# The Missing Models: A Data-driven Approach to Learning How Networks Grow

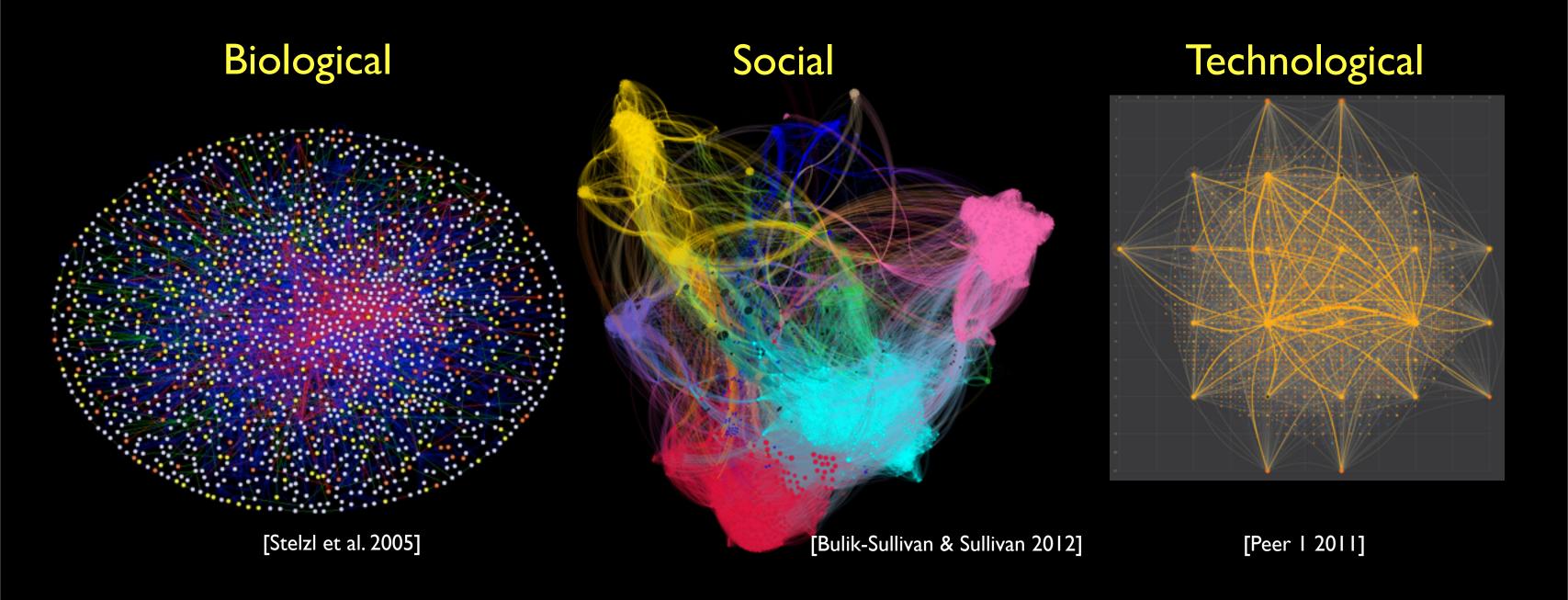
### **Carl Kingsford**

Professor
Computational Biology Department
School of Computer Science
Carnegie Mellon University

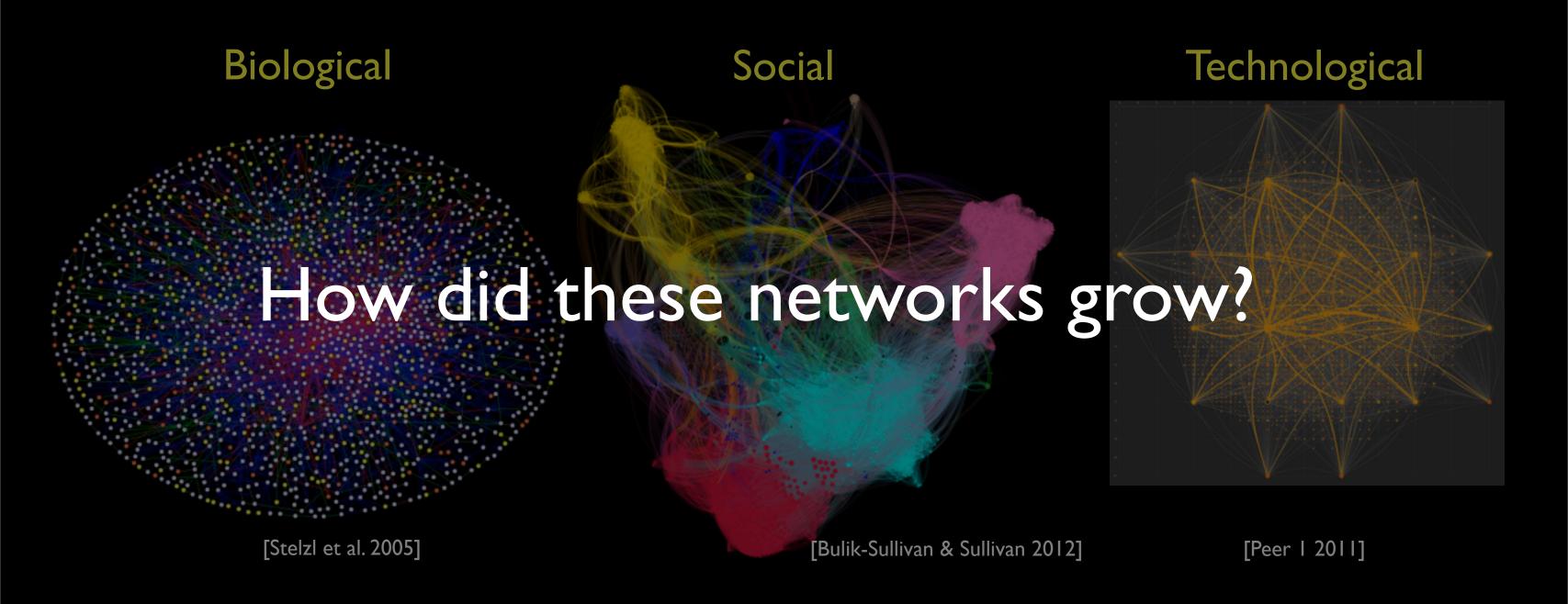
Robert Patro, Geet Duggal, Emre Sefer, Hao Wang, Darya Filippova, Carl Kingsford (2012). The missing models: A data-driven approach for learning how networks grow. *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 42-50.



