

Keyword-based Correlated Network Computation over Large Social Media

岡田莉奈(筑波大)

背景と研究目的

- 背景

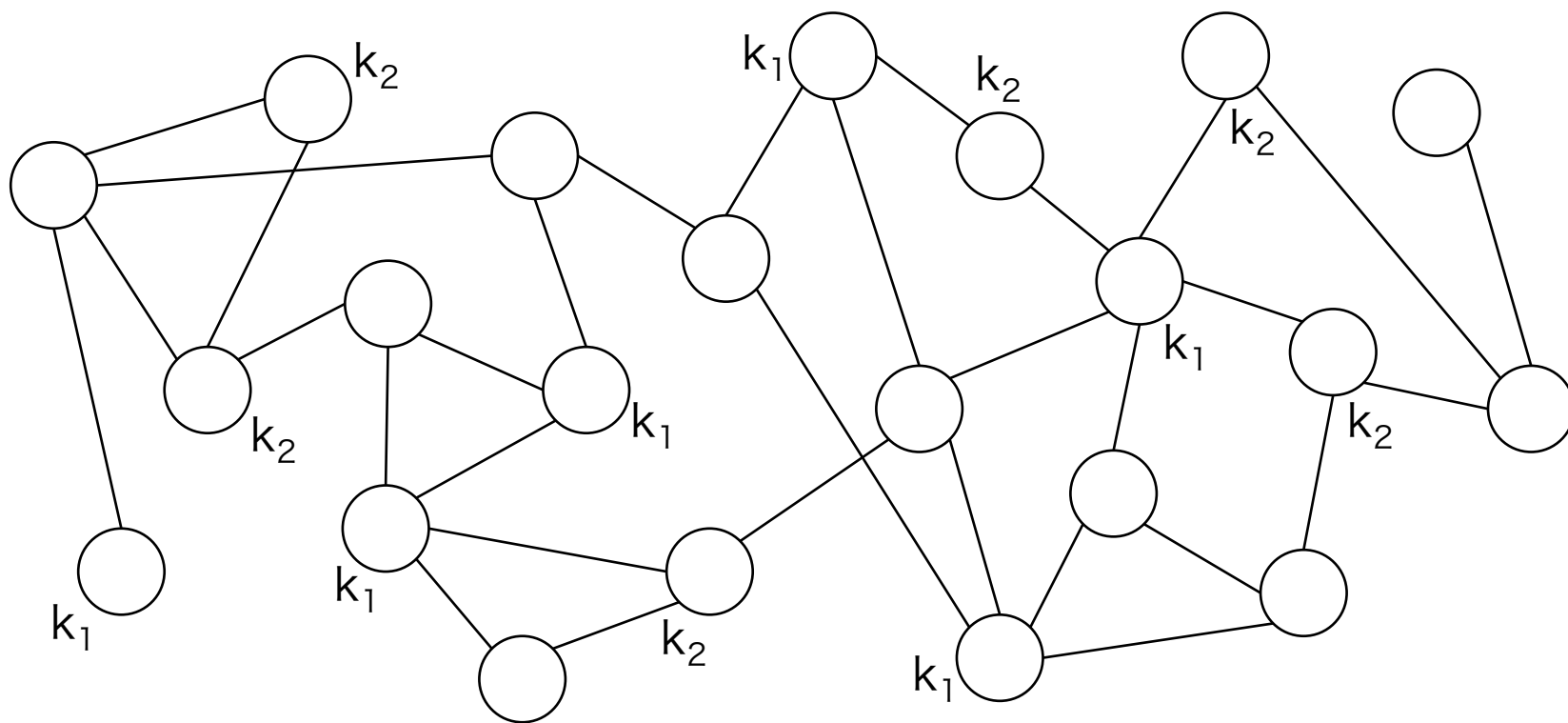
- ソーシャルメディアデータが増えており、これらを大きなグラフとみなすことができる。
 - ノード：実体， エッジ：実体間の関係
- 研究の動機
 - 何かイベントが起こったときに，キーワード検索をかける。欲しい情報は，そのキーワードを持つノードの情報だけでなく，そのノードの周辺情報も欲しい。

