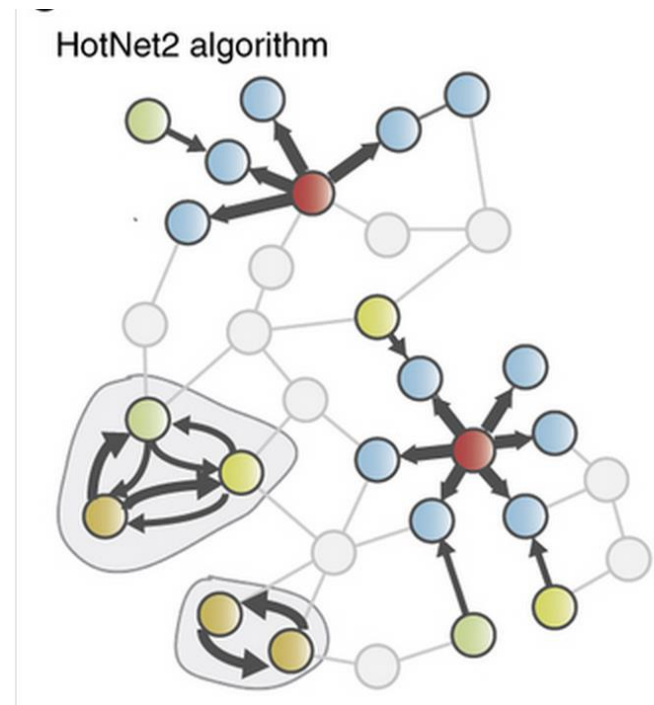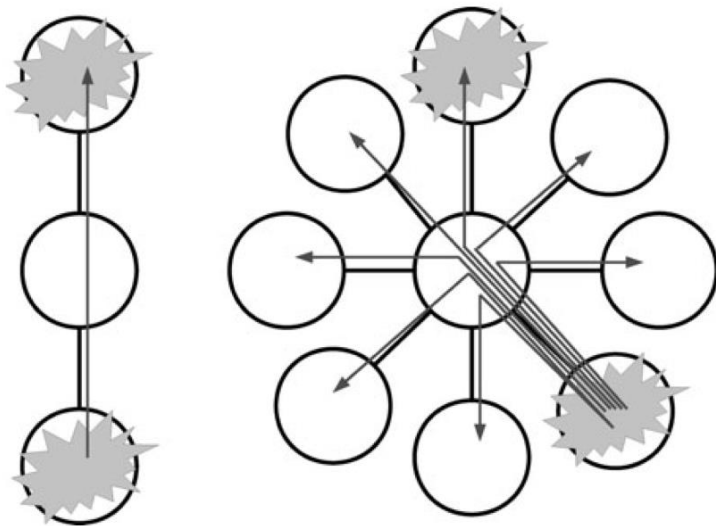# HotNet

# Background

- Determine significantly mutated subnetworks in a large gene interaction network

- Problems with current methods
  - Frequency doesn't always predict significance
  - Naïve subnetwork analysis
    - Enumeration prohibits subnetworks of reasonable size
    - Large number of hypotheses makes statistically significant difficult
    - Hub genes make for small gene diameters

# HotNet Overview

1. Formulate an influence measure between pairs of genes in the network

2. Identify subnetworks with *Combinatorial Model* or *Enhanced Influence Model*

3. Two-stage multiple hypothesis test to mitigate testing of large number of hypotheses

# Influence Graph

- Identify subnetworks that are significant with respect to a set of mutated genes



HotNet2 algorithm

# Diffusion

$f^{cs}_v(t)$    Amount of fluid @ node *V* at time *T*

$\mathbf{f}^s(t) = [f^s_1(t), \ldots, f^s_n(t)]^T$   Amount of fluid at all nodes

$L_\gamma = L + \gamma I.$   L is the laplacian matrix of the graph

$\dfrac{d\mathbf{f}^s(t)}{dt} = -L_\gamma \mathbf{f}^s(t) + \mathbf{b}^s u(t),$   Dynamics of the continuous process