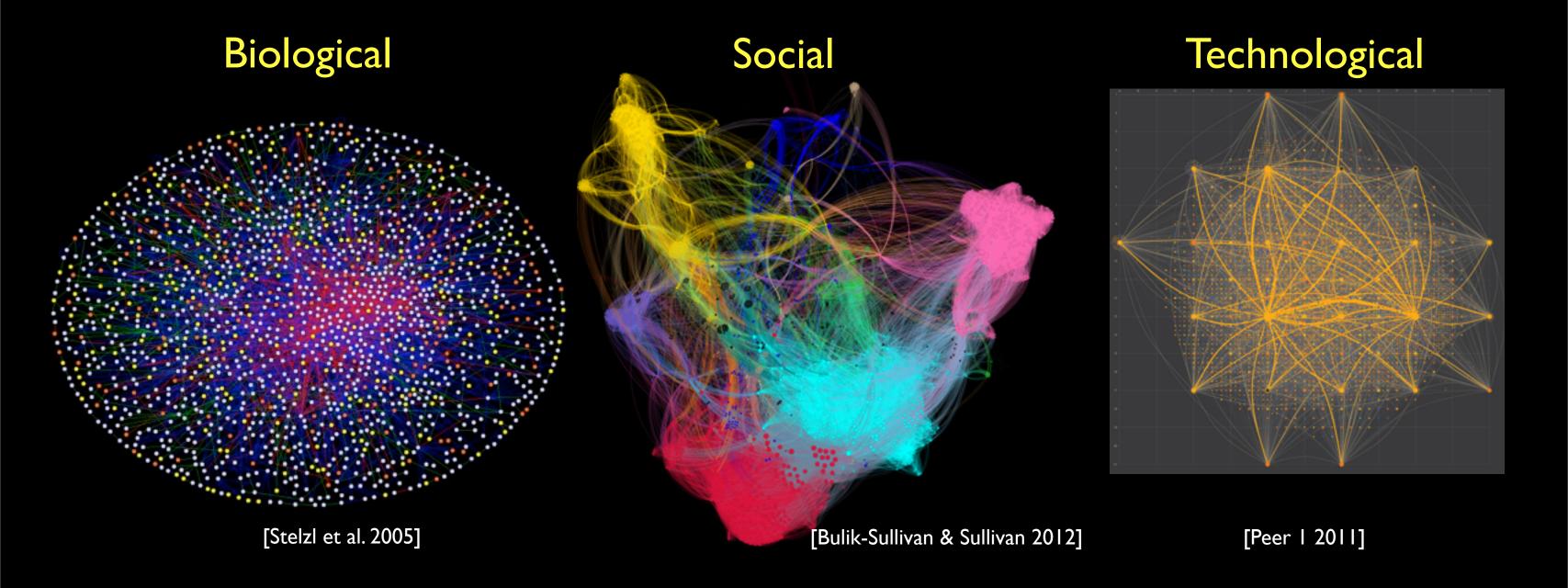
# Networks are everywhere



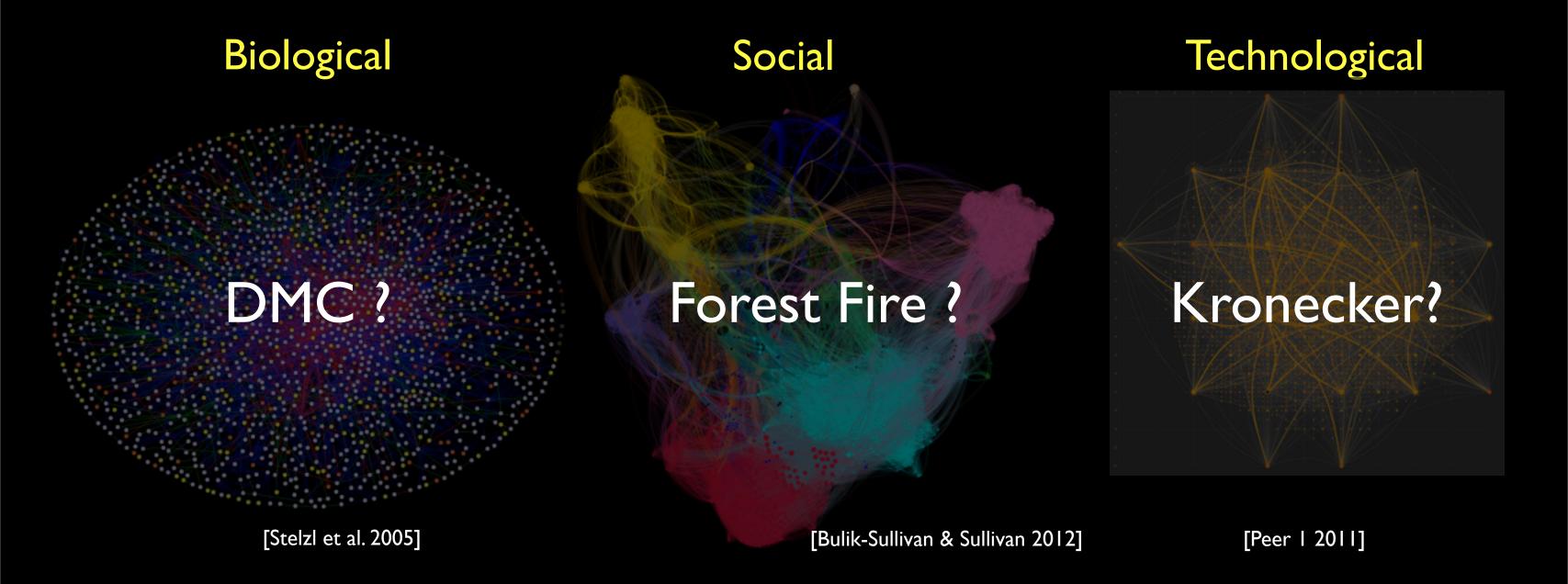
# Networks are everywhere



# Enter Network Growth Models



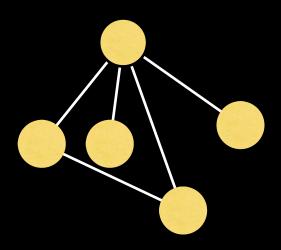
# Enter Network Growth Models



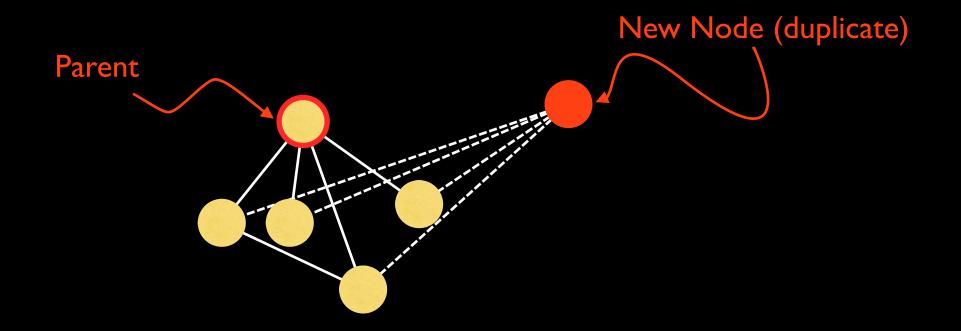
# Example: DMC Model

Plausible model of protein interaction network growth introduced by Vazquez et al. in 2001

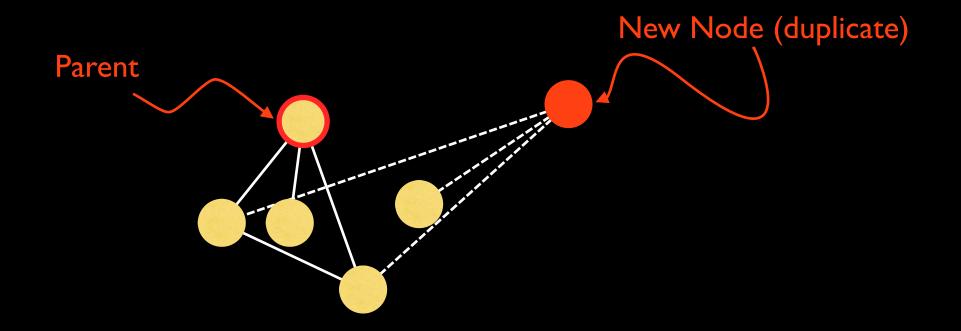
Based on gene duplication & divergence



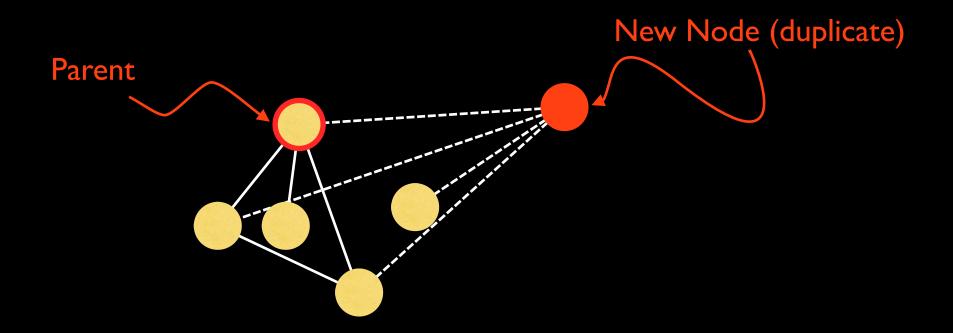
Network at time t



[Duplication, Mutation, Complementarity]

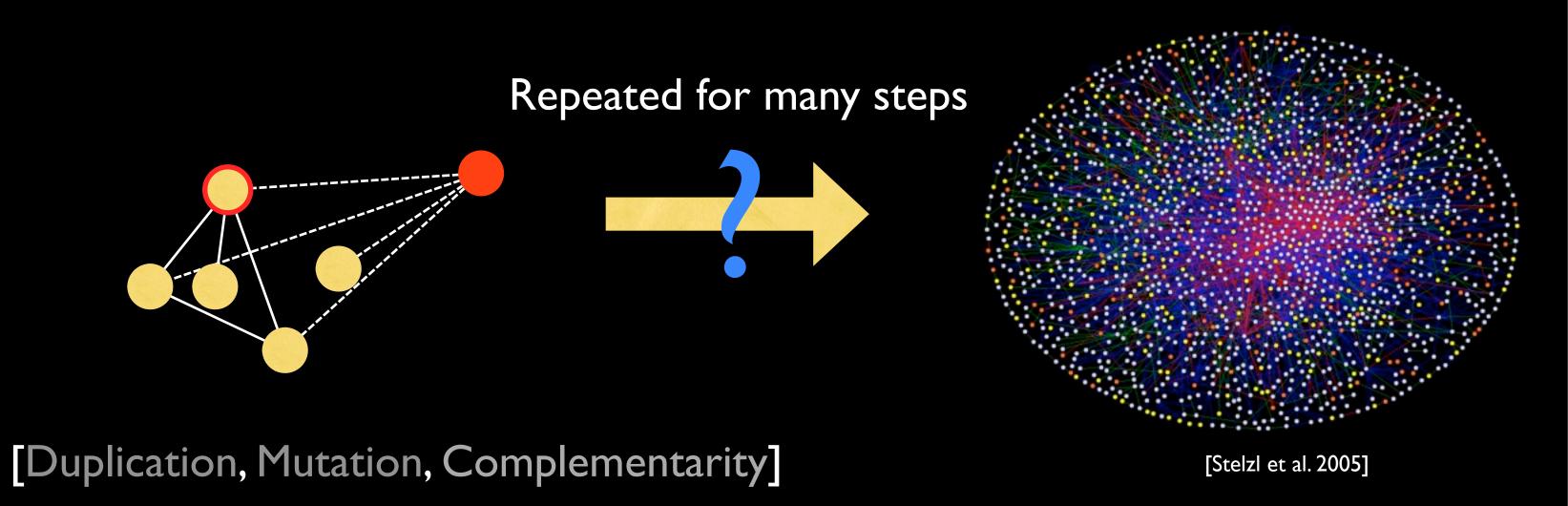


[Duplication, Mutation, Complementarity]



[Duplication, Mutation, Complementarity]

Network at time t+1



In addition to biologically plausible mechanism, can produce networks with similar degree distribution and clustering coeff. as real PPIs

# Network Growth Models (NGMs)

"What I Cannot Create, I Do Not Understand" -- Richard Feynman

Bottom-up generative model of network growth process

Creates "random" graphs with similar topological characteristics to the target

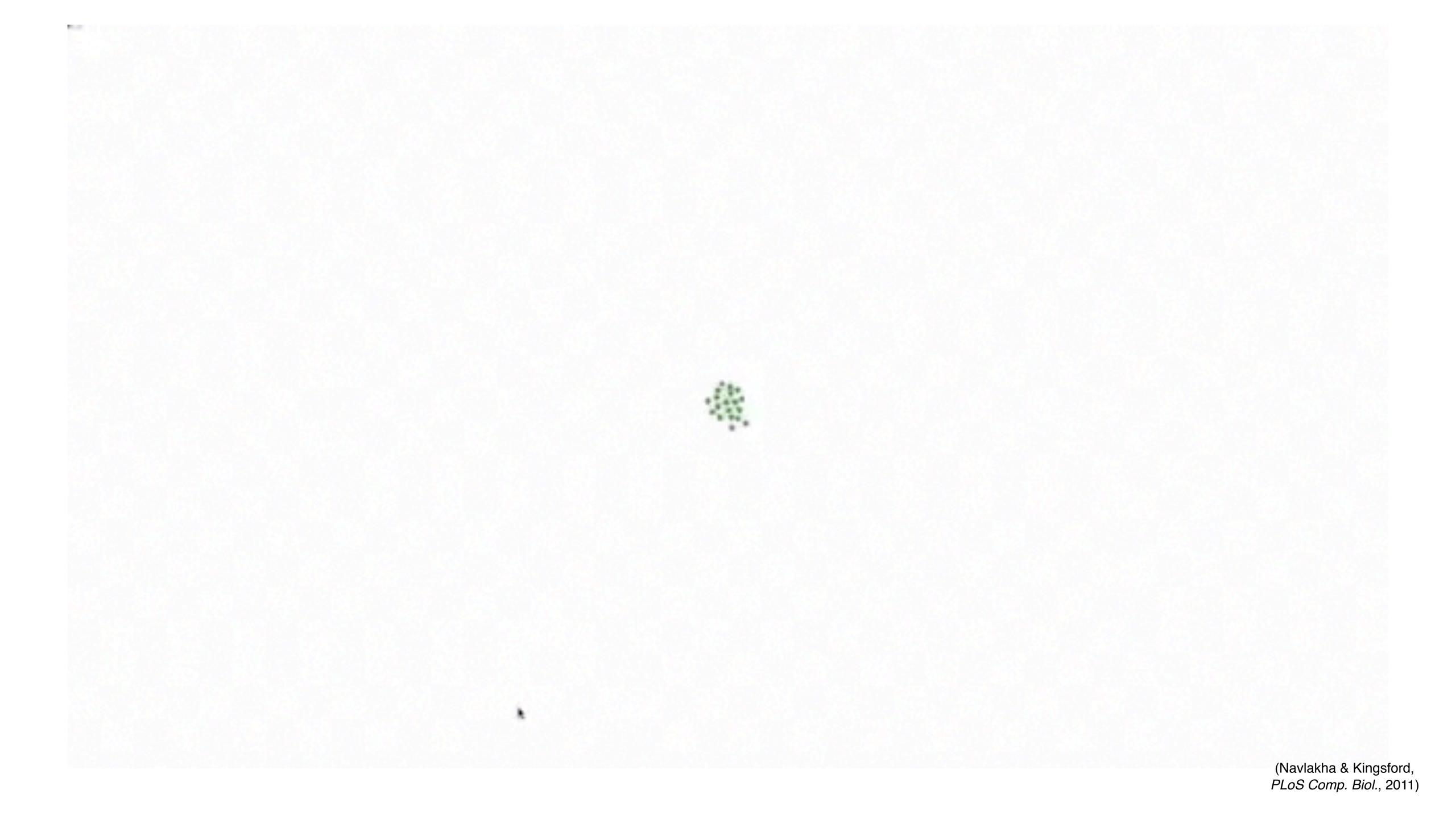
#### **Practical**

Theoretical

Evaluate statistical significance of observed features Test algorithms in different contexts