- 研究目的

- 大きなグラフからキーワードに基づいた相互関係のあるネットワークを見つけ出すこと。
 - サブグラフ(コンポーネント+ネイバー)の抽出。

相互関係のあるネットワークとは？

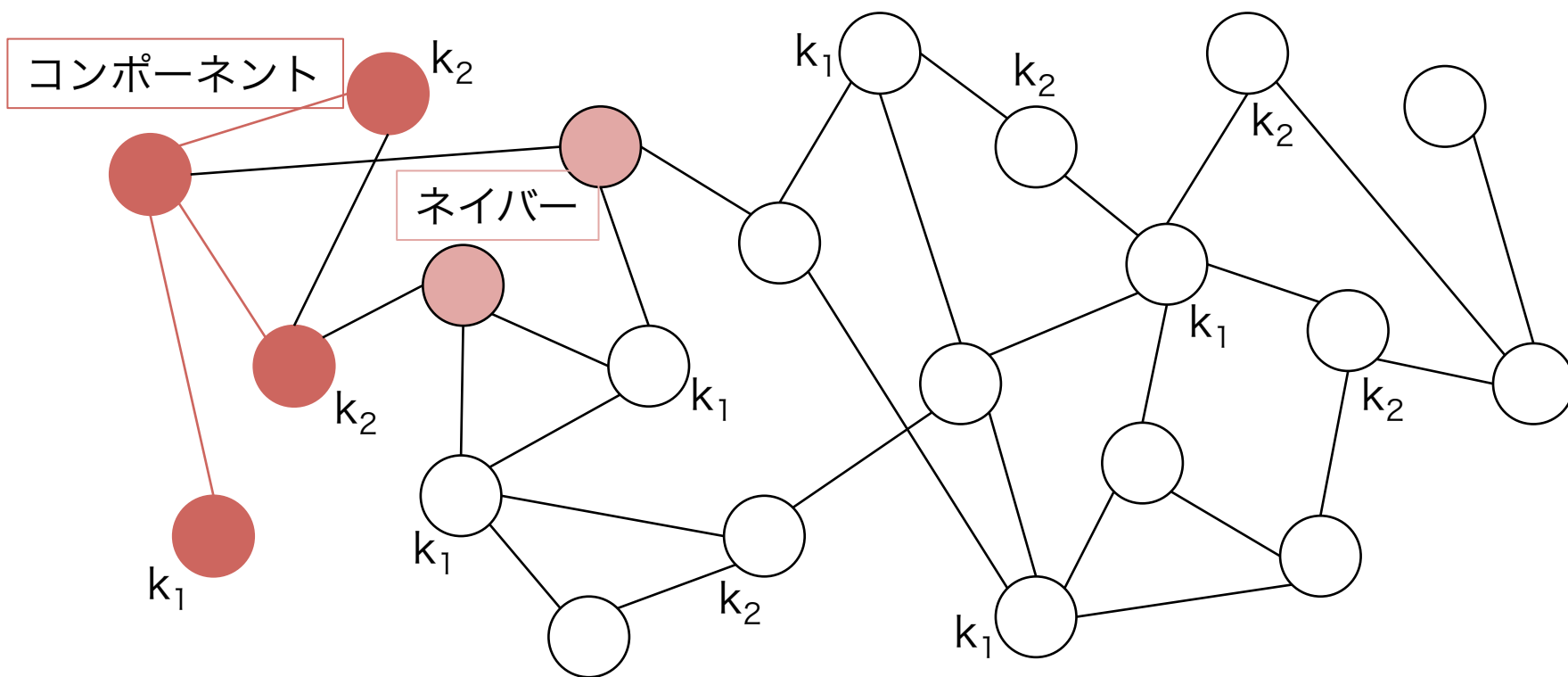
(例) 検索キーワード： $\{k_1, k_2\}$
コンポーネント内のキーワード間の最大距離 $r=2$



