

Алгоритм разрешения повторов по графу де Брюйна

Пржибельский Андрей, гр. 604 (SE)

Кафедра математических и информационных технологий

Научные руководители:

Банкевич А.В.

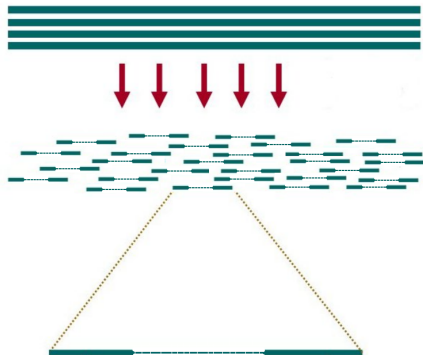
Нурк С.Ю.

Рецензент:

Алексеев М.А.

Геном и секвенирование

- Для биологов геном — это длинные молекулы ДНК
- Для нас — длинные строки над алфавитом из 4 букв (А,Т,Г,С)
- Секвенирование — процесс "чтения" генома
- Результат — короткие парные фрагменты (100 – 200 символов)



Граф де Брюина

- По полученным фрагментам строится граф де Брюина G
- Геном соответствует некоторому пути Γ в графе
- За счет парности исходных фрагментов между ребрами можно задать связи (парная информация)
- Вес парной информации $W(e_1, e_2, d)$ определяет степень уверенности в том, что ребра e_1 и e_2 находятся в пути генома Γ на расстоянии d

Сложности сборки

- Многие ребра повторяются в пути генома
- Граф де Брюйна является очень запутанным
- Из-за неравномерности исходных данных путь может Γ иметь разрывы

Постановка задачи

Дано:

- Граф де Брюйна \mathbf{G}
- Парная информация $W(e_1, e_2, d), e_1, e_2 \in E(\mathbf{G}), d \in \mathbb{Z}$

Найти:

- Множество путей P , таких что:
 - ▶ $|\{p \in P | p \notin \mathbf{\Gamma}\}| \rightarrow \min$
 - ▶ $LEN\{p \in P | p \in \mathbf{\Gamma}\} \rightarrow \max$

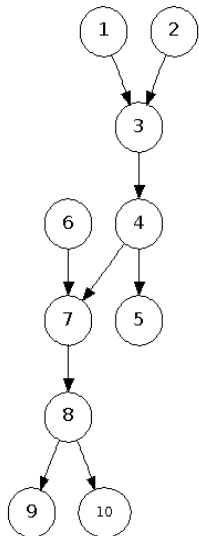
Предложенная стратегия

- Выберем в графе G достоверный путь p
- Пусть этот путь имеет несколько возможных продолжений h_1, \dots, h_n
- Добавим в p то, которое лучше остальных согласуется парной информацией с уже пройденным путем
- Продолжаем рост путей итеративно, пока это возможно

Выбор начальных путей

TRIVIAL_PATHS(G)

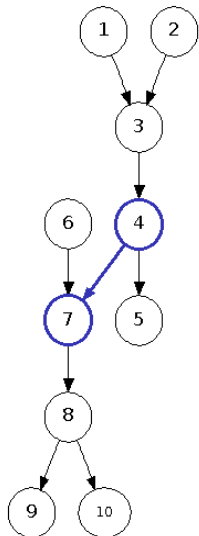
- 1 $P \leftarrow \{\}$
- 2 for each $e \in E(G)$:
- 3 $p \leftarrow (e)$
- 4 $(u, v) \leftarrow e$
- 5 while $\exists! w \in V(G) : (v, w) \in E(G)$:
- 6 $p \leftarrow (p, (v, w))$
- 7 $(u, v) \leftarrow (v, w)$
- 8 $(u, v) \leftarrow e$
- 9 while $\exists! w \in V(G) : (w, u) \in E(G)$:
- 10 $p \leftarrow ((w, u), p)$
- 11 $(u, v) \leftarrow (w, u)$
- 12 $P \leftarrow P \cup \{p\}$
- 13 return P



Выбор начальных путей

TRIVIAL_PATHS(G)