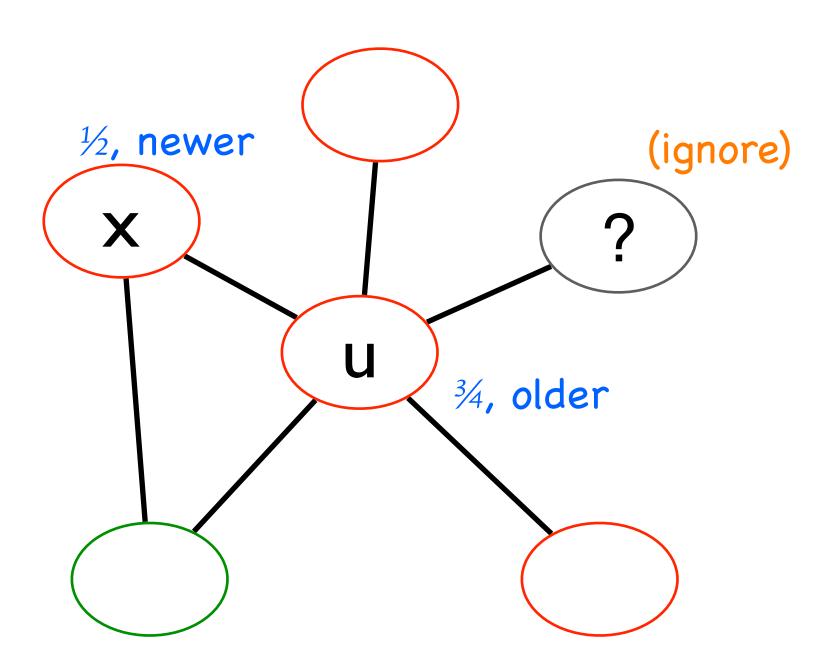
- Varying topological characteristics
- -Varying scales

Discover reasons for observed structure How does topology change over time? How did the network look in the past? How will it look in the future?



# Core vs. Peripheral Complex Members

Coreness of a protein = percentage of like-annotated neighbors

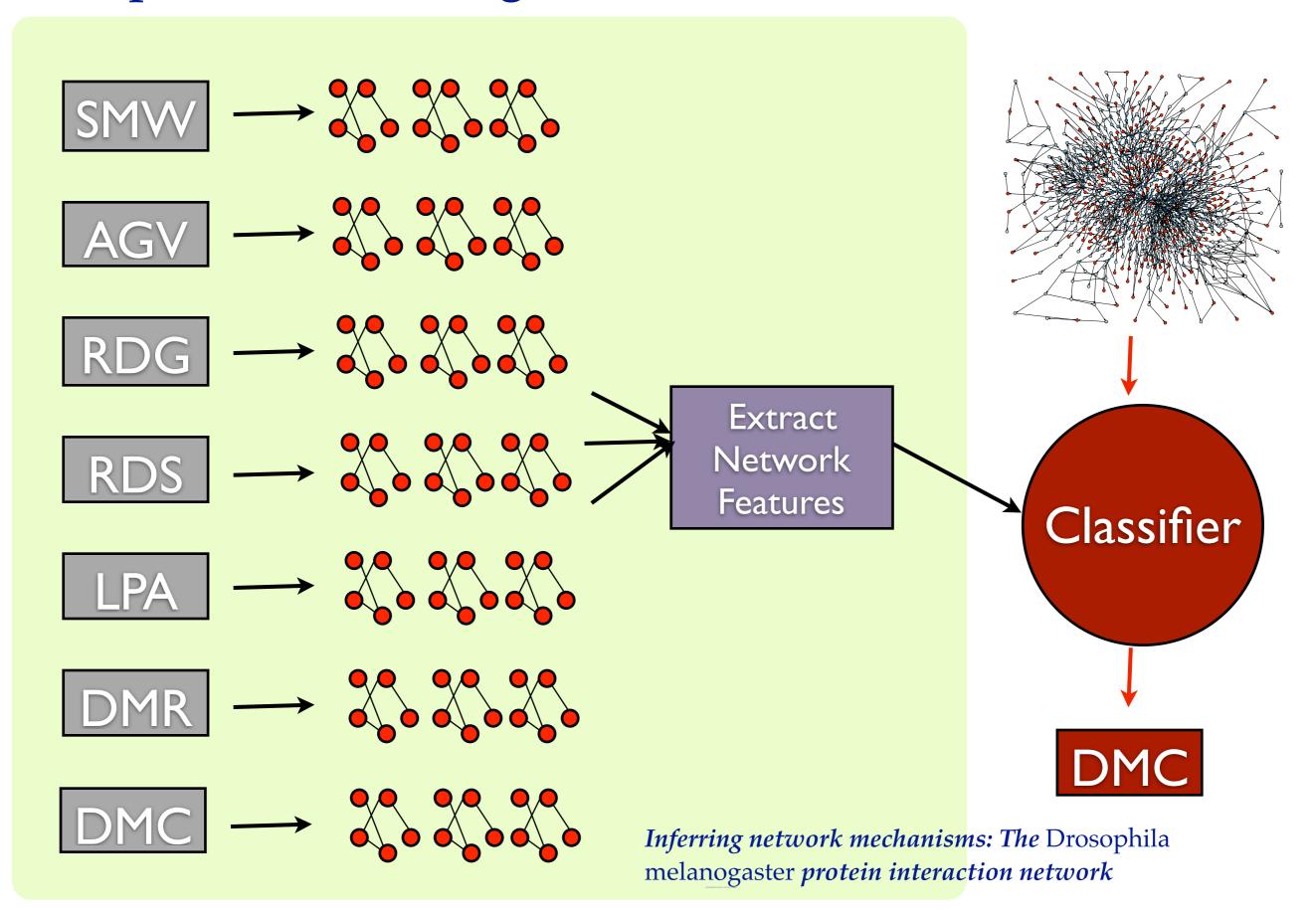


Are core members of a protein complex older than peripheral members?

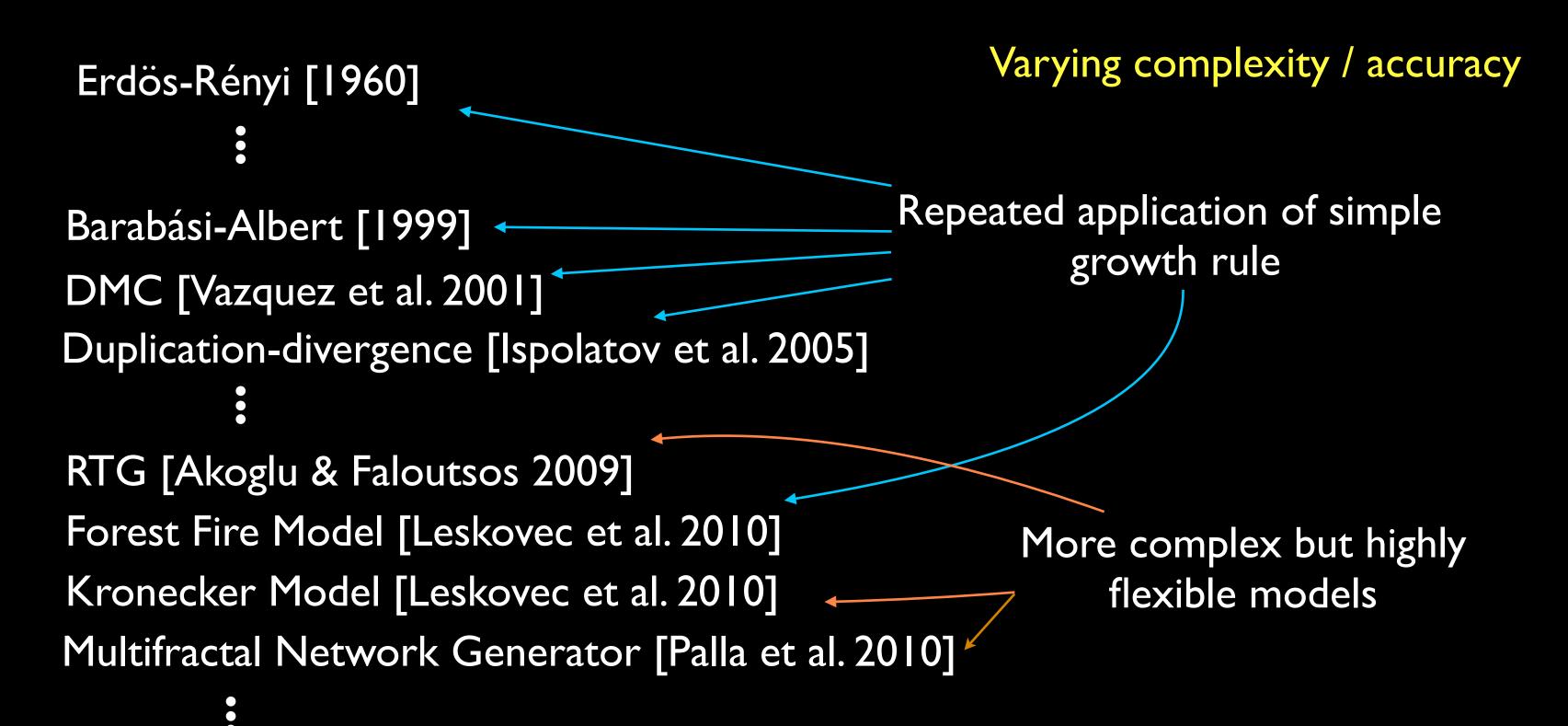
Yes, somewhat: R = 0.37, P < 0.01

Agrees with 3D protein structure analysis (Kim & Marcotte, 2008) looking at age distribution of domains among eukaryotic species.