- Interpret $f_i$ as the influence of gene $g_s$ on $g_i$

# Combinatorial Model

- Takes in influence measure between genes to discover significant subnetworks

**Combinatorial Algorithm**

**Input:** Influence graph $G_I$ and parameters $\delta$ and $k$

**Output:** Connected subgraph $\mathcal{C}$ of $G_I(\delta)$ with $k$ vertices

1. Construct $G_I(\delta)$ by removing from $G_I$ all edges with weight $< \delta$;
2. $\mathcal{C} \leftarrow \emptyset$;
3. **for** each node $v \in V$ **do**
4.     $C_v \leftarrow \{v\}$;
5.     **for** each $u \in V \setminus \{v\}$ **do** $p_v(u) \leftarrow$ shortest path from $v$ to $u$ in $G_I(\delta)$;
6.     **while** $|C_v| < k$ **do**

        // $\ell_v(u) = $ *set of nodes in* $p_v(u)$; $P_v(u) = $ *elements of* $I$ *covered by*
          $\ell_v(u)$; $P_{C_v} = $ *elements covered by* $C_v$; $P_{\mathcal{C}} = $ *elements covered by* $\mathcal{C}$

7.         $u \leftarrow \arg\max_{u \in V \setminus C_v : |\ell_v(u) \cup C_v| \leq k} \left\{ \dfrac{|P_v(u) \setminus P_{C_v}|}{|\ell_v(u) \setminus C_v|} \right\}$;
8.         $C_v \leftarrow \ell_v(u) \cup C_v$;
9.     **if** $|P_{C_v}| > |P_{\mathcal{C}}|$ **then** $\mathcal{C} \leftarrow C_v$;
10. **return** $\mathcal{C}$;

# Enhanced Influence Model

- Enhance the influence measure between genes by the number of mutations observed in each gene

**Enhanced Influence Algorithm**

**Input:** Influence graph $G_I$ and parameter $\delta$
**Output:** Connected components of $H(\delta)$

1  $V_H \leftarrow \{g_j : \mathcal{S}_j \neq \emptyset\}$;
2  $E \leftarrow \{g_j, g_k : g_j, g_k \in V_H, g_j \neq g_k\}$;
3  $H \leftarrow (V_H, E, h)$;
4  $E(\delta) \leftarrow \{(g_j, g_k) \in E : h(g_j, g_k) \geq \delta\}$;
5  $H(\delta) \leftarrow (V_H, E(\delta))$;
6  **return** connected components of $H(\delta)$;

# Statistics

- Calibrates with $H_o^{sample}$ and $H_0^{gene}$
  - Sample: mutations placed at random nodes
  - Gene: Move genes around…..?
- Compute significance of **number** of subnetworks
- Bound FDR

# Experimental Data

TABLE 1.  RESULTS OF THE COMBINATORIAL MODEL

| Dataset | k | Samples | p-value | | Pathway enrichment p-value | | |
|---------|---|---------|---------|---------|---------|---------|---------|
| | | | $H_0^{\text{sample}}$ | $H_0^{\text{gene}}$ | All | RTK/RAS/PI(3)K | p53 |
| GBM | 10 | 67 | $<10^{-10}$ | $4 \times 10^{-3}$ | $3 \times 10^{-4}$ | $8 \times 10^{-4}$ | 0.19 |
| | 20 | 78 | $<10^{-10}$ | $<10^{-3}$ | $10^{-5}$ | $8 \times 10^{-5}$ | 0.05 |
| Lung | 10 | 140 | $<10^{-10}$ | 0.02 | $8 \times 10^{-6}$ | / | |
| | 20 | 151 | $<10^{-10}$ | 0.03 | $3 \times 10^{-3}$ | / | |

$k$ is the number of genes in the subnetwork. *Samples* is the number of samples in which the subnetwork is mutated. *p-value* is the probability of observing a connected subgraph of size $k$ mutated in a number of samples $\geq$ *samples* under the random model $H_0^{\text{sample}}$ or $H_0^{\text{gene}}$. *enrichment p-value* is the $p$-value of the hypergeometric test for overlap between genes in the identified subgraph and genes reported significant pathways in TCGA (2008) or Ding et al. (2008). For GBM, *enrichment p-value* is the $p$-value of the hypergeometric test for RTK/RAS/PI(3)K and p53 pathways.

**a**

**b**

**c**