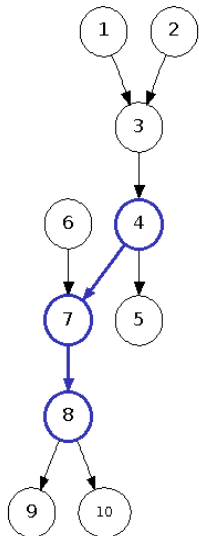
- 1 $P \leftarrow \{\}$
- 2 for each $e \in E(G)$:
- 3 $p \leftarrow (e)$
- 4 $(u, v) \leftarrow e$
- 5 while $\exists! w \in V(G) : (v, w) \in E(G)$:
- 6 $p \leftarrow (p, (v, w))$
- 7 $(u, v) \leftarrow (v, w)$
- 8 $(u, v) \leftarrow e$
- 9 while $\exists! w \in V(G) : (w, u) \in E(G)$:
- 10 $p \leftarrow ((w, u), p)$
- 11 $(u, v) \leftarrow (w, u)$
- 12 $P \leftarrow P \cup \{p\}$
- 13 return P



Выбор начальных путей

TRIVIAL_PATHS(G)

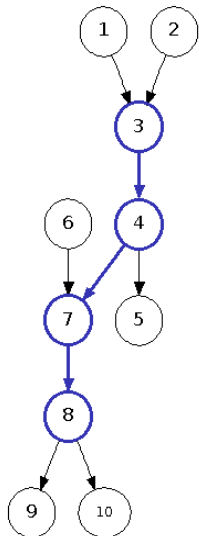
- 1 $P \leftarrow \{\}$
- 2 for each $e \in E(G)$:
- 3 $p \leftarrow (e)$
- 4 $(u, v) \leftarrow e$
- 5 while $\exists! w \in V(G) : (v, w) \in E(G)$:
- 6 $p \leftarrow (p, (v, w))$
- 7 $(u, v) \leftarrow (v, w)$
- 8 $(u, v) \leftarrow e$
- 9 while $\exists! w \in V(G) : (w, u) \in E(G)$:
- 10 $p \leftarrow ((w, u), p)$
- 11 $(u, v) \leftarrow (w, u)$
- 12 $P \leftarrow P \cup \{p\}$
- 13 return P



Выбор начальных путей

TRIVIAL_PATHS(G)

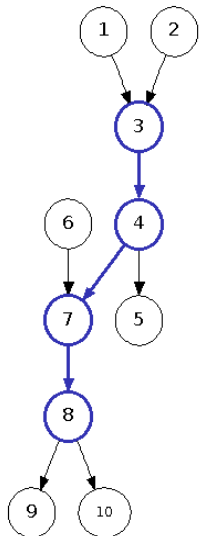
- 1 $P \leftarrow \{\}$
- 2 for each $e \in E(G)$:
- 3 $p \leftarrow (e)$
- 4 $(u, v) \leftarrow e$
- 5 while $\exists! w \in V(G) : (v, w) \in E(G)$:
- 6 $p \leftarrow (p, (v, w))$
- 7 $(u, v) \leftarrow (v, w)$
- 8 $(u, v) \leftarrow e$
- 9 while $\exists! w \in V(G) : (w, u) \in E(G)$:
- 10 $p \leftarrow ((w, u), p)$
- 11 $(u, v) \leftarrow (w, u)$
- 12 $P \leftarrow P \cup \{p\}$
- 13 return P



Выбор начальных путей

TRIVIAL_PATHS(G)