### Supervised Learning → Predict Network Models



# Many Existing Growth Models



# Many Existing Growth Models

```
Varying complexity / accuracy
Erdös-Rényi [1960]
                      Previous work focused on either
                                              Repeated application of simple
Barabási-Albert [1999]
DMC [Vazquez et al. 2 Manually designed growth models growth rule
Duplication-divergence [Ispolatov et al. 2005]
                                     or
RTG [Akoglu & Faloutsos 2009]
Forest Parameterized family of models (possibly with parameter learning) highly
Kronecker Model [Leskovec et al. 2010]
                                                         flexible models
Multifractal Network Generator [Palla et al. 2010]
```

## So What's New?

Method to automatically learn growth models which is nonparametric & data-driven

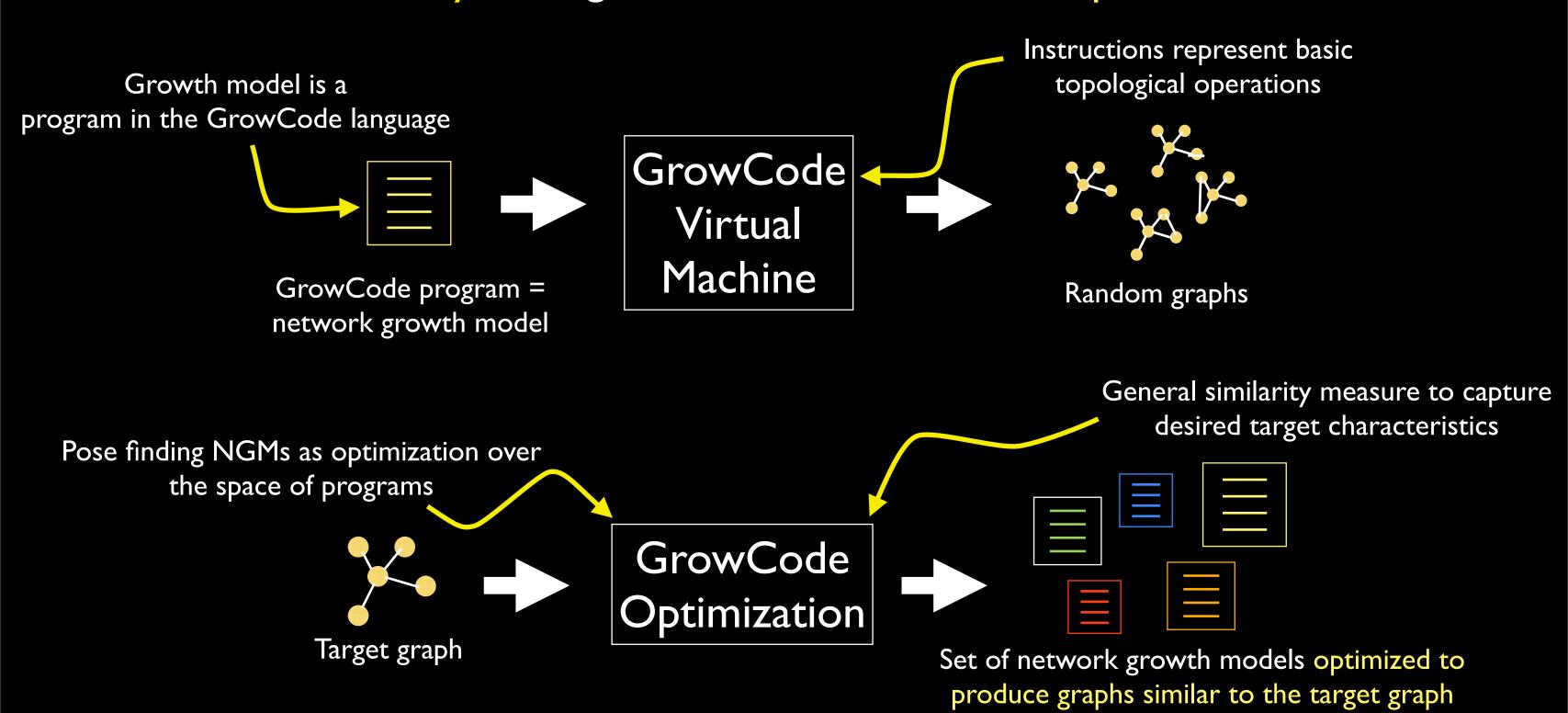




Set of network growth models optimized to produce graphs similar to the target graph

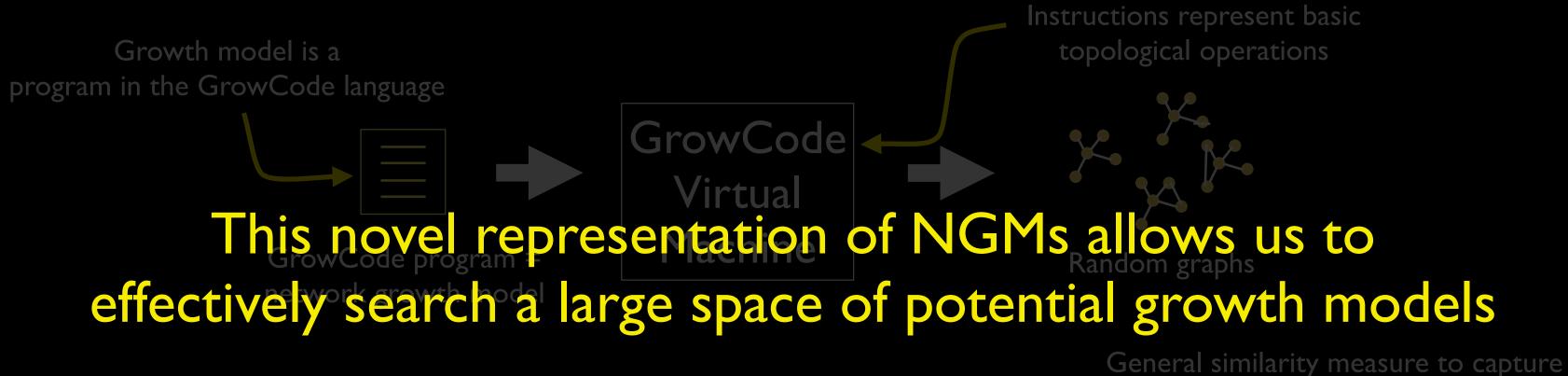
### So What's New?

Method to automatically learn growth models which is nonparametric & data-driven



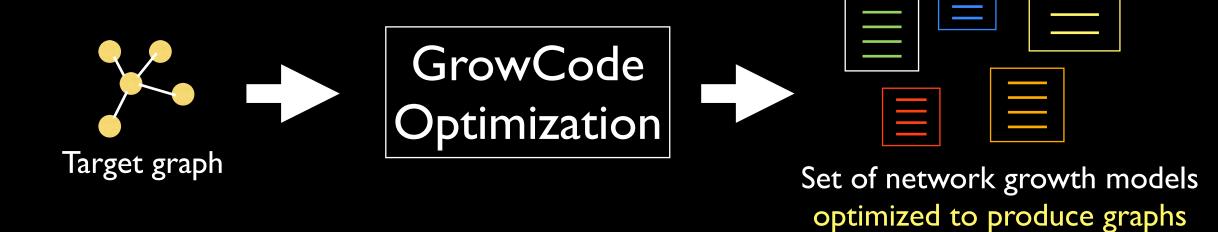
### So What's New?

Method to automatically learn growth models which is nonparametric & data-driven

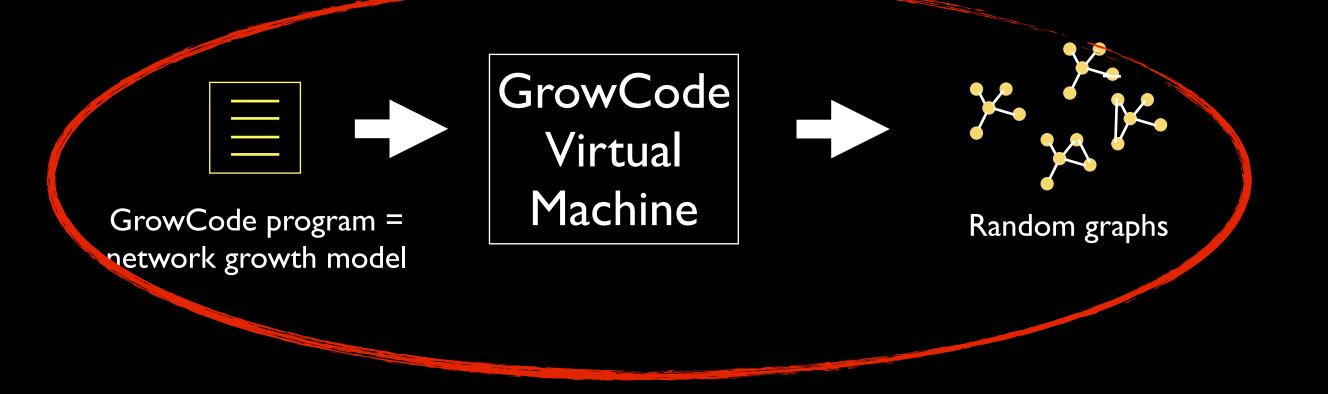


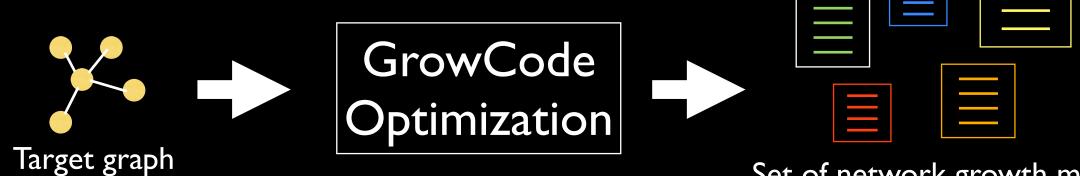






similar to the target graph





Set of network growth models optimized to produce graphs similar to the target graph

# GrowCode Virtual Machine

r0 r1 r2

Registers

Register-based virtual machine

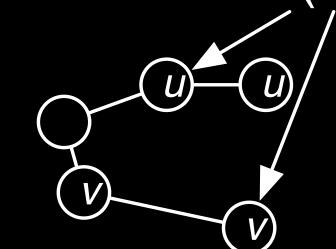
Node label memory  $L: V \rightarrow V$ 

Runs program iteratively to grow a graph

### Program

- 1: New node
- 2: RANDOM NODE
- 3: ATTACH TO INFLUENCED
- 4: CLEAR r2
- 5: **Set(1)**
- 6: RANDOM EDGE
- 7: DETACH FROM INFLUENCED
- 8: RANDOM NODE
- 9: Create edge
- Influence neighbors (0.692)

Node labels (act as memory)



