

Treffen 07.11.2014

Dorle Osterode, Stefan Dang & Lukas Götz

18. November 2014

1 Implementationsmöglichkeiten der benötigten Datentypen in Genometools

Erste Ideen zur Implementierung der Datentypen fuer Scaffold-Graphen in C unter Verwendung bestehender Datentypen aus Genome-Tools.

```
1  /* Scaffold-Graph */
3  typedef enum { VC_UNIQUE, VC_REPEAT, VC_POLYMORPHIC, VC_UNKNOWN }
      VertexClass;
5
6  /* Vertex */
7
8  struct GtScaffoldGraphVertex
9  {
10     /* eindeutige ID fuer den Knoten */
11     GtUword id;
12     /* Laenge der Sequenz, die der Contig darstellt */
13     GtUword seqlen;
14     /* Wert der A-Statistik, um Contigs als REPEAT oder UNIQUE
15        klassifizieren zu koennen;
16        in Genom-Tools vom Typ float */
17     float astat;
18     /* abgeschaetzte Anzahl an Vorkommen des Contigs im Genom */
19     float copynum;
20     /* zur Klassifikation des Knotens: REPEAT, UNIQUE, ... */
21     VertexClass vertexclass;
22     bool hasconflictinglink;
23     GtUword nofedges;
24     /* Sammlung von Kanten, die von dem Contig abgehen */
25     struct GtScaffoldGraphEdge **edges;
26     /* Markierung fuer Algorithmen; aus Genome-Tools entnommen siehe
27        match/rdj-contigs-graph.c */
28     GtContigsGraphMarks color;
29 };
31 /* Edge */
```

```

33 struct GtScaffoldGraphEdge
34 {
35     /* Knoten, zu dem die Kante fuehrt */
36     struct GtScaffoldGraphVertex *pend;
37     /* Kante, die genau in die andere Richtung fuehrt */
38     struct GtScaffoldGraphEdge *ptwin;
39     /* Markierung fuer Algorithmen */
40     GtContigsGraphMarks color;
41     /* Abschaetzung der Entfernung der verbundenen Contigs */
42     GtWord dist;
43     /* Standardabweichung von der abgeschaetzten Entfernung */
44     float stddev;
45     /* Anzahl der Distanzinformationen, die ein Anzeichen fuer die
46     Verbindung der Contigs geben */
47     GtUword numpairs;
48     /* enthaelt die Richtung (Sense, Antisense) und welche
49     Straenge die paired-Information enthalten (die gleiche
50     Richtung oder das Reverse) */
51     bool antisense;
52     bool reverse;
53 };
54
55 /* Graph */
56 struct GtScaffoldGraph
57 {
58     struct GtScaffoldGraphVertex **vertices;
59     GtUword novertices;
60     struct GtScaffoldGraphEdge **edges;
61     GtUword noedges;
62 };

```

2 Verfeinerung des Algorithmus zur Filterung der Knoten

- graph.visit() Funktion durchläuft alle Knoten (HashMap) und ruft auf jedem Knoten die übergebene Visitor-Funktion auf. Vorher wird previsit() und danach postvisit() aufgerufen. (Quelle ScaffoldGraph.h)
- Es dürfte keine Probleme beim zusammenlegen der Filterfunktionen geben, da die Bedingungen nacheinander für jeden Knoten lokal geprüft werden können.
- Es sollte beachtet werden, dass es gegebenenfalls andere Ergebnisse geben könnte, wenn in einem späteren Schritt die vorher schon herausgefilterten Knoten noch beachtet werden.
- Die repetitiven Knoten sollten schon bei der Konstruktion des Graphen heraus-

gefiltert werden.

- Die markierten Knoten und Kanten können nicht sofort gelöscht werden, da die restlichen Knoten noch nicht klassifiziert sind. Da die Klassifikation eines Knoten aber anhand aller Nachbarn durchgeführt wird, können die Knoten und Kanten noch nicht früher gelöscht werden.
- Bei SGA wird gezählt, wie viele Knoten anhand welches Kriteriums gelöscht werden. Dies könnte auch noch eingebaut werden.

Algorithm 1: Zusammengefasste Filterfunktion (Schritt 4a und 4b vereinigt)

```
1 foreach Knoten  $k_0$  im Graph  $G$  do
2   foreach Kantenrichtung  $dir$  in  $[ANTISENSE, SENSE]$  do
3     foreach Kantenpaar  $(A, B)$  in Richtung  $dir$  do
4        $k_1 = A.pend$ ;
5        $k_2 = B.pend$ ;
6       if  $AmbiguousOrdering(A, B, p\_cutoff)$  and
7          $k_1.estCopy + k_2.estCopy < cn\_cutoff$  then
8         if  $k_1.estCopy < k_2.estCopy$  then
9           markiere  $k_1$  als polymorph und alle aus- und eingehenden SENSE
10          und ANTISENSE Kanten von  $k_1$  schwarz, so dass sie im nächsten
11          Schritt nicht mitbeachtet werden.;
12         end
13         else
14           markiere  $k_2$  als polymorph und alle aus- und eingehenden SENSE
15          und ANTISENSE Kanten von  $k_2$  schwarz, so dass sie im nächsten
16          Schritt nicht mitbeachtet werden.;
17         end
18         // bei polymorphen Knoten wird nur das erste polymorphe
19         Kantenpaar markiert
20         if Knoten  $k_0$  ist polymorph markiert then
21           break;
22         end
23       end
24     end
25   end
26   // polymorphe Knoten müssen nicht mehr auf inkonsistente Kanten
27   überprüft werden
28   if Knoten  $k_0$  ist polymorph then
29     break;
30   end
31   foreach Kantenpaar  $(A, B)$  in Richtung  $dir$  do
32     if  $A$  ist nicht schwarz und  $B$  ist nicht schwarz then
33       Berechne Overlap von  $A$  und  $B$  und speichere längsten Overlap.;
34     end
35   end
36   if  $längster\ Overlap > 400$  then
37     Markiere alle ausgehenden Sense-/Antisensekanten von  $k_0$  rot;
38   end
39 end
40 end
41 Lösche alle markierten Knoten und Kanten;
```