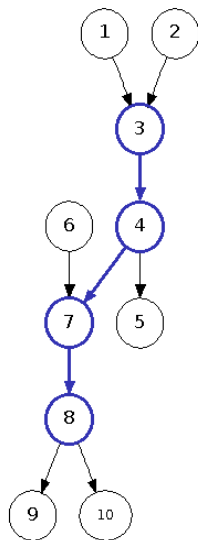
- 1 $P \leftarrow \{\}$
- 2 for each $e \in E(G)$:
- 3 $p \leftarrow (e)$
- 4 $(u, v) \leftarrow e$
- 5 while $\exists! w \in V(G) : (v, w) \in E(G)$:
- 6 $p \leftarrow (p, (v, w))$
- 7 $(u, v) \leftarrow (v, w)$
- 8 $(u, v) \leftarrow e$
- 9 while $\exists! w \in V(G) : (w, u) \in E(G)$:
- 10 $p \leftarrow ((w, u), p)$
- 11 $(u, v) \leftarrow (w, u)$
- 12 $P \leftarrow P \cup \{p\}$
- 13 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

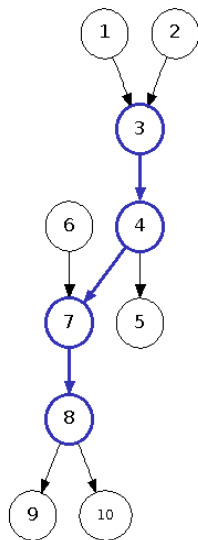
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow \text{LAST}(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(G)\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in p} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists ! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\}:$
- 9 $p \leftarrow (p, h_i)$
- 10 until $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

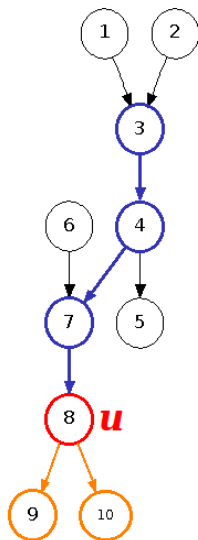
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow \text{LAST}(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(G)\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in P} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists ! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\}:$
- 9 $p \leftarrow (p, h_i)$
- 10 until $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

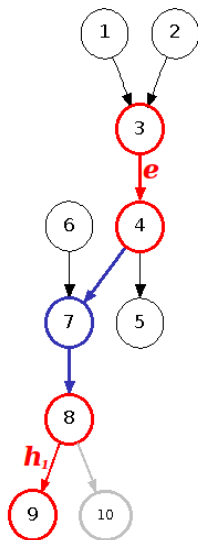
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) | (u, v) \in E(\mathbf{G})\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in P} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists ! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\} :$
- 9 $p \leftarrow (p, h_i)$
- 10 untill $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

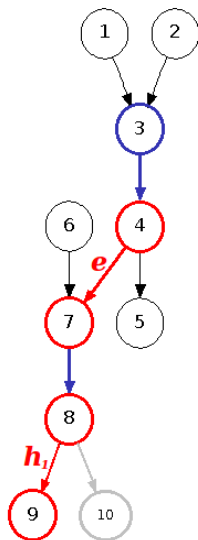
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(\mathbf{G})\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in p} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists ! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\} :$
- 9 $p \leftarrow (p, h_i)$
- 10 until $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

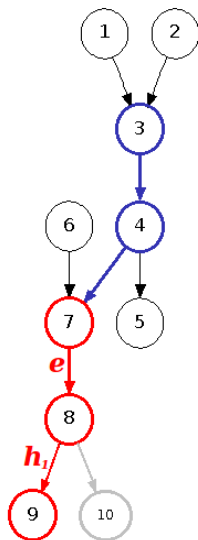
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(\mathbf{G})\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in P} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists ! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\} :$
- 9 $p \leftarrow (p, h_i)$
- 10 until $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

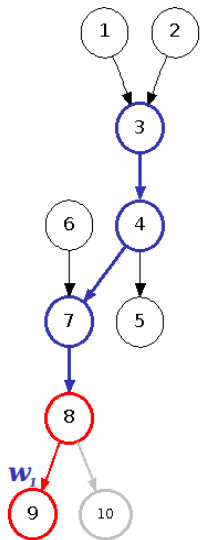
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(\mathbf{G})\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in p} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists ! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\}:$
- 9 $p \leftarrow (p, h_i)$
- 10 until $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