Current graph

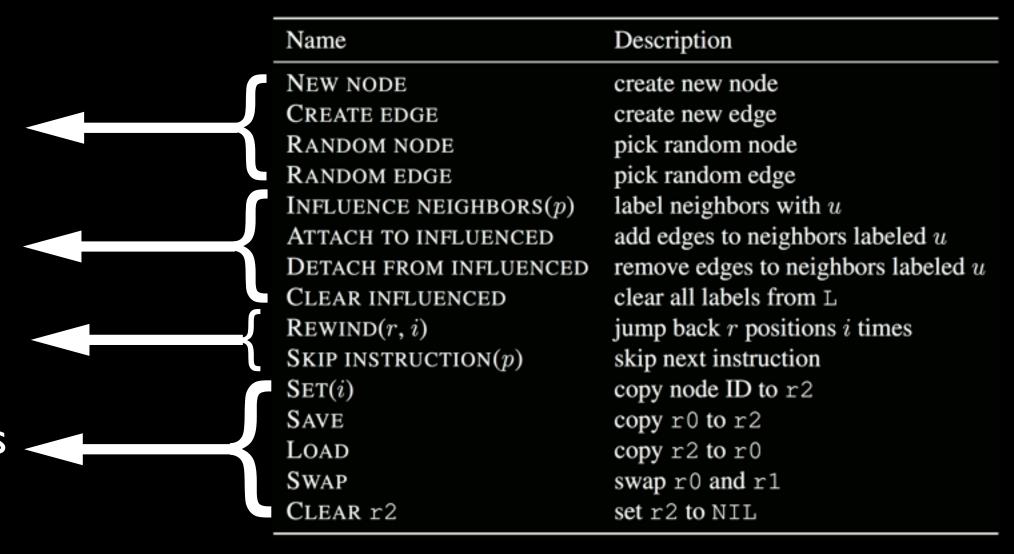
### Machine Instructions

Modify graph toplogy

Modify label memory

Program control flow

Manipulate machine registers



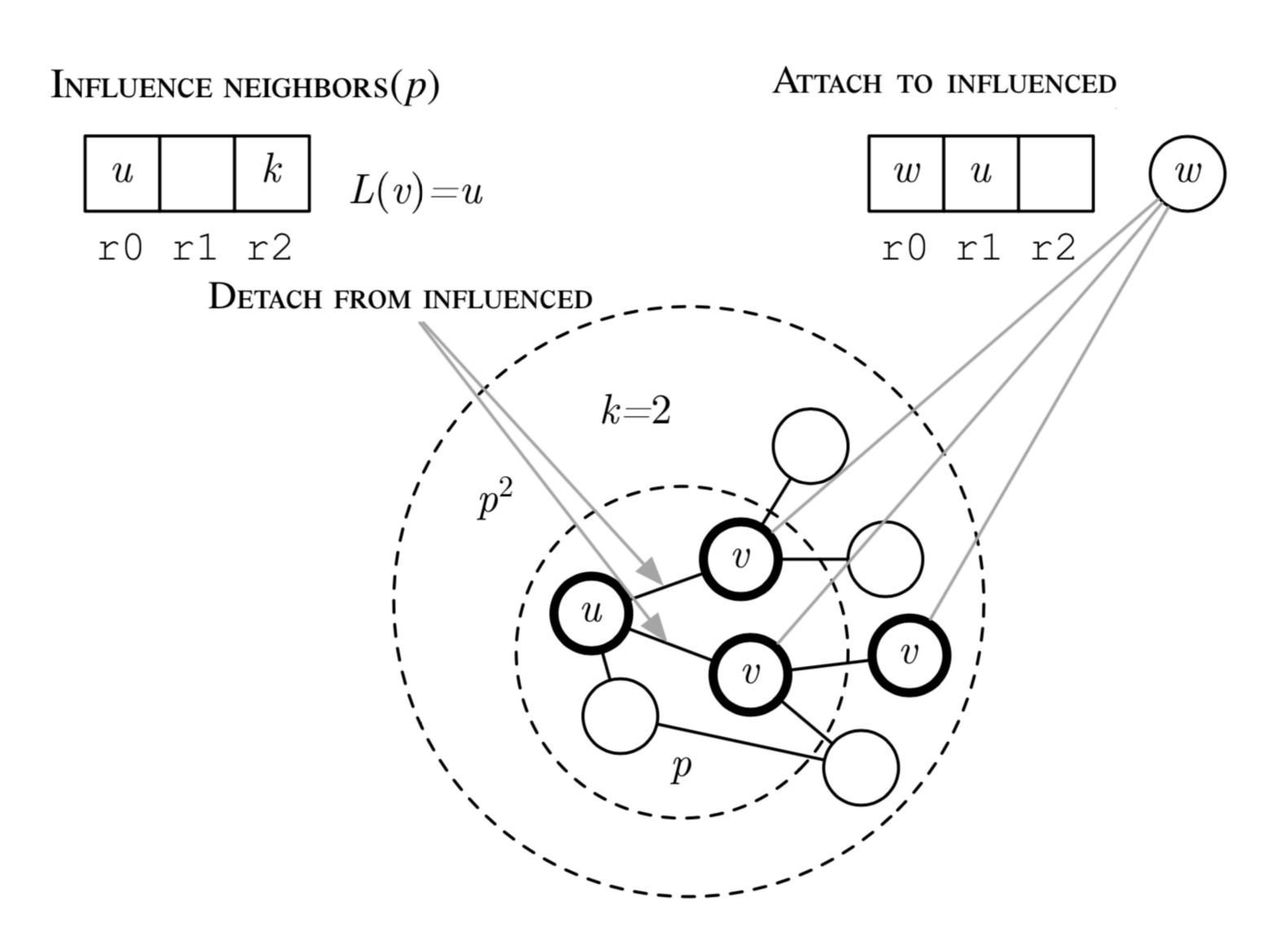
Mame

New Node

Create new mode

Create

# Influence Instructions



A new node duplicates an existing node u where u is selected proportional to its degree.

### Current graph:



### Program:

Set(1)
Random edge
New node
Swap
Influence neighbors(1.0)
Swap
Attach to influenced

Registers:

r0

r1 r2

A new node duplicates an existing node u where u is selected proportional to its degree.

### Current graph:



### Program:

#### Set (1)

Random edge
New node
Swap
Influence neighbors(1.0)
Swap
Attach to influenced



A new node duplicates an existing node u where u is selected proportional to its degree.

### Current graph:

Program:

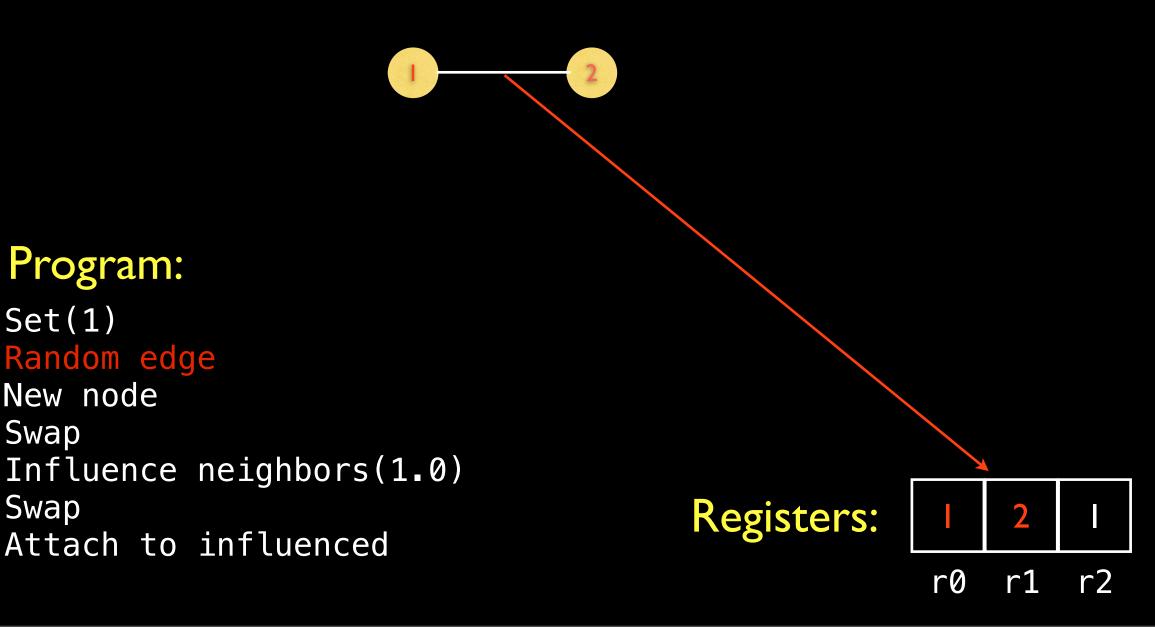
Random edge

Set(1)

Swap

Swap

New node



A new node duplicates an existing node u where u is selected proportional to its degree.

### Current graph:

Program:

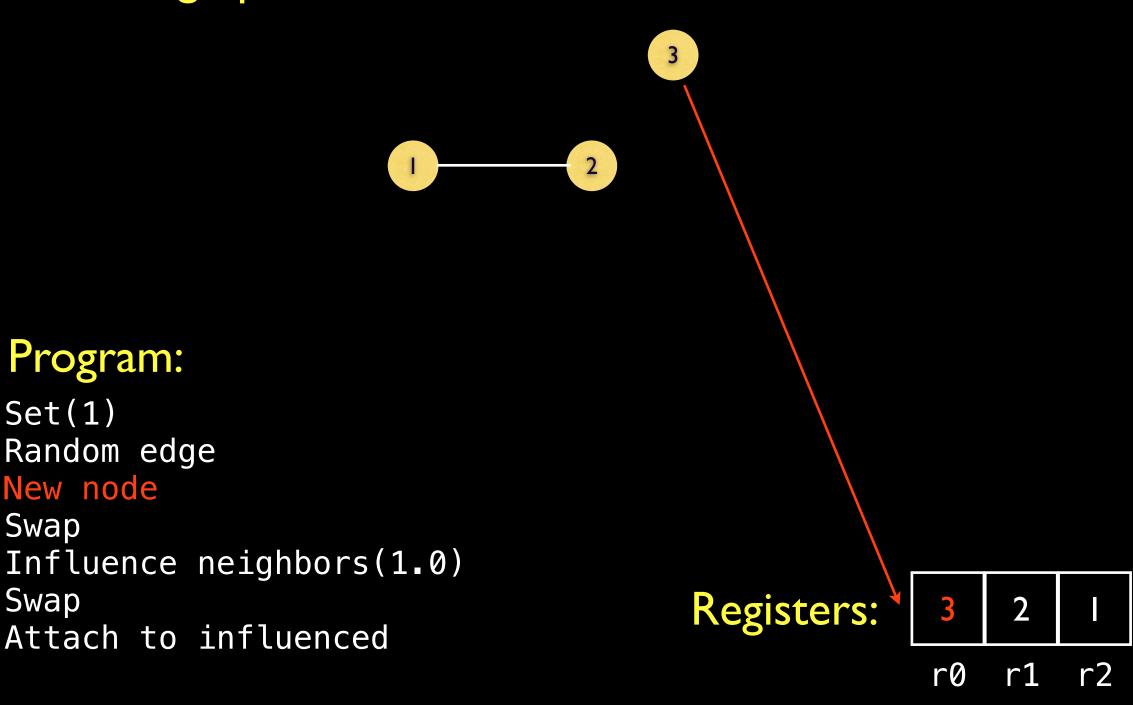
Random edge

Set(1)

Swap

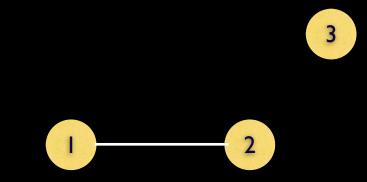
Swap

New node



A new node duplicates an existing node u where u is selected proportional to its degree.

### Current graph:



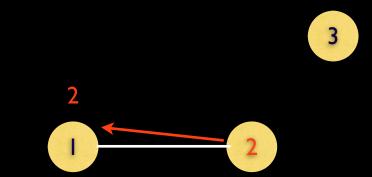
### Program:

Set(1)
Random edge
New node
Swap
Influence neighbors(1.0)
Swap
Attach to influenced

Registers: 2 3 I
r0 r1 r2

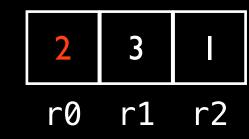
A new node duplicates an existing node u where u is selected proportional to its degree.