Algorithm 2: Funktion `AMBIGUOUSORDERING(A, B, p_cutoff)`

Data: Kante A und Kante B , die auf eindeutige Ordnung geprüft werden sollen.

Wahrscheinlichkeitsschwellenwert p_cutoff

Result: Ob die Kanten A und B nicht eindeutig geordnet werden können

```
1  $\mu = A.dist - B.dist$ ;  
2  $\sigma^2 = A.\sigma^2 + B.\sigma^2$ ;  
3  $t = \frac{-\mu}{\sigma \cdot \sqrt{2}}$ ;  
4  $P_{AB} = \frac{1}{2} \cdot \left(1 + \frac{2}{\sqrt{\pi}} \int_0^t \exp -x^2 dx\right)$ ;  
5  $P_{BA} = 1 - P_{AB}$ ;  
6 return  $\max\{P_{AB}, P_{BA}\} \leq p\_cutoff$ 
```

3 Notizen zu dem Layout-Algorithmus

- terminale Knoten: Knoten, die nur SENSE oder ANTISENES Kanten haben (Quelle: ScaffoldAlgorithms)
- Zusammenhangskomponenten-Funktion wird in StringGraph/GraphSearchTree.h definiert
- beim Walk für den Scaffold wird zuerst für jede Zusammenhangskomponente jeder Pfad zwischen terminalen Knoten mit einer Breitensuche (ohne heuristische Auswahl der Reihenfolge der Kindsknoten) berechnet. Dabei wird die Pfadlänge (Gap-Größe) minimiert. Als Layout für eine Zusammenhangskomponente wird dann der Pfad mit der längsten Sequenz (ohne Gaps) gewählt. (Quelle: ScaffoldGraph und ScaffoldWalk)

3.1 Layout-Algorithmus

Algorithm 3: Berechnung der Scaffolds (Schritt 6)

Data: Graph G

Result: Graph G ohne Knoten, die nicht zum bestem Walk gehören

```
1 Markiere alle Kanten aus  $G$  schwarz;  
2 Berechnung der Menge  $C$  aller Connected Components von  $G$ ;  
3 foreach Connected Component  $c_0$  aus der Menge  $C$  do  
4   Berechne Menge der terminalen Knoten  $T$  (mit ausschließlich SENSE oder  
   ANTISENSE Kanten) für die Connected Component  $c_0$ ;  
5   foreach Terminaler Knoten  $t_0$  aus der Menge  $T$  do  
6     Berechne die Menge  $W$  aller Walks durch die Connected Component  $c_0$  von  $t_0$   
     aus;  
7     foreach Walk  $w_0$  aus  $W$  do  
8       if Contig-Gesamtlänge  $>$  bislang beste Contig-Gesamtlänge then  
9         Setze aktuellen Walk  $w_0$  als bestWalk;  
10      end  
11    end  
12  end  
13 end  
14 Setze alle Kanten des bestWalk weiß;  
15 Lösche alle schwarzen Kanten;
```

Algorithm 4: Berechnung der Walks (Schritt 6.1)

Data: terminaler Startknoten t_0

Result: Alle von diesem Knoten möglichen Walks

```
1 Konstruktionsrichtung = Richtung der vom terminalen Knoten  $t_0$  ausgehenden
  Kanten;
2 foreach Kante  $A$  vom Knoten  $t_0$  ausgehend (in Konstruktionsrichtung) do
3    $k_0 = A.pend$ ;
4   Speichere Startkante  $A$  und Distanz  $A.dist$  in Map an Position  $k_0$  (für spätere
    Traversierung);
5   Schiebe Startkante  $A$  und Distanz  $A.dist$  in Queue;
6 end
7 while BFS über Queue nicht beendet do
8   Poppe Kante  $A$  und Distanz  $A.dist$  aus der Queue;
9    $k_0 = A.pend$ ;
10  foreach Kante  $B$  in Konstruktionsrichtung von  $k_0$  aus do
11     $k_1 = B.pend$ ;
12    if Distanz zu aktuell betrachtetem Knoten  $k_1 <$  bisher ermittelte Distanz zu
       $k_1$  OR Knoten  $k_1$  noch unbetrachtet then
13      Speichere Kante  $B$  und Distanz  $B.dist$  in Map an Position  $k_1$ ;
14      Schiebe Kante  $B$  und Distanz  $B.dist$  in Queue;
15    end
16  end
17  if  $k_0$  hat keine Kanten in Konstruktionsrichtung then
18    Schiebe Knoten  $k_0$  in terminalSet;
19  end
20 end
21 foreach Knoten  $k_0$  in terminalSet do
22   Erzeuge Walk mithilfe einer Traversierung über die Map
23 end
```