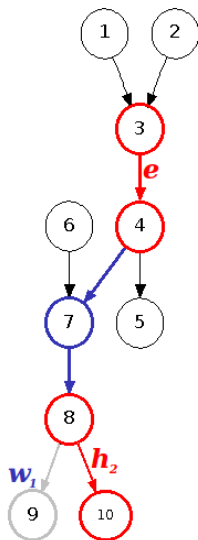
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(\mathbf{G})\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in p} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists ! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\} :$
- 9 $p \leftarrow (p, h_i)$
- 10 untill $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

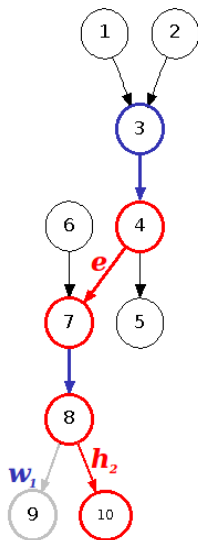
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(\mathbf{G})\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in P} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists ! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\}:$
- 9 $p \leftarrow (p, h_i)$
- 10 until $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

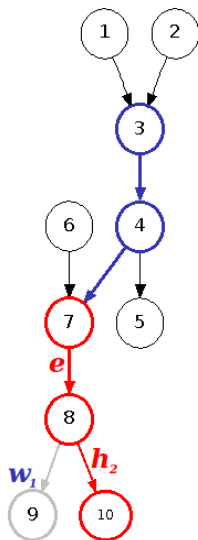
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) | (u, v) \in E(\mathbf{G})\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in P} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists ! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\}:$
- 9 $p \leftarrow (p, h_i)$
- 10 until $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

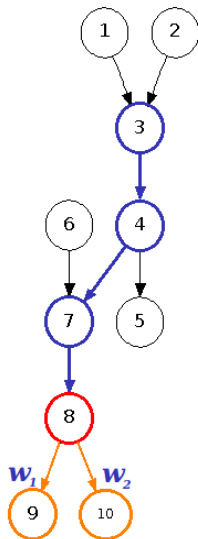
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(G)\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in P} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\}:$
- 9 $p \leftarrow (p, h_i)$
- 10 until $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

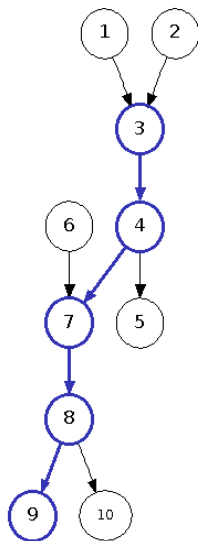
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(\mathbf{G})\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in p} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\}:$
- 9 $p \leftarrow (p, h_i)$
- 10 until $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

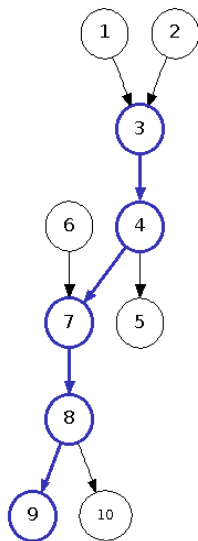
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(\mathbf{G})\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in p} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\}:$
- 9 $p \leftarrow (p, h_i)$
- 10 untill $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

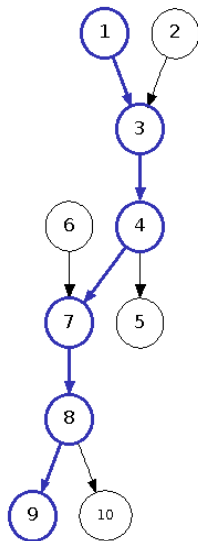
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(\mathbf{G})\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in p} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\}:$
- 9 $p \leftarrow (p, h_i)$
- 10 untill $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

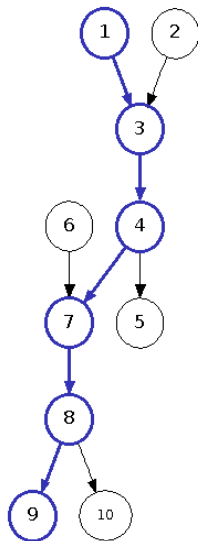
- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(\mathbf{G})\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in p} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\}:$
- 9 $p \leftarrow (p, h_i)$
- 10 until $q = p$
- 11 return P