### Current graph:



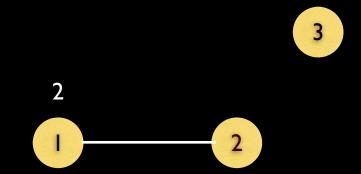
### Program:

Set(1)
Random edge
New node
Swap
Influence neighbors(1.0)
Swap
Attach to influenced



A new node duplicates an existing node u where u is selected proportional to its degree.

### Current graph:



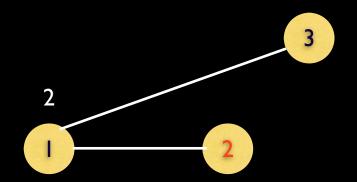
### Program:

Set(1)
Random edge
New node
Swap
Influence neighbors(1.0)
Swap
Attach to influenced



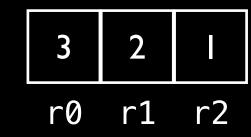
A new node duplicates an existing node u where u is selected proportional to its degree.

### Current graph:



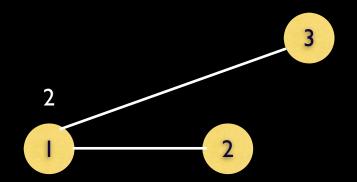
### Program:

Set(1)
Random edge
New node
Swap
Influence neighbors(1.0)
Swap
Attach to influenced



A new node duplicates an existing node u where u is selected proportional to its degree.

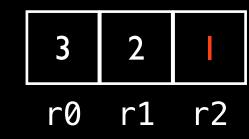
### Current graph:



### Program:

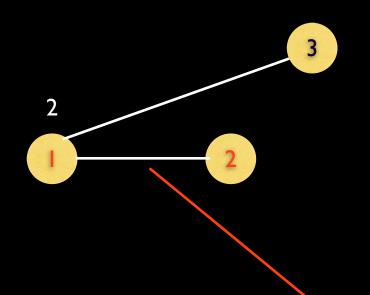
#### Set (1)

Random edge
New node
Swap
Influence neighbors(1.0)
Swap
Attach to influenced



A new node duplicates an existing node u where u is selected proportional to its degree.

### Current graph:



### Program:

Set(1)
Random edge
New node
Swap
Influence neighbors(1.0)
Swap
Attach to influenced

A new node duplicates an existing node u where u is selected proportional to its degree.

### Current graph:

Program:

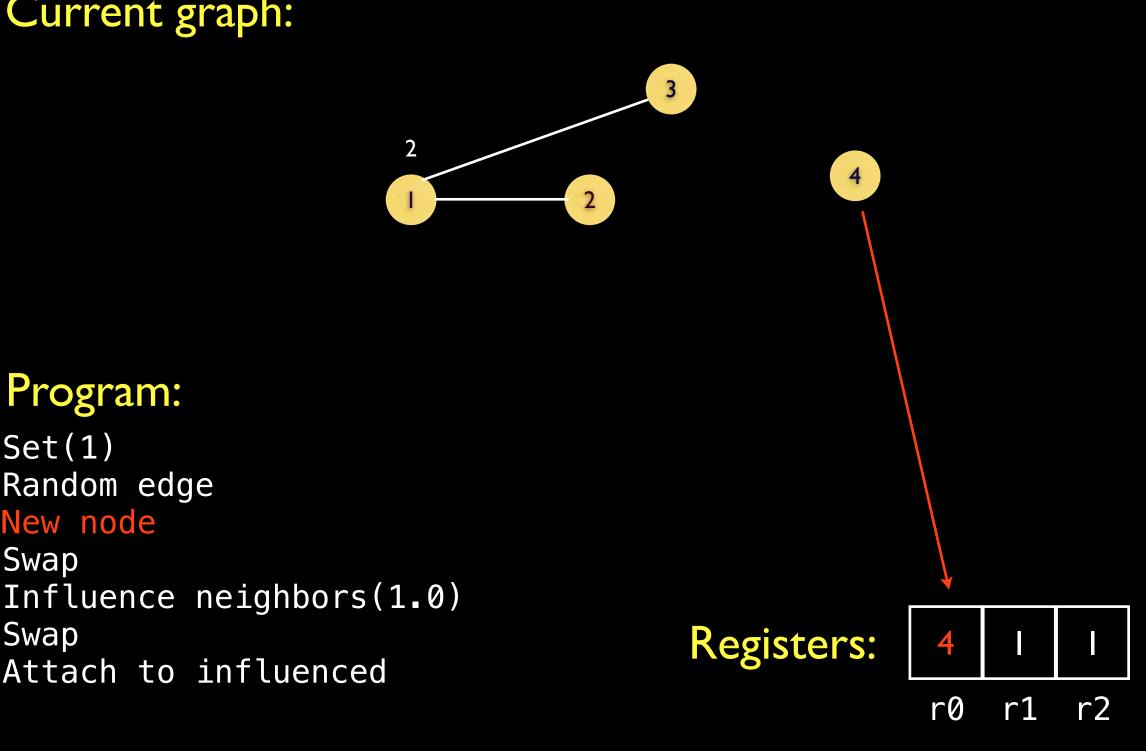
Random edge

Set(1)

Swap

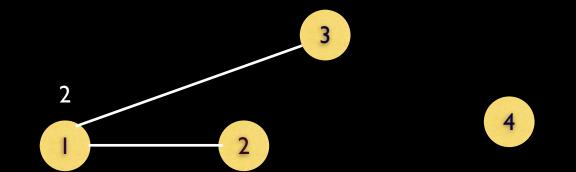
Swap

New node



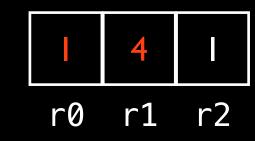
A new node duplicates an existing node u where u is selected proportional to its degree.

#### Current graph:



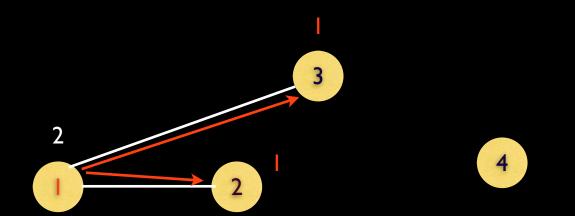
#### Program:

Set(1)
Random edge
New node
Swap
Influence neighbors(1.0)
Swap
Attach to influenced



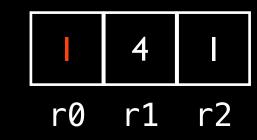
A new node duplicates an existing node u where u is selected proportional to its degree.

#### Current graph:



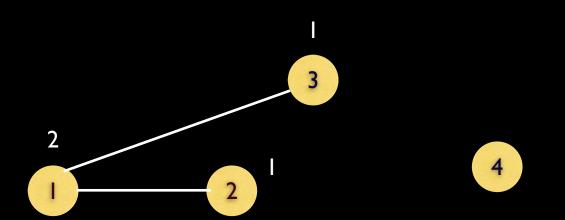
#### Program:

Set(1)
Random edge
New node
Swap
Influence neighbors(1.0)
Swap
Attach to influenced



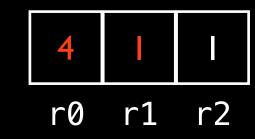
A new node duplicates an existing node u where u is selected proportional to its degree.

#### Current graph:



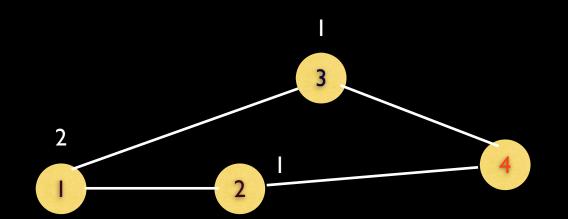
#### Program:

Set(1)
Random edge
New node
Swap
Influence neighbors(1.0)
Swap
Attach to influenced



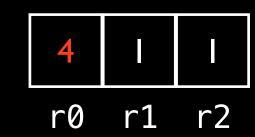
A new node duplicates an existing node u where u is selected proportional to its degree.

#### Current graph:



#### Program:

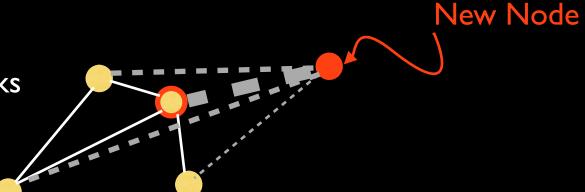
Set(1)
Random edge
New node
Swap
Influence neighbors(1.0)
Swap
Attach to influenced



# GrowCode Can Express Existing Models

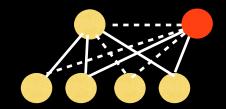
Barabási-Albert (models preferential attachment of new nodes)

Reproduces scale free distribution observed in large, real networks



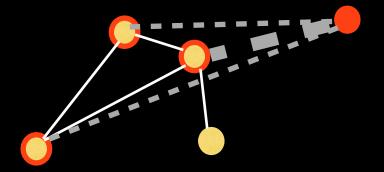
#### DMC (models protein duplication and functional divergence events)

Reproduces scale free distribution and clustering properties of protein-protein interaction networks (PPI)



#### Forest Fire (new nodes connect to a set of topologically close nodes)

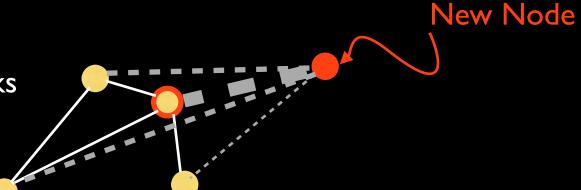
Reproduces scale free distribution, shrinking diameter, and densification over time



# GrowCode Can Express Existing Models

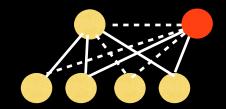
Barabási-Albert (models preferential attachment of new nodes)

Reproduces scale free distribution observed in large, real networks



#### DMC (models protein duplication and functional divergence events)

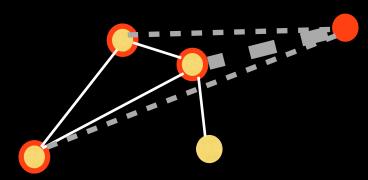
Reproduces scale free distribution and clustering properties of protein-protein interaction networks (PPI)



#### Forest Fire (new nodes connect to a set of topologically close nodes)

Reproduces scale free distribution, shrinking diameter, and densification over time

GrowCode instructions are reused throughout all three models.



