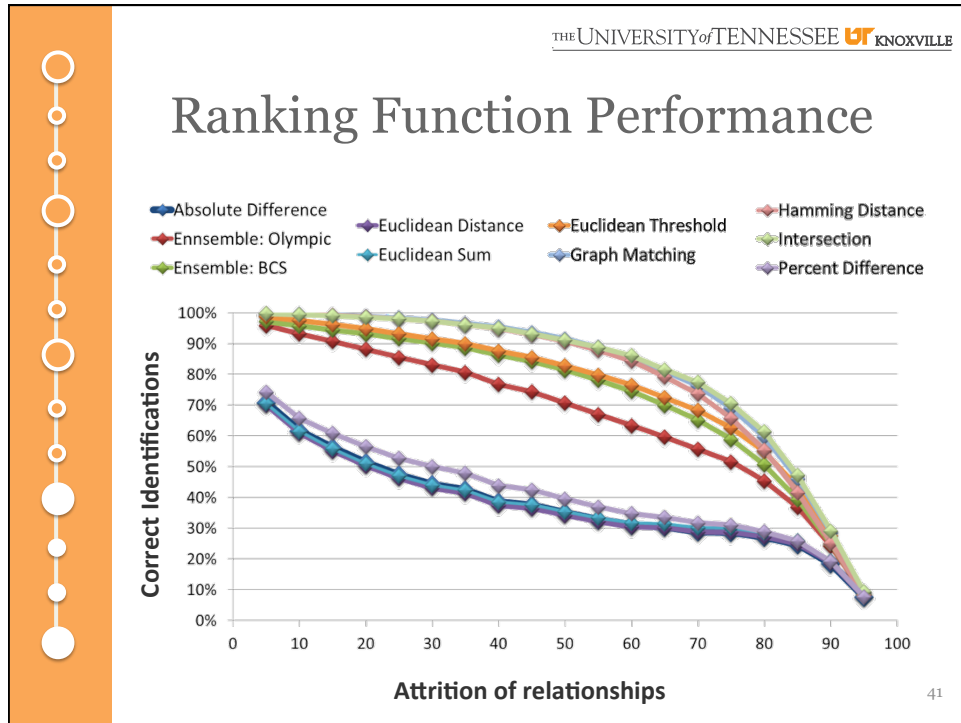
- Number of People: 1K, 10K, 100K
- Length of time: 6
- Attrition on relationships: varies
- Training/Testing Split: varies
- Model Types: Binary and Weighted
- Total trials: 57,000 example graphs