ソーシャルメディアデータ

相互関係のあるネットワークとは？

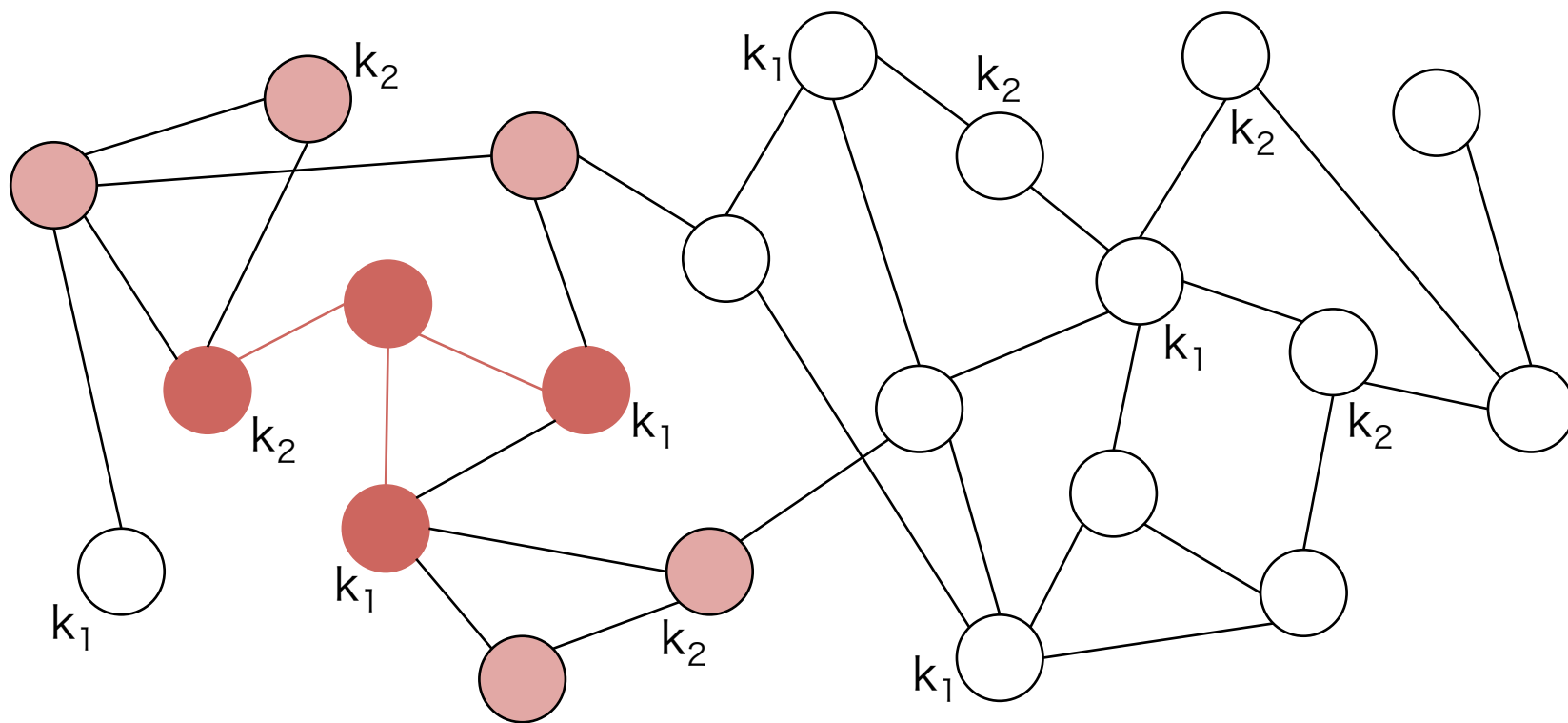
(例) 検索キーワード : $\{k_1, k_2\}$
コンポーネント内のキーワード間の最大距離 $r=2$