# Existing Models in GrowCode

#### Algorithm 1 Barabási-Albert

1: New node	$\triangleright$ Create a new node, $u$
2: Save	
3: Random edge	$\triangleright$ Choose a random edge, $e$
4: Skip instruction $(0.5)$	$\triangleright$ Choose random endpoint $v$ of $e$
5: SWAP	
6: Load	
7: Create edge	$\triangleright$ Create an edge between $v$ and $u$
8: Rewind(5, $i$ )	$\triangleright$ Attach it to <i>i</i> random nodes

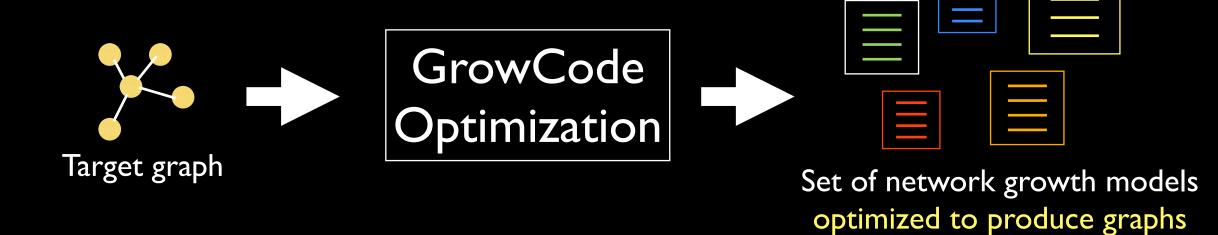
#### **Algorithm 3** FF

	,	
1:	RANDOM NODE	⊳ Put a random node in <b>r0</b>
2:	CLEAR r2	> Clear r2 to allow full graph influence
3:	Influence neighbors $(b)$	
4:	SWAP	⊳ Move the random node into r1
5:	New node	$\triangleright$ Create a new node, $u$
6:	Create edge	
7:	Attach to influenced	$\triangleright$ Connect <i>u</i> to influenced nodes

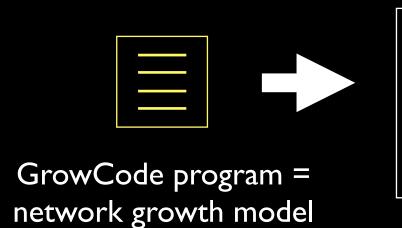
#### Algorithm 2 DMC

```
1: RANDOM NODE
                                              \triangleright Put a random node v in r0
 2: Set(1)
                                       \triangleright Set r2 (k-hop for influence) to 1
                                                  \triangleright Influence v's neighbors
 3: Influence neighbors (1.0)
 4: SWAP
                                                           \triangleright Swap r0 and r1
                                  \triangleright Create a new node u and put it in r0
 5: New Node
                                         \triangleright Connect u to influenced nodes
 6: ATTACH TO INFLUENCED
 7: CLEAR INFLUENCED
                                                  \triangleright Influence u's neighbors
 8: Influence neighbors (q_{\text{mod}}/2)
 9: SWAP
                                                  \triangleright Influence v's neighbors
10: Influence neighbors (q_{\text{mod}}/2)
                                           \triangleright Actually delete edges from v
11: Detach from influenced
12: SWAP
                                                         \triangleright Do the same for u
13: Detach from influenced
14: Clear influenced
                                                        \triangleright Skip adding \{u, v\}
15: Skip instruction (1.0 - q_{con})
16: Create edge
                                                      \triangleright with probability q_{\text{con}}
```

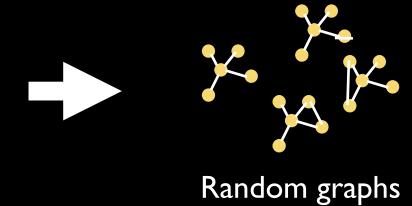


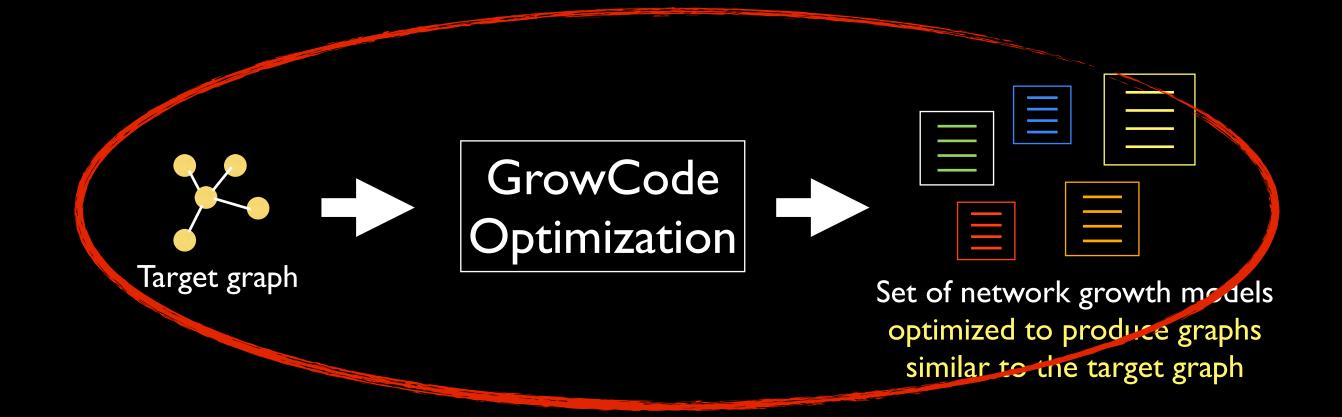


similar to the target graph

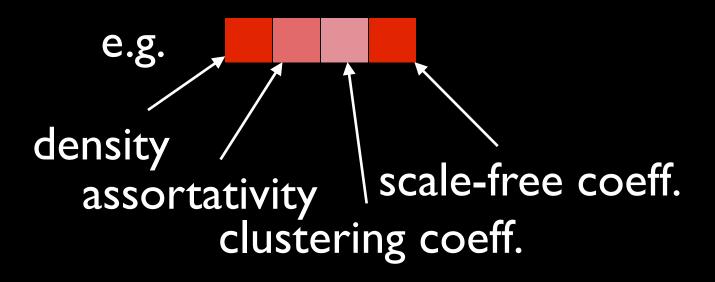


GrowCode Virtual Machine

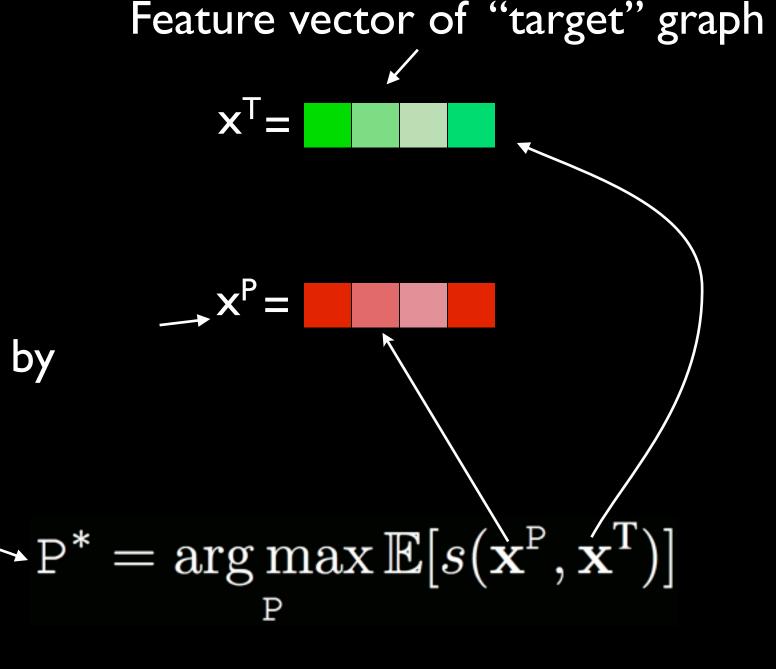




Characterize a graph with a user-defined feature vector



Feature vector for graph generated by GrowCode Prog.