39

## Validation of Approach



40



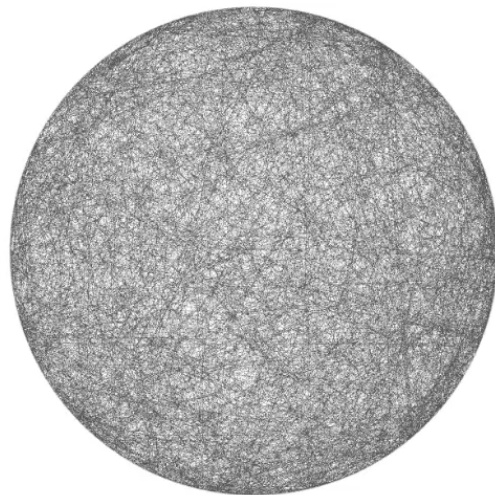


## Social Fingerprinting: Picking the individual out of crowd

43



## What does this look like?



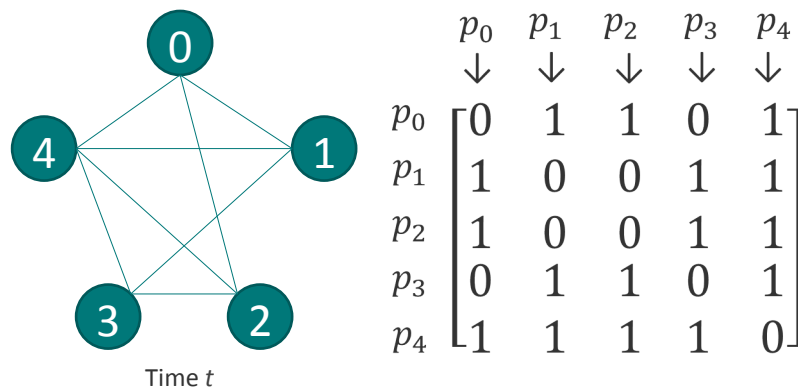
44

## Social Fingerprint Procedure:

1. Construct matrix A and query vector(s)
2. Semidiscrete Decomposition of matrix A to yield rank- $k$  approximation
3. Compute new query vectors
4. Rank the vectors w.r.t. cosine similarity
5. Evaluate

45

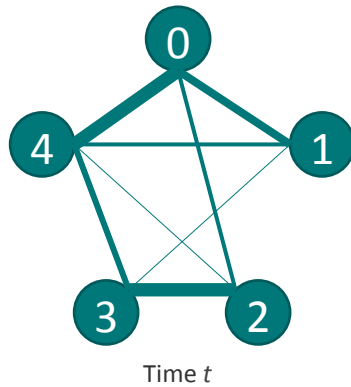
## Construction:



**Binary:**  $A[i,j] \rightarrow$  the presence of the relationship

46

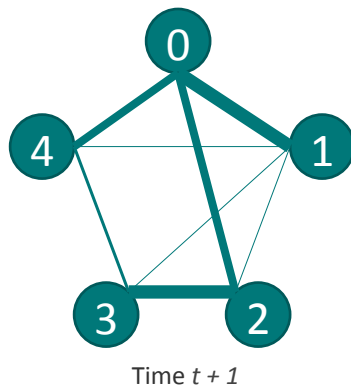
### Construction:



	$p_0$	$p_1$	$p_2$	$p_3$	$p_4$
	↓	↓	↓	↓	↓
$p_0$	0	5	2	0	9
$p_1$	5	0	0	1	2
$p_2$	2	0	0	9	1
$p_3$	0	1	9	0	4
$p_4$	9	2	1	4	0

**Term Frequency:**  $A[i,j] \rightarrow$  the *strength* of the relationship

### Construction: Query Vectors



	$q_0$	$q_1$	$q_2$	$q_3$	$q_4$
	↓	↓	↓	↓	↓
	0	9	1	0	3
	9	0	3	1	1
	4	1	0	9	0
	0	1	9	0	2
	3	1	0	2	0



## Semidiscrete Decomposition (SDD)

[Kolda and O'Leary 1998]

$$A_{n \times n} \approx U_{n \times k} \cdot \Sigma_{k \times k} \cdot V_{n \times k}$$

$$U_{i,j} \in \{-1, 0, 1\}$$

$$\Sigma_{i,i} = \sigma_i$$

$$V_{i,j} \in \{-1, 0, 1\}$$

49

## Validation: $q^{t+1}[j] * V^{(t)}[i]$

	V[0]	V[1]	V[2]	V[3]	V[4]
q[0]	<b>0.8467</b>	0	0	0.5319	0
q[1]	0.0704	<b>0.9859</b>	0.9859	0.1516	0.9859
q[2]	0.2095	<b>0.9778</b>	<b>0.9778</b>	0	<b>0.9778</b>
q[3]	0.2454	0	0	<b>0.9693</b>	0
q[4]	0.1414	<b>0.9899</b>	<b>0.9899</b>	0	<b>0.9899</b>