ソーシャルメディアデータ

相互関係のあるネットワークとは？

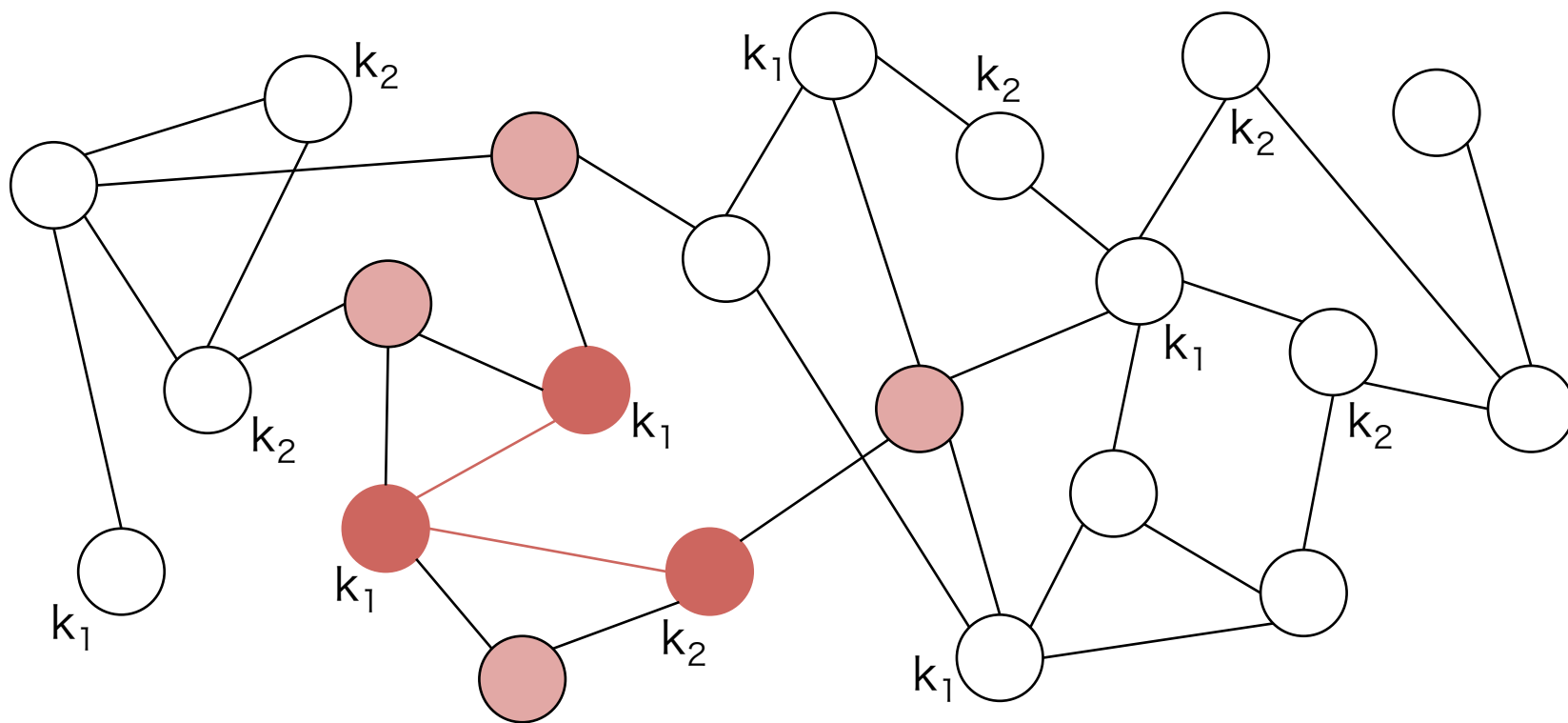
(例) 検索キーワード : $\{k_1, k_2\}$
コンポーネント内のキーワード間の最大距離 $r=2$



ソーシャルメディアデータ

相互関係のあるネットワークとは？

(例) 検索キーワード : $\{k_1, k_2\}$
コンポーネント内のキーワード間の最大距離 $r=2$



ソーシャルメディアデータ

相互関係のあるネットワークとは？

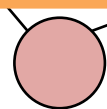
(例) 検索キーワード : $\{k_1, k_2\}$

$$\frac{\sum \{weight(v, G_1) * weight(v, G_2) | v \in G'_1 \cap G'_2\}}{|G'_1 \cup G'_2|}$$

ただし, $\begin{cases} weight(v', G_1) = 2^{-minDist(v', G_1)} & v' \text{がネイバーのとき} \\ 1 & \text{otherwise} \end{cases}$

という指標を用いて,
密に関連しているコンポーネントを探す.