Удлинение путей

EXTEND(G, P, W, θ, C)

- 1 for each $p \in P$:
- 2 repeat :
- 3 $q \leftarrow p$
- 4 $u \leftarrow LAST(p)$
- 5 $H \leftarrow \{(u, v) \mid (u, v) \in E(\mathbf{G})\}$
- 6 for each $h_i \in H$:
- 7 $w_i \leftarrow \sum_{e \in p} W(e, h_i, D_p(e, h_i))$
- 8 if $\exists! i : w_i > \theta \wedge w_i > C \cdot w_j,$
 $\forall j \in \{1, \dots, |H|\} / \{i\}:$
- 9 $p \leftarrow (p, h_i)$
- 10 untill $q = p$
- 11 return P

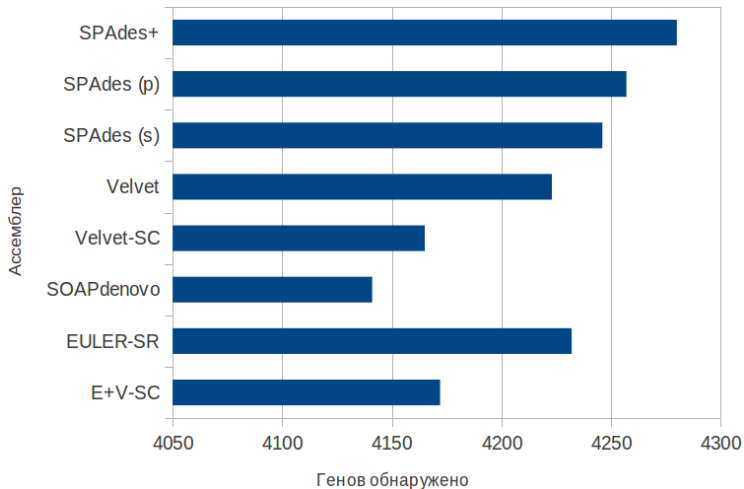


Реализация

- Модуль интегрирован в геномный ассемблер SPAdes
- Реализация включает в себя ряд дополнительных эвристик, учитывающих особенности строения графов и ошибки в начальных данных
- Архитектура проекта позволяет легко добавлять и изменять существующие
 - ▶ Процедуры роста путей
 - ▶ Стратегии выбора продолжений на каждом шаге
 - ▶ Алгоритмы подсчета парной информации

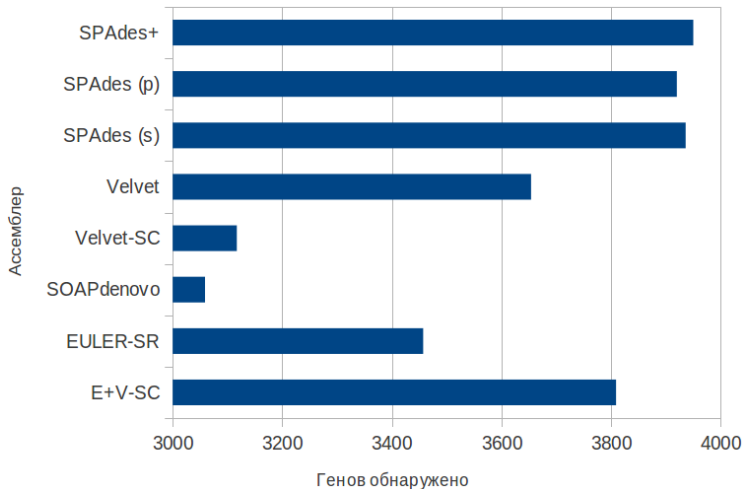
Результаты

- Сравнение по количеству обнаруженных генов
- Геном бактерии *Escherichia coli* (несколько клеток)



Результаты

- Сравнение по количеству обнаруженных генов
- Геном бактерии *Escherichia coli* (одна клетка)



Благодарности



Спасибо за внимание!