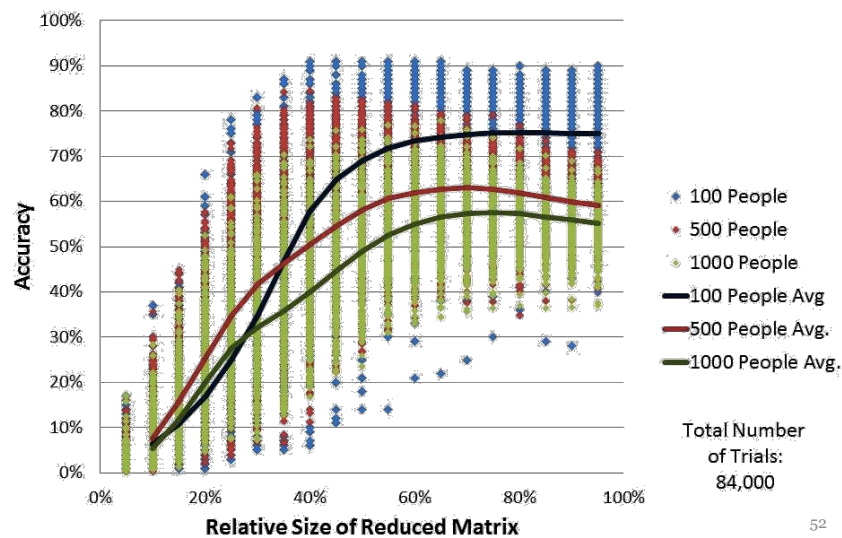
50

## SDD: Test Cases

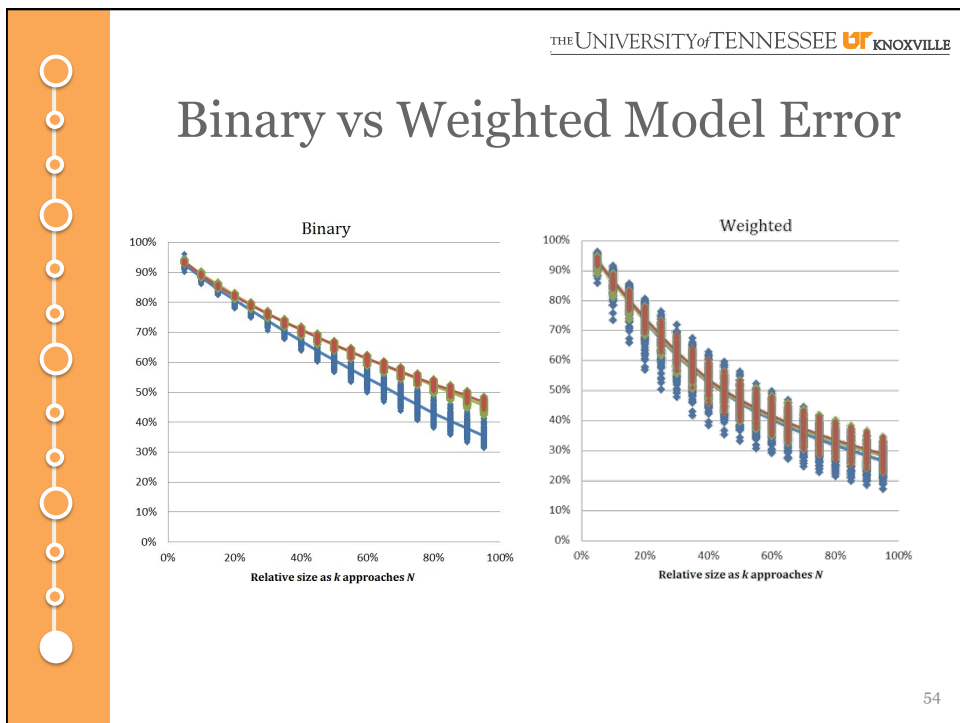
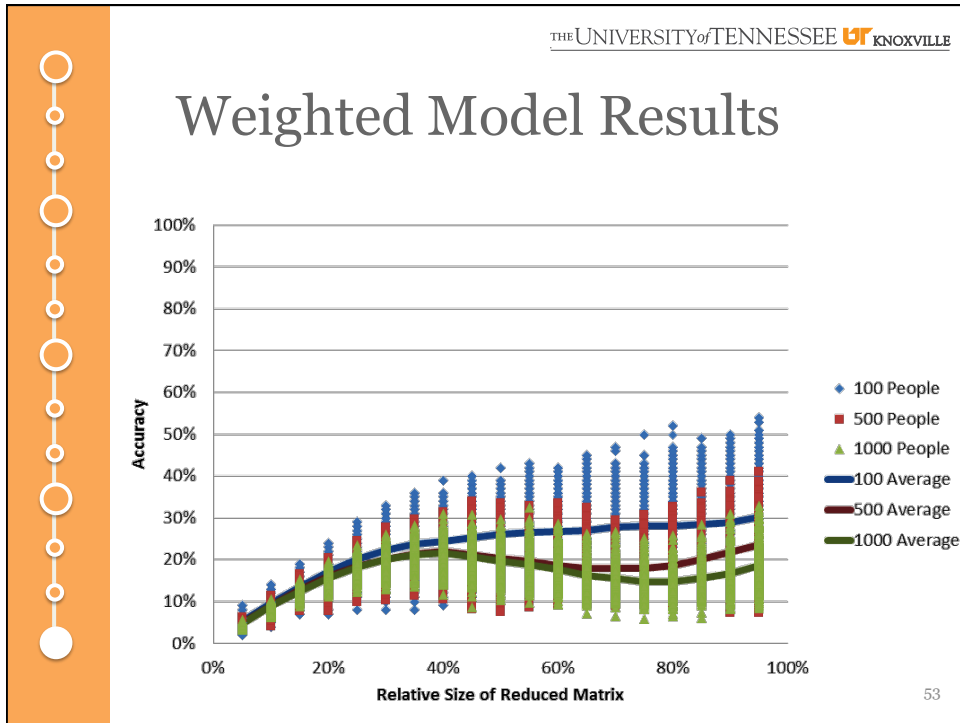
- Number of People: 100, 500 or 1,000
- Length of time: 12
- Attrition on relationships: 37%
- Low Rank Approximation rank: varies
- Training/Testing Split: varies
- Total trials: 84,000