G_1



k_2

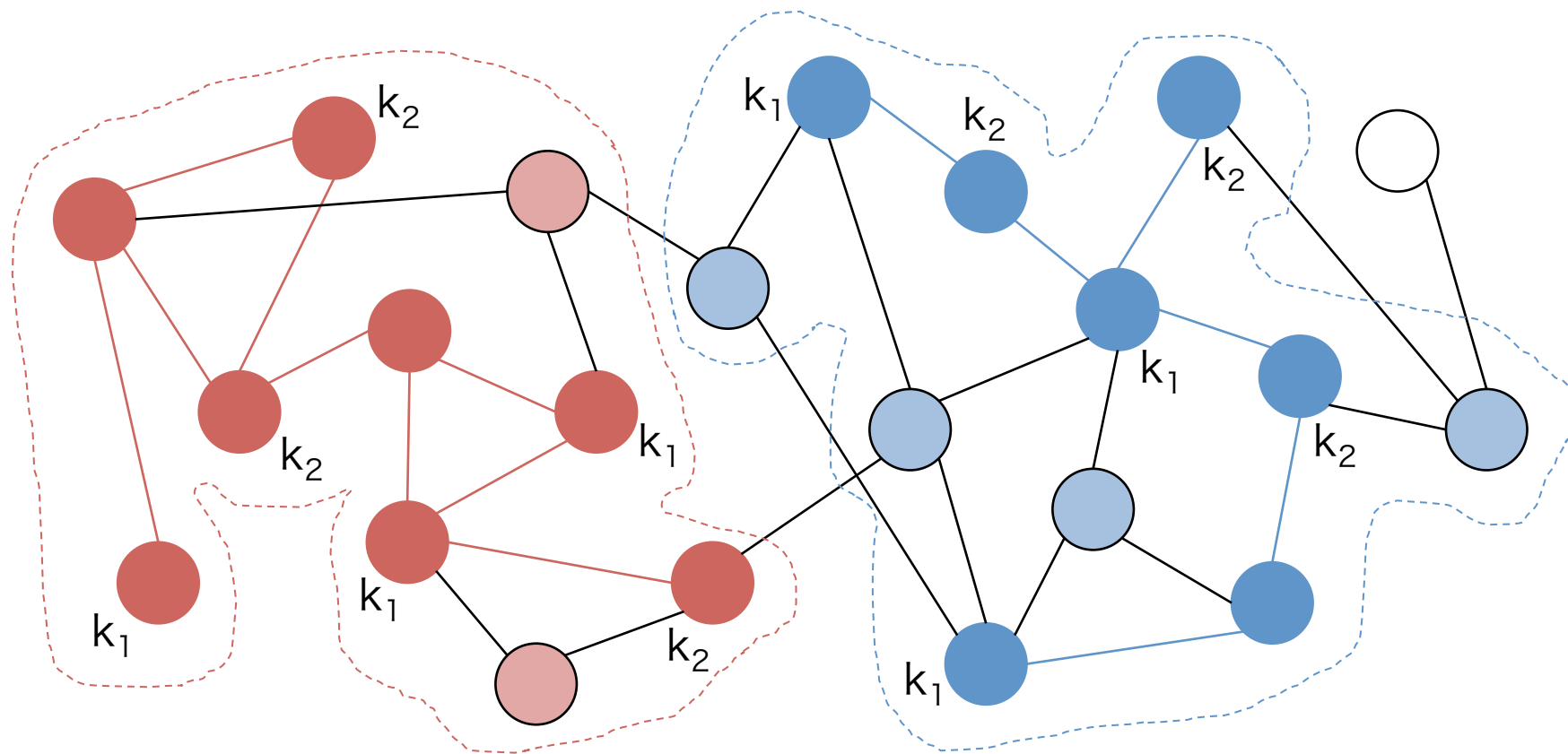


k_1

ソーシャルメディアデータ

相互関係のあるネットワークとは？

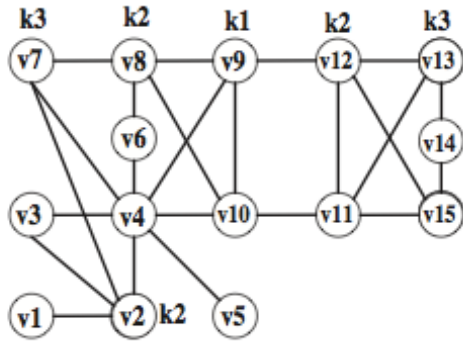
(例) 検索キーワード : $\{k_1, k_2\}$
コンポーネント内のキーワード間の最大距離 $r=2$