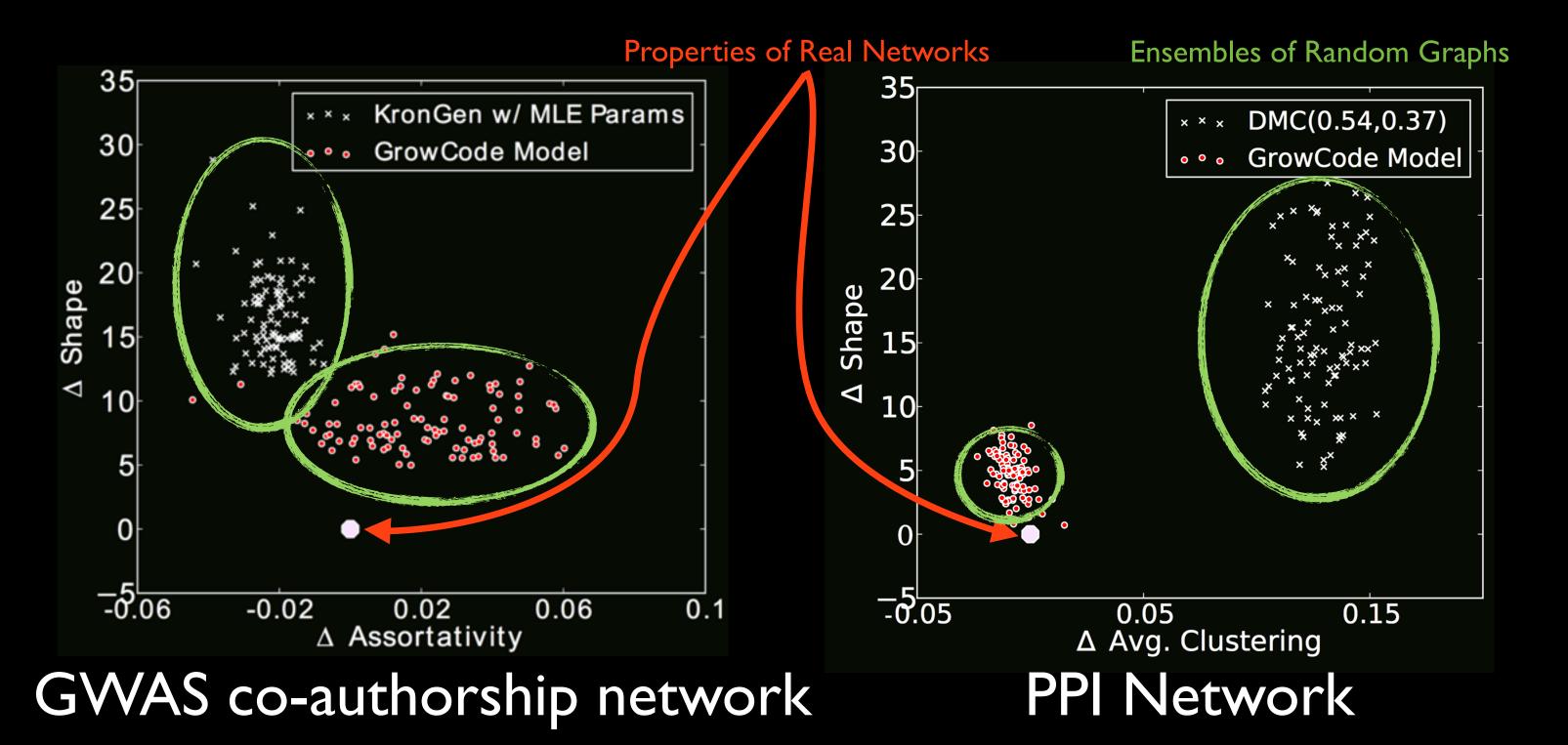


Finding a NGM becomes a search problem

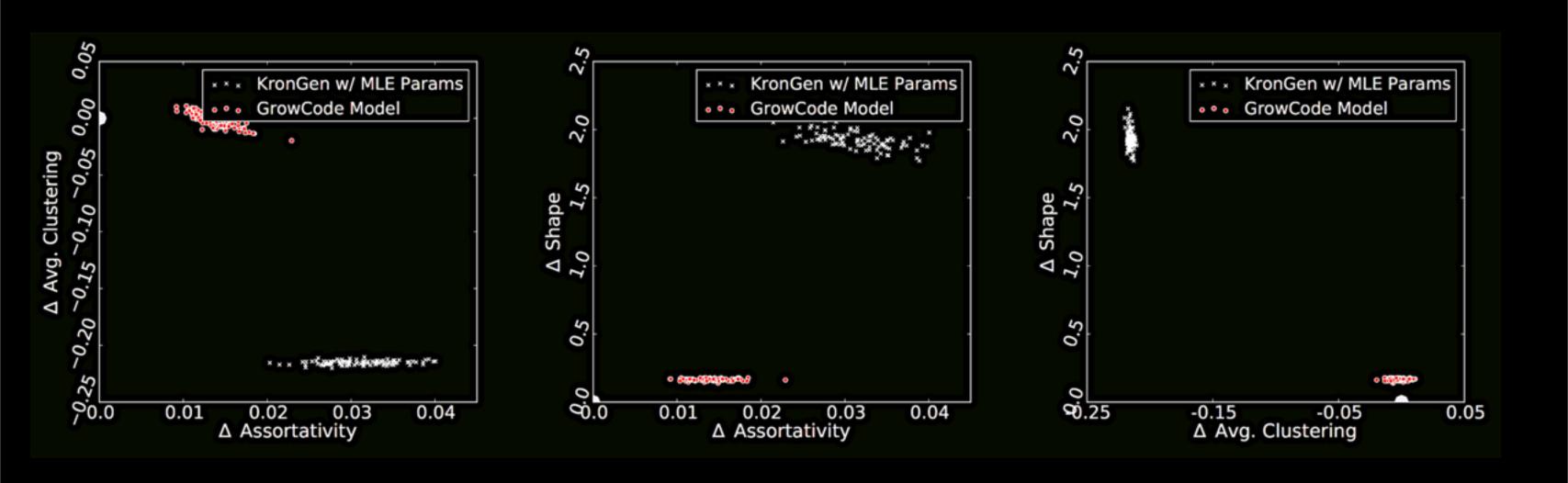
Want models (programs) that produce graphs whose feature vectors have high expected similarity with the target feature vector

# Learning Models for Real-World Networks

## GrowCode Better Matches Pairs of Features

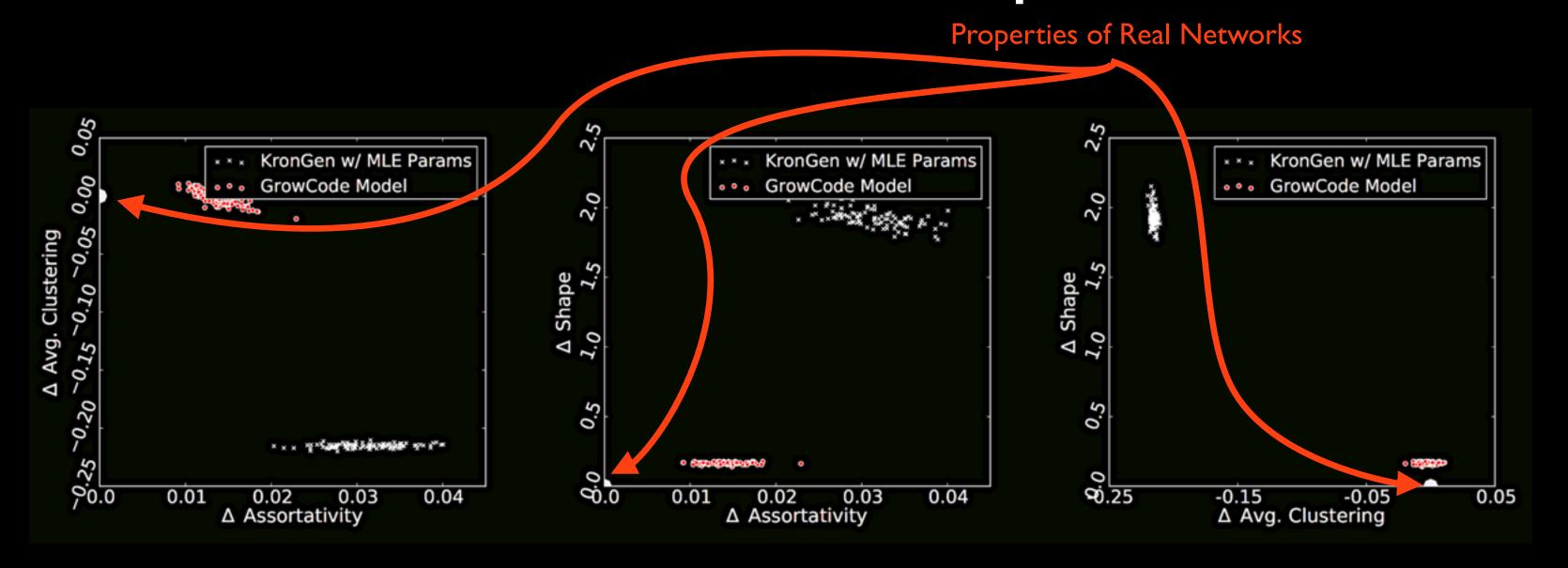


# GrowCode Better Matches Triplets of Features



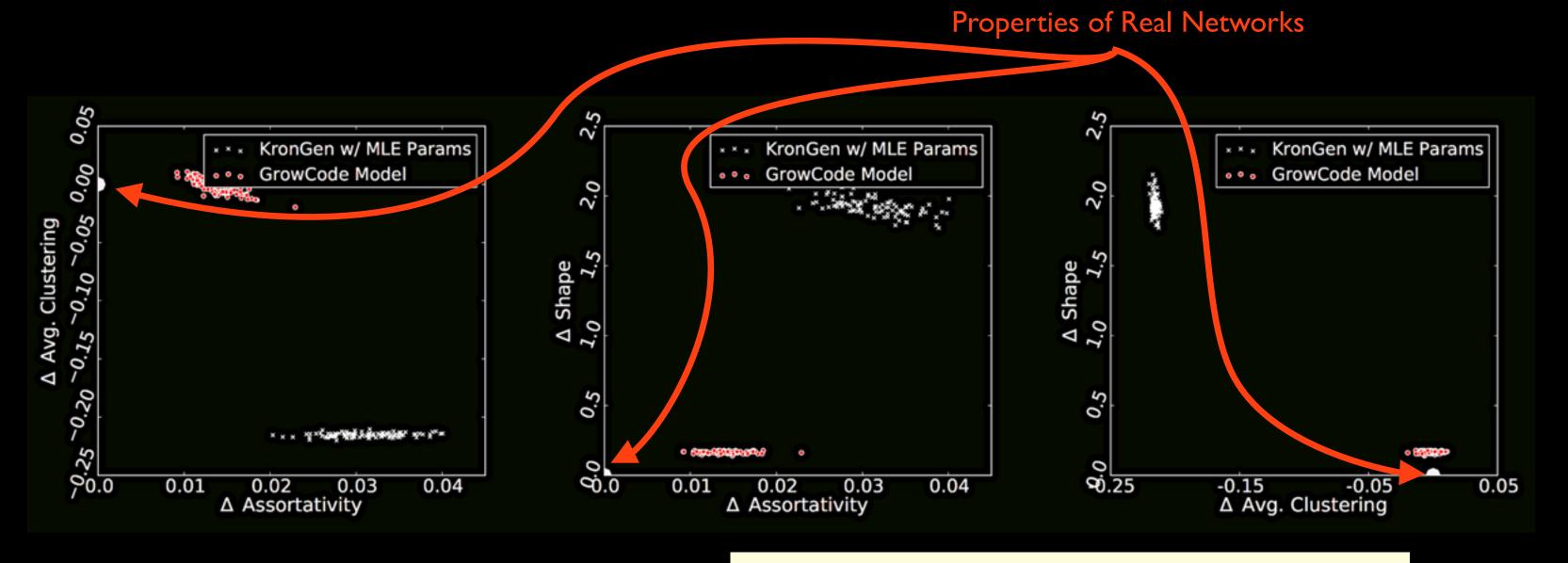
Internet Router Network

# GrowCode Better Matches Triplets of Features



Internet Router Network

# GrowCode Better Matches Triplets of Features



Social, technological, and biological!

Internet Router Network

# Random Walk with Restarts

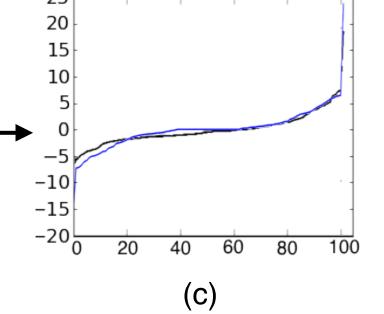
```
; RWR - Random Walk with Restart
; r3 is the start node; r2 is the new node; r1 is the current node
         ; r1 := random start node
SAVE
           ; r3 := r1
COPY; copy r1 from an existing node
ZERO; remove all of r1s copied edge
SWAP; r2 := r1
STEP 100
               ; do the random walk with restarts
 INC
               ; increment attribute for node r1
 SIFA != 5
               ; skip next statment if attribute for r1 != 50
 CREATE
               ; create an edge between r1 and r2
 NEIGH
               ; move to a random neighbor
  SCOIN 0.7
               ; skip next statment if coin < r
 LOAD
               r1 = r3 == start node
                    (a)
```

growcode model of network

evolution

ode 65 91 85 46 55 81 95 66 66 63 33 38 94 39 41 82 68 87 60 62 50 5 4 13 7 95 64 44 45 88 93 59 40 61 1 21 2 16 20 98 74 76 83 101 67 80 69 37 37 38 73 78 44 99 78

Spectra of adjacency matrix of (1) RWR model and (2) learned/optimized model



#### Conclusion

Novel representation of network growth models

Can express many existing models & define new ones

Automated process for learning models that grow graphs desired topological properties

Learned models can outperform (match better) than manually designed models with optimal parameters

Basic building blocks are interpretable -- entire model may not be

Direction for future work -- analyze programs to find interpretable motifs

# Questions

- **Accuracy:** Can something be proved about the deviation between the distribution of graphs produced by a GrowCode\* program and a hand-designed mechanism?
- **Completeness:** Can a set of properties and a GrowCode\* language be defined such that one could show that any graph growth model with these properties can be expressed with GrowCode\*?
- Efficiency: Is there a better optimization procedure than a GA?
- Regularization: How can one enforce "simple" models are determined?
- Interpretation: Can we mine learned programs to identify mechanism "motifs"?
- **Complexity**: Can we define a way to quantify how "easy" a property is to generate?

### Thanks



[CCF-1053918, EF-0849899, and IIS-0812111]



[IR2IAI085376]



Institute for Advanced Studies
New Frontiers Award



Research Fellowship to CK