51

## Binary Model Results



52



## SDD: Case Study

- Number of People: 100
- Length of time: 12
- Attrition on relationships: 37%
- Low Rank Approximation rank: 75%
- Training/Testing Split: 50/50

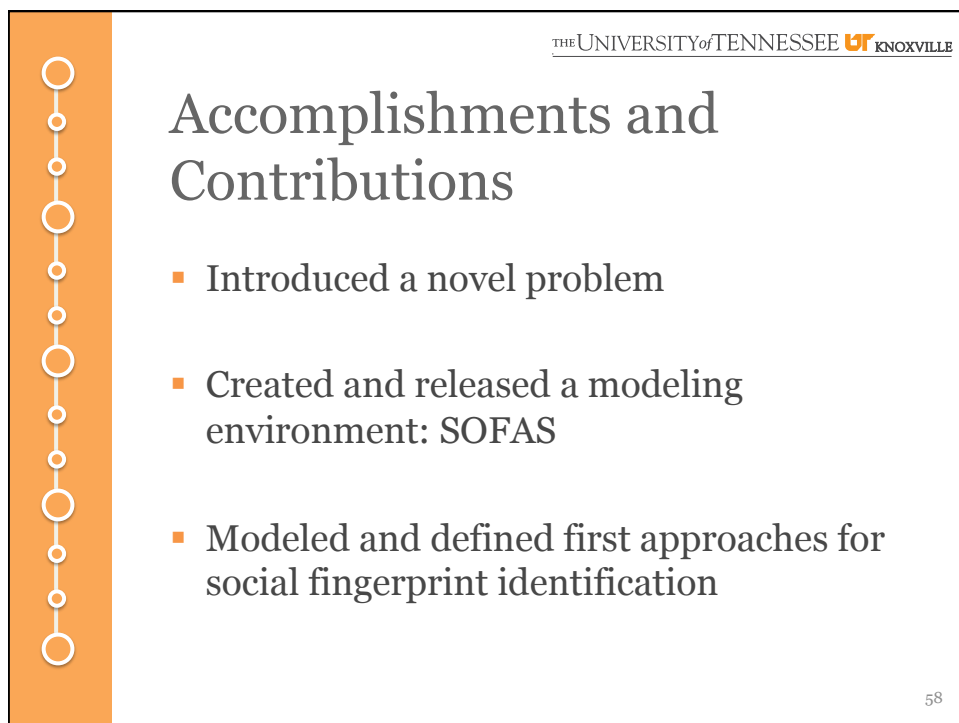
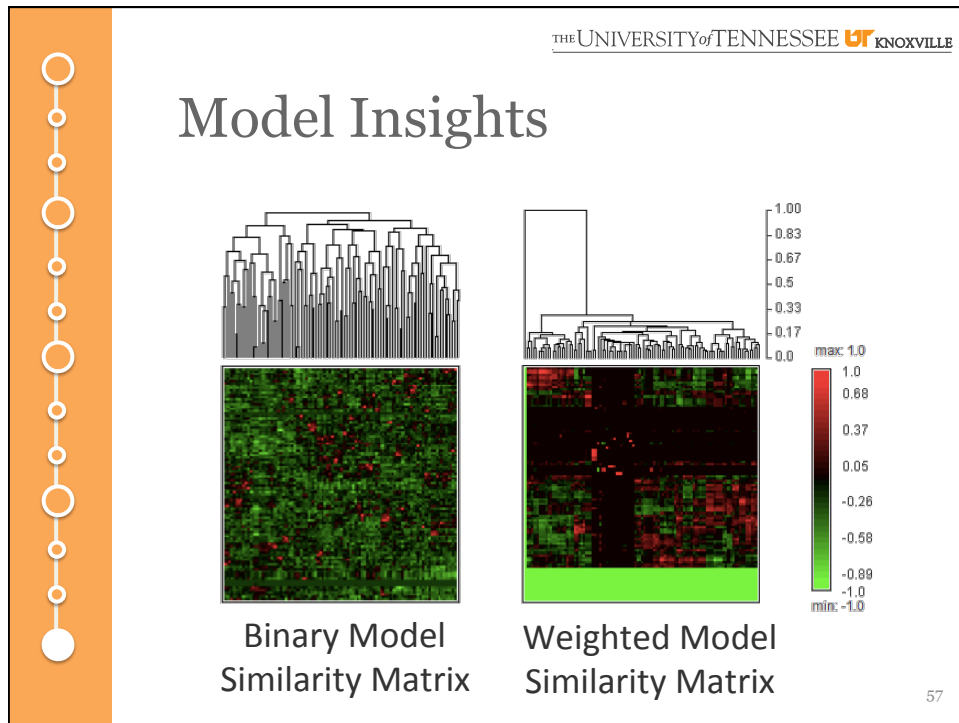
55

## SDD: Case Study Results

- Binary Model Accuracy: 81%
- Weighted Model Accuracy: 27%

	WM: Incorrect	WM: Correct
BM: Incorrect	12	7
BM: Correct	61	20

56



## Open Research

- New models for social fingerprint detection
- Additional features and graph statistics of the SOFA software
- Fingerprint Evasion

## Acknowledgements

Dr. Michael Berry

Dr. Judy Day

Dr. Jens Gregor

Dr. Bruce MacLennan



**SCALE-IT**

Scalable Computing and  
Leading Edge Innovative Technologies



THE UNIVERSITY of TENNESSEE  KNOXVILLE

# Social Fingerprinting: Identifying Users of Social Networks by their Data Footprint

The University of Tennessee  
College of Engineering Dissertation Defense  
Dr. Denise Koessler Gosnell, PhD Computer Science  
(Nerdery)