ソーシャルメディアデータ

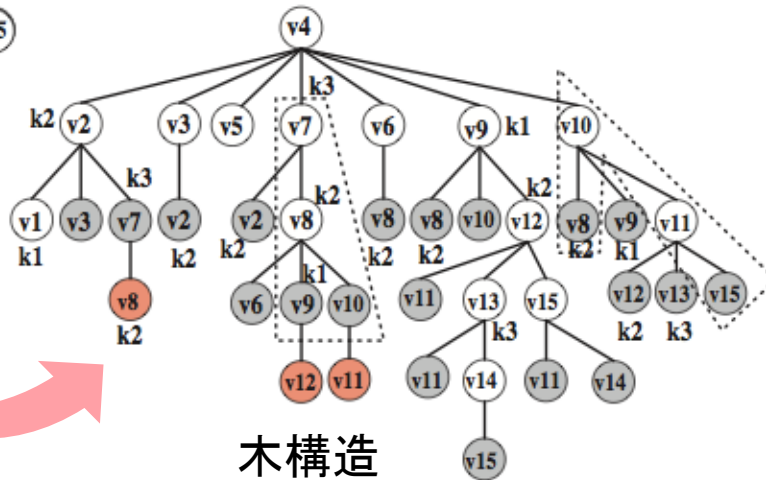
提案

- コンポーネントを効率的に見つけ出す方法を提案する。
 - 具体的には、新たな木構造を提案している。



グラフ

パラメータ γ によって、
木構造の深さをコントロール
ことができる。



木構造

- 利点
 - 距離を記憶した木構造をメインメモリ上に展開することができる。
 - 木構造を作らずにグラフから相互関係のあるコンポーネントを探すと、ペアワイズ比較 (N^2 の比較) を要する。しかし、この木構造を作れば比較はない。

実験

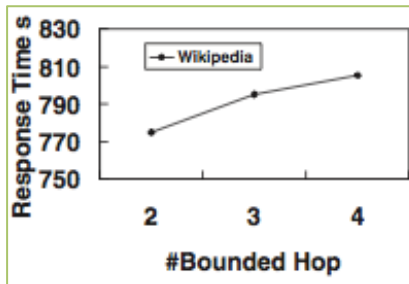
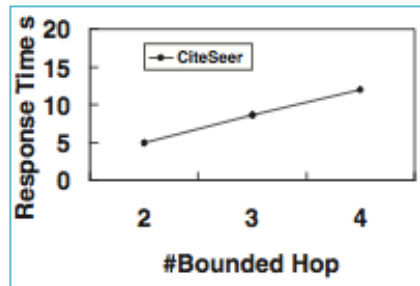
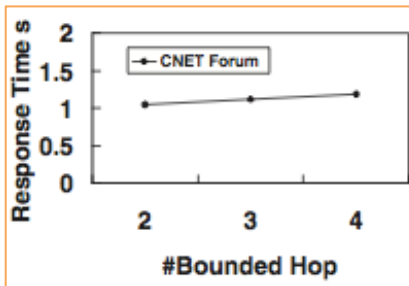
データセット

Dataset Name	Size	Number of Nodes	Number of Edges
CNET Forum	7MB	106,436	79,217
CiteSeer	89MB	230,470	406,784
Wikipedia	1,141MB	5,716,808	65,080,196

どちらも気にならない程度の
時間や使用メモリ領域
なので、これを使うのは
効果的である。

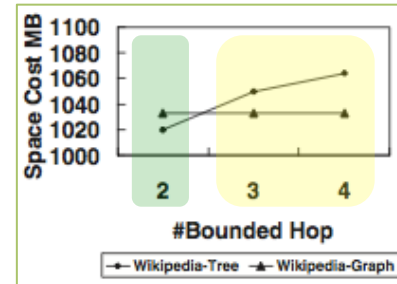
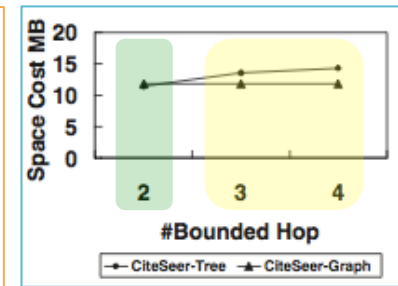
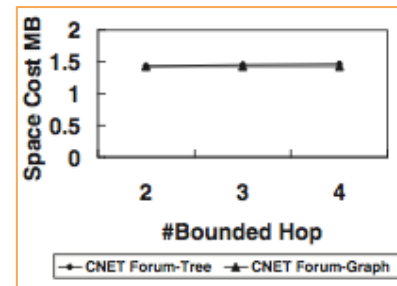
実験結果

① グラフから木への変換にかかる時間



平均次数が
10以上のときは
増加傾向にある。

② 木とグラフの使用メモリ領域



$r=2$ のときは木の方が
コストがかからない。
 $r \geq 3$ のときは木の方が
コストがかかるが、
気にならない程度。